

Linear Algebra II

Lectures Notes

MATH 425 Linear Algebra II, Spring 2012
LCD-undergrad 24908; LCD-grad 24909,
MWF 10:00-10:50, Addams Hall 303

Instructor: Shmuel Friedland
Office: 1215 SEO, phone: 413-2159, *e:mail*: friedlan@uic.edu,
web: <http://www.math.uic.edu/~friedlan>

Last update April 10, 2012

1 Basic facts on vector spaces, matrices and linear transformations

1.1 Fields and vector spaces

A commutative group, also called an abelian group, denoted by \mathbf{G} , is a set of elements with a binary operation \oplus , i.e. $a \oplus b \in \mathbf{G}$ for each $a, b \in \mathbf{G}$. This operation is

1. associative: $(a \oplus b) \oplus c = a \oplus (b \oplus c)$;
2. commutative: $a \oplus b = b \oplus a$;
3. there exist a *neutral* element \bigcirc such that $a + \bigcirc = a$ for each $a \in \mathbf{G}$;
4. for each $a \in \mathbf{G}$ there exists a unique $\ominus a$ such that $a + (\ominus a) = \bigcirc$.

Examples of commutative group:

1. The following subset of complex numbers where \oplus is the standard addition $+$, \ominus is $-$ and the neutral element is 0.
 - (a) The set of integers \mathbb{Z} .
 - (b) The set of rational numbers \mathbb{Q} .
 - (c) The set of real numbers \mathbb{R} .
 - (d) The set of complex numbers \mathbb{C}
2. The following subsets of $\mathbb{C}^* := \mathbb{C} \setminus \{0\}$, i.e. all nonzero complex numbers, where the operation \oplus is the *product*, the neutral element is 1, and $\ominus a$ is a^{-1} .
 - (a) $\mathbb{Q}^* := \mathbb{Q} \setminus \{0\}$.
 - (b) $\mathbb{R}^* := \mathbb{R} \setminus \{0\}$.

(c) $\mathbb{C}^* := \mathbb{C} \setminus \{0\}$.

From now on, we will denote the operation \oplus by $+$. An abelian group \mathbf{R} is called a *ring*, if there exist a second operation, called *product*, denoted by ab , which satisfies:

1. associativity: $(ab)c = a(bc)$;
2. distributivity: $a(b + c) = ab + ac$, $(b + c)a = ba + ca$;
3. existence of identity $1 \in \mathbf{R}$: $1a = a1 = a$ for all $a \in \mathbf{R}$.

\mathbf{R} is called a commutative ring if $ab = ba$ for all $a, b \in \mathbf{R}$.

Examples of rings:

1. The following subsets of $n \times n, n > 1$ complex valued matrices, denoted by $\mathbb{C}^{n \times n}$ with the addition $A + B$, multiplication AB , with the identity $I_n, n \times n$ identity matrix.
 - (a) $\mathbb{Z}^{n \times n}$, the ring of $n \times n$ matrices with integer entries. (Noncommutative ring.)
 - (b) $\mathbb{Q}^{n \times n}$, the ring of $n \times n$ matrices with rational entries. (Noncommutative ring.)
 - (c) $\mathbb{R}^{n \times n}$, the ring of $n \times n$ matrices with real entries. (Noncommutative ring.)
 - (d) $\mathbb{C}^{n \times n}$. (Noncommutative ring.)
 - (e) $\mathbf{D}(n, S)$, the set of $n \times n$ diagonal matrices with entries in $S = \mathbb{Z}, \mathbb{Q}, \mathbb{R}, \mathbb{C}$. (Commutative ring.)
2. $S[z]$, the commutative ring of polynomials in one variable z with coefficients in $S = \mathbb{Z}, \mathbb{Q}, \mathbb{R}, \mathbb{C}$.
3. $\mathbb{Z}_m = \mathbb{Z}/(m\mathbb{Z})$, all integers *modulo* a positive integer m , with the addition and multiplication modulo m . $\#\mathbb{Z}_m$, the number of elements in \mathbb{Z}_m , is m . \mathbb{Z}_m can be identified with $\{0, 1, \dots, m - 1\}$.

A commutative ring \mathbf{R} is called a *field* if each nonzero element $a \in \mathbf{R}$ has a unique inverse, denoted by a^{-1} such that $aa^{-1} = 1$. A field is denoted by \mathbb{F} .

Examples of fields:

1. $\mathbb{Q}, \mathbb{R}, \mathbb{C}$.
2. \mathbb{Z}_m , where m is a prime integer.

An abelian group \mathbf{V} is called a vector space over a field \mathbb{F} , if for any $a \in \mathbb{F}, \mathbf{v} \in \mathbf{V}$ the product $a\mathbf{v}$ is an element in \mathbf{V} . This operation satisfies the following properties:

$$a(\mathbf{u} + \mathbf{v}) = a\mathbf{u} + a\mathbf{v}, (a + b)\mathbf{v} = a\mathbf{v} + b\mathbf{v}, (ab)(\mathbf{u} + \mathbf{v}) = a(b(\mathbf{u} + \mathbf{v})), 1\mathbf{v} = \mathbf{v},$$

for all $a, b \in \mathbb{F}, \mathbf{u}, \mathbf{v} \in \mathbf{V}$.

A ring \mathbf{R} is called an algebra over \mathbb{F} if \mathbf{R} is a vector space over \mathbb{F} and $a(\mathbf{x}\mathbf{y}) = (a\mathbf{x})\mathbf{y} = \mathbf{x}(a\mathbf{y})$ for any $a \in \mathbb{F}, \mathbf{x}, \mathbf{y} \in \mathbf{R}$. A standard example of an algebra over \mathbb{F} is the set of $n \times n$ matrices with entries in \mathbb{F} , denoted by $\mathbb{F}^{n \times n}$. Another example is $\mathbb{F}[z]$, the ring of polynomials in one variable z with coefficients in \mathbb{F} .

A set \mathbf{H} is called a group, if there exists a binary operation $ab \in \mathbf{H}$ for each $a, b \in \mathbf{H}$ such that

1. associative: $(ab)c = a(bc)$;
2. identity: $1a = a1 = a$;
3. inverse: $a^{-1}a = aa^{-1} = 1$ for each $a \in \mathbf{H}$.

\mathbf{H} is abelian if $ab = ba$ for each $a, b \in \mathbf{H}$. This is the case discussed in the beginning of this section for \mathbf{G} , where the product operation is replaced by \oplus and the identity by the neutral element.

A standard example of noncommutative group is $\text{GL}(n, \mathbb{C}) \subset \mathbb{C}^{n \times n}, n > 1$, the group of $n \times n$ complex valued invertible matrices.

1.2 Dimension and basis

Let \mathbf{V} be a vector space over a field \mathbb{F} . (For simplicity of the exposition you may assume that $\mathbb{F} = \mathbb{R}$ or $\mathbb{F} = \mathbb{C}$.) $\mathbf{U} \subset \mathbf{V}$ is called a subspace of \mathbf{V} , if $a_1\mathbf{u}_1 + a_2\mathbf{u}_2 \in \mathbf{U}$ for each $\mathbf{u}_1, \mathbf{u}_2 \in \mathbf{U}, a_1, a_2 \in \mathbb{F}$. The subspace $\mathbf{U} := \{\mathbf{0}\}$ is called the zero, or trivial, subspace of \mathbf{V} . Let $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbf{V}$ be n vectors. Then for any $a_1, \dots, a_n \in \mathbb{F}$, called scalars, $a_1\mathbf{x}_1 + \dots + a_n\mathbf{x}_n$ is called a linear combination of $\mathbf{x}_1, \dots, \mathbf{x}_n$. $\mathbf{0} = \sum_{i=1}^n 0\mathbf{x}_i$ is the trivial combination. $\mathbf{x}_1, \dots, \mathbf{x}_n$ are *linearly independent*, abbreviated to l.i. or lin. ind., if the equality $\mathbf{0} = \sum_{i=1}^n a_i\mathbf{x}_i$ implies that $a_1 = \dots = a_n = 0$. Otherwise $\mathbf{x}_1, \dots, \mathbf{x}_n$ are *linearly dependent*, abbreviated to l.d. or lin. dep.. The set of all linear combination of $\mathbf{x}_1, \dots, \mathbf{x}_n$ is called the span of $\mathbf{x}_1, \dots, \mathbf{x}_n$, and denoted by $\text{span}(\mathbf{x}_1, \dots, \mathbf{x}_n)$. $\text{span}(\mathbf{x}_1, \dots, \mathbf{x}_n)$ is a subspace of \mathbf{V} .

1.2.1 More details on vector space - January 13, 2012

\mathbf{V} is called *finitely generated* if $\mathbf{V} = \text{span}(\mathbf{x}_1, \dots, \mathbf{x}_n)$ for some vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$. $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ is a *spanning set* of \mathbf{V} . (In these course we consider only finite generated vector spaces. We will show that every finitely generated vector space is finite dimensional). We will use the following notation for a positive integer n :

$$[n] := \{1, 2, \dots, n-1, n\} \quad (1.1)$$

Observe the following "obvious" relation for any $n \in \mathbb{N}$. (Here \mathbb{N} is the set of all positive integers $1, 2, \dots$.)

$$\text{span}(\mathbf{u}_1, \dots, \mathbf{u}_{i-1}, \mathbf{u}_{i+1}, \dots, \mathbf{u}_n) \subseteq \text{span}(\mathbf{u}_1, \dots, \mathbf{u}_n), \text{ for each } i \in [n]. \quad (1.2)$$

It is enough to consider the case $i = n$. (We assume here that for $n = 1$ the span of empty set is the zero subspace $\mathbf{U} = \{\mathbf{0}\}$.) So for $n = 1$ spaninc holds. Suppose that $n > 1$. Clearly, any linear combination of $\mathbf{u}_1, \dots, \mathbf{u}_{n-1}$, which is $\mathbf{u} = \sum_{i=1}^{n-1} b_i\mathbf{u}_i$ is a linear combination of $\mathbf{u}_1, \dots, \mathbf{u}_n$: $\mathbf{u} = 0\mathbf{u}_n + \sum_{i=1}^{n-1} b_i\mathbf{u}_i$. Hence (1.2) holds.

Theorem 1.1 Suppose that $n > 1$. Then there exists $i \in [n]$ such that

$$\text{span}(\mathbf{u}_1, \dots, \mathbf{u}_n) = \text{span}(\mathbf{u}_1, \dots, \mathbf{u}_{i-1}, \mathbf{u}_{i+1}, \dots, \mathbf{u}_n) \quad (1.3)$$

if and only if $\mathbf{u}_1, \dots, \mathbf{u}_n$ are linearly dependent vectors.

Proof. Suppose that (1.3) holds. So $\mathbf{u}_i = b_1 \mathbf{u}_1 + \dots + b_{i-1} \mathbf{u}_{i-1} + b_{i+1} \mathbf{u}_{i+1} + \dots + b_n \mathbf{u}_n$. Hence $\sum_{i=1}^n a_i \mathbf{u}_i = \mathbf{0}$, where $a_j = b_j$ for $j \neq i$ and $a_i = -1$. Hence $\mathbf{u}_1, \dots, \mathbf{u}_n$ are linearly dependent. Assume that $\mathbf{u}_1, \dots, \mathbf{u}_n$ are linearly dependent. Then there exist scalars $a_1, \dots, a_n \in \mathbb{F}$ not all of them are equal to zero so that $a_1 \mathbf{u}_1 + \dots + a_n \mathbf{u}_n = \mathbf{0}$. Assume that $a_i \neq 0$. For simplicity of notation, (or by renaming indices in $[n]$), we may assume that $i = n$. So $\mathbf{u}_n = -\frac{1}{a_n} \sum_{i=1}^{n-1} a_i \mathbf{u}_i$. Let $\mathbf{u} \in \text{span}(\mathbf{u}_1, \dots, \mathbf{u}_n)$. So

$$\mathbf{u} = \sum_{i=1}^n b_i \mathbf{u}_i = b_n \mathbf{u}_n + \sum_{i=1}^{n-1} b_i \mathbf{u}_i = -b_n \sum_{i=1}^{n-1} \frac{a_i}{a_n} \mathbf{u}_i + \sum_{i=1}^{n-1} b_i \mathbf{u}_i = \sum_{i=1}^{n-1} \frac{a_n b_i - a_i b_n}{a_n} \mathbf{u}_i.$$

That is $\mathbf{u} \in \text{span}(\mathbf{u}_1, \dots, \mathbf{u}_{n-1})$. This proves the theorem. \square

Corollary 1.2 Let $\mathbf{u}_1, \dots, \mathbf{u}_n \in \mathbf{V}$. Assume that not all \mathbf{u}_i are zero vectors. Then there exists $d \geq 1$ integers $1 \leq i_1 < i_2 < \dots < i_d \leq n$ such that $\mathbf{u}_{i_1}, \dots, \mathbf{u}_{i_d}$ are linearly independent and $\text{span}(\mathbf{u}_{i_1}, \dots, \mathbf{u}_{i_d}) = \text{span}(\mathbf{u}_1, \dots, \mathbf{u}_n)$.

Proof. Suppose that $\mathbf{u}_1, \dots, \mathbf{u}_n$ are linearly independent. Then $d = n$ and $i_k = k$ for $k = 1, \dots, n$ and we are done.

Assume that $\mathbf{u}_1, \dots, \mathbf{u}_n$ are linearly dependent. Apply Theorem 1.1. So consider now $n - 1$ vectors $\mathbf{u}_1, \dots, \mathbf{u}_{i-1}, \mathbf{u}_{i+1}, \dots, \mathbf{u}_n$ as given by Theorem 1.1. Note that it is not possible that all vectors in $\{\mathbf{u}_1, \dots, \mathbf{u}_{i-1}, \mathbf{u}_{i+1}, \dots, \mathbf{u}_n\}$ are zero since this will imply that $\mathbf{u}_i = \mathbf{0}$. This will contradict the assumption that not all \mathbf{u}_i are zero. Apply the previous arguments to $\mathbf{u}_1, \dots, \mathbf{u}_{i-1}, \mathbf{u}_{i+1}, \dots, \mathbf{u}_n$ and continue in this fashion until one gets d linearly independent vectors $\mathbf{u}_{i_1}, \dots, \mathbf{u}_{i_d}$. \square

Corollary 1.3 Let \mathbf{V} be a finitely generated nontrivial vectors space, i.e. contains more than one element. Then there exist $n \in \mathbb{N}$ and n linearly independent vectors $\mathbf{v}_1, \dots, \mathbf{v}_n$ such that $\mathbf{V} = \text{span}(\mathbf{v}_1, \dots, \mathbf{v}_n)$.

In the following lemma we use the fact that any m homogeneous linear equations with n variables has a nontrivial solution if $m < n$. (This will be proved later using REF of matrices.)

Lemma 1.4 Let $n > m \geq 1$ be integers. Then any $\mathbf{w}_1, \dots, \mathbf{w}_n \in \text{span}(\mathbf{u}_1, \dots, \mathbf{u}_m)$ are linearly dependent

Proof. Observe that $\mathbf{w}_j = \sum_{i=1}^m a_{ij} \mathbf{u}_i$. So

$$\sum_{j=1}^n x_j \mathbf{w}_j = \sum_{j=1}^n x_j \sum_{i=1}^m a_{ij} \mathbf{u}_i = \sum_{i=1}^m \left(\sum_{j=1}^n a_{ij} x_j \right) \mathbf{u}_i.$$

Consider the following m homogeneous equations in n unknowns: $\sum_{j=1}^n a_{ij} x_j = 0$ for $i = 1, \dots, m$. Since $n > m$ we have a nontrivial solution $(x_1, \dots, x_n)^\top$. So $\sum_{j=1}^n x_j \mathbf{w}_j = \mathbf{0}$. \square

Theorem 1.5 Let $\mathbf{V} = \text{span}(\mathbf{v}_1, \dots, \mathbf{v}_n)$ and assume that $\mathbf{v}_1, \dots, \mathbf{v}_n$ are linearly independent. Then the following holds.

1. Any vector \mathbf{u} can be expressed as a unique linear combination of $\mathbf{v}_1, \dots, \mathbf{v}_n$.
2. For an integer $N > n$ any N vectors in \mathbf{V} are linearly independent.
3. Assume that $\mathbf{u}_1, \dots, \mathbf{u}_n$ are linearly independent. Then $\mathbf{V} = \text{span}(\mathbf{u}_1, \dots, \mathbf{u}_n)$.

Proof. Assume that $\sum_{i=1}^n x_i \mathbf{v}_i = \sum_{i=1}^n y_i \mathbf{v}_i$. So $\sum_{i=1}^n (x_i - y_i) \mathbf{v}_i = \mathbf{0}$. As $\mathbf{v}_1, \dots, \mathbf{v}_n$ are linearly independent it follows that $x_i - y_i = 0$ for $i = 1, \dots, n$. Hence 1 holds.

Lemma 1.4 implies 2.

Suppose that $\mathbf{u}_1, \dots, \mathbf{u}_n$ are linearly independent. Let $\mathbf{v} \in \mathbf{V}$ and consider $n+1$ vectors $\mathbf{u}_1, \dots, \mathbf{u}_n, \mathbf{v}$. 2 implies that $\mathbf{u}_1, \dots, \mathbf{u}_n, \mathbf{v}$ are linearly dependent. Thus there exists $n+1$ scalars a_1, \dots, a_{n+1} , not all of them zero, such that $a_{n+1} \mathbf{v} + \sum_{i=1}^n a_i \mathbf{u}_i = \mathbf{0}$. Assume first that $a_{n+1} = 0$. Then $\sum_{i=1}^n a_i \mathbf{u}_i = \mathbf{0}$. Since not all a_i zero it follows that $\mathbf{u}_1, \dots, \mathbf{u}_n$ linearly dependent, which is a contradiction to our assumption. Hence $a_{n+1} \neq 0$ and $\mathbf{v} = \sum_{i=1}^n \frac{-a_i}{a_{n+1}} \mathbf{u}_i$. \square

Definition 1.6 $\mathbf{v}_1, \dots, \mathbf{v}_n$ is called a spanning set of a linear space \mathbf{V} if $\mathbf{V} = \text{span}(\mathbf{v}_1, \dots, \mathbf{v}_n)$.

1.2.2 Dimension

The dimension of a trivial vector space, consisting of the zero vector $\mathbf{V} = \{\mathbf{0}\}$ is zero. The dimension of a finite dimensional nonzero vector space \mathbf{V} is the number of vectors in any spanning linearly independent set. The dimension of \mathbf{V} is denoted by $\dim \mathbf{V}$. Assume that $\dim \mathbf{V} = n$. Suppose that $\mathbf{V} = \text{span}(\mathbf{x}_1, \dots, \mathbf{x}_n)$. Then $\mathbf{x}_1, \dots, \mathbf{x}_n$ is a *basis* of \mathbf{V} . I.e., each vector \mathbf{x} can be uniquely expressed as $\sum_{i=1}^n a_i \mathbf{x}_i$. Thus to each $\mathbf{x} \in \mathbf{V}$ one corresponds a unique column vector $\mathbf{a} = (a_1, \dots, a_n)^T \in \mathbb{F}^n$, where \mathbb{F}^n is the vector space on column vectors with n coordinates in \mathbb{F} . The correspondence $\mathbf{x} \mapsto \mathbf{a}$ is an isomorphism of $\iota : \mathbf{V} \rightarrow \mathbb{F}^n$. It will be convenient to denote $\mathbf{x} = \sum_{i=1}^n a_i \mathbf{x}_i$ be the equality $\mathbf{x} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n] \mathbf{a}$ ($= \sum_{i=1}^n \mathbf{x}_i a_i$). (Note that we use the standard way to multiply row by column.)

Assume now that $\mathbf{y}_1, \dots, \mathbf{y}_n$ be n vectors in \mathbf{V} . Then

$$\mathbf{y}_i = \sum_{j=1}^n y_{ji} \mathbf{x}_j, \quad i = 1, \dots, n.$$

Denote by $Y = [y_{ji}]_{i=j=1}^n$ the $n \times n$ matrix with the element y_{ji} in j -th row and i -th column. The above equalities are equivalent to the identity

$$[\mathbf{y}_1, \dots, \mathbf{y}_n] = [\mathbf{x}_1, \dots, \mathbf{x}_n] Y. \quad (1.4)$$

Then $\mathbf{y}_1, \dots, \mathbf{y}_n$ is a basis if and only if Y is invertible matrix. (See for details next subsection.) Y is called the matrix of the change of basis from $[\mathbf{y}_1, \dots, \mathbf{y}_n]$ to $[\mathbf{x}_1, \dots, \mathbf{x}_n]$. The matrix of the change of basis from $[\mathbf{x}_1, \dots, \mathbf{x}_n]$ to $[\mathbf{y}_1, \dots, \mathbf{y}_n]$ is given by Y^{-1} . (Just multiply the identity (1.4) by Y^{-1} from the right.) Suppose also that $[\mathbf{z}_1, \dots, \mathbf{z}_n]$ is a basis in \mathbf{V} . Let $[\mathbf{z}_1, \dots, \mathbf{z}_n] = [\mathbf{y}_1, \dots, \mathbf{y}_n] Z$. Then $[\mathbf{z}_1, \dots, \mathbf{z}_n] = [\mathbf{x}_1, \dots, \mathbf{x}_n] (YZ)$. (Use (1.4).)

1.3 Matrices

1.3.1 Basic properties of matrices

The set of $m \times n$ matrices with entries in \mathbb{F} is denoted by $\mathbb{F}^{m \times n}$. $\mathbb{F}^{m \times n}$ is a vector space over \mathbb{F} of dimension mn . A standard basis of $\mathbb{F}^{m \times n}$ is given by E_{ij} , $i = 1, \dots, m, j = 1, \dots, n$ where E_{ij} is a matrix with 1 in the entry (i, j) and 0 in all other entries. $\mathbb{F}^{m \times 1} = \mathbb{F}^m$ the set of column vectors with m coordinates. A standard basis of \mathbb{F}^m is $\mathbf{e}_i = (\delta_{1i}, \dots, \delta_{mi})^\top$, $i = 1, \dots, m$.

For $A \in \mathbb{F}^{m \times n}$, $B \in \mathbb{F}^{p \times q}$ the product AB is defined if and only if the number of columns in A is equal to the number of columns in B , i.e. $n = p$. In that case the resulting matrix $C = [c_{ij}]$ is $m \times q$. The entry c_{ij} is obtained by multiplying row i of A by the column j of B . Recall that for $\mathbf{x} = (x_1, \dots, x_n)^\top$, $\mathbf{y} = (y_1, \dots, y_n)^\top \in \mathbb{F}^n$ $\mathbf{x}^\top \mathbf{y} = \sum_{i=1}^n x_i y_i$. (This product can be regarded as a product of $1 \times n$ matrix \mathbf{x}^\top with $n \times 1$ product matrix \mathbf{y} .) The product of corresponding sizes of matrices satisfies the properties:

1. associativity $(AB)C = A(BC)$;
2. distributivity $(A_1 + A_2)B = A_1B + A_2B$, $A(B_1 + B_2)$;
3. $a(AB) = (aA)B = A(aB)$ for each $a \in \mathbb{F}$;
4. identities $I_m A = AI_n$, where $A \in \mathbb{F}^{m \times n}$ and $I_m = [\delta_{ij}]_{i=j=1}^m \in \mathbb{F}^{m \times m}$ is the identity matrix.

For $A = [a_{ij}] \in \mathbb{F}^{m \times n}$ denote by $A^\top \in \mathbb{F}^{n \times m}$ the transposed matrix of A . So the (i, j) entry of A^\top is the (j, i) entry of A . The following properties hold

$$(aA)^\top = aA, (A + B)^\top = A^\top + B^\top, (AC)^\top = C^\top A^\top.$$

$A = [a_{ij}] \in \mathbb{F}^{m \times n}$ is called diagonal if $a_{ij} = 0$ for $i \neq j$. Denote by $\mathbf{D}(m, n) \subset \mathbb{F}^{m \times n}$ the vector subspace of diagonal matrices. A square diagonal matrix with the diagonal entries d_1, \dots, d_n is denoted by $\text{diag}(d_1, \dots, d_n) \in \mathbb{F}^{n \times n}$. A is called a block matrix if $A = [A_{ij}]_{i=j=1}^{p,q}$ where each entry A_{ij} is an $m_i \times m_j$ matrix. So $A \in \mathbb{F}^{m \times n}$ where $m = \sum_{i=1}^p m_i$, $n = \sum_{j=1}^q n_j$. A block matrix A is called block diagonal if $A_{ij} = 0_{m_i \times m_j}$ for $i \neq j$. A block diagonal matrix with $p = q$ is denoted by $\text{diag}(A_{11}, \dots, A_{pp})$. A different notation is $\bigoplus_{l=1}^p A_{ll} := \text{diag}(A_{11}, \dots, A_{pp})$.

1.3.2 Elementary row operation

Elementary operations on the rows of $A \in \mathbb{F}^{m \times n}$ are defined as follows:

1. Multiply row i by a nonzero $\mathbf{a} \in \mathbb{F}$.
2. Interchange two distinct rows i and j in A .
3. Add to row j row i multiplied by a , where $i \neq j$.

$A \in \mathbb{F}^{m \times n}$ can be brought to a unique canonical form, called the *reduced row echelon form*, abbreviated as RREF by the elementary row operation. The RREF of zero matrix $0_{m \times n}$ is $0_{m \times n}$. For $A \neq 0_{m \times n}$, RREF of A given by B is of the following form. the first k rows of B are nonzero rows. The last $m - k$ rows are zero rows.

The number of nonzero rows is called the rank of A , and is denoted by $\text{rank } A$. The first nonzero element in each i -th row, $i \leq \text{rank } A$, located at the column j_i is 1. It is called a *pivot*. All other elements of B , in the same columns as the pivot, are zero elements. Furthermore the following condition hold $1 \leq j_1 < j_2 < \dots < j_k \leq n$.

1.3.3 Invertible matrices

$A \in \mathbb{F}^{m \times m}$ is called *invertible* if there exists $B \in \mathbb{F}^{m \times m}$ such that $AB = BA = I_m$. Note that B is unique. If $AC = CA = I_m$ then $B = BI_m = B(AC) = (BA)C = I_m C = C$. So we denote B , the inverse of A by A^{-1} . Denote by $\text{GL}(m, \mathbb{F}) \subset \mathbb{F}^{m \times m}$ the set of all invertible matrices. Note that $\text{GL}(m, \mathbb{F})$ is a group under the multiplication, with the unit I_n . (Observe that for $A, B \in \text{GL}(m, \mathbb{F})$ $(AB)^{-1} = B^{-1}A^{-1}$.)

A matrix $E \in \mathbb{F}^{m \times m}$ is called *elementary* if it is obtained from I_m by applying an appropriate elementary row operation. Note that applying an elementary row operation on A is equivalent to EA , where E is the corresponding elementary row operation. By reversing the corresponding elementary operation we see that E is invertible, and E^{-1} is also an elementary matrix.

Theorem 1.7 *Let $A \in \mathbb{F}^{m \times m}$. The following are equivalent.*

1. A is invertible.
2. $A\mathbf{x} = \mathbf{0}$ has only the trivial solution.
3. The RREF of A is I_m .
4. A is a product of elementary matrices. A matrix $A \in \text{GL}(m, \mathbb{F})$ if and only if I_m is its reduced row

Proof. $1 \Rightarrow 2$. $A\mathbf{x} = \mathbf{0}$ implies that $A^{-1}(A\mathbf{x}) = I_m\mathbf{x} = \mathbf{x} = A^{-1}\mathbf{0} = \mathbf{0}$.
 $2 \Rightarrow 3$. The system $A\mathbf{x} = \mathbf{0}$ does free variables, hence the number of pivots is the number of row, i.e. RREF of A is n .
 $3 \Rightarrow 4$. There exists a sequence of elementary matrices so that $E_p E_{p-1} \dots E_1 A = I_m$. Hence $A = E_1^{-1} E_2^{-1} \dots E_p^{-1}$, and each E_j^{-1} is also elementary.
 $4 \Rightarrow 1$. As each elementary matrix is invertible so is their product, which is equal to A . \square

Theorem 1.8 *Let $A \in \mathbb{F}^{n \times n}$. Define the matrix $B = [A \ I_n] \in \mathbb{F}^{n \times (2n)}$. Let $C = [C_1 \ C_2], C_1, C_2 \in \mathbb{F}^{n \times n}$ be the RREF of B . Then A is invertible if and only if $C_1 = I_n$. Furthermore, if $C_1 = I_n$ then $A^{-1} = C_2$.*

Proof. The fact that A is invertible if and only if $C_1 = I_n$ follows straightforward from Theorem 1.7. Note that for any matrix $F = [F_1 \ F_2], F_1, F_2 \in \mathbb{F}^{n \times p}, G \in \mathbb{F}^{l \times m}$ we have $GF = [(GF_1) \ (GF_2)]$. Let H be a product of elementary matrices such that $HB = [IC_2]$. So $HA = I_n, HI_n = C_1$. Hence $H = A^{-1}$ and $C_1 = A^{-1}$. \square

1.3.4 Row and column spaces of A

$A, B \in \mathbb{F}^{m \times n}$ are called *left equivalent* if A and B have the same reduced row echelon form, or B can be obtained from A by applying elementary row operations. Equivalently, $B = UA$ for some $U \in \text{GL}(m, \mathbb{F})$. $A, B \in \mathbb{F}^{m \times n}$ are called *equivalent* if $B = UAV$ for some $U \in \text{GL}(m, \mathbb{F}), V \in \text{GL}(n, \mathbb{F})$. Equivalently, B is equivalent to A , if B can be obtained from A by applying elementary row and column operation. A is equivalent to a block diagonal matrix $I_{\text{rank } A} \oplus 0$.

Let $A \in \mathbb{F}^{m \times n}$. Denote by $\mathbf{c}_1, \dots, \mathbf{c}_n$ the n column of A . We write $A = [\mathbf{c}_1 \ \mathbf{c}_2 \ \dots \ \mathbf{c}_n]$. Then $\text{span}(\mathbf{c}_1, \dots, \mathbf{c}_n)$ is called the column space of A . Its dimension is $\text{rank } A$. Similarly, the dimension of the columns space of A^\top , which is equal to the dimension of the row space of A , is $\text{rank } A$. I.e. $\text{rank } A^\top = \text{rank } A$. Let $\mathbf{x} = (x_1, \dots, x_n)^\top$. Then $A\mathbf{x} = \sum_{i=1}^n x_i \mathbf{c}_i$. Hence the system $A\mathbf{x} = \mathbf{b}$ is solvable if and only if \mathbf{b} is in the column space of A . The set of all $\mathbf{x} \in \mathbb{F}^n$ such that $A\mathbf{x} = \mathbf{0}$ is called the *null space* of A . It is a subspace of dimension of $n - \text{rank } A$, and is denoted as $N(A) \subset \mathbb{F}^n$. $\dim N(A)$ is called the *nullity* of A , and denoted by $\text{nul } A$. If $A \in \text{GL}(n, \mathbb{F})$ then $A\mathbf{x} = \mathbf{b}$ has the unique solution $\mathbf{x} = A^{-1}\mathbf{b}$.

Theorem 1.9 *Let $A \in \mathbb{F}^{m \times n}$. Assume that $B \in \mathbb{F}^{m \times n}$ is a row echelon form of A . Let $k = \text{rank } A$. Then*

1. *The nonzero rows of B form a basis in the row space of A .*
2. *Assume that the pivots of B are the columns $1 \leq j_1 < \dots < j_k \leq n$, i.e. x_{j_1}, \dots, x_{j_k} are the lead variables in the system $A\mathbf{x} = \mathbf{0}$. Then the columns $\mathbf{c}_{j_1}, \dots, \mathbf{c}_{j_k}$ of A form a basis of the column space of A .*
3. *If $n = \text{rank } A$ then the null space of A consists only of $\mathbf{0}$. Otherwise the null space of A has the following basis. For each free variable x_p let \mathbf{x}_p be the unique solution of $A\mathbf{x} = \mathbf{0}$, where $x_p = 1$ and all other free variables are zero. Then this $n - \text{rank } A$ vectors form a basis in $\text{nul } A$.*

Proof. 1. Note that if $E \in \mathbb{F}^{m \times m}$ then the row space of EA is contained in A , since any row in EA is a linear combination of the rows of A . Since $A = E^{-1}(EA)$ it follows that the row space of A is contained in the row space of EA . Hence the row space of A and EA are the same. Therefore the row space of A and B are the same. Since the last $m - k$ rows of B are zero rows, it follows that the row space of B spanned by the first k rows $\mathbf{b}_1^\top, \dots, \mathbf{b}_k^\top$. We claim that $\mathbf{b}_1^\top, \dots, \mathbf{b}_k^\top$ are linearly independent. Indeed, consider $x\mathbf{b}_1^\top + \dots + x_k\mathbf{b}_k^\top = \mathbf{0}^\top$. Since all the entries below the first pivot in B are zero we deduce that $x_1 = 0$. Continue in the same manner to obtain that $x_2 = 0, \dots, x_k = 0$. So $\mathbf{b}_1^\top, \dots, \mathbf{b}_k^\top$ form a basis in the row space of B and A .

2. We first show that $\mathbf{c}_{j_1}, \dots, \mathbf{c}_{j_k}$ linearly independent. Consider the equality $\sum_{j=1}^k x_j \mathbf{c}_{j_j} = \mathbf{0}$. This is equivalent to the system $A\mathbf{x} = \mathbf{0}$, where $\mathbf{x} = (x_1, \dots, x_n)^\top$ and $x_p = 0$ if x_p is a free variable. Since all free variables are zero, all lead variables are zero, i.e. $x_{j_i} = 0$ for $i = 1, \dots, k$. So $\mathbf{c}_{j_1}, \dots, \mathbf{c}_{j_k}$ linearly independent. It is left to show that \mathbf{c}_p , where x_p is a free variable is a linear combination of $\mathbf{c}_{j_1}, \dots, \mathbf{c}_{j_k}$. Again, consider $A\mathbf{x} = \mathbf{0}$, where each $x_p = 1$ and all other free variables are zero. We have a unique solution to $A\mathbf{x} = \mathbf{0}$, which states that $\mathbf{0} = \mathbf{c}_p + \sum_{i=1}^k x_{j_i} \mathbf{c}_{j_i}$, i.e. $\mathbf{c}_p = \sum_{i=1}^k -x_{j_i} \mathbf{c}_{j_i}$.

3. This follows for the way we write down the general solution of the system $A\mathbf{x} = \mathbf{0}$ in terms of free variables. \square

Remark: Let $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{F}^m$. To find a basis in $\mathbf{U} := \text{span}(\mathbf{x}_1, \dots, \mathbf{x}_n)$ consisting of k vectors in $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ apply Theorem 1.9 to the column space of the matrix $A = [\mathbf{x}_1 \ \mathbf{x}_2 \ \dots \ \mathbf{x}_n]$. To find a nice basis in \mathbf{U} apply Theorem 1.9 to the row space of A^\top .

1.3.5 Special types of matrices

The following are some special subsets of $\mathbb{F}^{n \times n}$.

1. A is symmetric if $A^\top = A$. Denote by $S(n, \mathbb{F})$ the subspace of $n \times n$ symmetric polynomials.
2. A is skew-symmetric, or anti-symmetric, if $A^\top = -A$. Denote by $\mathcal{A}(n, \mathbb{F})$ the subspace of the skew-symmetric matrices.
3. $A = [a_{ij}]$ is called upper triangular if $a_{ij} = 0$ for each $j < i$. Denote by $U(n, \mathbb{F})$ the subspace of upper triangular matrices.
4. $A = [a_{ij}]$ is called lower triangular if $a_{ij} = 0$ for each $j > i$. Denote by $L(n, \mathbb{F})$ the subspace of lower triangular matrices.

Denote $\langle m \rangle := \{1, \dots, m\}$ the set of consecutive integers from 1 to m . Let $\alpha = \{\alpha_1 < \alpha_2 < \dots < \alpha_p\} \subset \langle m \rangle$, $\beta = \{\beta_1 < \beta_2 < \dots < \beta_p\} \subset \langle n \rangle$. Then $A[\alpha, \beta]$ is an $p \times q$ submatrix of A , which is obtained by erasing in A all rows and columns which are not in α and β respectively. Assume that $p < m, q < n$. Denote by $A(\alpha, \beta)$ the $(m - p) \times (n - q)$ submatrix of A , which is obtained by erasing in A all rows and columns which are in α and β respectively. For $i \in \langle m \rangle, j \in \langle n \rangle$ denote by $A(i, j) := A(\{i\}, \{j\})$. Assume that $A \in \mathbb{F}^{n \times n}$. The $A[\alpha, \beta]$ is called a *principal* if and only if $\alpha = \beta$.

1.4 Determinants

1.4.1 The group of bijections on a set \mathcal{X}

Let \mathcal{X}, \mathcal{Y} be two sets of objects. Then $\phi : \mathcal{X} \rightarrow \mathcal{Y}$ is called a mapping. I.e. for each $x \in \mathcal{X}$ $\phi(x)$ is an element in \mathcal{Y} . $\phi : \mathcal{X} \rightarrow \mathcal{X}$ is called the identity map if $\phi(x) = x$ for each $x \in \mathcal{X}$. The identity map is denoted as id or $\text{id}_{\mathcal{X}}$. Let $\psi : \mathcal{Y} \rightarrow \mathcal{Z}$. Then one defines the composition map $\psi \circ \phi : \mathcal{X} \rightarrow \mathcal{Z}$ as follows $(\psi \circ \phi)(x) = \psi(\phi(x))$. $\phi : \mathcal{X} \rightarrow \mathcal{Y}$ is called bijection if there exists $\psi : \mathcal{Y} \rightarrow \mathcal{X}$ such that $\psi \circ \phi = \text{id}_{\mathcal{X}}, \phi \circ \psi = \text{id}_{\mathcal{Y}}$. ψ is denoted as ϕ^{-1} . Denote by $\text{Bi}(\mathcal{X})$ the set of all bijections of \mathcal{X} onto itself. It is easy to show that $\text{Bi}(\mathcal{X})$ forms a group under the composition, with the identity element $\text{id}_{\mathcal{X}}$. See the below problem.

Problems

1. Let \mathcal{X} be any nonempty set. Show that $\text{Bi}(\mathcal{X})$ is a group.
2. Show that if \mathcal{X} has one or two elements then $\text{Bi}(\mathcal{X})$ is a commutative group.

3. Show that if \mathcal{X} has at least three elements then $\text{Bi}(\mathcal{X})$ is not a commutative group.

1.4.2 The permutation group

Denote by \mathcal{S}_n is the group of the bijections of the set $\langle n \rangle := \{1, \dots, n\}$ onto itself. Introduce the following polynomials

$$P_\omega(\mathbf{x}) := \prod_{1 \leq i < j \leq n} (x_{\omega(i)} - x_{\omega(j)}), \mathbf{x} = (x_1, \dots, x_n), \text{ for each } \omega \in \mathcal{S}_n. \quad (1.5)$$

Define

$$\text{sign}(\omega) := \frac{P_\omega(\mathbf{x})}{P_{\text{id}}(\mathbf{x})}. \quad (1.6)$$

Theorem 1.10 *For each $\omega \in \mathcal{S}_n$ $\text{sign}(\omega) \in \{1, -1\}$. The map $\text{sign} : \mathcal{S}_n \rightarrow \{1, -1\}$ is a group homomorphism, i.e.*

$$\text{sign}(\omega \circ \sigma) = \text{sign}(\omega)\text{sign}(\sigma) \text{ for each pair } \omega, \sigma \in \mathcal{S}_n. \quad (1.7)$$

Proof. Clearly, if $\omega(i) < \omega(j)$ then the factor $(x_{\omega(i)} - x_{\omega(j)})$ appear in $P_{\text{id}}(\mathbf{x})$. If $\omega(i) > \omega(j)$ then the factor $-(x_{\omega(i)} - x_{\omega(j)})$ appear in $P_{\text{id}}(\mathbf{x})$. Hence $\text{sign}(\omega) \in \{1, -1\}$. Observe next that

$$\frac{P_{\omega \circ \sigma}(\mathbf{x})}{P_{\text{id}}(\mathbf{x})} = \frac{P_{\omega \circ \sigma}(\mathbf{x})}{P_\sigma(\mathbf{x})} \frac{P_\sigma(\mathbf{x})}{P_{\text{id}}(\mathbf{x})} = \text{sign}(\omega)\text{sign}(\sigma).$$

(To show the equality $\frac{P_{\omega \circ \sigma}(\mathbf{x})}{P_\sigma(\mathbf{x})} = \text{sign}(\omega)$ introduce new variables $y_i = x_{\sigma(i)}$ for $i \in \langle n \rangle$.)

□

$\tau \in \mathcal{S}_n$ is a transposition if there exists a pair of integers $1 \leq i < j \leq n$ such that $\tau(i) = j, \tau(j) = i$. Furthermore $\tau(k) = k$ for $k \neq i, j$. Note that $\tau \circ \tau = \text{id}$, i.e. $\tau^{-1} = \tau$.

Theorem 1.11 *For an integer $n \geq 2$ each $\omega \in \mathcal{S}_n$ is a product of transpositions*

$$\omega = \tau_1 \circ \tau_2 \circ \dots \circ \tau_m. \quad (1.8)$$

The product decomposition is not unique. However the parity of m is unique. More precisely $\text{sign}(\omega) = (-1)^m$.

Proof. We agree that in (1.8) $m = 0$ iff $\omega = \text{id}$. (This is true for any $n \in \mathbb{N}$.) We prove the theorem by induction on n . For $n = 2$ \mathcal{S}_2 consists of two elements id and a unique permutation τ , which satisfies $\tau(1) = 2, \tau(2) = 1$. So $\tau \circ \tau = \text{id}$. In this case the lemma follows straightforward.

Suppose that theorem holds for $n = N \geq 2$ and assume that $n = N + 1$. Let $\omega \in \mathcal{S}_n$. Suppose first that $\omega(n) = n$. So ω can be identified with the bijection $\omega' \in \mathcal{S}_{n-1}$, where $\omega'(i) = \omega(i)$ for $i = 1, \dots, n-1$. Use the induction hypothesis to express ω' as a product of m transposition $\tau'_1 \circ \tau'_2 \circ \dots \circ \tau'_m$ in $c\mathcal{S}_{n-1}$. Clearly each τ'_i extends to a transposition τ_i in \mathcal{S}_n . Hence (1.8) holds.

Suppose now that $\omega(n) = i < n$. Let τ be the transposition that interchange i and n . Let $\omega' = \tau \circ \omega$. Then $\omega'(n) = \tau(\omega(n)) = n$. The previous arguments show that $\omega' = \tau_1 \circ \dots \circ \tau_l$. So $\omega = \tau \circ \omega' = \tau \circ \tau_1 \circ \dots \circ \tau_l$.

It is left to show that the parity of m is unique. First observe that $\text{sign}(\tau) = -1$. See problem 4. Theorem 1.10 yields that $\text{sign}(\omega) = (-1)^m$. Hence the parity of m is fixed. \square

Problems

Show

1. $\text{sign}(\text{id}) = 1$.
2. $\text{sign}(\omega^{-1}) = \text{sign}(\omega)$.
3. τ is called an *elementary transposition* if τ interchanges two consecutive integers $\tau(i) = i + 1, \tau(i + 1) = i$. Show that any nonelementary transposition $\tau(i) = j, \tau(j) = i$, where $i < j - 1$ can be represented as a product of $2(j - i) - 1$ elementary transpositions.
4. $\text{sign}(\tau) = -1$ for any transposition τ .

1.5 Definition and properties of determinants

For $A \in \mathbb{F}^{n \times n}$ we define the determinant of A as

$$\det A = \sum_{\omega \in \mathcal{S}_n} \text{sign}(\omega) a_{1\omega(1)} a_{2\omega(2)} \dots a_{n\omega(n)}. \quad (1.9)$$

Theorem 1.12 *Let $A = [a_{ij}] \in \mathbb{F}^{n \times n}$. Then $\det A$ satisfies the following properties.*

1. $\det A = \det A^\top$.
2. $\det A$ is a multilinear function in rows or columns of A .
3. The determinant of lower triangular or upper triangular matrix is a product of the diagonal entries of A .
4. Let B obtained from A by permuting two rows or columns of A . Then $\det B = -\det A$.
5. If A has two equal rows or columns then $\det A = 0$.
6. A is invertible iff $\det A \neq 0$.
7. Let $A, B \in \mathbb{F}^{n \times n}$. Then $\det AB = \det A \det B$.
8. (Laplace row and column expansion for determinants) For $i, j \in [n]$ denote by $A(i, j) \in \mathbb{F}^{(n-1) \times (n-1)}$ the matrix obtained from A by deleting its i -th row and j -th column. Then

$$\det A = \sum_{j=1}^n a_{ij} (-1)^{i+j} \det A(i, j) = \sum_{i=1}^n a_{ji} (-1)^{i+j} \det A(i, j) \quad (1.10)$$

for $i = 1, \dots, n$.

Proof. 1. In (1.9) note that if $j = \omega(i)$ then $i = \omega^{-1}(j)$. since ω is a bijection, we deduce that when $i = 1, \dots, n$ then j takes each value in $[n]$. Since $\text{sign}(\omega) = \text{sign}(\omega^{-1})$ we obtain Hence

$$\begin{aligned} \det A &= \sum_{\omega \in \mathcal{S}_n} \text{sign}(\omega) a_{\omega^{-1}(1)1} a_{\omega^{-1}(2)2} \cdots a_{\omega^{-1}(n)n} = \\ & \sum_{\omega \in \mathcal{S}_n} \text{sign}(\omega^{-1}) a_{\omega^{-1}(1)1} a_{\omega^{-1}(2)2} \cdots a_{\omega^{-1}(n)n} = \det A^\top. \end{aligned}$$

(Note since \mathcal{S}_n is a group, when ω varies over \mathcal{S}_n so is ω^{-1} .)

2. Fix all rows of A except the row i . From (1.9) it follows that $\det A$ is a linear function in the row i of A . I.e. $\det A$ is a multilinear function in the rows of A . In view of the identity $\det A = \det A^\top$ we deduce that $\det A$ is a multilinear function of the columns of A .

3. Assume that A is upper triangular. Then $a_{n\omega(n)} = 0$ if $\omega(n) \neq n$. So all nonzero terms in (1.9) are zero unless $\omega(n) = n$. So assume that $\omega(n) = n$. Then $a_{(n-1)\omega(n-1)} = 0$ unless $\omega(n-1) = n-1$. So all nonzero terms in (1.9) must come from all ω satisfying $\omega(n) = n, \omega(n-1) = n-1$. Continuing in the same manner we deduce that the only nonzero term in (1.9) comes from $\omega = \text{id}$. As $\text{sign}(\text{id}) = 1$ it follows that $\det A = \prod_{i=1}^n a_{ii}$. Since $\det A = \det A^\top$ we deduce the claim for lower triangular matrix.

4. Note that a permutation of two elements $1 \leq i < j \leq n$ in $[n]$ is achieved by a transposition $\tau \in \mathcal{S}_n$. So $B = [b_{pq}]$ and $b_{pq} = a_{\tau(p)q}$. As in the proof of 1 we let $\tau(i) = j$, then $j = \tau^{-1}(i) = \tau(i)$. Hence

$$\begin{aligned} \det B &= \sum_{\omega \in \mathcal{S}_n} \text{sign}(\omega) a_{\tau(1)\omega(1)} a_{\tau(2)\omega(2)} \cdots a_{\tau(n)\omega(n)} = \\ & \sum_{\omega \in \mathcal{S}_n} \text{sign}(\omega) a_{1\omega(\tau(1))} a_{2\omega(\tau(2))} \cdots a_{n\omega(\tau(n))} = \\ & \sum_{\omega \in \mathcal{S}_n} -\text{sign}(\omega \circ \tau) a_{1(\omega \circ \tau)(1)} a_{2(\omega \circ \tau)(2)} \cdots a_{n(\omega \circ \tau)(n)} = -\det A. \end{aligned}$$

5. Suppose first that we are not in characteristic 2, i.e. $0 \neq 2$. Suppose A has two identical rows. Interchange these two rows to obtain $B = A$. So $\det A = \det B = -\det A$, where the last equality is established in 4. So $2 \det A = 0$. Hence $\det A = 0$.

Suppose now that $2 = 0$, i.e. $1 = -1$. In this case $\det A$ is equal to the *permanent* of A :

$$\text{perm } A = \sum_{\omega \in \mathcal{S}_n} a_{1\omega(1)} a_{2\omega(2)} \cdots a_{n\omega(n)}. \quad (1.11)$$

Since A has two identical rows each term appears twice. (For example if row one is equal to two two them $a_{1i}a_{2j} = a_{1j}a_{2i}$. So we have only $\frac{n!}{2}$ terms and each term is multiplied by $2 = 0$. Hence $\det A = 0$. Use $\det A = \det A^\top$ to deduce that $\det A = 0$ if A has two identical columns.

6. Use 2, 3 and 5 to deduce that $\det EA = \det E \det A$ if E is an elementary matrix. (Note that $\det E \neq 0$.) Hence if E_1, \dots, E_k are elementary matrices we deduce that $\det(E_k \cdots E_1) = \prod_{i=1}^k \det E_i$. Let B be the reduced row echelon for of A . So $B = E_k E_{k-1} \cdots E_1 A$. Hence $\det B = (\prod_{i=1}^k \det E_i) \det A$. If $I_n \neq B$, then the

last row of B is zero, so $\det B = 0$ which implies that $\det A = 0$. If $B = I_n$ then $\det A = \prod_{i=1}^n (\det E_i)^{-1}$.

7. Assume that either $\det A = 0$ or $\det B = 0$. We claim that $(AB)\mathbf{x} = \mathbf{0}$ has a nontrivial solution. Suppose $\det B = 0$. 6 and Theorem 1.7 yields that $B\mathbf{x} = \mathbf{0}$ has a nontrivial solution which satisfies $AB\mathbf{x} = \mathbf{0}$. Suppose that B is invertible and $A\mathbf{y} = \mathbf{0}$ for some $\mathbf{y} \neq \mathbf{0}$. Then $AB(B^{-1}\mathbf{x}) = \mathbf{0}$ which implies that $\det AB = 0$. Hence in these cases $\det AB = 0 = \det A \det B$. Suppose that A, B are invertible. Then each of them is a product of elementary matrices. Use the arguments in the proof of 6 to deduce that $\det AB = \det A \det B$.

8. First we prove the first part of the formula 1.10 for $i = n$. Clearly, (1.9) yield the equality

$$\det A = \sum_{j=1}^n a_{nj} \sum_{\omega \in \mathcal{S}_n, \omega(n)=j} \text{sign}(\omega) a_{1\omega(1)} \cdots a_{(n-1)\omega((n-1))}. \quad (1.12)$$

In the above sum let $j = n$. So $\omega(n) = n$. So ω can be viewed as $\omega' \in \mathcal{S}_{n-1}$. Also $\text{sign}(\omega) = \text{sign}(\omega')$. Hence $\sum_{\omega' \in \mathcal{S}_{n-1}} \text{sign}(\omega) a_{1\omega'(1)} \cdots a_{(n-1)\omega'(n-1)} = \det A(n, n)$. Note that $(-1)^{n+n} = 1$. This justify the form of the last term of the expansion 1.10 for $i = n$. To justify the sign of any term in 1.10 for $i = n$ we take the column row and interchange it first with column $j + 1$, then with column $j + 2$, and at last with the column n . The sign of $\det A(n, j)$ is $(-1)^{n+j}$. This proves the case $i = n$.

By interchanging any row $i < n$ with row $i + 1$, row $i + 2$, and finally with row n we deduce the first part of the formula 1.10 for any i . By considering $\det A^\top$ we deduce the second part of the formula 1.10.

□

Observe next that $\det I_n = 1$. Hence

$$1 = \det I_n = \det(AA^{-1}) = \det A \det A^{-1} \Rightarrow \det(A^{-1}) = \frac{1}{\det A}.$$

For $A = [a_{ij}] \in \mathbb{F}^{n \times n}$ denote by A_{ij} the determinant of the matrix obtained from A by deleting i -th row and j -th column and multiplied by $(-1)^{i+j}$. (This is called the (i, j) cofactor of A .)

$$A_{ij} := (-1)^{i+j} \det A(i, j). \quad (1.13)$$

Then the expansion of $\det A$ by the row i and the column j respectively, i.e. (1.12), is given by the equalities

$$\det A = \sum_{j=1}^n a_{ij} A_{ij} = \sum_{i=1}^n a_{ij} A_{ij}. \quad (1.14)$$

Then the *adjoint* matrix of A is defined as follows:

$$\text{adj } A := \begin{bmatrix} A_{11} & A_{21} & \cdots & A_{n1} \\ A_{12} & A_{22} & \cdots & A_{n2} \\ \vdots & \vdots & \vdots & \vdots \\ A_{1n} & A_{2n} & \cdots & A_{nn} \end{bmatrix}. \quad (1.15)$$

Proposition 1.13 Let $A \in \mathbb{F}^{n \times n}$. Then

$$A(\text{adj } A) = (\text{adj } A)A = (\det A)I_n. \quad (1.16)$$

Hence A is invertible if and only if $\det A \neq 0$. Furthermore, $A^{-1} = (\det A)^{-1} \text{adj } A$.

Proof. Consider an (i, k) entry of $A(\text{adj } A)$. It is given as $\sum_{j=1}^n a_{ij} A_{kj}$. For $i = k$ (1.14) yields that $\sum_{j=1}^n a_{ij} A_{ij} = \det A$. Suppose that $i \neq k$. Let B_k be the matrix obtained from A by replacing the row k in of A by the row i . So B_k has two identical rows, hence $\det B_k = 0$. On the other hand, expand B_k by the row k to obtain that $0 = \det B_k = \sum_{j=1}^n a_{ij} A_{kj}$. This shows that $A(\text{adj } A) = (\det A)I_n$. Similarly one shows that $(\text{adj } A)A = (\det A)I_n$.

Theorem 1.12 part 6 shows that A is invertible iff $\det A \neq 0$. Hence for invertible A , $A^{-1} = \frac{1}{\det A} \text{adj } A$. \square

Proposition 1.14 (Cramer's rule.) Let $A \in \text{GL}(n, \mathbb{F})$ and consider the system $A\mathbf{x} = \mathbf{b}$, where $\mathbf{x} = (x_1, \dots, x_n)^\top$. Denote by C_k the matrix obtained from A by replacing the column k in A by \mathbf{b} . Then $x_k = \frac{\det C_k}{\det A}$ for $k = 1, \dots, n$.

Proof. Clearly $\mathbf{x} = A^{-1}\mathbf{b} = \frac{1}{\det A}(\text{adj } A)\mathbf{b}$. Hence $x_k = (\det A)^{-1} \sum_{j=1}^n A_{jk} b_j$, where $\mathbf{b} = (b_1, \dots, b_n)^\top$. Expand B_k by the column k to deduce that $\det B_k = \sum_{j=1}^n b_j A_{jk}$. \square

Let B be a square submatrix of A . Then $\det B$ is called a *minor* of A . $\det B$ is called a principle minor of order m if B is an $m \times m$ principle submatrix of A .

1.6 Polynomials

For a field \mathbb{F} , (usually $\mathbb{F} = \mathbb{R}, \mathbb{C}$), denote by $\mathbb{F}[z]$, the *ring* of polynomials $p(z) = a_0 z^n + a_1 z^{n-1} + \dots + a_n$ with coefficients in \mathbb{F} . The *degree* of p , denoted by $\deg p$, is the maximal degree $n - j$ of a monomial $a_j x^{n-j}$ which is not identically zero, i.e. $a_j \neq 0$. So $\deg p = n$ if and only if $a_0 \neq 0$, the degree of a nonzero constant polynomial $p(z) = a_0$ is zero, and the degree of the zero polynomial is agreed to be equal to $-\infty$. For two polynomials $p, q \in \mathbb{F}[z]$ and two scalars $a, b \in \mathbb{F}$ $ap(z) + bq(z)$ is a well defined polynomial. Hence $\mathbb{F}[z]$ is a vector space over \mathbb{F} , whose dimension is infinite. The set of polynomials of degree n at most, is $n + 1$ dimensional subspace of $\mathbb{F}[z]$. Given two polynomials $p = \sum_{i=0}^n a_i z^{n-i}$, $q = \sum_{j=0}^m b_j z^{m-j} \in \mathbb{F}[z]$ one can form the product

$$p(z)q(z) = \sum_{k=0}^{n+m} \left(\sum_{i=0}^k a_i b_{k-i} \right) z^{n+m-k}, \text{ where } a_i = b_j = 0 \text{ for } i > n \text{ and } j > m.$$

Note that $pq = qp$ and $\deg pq = \deg p + \deg q$. The addition and the product in $\mathbb{F}[z]$ satisfies all the nice distribution identities as the addition and multiplication in \mathbb{F} . Here the constant polynomial $p \equiv 1$ is the identity element, and the zero polynomial as the zero element. (That is the reason for the name *ring* of polynomials in one variable over \mathbb{F} .)

Recall that given two polynomials $p, q \in \mathbb{F}[z]$ one can divide p by $q \neq 0$ with the residue r , i.e. $p = tq + r$ for some unique $t, r \in \mathbb{F}[z]$, where $\deg r < \deg q$. For

$p, q \in \mathbb{F}[z]$ let (p, q) be the *greatest common divisor* of p, q . If p and q are identically zero then (p, q) is the zero polynomial. Otherwise (p, q) is a polynomial s of the highest degree that divides p and q . s is determined up to a multiple of a nonzero scalar. s can be chosen as a unique *monic* polynomial:

$$s(z) = z^l + s_1 z^{l-1} + \dots + s_l \in \mathbb{F}[z]. \quad (1.17)$$

For $p, q \neq 0$ s can be found using the *Euclid* algorithm:

$$p_i(z) = t_i(z)p_{i+1}(z) + p_{i+2}(z), \quad \deg p_{i+2} < \deg p_{i+1} \quad i = 1, \dots \quad (1.18)$$

Start this algorithm with $p_1 = p, p_2 = q$. Continue it until $p_k = 0$ the first time. (Note that $k \geq 3$.) Then $p_{k-1} = (p, q)$. It is easy to show, for example by induction, that each p_i is of the form $u_i p + v_i q$ for some polynomials u_i, v_i . Hence the Euclid algorithm yields

$$(p(z), q(z)) = u(z)p(z) + v(z)q(z), \quad \text{for some } u(z), v(z) \in \mathbb{F}[z]. \quad (1.19)$$

(This formula holds for any $p, q \in \mathbb{F}[z]$.) $p, q \in \mathbb{F}[z]$ are called *coprime* if $(p, q) = 1$.

Recall that if we divide $p(z)$ by $z - a$ we get the residue $p(a)$, i.e. $p(z) = (z - a)q(z) + p(a)$. So $z - a$ divides $p(z)$ if and only if $p(a) = 0$, i.e. a is the root of p . A monic $p(z)$ *splits* to a product of linear factors if

$$p(z) = (z - z_1)(z - z_2) \dots (z - z_n) = \prod_{i=1}^n (z - z_i). \quad (1.20)$$

Note that z_1, \dots, z_n are the *roots* of p .

Let $\mathbf{z} = (z_1, \dots, z_n)^\top \in \mathbb{F}^n$. Denote

$$\sigma_k(\mathbf{z}) := \sum_{1 \leq i_1 < i_2 < \dots < i_k \leq n} z_{i_1} z_{i_2} \dots z_{i_k}, \quad k = 1, \dots, n. \quad (1.21)$$

$\sigma_k(\mathbf{z})$ is called the k -*th elementary symmetric polynomial* in z_1, \dots, z_n . Observe

$$\begin{aligned} \sigma_1(\mathbf{z}) &= z_1 + z_2 + \dots + z_n, \quad n \text{ summands} \\ \sigma_2(\mathbf{z}) &= z_1 z_2 + \dots + z_1 z_n + z_2 z_3 + \dots + z_{n-1} z_n, \quad \frac{n(n-1)}{2} \text{ summands,} \\ \sigma_n(\mathbf{z}) &= z_1 z_2 \dots z_n, \quad n \text{ terms in the product.} \end{aligned}$$

A straightforward calculation shows

$$\prod_{i=1}^n (z - z_i) = z^n + \sum_{i=1}^n (-1)^i \sigma_i(\mathbf{z}) z^{n-i}. \quad (1.22)$$

1.6.1 Finite extension of fields

A polynomial $p(z) \in \mathbb{F}[z]$ is called *irreducible*, if all polynomials q that divide p are either constant nonzero polynomials, or polynomials of the form $ap(z)$, where $a \in \mathbb{F} \setminus \{0\}$. \mathbb{F} is called an *algebraically closed field* if any monic polynomial $p(z) \in \mathbb{F}[z]$ splits to linear factors in $\mathbb{F}[z]$. It is easy to see that \mathbb{F} is algebraically closed if and only if the only irreducible monic polynomials are $z - a$ for all $a \in \mathbb{F}$. So \mathbb{F} is not algebraically closed if and only if there exists an irreducible monic polynomial in $\mathbb{F}[z]$ of degree greater than 1.

Theorem 1.15 *Let \mathbb{F} be a field. Assume that $p(z) = z^d + \sum_{i=1}^d a_i z^{d-i}$ be an irreducible polynomial, where $d > 1$. Denote by $\mathbb{F}[z]/(p(z)\mathbb{F}[z])$ the set of all polynomials modulo $p(z)$. That is $f(z) \equiv g(z)$ if the polynomial $f(z) - g(z)$ is divided by $p(z)$. Then this set is a field, denoted by $\mathbb{F}_{p(z)}$, under the addition and multiplication modulo $p(z)$. $\mathbb{F}_{p(z)}$ is a vector space over \mathbb{F} of dimension d . The set of all constant polynomials in $\mathbb{F}_{p(z)}$ is isomorphic to \mathbb{F} . ($\mathbb{F}_{p(z)}$ is called a finite extension of \mathbb{F} , and more precisely a an extension of degree d .)*

Problems

Show

1. Prove Theorem 1.15.
2. Show that there is only one monic irreducible polynomial of degree two over \mathbb{Z}_2 . Describe the extension of \mathbb{Z}_2 of degree 2.
3. Show that any finite extension of \mathbb{Z}_p , where $p \geq 2$ is prime, has p^d elements.
4. Recall that the characteristic polynomial of $A \in \mathbb{F}^{n \times n}$ is defined as $\det(zI_n - A)$. Show that the characteristic polynomial is a monic polynomial of degree n . Prove that the coefficient of z^{n-k} of the characteristic polynomial of A is the sum of all minors $\det A[\alpha, \alpha]$ where α runs over all subsets of $[n]$ of cardinality k .
5. Recall that $\lambda \in \mathbb{F}$ is an eigenvalue of $A \in \mathbb{F}^{n \times n}$ if there exists $\mathbf{0} \neq \mathbf{x} \in \mathbb{F}^n$ such that $A\mathbf{x} = \lambda\mathbf{x}$. This $\mathbf{x} \neq \mathbf{0}$ is called an eigenvector of A corresponding to λ . Show that λ is an eigenvalue of A if and only if λ is a zero of the characteristic polynomial of A .
6. Find a $A \in \mathbb{Z}_2^{2 \times 2}$ which does not have eigenvalues in \mathbb{Z}_2 .
7. Show that \mathbb{C} is a 2-extension of \mathbb{R} . What is the corresponding irreducible polynomial in $\mathbb{R}[z]$?
8. Let $p \geq 3$ be a prime and consider the polynomial $f_p = \sum_{i=0}^{p-1} z^i \in \mathbb{Q}[z]$. Show that this polynomial is irreducible over $\mathbb{Q}[z]$.

1.7 Complex numbers

\mathbb{C} is the field of complex numbers. A complex number z can be written in the form $z = x + \mathbf{i}y$, where $x, y \in \mathbb{R}$. Here $\mathbf{i}^2 = -1$. Sometimes in this notes we denote \mathbf{i} by $\sqrt{-1}$. So \mathbb{C} can be viewed as \mathbb{R}^2 , where the vector $(x, y)^\top$ represents z . Recall that $x = \Re z$, the real part of z , and $y = \Im z$, the imaginary part of z . $\bar{z} = x - \mathbf{i}y$ is the conjugate of z . Note that $|z| = \sqrt{x^2 + y^2}$, is the absolute value of z or the modulus of z . For $z \neq 0$, the argument of z is defined as $\arctan \frac{y}{x}$. The polar representation of z is $z = r e^{i\theta} = r(\cos \theta + \mathbf{i} \sin \theta)$. Here r and θ are the modulus and the argument of $z (\neq 0)$. Let $w = u + \mathbf{i}v = R(\cos \psi + \mathbf{i} \sin \psi)$, where $u, v \in \mathbb{R}$. Then

$$z + w = (x + u) + \mathbf{i}(y + v), \quad zw = (xu - yv) + \mathbf{i}(xv + yu) = rR e^{i(\theta + \psi)},$$

$$\frac{w}{z} = \frac{1}{z\bar{z}} w\bar{z} = \frac{R}{r} e^{i(\psi - \theta)} \text{ if } z \neq 0.$$

For a complex number $w = Re^{i\psi}$ and a positive integer $n \geq 2$, the equation $z^n - w = 0$ has n -complex distinct roots for $w \neq 0$, which are $R^{\frac{1}{n}} e^{i\frac{\psi+2k\pi}{n}}$ for $k = 0, 1, \dots, n-1$.

The fundamental theorem of algebra states that any monic polynomial $p(z) \in \mathbb{C}[z]$ of degree $n \geq 2$ splits to linear factors, i.e. $p(z) = \prod_{i=1}^n (z - z_i)$.

Problems

Show

1. Show that $p(z) = z^2 + bz + c \in \mathbb{R}[z]$ is irreducible over \mathbb{R} if and only if $b^2 < 4c$.
2. Show that any monic polynomial $p(z) \in \mathbb{R}[z]$ of degree 2 at least splits to a product of irreducible linear and quadratic monic polynomials over $\mathbb{R}[z]$.
3. Deduce from the previous problem that any $p(z) \in \mathbb{R}[z]$ of odd degree must have a real root.

1.8 Linear transformations

Let \mathbf{V}, \mathbf{U} be two finite dimensional subspaces, where $\dim \mathbf{V} = n, \dim \mathbf{U} = m$. $T : \mathbf{V} \rightarrow \mathbf{U}$ is called a linear transformation if $T(a\mathbf{u} + b\mathbf{v}) = aT(\mathbf{u}) + bT(\mathbf{v})$. Assume that $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}, \{\mathbf{u}_1, \dots, \mathbf{u}_m\}$ be two bases in \mathbf{V}, \mathbf{U} respectively. Then T is completely determined by $T(\mathbf{v}_j) = \sum_{i=1}^m a_{ij} \mathbf{u}_i, j = 1, \dots, n$. Let $A = [a_{ij}]_{j=1}^{m,n} \in \mathbb{F}^{m \times n}$. Then the above equality is equivalent to

$$T[\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n] = [T\mathbf{v}_1, T\mathbf{v}_2, \dots, T\mathbf{v}_n] = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_m]A. \quad (1.23)$$

A is called the representation matrix of T . Assume that

$$[\mathbf{y}_1, \dots, \mathbf{y}_n] = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n]Y, [\mathbf{x}_1, \dots, \mathbf{x}_m] = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_m]X, \quad (1.24)$$

where $Y \in \text{GL}(n, \mathbb{F}), X \in \text{GL}(m, \mathbb{F})$, are another bases in \mathbf{V} and \mathbf{U} respectively. Then

$$T[\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n] = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m]X^{-1}AY. \quad (1.25)$$

Denote by $L(\mathbf{V}, \mathbf{U})$ is the set of linear transformations $T : \mathbf{V} \rightarrow \mathbf{U}$. $L(\mathbf{V}, \mathbf{U})$ is a vector space by defining $(aT + bS)(\mathbf{v}) = aT(\mathbf{v}) + bS(\mathbf{v})$. Then $L(\mathbf{V}, \mathbf{U})$ is isomorphic to $\mathbb{F}^{m \times n}$. (Choose bases $[\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n], [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_m]$ in \mathbf{V}, \mathbf{U} respectively, and identify T with its representation matrix.)

Assume that $\dim \mathbf{U} = \dim \mathbf{V}$. $T \in L(\mathbf{V}, \mathbf{U})$ is called an *isomorphism* if it is 1-1, i.e. $T(\mathbf{x}) = T(\mathbf{y}) \Rightarrow \mathbf{x} = \mathbf{y}$ and onto, i.e. $T(\mathbf{V}) = \mathbf{U}$.

Denote by $L(\mathbf{V}) := L(\mathbf{V}, \mathbf{V})$. Any $T \in L(\mathbf{V})$ is represented by a matrix $A \in \mathbb{F}^{n \times n}$ in a basis $[\mathbf{v}_1, \dots, \mathbf{v}_n]$ of \mathbf{V} as follows:

$$T[\mathbf{v}_1, \dots, \mathbf{v}_n] = [v_1, \dots, \mathbf{v}_n], \text{ i.e. } T(\mathbf{v}_j) = \sum_{i=1}^n a_{ij} \mathbf{v}_i \text{ for } j \in [n]. \quad (1.26)$$

Denote by $\mathbf{V}^* := L(\mathbf{V}, \mathbb{F})$. \mathbf{V}^* is the set of all linear functions on \mathbf{V} . \mathbf{V}^* is called the *dual space*.

Problems

Show

1. Any n -dimensional space \mathbf{V} over \mathbb{F} is isomorphic to \mathbb{F}^n .
2. Any finite dimensional space \mathbf{V} over \mathbb{F} is isomorphic to $(\mathbf{V}^*)^*$. Define an explicit isomorphism from \mathbf{V} to $(\mathbf{V}^*)^*$.
3. Any $T \in L(\mathbf{V}, \mathbf{U})$ is an isomorphism if and only if T is represented by an invertible matrix in some bases of \mathbf{V} and \mathbf{U} .
4. Let $T \in L(\mathbf{V}, \mathbf{U})$. What are the conditions on the dimensions of \mathbf{V} and \mathbf{U} that the following holds: T is 1 – 1; T is onto; T is 1 – 1 and onto.
5. $A, B \in \mathbb{F}^{n \times n}$ are called similar if $B = Q^{-1}AQ$ for some $Q \in GL(n, \mathbb{F})$. Show that similarity is an equivalence relation on $\mathbb{F}^{n \times n}$.
6. $A, B \in \mathbb{F}^{n \times n}$ are similar if and only if they represent the same linear transformation $T \in L(\mathbf{V})$ in different bases, for a given n -dimensional vector space \mathbf{V} over \mathbb{F} .
7. If $A, B \in \mathbb{F}^{n \times n}$ are similar then A and B have the same characteristic polynomial. Give an example where the opposite claim does not hold.
8. Suppose that $A, B \in \mathbb{F}^n$ have the same characteristic polynomial, which has n distinct roots in \mathbb{F} . Show that A and B are similar.

2 Inner product spaces

2.1 Inner product

Definition 2.1 Let $\mathbb{F} = \mathbb{R}, \mathbb{C}$ and let \mathbf{V} be a vector space over \mathbb{F} . Then $\langle \cdot, \cdot \rangle : \mathbf{V} \times \mathbf{V} \rightarrow \mathbb{F}$ is called an inner product if the following conditions hold:

- (a) $\langle a\mathbf{x} + b\mathbf{y}, \mathbf{z} \rangle = a\langle \mathbf{x}, \mathbf{z} \rangle + b\langle \mathbf{y}, \mathbf{z} \rangle$, for all $a, b \in \mathbb{F}$, $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathbf{V}$,
- (br) for $\mathbb{F} = \mathbb{R}$ $\langle \mathbf{y}, \mathbf{x} \rangle = \langle \mathbf{x}, \mathbf{y} \rangle$, for all $\mathbf{x}, \mathbf{y} \in \mathbf{V}$;
- (bc) for $\mathbb{F} = \mathbb{C}$ $\langle \mathbf{y}, \mathbf{x} \rangle = \overline{\langle \mathbf{x}, \mathbf{y} \rangle}$, for all $\mathbf{x}, \mathbf{y} \in \mathbf{V}$;
- (c) $\langle \mathbf{x}, \mathbf{x} \rangle > 0$ for all $\mathbf{x} \in \mathbf{V} \setminus \{\mathbf{0}\}$.

$\|\mathbf{x}\| := \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}$ is called the norm (length) of $\mathbf{x} \in \mathbf{V}$.

Other standard properties of inner products are mentioned in Problems 2.2-2.3. We will use the abbreviation IPS for inner product space. In this chapter we assume that $\mathbb{F} = \mathbb{R}, \mathbb{C}$ unless stated otherwise.

Proposition 2.2 Let \mathbf{V} be a vector space over \mathbb{R} . Identify $\mathbf{V}_\mathbb{C}$ with the set of pairs (\mathbf{x}, \mathbf{y}) , $\mathbf{x}, \mathbf{y} \in \mathbf{V}$. Then $\mathbf{V}_\mathbb{C}$ is a vector space over \mathbb{C} with

$$(a + \sqrt{-1}b)(\mathbf{x}, \mathbf{y}) := a(\mathbf{x}, \mathbf{y}) + b(-\mathbf{y}, \mathbf{x}), \quad \text{for all } a, b \in \mathbb{R}, \mathbf{x}, \mathbf{y} \in \mathbf{V}.$$

If \mathbf{V} has a basis $\mathbf{e}_1, \dots, \mathbf{e}_n$ over \mathbb{F} then $(\mathbf{e}_1, \mathbf{0}), \dots, (\mathbf{e}_n, \mathbf{0})$ is a basis of $\mathbf{V}_\mathbb{C}$ over \mathbb{C} . Any inner product $\langle \cdot, \cdot \rangle$ on \mathbf{V} over \mathbb{R} induces the following inner product on $\mathbf{V}_\mathbb{C}$:

$$\langle (\mathbf{x}, \mathbf{y}), (\mathbf{u}, \mathbf{v}) \rangle = \langle \mathbf{x}, \mathbf{u} \rangle + \langle \mathbf{y}, \mathbf{v} \rangle + \sqrt{-1}(\langle \mathbf{y}, \mathbf{u} \rangle - \langle \mathbf{x}, \mathbf{v} \rangle), \quad \mathbf{x}, \mathbf{y}, \mathbf{u}, \mathbf{v} \in \mathbf{V}.$$

We leave the proof of this proposition to the reader (Problem 2.4).

- Definition 2.3** Let \mathbf{V} be an IPS. Then
- (a) $\mathbf{x}, \mathbf{y} \in \mathbf{V}$ are called orthogonal if $\langle \mathbf{x}, \mathbf{y} \rangle = 0$.
 - (b) $S, T \subset \mathbf{V}$ are called orthogonal if $\langle \mathbf{x}, \mathbf{y} \rangle = 0$ for any $\mathbf{x} \in S, \mathbf{y} \in T$.
 - (d) For any $S \subset \mathbf{V}$, $S^\perp \subset \mathbf{V}$ is the maximal orthogonal set to S .
 - (e) $\mathbf{x}_1, \dots, \mathbf{x}_m$ is called an orthonormal set if

$$\langle \mathbf{x}_i, \mathbf{x}_j \rangle = \delta_{ij}, \quad i, j = 1, \dots, m.$$

- (f) $\mathbf{x}_1, \dots, \mathbf{x}_n$ is called an orthonormal basis if it is an orthonormal set which is a basis in \mathbf{V} .

Definition 2.4 (Gram-Schmidt algorithm.) Let \mathbf{V} be an IPS and $S = \{\mathbf{x}_1, \dots, \mathbf{x}_m\} \subset \mathbf{V}$ a finite (possibly empty) set ($m \geq 0$). Then $\tilde{S} = \{\mathbf{e}_1, \dots, \mathbf{e}_p\}$ is the orthonormal set ($p \geq 1$) or the empty set ($p = 0$) obtained from S using the following recursive steps:

- (a) If $\mathbf{x}_1 = \mathbf{0}$ remove it from S . Otherwise replace \mathbf{x}_1 by $\|\mathbf{x}_1\|^{-1}\mathbf{x}_1$.
- (b) Assume that $\mathbf{x}_1, \dots, \mathbf{x}_k$ is an orthonormal set and $1 \leq k < m$. Let $\mathbf{y}_{k+1} = \mathbf{x}_{k+1} - \sum_{i=1}^k \langle \mathbf{x}_{k+1}, \mathbf{x}_i \rangle \mathbf{x}_i$. If $\mathbf{y}_{k+1} = \mathbf{0}$ remove \mathbf{x}_{k+1} from S . Otherwise replace \mathbf{x}_{k+1} by $\|\mathbf{y}_{k+1}\|^{-1}\mathbf{y}_{k+1}$.

Corollary 2.5 Let \mathbf{V} be an IPS and $S = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subset \mathbf{V}$ be n linearly independent vectors. Then the Gram-Schmidt algorithm on S is given as follows:

$$\begin{aligned} \mathbf{y}_1 &:= \mathbf{x}_1, \quad r_{11} := \|\mathbf{y}_1\|, \quad \mathbf{e}_1 := \frac{\mathbf{y}_1}{r_{11}}, \\ r_{ji} &:= \langle \mathbf{x}_i, \mathbf{e}_j \rangle, \quad j = 1, \dots, i-1, \\ \mathbf{y}_i &:= \mathbf{x}_i - \sum_{j=1}^{i-1} r_{ji} \mathbf{e}_j, \quad r_{ii} := \|\mathbf{y}_i\|, \quad \mathbf{e}_i := \frac{\mathbf{y}_i}{r_{ii}}, \quad i = 2, \dots, n. \end{aligned} \tag{2.1}$$

In particular, $\mathbf{e}_i \in S_i$ and $\|\mathbf{y}_i\| = \text{dist}(\mathbf{x}_i, S_{i-1})$, where $S_i = \text{span}(\mathbf{x}_1, \dots, \mathbf{x}_i)$ for $i = 1, \dots, n$ and $S_0 = \{\mathbf{0}\}$. (See Problem 2.5 for the definition of $\text{dist}(\mathbf{x}_i, S_{i-1})$.)

Corollary 2.6 Any (ordered) basis in a finite dimensional IPS \mathbf{V} induces an orthonormal basis by the Gram-Schmidt algorithm.

See Problem 2.5 for some known properties related to the above notions.

Remark 2.7 It is known, e.g. [1] that the Gram-Schmidt process as described in (2.1) is numerically unstable. That is, there is a severe loss of orthogonality of \mathbf{y}_1, \dots as we proceed to compute \mathbf{y}_i . In computations one uses either a modified GSP or Householder orthogonalization [1].

Definition 2.8 (Modified Gram-Schmidt algorithm.) Let \mathbf{V} be an IPS and $S = \{\mathbf{x}_1, \dots, \mathbf{x}_m\} \subset \mathbf{V}$ a finite (possibly empty) set ($m \geq 0$). Then $\tilde{S} = \{\mathbf{e}_1, \dots, \mathbf{e}_p\}$ is the orthonormal set ($p \geq 1$) or the empty set ($p = 0$) obtained from S using the following recursive steps:

- Initialize $j = 1$ and $p = m$.
- If $\mathbf{x}_j \neq \mathbf{0}$ let $\mathbf{e}_j := \frac{1}{\|\mathbf{x}_j\|}\mathbf{x}_j$. If $\mathbf{x}_j = \mathbf{0}$ replace p by $p - 1$ and \mathbf{x}_i by \mathbf{x}_{i+1} for $i = j, \dots, p$.
- $\mathbf{p}_i := \langle \mathbf{x}_i, \mathbf{e}_j \rangle \mathbf{e}_j$ and replace \mathbf{x}_i by $\mathbf{x}_i := \mathbf{x}_i - \mathbf{p}_i$ for $i = j + 1, \dots, p$.
- Let $j = j + 1$ and repeat the process.

MGS algorithm is stable, needs mn^2 flops, which is more time consuming than GS algorithm.

Problems

(2.2)

Let \mathbf{V} be an IPS over \mathbb{F} . Show

$$\begin{aligned} \langle \mathbf{0}, \mathbf{x} \rangle &= \langle \mathbf{x}, \mathbf{0} \rangle = 0, \\ \text{for } \mathbb{F} = \mathbb{R} \quad \langle \mathbf{z}, a\mathbf{x} + b\mathbf{y} \rangle &= a\langle \mathbf{z}, \mathbf{x} \rangle + b\langle \mathbf{z}, \mathbf{y} \rangle, \text{ for all } a, b \in \mathbb{R}, \mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathbf{V}, \\ \text{for } \mathbb{F} = \mathbb{C} \quad \langle \mathbf{z}, a\mathbf{x} + b\mathbf{y} \rangle &= \bar{a}\langle \mathbf{z}, \mathbf{x} \rangle + \bar{b}\langle \mathbf{z}, \mathbf{y} \rangle, \text{ for all } a, b \in \mathbb{C}, \mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathbf{V}. \end{aligned}$$

(2.3)

Let \mathbf{V} be an IPS. Show

- $\|a\mathbf{x}\| = |a| \|\mathbf{x}\|$ for $a \in \mathbb{F}$ and $\mathbf{x} \in \mathbf{V}$.
- The Cauchy-Schwarz inequality:

$$|\langle \mathbf{x}, \mathbf{y} \rangle| \leq \|\mathbf{x}\| \|\mathbf{y}\|,$$

and equality holds if and only if \mathbf{x}, \mathbf{y} are linearly dependent (collinear).

- The triangle inequality

$$\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|,$$

and equality holds if either $\mathbf{x} = \mathbf{0}$ or $\mathbf{y} = a\mathbf{x}$ for $a \in \mathbb{R}_+$.

(2.4)

Prove Proposition 2.2.

(2.5)

Let \mathbf{V} be a finite dimensional IPS of dimension n . Assume that $S \subset \mathbf{V}$. Show

- If $\mathbf{x}_1, \dots, \mathbf{x}_m$ is an orthonormal set then $\mathbf{x}_1, \dots, \mathbf{x}_m$ are linearly independent.
- Assume that $\mathbf{e}_1, \dots, \mathbf{e}_n$ is an orthonormal basis in \mathbf{V} . Show that for any $\mathbf{x} \in \mathbf{V}$ the orthonormal expansion holds

$$\mathbf{x} = \sum_{i=1}^n \langle \mathbf{x}, \mathbf{e}_i \rangle \mathbf{e}_i. \quad (2.6)$$

Furthermore for any $\mathbf{x}, \mathbf{y} \in \mathbf{V}$

$$\langle \mathbf{x}, \mathbf{y} \rangle = \sum_{i=1}^n \langle \mathbf{x}, \mathbf{e}_i \rangle \overline{\langle \mathbf{y}, \mathbf{e}_i \rangle}. \quad (2.7)$$

(c) Assume that S is a finite set. Let \tilde{S} be the set obtained by the Gram-Schmidt process. Show that $\tilde{S} = \emptyset \iff \text{span} S = \{\mathbf{0}\}$. Show that if $\tilde{S} \neq \emptyset$ then $\mathbf{e}_1, \dots, \mathbf{e}_p$ is an orthonormal basis in $\text{span} S$.

(d) There exists an orthonormal basis $\mathbf{e}_1, \dots, \mathbf{e}_n$ in \mathbf{V} and $0 \leq m \leq n$ such that

$$\begin{aligned}\mathbf{e}_1, \dots, \mathbf{e}_m &\in S, & \text{span} S &= \text{span}(\mathbf{e}_1, \dots, \mathbf{e}_m), \\ S^\perp &= \text{span}(\mathbf{e}_{m+1}, \dots, \mathbf{e}_n), \\ (S^\perp)^\perp &= \text{span} S.\end{aligned}$$

(e) Assume from here to the end of the problem that S is a subspace. Show $\mathbf{V} = S \oplus S^\perp$.

(f) Let $\mathbf{x} \in \mathbf{V}$ and let $\mathbf{x} = \mathbf{u} + \mathbf{v}$ for unique $\mathbf{u} \in S$, $\mathbf{v} \in S^\perp$. Let $P(\mathbf{x}) := \mathbf{u}$ be the projection of \mathbf{x} on S . Show that $P : \mathbf{V} \rightarrow \mathbf{V}$ is a linear transformation satisfying

$$P^2 = P, \quad \text{Range } P = S, \quad \text{Ker } P = S^\perp.$$

(g) Show

$$\begin{aligned}\text{dist}(\mathbf{x}, S) &:= \|\mathbf{x} - P\mathbf{x}\| \leq \|\mathbf{x} - \mathbf{w}\| \text{ for any } \mathbf{w} \in S \\ \text{and equality} &\iff \mathbf{w} = P\mathbf{x}.\end{aligned}\tag{2.8}$$

(h) Show that $\text{dist}(\mathbf{x}, S) = \|\mathbf{x} - \mathbf{w}\|$ for some $\mathbf{w} \in S$ if and only if $\mathbf{x} - \mathbf{w}$ is orthogonal to S .

(i) Let $\mathbf{e}_1, \dots, \mathbf{e}_m$ be an orthonormal basis of S . Show that for each $\mathbf{x} \in \mathbf{V}$ $P\mathbf{x} = \sum_{i=1}^m \langle \mathbf{x}, \mathbf{e}_i \rangle \mathbf{e}_i$.

(Note: $P\mathbf{x}$ is called *the least square approximation* to \mathbf{x} in the subspace S .)

(2.9)

Let $X \in \mathbb{C}^{m \times n}$ and assume that $m \geq n$ and $\text{rank } X = n$. Let $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{C}^m$ be the columns of X , i.e. $X = (\mathbf{x}_1, \dots, \mathbf{x}_n)$. Assume that \mathbb{C}^m is an IPS with the standard inner product $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{y}^* \mathbf{x}$. Perform the Gram-Schmidt algorithm (2.5) to obtain the matrix $Q = (\mathbf{e}_1, \dots, \mathbf{e}_n) \in \mathbb{C}^{m \times n}$. Let $R = (r_{ji})_1^n \in \mathbb{C}^{n \times n}$ be the upper triangular matrix with r_{ji} , $j \leq i$ given by (2.1). Show that $\bar{Q}^T Q = I_n$ and $X = QR$. (This is the QR algorithm.) Show that if in addition $X \in \mathbb{R}^{m \times n}$ then Q and R are real valued matrices.

(2.10)

Let $C \in \mathbb{C}^{n \times n}$ and assume that $\{\lambda_1, \dots, \lambda_n\}$ are n eigenvalues of C counted with their multiplicities. View C as an operator $C : \mathbb{C}^n \rightarrow \mathbb{C}^n$. View \mathbb{C}^n as $2n$ -dimensional vector space over \mathbb{R}^{2n} . Let $C = A + \sqrt{-1}B$, $A, B \in M_n(\mathbb{R})$.

a. Then $\hat{C} := \begin{bmatrix} A & -B \\ B & A \end{bmatrix} \in M_{2n}(\mathbb{R})$ represents the operator $C : \mathbb{C}^n \rightarrow \mathbb{C}^n$ as an operator over \mathbb{R} in suitably chosen basis.

b. Show that $\{\lambda_1, \bar{\lambda}_1, \dots, \lambda_n, \bar{\lambda}_n\}$ are the $2n$ eigenvalues of \hat{C} counting with multiplicities.

c. Show that the Jordan canonical form of \hat{C} , is obtained by replacing each Jordan block $\lambda I + H$ in C by two Jordan blocks $\lambda I + H$ and $\bar{\lambda} I + H$.

2.2 Geometric interpretation of the determinant

Definition 2.9 Let $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^n$ be m given vectors. Then the parallelepiped $P(\mathbf{x}_1, \dots, \mathbf{x}_m)$ is defined as follows. The 2^m vertices of $P(\mathbf{x}_1, \dots, \mathbf{x}_m)$ are of the form $\mathbf{v} := \sum_{i=1}^m a_i \mathbf{x}_i$, where $a_i = 0, 1$ for $i = 1, \dots, m$. Two vertices $\mathbf{v} = \sum_{i=1}^m a_i \mathbf{x}_i$ and $\mathbf{w} = \sum_{i=1}^m b_i \mathbf{x}_i$ of $P(\mathbf{x}_1, \dots, \mathbf{x}_m)$ are adjacent, i.e. connected by an edge in $P(\mathbf{x}_1, \dots, \mathbf{x}_m)$, if $\|(a_1, \dots, a_m)^\top - (b_1, \dots, b_m)^\top\| = 1$, i.e. the 0–1 coordinates of $(a_1, \dots, a_m)^\top$ and $(b_1, \dots, b_m)^\top$ differ only at one coordinate k , for some $k \in [1, m]$.

Note that if $\mathbf{e}_1, \dots, \mathbf{e}_n$ is the standard basis in \mathbb{R}^n , i.e. $\mathbf{e}_i = (\delta_{1i}, \dots, \delta_{ni})^\top$, $i = 1, \dots, n$, then $P(\mathbf{e}_1, \dots, \mathbf{e}_m)$ is the m -dimensional unit cube, whose edges are parallel to $\mathbf{e}_1, \dots, \mathbf{e}_m$ and its center (of gravity) is $\frac{1}{2}(\underbrace{1, \dots, 1}_m, 0, \dots, 0)^\top$, where 1 appears m times for $1 \leq m \leq n$.

For $m > n$ $P(\mathbf{x}_1, \dots, \mathbf{x}_m)$ is "flattened" parallelepiped, since $\mathbf{x}_1, \dots, \mathbf{x}_m$ are always linearly dependent in \mathbb{R}^n for $m > n$.

Proposition 2.10 Let $A \in \mathbb{R}^{n \times n}$ and view $A = [\mathbf{c}_1 \ \mathbf{c}_2 \ \dots \ \mathbf{c}_n]$ as an ordered set of n vectors, (columns), $\mathbf{c}_1, \dots, \mathbf{c}_n$. Then $|\det A|$ is the n -dimensional volume of the parallelepiped $P(\mathbf{c}_1, \dots, \mathbf{c}_n)$. If $\mathbf{c}_1, \dots, \mathbf{c}_n$ are linearly independent then the orientation in \mathbb{R}^n induced by $\mathbf{c}_1, \dots, \mathbf{c}_n$ is the same as the orientation induced by $\mathbf{e}_1, \dots, \mathbf{e}_n$ if $\det A > 0$, and is the opposite orientation if $\det A < 0$.

Proof. $\det A = 0$ if and only if the columns of A are linearly dependent. If $\mathbf{c}_1, \dots, \mathbf{c}_n$ are linearly dependent, then $P(\mathbf{c}_1, \dots, \mathbf{c}_n)$ lies in a subspace of \mathbb{R}^n , i.e. some $n - 1$ dimensional subspace, and hence the n -dimensional volume of $P(\mathbf{c}_1, \dots, \mathbf{c}_n)$ is zero.

Assume now that $\det A \neq 0$, i.e. $\mathbf{c}_1, \dots, \mathbf{c}_n$ are linearly independent. Perform that Gram-Schmidt process 2.4. Then $A = QR$, where $Q = [\mathbf{e}_1 \ \mathbf{e}_2 \ \dots \ \mathbf{e}_n]$ is an orthogonal matrix and $R = (r_{ji}) \in \mathbb{R}^{n \times n}$ is an upper diagonal matrix. (See Problem 2.9.) So $\det A = \det Q \det R$. Since $Q^\top Q = I_n$ we deduce that $1 = \det I_n = \det Q^\top \det Q = \det Q \det Q = (\det Q)^2$. So $\det Q = \pm 1$ and the sign of $\det Q$ is the sign of $\det A$.

Hence $|\det A| = \det R = r_{11}r_{22} \dots r_{nn}$. Recall that r_{11} is the length of the vector \mathbf{c}_1 , and r_{ii} is the distance of the vector \mathbf{e}_i to the subspace spanned by $\mathbf{e}_1, \dots, \mathbf{e}_{i-1}$ for $i = 2, \dots, n$. (See Problem 2.5 parts (f-i).) Thus the length of $P(\mathbf{c}_1)$ is r_{11} . The distance of \mathbf{c}_2 to $P(\mathbf{c}_1)$ is r_{22} . Hence the area, i.e 2-dimensional volume of $P(\mathbf{c}_1, \mathbf{c}_2)$ is $r_{11}r_{22}$. Continuing in this manner we deduce that the $i - 1$ dimensional volume of $P(\mathbf{c}_1, \dots, \mathbf{c}_{i-1})$ is $r_{11} \dots r_{(i-1)(i-1)}$. As the distance of \mathbf{c}_i to $P(\mathbf{c}_1, \dots, \mathbf{c}_{i-1})$ is r_{ii} it follows that the i -dimensional volume of $P(\mathbf{c}_1, \dots, \mathbf{c}_i)$ is $r_{11} \dots r_{ii}$. For $i = n$ we get that $|\det A| = r_{11} \dots r_{nn}$ which is equal to the n -dimensional volume of $P(\mathbf{c}_1, \dots, \mathbf{c}_n)$.

As we already pointed out the sign of $\det A$ is equal to the sign of $\det Q = \pm 1$. If $\det Q = 1$ it is possible to "rotate" the standard basis in \mathbb{R}^n to the basis given by the columns of an orthogonal matrix Q with $\det Q = 1$. If $\det Q = -1$, we need one reflection, i.e. replace the standard basis $\mathbf{e}_1, \dots, \mathbf{e}_n$ by the new basis $-\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n$ and then rotate the new basis $-\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n$ to the basis consisting of the columns of an orthogonal matrix Q , where $\det Q = -1$. \square

Theorem 2.11 (The Hadamard determinantal inequality) Let $A = [\mathbf{c}_1, \dots, \mathbf{c}_n] \in \mathbb{C}^{n \times n}$. Then $|\det A| \leq \|\mathbf{c}_1\| \|\mathbf{c}_2\| \dots \|\mathbf{c}_n\|$. Equality holds if and only if either $\mathbf{c}_i = \mathbf{0}$ for some i or $\langle \mathbf{c}_i, \mathbf{c}_j \rangle = 0$ for all $i \neq j$, i.e. $\mathbf{c}_1, \dots, \mathbf{c}_n$ is an orthogonal system.

Proof. Assume first that $\det A = 0$. Clearly the Hadamard inequality holds. Equality in Hadamard inequality if and only if $\mathbf{c}_i = \mathbf{0}$ for some i .

Assume now that $\det A \neq 0$ and perform the Gram-Schmidt process. From (2.1) it follows that $A = QR$ where Q is a unitary matrix, i.e. $Q^*Q = I_n$ and $R = (r_{ji}) \in \mathbb{C}^{n \times n}$ upper triangular with r_{ii} real and positive numbers. So $\det A = \det Q \det R$. Thus

$$1 = \det I_n = \det Q^*Q = \det Q^* \det Q = \overline{\det Q} \det Q = |\det Q|^2 \Rightarrow |\det Q| = 1.$$

Hence $|\det A| = \det R = r_{11}r_{22} \dots r_{nn}$. According to Problem 2.5 and the proof of Proposition 2.10 we know that $\|\mathbf{c}_i\| \geq \text{dist}(\mathbf{c}_i, \text{span}(\mathbf{c}_1, \dots, \mathbf{c}_{i-1})) = r_{ii}$ for $i = 2, \dots, n$. Hence $|\det A| = \det R \leq \|\mathbf{c}_1\| \|\mathbf{c}_2\| \dots \|\mathbf{c}_n\|$. Equality holds if $\|\mathbf{c}_i\| = \text{dist}(\mathbf{c}_i, \text{span}(\mathbf{c}_1, \dots, \mathbf{c}_{i-1}))$ for $i = 2, \dots, n$. Use Problem 2.5 to deduce that $\|\mathbf{c}_i\| = \text{dist}(\mathbf{c}_i, \text{span}(\mathbf{c}_1, \dots, \mathbf{c}_{i-1}))$ if and only if $\langle \mathbf{c}_i, \mathbf{c}_j \rangle = 0$ for $j = 1, \dots, i-1$. Use these conditions for $i = 2, \dots$ to deduce that equality in Hadamard inequality holds if and only if $\mathbf{c}_1, \dots, \mathbf{c}_n$ is an orthogonal system. \square

Problems

1. Let $A = (a_{ij})_{i,j} \in \mathbb{C}^{n \times n}$. Assume that $|a_{ij}| \leq K$ for all $i, j = 1, \dots, n$. Show that $|\det A| \leq K^n n^{\frac{n}{2}}$.
2. Let $A = (a_{ij})_{i,j=1}^n \in \mathbb{C}^{n \times n}$ such that $|a_{ij}| \leq 1$ for $i, j = 1, \dots, n$. Show that $|\det A| = n^{\frac{n}{2}}$ if and only if $A^*A = AA^* = nI_n$. In particular, if $|\det A| = n^{\frac{n}{2}}$ then $|a_{ij}| = 1$ for $i, j = 1, \dots, n$.
3. Show that for each n there exists a matrix $A = (a_{ij})_{i,j=1}^n \in \mathbb{C}^{n \times n}$ such that $|a_{ij}| = 1$ for $i, j = 1, \dots, n$ and $|\det A| = n^{\frac{n}{2}}$.
4. Let $A = (a_{ij}) \in \mathbb{R}^{n \times n}$ and assume that $a_{ij} = \pm 1, i, j = 1, \dots, n$. Show that if $n > 2$ then the assumption that $|\det A| = n^{\frac{n}{2}}$ yields that n is divisible by 4.
5. Show that for any $n = 2^m, m = 0, 1, \dots$ there exists $A = (a_{ij}) \in \mathbb{R}^{n \times n}$ such that $a_{ij} = \pm 1, i, j = 1, \dots, n$ and $|\det A| = n^{\frac{n}{2}}$. (*Hint:* Try to prove by induction on m that $A \in \mathbb{R}^{2^m \times 2^m}$ can be chosen symmetric, and then construct $B \in \mathbb{R}^{2^{m+1} \times 2^{m+1}}$ using A .)

Note: A matrix $A = (a_{ij})_{i,j=1}^n \in \mathbb{R}^{n \times n}$ such that $a_{ij} = \pm 1$ for $i, j = 1, \dots, n$ and $|\det A| = n^{\frac{n}{2}}$ is called a *Hadamard matrix*. It is conjectured that for each n divisible by 4 there exists a Hadamard matrix.

2.3 Special transformations in IPS

Proposition 2.12 Let \mathbf{V} be an IPS and $T : \mathbf{V} \rightarrow \mathbf{V}$ a linear transformation. Then there exists a unique linear transformation $T^* : \mathbf{V} \rightarrow \mathbf{V}$ such that $\langle T\mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{x}, T^*\mathbf{y} \rangle$ for all $\mathbf{x}, \mathbf{y} \in \mathbf{V}$.

See Problems 1-2.

Definition 2.13 Let \mathbf{V} be an IPS and let $T : \mathbf{V} \rightarrow \mathbf{V}$ be a linear transformation. Then

- (a) T is called self-adjoint if $T^* = T$;
- (b) T is called anti self-adjoint if $T^* = -T$;
- (c) T is called unitary if $T^*T = TT^* = I$;
- (d) T is called normal if $T^*T = TT^*$.

Denote by $\mathbf{S}(\mathbf{V})$, $\mathbf{AS}(\mathbf{V})$, $\mathbf{U}(\mathbf{V})$, $\mathbf{N}(\mathbf{V})$ the sets of self-adjoint, anti self-adjoint, unitary and normal operators on \mathbf{V} respectively.

Proposition 2.14 Let \mathbf{V} be an IPS over $\mathbb{F} = \mathbb{R}, \mathbb{C}$ with an orthonormal basis $E = \{\mathbf{e}_1, \dots, \mathbf{e}_n\}$. Let $T : \mathbf{V} \rightarrow \mathbf{V}$ be a linear transformation. Let $A = (a_{ij}) \in \mathbb{F}^{n \times n}$ be the representation matrix of T in the basis E :

$$a_{ij} = \langle T\mathbf{e}_j, \mathbf{e}_i \rangle, \quad i, j = 1, \dots, n. \quad (2.1)$$

Then for $\mathbb{F} = \mathbb{R}$:

- (a) T^* is represented by A^\top ,
- (b) T is selfadjoint $\iff A = A^\top$,
- (c) T is anti selfadjoint $\iff A = -A^\top$,
- (d) T is unitary $\iff A$ is orthogonal $\iff AA^\top = A^\top A = I$,
- (e) T is normal $\iff A$ is normal $\iff AA^\top = A^\top A$,

and for $\mathbb{F} = \mathbb{C}$:

- (a) T^* is represented by $A^* (:= \bar{A}^\top)$,
- (b) T is selfadjoint $\iff A$ is hermitian $\iff A = A^*$,
- (c) T is anti selfadjoint $\iff A$ is anti hermitian $\iff A = -A^*$,
- (d) T is unitary $\iff A$ is unitary $\iff AA^* = A^*A = I$,
- (e) T is normal $\iff A$ is normal $\iff AA^* = A^*A$.

See Problem 3.

Proposition 2.15 Let \mathbf{V} be an IPS over \mathbb{R} , and let $T \in \text{Hom}(\mathbf{V})$. Let \mathbf{V}_c be the complexification of \mathbf{V} . Show that there exists a unique $T_c \in \text{Hom}(\mathbf{V}_c)$ such that $T_c|_{\mathbf{V}} = T$. Furthermore T is self-adjoint, unitary or normal if and only if T_c is self-adjoint, unitary or normal respectively.

See Problem 4

Definition 2.16 For a field \mathbb{F} let

$$\begin{aligned}
\mathbf{S}(n, \mathbb{F}) &:= \{A \in \mathbb{F}^{n \times n} : A = A^\top\}, \\
\mathbf{AS}(n, \mathbb{F}) &:= \{A \in \mathbb{F}^{n \times n} : A = -A^\top\}, \\
\mathbf{O}(n, \mathbb{F}) &:= \{A \in \mathbb{F}^{n \times n} : AA^\top = A^\top A = I\}, \\
\mathbf{SO}(n, \mathbb{F}) &:= \{A \in \mathbf{O}(n, \mathbb{F}) : \det A = 1\}, \\
\mathbf{DO}(n, \mathbb{F}) &:= \mathbf{D}(n, \mathbb{F}) \cap \mathbf{O}(n, \mathbb{F}), \\
\mathbf{N}(n, \mathbb{R}) &:= \{A \in \mathbb{R}^{n \times n} : AA^\top = A^\top A\}, \\
\mathbf{N}(n, \mathbb{C}) &:= \{A \in \mathbb{C}^{n \times n} : AA^* = A^*A\}, \\
\mathbf{H}_n &:= \{A \in \mathbb{C}^{n \times n} : A = A^*\}, \\
\mathbf{AH}_n &:= \{A \in \mathbb{C}^{n \times n} : A = -A^*\}, \\
\mathbf{U}_n &:= \{A \in \mathbb{C}^{n \times n} : AA^* = A^*A = I\}, \\
\mathbf{SU}_n &:= \{A \in \mathbf{U}_n : \det A = 1\}, \\
\mathbf{DU}_n &:= \mathbf{D}(n, \mathbb{C}) \cap \mathbf{U}_n.
\end{aligned}$$

See Problem 5 for relations between these classes.

Theorem 2.17 Let \mathbf{V} be an IPS over \mathbb{C} of dimension n . Then a linear transformation $T : \mathbf{V} \rightarrow \mathbf{V}$ is normal if and only if \mathbf{V} has an orthonormal basis consisting of eigenvectors of T .

Proof. Suppose first that \mathbf{V} has an orthonormal basis $\mathbf{e}_1, \dots, \mathbf{e}_n$ such that $T\mathbf{e}_i = \lambda_i \mathbf{e}_i$, $i = 1, \dots, n$. From the definition of T^* it follows that $T^*\mathbf{e}_i = \bar{\lambda}_i \mathbf{e}_i$, $i = 1, \dots, n$. Hence $TT^* = T^*T$.

Assume now T is normal. Since \mathbb{C} is algebraically closed T has an eigenvalue λ_1 . Let \mathbf{V}_1 be the subspace of \mathbf{V} spanned by all eigenvectors of T corresponding to the eigenvalue λ_1 . Clearly $T\mathbf{V}_1 \subset \mathbf{V}_1$. Let $\mathbf{x} \in \mathbf{V}_1$. Then $T\mathbf{x} = \lambda_1 \mathbf{x}$. Thus

$$T(T^*\mathbf{x}) = (TT^*)\mathbf{x} = (T^*T)\mathbf{x} = T^*(T\mathbf{x}) = \lambda_1 T^*\mathbf{x} \Rightarrow T^*\mathbf{V}_1 \subset \mathbf{V}_1.$$

Hence $T\mathbf{V}_1^\perp, T^*\mathbf{V}_1^\perp \subset \mathbf{V}_1^\perp$. Since $\mathbf{V} = \mathbf{V}_1 \oplus \mathbf{V}_1^\perp$ it is enough to prove the theorem for $T|_{\mathbf{V}_1}$ and $T|_{\mathbf{V}_1^\perp}$.

As $T|_{\mathbf{V}_1} = \lambda_1 I_{\mathbf{V}_1}$ it is straightforward to show $T^*|_{\mathbf{V}_1} = \bar{\lambda}_1 I_{\mathbf{V}_1}$ (see Problem 2). Hence for $T|_{\mathbf{V}_1}$ the theorem trivially holds. For $T|_{\mathbf{V}_1^\perp}$ the theorem follows by induction. \square

The proof of Theorem 2.17 yields:

Corollary 2.18 Let \mathbf{V} be an IPS over \mathbb{R} of dimension n . Then the linear transformation $T : \mathbf{V} \rightarrow \mathbf{V}$ with a real spectrum is normal if and only if \mathbf{V} has an orthonormal basis consisting of eigenvectors of T .

Proposition 2.19 Let \mathbf{V} be an IPS over \mathbb{C} . Let $T \in \mathbf{N}(\mathbf{V})$. Then

$$\begin{aligned}
T \text{ is self-adjoint} &\iff \text{spec}(T) \subset \mathbb{R}, \\
T \text{ is unitary} &\iff \text{spec}(T) \subset S^1 = \{z \in \mathbb{C} : |z| = 1\}.
\end{aligned}$$

Proof. Since T is normal there exists an orthonormal basis $\mathbf{e}_1, \dots, \mathbf{e}_n$ such that $T\mathbf{e}_i = \lambda_i\mathbf{e}_i$, $i = 1, \dots, n$. Hence $T^*\mathbf{e}_i = \bar{\lambda}_i\mathbf{e}_i$. Then

$$\begin{aligned} T = T^* &\iff \lambda_i = \bar{\lambda}_i, \quad i = 1, \dots, n, \\ TT^* = T^*T = I &\iff |\lambda_i| = 1, \quad i = 1, \dots, n. \end{aligned}$$

□

Combine Proposition 2.15 and Corollary 2.18 with the above proposition to deduce:

Corollary 2.20 *Let \mathbf{V} be an IPS over \mathbb{R} and let $T \in \mathbf{S}(\mathbf{V})$. Then $\text{spec}(T) \subset \mathbb{R}$ and \mathbf{V} has an orthonormal basis consisting of the eigenvectors of T .*

Proposition 2.21 *Let \mathbf{V} be an IPS over \mathbb{R} and let $T \in \mathbf{U}(\mathbf{V})$. Then $\mathbf{V} = \bigoplus_{i \in \{-1, 1, 2, \dots, k\}} \mathbf{V}_i$, where $k \geq 1$, \mathbf{V}_i and \mathbf{V}_j are orthogonal for $i \neq j$, such that*

- (a) $T|_{\mathbf{V}_{-1}} = -I_{\mathbf{V}_{-1}}$ $\dim \mathbf{V}_{-1} \geq 0$,
- (b) $T|_{\mathbf{V}_1} = I_{\mathbf{V}_1}$ $\dim \mathbf{V}_1 \geq 0$,
- (c) $T\mathbf{V}_i = \mathbf{V}_i$, $\dim \mathbf{V}_i = 2$, $\text{spec}(T|_{\mathbf{V}_i}) \subset S^1 \setminus \{-1, 1\}$ for $i = 2, \dots, k$.

See Problem 7.

Proposition 2.22 *Let \mathbf{V} be an IPS over \mathbb{R} and let $T \in \mathbf{AS}(\mathbf{V})$. Then $\mathbf{V} = \bigoplus_{i \in \{1, 2, \dots, k\}} \mathbf{V}_i$, where $k \geq 1$, \mathbf{V}_i and \mathbf{V}_j are orthogonal for $i \neq j$, such that*

- (a) $T|_{\mathbf{V}_1} = 0_{\mathbf{V}_1}$ $\dim \mathbf{V}_1 \geq 0$,
- (b) $T\mathbf{V}_i = \mathbf{V}_i$, $\dim \mathbf{V}_i = 2$, $\text{spec}(T|_{\mathbf{V}_i}) \subset \sqrt{-1}\mathbb{R} \setminus \{0\}$ for $i = 2, \dots, k$.

See Problem 8.

Theorem 2.23 *Let \mathbf{V} be an IPS over \mathbb{C} of dimension n . Let $T \in \text{Hom}(\mathbf{V})$. Let $\lambda_1, \dots, \lambda_n \in \mathbb{C}$ be n eigenvalues of T counted with their multiplicities. Then there exists an orthonormal basis $\mathbf{g}_1, \dots, \mathbf{g}_n$ of \mathbf{V} with the following properties:*

$$T\text{span}(\mathbf{g}_1, \dots, \mathbf{g}_i) \subset \text{span}(\mathbf{g}_1, \dots, \mathbf{g}_i), \quad \langle T\mathbf{g}_i, \mathbf{g}_i \rangle = \lambda_i, \quad i = 1, \dots, n. \quad (2.2)$$

Let \mathbf{V} be an IPS over \mathbb{R} of dimension n . Let $T \in \text{Hom}(\mathbf{V})$ and assume that $\text{spec}(T) \subset \mathbb{R}$. Let $\lambda_1, \dots, \lambda_n \in \mathbb{R}$ be n eigenvalues of T counted with their multiplicities. Then there exists an orthonormal basis $\mathbf{g}_1, \dots, \mathbf{g}_n$ of \mathbf{V} such that (2.2) holds.

Proof. Assume first that \mathbf{V} is IPS over \mathbb{C} of dimension n . The proof is by induction on n . For $n = 1$ the theorem is trivial. Assume that $n > 1$. Since $\lambda_1 \in \text{spec}(T)$ it follows that there exists $\mathbf{g}_1 \in \mathbf{V}$, $\langle \mathbf{g}_1, \mathbf{g}_1 \rangle = 1$ such that $T\mathbf{g}_1 = \lambda_1\mathbf{g}_1$. Let $\mathbf{U} := \text{span}(\mathbf{g}_1)^\perp$. Let P be the orthogonal projection on \mathbf{U} . Let $T_1 := PT|_{\mathbf{U}}$. Then $T_1 \in \text{Hom}(\mathbf{U})$. Let $\tilde{\lambda}_2, \dots, \tilde{\lambda}_n$ be the eigenvalues of T_1 counted with their multiplicities. The induction hypothesis yields the existence of an orthonormal basis $\mathbf{g}_2, \dots, \mathbf{g}_n$ of \mathbf{U} such that

$$T_1\text{span}(\mathbf{g}_2, \dots, \mathbf{g}_i) \subset \text{span}(\mathbf{g}_2, \dots, \mathbf{g}_i), \quad \langle T_1\mathbf{g}_i, \mathbf{g}_i \rangle = \tilde{\lambda}_i, \quad i = 1, \dots, n.$$

It is straightforward to show that $T\text{span}(\mathbf{g}_1, \dots, \mathbf{g}_i) \subset \text{span}(\mathbf{g}_1, \dots, \mathbf{g}_i)$ for $i = 1, \dots, n$. Hence in the orthonormal basis $\mathbf{g}_1, \dots, \mathbf{g}_n$ T is presented by an upper diagonal matrix

$B = (b_{ij})_1^n$, with $b_{11} = \lambda_1$ and $b_{ii} = \tilde{\lambda}_i$, $i = 2, \dots, n$. Hence $\lambda_1, \tilde{\lambda}_2, \dots, \tilde{\lambda}_n$ are the eigenvalues of T counted with their multiplicities. This establishes the theorem in this case. The real case is treated similarly. \square

Combine the above results with Problems 6 and 12 to deduce:

Corollary 2.24 *Let $A \in \mathbb{C}^{n \times n}$. Let $\lambda_1, \dots, \lambda_n \in \mathbb{C}$ be n eigenvalues of A counted with their multiplicities. Then there exist an upper triangular matrix $B = (b_{ij})_1^n \in \mathbb{C}^{n \times n}$, such that $b_{ii} = \lambda_i$, $i = 1, \dots, n$, and a unitary matrix $U \in \mathbf{U}_n$ such that $A = UBU^{-1}$. If $A \in \mathbf{N}(n, \mathbb{C})$ then B is a diagonal matrix.*

Let $A \in \mathbb{R}^{n \times n}$ and assume that $\text{spec}(T) \subset \mathbb{R}$. Then $A = UBU^{-1}$ where U can be chosen a real orthogonal matrix and B a real upper triangular matrix. If $A \in \mathbf{N}(n, \mathbb{R})$ and $\text{spec}(A) \subset \mathbb{R}$ then B is a diagonal matrix.

It is easy to show that U in the above Corollary can be chosen in \mathbf{SU}_n or $\mathbf{SO}(n, \mathbb{R})$ respectively (Problem 11).

Definition 2.25 *Let \mathbf{V} be a vector space and assume that $T : \mathbf{V} \rightarrow \mathbf{V}$ is a linear operator. Let $0 \neq \mathbf{v} \in \mathbf{V}$. Then $\mathbf{W} = \text{span}(\mathbf{v}, T\mathbf{v}, T^2\mathbf{v}, \dots)$ is called a cyclic invariant subspace of T generated by \mathbf{v} . (It is also referred as a Krylov subspace of T generated by \mathbf{v} .) Sometimes we will call \mathbf{W} just a cyclic subspace, or Krylov subspace.*

Theorem 2.26 *Let \mathbf{V} be a finite dimensional IPS. Let $T : \mathbf{V} \rightarrow \mathbf{V}$ be a linear operator. For $0 \neq \mathbf{v} \in \mathbf{V}$ let $\mathbf{W} = \text{span}(\mathbf{v}, T\mathbf{v}, \dots, T^{r-1}\mathbf{v})$ be a cyclic T -invariant subspace of dimension r generated by \mathbf{v} . Let $\mathbf{u}_1, \dots, \mathbf{u}_r$ be an orthonormal basis of \mathbf{W} obtained by the Gram-Schmidt process from the basis $[\mathbf{v}, T\mathbf{v}, \dots, T^{r-1}\mathbf{v}]$ of \mathbf{W} . Then $\langle T\mathbf{u}_i, \mathbf{u}_j \rangle = 0$ for $1 \leq i \leq j - 2$, i.e. the representation matrix of $T|_{\mathbf{W}}$ in the basis $[\mathbf{u}_1, \dots, \mathbf{u}_r]$ is upper Hessenberg. If T is self-adjoint then the representation matrix of $T|_{\mathbf{W}}$ in the basis $[\mathbf{u}_1, \dots, \mathbf{u}_r]$ is a tridiagonal hermitian matrix.*

Proof. Let $\mathbf{W}_j = \text{span}(\mathbf{v}, \dots, T^{j-1}\mathbf{v})$ for $j = 1, \dots, r + 1$. Clearly $T\mathbf{W}_j \subset \mathbf{W}_{j+1}$ for $j = 1, \dots, r$. The assumption that \mathbf{W} is T -invariant subspace yields $\mathbf{W} = \mathbf{W}_r = \mathbf{W}_{r+1}$. Since $\dim \mathbf{W} = r$ it follows that $\mathbf{v}, \dots, T^{r-1}\mathbf{v}$ are linearly independent. Hence $[\mathbf{v}, \dots, T^{r-1}\mathbf{v}]$ is a basis for \mathbf{W} . Recall that $\text{span}(\mathbf{u}_1, \dots, \mathbf{u}_j) = \mathbf{W}_j$ for $j = 1, \dots, r$. Let $r \geq j \geq i + 2$. Then $T\mathbf{u}_i \in T\mathbf{W}_i \subset \mathbf{W}_{i+1}$. As $\mathbf{u}_j \perp \mathbf{W}_{i+1}$ it follows that $\langle T\mathbf{u}_i, \mathbf{u}_j \rangle = 0$. Assume that $T^* = T$. Let $r \geq i \geq j + 2$. Then $\langle T\mathbf{u}_i, \mathbf{u}_j \rangle = \langle \mathbf{u}_i, T\mathbf{u}_j \rangle = 0$. Hence the representation matrix of $T|_{\mathbf{W}}$ in the basis $[\mathbf{u}_1, \dots, \mathbf{u}_r]$ is a tridiagonal hermitian matrix. \square

Problems

1. Prove Proposition 2.12.
2. Let $P, Q \in \text{Hom}(\mathbf{V})$, $\mathbf{a}, b \in \mathbb{F}$. Show that $(aP + bQ)^* = \bar{a}P^* + \bar{b}Q^*$.
3. Prove Proposition 2.14.

4. Prove Proposition 2.15 for finite dimensional \mathbf{V} . (*Hint*: Choose an orthonormal basis in \mathbf{V} .)
5. Show the following

$$\mathbf{SO}(n, \mathbb{D}) \subset \mathbf{O}(n, \mathbb{D}) \subset \text{GL}(n, \mathbb{D}),$$

$$\mathbf{S}(n, \mathbb{R}) \subset \mathbf{H}_n \subset \mathbf{N}(n, \mathbb{C}),$$

$$\mathbf{AS}(n, \mathbb{R}) \subset \mathbf{AH}_n \subset \mathbf{N}(n, \mathbb{C}),$$

$$\mathbf{S}(n, \mathbb{R}), \mathbf{AS}(n, \mathbb{R}) \subset \mathbf{N}(n, \mathbb{R}) \subset \mathbf{N}(n, \mathbb{C}),$$

$$\mathbf{O}(n, \mathbb{R}) \subset \mathbf{U}_n \subset \mathbf{N}(n, \mathbb{C}),$$

$$\mathbf{SO}(n, \mathbb{D}), \mathbf{O}(n, \mathbb{D}), \mathbf{SU}_n, \mathbf{U}_n \text{ are groups}$$

$$\mathbf{S}(n, \mathbb{D}) \text{ is a } \mathbb{D} \text{ - module of dimension } \binom{n+1}{2},$$

$$\mathbf{AS}(n, \mathbb{D}) \text{ is a } \mathbb{D} \text{ - module of dimension } \binom{n}{2},$$

$$\mathbf{H}_n \text{ is an } \mathbb{R} \text{ - vector space of dimension } n^2.$$

$$\mathbf{AH}_n = \sqrt{-1} \mathbf{H}_n$$

6. Let $E = \{\mathbf{e}_1, \dots, \mathbf{e}_n\}$ be an orthonormal basis in IPS \mathbf{V} over \mathbb{F} . Let $G = \{\mathbf{g}_1, \dots, \mathbf{g}_n\}$ be another basis in \mathbf{V} . Show that F is an orthonormal basis if and only if the transfer matrix either from E to G or from G to E is a unitary matrix.
7. Prove Proposition 2.21
8. Prove Proposition 2.22
9. a. Show that $A \in \mathbf{SO}(2, \mathbb{R})$ is of the form $A = \begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix}, \theta \in \mathbb{R}$.
- b. Show that $\mathbf{SO}(2, \mathbb{R}) = e^{\mathbf{AS}(2, \mathbb{R})}$. That is for any $B \in \mathbf{AS}(2, \mathbb{R})$ $e^B \in \mathbf{SO}(2, \mathbb{R})$ and any $A \in \mathbf{SO}(2, \mathbb{R})$ is e^B for some $B \in \mathbf{AS}(2, \mathbb{R})$. (*Hint*: Consider the power series for e^B , $B = \begin{bmatrix} 0 & \theta \\ -\theta & 0 \end{bmatrix}$.)
- c. Show that $\mathbf{SO}(n, \mathbb{R}) = e^{\mathbf{AS}(n, \mathbb{R})}$. (*Hint*: Use Propositions 2.21 and 2.22 and part b.)
- d. Show that $\mathbf{SO}(n, \mathbb{R})$ is a path connected space. (See part e.)
- e. Let \mathbf{V} be an $n (> 1)$ -dimensional IPS over $\mathbb{F} = \mathbb{R}$. Let $p \in \langle n-1 \rangle$. Assume that $\mathbf{x}_1, \dots, \mathbf{x}_p$ and $\mathbf{y}_1, \dots, \mathbf{y}_p$ be two orthonormal systems in \mathbf{V} . Show that these two o.n.s. are path connected. That is there are p continuous mappings $\mathbf{z}_i(t) : [0, 1] \rightarrow \mathbf{V}$, $i = 1, \dots, p$ such that for each $t \in [0, 1]$ $\mathbf{z}_1(t), \dots, \mathbf{z}_p(t)$ is an o.n.s. and $\mathbf{z}_i(0) = \mathbf{x}_i, \mathbf{z}_i(1) = \mathbf{y}_i, i = 1, \dots, p$.
10. a. Show that $\mathbf{U}_n = e^{\mathbf{AH}_n}$. (*Hint*: Use Proposition 2.19 and its proof.)
- b. Show that \mathbf{U}_n is path connected.
- c. Prove Problem 9e for $\mathbb{F} = \mathbb{C}$.

11. Show
- (a) $D_1 D D_1^* = D$ for any $D \in \mathbf{D}(n, \mathbb{C})$, $D_1 \in \mathbf{DU}_n$.
 - (b) $A \in \mathbf{N}(n, \mathbb{C}) \iff A = U D U^*$, $U \in \mathbf{SU}_n$, $D \in \mathbf{D}(n, \mathbb{C})$.
 - (c) $A \in \mathbf{N}(n, \mathbb{R})$, $\sigma(A) \subset \mathbb{R} \iff A = U D U^\top$, $U \in \mathbf{SO}_n$, $D \in \mathbf{D}(n, \mathbb{R})$.
12. Show that an upper triangular or a lower triangular matrix $B \in \mathbb{C}^{n \times n}$ is normal if and only if B is diagonal. (**Hint**: consider the equality $(B B^*)_{11} = (B^* B)_{11}$.)
13. Let the assumptions of Theorem 2.26 hold. Show that instead of performing the Gram-Schmidt process on $\mathbf{v}, T\mathbf{v}, \dots, T^{r-1}\mathbf{v}$ one can perform the following process. Let $\mathbf{w}_1 := \frac{1}{\|\mathbf{v}\|}\mathbf{v}$. Assume that one already obtained i orthonormal vectors $\mathbf{w}_1, \dots, \mathbf{w}_i$. Let $\tilde{\mathbf{w}}_{i+1} := T\mathbf{w}_i - \sum_{j=1}^i \langle T\mathbf{w}_i, \mathbf{w}_j \rangle \mathbf{w}_j$. If $\tilde{\mathbf{w}}_{i+1} = 0$ then stop the process, i.e. one is left with i orthonormal vectors. If $\tilde{\mathbf{w}}_{i+1} \neq 0$ then $\mathbf{w}_{i+1} := \frac{1}{\|\tilde{\mathbf{w}}_{i+1}\|}\tilde{\mathbf{w}}_{i+1}$ and continue the process. Show that the process ends after obtaining r orthonormal vectors $\mathbf{w}_1, \dots, \mathbf{w}_r$ and $\mathbf{u}_i = \mathbf{w}_i$ for $i = 1, \dots, r$. (This is a version of *Lanczos tridiagonalization* process.)

2.4 Symmetric bilinear and hermitian forms

Definition 2.27 Let \mathbf{V} be a module over \mathbb{D} and $Q : \mathbf{V} \times \mathbf{V} \rightarrow \mathbb{D}$. Q is called a symmetric bilinear form (on \mathbf{V}) if the following conditions are satisfied:

- (a) $Q(\mathbf{x}, \mathbf{y}) = Q(\mathbf{y}, \mathbf{x})$ for all $\mathbf{x}, \mathbf{y} \in \mathbf{V}$ (symmetry);
- (b) $Q(a\mathbf{x} + b\mathbf{z}, \mathbf{y}) = aQ(\mathbf{x}, \mathbf{y}) + bQ(\mathbf{z}, \mathbf{y})$ for all $a, b \in \mathbb{D}$ and $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathbf{V}$ (bilinearity).

For $\mathbb{D} = \mathbb{C}$ Q is called hermitian form (on \mathbf{V}) if Q satisfies the conditions (a') and (b) where

- (a') $Q(\mathbf{x}, \mathbf{y}) = \overline{Q(\mathbf{y}, \mathbf{x})}$ for all $\mathbf{x}, \mathbf{y} \in \mathbf{V}$ (conjugate symmetry).

The following results are elementary (see Problems 1-2):

Proposition 2.28 Let \mathbf{V} be a module over \mathbb{D} with a basis $E = \{\mathbf{e}_1, \dots, \mathbf{e}_n\}$. Then there is 1 – 1 correspondence between a symmetric bilinear form Q on \mathbf{V} and $A \in \mathbf{S}(n, \mathbb{D})$:

$$Q(\mathbf{x}, \mathbf{y}) = \eta^\top A \xi,$$

$$\mathbf{x} = \sum_{i=1}^n \xi_i \mathbf{e}_i, \mathbf{y} = \sum_{i=1}^n \eta_i \mathbf{e}_i, \xi = (\xi_1, \dots, \xi_n)^\top, \eta = (\eta_1, \dots, \eta_n)^\top \in \mathbb{D}^n.$$

Let \mathbf{V} be a vector space over \mathbb{C} with a basis $E = \{\mathbf{e}_1, \dots, \mathbf{e}_n\}$. Then there is 1 – 1 correspondence between a hermitian form Q on \mathbf{V} and $A \in \mathbf{H}_n$:

$$Q(\mathbf{x}, \mathbf{y}) = \eta^* A \xi,$$

$$\mathbf{x} = \sum_{i=1}^n \xi_i \mathbf{e}_i, \mathbf{y} = \sum_{i=1}^n \eta_i \mathbf{e}_i, \xi = (\xi_1, \dots, \xi_n)^\top, \eta = (\eta_1, \dots, \eta_n)^\top \in \mathbb{C}^n.$$

Definition 2.29 Let the assumptions of Proposition 2.28 hold. Then A is called the representation matrix of Q in the basis E .

Proposition 2.30 *Let the assumptions of Proposition 2.28. Let $F = \{\mathbf{f}_1, \dots, \mathbf{f}_n\}$ be another basis of the \mathbb{D} module \mathbf{V} . Then the symmetric bilinear form Q is represented by $B \in \mathbf{S}(n, \mathbb{D})$ in the basis F , where B is congruent A :*

$$B = U^\top AU, \quad U \in \text{GL}(n, \mathbb{D})$$

and U is the matrix corresponding to the basis change from F to E . For $\mathbb{D} = \mathbb{C}$ the hermitian form Q is presented by $B \in \mathbf{H}_n$ in the basis F , where B hermiticongruent to A :

$$B = U^*AU, \quad U \in \text{GL}(n, \mathbb{C})$$

and U is the matrix corresponding to the basis change from F to E .

In what follows we assume that $\mathbb{D} = \mathbb{F} = \mathbb{R}, \mathbb{C}$.

Proposition 2.31 *Let \mathbf{V} be an n dimensional vector space over \mathbb{R} . Let $Q : \mathbf{V} \times \mathbf{V} \rightarrow \mathbb{R}$ be a symmetric bilinear form. Let $A \in \mathbf{S}(n, \mathbb{R})$ the representation matrix of Q with respect to a basis E in \mathbf{V} . Let \mathbf{V}_c be the extension of \mathbf{V} over \mathbb{C} . Then there exists a unique hermitian form $Q_c : \mathbf{V}_c \times \mathbf{V}_c \rightarrow \mathbb{C}$ such that $Q_c|_{\mathbf{V} \times \mathbf{V}} = Q$ and Q_c is presented by A with respect to the basis E in \mathbf{V}_c .*

See Problem 3

Normalization 2.32 *Let \mathbf{V} is a finite dimensional IPS over \mathbb{F} . Let $Q : \mathbf{V} \times \mathbf{V} \rightarrow \mathbb{F}$ be either a symmetric bilinear form for $\mathbb{F} = \mathbb{R}$ or a hermitian form for $\mathbb{F} = \mathbb{C}$. Then a representation matrix A of Q is chosen with respect to an orthonormal basis E .*

The following proposition is straightforward (see Problem 4).

Proposition 2.33 *Let \mathbf{V} is an n -dimensional IPS over \mathbb{F} . Let $Q : \mathbf{V} \times \mathbf{V} \rightarrow \mathbb{F}$ be either a symmetric bilinear form for $\mathbb{F} = \mathbb{R}$ or a hermitian form for $\mathbb{F} = \mathbb{C}$. Then there exists a unique $T \in \mathbf{S}(\mathbf{V})$ such that $Q(\mathbf{x}, \mathbf{y}) = \langle T\mathbf{x}, \mathbf{y} \rangle$ for any $\mathbf{x}, \mathbf{y} \in \mathbf{V}$. In any orthonormal basis of \mathbf{V} Q and T represented by the same matrix A . In particular the characteristic polynomial $p(\lambda)$ of T is called the characteristic polynomial of Q . Q has only real roots:*

$$\lambda_1(Q) \geq \dots \geq \lambda_n(Q),$$

which are called the eigenvalues of Q . Furthermore there exists an orthonormal basis $F = \{\mathbf{f}_1, \dots, \mathbf{f}_n\}$ in \mathbf{V} such that $D = \text{diag}(\lambda_1(Q), \dots, \lambda_n(Q))$ is the representation matrix of Q in F .

Vice versa, for any $T \in \mathbf{S}(\mathbf{V})$ and any subspace $\mathbf{U} \subset \mathbf{V}$ the form $Q(T, \mathbf{U})$ defined by

$$Q(T, \mathbf{U})(\mathbf{x}, \mathbf{y}) := \langle T\mathbf{x}, \mathbf{y} \rangle \quad \text{for } \mathbf{x}, \mathbf{y} \in \mathbf{U}$$

is either a symmetric bilinear form for $\mathbb{F} = \mathbb{R}$ or a hermitian form for $\mathbb{F} = \mathbb{C}$.

In the rest of the book we use the following normalization unless stated otherwise.

Normalization 2.34 *Let \mathbf{V} is an n -dimensional IPS over \mathbb{F} . Assume that $T \in \mathbf{S}(\mathbf{V})$. Then arrange the eigenvalues of T counted with their multiplicities in the decreasing order*

$$\lambda_1(T) \geq \dots \geq \lambda_n(T).$$

Same normalization applies to real symmetric matrices and complex hermitian matrices.

Problems

1. Prove Proposition 2.28.
2. Prove Proposition 2.30.
3. Prove Proposition 2.31.
4. Prove Proposition 2.33.

2.5 Max-min characterizations of eigenvalues

Definition 2.35 Let \mathbf{V} be a finite dimensional space over the field \mathbb{F} . Denote by $\text{Gr}(m, \mathbf{V})$ be the space of all m -dimensional subspaces in \mathbf{U} of dimension $m \in [0, n] \cap \mathbb{Z}_+$.

Theorem 2.36 (The convoy principle) Let \mathbf{V} be an n -dimensional IPS. Let $T \in \mathbf{S}(\mathbf{V})$. Then

$$\lambda_k(T) = \max_{\mathbf{U} \in \text{Gr}(k, \mathbf{V})} \min_{\mathbf{0} \neq \mathbf{x} \in \mathbf{U}} \frac{\langle T\mathbf{x}, \mathbf{x} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle} = \max_{\mathbf{U} \in \text{Gr}(k, \mathbf{V})} \lambda_k(Q(T, \mathbf{U})), \quad k = 1, \dots, n, \quad (2.1)$$

where the quadratic form $Q(T, \mathbf{U})$ is defined in Proposition 2.33. For $k \in [1, n] \cap \mathbb{N}$ let \mathbf{U} be an invariant subspace of T spanned by eigenvectors $\mathbf{e}_1, \dots, \mathbf{e}_k$ corresponding to the eigenvalues $\lambda_1(T), \dots, \lambda_k(T)$. Then $\lambda_k(T) = \lambda_k(Q(T, \mathbf{U}))$. Let $\mathbf{U} \in \text{Gr}(k, \mathbf{V})$ and assume that $\lambda_k(T) = \lambda_k(Q(T, \mathbf{U}))$. Then \mathbf{U} contains an eigenvector of T corresponding to $\lambda_k(T)$.

In particular

$$\lambda_1(T) = \max_{\mathbf{0} \neq \mathbf{x} \in \mathbf{V}} \frac{\langle T\mathbf{x}, \mathbf{x} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle}, \quad \lambda_n(T) = \min_{\mathbf{0} \neq \mathbf{x} \in \mathbf{V}} \frac{\langle T\mathbf{x}, \mathbf{x} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle} \quad (2.2)$$

Moreover for any $\mathbf{x} \neq \mathbf{0}$

$$\begin{aligned} \lambda_1(T) = \frac{\langle T\mathbf{x}, \mathbf{x} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle} &\iff T\mathbf{x} = \lambda_1(T)\mathbf{x}, \\ \lambda_n(T) = \frac{\langle T\mathbf{x}, \mathbf{x} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle} &\iff T\mathbf{x} = \lambda_n(T)\mathbf{x}, \end{aligned}$$

The quotient $\frac{\langle T\mathbf{x}, \mathbf{x} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle}$, $\mathbf{0} \neq \mathbf{x} \in \mathbf{V}$ is called *Rayleigh quotient*. The characterization (2.2) is called *convoy principle*.

Proof. Choose an orthonormal basis $E = \{\mathbf{e}_1, \dots, \mathbf{e}_n\}$ such that

$$T\mathbf{e}_i = \lambda_i(T)\mathbf{e}_i, \quad \langle \mathbf{e}_i, \mathbf{e}_j \rangle = \delta_{ij} \quad i, j = 1, \dots, n. \quad (2.3)$$

Then

$$\frac{\langle T\mathbf{x}, \mathbf{x} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle} = \frac{\sum_{i=1}^n \lambda_i(T) |x_i|^2}{\sum_{i=1}^n |x_i|^2}, \quad \mathbf{x} = \sum_{i=1}^n x_i \mathbf{e}_i \neq \mathbf{0}. \quad (2.4)$$

The above equality yields straightforward (2.2) and the equality cases in these characterizations. Let $\mathbf{U} \in \text{Gr}(k, \mathbf{V})$. Then the minimal characterization of $\lambda_k(Q(T, \mathbf{U}))$ yields the equality

$$\lambda_k(Q(T, \mathbf{U})) = \min_{\mathbf{0} \neq \mathbf{x} \in \mathbf{U}} \frac{\langle T\mathbf{x}, \mathbf{x} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle} \quad \text{for any } \mathbf{U} \in \text{Gr}(k, \mathbf{U}). \quad (2.5)$$

Next there exists $\mathbf{0} \neq \mathbf{x} \in \mathbf{U}$ such that $\langle \mathbf{x}, \mathbf{e}_i \rangle = 0$ for $i = 1, \dots, k-1$. (For $k = 1$ this condition is void.) Hence

$$\frac{\langle T\mathbf{x}, \mathbf{x} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle} = \frac{\sum_{i=k}^n \lambda_i(T) |x_i|^2}{\sum_{i=k}^n |x_i|^2} \leq \lambda_k(T) \Rightarrow \lambda_k(T) \geq \lambda_k(Q(T, \mathbf{U})).$$

Let

$$\begin{aligned} \lambda_1(T) = \dots = \lambda_{n_1}(T) &> \lambda_{(n_1+1)}(T) = \dots = \lambda_{n_2}(T) > \dots > \\ \lambda_{(n_{r-1}+1)}(T) = \dots = \lambda_{n_r}(T) &= \lambda_n(T), \quad n_0 = 0 < n_1 < \dots < n_r = n. \end{aligned} \quad (2.6)$$

Assume that $n_{j-1} < k \leq n_j$. Suppose that $\lambda_k(Q(T, \mathbf{U})) = \lambda_k(T)$. Then for $\mathbf{x} \in \mathbf{U}$ such that $\langle \mathbf{x}, \mathbf{e}_i \rangle = 0$ we have equality $\lambda_k(Q(T, \mathbf{U})) = \lambda_k(T)$ if and only if $\mathbf{x} = \sum_{i=k}^{n_j} x_i \mathbf{e}_i$. Thus $T\mathbf{x} = \lambda_k(T)\mathbf{x}$.

Let $\mathbf{U}_k = \text{span}(\mathbf{e}_1, \dots, \mathbf{e}_k)$. Let $\mathbf{0} \neq \mathbf{x} = \sum_{i=1}^k x_i \mathbf{e}_i \in \mathbf{U}_k$. Then

$$\frac{\langle T\mathbf{x}, \mathbf{x} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle} = \frac{\sum_{i=1}^k \lambda_i(T) |x_i|^2}{\sum_{i=1}^k |x_i|^2} \geq \lambda_k(T) \Rightarrow \lambda_k(Q(T, \mathbf{U}_k)) \geq \lambda_k(T).$$

Hence $\lambda_k(Q(T, \mathbf{U}_k)) = \lambda_k(T)$. □

It can be shown that for $k > 1$ and $\lambda_1(T) > \lambda_k(T)$ there exist $\mathbf{U} \in \text{Gr}(k, \mathbf{V})$ such that $\lambda_k(T) = \lambda_k(Q(T, \mathbf{U}))$ and \mathbf{U} is not an invariant subspace of T , in particular \mathbf{U} does not contain all $\mathbf{e}_1, \dots, \mathbf{e}_k$ satisfying (2.3). (See Problem 1.)

Corollary 2.37 *Let the assumptions of Theorem 2.36 hold. Let $1 \leq \ell \leq n$. Then*

$$\lambda_k(T) = \max_{\mathbf{W} \in \text{Gr}(\ell, \mathbf{V})} \lambda_k(Q(T, \mathbf{W})), \quad k = 1, \dots, \ell. \quad (2.7)$$

Proof. For $k \leq \ell$ apply Theorem 2.36 to $\lambda_k(Q(T, \mathbf{W}))$ to deduce that $\lambda_k(Q(T, \mathbf{W})) \leq \lambda_k(T)$. Let $\mathbf{U}_\ell = \text{span}(\mathbf{e}_1, \dots, \mathbf{e}_\ell)$. Then

$$\lambda_k(Q(T, \mathbf{U}_\ell)) = \lambda_k(T), \quad k = 1, \dots, \ell.$$

□

Theorem 2.38 (*Courant-Fisher principle*) *Let \mathbf{V} be an n -dimensional IPS and $T \in \mathbf{S}(\mathbf{V})$. Then*

$$\lambda_k(T) = \min_{\mathbf{W} \in \text{Gr}(k-1, \mathbf{V})} \max_{\mathbf{0} \neq \mathbf{x} \in \mathbf{W}^\perp} \frac{\langle T\mathbf{x}, \mathbf{x} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle}, \quad k = 1, \dots, n.$$

See Problem 2 for the proof of the theorem and the following corollary.

Corollary 2.39 *Let \mathbf{V} be an n -dimensional IPS and $T \in \mathbf{S}(\mathbf{V})$. Let $k, \ell \in [1, n]$ be integers satisfying $k \leq \ell$. Then*

$$\lambda_{n-\ell+k}(T) \leq \lambda_k(Q(T, \mathbf{W})) \leq \lambda_k(T), \quad \text{for any } \mathbf{W} \in \text{Gr}(\ell, \mathbf{V}).$$

Theorem 2.40 *Let \mathbf{V} be an n -dimensional IPS and $S, T \in \mathbf{S}(\mathbf{V})$. Then for any $i, j \in \mathbb{N}, i + j - 1 \leq n$ the inequality $\lambda_{i+j-1}(S + T) \leq \lambda_i(S) + \lambda_j(T)$ holds.*

Proof. Let $\mathbf{U}_{i-1}, \mathbf{V}_{j-1} \subset \mathbf{V}$ be eigenspaces of S, T spanned by the first $i - 1, j - 1$ eigenvectors of S, T respectively. So

$$\langle S\mathbf{x}, \mathbf{x} \rangle \leq \lambda_i(S)\langle \mathbf{x}, \mathbf{x} \rangle, \quad \langle T\mathbf{y}, \mathbf{y} \rangle \leq \lambda_j(T)\langle \mathbf{y}, \mathbf{y} \rangle \quad \text{for all } \mathbf{x} \in \mathbf{U}_{i-1}^\perp, \mathbf{y} \in \mathbf{V}_{j-1}^\perp.$$

Note that $\dim \mathbf{U}_{i-1} = i - 1, \dim \mathbf{V}_{j-1} = j - 1$. Let $\mathbf{W} = \mathbf{U}_{i-1} + \mathbf{V}_{j-1}$. Then $\dim \mathbf{W} = i + j - 2$. Assume that $\mathbf{z} \in \mathbf{W}^\perp$. Then $\langle (S + T)\mathbf{z}, \mathbf{z} \rangle = \langle S\mathbf{z}, \mathbf{z} \rangle + \langle T\mathbf{z}, \mathbf{z} \rangle \leq (\lambda_i(S) + \lambda_j(T))\langle \mathbf{z}, \mathbf{z} \rangle$. Hence $\max_{\mathbf{0} \neq \mathbf{z} \in \mathbf{W}^\perp} \frac{\langle (S+T)\mathbf{z}, \mathbf{z} \rangle}{\langle \mathbf{z}, \mathbf{z} \rangle} \leq \lambda_i(S) + \lambda_j(T)$. Use Theorem 2.38 to deduce that $\lambda_{i+j-1}(S + T) \leq \lambda_l(S + T) \leq \lambda_i(S) + \lambda_j(T)$. \square

Definition 2.41 *Let \mathbf{V} be an n -dimensional IPS. Fix an integer $k \in [1, n]$. Then $F_k = \{\mathbf{f}_1, \dots, \mathbf{f}_k\}$ is called an orthonormal k -frame if $\langle \mathbf{f}_i, \mathbf{f}_j \rangle = \delta_{ij}$ for $i, j = 1, \dots, k$. Denote by $\text{Fr}(k, \mathbf{V})$ the set of all orthonormal k -frames in \mathbf{V} .*

Note that each $F_k \in \text{Fr}(k, \mathbf{V})$ induces $\mathbf{U} = \text{span}F_k \in \text{Gr}(k, \mathbf{V})$. Vice versa, any $\mathbf{U} \in \text{Gr}(k, \mathbf{V})$ induces the set $\text{Fr}(k, \mathbf{U})$ of all orthonormal k -frames which span \mathbf{U} .

Theorem 2.42 *Let \mathbf{V} be an n -dimensional IPS and $T \in \mathbf{S}(\mathbf{V})$. Then for any integer $k \in [1, n]$*

$$\sum_{i=1}^k \lambda_i(T) = \max_{\{\mathbf{f}_1, \dots, \mathbf{f}_k\} \in \text{Fr}(k, \mathbf{V})} \sum_{i=1}^k \langle T\mathbf{f}_i, \mathbf{f}_i \rangle.$$

Furthermore

$$\sum_{i=1}^k \lambda_i(T) = \sum_{i=1}^k \langle T\mathbf{f}_i, \mathbf{f}_i \rangle$$

for some k -orthonormal frame $F_k = \{\mathbf{f}_1, \dots, \mathbf{f}_k\}$ if and only if $\text{span}F_k$ is spanned by $\mathbf{e}_1, \dots, \mathbf{e}_k$ satisfying (2.3).

Proof. Define

$$\text{tr} Q(T, \mathbf{U}) := \sum_{i=1}^k \lambda_i(Q(T, \mathbf{U})) \quad \text{for } \mathbf{U} \in \text{Gr}(k, \mathbf{V}), \tag{2.8}$$

$$\text{tr}_k T := \sum_{i=1}^k \lambda_i(T).$$

Let $F_k = \{\mathbf{f}_1, \dots, \mathbf{f}_k\} \in \text{Fr}(k, \mathbf{V})$. Set $\mathbf{U} = \text{span}F_k$. Then in view of Corollary 2.37

$$\sum_{i=1}^k \langle T\mathbf{f}_i, \mathbf{f}_i \rangle = \text{tr} Q(T, \mathbf{U}) \leq \sum_{i=1}^k \lambda_i(T).$$

Let $E_k := \{\mathbf{e}_1, \dots, \mathbf{e}_k\}$ where $\mathbf{e}_1, \dots, \mathbf{e}_n$ are given by (2.3). Clearly $\text{tr}_k T = \text{tr } Q(T, \text{span} E_k)$. This shows the maximal characterization of $\text{tr}_k T$.

Let $\mathbf{U} \in \text{Gr}(k, \mathbf{V})$ and assume that $\text{tr}_k T = \text{tr } Q(T, \mathbf{U})$. Hence $\lambda_i(T) = \lambda_i(Q(T, \mathbf{U}))$ for $i = 1, \dots, k$. Then there exists $G_k = \{\mathbf{g}_1, \dots, \mathbf{g}_k\} \in \text{Fr}(k, \mathbf{U})$ such that

$$\min_{\mathbf{0} \neq \mathbf{x} \in \text{span}\{\mathbf{g}_1, \dots, \mathbf{g}_k\}} \frac{\langle T\mathbf{x}, \mathbf{x} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle} = \lambda_i(Q(T, \mathbf{U})) = \lambda_i(T), \quad i = 1, \dots, k.$$

Use Theorem 2.36 to deduce that $T\mathbf{g}_i = \lambda_i(T)\mathbf{g}_i$ for $i = 1, \dots, k$. \square

Theorem 2.43 *Let \mathbf{V} be an n -dimensional IPS and $T \in \mathbf{S}(\mathbf{V})$. Then for any integer $k, l \in [1, n]$, such that $k + l \leq n$*

$$\sum_{i=l+1}^{l+k} \lambda_i(T) = \min_{\mathbf{W} \in \text{Gr}(l, \mathbf{V})} \max_{\{\mathbf{f}_1, \dots, \mathbf{f}_k\} \in \text{Fr}(k, \mathbf{V} \cap \mathbf{W}^\perp)} \sum_{i=1}^k \langle T\mathbf{f}_i, \mathbf{f}_i \rangle.$$

Proof. Let $\mathbf{W}_j := \text{span}\{\mathbf{e}_1, \dots, \mathbf{e}_j\}$, $j = 1, \dots, n$, where $\mathbf{e}_1, \dots, \mathbf{e}_n$ are given by (2.3). Then $\mathbf{V}_1 := \mathbf{V} \cap \mathbf{W}_l$ is an invariant subspace of T . Let $T_1 := T|_{\mathbf{V}_1}$. Then $\lambda_i(T_1) = \lambda_{l+i}(T)$ for $i = 1, \dots, n - l$. Theorem 2.42 for T_1 yields

$$\max_{\{\mathbf{f}_1, \dots, \mathbf{f}_k\} \in \text{Fr}(k, \mathbf{V} \cap \mathbf{W}_l^\perp)} \sum_{i=1}^k \langle T\mathbf{f}_i, \mathbf{f}_i \rangle = \sum_{i=l+1}^{l+k} \lambda_i(T).$$

Let $T_2 := T|_{\mathbf{W}_{l+k}}$ and $\mathbf{W} \in \text{Gr}(l, \mathbf{V})$. Set $\mathbf{U} := \mathbf{W}_{l+k} \cap \mathbf{W}^\perp$. Then $\dim \mathbf{U} \geq k$. Apply Theorem 2.42 to $-T_2$ to deduce

$$\sum_{i=1}^k \lambda_i(-T_2) \geq \sum_{i=1}^k \langle -T\mathbf{f}_i, \mathbf{f}_i \rangle \text{ for } \{\mathbf{f}_1, \dots, \mathbf{f}_k\} \in \text{Fr}(k, \mathbf{U}).$$

The above inequality is equal to the inequality

$$\sum_{i=l+1}^{l+k} \lambda_i(T) \leq \sum_{i=1}^k \langle T\mathbf{f}_i, \mathbf{f}_i \rangle \text{ for } \{\mathbf{f}_1, \dots, \mathbf{f}_k\} \in \text{Fr}(k, \mathbf{U}) \leq \max_{\{\mathbf{f}_1, \dots, \mathbf{f}_k\} \in \text{Fr}(k, \mathbf{V} \cap \mathbf{W}^\perp)} \sum_{i=1}^k \langle T\mathbf{f}_i, \mathbf{f}_i \rangle.$$

The above inequalities yield the theorem. \square

Problems

1. Let \mathbf{V} be 3 dimensional IPS and $T \in \text{Hom}(\mathbf{V})$ be self-adjoint. Assume that

$$\lambda_1(T) > \lambda_2(T) > \lambda_3(T), \quad T\mathbf{e}_i = \lambda_i(T)\mathbf{e}_i, \quad i = 1, 2, 3.$$

Let $\mathbf{W} = \text{span}\{\mathbf{e}_1, \mathbf{e}_3\}$.

(a) Show that for each $t \in (\lambda_3(T), \lambda_1(T))$ there exist two $\mathbf{W}(t) \in \text{Gr}(1, \mathbf{W})$ such that $\lambda_1(Q(T, \mathbf{W}(t))) = t$.

(b) Let $t \in [\lambda_2(T), \lambda_1(T)]$. Let $\mathbf{U}(t) = \text{span}\{\mathbf{W}(t), \mathbf{e}_2\} \in \text{Gr}(2, \mathbf{V})$. Show that $\lambda_2(T) = \lambda_2(Q(T, \mathbf{U}(t)))$.

2. (a) Let the assumptions of Theorem 2.38 hold. Let $\mathbf{W} \in \text{Gr}(k-1, \mathbf{V})$. Show that there exists $\mathbf{0} \neq \mathbf{x} \in \mathbf{W}^\perp$ such that $\langle \mathbf{x}, \mathbf{e}_i \rangle = 0$ for $k+1, \dots, n$, where $\mathbf{e}_1, \dots, \mathbf{e}_n$ satisfy (2.3). Conclude that $\lambda_1(Q(T, \mathbf{W}^\perp)) \geq \frac{\langle T\mathbf{x}, \mathbf{x} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle} \geq \lambda_k(T)$.
- (b) Let $\mathbf{U}_\ell = \text{span}(\mathbf{e}_1, \dots, \mathbf{e}_\ell)$. Show that $\lambda_1(Q(T, \mathbf{U}_\ell^\perp)) = \lambda_{\ell+1}(T)$ for $\ell = 1, \dots, n-1$.
- (c) Prove Theorem 2.38.
- (d) Prove Corollary 2.39. (**Hint:** Choose $\mathbf{U} \in \text{Gr}(k, \mathbf{W})$ such that $\mathbf{U} \subset \mathbf{W} \cap \text{span}(\mathbf{e}_{n-\ell+k+1}, \dots, \mathbf{e}_n)^\perp$. Then $\lambda_{n-\ell+k}(T) \leq \lambda_k(Q(T, \mathbf{U})) \leq \lambda_k(Q(T, \mathbf{W}))$.)
3. Let $B = [b_{ij}]_{i,j=1}^n \in \mathbf{H}_n$ and denote by $A \in \mathbf{H}_{n-1}$ the matrix obtained from B by deleting the j -th row and column.

- (a) Show the Cauchy interlacing inequalities

$$\lambda_i(B) \geq \lambda_i(A) \geq \lambda_{i+1}(B), \text{ for } i = 1, \dots, n-1.$$

- (b) Show that inequality $\lambda_1(B) + \lambda_n(B) \leq \lambda_1(A) + b_{ii}$.

Hint. Express the traces of B and A respectively in terms of eigenvalues to obtain

$$\lambda_1(B) + \lambda_n(B) = b_{ii} + \lambda_1(A) + \sum_{i=2}^{n-1} (\lambda_i(A) - \lambda_i(B)).$$

Then use the Cauchy interlacing inequalities.

4. Show the following generalization of Problem 3.b ([?, p.56]). Let $B \in \mathbf{H}_n$ be the following 2×2 block matrix $B = \begin{bmatrix} B_{11} & B_{12} \\ B_{12}^* & B_{22} \end{bmatrix}$. Show that

$$\lambda_1(B) + \lambda_n(B) \leq \lambda_1(B_{11}) + \lambda_1(B_{22}).$$

Hint. Assume that $B\mathbf{x} = \lambda_1(B)\mathbf{x}$, $\mathbf{x}^\top = (\mathbf{x}_1^\top, \mathbf{x}_2^\top)$, partitioned as B . Consider $\mathbf{U} = \text{span}((\mathbf{x}_1^\top, \mathbf{0})^\top, (\mathbf{0}, \mathbf{x}_2^\top)^\top)$. Analyze $\lambda_1(Q(T, \mathbf{U})) + \lambda_2(Q(T, \mathbf{U}))$.

5. Let $B = (b_{ij})_1^n \in \mathbf{H}_n$. Show that $B > 0$ if and only if $\det(b_{ij})_1^k > 0$ for $k = 1, \dots, n$.
6. Let $T \in \mathbf{S}(\mathbf{V})$. Denote by $\iota_+(T), \iota_0(T), \iota_-(T)$ the number of positive, negative and zero eigenvalues among $\lambda_1(T) \geq \dots \geq \lambda_n(T)$. The triple $\iota(T) := (\iota_+(T), \iota_0(T), \iota_-(T))$ is called the inertia of T . For $B \in \mathbf{H}_n$ let $\iota(B) := (\iota_+(B), \iota_0(B), \iota_-(B))$ be the inertia of B , where $\iota_+(B), \iota_0(B), \iota_-(B)$ is the number of positive, negative and zero eigenvalues of B respectively. Let $\mathbf{U} \in \text{Gr}(k, \mathbf{V})$. Show
- (a) Assume that $\lambda_k(Q(T, \mathbf{U})) > 0$, i.e. $Q(T, \mathbf{U}) > 0$. Then $k \leq \iota_+(T)$. If $k = \iota_+(T)$ then one can choose \mathbf{U} to be an invariant subspace of \mathbf{V} spanned by the eigenvectors of T corresponding to positive eigenvalues of T . (Usually such a subspace is not unique.)
- (b) Assume that $\lambda_k(Q(T, \mathbf{U})) \geq 0$, i.e. $Q(T, \mathbf{U}) \geq 0$. Then $k \leq \iota_+(T) + \iota_0(T)$. If $k = \iota_+(T) + \iota_0(T)$ then \mathbf{U} is the unique invariant subspace of \mathbf{V} spanned by the eigenvectors of T corresponding to nonnegative eigenvalues of T .

(c) Assume that $\lambda_1(Q(T, \mathbf{U})) < 0$, i.e. $Q(T, \mathbf{U}) < 0$. Then $k \leq \iota_-(T)$. If $k = \iota_-(T)$ then \mathbf{U} can be chosen to be an invariant subspace of \mathbf{V} spanned by the eigenvectors of T , corresponding to negative eigenvalues of T . (Usually such a subspace may not be unique.)

(d) Assume that $\lambda_1(Q(T, \mathbf{U})) \leq 0$, i.e. $Q(T, \mathbf{U}) \leq 0$. Then $k \leq \iota_-(T) + \iota_0(T)$. If $k = \iota_-(T) + \iota_0(T)$ then \mathbf{U} is a unique invariant subspace of \mathbf{V} spanned by the eigenvectors of T corresponding to nonpositive eigenvalues of T .

7. Let $B \in \mathbf{H}_n$ and assume that $A = PBP^*$ for some $P \in \text{GL}(n, \mathbb{C})$. Then $\iota(A) = \iota(B)$.

2.6 Positive definite operators and matrices

Definition 2.44 Let \mathbf{V} be a finite dimensional IPS over $\mathbb{F} = \mathbb{C}, \mathbb{R}$. Let $S, T \in \mathbf{S}(\mathbf{V})$. Then $T > S$, ($T \geq S$) if $\langle T\mathbf{x}, \mathbf{x} \rangle > \langle S\mathbf{x}, \mathbf{x} \rangle$, ($\langle T\mathbf{x}, \mathbf{x} \rangle \geq \langle S\mathbf{x}, \mathbf{x} \rangle$) for all $\mathbf{0} \neq \mathbf{x} \in \mathbf{V}$. T is called positive (nonnegative) definite if $T > 0$ ($T \geq 0$), where 0 is the zero operator in $\text{Hom}(\mathbf{V})$.

Denote by $\mathbf{S}_+(\mathbf{V})^o \subset \mathbf{S}_+(\mathbf{V}) \subset \mathbf{S}(\mathbf{V})$ the open set of positive definite self adjoint operators and the closed set of nonnegative selfadjoint operators respectively.

Let P, Q be either quadratic forms if $\mathbb{F} = \mathbb{R}$ or hermitian forms if $\mathbb{F} = \mathbb{C}$. Then $Q > P$, ($Q \geq P$) if $Q(\mathbf{x}, \mathbf{x}) > P(\mathbf{x}, \mathbf{x})$, ($Q(\mathbf{x}, \mathbf{x}) \geq P(\mathbf{x}, \mathbf{x})$) for all $\mathbf{0} \neq \mathbf{x} \in \mathbf{V}$. Q is called positive (nonnegative) definite if $Q > 0$ ($Q \geq 0$), where 0 is the zero operator in $\text{Hom}(\mathbf{V})$.

For $A, B \in \mathbf{H}_n$ $B > A$ ($B \geq A$) if $\mathbf{x}^*B\mathbf{x} > \mathbf{x}^*A\mathbf{x}$ ($\mathbf{x}^*B\mathbf{x} \geq \mathbf{x}^*A\mathbf{x}$) for all $\mathbf{0} \neq \mathbf{x} \in \mathbb{C}^n$. $B \in \mathbf{H}_n$ is called is called positive (nonnegative) definite if $B > 0$ ($B \geq 0$). Denote by $\mathbf{H}_{n,+}^o \subset \mathbf{H}_{n,+} \subset \mathbf{H}_n$ the open set of positive definite $n \times n$ hermitian matrices and the closed set of $n \times n$ nonnegative hermitian matrices respectively. Let $\mathbf{S}_+(n, \mathbb{R}) := \mathbf{S}(n, \mathbb{R}) \cap \mathbf{H}_{n,+}$, $\mathbf{S}_+(n, \mathbb{R})^o := \mathbf{S}(n, \mathbb{R}) \cap \mathbf{H}_{n,+}^o$.

Use (2.1) to deduce.

Corollary 2.45 Let \mathbf{V} be n -dimensional IPS. Let $T \in \mathbf{S}(\mathbf{V})$. Then $T > 0$ ($T \geq 0$) if and only if $\lambda_n(T) > 0$ ($\lambda_n(T) \geq 0$). Let $S \in \mathbf{S}(\mathbf{V})$ and assume that $T > S$ ($T \geq S$). Then $\lambda_i(T) > \lambda_i(S)$ ($\lambda_i(T) \geq \lambda_i(S)$) for $i = 1, \dots, n$.

Proposition 2.46 Let \mathbf{V} be a finite dimensional IPS. Assume that $T \in \mathbf{S}(\mathbf{V})$. Then $T \geq 0$ if and only if there exists $S \in \mathbf{S}(\mathbf{V})$ such that $T = S^2$. Furthermore $T > 0$ if and only if S is invertible. For $0 \leq T \in \mathbf{S}(\mathbf{V})$ there exists a unique $0 \leq S \in \mathbf{S}(\mathbf{V})$ such that $T = S^2$. This S is called the square root of T and is denoted by $T^{\frac{1}{2}}$.

Proof. Assume first that $T \geq 0$. Let $\mathbf{e}_1, \dots, \mathbf{e}_n$ be an orthonormal basis consisting of eigenvectors of T as in (2.3). Since $\lambda_i(T) \geq 0$, $i = 1, \dots, n$ we can define $P \in \text{Hom}(\mathbf{V})$ as follows

$$P\mathbf{e}_i = \sqrt{\lambda_i(T)}\mathbf{e}_i, \quad i = 1, \dots, n.$$

Clearly P is self-adjoint nonnegative and $T = P^2$.

Suppose now that $T = S^2$ for some $S \in \mathbf{S}(\mathbf{V})$. Then $T \in \mathbf{S}(\mathbf{V})$ and $\langle T\mathbf{x}, \mathbf{x} \rangle = \langle S\mathbf{x}, S\mathbf{x} \rangle \geq 0$. Hence $T \geq 0$. Clearly $\langle T\mathbf{x}, \mathbf{x} \rangle = 0 \iff S\mathbf{x} = 0$. Hence

$T > 0 \iff S \in \text{GL}(\mathbf{V})$. Suppose that $S \geq 0$. Then $\lambda_i(S) = \sqrt{\lambda_i(T)}$, $i = 1, \dots, n$. Furthermore each eigenvector of S is an eigenvector of T . It is straightforward to show that $S = P$, where P is defined above. Clearly $T > 0$ if and only if $\sqrt{\lambda_n(T)} > 0$, i.e. if and only if S is invertible. \square

Corollary 2.47 *Let $B \in \mathbf{H}_n(\mathbf{S}(n, \mathbb{R}))$. Then $B \geq 0$ if and only there exists $A \in \mathbf{H}_n(\mathbf{S}(n, \mathbb{R}))$ such that $B = A^2$. Furthermore $B > 0$ if and only if A is invertible. For $B \geq 0$ there exists a unique $A \geq 0$ such that $B = A^2$. This A is denoted by $B^{\frac{1}{2}}$.*

Theorem 2.48 *Let \mathbf{V} be an IPS over $\mathbb{F} = \mathbb{C}, \mathbb{R}$. Let $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbf{V}$. Then the grammian matrix $G(\mathbf{x}_1, \dots, \mathbf{x}_n) := (\langle \mathbf{x}_i, \mathbf{x}_j \rangle)_{i,j=1}^n$ is a hermitian nonnegative definite matrix. (If $\mathbb{F} = \mathbb{R}$ then $G(\mathbf{x}_1, \dots, \mathbf{x}_n)$ is real symmetric nonnegative definite.) $G(\mathbf{x}_1, \dots, \mathbf{x}_n) > 0$ if and only $\mathbf{x}_1, \dots, \mathbf{x}_n$ are linearly independent. Furthermore for any integer $k \in [1, n-1]$*

$$\det G(\mathbf{x}_1, \dots, \mathbf{x}_n) \leq \det G(\mathbf{x}_1, \dots, \mathbf{x}_k) \det G(\mathbf{x}_{k+1}, \dots, \mathbf{x}_n). \quad (2.1)$$

Equality holds if and only if either $\det G(\mathbf{x}_1, \dots, \mathbf{x}_k) = 0$ or $\langle \mathbf{x}_i, \mathbf{x}_j \rangle = 0$ for $i = 1, \dots, k$ and $j = k+1, \dots, n$.

Proof. Clearly $G(\mathbf{x}_1, \dots, \mathbf{x}_n) \in \mathbf{H}_n$. If \mathbf{V} is an IPS over \mathbb{R} then $G(\mathbf{x}_1, \dots, \mathbf{x}_n) \in \mathbf{S}(n, \mathbb{R})$. Let $\mathbf{a} = (a_1, \dots, a_n)^\top \in \mathbb{F}^n$. Then

$$\mathbf{a}^* G(\mathbf{x}_1, \dots, \mathbf{x}_n) \mathbf{a} = \left\langle \sum_{i=1}^n a_i \mathbf{x}_i, \sum_{j=1}^n a_j \mathbf{x}_j \right\rangle \geq 0.$$

Equality holds if and only if $\sum_{i=1}^n a_i \mathbf{x}_i = 0$. Hence $G(\mathbf{x}_1, \dots, \mathbf{x}_n) \geq 0$ and $G(\mathbf{x}_1, \dots, \mathbf{x}_n) > 0$ if and only if $\mathbf{x}_1, \dots, \mathbf{x}_n$ are linearly independent. In particular $\det G(\mathbf{x}_1, \dots, \mathbf{x}_n) \geq 0$ and $\det G(\mathbf{x}_1, \dots, \mathbf{x}_n) > 0$ if and only if $\mathbf{x}_1, \dots, \mathbf{x}_n$ are linearly independent.

We now prove the inequality (2.1). Assume first that the right-hand side of (2.1) is zero. Then either $\mathbf{x}_1, \dots, \mathbf{x}_k$ or $\mathbf{x}_{k+1}, \dots, \mathbf{x}_n$ are linearly dependent. Hence $\mathbf{x}_1, \dots, \mathbf{x}_n$ are linearly dependent and $\det G = 0$.

Assume now that the right-hand side of (2.1) is positive. Hence $\mathbf{x}_1, \dots, \mathbf{x}_k$ and $\mathbf{x}_{k+1}, \dots, \mathbf{x}_n$ are linearly independent. If $\mathbf{x}_1, \dots, \mathbf{x}_n$ are linearly dependent then $\det G = 0$ and strict inequality holds in (2.1). It is left to show the inequality (2.1) and the equality case when $\mathbf{x}_1, \dots, \mathbf{x}_n$ are linearly independent. Perform the Gram-Schmidt algorithm on $\mathbf{x}_1, \dots, \mathbf{x}_n$ as given in (2.1). Let $S_j = \text{span}(\mathbf{x}_1, \dots, \mathbf{x}_j)$ for $j = 1, \dots, n$. Corollary 2.1 yields that $\text{span}(\mathbf{e}_1, \dots, \mathbf{e}_{n-1}) = S_{n-1}$. Hence $\mathbf{y}_n = \mathbf{x}_n - \sum_{j=1}^{n-1} b_j \mathbf{x}_j$ for some $b_1, \dots, b_{n-1} \in \mathbb{F}$. Let G' be the matrix obtained from $G(\mathbf{x}_1, \dots, \mathbf{x}_n)$ by subtracting from the n -th row b_j times j -th row. Thus the last row of G' is $(\langle \mathbf{y}_n, \mathbf{x}_1 \rangle, \dots, \langle \mathbf{y}_n, \mathbf{x}_n \rangle) = (0, \dots, 0, \|\mathbf{y}_n\|^2)$. Clearly $\det G(\mathbf{x}_1, \dots, \mathbf{x}_n) = \det G'$. Expand $\det G'$ by the last row to deduce

$$\begin{aligned} \det G(\mathbf{x}_1, \dots, \mathbf{x}_n) &= \det G(\mathbf{x}_1, \dots, \mathbf{x}_{n-1}) \|\mathbf{y}_n\|^2 = \dots = \\ \det G(\mathbf{x}_1, \dots, \mathbf{x}_k) &\prod_{i=k+1}^n \|\mathbf{y}_i\|^2 = \\ \det G(\mathbf{x}_1, \dots, \mathbf{x}_k) &\prod_{i=k+1}^n \text{dist}(\mathbf{x}_i, S_{i-1})^2, \quad k = n-1, \dots, 1. \end{aligned} \quad (2.2)$$

Perform the Gram-Schmidt process on $\mathbf{x}_{k+1}, \dots, \mathbf{x}_n$ to obtain the orthogonal set of vectors $\hat{\mathbf{y}}_{k+1}, \dots, \hat{\mathbf{y}}_n$ such that

$$\hat{S}_j := \text{span}(\mathbf{x}_{k+1}, \dots, \mathbf{x}_j) = \text{span}(\hat{\mathbf{y}}_{k+1}, \dots, \hat{\mathbf{y}}_j), \quad \text{dist}(\mathbf{x}_j, \hat{S}_{j-1}) = \|\hat{\mathbf{y}}_j\|,$$

for $j = k + 1, \dots, n$, where $\hat{S}_k = \{\mathbf{0}\}$. Use (2.2) to deduce that $\det G(\mathbf{x}_{k+1}, \dots, \mathbf{x}_n) = \prod_{j=k+1}^n \|\hat{\mathbf{y}}_j\|^2$. As $\hat{S}_{j-1} \subset S_{j-1}$ for $j > k$ it follows that

$$\|\mathbf{y}_j\| = \text{dist}(\mathbf{x}_j, S_{j-1}) \leq \text{dist}(\mathbf{x}_j, \hat{S}_{j-1}) = \|\hat{\mathbf{y}}_j\|, \quad j = k + 1, \dots, n.$$

This shows (2.1). Assume now equality holds in (2.1). Then $\|\mathbf{y}_j\| = \|\hat{\mathbf{y}}_j\|$ for $j = k + 1, \dots, n$. Since $\hat{S}_{j-1} \subset S_{j-1}$ and $\hat{\mathbf{y}}_j - \mathbf{x}_j \in \hat{S}_{j-1} \subset S_{j-1}$ it follows that $\text{dist}(\mathbf{x}_j, S_{j-1}) = \text{dist}(\hat{\mathbf{y}}_j, S_{j-1}) = \|\mathbf{y}_j\|$. Hence $\|\hat{\mathbf{y}}_j\| = \text{dist}(\hat{\mathbf{y}}_j, S_{j-1})$. Part (h) of Problem 2.1.2.5 yields that $\hat{\mathbf{y}}_j$ is orthogonal on S_{j-1} . In particular each $\hat{\mathbf{y}}_j$ is orthogonal to S_k for $j = k + 1, \dots, n$. Hence $\mathbf{x}_j \perp S_k$ for $j = k + 1, \dots, n$, i.e. $\langle \mathbf{x}_j, \mathbf{x}_i \rangle = 0$ for $j > k$ and $i \leq k$. Clearly, if the last condition holds then $\det G(\mathbf{x}_1, \dots, \mathbf{x}_n) = \det G(\mathbf{x}_1, \dots, \mathbf{x}_k) \det G(\mathbf{x}_{k+1}, \dots, \mathbf{x}_n)$. \square

$\det G(\mathbf{x}_1, \dots, \mathbf{x}_n)$ has the following geometric meaning. Consider a parallelepiped Π in \mathbf{V} spanned by $\mathbf{x}_1, \dots, \mathbf{x}_n$ starting from the origin $\mathbf{0}$. That is Π is a convex hull spanned by the vectors $\mathbf{0}$ and $\sum_{i \in S} \mathbf{x}_i$ for all nonempty subsets $S \subset \{1, \dots, n\}$. Then $\sqrt{\det G(\mathbf{x}_1, \dots, \mathbf{x}_n)}$ is the n -volume of Π . The inequality (2.1) and equalities (2.2) are "obvious" from this geometrical point of view.

Corollary 2.49 *Let $0 \leq B = (b_{ij})_1^n \in \mathbf{H}_{n,+}$. Then*

$$\det B \leq \det(b_{ij})_1^k \det(b_{ij})_{k+1}^n, \quad \text{for } k = 1, \dots, n - 1.$$

For a fixed k equality holds if and only if either the right-hand side of the above inequality is zero or $b_{ij} = 0$ for $i = 1, \dots, k$ and $j = k + 1, \dots, n$.

Proof. From Corollary 2.47 it follows that $B = X^2$ for some $X \in \mathbf{H}_n$. Let $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{C}^n$ be the n -columns of $X^T = (\mathbf{x}_1, \dots, \mathbf{x}_n)$. Let $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{y}^* \mathbf{x}$. Since $X \in \mathbf{H}_n$ we deduce that $B = G(\mathbf{x}_1, \dots, \mathbf{x}_n)$. \square

Theorem 2.50 *Let \mathbf{V} be an n -dimensional IPS. Let $T \in \mathbf{S}$. TFAE:*

- (a) $T > 0$.
- (b) *Let $\mathbf{g}_1, \dots, \mathbf{g}_n$ be a basis of \mathbf{V} . Then $\det(\langle T \mathbf{g}_i, \mathbf{g}_j \rangle)_{i,j=1}^k > 0$, $k = 1, \dots, n$.*

Proof. (a) \Rightarrow (b). According to Proposition 2.46 $T = S^2$ for some $S \in \mathbf{S}(\mathbf{V}) \cap \text{GL}(\mathbf{V})$. Then $\langle T \mathbf{g}_i, \mathbf{g}_j \rangle = \langle S \mathbf{g}_i, S \mathbf{g}_j \rangle$. Hence $\det(\langle T \mathbf{g}_i, \mathbf{g}_j \rangle)_{i,j=1}^k = \det G(S \mathbf{g}_1, \dots, S \mathbf{g}_k)$. Since S is invertible and $\mathbf{g}_1, \dots, \mathbf{g}_k$ linearly independent it follows that $S \mathbf{g}_1, \dots, S \mathbf{g}_k$ are linearly independent. Theorem 2.1 implies that $\det G(S \mathbf{g}_1, \dots, S \mathbf{g}_k) > 0$ for $k = 1, \dots, n$.

(b) \Rightarrow (a). The proof is by induction on n . For $n = 1$ (a) is obvious. Assume that (a) holds for $n = m - 1$. Let $\mathbf{U} := \text{span}(\mathbf{g}_1, \dots, \mathbf{g}_{n-1})$ and $Q := Q(T, \mathbf{U})$. Then there exists $P \in \mathbf{S}(\mathbf{U})$ such that $\langle P \mathbf{x}, \mathbf{y} \rangle = Q(\mathbf{x}, \mathbf{y}) = \langle T \mathbf{x}, \mathbf{y} \rangle$ for any $\mathbf{x}, \mathbf{y} \in \mathbf{U}$. By induction $P > 0$. Corollary 2.37 yields that $\lambda_{n-1}(T) \geq \lambda_{n-1}(P) > 0$. Hence T has at least $n - 1$ positive eigenvalues. Let $\mathbf{e}_1, \dots, \mathbf{e}_n$ be given by (2.3).

Then $\det(\langle T\mathbf{e}_i, \mathbf{e}_j \rangle)_{i,j=1}^n = \prod_{i=1}^n \lambda_i(T) > 0$. Let $A = (a_{pq})_1^n \in \text{GL}(n, \mathbb{C})$ be the transformation matrix from the basis $\mathbf{g}_1, \dots, \mathbf{g}_n$ to $\mathbf{e}_1, \dots, \mathbf{e}_n$, i.e.

$$\mathbf{g}_i = \sum_{p=1}^n a_{pi} \mathbf{e}_p, \quad i = 1, \dots, n.$$

It is straightforward to show that

$$(\langle T\mathbf{g}_i, \mathbf{g}_j \rangle)_1^n = A^T (\langle T\mathbf{e}_p, \mathbf{e}_q \rangle) \bar{A} \Rightarrow \quad (2.3)$$

$$\det(\langle T\mathbf{g}_i, \mathbf{g}_j \rangle)_1^n = \det(\langle T\mathbf{e}_i, \mathbf{e}_j \rangle)_1^n |\det A|^2 = |\det A|^2 \prod_{i=1}^n \lambda_i(T).$$

Since $\det(\langle T\mathbf{g}_i, \mathbf{g}_j \rangle)_1^n > 0$ and $\lambda_1(T) \geq \dots \geq \lambda_{n-1}(T) > 0$ it follows that $\lambda_n(T) > 0$. \square

Corollary 2.51 *Let $B = (b_{ij})_1^n \in \mathbf{H}_n$. Then $B > 0$ if and only if $\det(b_{ij})_1^k > 0$ for $k = 1, \dots, n$.*

The following result is straightforward (see Problem 1):

Proposition 2.52 *Let \mathbf{V} be a finite dimensional IPS over $\mathbb{F} = \mathbb{R}, \mathbb{C}$ with the inner product $\langle \cdot, \cdot \rangle$. Assume that $T \in \mathbf{S}(\mathbf{V})$. Then $T > 0$ if and only if $(\mathbf{x}, \mathbf{y}) := \langle T\mathbf{x}, \mathbf{y} \rangle$ is an inner product on \mathbf{V} . Vice versa any inner product $(\cdot, \cdot) : \mathbf{V} \times \mathbf{V} \rightarrow \mathbb{R}$ is of the form $(\mathbf{x}, \mathbf{y}) = \langle T\mathbf{x}, \mathbf{y} \rangle$ for a unique self-adjoint positive definite operator $T \in \text{Hom}(\mathbf{V})$.*

Example 2.53 *Each $0 < B \in \mathbf{H}_n$ induces an inner product on \mathbb{C}^n : $(\mathbf{x}, \mathbf{y}) = \mathbf{y}^* B \mathbf{x}$. Each $0 < B \in \mathbf{S}(n, \mathbb{R})$ induces an inner product on \mathbb{R}^n : $(\mathbf{x}, \mathbf{y}) = \mathbf{y}^T B \mathbf{x}$. Furthermore any inner product on \mathbb{C}^n or \mathbb{R}^n is of the above form. In particular, the standard inner products on \mathbb{C}^n and \mathbb{R}^n are induced by the identity matrix I .*

Definition 2.54 *Let \mathbf{V} be a finite dimensional IPS with the inner product $\langle \cdot, \cdot \rangle$. Let $S \in \text{Hom}(\mathbf{V})$. Then S is called symmetrizable if there exists an inner product (\cdot, \cdot) on \mathbf{V} such that S is self-adjoint with respect to (\cdot, \cdot) .*

Problems

1. Show Proposition 2.52.
2. Recall the Hölder inequality

$$\sum_{l=1}^n x_l y_l a_l \leq \left(\sum_{l=1}^n x_l^p a_l \right)^{\frac{1}{p}} \left(\sum_{l=1}^n y_l^q a_l \right)^{\frac{1}{q}} \quad (2.4)$$

for any $\mathbf{x} = (x_1, \dots, x_n)^\top, \mathbf{y} = (y_1, \dots, y_n)^\top, \mathbf{a} = (a_1, \dots, a_n) \in \mathbb{R}_+^n$ and $p, q \in (1, \infty)$ such that $\frac{1}{p} + \frac{1}{q} = 1$. Show

(a) Let $A \in \mathbf{H}_{n,+}$, $\mathbf{x} \in \mathbb{C}^n$ and $0 \leq i < j < k$ be three integers. Then

$$\mathbf{x}^* A^j \mathbf{x} \leq (\mathbf{x}^* A^i \mathbf{x})^{\frac{k-j}{k-i}} (\mathbf{x}^* A^k \mathbf{x})^{\frac{j-i}{k-i}}. \quad (2.5)$$

Hint: Diagonalize A .

(b) Assume that $A = e^B$ for some $B \in \mathbf{H}_n$. Show that (2.5) holds for any three real numbers $i < j < k$.

2.7 Inequalities for traces

Let \mathbf{V} be a finite dimensional IPS over $\mathbb{F} = \mathbb{R}, \mathbb{C}$. Let $T : \mathbf{V} \rightarrow \mathbf{V}$ be a linear operator. Then $\text{tr } T$ is the trace of the representation matrix A of with respect to any orthonormal basis of \mathbf{V} . See Problem 1.

Theorem 2.55 *Let \mathbf{V} be an n -dimensional IPS over $\mathbb{F} = \mathbb{R}, \mathbb{C}$. Assume that $S, T \in \mathbf{S}(\mathbf{V})$. Then $\text{tr } ST$ is bounded below and above by*

$$\sum_{i=1}^n \lambda_i(S) \lambda_{n-i+1}(T) \leq \text{tr } ST \leq \sum_{i=1}^n \lambda_i(S) \lambda_i(T). \quad (2.1)$$

Equality for the upper bound holds if and only if $ST = TS$ and there exists an orthonormal basis $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbf{V}$ such that

$$S\mathbf{x}_i = \lambda_i(S)\mathbf{x}_i, \quad T\mathbf{x}_i = \lambda_i(T)\mathbf{x}_i, \quad i = 1, \dots, n. \quad (2.2)$$

Equality for the lower bound holds if and only if $ST = TS$ and there exists an orthonormal basis $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbf{V}$ such that

$$S\mathbf{x}_i = \lambda_i(S)\mathbf{x}_i, \quad T\mathbf{x}_i = \lambda_{n-i+1}(T)\mathbf{x}_i, \quad i = 1, \dots, n. \quad (2.3)$$

Proof. Let $\mathbf{y}_1, \dots, \mathbf{y}_n$ be an orthonormal basis of \mathbf{V} such that

$$\begin{aligned} T\mathbf{y}_i &= \lambda_i(T)\mathbf{y}_i, \quad i = 1, \dots, n, \\ \lambda_1(T) &= \dots = \lambda_{i_1}(T) > \lambda_{i_1+1}(T) = \dots = \lambda_{i_2}(T) > \dots > \\ \lambda_{i_{k-1}+1}(T) &= \dots = \lambda_{i_k}(T) = \lambda_n(T), \quad 1 \leq i_1 < \dots < i_k = n. \end{aligned}$$

If $k = 1 \iff i_1 = n$ it follows that $T = \lambda_1 I$ and the theorem is trivial in this case. Assume that $k > 1$. Then

$$\begin{aligned} \text{tr } ST &= \sum_{i=1}^n \lambda_i(T) \langle S\mathbf{y}_i, \mathbf{y}_i \rangle = \\ &= \sum_{i=1}^{n-1} (\lambda_i(T) - \lambda_{i+1}(T)) \left(\sum_{l=1}^i \langle S\mathbf{y}_l, \mathbf{y}_l \rangle \right) + \lambda_n(T) \left(\sum_{l=1}^n \langle S\mathbf{y}_l, \mathbf{y}_l \rangle \right) = \\ &= \sum_{j=1}^{k-1} (\lambda_{i_j}(T) - \lambda_{i_{j+1}}(T)) \sum_{l=1}^{i_j} \langle S\mathbf{y}_l, \mathbf{y}_l \rangle + \lambda_n(T) \text{tr } S. \end{aligned}$$

Theorem 2.42 yields that $\sum_{l=1}^{i_j} \langle S\mathbf{y}_l, \mathbf{y}_l \rangle \leq \sum_{l=1}^{i_j} \lambda_l(S)$. Substitute these inequalities for $j = 1, \dots, k-1$ in the above identity to deduce the upper bound in (2.1).

Clearly the condition (2.2) implies that $\text{tr } ST$ is equal to the upper bound in (2.1). Assume now that $\text{tr } ST$ is equal to the upper bound in (2.1). Then $\sum_{l=1}^{i_j} \langle S\mathbf{y}_l, \mathbf{y}_l \rangle = \sum_{l=1}^{i_j} \lambda_l(S)$ for $j = 1, \dots, k-1$. Theorem 2.42 yields that $\text{span}(\mathbf{y}_1, \dots, \mathbf{y}_{i_j})$ is spanned by some i_j eigenvectors of S corresponding to the first i_j eigenvalues of S for $j = 1, \dots, k-1$. Let $\mathbf{x}_1, \dots, \mathbf{x}_{i_1}$ be an orthonormal basis of $\text{span}(\mathbf{y}_1, \dots, \mathbf{y}_{i_1})$ consisting of the eigenvectors of S corresponding to the eigenvalues of $\lambda_1(S), \dots, \lambda_{i_1}(S)$. Since any $0 \neq \mathbf{x} \in \text{span}(\mathbf{y}_1, \dots, \mathbf{y}_{i_1})$ is an eigenvector of T corresponding to the eigenvalue $\lambda_{i_1}(T)$ it follows that (2.2) holds for $i = 1, \dots, i_1$. Consider $\text{span}(\mathbf{y}_1, \dots, \mathbf{y}_{i_2})$. The above arguments imply that this subspace contains i_2 eigenvectors of S and T corresponding to the first i_2 eigenvalues of S and T . Hence \mathbf{U}_2 , the orthogonal complement of $\text{span}(\mathbf{x}_1, \dots, \mathbf{x}_{i_1})$ in $\text{span}(\mathbf{y}_1, \dots, \mathbf{y}_{i_2})$, spanned by $\mathbf{x}_{i_1+1}, \dots, \mathbf{x}_{i_2}$, which are $i_2 - i_1$ orthonormal eigenvectors of S corresponding to the eigenvalues $\lambda_{i_1+1}(S), \dots, \lambda_{i_2}(S)$. Since any nonzero vector in \mathbf{U}_2 is an eigenvector of T corresponding to the eigenvalue $\lambda_{i_2}(T)$ we deduce that (2.2) holds for $i = 1, \dots, i_2$. Continuing in the same manner we obtain (2.2).

To prove the equality case in the lower bound consider the equality in the upper bound for $\text{tr } S(-T)$. \square

Corollary 2.56 *Let \mathbf{V} be an n -dimensional IPS over $\mathbb{F} = \mathbb{R}, \mathbb{C}$. Assume that $S, T \in \mathbf{S}(\mathbf{V})$. Then*

$$\sum_{i=1}^n (\lambda_i(S) - \lambda_i(T))^2 \leq \text{tr}(S - T)^2. \quad (2.4)$$

Equality holds if and only if $ST = TS$ and \mathbf{V} has an orthonormal basis $\mathbf{x}_1, \dots, \mathbf{x}_n$ satisfying (2.2).

Proof. Note

$$\sum_{i=1}^n (\lambda_i(S) - \lambda_i(T))^2 = \text{tr } S^2 + \text{tr } T^2 - 2 \sum_{i=1}^n \lambda_i(S)\lambda_i(T).$$

\square

Corollary 2.57 *Let $S, T \in \mathbf{H}_n$. Then the inequalities (2.1) and (2.4) hold. Equalities in the upper bounds hold if and only if there exists $U \in \mathbf{U}_n$ such that $S = U \text{diag } \lambda(S)U^*, T = U \text{diag } \lambda(T)U^*$. Equality in the lower bound of (2.1) if and only if there exists $V \in \mathbf{U}_n$ such that $S = V \text{diag } \lambda(S)V^*, -T = V \text{diag } \lambda(-T)V^*$.*

Problems

1. Let \mathbf{V} be a n -dimensional IPS over $\mathbb{F} = \mathbb{R}, \mathbb{C}$.

- (a) Assume that $T : \mathbf{V} \rightarrow \mathbf{V}$ be a linear transformation. Show that for any o.n. basis $\mathbf{x}_1, \dots, \mathbf{x}_n$

$$\text{tr } T = \sum_{i=1}^n \langle T\mathbf{x}_i, \mathbf{x}_i \rangle.$$

Furthermore, if $\mathbb{F} = \mathbb{C}$ then $\text{tr } T$ is the sum of the n eigenvalues of T .

- (b) Let $S, T \in \mathbf{S}(\mathbf{V})$. Show that $\text{tr } ST = \text{tr } TS \in \mathbb{R}$.

2.8 Singular Value Decomposition

Let \mathbf{U}, \mathbf{V} , be finite dimensional IPS over $\mathbb{F} = \mathbb{R}, \mathbb{C}$, with the inner products $\langle \cdot, \cdot \rangle_{\mathbf{U}}, \langle \cdot, \cdot \rangle_{\mathbf{V}}$ respectively. Let $\mathbf{u}_1, \dots, \mathbf{u}_m$ and $\mathbf{v}_1, \dots, \mathbf{v}_n$ be bases in \mathbf{U} and \mathbf{V} respectively. Let $T : \mathbf{V} \rightarrow \mathbf{U}$ be a linear operator. In these bases T is represented by a matrix $A = [a_{ij}] \in \mathbb{F}^{m \times n}$ as given by

$$T\mathbf{v}_j = \sum_{i=1}^m a_{ij}\mathbf{u}_i, \quad j = 1, \dots, n.$$

Let $T^* : \mathbf{U}^* = \mathbf{U} \rightarrow \mathbf{V}^* = \mathbf{V}$. Then $T^*T : \mathbf{V} \rightarrow \mathbf{V}$ and $TT^* : \mathbf{U} \rightarrow \mathbf{U}$ are selfadjoint operators. As

$$\langle T^*T\mathbf{v}, \mathbf{v} \rangle_{\mathbf{V}} = \langle T\mathbf{v}, T\mathbf{v} \rangle_{\mathbf{U}} \geq 0, \quad \langle TT^*\mathbf{u}, \mathbf{u} \rangle_{\mathbf{U}} = \langle T^*\mathbf{u}, T^*\mathbf{u} \rangle_{\mathbf{V}} \geq 0$$

it follows that $T^*T \geq 0, TT^* \geq 0$. Let

$$T^*T\mathbf{c}_i = \lambda_i(T^*T)\mathbf{c}_i, \quad \langle \mathbf{c}_i, \mathbf{c}_k \rangle_{\mathbf{V}} = \delta_{ik}, \quad i, k = 1, \dots, n, \quad (2.1)$$

$$\lambda_1(T^*T) \geq \dots \geq \lambda_n(T^*T) \geq 0,$$

$$TT^*\mathbf{d}_j = \lambda_j(TT^*)\mathbf{d}_j, \quad \langle \mathbf{d}_j, \mathbf{d}_l \rangle_{\mathbf{U}} = \delta_{jl}, \quad j, l = 1, \dots, m, \quad (2.2)$$

$$\lambda_1(TT^*) \geq \dots \geq \lambda_m(TT^*) \geq 0,$$

Proposition 2.58 *Let \mathbf{U}, \mathbf{V} , be finite dimensional IPS over $\mathbb{F} = \mathbb{R}, \mathbb{C}$. Let $T : \mathbf{V} \rightarrow \mathbf{U}$. Then $\text{rank } T = \text{rank } T^* = \text{rank } T^*T = \text{rank } TT^* = r$. Furthermore the selfadjoint nonnegative definite operators T^*T and TT^* have exactly r positive eigenvalues, and*

$$\lambda_i(T^*T) = \lambda_i(TT^*) > 0, \quad i = 1, \dots, \text{rank } T. \quad (2.3)$$

Moreover for $i \in [1, r]$ $T\mathbf{c}_i$ and $T^*\mathbf{d}_i$ are eigenvectors of TT^* and T^*T corresponding to the eigenvalue $\lambda_i(TT^*) = \lambda_i(T^*T)$ respectively. Furthermore if $\mathbf{c}_1, \dots, \mathbf{c}_r$ satisfy (2.1) then $\tilde{\mathbf{d}}_i := \frac{T\mathbf{c}_i}{\|T\mathbf{c}_i\|}, i = 1, \dots, r$ satisfy (2.2) for $i = 1, \dots, r$. Similar result holds for $\mathbf{d}_1, \dots, \mathbf{d}_r$.

Proof. Clearly $T\mathbf{x} = 0 \iff \langle T\mathbf{x}, T\mathbf{x} \rangle = 0 \iff T^*T\mathbf{x} = 0$. Hence

$$\text{rank } T^*T = \text{rank } T = \text{rank } T^* = \text{rank } TT^* = r.$$

Thus T^*T and TT^* have exactly r positive eigenvalues. Let $i \in [1, r]$. Then $T^*T\mathbf{c}_i \neq 0$. Hence $T\mathbf{c}_i \neq 0$. (2.1) yields that $TT^*(T\mathbf{c}_i) = \lambda_i(T^*T)(T\mathbf{c}_i)$. Similarly $T^*T(T^*\mathbf{d}_i) = \lambda_i(TT^*)(T^*\mathbf{d}_i) \neq 0$. Hence (2.3) holds. Assume that $\mathbf{c}_1, \dots, \mathbf{c}_r$ satisfy (2.1). Let $\tilde{\mathbf{d}}_1, \dots, \tilde{\mathbf{d}}_r$ be defined as above. By the definition $\|\tilde{\mathbf{d}}_i\| = 1, i = 1, \dots, r$. Let $1 \leq i < j \leq r$. Then

$$0 = \langle \mathbf{c}_i, \mathbf{c}_j \rangle = \lambda_i(T^*T)\langle \mathbf{c}_i, \mathbf{c}_j \rangle = \langle T^*T\mathbf{c}_i, \mathbf{c}_j \rangle = \langle T\mathbf{c}_i, T\mathbf{c}_j \rangle \Rightarrow \langle \tilde{\mathbf{d}}_i, \tilde{\mathbf{d}}_j \rangle = 0.$$

Hence $\tilde{\mathbf{d}}_1, \dots, \tilde{\mathbf{d}}_r$ is an orthonormal system. \square

Let

$$\sigma_i(T) = \sqrt{\lambda_i(T^*T)} \text{ for } i = 1, \dots, r, \quad \sigma_i(T) = 0 \text{ for } i > r, \quad (2.4)$$

$$\sigma_{(p)}(T) := (\sigma_1(T), \dots, \sigma_p(T))^{\top} \in \mathbb{R}_{\geq}^p, \quad p \in \mathbb{N}.$$

Then $\sigma_i(T) = \sigma_i(T^*)$, $i = 1, \dots, \min(m, n)$ are called the singular values of T and T^* respectively. Note that the singular values are arranged in a decreasing order. The positive singular values are called principal singular values of T and T^* respectively. Note that

$$\begin{aligned} \|T\mathbf{c}_i\|^2 &= \langle T\mathbf{c}_i, T\mathbf{c}_i \rangle = \langle T^*T\mathbf{c}_i, \mathbf{c}_i \rangle = \lambda_i(T^*T) = \sigma_i^2 \Rightarrow \\ \|T\mathbf{c}_i\| &= \sigma_i, \quad i = 1, \dots, n, \\ \|T^*\mathbf{d}_j\|^2 &= \langle T^*\mathbf{d}_j, T^*\mathbf{d}_j \rangle = \langle TT^*\mathbf{d}_j, \mathbf{d}_j \rangle = \lambda_j(TT^*) = \sigma_j^2 \Rightarrow \\ \|T\mathbf{d}_j\| &= \sigma_j, \quad j = 1, \dots, m. \end{aligned}$$

Let $\mathbf{c}_1, \dots, \mathbf{c}_n$ be an orthonormal basis of \mathbf{V} satisfying (2.1). Choose an orthonormal basis $\mathbf{d}_1, \dots, \mathbf{d}_m$ as follows. Set $\mathbf{d}_i := \frac{T\mathbf{c}_i}{\sigma_i}$, $i = 1, \dots, r$. Then complete the orthonormal set $\{\mathbf{d}_1, \dots, \mathbf{d}_r\}$ to an orthonormal basis of \mathbf{U} . Since $\text{span}(\mathbf{d}_1, \dots, \mathbf{d}_r)$ is spanned by all eigenvectors of TT^* corresponding to nonzero eigenvalues of TT^* it follows that $\ker T^* = \text{span}(\mathbf{d}_{r+1}, \dots, \mathbf{d}_m)$. Hence (2.2) holds. In these orthonormal bases of \mathbf{U} and \mathbf{V} the operators T and T^* represented quite simply:

$$T\mathbf{c}_i = \sigma_i(T)\mathbf{d}_i, \quad i = 1, \dots, n, \quad \text{where } \mathbf{d}_i = 0 \text{ for } i > m, \quad (2.5)$$

$$T^*\mathbf{d}_j = \sigma_j(T)\mathbf{c}_j, \quad j = 1, \dots, m, \quad \text{where } \mathbf{c}_j = 0 \text{ for } j > n..$$

Let

$$\Sigma = (s_{ij})_{i,j=1}^{m,n}, \quad s_{ij} = 0 \text{ for } i \neq j, \quad s_{ii} = \sigma_i \text{ for } i = 1, \dots, \min(m, n). \quad (2.6)$$

In the case $m \neq n$ we call Σ a diagonal matrix with the diagonal $\sigma_1, \dots, \sigma_{\min(m, n)}$. Then in the bases $[\mathbf{d}_1, \dots, \mathbf{d}_m]$ and $[\mathbf{c}_1, \dots, \mathbf{c}_n]$ T and T^* represented by the matrices Σ and Σ^\top respectively.

Lemma 2.59 *Let $[\mathbf{u}_1, \dots, \mathbf{u}_m], [\mathbf{v}_1, \dots, \mathbf{v}_n]$ be orthonormal bases in the vector spaces \mathbf{U}, \mathbf{V} over $\mathbb{F} = \mathbb{R}, \mathbb{C}$ respectively. Then T and T^* are presented by the matrices $A \in \mathbb{F}^{m \times n}$ and $A^* \in \mathbb{F}^{n \times m}$ respectively. Let $U \in \mathbf{U}(m)$ and $V \in \mathbf{U}(n)$ be the unitary matrices representing the change of base $[\mathbf{d}_1, \dots, \mathbf{d}_m]$ to $[\mathbf{u}_1, \dots, \mathbf{u}_m]$ and $[\mathbf{c}_1, \dots, \mathbf{c}_n]$ to $[\mathbf{v}_1, \dots, \mathbf{v}_n]$ respectively. (If $\mathbb{F} = \mathbb{R}$ then U and V are orthogonal matrices.) Then*

$$A = U\Sigma V^* \in \mathbb{F}^{m \times n}, \quad U \in \mathbf{U}(m), \quad V \in \mathbf{U}(n). \quad (2.7)$$

Proof. By the definition $T\mathbf{v}_j = \sum_{i=1}^m a_{ij}\mathbf{u}_i$. Let $U = (u_{ip})_{i,p=1}^m, V = (v_{jq})_{j,q=1}^n$. Then

$$T\mathbf{c}_q = \sum_{j=1}^n v_{jq}T\mathbf{v}_j = \sum_{j=1}^n v_{jq} \sum_{i=1}^m a_{ij}\mathbf{u}_i = \sum_{j=1}^n v_{jq} \sum_{i=1}^m a_{ij} \sum_{p=1}^m \bar{u}_{ip}\mathbf{d}_p.$$

Use the first equality of (2.5) to deduce that $U^*AV = \Sigma$. □

Definition 2.60 (2.7) is called the singular value decomposition (SVD) of A .

Proposition 2.61 *Let $\mathbb{F} = \mathbb{R}, \mathbb{C}$ and denote by $\mathcal{R}_{m,n,k}(\mathbb{F}) \subset \mathbb{F}^{m \times n}$ the set of all matrices of rank $k \in [1, \min(m, n)]$ at most. Then $A \in \mathcal{R}_{m,n,k}(\mathbb{F})$ if and only if A can be expressed as a sum of at most k matrices of rank 1. Furthermore $\mathcal{R}_{m,n,k}(\mathbb{F})$ is a variety in $\mathbb{F}^{m \times n}$ given by the polynomial conditions: Each $(k+1) \times (k+1)$ minor of A is equal to zero.*

For the proof see Problem 2

Definition 2.62 Let $A \in \mathbb{C}^{m \times n}$ and assume that A has the SVD given by (2.7), where $U = [\mathbf{u}_1, \dots, \mathbf{u}_m]$, $V = [\mathbf{v}_1, \dots, \mathbf{v}_n]$. Denote by $A_k := \sum_{i=1}^k \sigma_i \mathbf{u}_i \mathbf{v}_i^* \in \mathbb{C}^{m \times n}$ for $k = 1, \dots, \text{rank } A$. For $k > \text{rank } A$ we define $A_k := A (= A_{\text{rank } A})$.

Note that for $1 \leq k < \text{rank } A$, the matrix A_k is uniquely defined if and only if $\sigma_k > \sigma_{k+1}$. (See Problem 1.)

Theorem 2.63 For $\mathbb{F} = \mathbb{R}, \mathbb{C}$ and $A = (a_{ij}) \in \mathbb{F}^{m \times n}$ the following conditions hold:

$$\|A\|_F := \sqrt{\text{tr } A^* A} = \sqrt{\text{tr } A A^*} = \sqrt{\sum_{i=1}^{\text{rank } A} \sigma_i(A)^2}. \quad (2.8)$$

$$\|A\|_2 := \max_{\mathbf{x} \in \mathbb{F}^n, \|\mathbf{x}\|_2=1} \|A\mathbf{x}\|_2 = \sigma_1(A). \quad (2.9)$$

$$\min_{B \in \mathcal{R}_{m,n,k}(\mathbb{F})} \|A - B\|_2 = \|A - A_k\| = \sigma_{k+1}(A), k = 1, \dots, \text{rank } A - 1. \quad (2.10)$$

$$\sigma_i(A) \geq \sigma_i((a_{i_p j_q})_{p=1, q=1}^{m', n'}) \geq \sigma_{i+(m-m')+(n-n')}(A), \quad (2.11)$$

$$m' \in [1, m], n' \in [1, n], 1 \leq i_1 < \dots < i_{m'} \leq m, 1 \leq j_1 < \dots < j_{n'} \leq n.$$

Proof. The proof of (2.8) is left a Problem 7. We now show the equality in (2.9). View A as an operator $A : \mathbb{C}^n \rightarrow \mathbb{C}^m$. From the definition of $\|A\|_2$ it follows

$$\|A\|_2^2 = \max_{\mathbf{x} \in \mathbb{R}^n, \mathbf{x} \neq 0} \frac{\mathbf{x}^* A^* A \mathbf{x}}{\mathbf{x}^* \mathbf{x}} = \lambda_1(A^* A) = \sigma_1(A)^2,$$

which proves (2.9).

We now prove (2.10). In the SVD decomposition of A (2.7) assume that $U = (\mathbf{u}_1, \dots, \mathbf{u}_m)$ and $V = (\mathbf{v}_1, \dots, \mathbf{v}_n)$. Then (2.7) is equivalent to the following representation of A :

$$A = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^*, \mathbf{u}_1, \dots, \mathbf{u}_r \in \mathbb{R}^m, \mathbf{v}_1, \dots, \mathbf{v}_r \in \mathbb{R}^n, \mathbf{u}_i^* \mathbf{u}_j = \mathbf{v}_i^* \mathbf{v}_j = \delta_{ij}, i, j = 1, \dots, r, \quad (2.12)$$

where $r = \text{rank } A$. Let $B = \sum_{i=1}^k \sigma_i \mathbf{u}_i \mathbf{v}_i^* \in \mathcal{R}_{m,n,k}$. Then in view of (2.9)

$$\|A - B\|_2 = \left\| \sum_{k+1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^* \right\|_2 = \sigma_{k+1}.$$

Let $B \in \mathcal{R}_{m,n,k}$. To show (2.10) it is enough to show that $\|A - B\|_2 \geq \sigma_{k+1}$. Let

$$\mathbf{W} := \{\mathbf{x} \in \mathbb{R}^n : B\mathbf{x} = 0\}.$$

Then $\text{codim } \mathbf{W} \geq k$. Furthermore

$$\|A - B\|_2^2 \geq \max_{\|\mathbf{x}\|_2=1, \mathbf{x} \in \mathbf{W}} \|(A - B)\mathbf{x}\|^2 = \max_{\|\mathbf{x}\|_2=1, \mathbf{x} \in \mathbf{W}} \mathbf{x}^* A^* A \mathbf{x} \geq \lambda_{k+1}(A^* A) = \sigma_{k+1}^2,$$

where the last inequality follows from the min-max characterization of $\lambda_{k+1}(A^*A)$.

Let $C = (a_{ijq})_{i,q=1}^{m,n'}$. Then C^*C is a principal submatrix of A^*A of dimension n' . The interlacing inequalities between the eigenvalues of A^*A and C^*C yields (2.11) for $m' = m$. Let $D = (a_{ipjq})_{p,q=1}^{m',n'}$. Then DD^* is a principle submatrix of CC^* . Use the interlacing properties of the eigenvalues of CC^* and DD^* to deduce (2.11). \square

We now restate the above results for linear operators.

Definition 2.64 Let \mathbf{U}, \mathbf{V} be finite dimensional vector spaces over $\mathbb{F} = \mathbb{R}, \mathbb{C}$. For $k \in \mathbb{Z}_+$ denote $L_k(\mathbf{V}, \mathbf{U}) := \{T \in L(\mathbf{V}, \mathbf{U}) : \text{rank } T \leq k\}$. Assume furthermore that \mathbf{U}, \mathbf{V} are IPS. Let $T \in L(\mathbf{V}, \mathbf{U})$ and assume that the orthonormal bases of $[\mathbf{d}_1, \dots, \mathbf{d}_m], [\mathbf{c}_1, \dots, \mathbf{c}_n]$ of \mathbf{U}, \mathbf{V} respectively satisfy (2.5). Define $T_0 := \mathbf{0}$ and $T_k := T$ for an integer $k \geq \text{rank } T$. Let $k \in [1, \text{rank } T - 1] \cap \mathbb{N}$. Define $T_k \in L(\mathbf{V}, \mathbf{U})$ by the equality $T_k(\mathbf{v}) = \sum_{i=1}^k \sigma_i(T) \langle \mathbf{v}, \mathbf{c}_i \rangle \mathbf{d}_i$ for any $\mathbf{v} \in \mathbf{V}$.

It is straightforward to show that $T_k \in L_k(\mathbf{V}, \mathbf{U})$ and T_k is unique if and only if $\sigma_k(T) > \sigma_{k+1}(T)$. See Problem 8. Theorem 2.63 yields:

Corollary 2.65 Let \mathbf{U} and \mathbf{V} be finite dimensional IPS over $\mathbb{F} = \mathbb{R}, \mathbb{C}$. Let $T : \mathbf{V} \rightarrow \mathbf{U}$ be a linear operator. Then

$$\|T\|_F := \sqrt{\text{tr } T^*T} = \sqrt{\text{tr } TT^*} = \sqrt{\sum_{i=1}^{\text{rank } T} \sigma_i(T)^2}. \quad (2.13)$$

$$\|T\|_2 := \max_{\mathbf{x} \in \mathbf{V}, \|\mathbf{x}\|_2=1} \|T\mathbf{x}\|_2 = \sigma_1(T). \quad (2.14)$$

$$\min_{Q \in L_k(\mathbf{V}, \mathbf{U})} \|T - Q\|_2 = \sigma_{k+1}(T), \quad k = 1, \dots, \text{rank } T - 1. \quad (2.15)$$

Problems

1. Let \mathbf{U}, \mathbf{V} be finite dimensional inner product spaces. Assume that $T \in L(\mathbf{U}, \mathbf{V})$. Show that for any complex number $t \in \mathbb{C}$ $\sigma_i(tT) = |t|\sigma_i(T)$ for all i .
2. Prove Proposition 2.61. (Use SVD to prove the nontrivial part of the Proposition.)
3. Let $A \in \mathbb{C}^{m \times n}$ and assume that $U \in \mathbf{U}_m, V \in \mathbf{V}_n$. Show that $\sigma_i(UAV) = \sigma_i(A)$ for all i .
4. Let $A \in \text{GL}(n, \mathbb{C})$. Show that $\sigma_1(A^{-1}) = \sigma_n(A)^{-1}$.
5. Let \mathbf{U}, \mathbf{V} be IPS inner product space of dimensions m and n respectively. Assume that

$$\begin{aligned} \mathbf{U} &= \mathbf{U}_1 \oplus \mathbf{U}_2, \dim \mathbf{U}_1 = m_1, \dim \mathbf{U}_2 = m_2, \mathbf{U}_1 \perp \mathbf{U}_2, \\ \mathbf{V} &= \mathbf{V}_1 \oplus \mathbf{V}_2, \dim \mathbf{V}_1 = n_1, \dim \mathbf{V}_2 = n_2, \mathbf{V}_1 \perp \mathbf{V}_2. \end{aligned}$$

Assume that $T \in L(\mathbf{V}, \mathbf{U})$. Suppose furthermore that $T\mathbf{V}_1 \subseteq \mathbf{U}_1, T\mathbf{V}_2 \subseteq \mathbf{U}_2$. Let $T_i \in L(\mathbf{V}_i, \mathbf{U}_i)$ be the restriction of T to \mathbf{V}_i for $i = 1, 2$. Then $\text{rank } T = \text{rank } T_1 + \text{rank } T_2$ and $\{\sigma_1(T), \dots, \sigma_{\text{rank } T}(T)\} = \{\sigma_1(T_1), \dots, \sigma_{\text{rank } T_1}(T_1)\} \cup \{\sigma_1(T_2), \dots, \sigma_{\text{rank } T_2}(T_2)\}$.

6. Let the assumptions of the Definition 2.62 hold. Show that for $1 \leq k < \text{rank } A$ A_k is uniquely defined if and only if $\sigma_k > \sigma_{k+1}$.
7. Prove the equalities in (2.8).
8. Let the assumptions of Definition 2.64 hold. Show that for $k \in [1, \text{rank } T - 1] \cap \mathbb{N}$ $\text{rank } T_k = k$ and T_k is unique if and only if $\sigma_k(T) > \sigma_{k+1}(T)$.
9. Let \mathbf{V} be an n -dimensional IPS. Assume that $T \in L(\mathbf{V})$ is a normal operator. Let $\lambda_1(T), \dots, \lambda_n(T)$ be the eigenvalues of T arranged in the order $|\lambda_1(T)| \geq \dots \geq |\lambda_n(T)|$. Show that $\sigma_i(T) = |\lambda_i(T)|$ for $i = 1, \dots, n$.

2.9 Characterizations of singular values

Theorem 2.66 *Let $\mathbb{F} = \mathbb{R}, \mathbb{C}$ and assume that $A \in \mathbb{F}^{m \times n}$. Define*

$$H(A) = \begin{bmatrix} 0 & A \\ A^* & 0 \end{bmatrix} \in \mathbb{H}_{m+n}(\mathbb{F}). \quad (2.1)$$

Then

$$\begin{aligned} \lambda_i(H(A)) &= \sigma_i(A), \quad \lambda_{m+n+1-i}(H(A)) = -\sigma_i(A), \quad i = 1, \dots, \text{rank } A, \\ \lambda_j(H(A)) &= 0, \quad j = \text{rank } A + 1, \dots, n + m - \text{rank } A. \end{aligned} \quad (2.2)$$

View A as an operator $A : \mathbb{F}^n \rightarrow \mathbb{F}^m$. Choose orthonormal bases $[\mathbf{d}_1, \dots, \mathbf{d}_m], [\mathbf{c}_1, \dots, \mathbf{c}_n]$ in $\mathbb{F}^m, \mathbb{F}^n$ respectively satisfying (2.5). Then

$$\begin{aligned} \begin{bmatrix} 0 & A \\ A^* & 0 \end{bmatrix} \begin{bmatrix} \mathbf{d}_i \\ \mathbf{c}_i \end{bmatrix} &= \sigma_i(A) \begin{bmatrix} \mathbf{d}_i \\ \mathbf{c}_i \end{bmatrix}, \quad \begin{bmatrix} 0 & A \\ A^* & 0 \end{bmatrix} \begin{bmatrix} \mathbf{d}_i \\ -\mathbf{c}_i \end{bmatrix} = -\sigma_i(A) \begin{bmatrix} \mathbf{d}_i \\ -\mathbf{c}_i \end{bmatrix}, \\ & i = 1, \dots, \text{rank } A, \end{aligned} \quad (2.3)$$

$$\ker H(A) = \text{span}((\mathbf{d}_{r+1}^*, 0)^*, \dots, (\mathbf{d}_m^*, 0)^*, (0, \mathbf{c}_{r+1}^*)^*, \dots, (0, \mathbf{c}_n^*)^*), \quad r = \text{rank } A.$$

Proof. It is straightforward to show the equalities (2.3). Since all the eigenvectors appearing in (2.3) are linearly independent we deduce (2.2). \square

Corollary 2.67 *Let $\mathbb{F} = \mathbb{R}, \mathbb{C}$ and assume that $A \in \mathbb{F}^{m \times n}$. Let $\hat{A} := A[\alpha, \beta] \in \mathbb{F}^{p \times q}$ be a submatrix of A , formed by the set of rows and columns $\alpha \in Q_{p,m}, \beta \in Q_{q,n}$ respectively. Then*

$$\sigma_i(\hat{A}) \leq \sigma_i(A) \text{ for } i = 1, \dots \quad (2.4)$$

For $l \in [1, \text{rank } A] \cap \mathbb{N}$ the equalities $\sigma_i(\hat{A}) = \sigma_i(A), i = 1, \dots, l$ hold if and only if there exists two orthonormal systems of l right and left singular vectors $\mathbf{c}_1, \dots, \mathbf{c}_l \in \mathbb{F}^n, \mathbf{d}_1, \dots, \mathbf{d}_l \in \mathbb{F}^m$ satisfying (2.3) for $i = 1, \dots, l$ such that the nonzero coordinates vectors $\mathbf{c}_1, \dots, \mathbf{c}_l$ and $\mathbf{d}_1, \dots, \mathbf{d}_l$ are located at the indices β, α respectively.

See Problem 1.

Corollary 2.68 *Let \mathbf{V}, \mathbf{U} be IPS over $\mathbb{F} = \mathbb{R}, \mathbb{C}$. Assume that \mathbf{W} is a subspace of \mathbf{V} . Let $T \in L(\mathbf{V}, \mathbf{U})$ and denote by $\hat{T} \in L(\mathbf{W}, \mathbf{U})$ the restriction of T to \mathbf{W} . Then $\sigma_i(\hat{T}) \leq \sigma_i(T)$ for any $i \in \mathbb{N}$. Furthermore $\sigma_i(\hat{T}) = \sigma_i(T)$ for $i = 1, \dots, l \leq \text{rank } T$ if and only if \mathbf{U} contains a subspace spanned by the first l right singular vectors of T .*

See Problem 2.

Define by $\mathbb{R}_{+, \searrow}^n := \mathbb{R}_{\searrow}^n \cap \mathbb{R}_+^n$. Then $D \subset \mathbb{R}_{+, \searrow}^n$ is called a strong Schur set if for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}_{+, \searrow}^n$, $\mathbf{x} \preceq \mathbf{y}$ we have the implication $\mathbf{y} \in D \Rightarrow \mathbf{x} \in D$.

Theorem 2.69 *Let $p \in \mathbb{N}$ and $D \subset \mathbb{R}_{\searrow}^p \cap \mathbb{R}_+^p$ be a regular convex strong Schur domain. Fix $m, n \in \mathbb{N}$ and let $\sigma_{(p)}(D) := \{A \in \mathbb{F}^{m \times n} : \sigma_{(p)}(A) \in D\}$. Let $h : D \rightarrow \mathbb{R}$ be a convex and strongly Schur's order preserving on D . Let $f : \sigma_{(p)} \rightarrow \mathbb{R}$ be given as $h \circ \sigma_{(p)}$. Then f is a convex function.*

See Problem 3.

Corollary 2.70 *Let $\mathbb{F} = \mathbb{R}, \mathbb{C}$, $m, n, p \in \mathbb{N}$, $q \in [1, \infty)$ and $w_1 \geq w_2 \geq \dots \geq w_p > 0$. Then the following function*

$$f : \mathbb{F}^{m \times n} \rightarrow \mathbb{R}, \quad f(A) := \left(\sum_{i=1}^p w_i \sigma_i(A)^q \right)^{\frac{1}{q}}, \quad A \in \mathbb{F}^{m \times n}$$

is a convex function.

See Problem 4

We now translate Theorem 2.66 to the operator setting.

Lemma 2.71 *Let \mathbf{U}, \mathbf{V} be finite dimensional IPS spaces with the inner products $\langle \cdot, \cdot \rangle_{\mathbf{U}}, \langle \cdot, \cdot \rangle_{\mathbf{V}}$ respectively. Define $\mathbf{W} := \mathbf{V} \oplus \mathbf{U}$ be the induced IPS with*

$$\langle (\mathbf{y}, \mathbf{x}), (\mathbf{v}, \mathbf{u}) \rangle_{\mathbf{W}} := \langle \mathbf{y}, \mathbf{v} \rangle_{\mathbf{V}} + \langle \mathbf{x}, \mathbf{u} \rangle_{\mathbf{U}}.$$

Let $T : \mathbf{V} \rightarrow \mathbf{U}$ be a linear operator, and $T^* : \mathbf{U} \rightarrow \mathbf{V}$ be the adjoint of T . Define the operator

$$\hat{T} : \mathbf{W} \rightarrow \mathbf{W}, \quad \hat{T}(\mathbf{y}, \mathbf{x}) := (T^* \mathbf{x}, T \mathbf{y}). \quad (2.5)$$

Then \hat{T} is self-adjoint operator and $\hat{T}^2 = T^* T \oplus T T^*$. Hence the spectrum of \hat{T} is symmetric with respect to the origin and \hat{T} has exactly $2 \text{rank } T$ nonzero eigenvalues. More precisely, if $\dim \mathbf{U} = m, \dim \mathbf{V} = n$ then:

$$\begin{aligned} \lambda_i(\hat{T}) &= -\lambda_{m+n-i+1}(\hat{T}) = \sigma_i(T), \quad \text{for } i = 1, \dots, \text{rank } T, \\ \lambda_j(\hat{T}) &= 0, \quad \text{for } j = \text{rank } T + 1, \dots, m + n - \text{rank } T. \end{aligned} \quad (2.6)$$

Let $\{\mathbf{d}_1, \dots, \mathbf{d}_{\min(m,n)}\} \in \text{Fr}(\min(m,n), \mathbf{U}), \{\mathbf{c}_1, \dots, \mathbf{c}_{\min(m,n)}\} \in \text{Fr}(\min(m,n), \mathbf{V})$ be the set of vectors satisfying (2.5). Define

$$\mathbf{z}_i := \frac{1}{\sqrt{2}}(\mathbf{c}_i, \mathbf{d}_i), \mathbf{z}_{m+n-i+1} := \frac{1}{\sqrt{2}}(\mathbf{c}_i, -\mathbf{d}_i), i = 1, \dots, \min(m,n). \quad (2.7)$$

Then $\{\mathbf{z}_1, \mathbf{z}_{m+n}, \dots, \mathbf{z}_{\min(m,n)}, \mathbf{z}_{m+n-\min(m,n)+1}\} \in \text{Fr}(2 \min(m,n), \mathbf{W})$. Furthermore $\hat{T} \mathbf{z}_i = \sigma_i(T) \mathbf{z}_i, \hat{T} \mathbf{z}_{m+n-i+1} = -\sigma_i(T) \mathbf{z}_{m+n-i+1}$ for $i = 1, \dots, \min(m,n)$.

See Problem 5.

Theorem 2.72 *Let \mathbf{U}, \mathbf{V} be m and n -dimensional IPS over \mathbb{C} respectively. Let $T : \mathbf{V} \rightarrow \mathbf{U}$ be a linear operator. Then for each $k \in [1, \min(m, n)] \cap \mathbb{Z}$*

$$\begin{aligned} \sum_{i=1}^k \sigma_i(T) &= \max_{\{\mathbf{f}_1, \dots, \mathbf{f}_k\} \in \text{Fr}(k, \mathbf{U}), \{\mathbf{g}_1, \dots, \mathbf{g}_k\} \in \text{Fr}(k, \mathbf{V})} \sum_{i=1}^k \Re \langle T \mathbf{g}_i, \mathbf{f}_i \rangle_{\mathbf{U}} = \\ & \max_{\{\mathbf{f}_1, \dots, \mathbf{f}_k\} \in \text{Fr}(k, \mathbf{U}), \{\mathbf{g}_1, \dots, \mathbf{g}_k\} \in \text{Fr}(k, \mathbf{V})} \sum_{i=1}^k |\langle T \mathbf{g}_i, \mathbf{f}_i \rangle_{\mathbf{U}}|. \end{aligned} \quad (2.8)$$

Furthermore $\sum_{i=1}^k \sigma_i(T) = \sum_{i=1}^k \Re \langle T \mathbf{g}_i, \mathbf{f}_i \rangle_{\mathbf{U}}$ for some two k -orthonormal frames $F_k = \{\mathbf{f}_1, \dots, \mathbf{f}_k\}, G_k = \{\mathbf{g}_1, \dots, \mathbf{g}_k\}$ if and only if $\text{span}((\mathbf{g}_1, \mathbf{f}_1), \dots, (\mathbf{g}_k, \mathbf{f}_k))$ is spanned by k eigenvectors of \hat{T} corresponding to the first k eigenvalues of \hat{T} .

Proof. Assume that $\{\mathbf{f}_1, \dots, \mathbf{f}_k\} \in \text{Fr}(k, \mathbf{U}), \{\mathbf{g}_1, \dots, \mathbf{g}_k\} \in \text{Fr}(k, \mathbf{V})$. Let $\mathbf{w}_i := \frac{1}{\sqrt{2}}(\mathbf{g}_i, \mathbf{f}_i), i = 1, \dots, k$. Then $\{\mathbf{w}_1, \dots, \mathbf{w}_k\} \in \text{Fr}(k, \mathbf{W})$. A straightforward calculation shows $\sum_{i=1}^k \langle \hat{T} \mathbf{w}_i, \mathbf{w}_i \rangle_{\mathbf{W}} = \sum_{i=1}^k \Re \langle T \mathbf{g}_i, \mathbf{f}_i \rangle_{\mathbf{U}}$. The maximal characterization of $\sum_{i=1}^k \lambda_i(\hat{T})$, (Theorem 2.42), and (2.6) yield the inequality $\sum_{i=1}^k \sigma_i(\hat{T}) \geq \sum_{i=1}^k \Re \langle T \mathbf{g}_i, \mathbf{f}_i \rangle_{\mathbf{U}}$ for $k \in [\min(m, n) \cap \mathbb{Z}]$. Let $\mathbf{c}_1, \dots, \mathbf{c}_{\min(m, n)}, \mathbf{d}_1, \dots, \mathbf{d}_{\min(m, n)}$ satisfy (2.5). Then Lemma 2.71 yields that $\sum_{i=1}^k \sigma_i(\hat{T}) = \sum_{i=1}^k \Re \langle T \mathbf{c}_i, \mathbf{d}_i \rangle_{\mathbf{U}}$ for $k \in [\min(m, n) \cap \mathbb{Z}]$. This proves the first equality of (2.8). The second equality of (2.8) is straightforward. (See Problem 6.)

Assume now that $\sum_{i=1}^k \sigma_i(T) = \sum_{i=1}^k \Re \langle T \mathbf{g}_i, \mathbf{f}_i \rangle_{\mathbf{U}}$ for some two k -orthonormal frames $F_k = \{\mathbf{f}_1, \dots, \mathbf{f}_k\}, G_k = \{\mathbf{g}_1, \dots, \mathbf{g}_k\}$. Define $\mathbf{w}_1, \dots, \mathbf{w}_k$ as above. The above arguments yield that $\sum_{i=1}^k \langle \hat{T} \mathbf{w}_i, \mathbf{w}_i \rangle_{\mathbf{W}} = \sum_{i=1}^k \lambda_i(\hat{T})$. Theorem 2.42 yields that $\text{span}((\mathbf{g}_1, \mathbf{f}_1), \dots, (\mathbf{g}_k, \mathbf{f}_k))$ is spanned by k eigenvectors of \hat{T} corresponding to the first k eigenvalues of \hat{T} . Vice versa, assume that $\{\mathbf{f}_1, \dots, \mathbf{f}_k\} \in \text{Fr}(k, \mathbf{U}), \{\mathbf{g}_1, \dots, \mathbf{g}_k\} \in \text{Fr}(k, \mathbf{V})$ and $\text{span}((\mathbf{g}_1, \mathbf{f}_1), \dots, (\mathbf{g}_k, \mathbf{f}_k))$ is spanned by k eigenvectors of \hat{T} corresponding to the first k eigenvalues of \hat{T} . Define $\{\mathbf{w}_1, \dots, \mathbf{w}_k\} \in \text{Fr}(\mathbf{W})$ as above. Then $\text{span}(\mathbf{w}_1, \dots, \mathbf{w}_k)$ contains k linearly independent eigenvectors corresponding to the the first k eigenvalues of \hat{T} . Theorem 2.42 and Lemma 2.71 yield that $\sigma_i(T) = \sum_{i=1}^k \langle \hat{T} \mathbf{w}_i, \mathbf{w}_i \rangle_{\mathbf{W}} = \sum_{i=1}^k \Re \langle T \mathbf{g}_i, \mathbf{f}_i \rangle_{\mathbf{U}}$. \square

Theorem 2.73 *\mathbf{U}, \mathbf{V} be m and n dimensional IPS spaces. Assume that Let $S, T : \mathbf{V} \rightarrow \mathbf{U}$ be linear operators. Then*

$$\Re \text{tr}(S^* T) \leq \sum_{i=1}^{\min(m, n)} \sigma_i(S) \sigma_i(T). \quad (2.9)$$

Equality holds if and only if there exists two orthonormal set $\{\mathbf{d}_1, \dots, \mathbf{d}_{\min(m, n)}\} \in \text{Fr}(\min(m, n), \mathbf{U}), \{\mathbf{c}_1, \dots, \mathbf{c}_{\min(m, n)}\} \in \text{Fr}(\min(m, n), \mathbf{V})$, such that

$$S \mathbf{c}_i = \sigma_i(S) \mathbf{d}_i, T \mathbf{c}_i = \sigma_i(T) \mathbf{d}_i, S^* \mathbf{d}_i = \sigma_i(S) \mathbf{c}_i, T^* \mathbf{d}_i = \sigma_i(T) \mathbf{c}_i, i = 1, \dots, \min(m, n). \quad (2.10)$$

Proof. Let $A, B \in \mathbb{C}^{n \times m}$. Then

$\text{tr } B^*A = \overline{\text{tr } AB^*}$. Hence $2\Re \text{tr } AB^* = \text{tr } H(A)H(B)$. Therefore $2\Re \text{tr } S^*T = \text{tr } \hat{S}\hat{T}$. Use Theorem 2.55 for \hat{S}, \hat{T} and Lemma 2.71 to deduce (2.9). Equality in (2.9) if and only if $\text{tr } \hat{S}\hat{T} = \sum_{i=1}^{m+n} \lambda_i(\hat{S})\lambda_i(\hat{T})$.

Clearly, the assumptions that $\{\mathbf{d}_1, \dots, \mathbf{d}_{\min(m,n)}\} \in \text{Fr}(\min(m,n), \mathbf{U})$, $\{\mathbf{c}_1, \dots, \mathbf{c}_{\min(m,n)}\} \in \text{Fr}(\min(m,n), \mathbf{V})$, and the equalities (2.10) imply equality in (2.9).

Assume equality in (2.9). Theorem 2.55 and the definitions of \hat{S}, \hat{T} yields the existence $\{\mathbf{d}_1, \dots, \mathbf{d}_{\min(m,n)}\} \in \text{Fr}(\min(m,n), \mathbf{U})$, $\{\mathbf{c}_1, \dots, \mathbf{c}_{\min(m,n)}\} \in \text{Fr}(\min(m,n), \mathbf{V})$, such that (2.10) hold. \square

Theorem 2.74 *Let \mathbf{U} and \mathbf{V} be finite dimensional IPS over $\mathbb{F} = \mathbb{R}, \mathbb{C}$. Let $T : \mathbf{V} \rightarrow \mathbf{U}$ be a linear operator. Then*

$$\min_{Q \in L_k(\mathbf{V}, \mathbf{U})} \|T - Q\|_F = \sqrt{\sum_{i=k+1}^{\text{rank } T} \sigma_i^2(T)}, \quad k = 1, \dots, \text{rank } T - 1. \quad (2.11)$$

Furthermore $\|T - Q\|_F = \sqrt{\sum_{i=k+1}^{\text{rank } T} \sigma_i^2(T)}$ for some $Q \in L_k(\mathbf{V}, \mathbf{U})$, $k < \text{rank } T$, if and only there $Q = T_k$, where T_k is defined in Definition 2.64.

Proof. Use Theorem 2.73 to deduce that for any $Q \in L(\mathbf{V}, \mathbf{U})$ one has

$$\begin{aligned} \|T - Q\|_F^2 &= \text{tr } T^*T - 2\Re \text{tr } Q^*T + \text{tr } Q^*Q \geq \\ &= \sum_{i=1}^{\text{rank } T} \sigma_i^2(T) - 2 \sum_{i=1}^k \sigma_i(T)\sigma_i(Q) + \sum_{i=1}^k \sigma_i^2(Q) = \\ &= \sum_{i=1}^k (\sigma_i(T) - \sigma_i(Q))^2 + \sum_{i=k+1}^{\text{rank } T} \sigma_i^2(T) \geq \sum_{i=k+1}^{\text{rank } T} \sigma_i^2(T). \end{aligned}$$

Clearly $\|T - T_k\|_F^2 = \sum_{i=k+1}^{\text{rank } T} \sigma_i^2(T)$. Hence (2.11) holds. Vice versa if $Q \in L_k(\mathbf{V}, \mathbf{U})$ and $\|T - Q\|_F^2 = \sum_{i=k+1}^{\text{rank } T} \sigma_i^2(T)$ then the equality case in Theorem 2.73 yields that $Q = T_k$. \square

Corollary 2.75 *Let $F = \mathbb{R}, \mathbb{C}$ and $A \in \mathbb{F}^{m \times n}$. Then*

$$\min_{B \in \mathcal{R}_{m,n,k}(F)} \|A - B\|_F = \sqrt{\sum_{i=k+1}^{\text{rank } A} \sigma_i^2(A)}, \quad k = 1, \dots, \text{rank } A - 1. \quad (2.12)$$

Furthermore $\|A - B\|_F = \sqrt{\sum_{i=k+1}^{\text{rank } A} \sigma_i^2(A)}$ for some $B \in \mathcal{R}_{m,n,k}(F)$, $k < \text{rank } A$, if and only there $B = A_k$, where A_k is defined in Definition 2.62.

Theorem 2.76 *Let $F = \mathbb{R}, \mathbb{C}$ and $A \in \mathbb{F}^{m \times n}$. Then*

$$\min_{B \in \mathcal{R}_{m,n,k}(F)} \sum_{i=1}^j \sigma_i(A-B) = \sum_{i=k+1}^{k+j} \sigma_i(A), \quad j = 1, \dots, \min(m,n) - k, \quad k = 1, \dots, \min(m,n) - 1. \quad (2.13)$$

Proof. Clearly, for $B = A_k$ we have the equality $\sum_{i=1}^j \sigma_i(A-B) = \sum_{i=k+1}^{k+j} \sigma_i(A)$. Let $B \in \mathcal{R}_{m,n,k}(\mathbb{F})$. Let $\mathbf{X} \in \text{Gr}(k, \mathbb{C}^n)$ be an subspace which contains the columns of B . Let $\mathbf{W} = \{(\mathbf{0}^\top, \mathbf{x}^\top)^\top \in \mathbb{F}^{m+n}, \mathbf{x} \in \mathbf{X}\}$. Observe that for any $\mathbf{z} \in \mathbf{W}^\perp$ one has the equality $\mathbf{z}^* H((A-B)) \mathbf{z} = \mathbf{z}^* H(A) \mathbf{z}$. Combine Theorems 2.43 and 2.66 to deduce $\sum_{i=1}^j \sigma_i(B-A) \geq \sum_{i=k+1}^{k+j} \sigma_i(A)$. \square

Theorem 2.77 *Let \mathbf{V} be an n -dimensional IPS over \mathbb{C} . Let $T : \mathbf{V} \rightarrow \mathbf{V}$ be a linear operator. Assume the n eigenvalues of T $\lambda_1(T), \dots, \lambda_n(T)$ are arranged the order $|\lambda_1(T)| \geq \dots \geq |\lambda_n(T)|$. Let $\boldsymbol{\lambda}_a(T) := (|\lambda_1(T)|, \dots, |\lambda_n(T)|)$, $\boldsymbol{\sigma}(T) := (\sigma_1(T), \dots, \sigma_n(T))$. Then $\boldsymbol{\lambda}_a(T) \preceq \boldsymbol{\sigma}(T)$. That is*

$$\sum_{i=1}^k |\lambda_i(T)| \leq \sum_{i=1}^k \sigma_i(T), \quad i = 1, \dots, n. \quad (2.14)$$

Furthermore, $\sum_{i=1}^k |\lambda_i(T)| = \sum_{i=1}^k \sigma_i(T)$ for some $k \in [1, n] \cap \mathbb{Z}$ if and only if the following conditions are satisfied. There exists an orthonormal basis $\mathbf{x}_1, \dots, \mathbf{x}_n$ of \mathbf{V} such that:

1. $T\mathbf{x}_i = \lambda_i(T)\mathbf{x}_i, T^*\mathbf{x}_i = \overline{\lambda_i(T)}\mathbf{x}_i$ for $i = 1, \dots, k$.
2. Denote by $S : \mathbf{U} \rightarrow \mathbf{U}$ the restriction of T to the invariant subspace $\mathbf{U} = \text{span}(\mathbf{x}_{k+1}, \dots, \mathbf{x}_n)$. Then $\|S\|_2 \leq |\lambda_k(T)|$.

Proof. Use Theorem 2.23 to choose an orthonormal basis $\mathbf{g}_1, \dots, \mathbf{g}_n$ of \mathbf{V} , such that T is represented by an upper diagonal matrix $A = [a_{ij}] \in \mathbb{C}^{n \times n}$ such that $a_{ii} = \lambda_i(T), i = 1, \dots, n$. Let $\epsilon_i \in \mathbb{C}, |\epsilon_i| = 1$ such that $\bar{\epsilon}_i \lambda_i(T) = |\lambda_i(T)|$ for $i = 1, \dots, n$. Let $S \in L(\mathbf{V})$ be presented in the basis $\mathbf{g}_1, \dots, \mathbf{g}_n$ by a diagonal matrix $\text{diag}(\epsilon_1, \dots, \epsilon_k, 0, \dots, 0)$. Clearly, $\sigma_i(S) = 1$ for $i = 1, \dots, k$ and $\sigma_i(S) = 0$ for $i = k+1, \dots, n$. Furthermore, $\Re \text{tr } S^* C = \sum_{i=1}^k |\lambda_i(T)|$. Hence Theorem 2.73 yields (2.14).

Assume now that $\sum_{i=1}^k |\lambda_i(T)| = \sum_{i=1}^k \sigma_i(T)$. Hence equality sign holds in (2.9). Hence there exists two orthonormal bases $\{\mathbf{c}_1, \dots, \mathbf{c}_n\}, \{\mathbf{d}_1, \dots, \mathbf{d}_n\}$ in \mathbf{V} such that (2.10) holds. It easily follows that $\{\mathbf{c}_1, \dots, \mathbf{c}_k\}, \{\mathbf{d}_1, \dots, \mathbf{d}_k\}$ are orthonormal bases of $\mathbf{W} := \text{span}(\mathbf{g}_1, \dots, \mathbf{g}_k)$. Hence \mathbf{W} is an invariant subspace of T and T^* . Hence $A = A_1 \oplus A_2$, i.e. A is a block diagonal matrix. Thus $A_1 = (a_{ij})_{i,j=1}^k \in \mathbb{C}^{k \times k}, A_2 = (a_{ij})_{i,j=k+1}^n \in \mathbb{C}^{(n-k) \times (n-k)}$ represent the restriction of T to $\mathbf{W}, \mathbf{U} := \mathbf{W}^\perp$, denoted by T_1 and T_2 respectively. Hence $\sigma_i(T_1) = \sigma_i(T)$ for $i = 1, \dots, k$. Note that the restriction of S to \mathbf{W} , denoted by S_1 is given by the diagonal matrix $D_1 := \text{diag}(\epsilon_1, \dots, \epsilon_k) \in \mathbf{U}(k)$. (2.10) yield that $S_1^{-1} T_1 \mathbf{c}_i = \sigma_i(T) \mathbf{c}_i$ for $i = 1, \dots, k$, i.e. $\sigma_1(T), \dots, \sigma_k(T)$ are the eigenvalues of $S_1^{-1} T_1$. Clearly $S_1^{-1} T_1$ is presented in the basis $[\mathbf{g}_1, \dots, \mathbf{g}_k]$ by the matrix $D_1^{-1} A_1$, which is a diagonal matrix with $|\lambda_1(T)|, \dots, |\lambda_k(T)|$ on the main diagonal. That is $S_1^{-1} T_1$ has eigenvalues $|\lambda_1(T)|, \dots, |\lambda_k(T)|$. Therefore $\sigma_i(T) = |\lambda_i(T)|$ for $i = 1, \dots, k$. Theorem 2.63 yields that

$$\text{tr } A_1^* A_1 = \sum_{i,j=1}^k |a_{ij}|^2 = \sum_{i=1}^k \sigma_i^2(A_1) = \sum_{i=1}^k \sigma_i^2(T_1) = \sum_{i=1}^k |\lambda_i(T)|^2.$$

As $\lambda_1(T), \dots, \lambda_k(T)$ are the diagonal elements of A_1 it follows from the above equality that A_1 is a diagonal matrix. Hence we can choose $\mathbf{x}_i = \mathbf{g}_i$ for $i = 1, \dots, n$ to obtain the part 1 of the equality case.

Let $T\mathbf{x} = \lambda\mathbf{x}$ where $\|\mathbf{x}\| = 1$ and $\rho(T) = |\lambda|$. Recall $\|T\|_2 = \sigma_1(T)$, where $\sigma_1(T)^2 = \lambda_1(T^*T)$ is the maximal eigenvalue of the self-adjoint operator T^*T . The maximum characterization of $\lambda_1(T^*T)$ yields that $|\lambda|^2 = \langle T\mathbf{x}, T\mathbf{x} \rangle = \langle T^*T\mathbf{x}, \mathbf{x} \rangle \leq \lambda_1(T^*T) = \|T\|_2^2$. Hence $\rho(T) \leq \|T\|_2$.

Assume now that $\rho(T) = \|T\|_2$. $\rho(T) = 0$ then $\|T\|_2 = 0 \Rightarrow T = 0$, and theorem holds trivially in this case. Assume that $\rho(T) > 0$. Hence the eigenvector $\mathbf{x}_1 := \mathbf{x}$ is also the eigenvector of T^*T corresponding to $\lambda_1(T^*T) = |\lambda|^2$. Hence $|\lambda|^2\mathbf{x} = T^*T\mathbf{x} = T^*(\lambda\mathbf{x})$, which implies that $T^*\mathbf{x} = \bar{\lambda}\mathbf{x}$. Let $\mathbf{U} = \text{span}(\mathbf{x})^\perp$ be the orthogonal complement of $\text{span}(\mathbf{x})$. Since $T\text{span}(\mathbf{x}) = \text{span}(\mathbf{x})$ it follows that $T^*\mathbf{U} \subseteq \mathbf{U}$. Similarly, since $T^*\text{span}(\mathbf{x}) = \text{span}(\mathbf{x})$ $T\mathbf{U} \subseteq \mathbf{U}$. Thus $\mathbf{V} = \text{span}(\mathbf{x}) \oplus \mathbf{U}$ and $\text{span}(\mathbf{x}), \mathbf{U}$ are invariant subspaces of T and T^* . Hence $\text{span}(\mathbf{x}), \mathbf{U}$ are invariant subspaces of T^*T and TT^* . Let T_1 be the restriction of T to \mathbf{U} . Then $T_1^*T_1$ is the restriction of T^*T . Therefore $\|T_1\|_2^2 = \lambda_1(T_1^*T_1) \geq \lambda_1(T^*T) = \|T\|_2^2$. This establishes the second part of theorem, labeled (a) and (b).

The above result imply that the conditions (a) and (b) of the theorem yield the equality $\rho(T) = \|T\|_2$. \square

Corollary 2.78 *Let \mathbf{U} be an n -dimensional IPS over \mathbb{C} . Let $T : \mathbf{U} \rightarrow \mathbf{U}$ be a linear operator. Denote by $|\lambda(T)| = (|\lambda_1(T)|, \dots, |\lambda_n(T)|)^T$ the absolute eigenvalues of T , (counting with their multiplicities), arranged in a decreasing order. Then $|\lambda(T)| = (\sigma_1(T), \dots, \sigma_n(T))^T$ if and only if T is a normal operator.*

Problems

1. Let the assumptions of Corollary 2.67 hold.

- (a) Since $\text{rank } \hat{A} \leq \text{rank } A$ show that the inequalities (2.4) reduce to $\sigma_i(\hat{A}) = \sigma_i(A) = 0$ for $i > \text{rank } A$.
- (b) Since $H(\hat{A})$ is a submatrix of $H(A)$ use the Cauchy interlacing principle to deduce the inequalities (2.4) for $i = 1, \dots, \text{rank } A$. Furthermore, if $p' := m - \#\alpha, q' = n - \#\beta$ then the Cauchy interlacing principle gives the complementary inequalities $\sigma_i(\hat{A}) \geq \sigma_{i+p'+q'}(A)$ for any $i \in \mathbb{N}$.
- (c) Assume that $\sigma_i(\hat{A}) = \sigma_i(A)$ for $i = 1, \dots, l \leq \text{rank } A$. Compare the maximal characterization of the sum of the first k eigenvalues of $H(\hat{A})$ and $H(A)$ given by Theorem 2.42 for $k = 1, \dots, l$ to deduce the last part of Corollary (2.67).

2. Prove Corollary 2.68 by choosing any orthonormal basis in \mathbf{U} , an orthonormal basis in \mathbf{V} whose first $\dim \mathbf{W}$ elements span \mathbf{W} , and using Problem 1.

3. Combine Theorems ?? and 2.66 to deduce Theorem 2.69.

4. (a) Prove Corollary 2.70

- (b) Recall the definition of a norm on a vector space over $\mathbb{F} = \mathbb{R}, \mathbb{C}$???. Show that the function f defined in Corollary 2.70 is a norm. For $p = \min(m, n)$ and $w_1 = \dots = w_p = 1$ this norm is called the q -Schatten norm.

5. Prove Lemma 2.71.
6. Under the assumptions of Theorem 2.72 show the equalities.

$$\begin{aligned} \max_{\{\mathbf{f}_1, \dots, \mathbf{f}_k\} \in \text{Fr}(k, \mathbf{U}), \{\mathbf{g}_1, \dots, \mathbf{g}_k\} \in \text{Fr}(k, \mathbf{V})} \sum_{i=1}^k \Re \langle T \mathbf{g}_i, \mathbf{f}_i \rangle_{\mathbf{U}} = \\ \max_{\{\mathbf{f}_1, \dots, \mathbf{f}_k\} \in \text{Fr}(k, \mathbf{U}), \{\mathbf{g}_1, \dots, \mathbf{g}_k\} \in \text{Fr}(k, \mathbf{V})} \sum_{i=1}^k |\langle T \mathbf{g}_i, \mathbf{f}_i \rangle_{\mathbf{U}}|. \end{aligned}$$

7. Under the assumptions of Theorem 2.72 is it *true* that that for $k > 1$

$$\sum_{i=1}^k \sigma_i(T) = \max_{\{\mathbf{f}_1, \dots, \mathbf{f}_k\} \in \text{Fr}(k, \mathbf{U})} \sum_{i=1}^k \|T \mathbf{f}_i\|_{\mathbf{V}}.$$

I doubt it.

8. Let \mathbf{U}, \mathbf{V} be finite dimensional IPS. Assume that $P, T \in L(\mathbf{U}, \mathbf{V})$. Show that $\Re \text{tr}(P^*T) \geq -\sum_{i=1}^{\min(m,n)} \sigma_i(S) \sigma_i(T)$. Equality holds if and only if $S = -P$ and T satisfy the conditions of Theorem 2.73.

2.10 Moore-Penrose generalized inverse

Let $A \in \mathbb{C}^{m \times n}$. Then (2.12) is called the *reduced SVD* of A . It can be written as

$$A = U_r \Sigma_r V_r^*, \quad r = \text{rank } A, \quad \Sigma_r := \text{diag}(\sigma_1(A), \dots, \sigma_r(A)) \in S_r(\mathbb{R}), \quad (2.1)$$

$$U_r = [\mathbf{u}_1, \dots, \mathbf{u}_r] \in \mathbb{C}^{m \times r}, V_r = [\mathbf{v}_1, \dots, \mathbf{v}_r] \in \mathbb{C}^{n \times r}, U_r^* U_r = V_r^* V_r = I_r, .$$

Recall that

$$\begin{aligned} AA^* \mathbf{u}_i &= \sigma_i(A)^2 \mathbf{u}_i, A^* A \mathbf{v}_i = \sigma_i(A)^2 \mathbf{v}_i, \\ \mathbf{v}_i &= \frac{1}{\sigma_i(A)} A^* \mathbf{u}_i, \mathbf{u}_i = \frac{1}{\sigma_i(A)} A \mathbf{v}_i, i = 1, \dots, r. \end{aligned}$$

Then

$$A^\dagger := V_r \Sigma_r^{-1} U_r^* \in \mathbb{C}^{n \times m} \quad (2.2)$$

is the *Moore-Penrose* generalized inverse of A . If $A \in \mathbb{R}^{m \times n}$ then we assume that $U \in \mathbb{R}^{m \times r}$ and $V \in \mathbb{R}^{n \times r}$, i.e. U, V are real values matrices over the real numbers \mathbb{R} .

Theorem 2.79 *Let $A \in \mathbb{C}^{m \times n}$ matrix. Then the Moore-Penrose generalized inverse $A^\dagger \in \mathbb{C}^{n \times m}$ satisfies the following properties.*

1. $\text{rank } A = \text{rank } A^\dagger$.
2. $A^\dagger A A^\dagger = A^\dagger$, $A A^\dagger A = A$, $A^* A A^\dagger = A^\dagger A A^* = A^*$.
3. $A^\dagger A$ and $A A^\dagger$ are Hermitian nonnegative definite idempotent matrices, i.e. $(A^\dagger A)^2 = A^\dagger A$ and $(A A^\dagger)^2 = A A^\dagger$, having the same rank as A .

4. The least square solution of $A\mathbf{x} = \mathbf{b}$, i.e. the solution of the system $A^*A\mathbf{x} = A^*\mathbf{b}$, has a solution $\mathbf{y} = A^\dagger\mathbf{b}$. This solution has the minimal norm $\|\mathbf{y}\|$, for all possible solutions of $A^*A\mathbf{x} = A^*\mathbf{b}$.
5. If $\text{rank } A = n$ then $A^\dagger = (A^*A)^{-1}A^*$. In particular, if $A \in \mathbb{C}^{n \times n}$ is invertible then $A^\dagger = A^{-1}$.

To prove the above theorem we need the following proposition.

Proposition 2.80 *Let $E \in \mathbb{C}^{l \times m}, G \in \mathbb{C}^{m \times n}$. Then $\text{rank } EG \leq \min(\text{rank } E, \text{rank } G)$. If $l = m$ and E is invertible then $\text{rank } EG = \text{rank } G$. If $m = n$ and G is invertible then $\text{rank } EG = \text{rank } E$.*

Proof. Let $\mathbf{e}_1, \dots, \mathbf{e}_m \in \mathbb{C}^l, \mathbf{g}_1, \dots, \mathbf{g}_n \in \mathbb{C}^m$ be the columns of E and G respectively. Then $\text{rank } E = \dim \text{span}(\mathbf{e}_1, \dots, \mathbf{e}_l)$. Observe that $EG = [E\mathbf{g}_1, \dots, E\mathbf{g}_n] \in \mathbb{C}^{l \times n}$. Clearly $E\mathbf{g}_i$ is a linear combination of the columns of E . Hence $E\mathbf{g}_i \in \text{span}(\mathbf{e}_1, \dots, \mathbf{e}_l)$. Therefore $\text{span}(E\mathbf{g}_1, \dots, E\mathbf{g}_n) \subseteq \text{span}(\mathbf{e}_1, \dots, \mathbf{e}_l)$, which implies that $\text{rank } EG \leq \text{rank } E$. Note that $(EG)^T = G^T E^T$. Hence $\text{rank } EG = \text{rank } (EG)^T \leq \text{rank } G^T = \text{rank } G$. Thus $\text{rank } EG \leq \min(\text{rank } E, \text{rank } G)$. Suppose E is invertible. Then $\text{rank } EG \leq \text{rank } G = \text{rank } E^{-1}(EG) \leq \text{rank } EG$. Hence $\text{rank } EG = \text{rank } G$. Similarly $\text{rank } EG = \text{rank } E$ if G is invertible. \square

Proof of Theorem 2.79.

1. Proposition 2.80 yields that $\text{rank } A^\dagger = \text{rank } V_r \Sigma_r^{-1} U_r^* \leq \text{rank } \Sigma_r^{-1} U_r^* \leq \text{rank } \Sigma_r^{-1} = r = \text{rank } A$. Since $\Sigma_r = V_r^* A^\dagger U_r$ Proposition 2.80 yields that $\text{rank } A^\dagger \geq \text{rank } \Sigma_r^{-1} = r$. Hence $\text{rank } A = \text{rank } A^\dagger$.

2. $AA^\dagger = (U_r \Sigma_r V_r^*)(V_r \Sigma_r^{-1} U_r^*) = U_r \Sigma_r \Sigma_r^{-1} U_r^* = U_r U_r^*$. Hence

$$AA^\dagger A = (U_r U_r^*)(U_r \Sigma_r V_r^*) = U_r \Sigma_r V_r^* = A.$$

Hence $A^*AA^\dagger = (V_r \Sigma_r U_r^*)(U_r U_r^*) = A^*$. Similarly $A^\dagger A = V_r V_r^*$ and $A^\dagger AA^\dagger = A^\dagger, A^\dagger AA^* = A^*$.

3. Since $AA^\dagger = U_r U_r^*$ we deduce that $(AA^\dagger)^* = (U_r U_r^*)^* = (U_r^*)^* U_r^* = AA^\dagger$, i.e. AA^\dagger is Hermitian. Next $(AA^\dagger)^2 = (U_r U_r^*)^2 = (U_r U_r^*)(U_r U_r^*) = (U_r U_r^*) = AA^\dagger$, i.e. AA^\dagger is idempotent. Hence AA^\dagger is nonnegative definite. As $AA^\dagger = U_r I_r U_r^*$, the arguments of part 1 yield that $\text{rank } AA^\dagger = r$. Similar arguments apply to $A^\dagger A = V_r V_r^*$.

4. Since $A^*AA^\dagger = A^*$ it follows that $A^*A(A^\dagger\mathbf{b}) = A^*\mathbf{b}$, i.e. $\mathbf{y} = A^\dagger\mathbf{b}$ is a least square solution. It is left to show that if $A^*A\mathbf{x} = A^*\mathbf{b}$ then $\|\mathbf{x}\| \geq \|A^\dagger\mathbf{b}\|$ and equality holds if and only if $\mathbf{x} = A^\dagger\mathbf{b}$.

We now consider the system $A^*A\mathbf{x} = A^*\mathbf{b}$. To analyze this system we use the full form of SVD given in (2.7). It is equivalent to

$(V\Sigma^T U^*)(U\Sigma V^*)\mathbf{x} = V\Sigma^T U^*\mathbf{b}$. Multiplying by V^* we obtain the system $\Sigma^T \Sigma(V^*\mathbf{x}) = \Sigma^T(U^*\mathbf{b})$. Let $\mathbf{z} = (z_1, \dots, z_n)^T := V^*\mathbf{x}$,

$\mathbf{c} = (c_1, \dots, c_m)^T := U^* \mathbf{b}$. Note that $\mathbf{z}^* \mathbf{z} = \mathbf{x}^* V V \mathbf{x} = \mathbf{x}^* \mathbf{x}$, i.e. $\|\mathbf{z}\| = \|\mathbf{x}\|$. After these substitutions the least square system in z_1, \dots, z_n variables is given in the form $\sigma_i(A)^2 z_i = \sigma_i(A) c_i$ for $i = 1, \dots, n$. Since $\sigma_i(A) = 0$ for $i > r$ we obtain that $z_i = \frac{1}{\sigma_i(A)} c_i$ for $i = 1, \dots, r$ while z_{r+1}, \dots, z_n are free variables. Thus $\|\mathbf{z}\|^2 = \sum_{i=1}^r \frac{1}{\sigma_i(A)^2} + \sum_{i=r+1}^n |z_i|^2$. Hence the least square solution with the minimal length $\|\mathbf{z}\|$ is the solution with $z_i = 0$ for $i = r + 1, \dots, n$. This solution corresponds the $\mathbf{x} = A^\dagger \mathbf{b}$.

5. Since $\text{rank } A^* A = \text{rank } A = n$ it follows that $A^* A$ is an invertible matrix. Hence the least square solution is unique and is given by $\mathbf{x} = (A^* A)^{-1} A^* \mathbf{b}$. Thus for each \mathbf{b} one has $(A^* A)^{-1} A^* \mathbf{b} = A^\dagger \mathbf{b}$, hence $A^\dagger = (A^* A)^{-1} A^*$.

If A is an $n \times n$ matrix and is invertible it follows that $(A^* A)^{-1} A^* = A^{-1} (A^*)^{-1} A^* = A^{-1}$. \square

Problems

- $P \in \mathbb{C}^{n \times n}$ is called a *projection* if $P^2 = P$. Show that P is a projection if and only if the following two conditions are satisfied:
 - Each eigenvalue of P is either 0 or 1.
 - P is a diagonalizable matrix.
- $P \in \mathbb{R}^{n \times n}$ is called an *orthogonal projection* if P is a projection and a symmetric matrix. Let $\mathbf{V} \subseteq \mathbb{R}^n$ be the subspace spanned by the columns of P . Show that for any $\mathbf{a} \in \mathbb{R}^n, \mathbf{b} \in P\mathbf{V}$ the following inequality holds: $\|\mathbf{a} - \mathbf{b}\| \geq \|\mathbf{a} - P\mathbf{a}\|$. Furthermore equality holds if and only if $\mathbf{b} = P\mathbf{a}$. That is, $P\mathbf{a}$ is the orthogonal projection of \mathbf{a} on the column space of P .
- Let $A \in \mathbb{R}^{m \times n}$ and assume that the SVD of A is given by (2.7), where $U \in \mathbf{O}(m, \mathbb{R}), V \in \mathbf{O}(n, \mathbb{R})$.
 - What is the SVD of A^T ?
 - Show that $(A^T)^\dagger = (A^\dagger)^T$.
 - Suppose that $B \in \mathbb{R}^{l \times m}$. Is it true that $(BA)^\dagger = A^\dagger B^\dagger$? Justify!

3 Jordan canonical form

3.1 Eigenvalues and eigenvectors

Let $A \in \mathbb{F}^{n \times n}$. $p(z) := \det(zI_n - A)$ is called the characteristic polynomial of A .

Lemma 3.1 *Let $A \in \mathbb{F}^{n \times n}$. Then $\det(zI_n - A) = z^n + \sum_{i=1}^n a_i z^{n-i}$. $(-1)^i a_i$ is the sum of all $i \times i$ principle minors of A . Assume that $\det(zI_n - A) = (z - z_1) \dots (z - z_n)$, where $z_1, \dots, z_n \in \mathbb{F}$. Then $(-1)^i a_i = \sigma_i(\mathbf{z})$ for $i = 1, \dots, n$.*

Proof. Consider $\det(zI - A)$. To obtain the coefficient of z^{n-i} we need to take the product of some $n - i$ diagonal elements of $zI_n - A$: $(z - a_{j_1 j_1}) \dots (z - a_{j_{n-i} j_{n-i}})$. We take z^{n-i} in this product. Then this product is multiplied by the $\det(-A[\boldsymbol{\alpha}, \boldsymbol{\alpha}])$,

where α is the complement of $\{j_1, \dots, j_{n-i}\}$ in the set $\langle n \rangle$. This shows that $(-1)^i a_i$ is the sum of all principal minors of A of order i .

Suppose that $\det(zI_n - A)$ splits to linear factors in $\mathbb{F}[z]$. Then (1.22) implies that $(-1)^i a_i = \sigma_i(\mathbf{z})$. \square

Corollary 3.2 *Let $A \in \mathbb{F}^{n \times n}$ and assume that $\det(zI_n - A) = \prod_{i=1}^n (z - z_i)$. Then*

$$\operatorname{tr} A := \sum_{i=1}^n a_{ii} = \sum_{i=1}^n z_i, \quad \det A = \prod_{i=1}^n z_i.$$

Definition 3.3 *Let $\operatorname{GL}(n, \mathbb{F}) \subset \mathbb{F}^{n \times n}$ denote the set (group) of all $n \times n$ invertible matrices with entries in a given field \mathbb{F} . $A, B \in \mathbb{F}^{n \times n}$ are called similar, and this is denoted by $A \sim B$, if $B = UAU^{-1}$ for some $U \in \operatorname{GL}(n, \mathbb{F})$. The set of all $B \in \mathbb{F}^{n \times n}$ similar to a fixed $A \in \mathbb{F}^{n \times n}$ is called the similarity class corresponding to A , or simply a similarity class.*

The following proposition is straightforward:

Proposition 3.4 *Let \mathbb{F} be a field, ($\mathbb{F} = \mathbb{R}, \mathbb{C}$). Then the similarity relation on $\mathbb{F}^{n \times n}$ is an equivalence relation:*

$$A \sim A, \quad A \sim B \iff B \sim A, \quad A \sim B \text{ and } B \sim C \Rightarrow A \sim B.$$

Furthermore if $B = UAU^{-1}$ then

1.

$$\det(zI_n - B) = \det(U(zI_n - A)U^{-1}) = \det U \det(zI_n - A) \det(U^{-1}) = \det(zI_n - A),$$

i.e. A and B have the same characteristic polynomial.

2. *For any integer $m \geq 2$ $B^m = UA^mU^{-1}$.*

3. *If in addition A is invertible, then B is invertible and $B^m = UA^mU^{-1}$ for any integer m .*

Corollary 3.5 *Let \mathbf{V} be n -dimensional vector space over \mathbb{F} . Assume that $T : \mathbf{V} \rightarrow \mathbf{V}$ is a linear transformation. Then the set of all representation matrices of T in different bases $\mathbf{v}_1, \dots, \mathbf{v}_n$ is a similarity class. (Use (1.24) and (1.25), where $m = n$, $\mathbf{x}_i = \mathbf{y}_i, i = 1, \dots, n, X = Y$.) Hence, the characteristic polynomial of T is defined as $\det(zI_n - A) = z^n + \sum_{i=0}^{n-1} a_i z^i$, where A is the representation matrix of T in any basis $[\mathbf{u}_1, \dots, \mathbf{u}_n]$, and this definition is independent of the choice of a basis. In particular $\det T := \det A$, and $\operatorname{trace} T^m = \operatorname{trace} A^m$ for any nonnegative integer. (T^0 is the identity operator, *i.e.* $T^0 \mathbf{v} = \mathbf{v}$ for all $\mathbf{v} \in \mathbf{V}$, and $A^0 = I$. Here by the trace of $B \in \mathbb{F}^{n \times n}$, denoted by $\operatorname{trace} B$, we mean the sum of all diagonal elements of B .)*

Problem 3.6 *(The representation problem.) Let \mathbf{V} be n -dimensional vector space over \mathbb{F} . Assume that $T : \mathbf{V} \rightarrow \mathbf{V}$ is a linear transformation. Find a basis $[\mathbf{v}_1, \dots, \mathbf{v}_n]$ in which T has the simplest form. Equivalently, given $A \in \mathbb{F}^{n \times n}$ find $B \sim A$ of the simplest form.*

In the following case the answer is well known. Recall that $\mathbf{v} \in \mathbf{V}$ is called an *eigenvector* of T corresponding to the *eigenvalue* $\lambda \in \mathbb{F}$, if $\mathbf{v} \neq \mathbf{0}$ and $T\mathbf{v} = \lambda\mathbf{v}$. This is equivalent to the existence $\mathbf{0} \neq \mathbf{x} \in \mathbb{F}^n$ such that $A\mathbf{x} = \lambda\mathbf{x}$. Hence $(\lambda I - A)\mathbf{x} = \mathbf{0}$ which implies that $\det(\lambda I - A) = 0$. Hence λ is the zero of the characteristic polynomial of A and T . The assumption λ is a zero of the characteristic polynomial yields that the system $(\lambda I - A)\mathbf{x}$ has a nontrivial solution $\mathbf{x} \neq \mathbf{0}$.

Corollary 3.7 *Let $A \in \mathbb{F}^{n \times n}$. Then λ is an eigenvalue of A if and only if λ is a zero of the characteristic polynomial of A : $\det(zI - A)$. Let \mathbf{V} be n -dimensional vector space over \mathbb{F} . Assume that $T : \mathbf{V} \rightarrow \mathbf{V}$ is a linear transformation. Then λ is an eigenvalue of T if and only if λ is a zero of the characteristic polynomial of T .*

Proposition 3.8 *Let \mathbf{V} be n -dimensional vector space over \mathbb{F} . Assume that $T : \mathbf{V} \rightarrow \mathbf{V}$ is a linear transformation. Then there exists a basis in V such that T is represented in this basis by a diagonal matrix*

$$\text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n) := \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \lambda_n \end{bmatrix},$$

if and only if the characteristic polynomial of T is $(z - \lambda_1)(z - \lambda_2) \dots (z - \lambda_n)$, and \mathbf{V} has a basis consisting of eigenvectors of T .

Equivalently, $A \in \mathbb{F}^{n \times n}$ is similar to a diagonal matrix $\text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$ if and only if $\det(zI - A) = (z - \lambda_1)(z - \lambda_2) \dots (z - \lambda_n)$, and A has n -linearly independent eigenvectors.

Proof. Assume that there exists a basis $[\mathbf{u}_1, \dots, \mathbf{u}_n]$ in V such that T is represented in this basis by a diagonal matrix $\Lambda := \text{diag}(\lambda_1, \dots, \lambda_n)$. Then the characteristic polynomial of T is $\det(zI - \Lambda) = \prod_{i=1}^n (z - \lambda_i)$. From the definition of the representation matrix of T , it follows that $T\mathbf{u}_i = \lambda_i\mathbf{u}_i$ for $i = 1, \dots, n$. Since each $\mathbf{u}_i \neq \mathbf{0}$, we deduce that each \mathbf{u}_i is an eigenvector of T . By our assumption $\mathbf{u}_1, \dots, \mathbf{u}_n$ for a basis in \mathbf{V} .

Assume now that \mathbf{V} has a basis $[\mathbf{u}_1, \dots, \mathbf{u}_n]$ consisting eigenvectors of T . So $T\mathbf{u}_i = \lambda_i\mathbf{u}_i$ for $i = 1, \dots, n$. Hence Λ is the representation matrix of T in the basis $[\mathbf{u}_1, \dots, \mathbf{u}_n]$.

To prove the corresponding results for $A \in \mathbb{F}^{n \times n}$, let $\mathbf{V} := \mathbb{F}^n$ and define the linear operator $T\mathbf{x} := A\mathbf{x}$ for all $\mathbf{x} \in \mathbb{F}^n$. □

Lemma 3.9 *Let $A \in \mathbb{F}^{n \times n}$ and assume that $\mathbf{x}_1, \dots, \mathbf{x}_k$ be k eigenvectors corresponding to k eigenvalues $\lambda_1, \dots, \lambda_k$. Suppose that $\lambda_i \neq \lambda_j$ for $i \neq j$. Then $\mathbf{x}_1, \dots, \mathbf{x}_k$ are linearly independent.*

Proof. We prove by induction on k . For $k = 1$ $\mathbf{x}_1 \neq \mathbf{0}$, hence \mathbf{x}_1 is linearly independent. Assume that the lemma holds for $k = m - 1$. Suppose that $k = m$. Assume that $\sum_{i=1}^m a_i \mathbf{x}_i = \mathbf{0}$. So

$$\mathbf{0} = A\mathbf{0} = A \sum_{i=1}^m a_i \mathbf{x}_i = \sum_{i=1}^m a_i A\mathbf{x}_i = \sum_{i=1}^m a_i \lambda_i \mathbf{x}_i.$$

Multiply the equality $\sum_{i=1}^m a_i \mathbf{x}_i = \mathbf{0}$ by λ_m and subtract it from the above inequality to deduce that $\sum_{i=1}^{m-1} a_i (\lambda_i - \lambda_m) \mathbf{x}_i = \mathbf{0}$. Since $\mathbf{x}_1, \dots, \mathbf{x}_{m-1}$ are linearly independent by the induction hypothesis, we deduce that $a_i (\lambda_i - \lambda_m) = 0$ for $i = 1, \dots, m-1$. As $\lambda_i - \lambda_m \neq 0$ for $i < m$ we get that $a_i = 0$ for $i = 1, \dots, m-1$. The assumption that $\sum_{i=1}^m a_i \mathbf{x}_i = \mathbf{0}$ yields that $a_m \mathbf{x}_m = \mathbf{0}$. Since $\mathbf{x}_m \neq \mathbf{0}$ we obtain that $a_m = 0$. Hence $a_1 = \dots = a_m = 0$. \square

Theorem 3.10 *Let \mathbf{V} be n -dimensional vector space over \mathbb{F} . Assume that $T : \mathbf{V} \rightarrow \mathbf{V}$ is a linear transformation. Assume that the characteristic polynomial of T $p(z)$ has n distinct roots over \mathbb{F} , i.e. $p(z) = \prod_{i=1}^n (z - \lambda_i)$ where $\lambda_1, \dots, \lambda_n \in \mathbb{F}$, and $\lambda_i \neq \lambda_j$ for each $i \neq j$. Then there exists a basis in \mathbf{V} in which T is represented by a diagonal matrix.*

Similarly, let $A \in \mathbb{F}^{n \times n}$ and assume that $\det(zI - A)$ has n distinct roots in \mathbb{F} . Then A is similar to a diagonal matrix.

Proof. It is enough to consider the case of the linear transformation T . Recall that each root of the characteristic polynomial of T is an eigenvalue of T (Corollary 3.7). Hence to each λ_i corresponds an eigenvector \mathbf{u}_i : $T\mathbf{u}_i = \lambda_i \mathbf{u}_i$. Then the proof of the theorem follows Lemma 3.9 and Proposition 3.8. \square

Given $A \in \mathbb{F}^{n \times n}$ it may happen that $\det(zI - A)$ does not have n roots in \mathbb{F} . (See for example Problem 2 of this section.) Hence we can not *diagonalize* A , i.e. A is not similar to a diagonal matrix. If \mathbb{F} is *algebraically closed*, i.e. any $\det(zI - A)$ has n roots in \mathbb{F} we can apply Proposition 3.8 in general and Theorem 3.10 in particular to see if A is diagonalizable.

Since \mathbb{R} is not algebraically closed and \mathbb{C} is, that is the reason that we sometimes view a real valued matrix $A \in \mathbb{R}^{n \times n}$ as a complex valued matrix $A \in \mathbb{C}^{n \times n}$. (See Problem 2 of this section.)

Corollary 3.11 *Let $A \in \mathbb{C}^{n \times n}$ be nondiagonalizable. Then its characteristic polynomial must have a multiple root.*

Definition 3.12 1. *Let k be a positive integer and $\lambda \in \mathbb{F}$. Then $J_k(\lambda) \in \mathbb{F}^{k \times k}$ be a $k \times k$ upper triangular matrix, with λ on the main diagonal, 1 on the next sub-diagonal and other entries are equal to 0 for $k > 1$:*

$$J_k(\lambda) := \begin{bmatrix} \lambda & 1 & 0 & \dots & 0 & 0 \\ 0 & \lambda & 1 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & \lambda & 1 \\ 0 & 0 & 0 & \dots & 0 & \lambda \end{bmatrix},$$

$$(J_1(\lambda) = [\lambda].)$$

2. *Let $A_i \in \mathbb{F}^{n_i \times n_i}$ for $i = 1, \dots, k$. Denote by*

$$\oplus_{i=1}^k A_i = A_1 \oplus A_2 \oplus \dots \oplus A_k = \text{diag}(A_1, A_2, \dots, A_k) := \begin{bmatrix} A_1 & 0 & \dots & 0 \\ 0 & A_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & A_k \end{bmatrix} \in \mathbb{F}^{n \times n}, \quad n = n_1 + n_2 + \dots + n_k,$$

the $n \times n$ block diagonal matrix, whose blocks are A_1, A_2, \dots, A_k .

Theorem 3.13 (The Jordan Canonical Form) Let $A \in \mathbb{C}^{n \times n}$, ($A \in \mathbb{F}^{n \times n}$, where \mathbb{F} is an algebraically closed field.) Then A is similar to its Jordan canonical form $\bigoplus_{i=1}^k J_{n_i}(\lambda_i)$ for some $\lambda_1, \dots, \lambda_k \in \mathbb{C}$, ($\lambda_1, \dots, \lambda_k \in \mathbb{F}$), and positive integers n_1, \dots, n_k . The Jordan canonical form is unique up to the permutations of the Jordan blocks $J_{n_1}(\lambda_1), \dots, J_{n_k}(\lambda_k)$.

Equivalently, let $T : \mathbf{V} \rightarrow \mathbf{V}$ be a linear transformation of an n -dimensional space over \mathbb{C} , or any other algebraically closed field. Then there exists a basis in \mathbf{V} , such that $\bigoplus_{i=1}^k J_{n_i}(\lambda_i)$ is the representation matrix of T in this basis. The blocks $J_{n_i}(\lambda_i), i = 1, \dots, k$ are unique.

Note that $A \in \mathbb{C}^{n \times n}$ is diagonalizable if and only in its Jordan canonical form $k = n$, i.e. $n_1 = \dots = n_n = 1$. For $k < n$, the Jordan canonical form is the simplest form of the similarity class of a nondiagonalizable $A \in \mathbb{C}^{n \times n}$.

We will prove Theorem 3.13 in the next several sections.

Problems

- Let \mathbf{V} be a vector space over \mathbb{F} . (You may assume that $\mathbb{F} = \mathbb{C}$.) Let $T : \mathbf{V} \rightarrow \mathbf{V}$ be a linear transformation. Suppose that \mathbf{u}_i is an eigenvector of T with the corresponding eigenvalue λ_i for $i = 1, \dots, m$. Show by induction on m that if $\lambda_1, \dots, \lambda_m$ are m distinct scalars then $\mathbf{u}_1, \dots, \mathbf{u}_m$ are linearly independent.
- Let $A = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix} \in \mathbb{R}^{2 \times 2}$.
 - Show that A is not diagonalizable over the real numbers \mathbb{R} .
 - Show that A is diagonalizable over the complex numbers \mathbb{C} . Find $U \in \mathbb{C}^{2 \times 2}$ and a diagonal $\Lambda \in \mathbb{C}^{2 \times 2}$ such that $A = U\Lambda U^{-1}$.
- Let $A = \bigoplus_{i=1}^k J_{n_i}(\lambda_i)$. Show that $\det(zI - A) = \prod_{i=1}^k (z - \lambda_i)^{n_i}$. (You may use the fact that the determinant of an upper triangular matrix is the product of its diagonal entries.)
- Let $A = \bigoplus_{i=1}^k A_i$ where $A_i \in \mathbb{C}^{n_i \times n_i}, i = 1, \dots, k$. Show that $\det(zI_n - A) = \prod_{i=1}^k \det(zI_{n_i} - A_i)$. (First show the identity for $k = 2$ using the determinant expansion by rows. Then use induction for $k > 2$.)
 - Show that any eigenvector of $J_n(\lambda) \in \mathbb{C}^{n \times n}$ is in the subspace spanned by \mathbf{e}_1 . Conclude that $J_n(\lambda)$ is not diagonalizable unless $n = 1$.
 - What is the rank of $zI_n - J_n(\lambda)$ for a fixed $\lambda \in \mathbb{C}$ and for each $z \in \mathbb{C}$?
 - What is the rank of $zI - \bigoplus_{i=1}^k J_{n_i}(\lambda_i)$ for fixed $\lambda_1, \dots, \lambda_k \in \mathbb{C}$ and for each $z \in \mathbb{C}$?
- Let $A \in \mathbb{C}^{n \times n}$ and assume that $\det(zI_n - A) = z^n + a_1 z^{n-1} + \dots + a_{n-1} z + a_n$ has n distinct complex roots. Show that $A^n + a_1 A^{n-1} + \dots + a_{n-1} A + a_n I_n = \mathbf{0}$, where $\mathbf{0} \in \mathbb{C}^{n \times n}$ denotes the zero matrix, i.e. the matrix whose all entries are

0. (This is a special case of the Cayley-Hamilton theorem, which claims that the above identity holds for *any* $A \in \mathbb{C}^{n \times n}$.) **Hint:** Use the fact that A is diagonalizable.

3.2 Matrix polynomials

Let $P(z) = (p_{ij}(z))_{i,j=1}^{m,n}$ be an $m \times n$ matrix whose entries are polynomials in $\mathbb{F}[z]$. The set of all such $m \times n$ matrices is denoted by $\mathbb{F}[z]^{m \times n}$. Clearly $\mathbb{F}[z]^{m \times n}$ is a vector space over \mathbb{F} , of infinite dimension. Given $p(z) \in \mathbb{F}[z]$ and $P(z) \in \mathbb{F}[z]^{m \times n}$ one can define $p(z)P(z) := (p(z)p_{ij}) \in \mathbb{F}[z]$. Again, this product satisfies nice distribution properties. Thus $\mathbb{F}[z]$ is a *module* over the ring $\mathbb{F}[z]$. (Note $\mathbb{F}[z]$ is not a field!)

Let $P(z) = (p_{ij}(z)) \in \mathbb{F}[z]^{m \times n}$. Then $\deg P(z) := \max_{i,j} \deg p_{ij}(z) = l$. Write

$$p_{ij}(z) = \sum_{k=0}^l p_{ij,k} z^{l-k}, \quad P_k := (p_{ij,k})_{i,j=1}^{m,n} \in \mathbb{F}^{m \times n} \text{ for } k = 0, \dots, l.$$

Then

$$P(z) = P_0 z^l + P_1 z^{l-1} + \dots + P_l, \quad P_i \in \mathbb{F}^{m \times n}, \quad i = 0, \dots, l, \quad (3.1)$$

is a matrix polynomial with coefficients in $\mathbb{F}^{m \times n}$.

Assume that $P(z), Q(z) \in \mathbb{F}[z]^{n \times n}$. Then we can define $P(z)Q(z) \in \mathbb{F}[z]$. Note that in general $P(z)Q(z) \neq Q(z)P(z)$. Hence $\mathbb{F}[z]^{n \times n}$ is a *noncommutative* ring. For $P(z) \in \mathbb{F}^{n \times n}$ of the form (3.1) and any $A \in \mathbb{F}^{n \times n}$ we define

$$P(A) = \sum_{i=0}^l P_i A^{l-i} = P_0 A^l + P_1 A^{l-1} + \dots + P_l, \quad \text{where } A^0 = I_n.$$

Recall that given two polynomials $p, q \in \mathbb{F}[z]$ one can divide p by $q \neq 0$ with the residue r , i.e. $p = tq + r$ for some unique $t, r \in \mathbb{F}[z]$, where $\deg r < \deg q$. One can trivially generalize that to polynomial matrices:

Proposition 3.14 *Let $p(z), q(z) \in \mathbb{F}[z]$ and assume that $q(z) \neq 0$. Let $p(z) = t(z)q(z) + r(z)$, where $t(z), r(z) \in \mathbb{F}[z]$ are unique polynomials with $\deg r(z) < \deg q(z)$. Let $n > 1$ be an integer, and define the following scalar polynomials: $P(z) := p(z)I_n, Q(z) := q(z)I_n, T(z) := t(z)I_n, R(z) := r(z)I_n \in \mathbb{F}[z]^{n \times n}$. Then $P(A) = T(A)Q(A) + R(A)$ for any $A \in \mathbb{F}^{n \times n}$.*

Proof. Since $A^i A^j = A^{i+j}$ for any nonnegative integer, with $A^0 = I_n$, the equality $P(A) = T(A)Q(A) + R(A)$ follows trivially from the equality $p(z) = t(z)q(z) + r(z)$. \square

Recall that p is divisible by q , denoted as $q|p$, if $p = tq$, i.e. r is the zero polynomial. Note that if $q(z) = (z - a)$ then $p(z) = t(z)(z - a) + p(a)$. Thus $(z - a)|p$ if and only if $p(a) = 0$. Similar results hold for square polynomial matrices, which are not scalar.

Lemma 3.15 *Let $P(z) \in \mathbb{F}[z]^{n \times n}, A \in \mathbb{F}^{n \times n}$. Then there exists a unique $T_{left}(z)$, of degree $\deg P - 1$ if $\deg P > 0$ or degree $-\infty$ if $\deg P \leq 0$, such that*

$$P(z) = T_{left}(z)(zI - A) + P(A). \quad (3.2)$$

In particular, $P(z)$ is divisible from the right by $zI - A$ if and only if $P(A) = 0$.

Proof. We prove the lemma by induction on $\deg P$. If $\deg P \leq 0$, i.e. $P(z) = P_0 \in \mathbb{F}^{n \times n}$ then $T_{left} = \mathbf{0}$, $P(A) = P_0$ and the lemma trivially holds. Suppose that the lemma holds for all P with $\deg P \leq l-1$, where $l \geq 1$. Let $P(z)$ be of degree $l \geq 1$ of the form (3.1). Then $P(z) = P_0 z^l + \tilde{P}(z)$, where $\tilde{P}(z) = \sum_{i=1}^l P_i z^{l-i}$. By the induction assumption $\tilde{P}(z) = \hat{T}_{left}(z)(zI_n - A) + \tilde{P}(A)$, where $\hat{T}_{left}(z)$ is unique. A straightforward calculation shows that

$$P_0 z^l = \hat{T}_{left}(z)(zI_n - A) + P_0 A^l, \text{ where } \hat{T}_{left}(z) = \sum_{i=0}^{l-1} P_0 A^i z^{l-i-1},$$

and \hat{T}_{left} is unique. Hence $T_{left}(z) = \hat{T}_{left}(z) + \tilde{T}_{left}$ is unique, $P(A) = P_0 A^l + \tilde{P}(A)$ and (3.2) follows.

Suppose that $P(A) = \mathbf{0}$. Then $P(z) = T_{left}(z)(zI - A)$, i.e. $P(z)$ is divisible by $zI_n - A$ from the right. Assume that $P(z)$ is divisible by $(zI_n - A)$ from the right, i.e. there exists $T(z) \in \mathbb{F}[z]^{n \times n}$ such that $P(z) = T(z)(zI_n - A)$. Subtract (3.2) from $P(z) = T(z)(zI_n - A)$ to deduce that $\mathbf{0} = (T(z) - T_{left}(z))(zI_n - A) - P(A)$. Hence $T(z) = T_{left}(z)$ and $P(A) = \mathbf{0}$. \square

The above lemma can be generalized to any $Q(z) = Q_0 z^l + Q_1 z^{l-1} + \dots + Q_l \in \mathbb{F}[z]$, where $Q_0 \in \text{GL}(n, \mathbb{F})$: There exists unique $T_{left}(z), R_{left}(z) \in \mathbb{F}[z]$ such that

$$P(z) = T_{left}(z)Q(z) + R_{left}(z), \deg R_{left} < \deg Q, Q(z) = \sum_{i=0}^l Q_i z^{l-i}, Q_0 \in \text{GL}(n, \mathbb{F}). \quad (3.3)$$

Here we agree that $(Az^i)(Bz^j) = (AB)z^{i+j}$ for any $A, B \in \mathbb{F}^{n \times n}$ and nonnegative integers i, j .

Theorem 3.16 (*Cayley-Hamilton theorem.*) *Let $A \in \mathbb{F}^{n \times n}$ and $p(z) = \det(zI_n - A)$ be the characteristic polynomial of A . Let $P(z) = p(z)I_n \in \mathbb{F}[z]^{n \times n}$. Then $P(A) = \mathbf{0}$.*

Proof. Let $A(z) = zI_n - A$. Fix $z \in \mathbb{F}$ and let $B(z) = (b_{ij}(z))$ be the adjoint matrix of $A(z)$, whose entries are the cofactors of $A(z)$. That is $b_{ij}(z)$ is $(-1)^{i+j}$ times the determinant of the matrix obtained from $A(z)$ by deleting row j and column i . If one views z as indeterminate then $B(z) \in \mathbb{F}[z]^{n \times n}$. Recall the identity

$$A(z)B(z) = B(z)A(z) = \det A(z)I_n = p(z)I_n = P(z).$$

Hence $(zI_n - A)$ divides from the right $P(z)$. Lemma 3.15 yields that $P(A) = \mathbf{0}$. \square

For $p, q \in \mathbb{F}[z]$ let (p, q) be the *greatest common divisor* of p, q . If p and q are identically zero then (p, q) is the zero polynomial. Otherwise (p, q) is a polynomial s of the highest degree that divides p and q . s is determined up to a multiple of a nonzero scalar. s can be chosen as a unique *monic* polynomial:

$$s(z) = z^l + s_1 z^{l-1} + \dots + s_l \in \mathbb{F}[z]. \quad (3.4)$$

Equality (1.19) yields.

Corollary 3.17 *Let $p, q \in \mathbb{F}[z]$ be coprime. Then there exists $u, v \in \mathbb{F}[z]$ such that $1 = up + vq$. Let $n > 1$ be an integer and define $P(z) := p(z)I_n, Q(z) := q(z)I_n, U(z) := u(z)I_n, V(z) := v(z)I_n \in \mathbb{F}[z]^{n \times n}$. Then for any $A \in \mathbb{F}^{n \times n}$ we have the identity $I_n = U(A)P(A) + V(A)Q(A)$, where $U(A)P(A) = P(A)U(A)$ and $V(A)Q(A) = Q(A)V(A)$.*

Let us consider that case where $p, q \in \mathbb{F}[z]$ are both nonzero polynomials that split (to linear factors) over \mathbb{F} . So

$$p(z) = p_0(z - \alpha_1) \dots (z - \alpha_i), p_0 \neq 0, \quad q(z) = q_0(z - \beta_1) \dots (z - \beta_j), q_0 \neq 0.$$

In that case $(p, q) = 1$, if p and q do not have a common root. If p and q have a common zero then (p, q) is a nonzero polynomial that has the maximal number of common roots of p and q counting with multiplicities.

From now on for any $p \in \mathbb{F}[z]$ and $A \in \mathbb{F}^{n \times n}$ we identify $p(A)$ with $P(A)$, where $P(z) = p(z)I_n$.

3.3 Minimal polynomial and decomposition to invariant subspaces

Recall that $\mathbb{F}^{n \times n}$ is a vector space over \mathbb{F} of dimension n^2 . Let $A \in \mathbb{F}^{n \times n}$ and consider the powers $A^0 = I_n, A, A^2, \dots, A^m$. Let m be the smallest positive integer such that these $m+1$ matrices are linearly dependent as vectors in $\mathbb{F}^{n \times n}$. (Note that $A^0 \neq \mathbf{0}$.) So $\sum_{i=0}^m b_i A^{m-i} = \mathbf{0}$, and $(b_0, \dots, b_m)^\top \neq \mathbf{0}$. If $b_0 = 0$ then A^0, \dots, A^{m-1} are linearly dependent, which contradicts the definition of m . Hence $b_0 \neq 0$. Divide the linear dependence by b_0 to obtain.

$$\psi(A) = 0, \quad \psi(z) = z^m + \sum_{i=1}^m a_i z^{m-i} \in \mathbb{F}[z], \quad a_i = \frac{b_i}{b_0} \text{ for } i = 1, \dots, m. \quad (3.5)$$

ψ is called the *minimal polynomial* of A . In principle $m \leq n^2$, but in reality $m \leq n$:

Theorem 3.18 *Let $A \in \mathbb{F}^{n \times n}$ and $\psi(z)$ be its characteristic polynomial. Assume that $p(z) \in \mathbb{F}[z]$ is an annihilated polynomial of A , i.e. $p(A) = \mathbf{0}$. Then ψ divides p . In particular, the characteristic polynomial $p(z) = \det(zI_n - A)$ is divisible by $\psi(z)$. Hence $\deg \psi \leq \deg p = n$.*

Proof. Divide the annihilating polynomial p by ψ to obtain $p(z) = t(z)\psi(z) + r(z)$, where $\deg r < \deg \psi = m$. Proposition 3.14 yields that $p(A) = t(A)\psi(A) + r(A)$ which implies that $r(A) = \mathbf{0}$. Assume that $l = \deg r(z) \geq 0$, i.e. r is not identically the zero polynomial. So A^0, \dots, A^l are linearly dependent, which contradicts the definition of m . Hence $r(z) \equiv 0$.

The Cayley-Hamilton theorem yields that the characteristic polynomial $p(z)$ of A annihilates A . Hence $\psi|p$ and $\deg \psi \leq \deg p = n$. \square

Theorem 3.19 *Let $A \in \mathbb{F}^{n \times n}$. Denote by $q(z)$ the g.c.d., (monic polynomial), of all the entries of $\text{adj}(zI_n - A) \in \mathbb{F}[z]^{n \times n}$, which is the g.c.d of all $(n-1) \times (n-1)$ minors of $(zI_n - A)$. Then $q(z)$ divides the characteristic polynomial $p(z) = \det(zI_n - A)$ of A . Furthermore, the minimal polynomial of A is equal to $\frac{p(z)}{q(z)}$.*

Proof. Expand $\det(zI_n - A)$ by the first column. Since $q(z)$ divides each $(n - 1) \times (n - 1)$ minor of $zI_n - A$ we deduce that $q(z)$ divides $p(z)$. Let $\phi(z) = \frac{p(z)}{q(z)}$. As $q(z)$ divides each entry of $\text{adj}(zI_n - A)$ we deduce that $\text{adj}(zI_n - A) = q(z)C(z)$ for some $C(z) \in \mathbb{F}[z]^{n \times n}$. Divide the equality $p(z)I_n = \text{adj}(zI_n - A)(zI_n - A)$ by $q(z)$ to deduce that $\phi(z)I_n = C(z)(zI_n - A)$. Lemma 3.15 yields that $\phi(A) = 0$.

Let $\psi(z)$ be the minimal polynomial of A . Theorem 3.18 yields that ψ divides ϕ . We now show that ϕ divides ψ . Theorem 3.18 implies that $p(z) = s(z)\psi(z)$ for some monic polynomial $s(z)$. Since $\psi(A) = 0$ Lemma 3.15 yields that $\psi(z)I_n = D(z)(zI_n - A)$ for some $D(z) \in \mathbb{F}^{n \times n}[z]$. So $p(z)I_n = s(z)\psi(z)I_n = s(z)D(z)(zI_n - A)$. Since $p(z)I_n = \text{adj}(zI_n - A)(zI_n - A)$ we deduce that $s(z)D(z) = \text{adj}(zI_n - A)$. As all the entries of $D(z)$ are polynomials, it follows that $s(z)$ divides all the entries of $\text{adj}(zI_n - A)$. Since $q(z)$ is the g.c.d. of all entries of $\text{adj}(zI_n - A)$ we deduce that $s(z)$ divides $q(z)$. Consider the equality $p(z) = s(z)\psi(z) = q(z)\phi(z)$. Thus $\psi(z) = \frac{q(z)}{s(z)}\phi(z)$. Hence $\phi(z)$ divides $\psi(z)$. As $\psi(z)$ and $\phi(z)$ are monic we deduce that $\psi(z) = \phi(z)$. \square

Proposition 3.20 *Let $A \in \mathbb{F}^{n \times n}$ and assume that $\lambda \in \mathbb{F}$ is an eigenvalue of A with the corresponding eigenvector $\mathbf{x} \in \mathbb{F}^n$. Then for any $h(z) \in \mathbb{F}[z]$, $h(A)\mathbf{x} = h(\lambda)\mathbf{x}$. In particular, λ is a root of the minimal polynomial $\psi(z)$ of A , i.e. $\psi(\lambda) = 0$.*

Proof. Clearly, $A^m\mathbf{x} = \lambda^m\mathbf{x}$. Hence, $h(A)\mathbf{x} = h(\lambda)\mathbf{x}$. Assume that $h(A) = 0$. As $\mathbf{x} \neq \mathbf{0}$ we deduce that $h(\lambda) = 0$. Hence $\psi(\lambda) = 0$. \square

Definition 3.21 *A matrix $A \in \mathbb{F}^{n \times n}$ is called nonderogatory if the minimal polynomial of A is equal to its characteristic polynomial.*

Definition 3.22 *Let \mathbf{V} be a finite dimensional vector space over \mathbb{F} , and assume that $\mathbf{V}_1, \dots, \mathbf{V}_i$ nonzero subspaces of \mathbf{V} . Then \mathbf{V} is a direct sum of $\mathbf{V}_1, \dots, \mathbf{V}_i$, denoted as $\mathbf{V} = \bigoplus_{j=1}^i \mathbf{V}_j$ if any vector $\mathbf{v} \in \mathbf{V}$ has a unique representation as $\mathbf{v} = \mathbf{v}_1 + \dots + \mathbf{v}_i$, where $\mathbf{v}_j \in \mathbf{V}_j$ for $j = 1, \dots, i$. Equivalently, let $[\mathbf{v}_{j1}, \dots, \mathbf{v}_{jl_j}]$ be a basis of \mathbf{V}_j for $j = 1, \dots, i$. Then $\dim \mathbf{V} = \sum_{j=1}^i \dim \mathbf{V}_j = \sum_{j=1}^i l_j$ and the $\dim \mathbf{V}$ vectors $\mathbf{v}_{11}, \dots, \mathbf{v}_{1l_1}, \dots, \mathbf{v}_{i1}, \dots, \mathbf{v}_{il_i}$ are linearly independent.*

Let $T : \mathbf{V} \rightarrow \mathbf{V}$ be a linear operator. A subspace \mathbf{U} of \mathbf{V} is called a T -invariant subspace, or simply an invariant subspace when there is no ambiguity about T , if $T\mathbf{u} \in \mathbf{U}$ for each $\mathbf{u} \in \mathbf{U}$. We denote this fact by $T\mathbf{U} \subseteq \mathbf{U}$. Denote by $T|_{\mathbf{U}}$ the restriction of T to the invariant subspace of T . Clearly, $T|_{\mathbf{U}}$ is a linear operator on \mathbf{U} .

Note \mathbf{V} and the zero subspace $\{\mathbf{0}\}$, (which consist only of the zero element), are invariant subspaces. Those are called *trivial* invariant subspaces. \mathbf{U} is called a *nontrivial* invariant subspace if \mathbf{U} is an invariant subspace such that $0 < \dim \mathbf{U} < \dim \mathbf{V}$.

Since the representation matrices of T in different bases form a similarity class we can define the *minimal polynomial* $\psi(z) \in \mathbb{F}[z]$ of T , as the minimal polynomial of any representation matrix of T . (See Problem 1 in the end of this section.) Equivalently $\psi(z)$ is the monic polynomial of the minimal degree which annihilates T : $\psi(T) = \mathbf{0}$.

Theorem 3.23 Let $T : \mathbf{V} \rightarrow \mathbf{V}$ be a linear operator on a finite dimensional space $\dim \mathbf{V} > 0$. Let $\psi(z)$ be the minimal polynomial of T . Assume that $\psi(z)$ decomposes to $\psi(z) = \psi_1(z) \dots \psi_k(z)$, where each $\psi_i(z)$ is a monic polynomial of degree at least 1. Suppose furthermore that for each pair $i \neq j$ $\psi_i(z)$ and $\psi_j(z)$ are coprime. Then \mathbf{V} is a direct sum of $\mathbf{V}_1, \dots, \mathbf{V}_k$, where each \mathbf{V}_i is a nontrivial invariant subspace of T . Furthermore the minimal polynomial of $T|_{\mathbf{V}_i}$ is equal to $\psi_i(z)$ for $i = 1, \dots, k$. Moreover, each \mathbf{V}_i is uniquely determined by $\psi_i(z)$ for $i = 1, \dots, k$.

Proof. We prove the theorem by induction on $k \geq 2$. Let $k = 2$. So $\psi(z) = \psi_1(z)\psi_2(z)$. Let $\mathbf{V}_1 := \psi_2(T)\mathbf{V}, \mathbf{V}_2 = \psi_1(T)\mathbf{V}$ be the ranges of the operators $\psi_2(T), \psi_1(T)$ respectively. Observe that

$$T\mathbf{V}_1 = T(\psi_2(T)\mathbf{V}) = (T\psi_2(T))\mathbf{V} = (\psi_2(T)T)\mathbf{V} = \psi_2(T)(T\mathbf{V}) \subseteq \psi_2(T)\mathbf{V} = \mathbf{V}_1.$$

Thus \mathbf{V}_1 is a T -invariant subspace. Assume that $\mathbf{V}_1 = \{\mathbf{0}\}$. This is equivalent to that $\psi_2(T) = 0$. By Theorem 3.18 ψ divides ψ_2 which is impossible since $\deg \psi = \deg \psi_1 + \deg \psi_2 > \deg \psi_1$. Thus $\dim \mathbf{V}_1 > 0$. Similarly \mathbf{V}_2 is a nonzero T -invariant subspace. Let $T_i = T|_{\mathbf{V}_i}$ for $i = 1, 2$. Clearly

$$\psi_1(T_1)\mathbf{V}_1 = \psi_1(T)\mathbf{V}_1 = \psi_1(T)(\psi_2(T)\mathbf{V}) = (\psi_1(T)\psi_2(T))\mathbf{V} = \{\mathbf{0}\},$$

since ψ is the minimal polynomial of T . Hence $\psi_1(T_1) = \mathbf{0}$, i.e. ψ_1 is an annihilating polynomial of T_1 . Similarly, $\psi_2(T_2) = \mathbf{0}$.

Let $\mathbf{U} = \mathbf{V}_1 \cap \mathbf{V}_2$. Then \mathbf{U} is an invariant subspace of T . We claim that $\mathbf{U} = \{\mathbf{0}\}$, i.e. $\dim \mathbf{U} = 0$. Assume to the contrary that $\dim \mathbf{U} \geq 1$. Let $Q := T|_{\mathbf{U}}$ and denote by $\phi \in \mathbb{F}[z]$ the minimal polynomial of Q . Clearly $\deg \phi \geq 1$. Since $\mathbf{U} \subseteq \mathbf{V}_i$ it follows that ψ_i is an annihilating polynomial of Q for $i = 1, 2$. Hence $\phi|\psi_1$ and $\phi|\psi_2$, i.e. ϕ is a nontrivial factor of ψ_1 and ψ_2 . This contradicts the assumption that ψ_1 and ψ_2 are coprime. Hence $\mathbf{V}_1 \cap \mathbf{V}_2 = \{\mathbf{0}\}$.

Since $(\psi_1, \psi_2) = 1$ there exists polynomials $f, g \in \mathbb{F}[z]$ such that $\psi_1 f + \psi_2 g = 1$. Hence $I = \psi_1(T)f(T) + \psi_2(T)g(T)$, where I is the identity operator $I\mathbf{v} = \mathbf{v}$ on \mathbf{V} . In particular for any $\mathbf{v} \in \mathbf{V}$ we have $\mathbf{v} = \mathbf{v}_2 + \mathbf{v}_1$, where $\mathbf{v}_1 = \psi_2(T)(g(T)\mathbf{v}) \in \mathbf{V}_1, \mathbf{v}_2 = \psi_1(T)(f(T)\mathbf{v}) \in \mathbf{V}_2$. Since $\mathbf{V}_1 \cap \mathbf{V}_2 = \{\mathbf{0}\}$ it follows that $\mathbf{V} = \mathbf{V}_1 \oplus \mathbf{V}_2$. Let $\tilde{\psi}_i$ be the minimal polynomial of T_i . Then $\tilde{\psi}_i|\psi_i$ for $i = 1, 2$. Hence $\tilde{\psi}_1\tilde{\psi}_2|\psi_1\psi_2$. Let $\mathbf{v} \in \mathbf{V}$. Then $\mathbf{v} = \mathbf{v}_1 + \mathbf{v}_2$, where $\mathbf{v}_i \in \mathbf{V}_i, i = 1, 2$. Using the facts that $\tilde{\psi}_1(T)\tilde{\psi}_2(T) = \tilde{\psi}_2(T)\tilde{\psi}_1(T)$, $\tilde{\psi}_i$ is the minimal polynomial of T_i , and the definition of T_i we deduce

$$\tilde{\psi}_1(T)\tilde{\psi}_2(T)\mathbf{v} = \tilde{\psi}_2(T)\tilde{\psi}_1(T)\mathbf{v}_1 + \tilde{\psi}_1(T)\tilde{\psi}_2(T)\mathbf{v}_2 = \mathbf{0}.$$

Hence the monic polynomial $\theta(z) := \tilde{\psi}_1(z)\tilde{\psi}_2(z)$ is an annihilating polynomial of T . Thus $\psi(z)|\theta(z)$ which implies that $\psi(z) = \theta(z)$, hence $\tilde{\psi}_i = \tilde{\psi}$ for $i = 1, 2$.

It is left to show that \mathbf{V}_1 and \mathbf{V}_2 are unique. Let $\bar{\mathbf{V}}_i := \{\mathbf{v} \in \mathbf{V} : \psi_i(T)\mathbf{v} = \mathbf{0}\}$ for $i = 1, 2$. So $\bar{\mathbf{V}}_i$ is a subspace that contains \mathbf{V}_i for $i = 1, 2$. If $\psi_i(T)\mathbf{v} = \mathbf{0}$ then

$$\psi_i(T)(T\mathbf{v}) = (\psi_i(T)T)\mathbf{v} = (T\psi_i(T))\mathbf{v} = T(\psi_i(T)\mathbf{v}) = T\mathbf{0} = \mathbf{0}.$$

Hence $\bar{\mathbf{V}}_i$ is T -invariant subspace. We claim that $\bar{\mathbf{V}}_i = \mathbf{V}_i$. Suppose to the contrary that $\dim \bar{\mathbf{V}}_i > \dim \mathbf{V}_i$ for some $i \in \{1, 2\}$. Let $j \in \{1, 2\}$ and $j \neq i$. Then

$\dim(\bar{\mathbf{V}}_i \cap \mathbf{V}_j) \geq 0$. As before we conclude that $\mathbf{U} := \bar{\mathbf{V}}_i \cap \mathbf{V}_j$ is T -invariant subspace. As above, the minimal polynomial of $T|_{\mathbf{U}}$ must divide $\psi_1(z)$ and $\psi_2(z)$, which contradicts the assumption that $(\psi_1, \psi_2) = 1$. This concludes the proof of the theorem for $k = 2$.

Assume that $k \geq 3$. Let $\hat{\psi}_2 := \psi_2 \dots \psi_k$. Then $(\psi_1, \hat{\psi}_2) = 1$ and $\psi = \psi_1 \hat{\psi}_2$. Then $\mathbf{V} = \mathbf{V}_1 \oplus \hat{\mathbf{V}}_2$, where $T : \mathbf{V}_1 \rightarrow \mathbf{V}_1$, has the minimal polynomial ψ_1 , and $T : \hat{\mathbf{V}}_2 \rightarrow \hat{\mathbf{V}}_2$ has the minimal polynomial $\hat{\psi}_2$. Note that \mathbf{V}_1 and $\hat{\mathbf{V}}_2$ are unique. Apply the induction hypothesis to $T|_{\hat{\mathbf{V}}_2}$ to deduce the theorem. \square

Problems

1. Let $A, B \in \mathbb{F}^{n \times n}$ and $p(z) \in \mathbb{F}[z]$. Show
 - (a) If $B = UAU^{-1}$, for some $U \in \text{GL}(n, \mathbb{F})$, then $p(B) = Up(A)U^{-1}$.
 - (b) If $A \sim B$ then A and B have the same minimal polynomial.
 - (c) Let $A\mathbf{x} = \lambda\mathbf{x}$. Then $p(A)\mathbf{x} = p(\lambda)\mathbf{x}$. Deduce that each eigenvalue of A is a root of the minimal polynomial of A .
 - (d) Assume that A has n distinct eigenvalues. Then A is nonderogatory.
2. (a) Show that the Jordan block $J_k(\lambda) \in \mathbb{F}^{k \times k}$ is nonderogatory.
 (b) Let $\lambda_1, \dots, \lambda_k \in \mathbb{F}$ be k distinct elements. Let

$$A = \bigoplus_{i=1}^k J_{m_i}^{l_i}(\lambda_i), \text{ where } m_i = m_{i1} \geq \dots \geq m_{il_i} \geq 1, \text{ for } i = 1, \dots, k. \quad (3.6)$$

Here m_{ij} and l_i are positive integers be integers. Find the minimal polynomial of A . When A is nonderogatory?

3. Find the characteristic and the minimal polynomials of

$$C := \begin{bmatrix} 2 & 2 & -2 & 4 \\ -4 & -3 & 4 & -6 \\ 1 & 1 & -1 & 2 \\ 2 & 2 & -2 & 4 \end{bmatrix},$$

4. Let $A := \begin{bmatrix} x & y \\ u & v \end{bmatrix}$. Then A is a point in four dimensional space \mathbb{R}^4 .
 - (a) What is the condition that A has a multiple eigenvalue ($\det(zI_2 - A) = (z - \lambda)^2$) ? Conclude that the set (variety) all 2×2 matrices with a multiple eigenvalue is a quadratic hypersurface in \mathbb{R}^4 , i.e. it satisfies a polynomial equation in (x, y, u, v) of degree 2. Hence its dimension is 3.
 - (b) What is the condition that A has a multiple eigenvalue and it is a diagonalizable matrix, i.e. similar to a diagonal matrix? Show that this is a line in \mathbb{R}^4 . Hence its dimension is 1.
 - (c) Conclude that the set (variety) of 2×2 matrices which have multiple eigenvalues and diagonalizable is "much smaller" than the variety of matrices with multiple eigenvalue.

This fact holds for any $n \times n$ matrices $\mathbb{R}^{n \times n}$ or $\mathbb{C}^{n \times n}$.

5. Programming Problem

Spectrum and pseudo spectrum: Let $A = (a_{ij})_{i,j=1}^n \in \mathbb{C}^{n \times n}$. Then $\det(zI_n - A) = (z - \lambda_1) \dots (z - \lambda_n)$ and the *spectrum* of A is given as $\text{spec } A := \{\lambda_1, \dots, \lambda_n\}$. In computations, the entries of A are known or given up to a certain precision. Say, in regular precision each a_{ij} is known with precision to eight digits: $a_1.a_2 \dots a_8 \times 10^m$ for some integer m , e.g. $1.2345678 \times 10^{-12}$, in floating point notation. Thus, with a given matrix A , we associate a whole class of matrices $\mathcal{C}(A) \subset \mathbb{C}^{n \times n}$ of matrices $B \in \mathbb{C}^{n \times n}$ that are represented by A . For each $B \in \mathcal{C}(A)$ we have the spectrum $\text{spec } B$. Then the *pseudo spectrum* of A is the union of all the spectra of $B \in \mathcal{C}(A)$: $\text{pspec } A := \cup_{B \in \mathcal{C}(A)} \text{spec } (B)$. $\text{spec } A$ and $\text{pspec } A$ are subsets of the complex plane \mathbb{C} and can be easily plotted by computer. The shape of $\text{pspec } A$ gives an idea of our real knowledge of the spectrum of A , and to changes of the spectrum of A under perturbations. The purpose of this programming problems to give the student a taste of this subject.

In all the computations use double precision.

- (a) Choose at random $A = (a_{ij}) \in \mathbb{R}^{5 \times 5}$ as follows: each entry a_{ij} is chosen at random from the interval $[-1, 1]$, using uniform distribution. Find the spectrum of A and plot the eigenvalues of A on the $X - Y$ axis as complex numbers, marked say as $+$, where the center of $+$ is at each eigenvalue.
 - i. For each $\epsilon = 0.1, 0.01, 0.0001, 0.000001$ do the following:
For $i = 1, \dots, 100$ choose $B_i \in \mathbb{R}^{5 \times 5}$ at random as A in the item (a) and find the spectrum of $A + \epsilon B_i$. Plot these spectra, each eigenvalue of $A + \epsilon B_i$ plotted as \cdot on the $X - Y$ axis, together with the plot of the spectrum of A . (Altogether you will have 4 graphs.)
- (b) Let $A := \text{diag}(0.1C, [-0.5])$, i.e. $A \in \mathbb{R}^{5 \times 5}$ be a block diagonal matrix where the first 4×4 block is $0.1C$, where the matrix C is given in Problem 3 above, and the second block is 1×1 matrix with the entry -0.5 . Repeat part (i) of part (a) above with this specific A . (Again you will have 4 graphs.)
- (c) Repeat (a) by choosing at random a symmetric matrix $A = (a_{ij}) \in \mathbb{R}^{5 \times 5}$. That is choose at random a_{ij} for $1 \leq i \leq j$, and let $a_{ji} = a_{ij}$ for $i < j$.
 - i. Repeat the part (i) of (a). (B_j are not symmetric!) You will have 4 graphs.
 - ii. Repeat part (i) of (a), with the restriction that each B_j is a random symmetric matrix, as explained in (c). You will have 4 graphs.
- (d) Can you draw some conclusions about these numerical experiments?

3.4 Existence and uniqueness of the Jordan canonical form

Definition 3.24 $A \in \mathbb{F}^{n \times n}$ or a linear transformation $T : \mathbf{V} \rightarrow \mathbf{V}$ is called *nilpotent* respectively, if $A^m = \mathbf{0}$ or $T^m = \mathbf{0}$. The minimal $m \geq 1$ for which $A^m = \mathbf{0}$ or $T^m = \mathbf{0}$ is called the *index of nilpotency* of A and T respectively, and denoted by $\text{index } A$ or $\text{index } T$ respectively.

Assume that A or T are nilpotent, then the s -numbers are defined as

$$\begin{aligned} s_i(A) &:= \text{rank } A^{i-1} - 2\text{rank } A^i + \text{rank } A^{i+1}, \\ s_i(T) &:= \text{rank } T^{i-1} - 2\text{rank } T^i + \text{rank } T^{i+1}, \quad i = 1, \dots \end{aligned} \quad (3.7)$$

Note that A or T are nilpotent with the index of nilpotency m if and only if z^m is the minimal polynomial of A or T respectively. Furthermore if A or T are nilpotent then the maximal l for which $s_l > 0$ is equal to the index of nilpotency of A or T respectively.

Proposition 3.25 *Let $T : \mathbf{V} \rightarrow \mathbf{V}$ be a nilpotent operator, with the index of nilpotency m , on the finite dimensional vector \mathbf{V} . Then*

$$\text{rank } T^i = \sum_{j=i+1}^m (j-i)s_j = (m-i)s_m + (m-i-1)s_{m-1} + \dots + s_{i+1}, \quad i = 0, \dots, m-1. \quad (3.8)$$

Proof. Since $T^l = \mathbf{0}$ for $l \geq m$ it follows that $s_m(T) = \text{rank } T^{m-1}$ and $s_{m-1} = \text{rank } T^{m-2} - 2\text{rank } T^{m-1}$ if $m > 1$. This proves (3.8) for $i = m-1, m-2$. For other values of i (3.8) follows straightforward from (3.7) by induction on $m-i \geq 2$. \square

Theorem 3.26 *Let $T : \mathbf{V} \rightarrow \mathbf{V}$ be a linear transformation on a finite dimensional space. Assume that T is nilpotent with the index of nilpotency m . Then \mathbf{V} has a basis of the form*

$$\mathbf{x}_j, T\mathbf{x}_j, \dots, T^{l_j-1}\mathbf{x}_j, \quad j = 1, \dots, i, \quad \text{where } l_1 = m \geq \dots \geq l_i \geq 1, \quad \text{and } T^{l_j}\mathbf{x}_j = \mathbf{0}, \quad j = 1, \dots, i. \quad (3.9)$$

More precisely, the number of \mathbf{x}_j , which are equal to an integer $l \in [1, m]$, is equal to $s_l(T)$ given in (3.7).

Proof. Let $s_i := s_i(T), i = 1, \dots, m$ be given by (3.7). Since $T^l = \mathbf{0}$ for $l \geq m$ it follows that $s_m = \text{rank } T^{m-1} = \dim \text{range } T^{m-1}$. So $[\mathbf{y}_1, \dots, \mathbf{y}_{s_m}]$ is a basis for $T^{m-1}\mathbf{V}$. Clearly $\mathbf{y}_i = T^{m-1}\mathbf{x}_i$ for some $\mathbf{x}_1, \dots, \mathbf{x}_{s_m} \in \mathbf{V}$. We claim that the ms_m vectors

$$\mathbf{x}_1, T\mathbf{x}_1, \dots, T^{m-1}\mathbf{x}_1, \dots, \mathbf{x}_{s_m}, T\mathbf{x}_{s_m}, \dots, T^{m-1}\mathbf{x}_{s_m} \quad (3.10)$$

are linearly independent. Suppose that there exists a linear combination of these vectors that is equal to $\mathbf{0}$:

$$\sum_{j=0}^{m-1} \sum_{k=1}^{s_m} \alpha_{jk} T^j \mathbf{x}_k = \mathbf{0}. \quad (3.11)$$

Multiply this equality by T^{m-1} . Thus we obtain $\sum_{j=0}^{m-1} \sum_{k=1}^{s_m} \alpha_{jk} T^{m-1+j} \mathbf{x}_k = \mathbf{0}$. Recall that $T^l = \mathbf{0}$ for any $l \geq m$. Hence this equality reduces to $\sum_{k=1}^{s_m} \alpha_{0k} T^{m-1} \mathbf{x}_k = \mathbf{0}$. Since $T^{m-1}\mathbf{x}_1, \dots, T^{m-1}\mathbf{x}_{s_m}$ form a basis in $T^{m-1}\mathbf{V}$ it follows that $\alpha_{0k} = 0$ for $k = 1, \dots, s_m$. If $m = 1$ we deduce that the vectors in (3.10) are linearly independent. Assume that $m > 1$. Suppose that we already proved that $\alpha_{jk} = 0$ for $k = 1, \dots, s_m$ and $j = 0, \dots, l-1$, where $1 \leq l \leq m-1$. Hence in (3.11) we

can assume that the summation on j starts from $j = l$. Multiply (3.11) by T^{m-l+1} and use the above arguments to deduce that $\alpha_{lk} = 0$ for $k = 1, \dots, s_m$. Use this argument iteratively for $l = 1, \dots, m - 1$ to deduce the linear independence of the vectors in (3.10).

Note that for $m = 1$ we proved the theorem. Assume that $m > 1$. Let $p \in [1, m]$ be an integer. We claim that the vectors

$$x_j, T\mathbf{x}_j, \dots, T^{l_j-1}\mathbf{x}_j, \quad \text{for all } j \text{ such that } l_j \geq p \quad (3.12)$$

are linearly independent and satisfy the condition $T^{l_j}\mathbf{x}_j = \mathbf{0}$ for all $l_j \geq p$. Moreover, the vectors

$$T^{p-1}\mathbf{x}_j, \dots, T^{l_j-1}\mathbf{x}_j, \quad \text{for all } j \text{ such that } l_j \geq p \quad (3.13)$$

is a basis for range T^{p-1} . Furthermore for each integer $l \in [p, m]$ the number of l_j , which are equal to l , is equal to $s_l(T)$.

We prove this claim by the induction on $m - p + 1$. For $p = m$ our previous argument give this claim. Assume that the claim holds for $p = q \geq m$ and let $p = q - 1$. By the induction assumption the vectors in (3.12) are linearly independent for $l_j \geq q$. Hence that vectors $T^{q-2}\mathbf{x}_j, \dots, T^{l_j-1}\mathbf{x}_j$ for all $l_j \geq q$ are linearly independent. Use the induction assumption that the number of $l_j = l \in [q, m]$ is equal to $s_l(T)$ to deduce that the number of this vectors is equal to $t_{q-2} := (m - q + 2)s_m + (m - q + 1)s_{m-1} + \dots + 2s_q$. Also the number of $l_j \geq q$ is $L_q = s_m + s_{m-1} + \dots + s_q$. Use the formula for rank T^{q-2} in (3.8) to deduce that $\text{rank } T^{q-2} - t_{q-2} = s_{q-1}$.

Suppose first that $s_{q-1} = 0$. Hence the vectors $T^{q-2}\mathbf{x}_j, \dots, T^{l_j-1}\mathbf{x}_j$ for all $l_j \geq q$ form a basis in range T^{q-2} . In this case we assume that there is no l_j that is equal to $q - 1$. This concludes the proof of the induction step and the proof of the theorem in this case.

Assume now that $s_{q-1} > 0$. Then there exist vectors $\mathbf{z}_1, \dots, \mathbf{z}_{s_{q-1}}$ that together with the vectors $T^{q-2}\mathbf{x}_j, \dots, T^{l_j-1}\mathbf{x}_j$ for all $l_j \geq q$ form a basis in $T^{q-2}\mathbf{V}$. Let $\mathbf{z}_k = T^{q-2}\mathbf{u}_k, k = 1, \dots, s_{q-1}$. Observe next that by induction hypothesis the vectors given in (3.13) form a basis in range T^{p-1} for $p = q$. Hence $T^{q-1}\mathbf{u}_k = \sum_{j:l_j \geq q} \sum_{r=q-1}^{l_j-1} \beta_{k,r,j} T^r \mathbf{x}_j$. Let $\mathbf{v}_k := \mathbf{u}_k - \sum_{j:l_j \geq q} \sum_{r=q-1}^{l_j-1} \beta_{k,r,j} T^{r-q+1} \mathbf{x}_j$. Clearly $T^{q-1}\mathbf{v}_k = 0$ for $k = 1, \dots, s_{q-1}$. Also $T^{q-2}\mathbf{v}_k = \mathbf{z}_k - \sum_{j:l_j \geq q} \sum_{r=q-1}^{l_j-1} \beta_{k,r,j} T^{r-1} \mathbf{x}_j$. Hence $T^{q-2}\mathbf{v}_1, \dots, T^{q-2}\mathbf{v}_{s_{q-1}}$ and the vectors $T^{q-2}\mathbf{x}_j, \dots, T^{l_j-1}\mathbf{x}_j$ for all $l_j \geq q$ form a basis in $T^{q-2}\mathbf{V}$. From the above definition of L_q $l_j \geq q$ if and only if $j = [1, L_q]$. Let $\mathbf{x}_j = \mathbf{v}_{j-L_q}$ and $l_j = s_{q-1}$ for $j = L_q + 1, \dots, L_{q-1} := L_q + s_{q-1}$.

It is left to show that the vectors given in (3.12) are linearly independent for $p = q - 1$. This is done as in the beginning of the proof of the theorem. (Assume that a linear combination of these vectors is equal to $\mathbf{0}$. Then apply T^{q-2} and use the fact that $T^{l_j}\mathbf{x}_j = \mathbf{0}$ for $j = 1, \dots, L_{q-1}$. Then continue as in the beginning of the proof of this theorem.) This concludes the proof of this theorem by induction. \square

Corollary 3.27 *Let T satisfies the assumption of Theorem 3.26 hold. Denote $\mathbf{V}_j := \text{span}(T^{l_j-1}\mathbf{x}_j, \dots, T\mathbf{x}_j, \mathbf{x}_j)$ for $j = 1, \dots, i$. Then each \mathbf{V}_j is a T -invariant subspace, $T|_{\mathbf{V}_j}$ is represented by $J_{l_j}(0) \in \mathbb{C}^{l_j \times l_j}$ in the basis $[T^{l_j-1}\mathbf{x}_j, \dots, T\mathbf{x}_j, \mathbf{x}_j]$, and $\mathbf{V} = \bigoplus_{j=1}^i \mathbf{V}_j$. Each l_j is uniquely determined by the sequence $s_i(T), i = 1, \dots, m$. Namely, the index m of the nilpotent T is the largest $i \geq 1$ such that $s_i(T) \geq 1$. Let*

$k_1 = s_m(T), l_1 = \dots = l_{k_1} = p_1 = m$ and define recursively $k_r := k_{r-1} + s_{p_r}(T)$, $l_{k_{r-1}+1} = \dots = l_{k_r} = p_r$, where $2 \leq r, p_r \in [1, m-1], s_{p_r}(T) > 0$ and $k_{r-1} = \sum_{j=1}^{m-p_r} s_{m-j+1}(T)$.

Definition 3.28 $T : \mathbf{V} \rightarrow \mathbf{V}$ be a nilpotent operator. Then the sequence (l_1, \dots, l_i) defined in Theorem 3.26, which gives the lengths of the corresponding Jordan blocks of T in a decreasing order, is called the Segré characteristic of T . The Weyr characteristic of T is the dual to Segre's characteristic. That is consider an $m \times i$ 0-1 matrix $B = (b_{pq}) \in \{0, 1\}^{m \times i}$. The j -th column of B has 1 in the rows $1, \dots, l_j$ and 0 in the rest of the rows. Let ω_p be the p -th row sum of B for $p = 1, \dots, m$. Then $\omega_1 \geq \dots \geq \omega_m \geq 1$ is the Weyr characteristic.

Proof of Theorem 3.13 (The Jordan Canonical Form)

Let $p(z) = \det(zI_n - A)$ be the characteristic polynomial of $A \in \mathbb{C}^{n \times n}$. Since \mathbb{C} is algebraically closed $p(z) = \prod_{j=1}^k (z - \lambda_j)^{n_j}$. Here $\lambda_1, \dots, \lambda_k$ are k distinct roots, (eigenvalues of A), where $n_j \geq 1$ is the multiplicity of λ_j in $p(z)$. Note that $\sum_{j=1}^k n_j = n$. Let $\psi(z)$ be the minimal polynomial of A . By Theorem 3.18 $\psi(z)|p(z)$. Problem 1(c) of §3.3 we deduce that $\psi(\lambda_j) = 0$ for $j = 1, \dots, k$. Hence

$$\det(zI_n - A) = \prod_{j=1}^k (z - \lambda_j)^{n_j}, \quad \psi(z) = \prod_{j=1}^k (z - \lambda_j)^{m_j}, \quad (3.14)$$

$$1 \leq m_j \leq n_j, \quad \lambda_j \neq \lambda_i \text{ for } j \neq i, \quad i, j = 1, \dots, k.$$

Let $\psi_j := (z - \lambda_j)^{m_j}$ for $j = 1, \dots, k$. Then $(\psi_j, \psi_i) = 1$ for $j \neq i$. Let $\mathbf{V} := \mathbb{C}^n$ and $T : \mathbf{V} \rightarrow \mathbf{V}$ be given by $T\mathbf{x} := A\mathbf{x}$ for any $\mathbf{x} \in \mathbb{C}^n$. Then $\det(zI_n - A)$ and $\psi(z)$ are the characteristic and the minimal polynomial of T respectively. Use Theorem 3.23 to obtain the decomposition $\mathbf{V} = \bigoplus_{i=1}^k \mathbf{V}_i$, where each \mathbf{V}_i is a nontrivial T -invariant subspace such that the minimal polynomial of $T_i := T|_{\mathbf{V}_i}$ is ψ_i for $i = 1, \dots, k$. That is $T_i - \lambda_i I_i$, where I_i is the identity operator, i.e. $I_i \mathbf{v} = \mathbf{v}$ for all $\mathbf{v} \in \mathbf{V}_i$, is a nilpotent operator on \mathbf{V}_i and $\text{index}(T_i - \lambda_i I_i) = m_i$. Let $Q_i := T_i - \lambda_i I_i$. Then Q_i is nilpotent and $\text{index } Q_i = m_i$. Apply Theorem 3.26 and Corollary 3.27 to deduce that $\mathbf{V}_i = \bigoplus_{j=1}^{q_j} \mathbf{V}_{i,j}$, where each $\mathbf{V}_{i,j}$ is Q_i -invariant subspace, and each $\mathbf{V}_{i,j}$ has a basis in which Q_i is represented by a Jordan block $J_{m_{i,j}}(0)$ for $j = 1, \dots, q_j$. According to Corollary 3.27

$$m_i = m_{i1} \geq \dots \geq m_{iq_i} \geq 1, \quad i = 1, \dots, k. \quad (3.15)$$

Furthermore, the above sequence is completely determined by $\text{rank } Q_i^j, j = 0, 1, \dots$ for $i = 1, \dots, k$. Noting that $T_i = Q_i + \lambda_i I_i$ it easily follows that each $\mathbf{V}_{i,j}$ is a T_i -invariant subspace, hence T -invariant subspace. Moreover, in the same basis of $\mathbf{V}_{i,j}$ that Q_i is represented by $J_{m_{i,j}}(0)$ T_i is represented by $J_{m_{i,j}}(\lambda_i)$ for $j = 1, \dots, q_i$ and $i = 1, \dots, k$. This shows the existence of the Jordan canonical form.

We now show that the Jordan canonical form is unique, up to a permutation of factors. Note that the minimal polynomial of A is completely determined by its Jordan canonical form. Namely $\psi(z) = \prod_{i=1}^k (z - \lambda_i)^{m_{i1}}$, where m_{i1} is the biggest Jordan block with the eigenvalues λ_i . (See Problems 1,2 in §3.3.) Thus $m_{i1} = m_i$ for $i = 1, \dots, k$. Theorem 3.23 yields that the subspaces $\mathbf{V}_1, \dots, \mathbf{V}_k$ are uniquely determined by ψ . So each T_i and $Q_i = T - \lambda_i I_i$ are uniquely determined. Theorem

3.26 yields that $\text{rank } Q_i^j, j = 0, 1, \dots$ determines the sizes of the Jordan blocks of Q_i . Hence all the Jordan blocks corresponding to λ_i are uniquely determined for each $i \in [1, k]$. \square

Corollary 3.29 *Let $A \in \mathbb{F}^{n \times n}$ and assume that the characteristic polynomial of $p(z) = \det(zI_n - A)$ splits to linear factors, i.e. (3.14) holds. Let B be the Jordan canonical form of A . Then*

1. *The multiplicity of the eigenvalue λ_i in the minimal polynomial $\psi(z)$ of A is the size of the biggest Jordan block corresponding to λ_i in B .*
2. *The number of Jordan blocks in B corresponding to λ_i is the nullity of $A - \lambda_i I_n$, i.e. the number of Jordan block in B corresponding to $A - \lambda_i I_n$ is the number of linearly independent eigenvectors of A corresponding to the eigenvalue λ_i .*
3. *Let λ_i be an eigenvalue of A . Then the number of the Jordan blocks of order i corresponding to λ_i in B is given in (3.16).*

Proof. 1. Since $J_n(0)^n = 0$ and $J_n(0)^{n-1} \neq 0$, it follows that the minimal polynomial of $J_n(\lambda)$ is $(z - \lambda)^n = \det(zI_n - J_n(\lambda))$. Use Problem 4a to deduce the first part of the corollary. 2. Since $J_n(0)$ has one independent eigenvector, use Problem 4b to deduce the second part of the corollary. 3. Use 3a to establish the last part of the corollary. \square

Problems

1. Let $T : \mathbf{V} \rightarrow \mathbf{V}$ be nilpotent with $m = \text{index } T$. Let $(\omega_1, \dots, \omega_m)$ be the Weyr characteristic. Show that $\text{rank } T^j = \sum_{p=1}^j \omega_p$ for $j = 1, \dots, m$.
2. Let $A \in \mathbb{C}^{n \times n}$. Show that A is diagonalizable if and only if all the zeros of the minimal polynomial ψ of A are simple, i.e. ψ does not have multiple roots.
3. Let $A \in \mathbb{C}^{n \times n}$ and assume that $\det(zI_n - A) = \prod_{i=1}^k (z - \lambda_i)^{n_i}$, where $\lambda_1, \dots, \lambda_k$ are k distinct eigenvalues of A . Let

$$s_i(A, \lambda_j) := \text{rank}(A - \lambda_j I_n)^{i-1} - 2\text{rank}(A - \lambda_j I_n)^i + \text{rank}(A - \lambda_j I_n)^{i+1} \quad (3.16)$$

$$i = 1, \dots, n_j, j = 1, \dots, k.$$

- (a) Show that $s_i(A, \lambda_j)$ is the number of Jordan blocks of order i corresponding to λ_j for $i = 1, \dots, n_j$.
 - (b) Show that in order to find all Jordan blocks of A corresponding to λ_j one can stop computing $s_i(A, \lambda_j)$ at the smallest $i \in [1, n_j]$ such that $1s_1(A, \lambda_j) + 2s_2(A, \lambda_j) \dots + is_i(A, \lambda_j) = n_j$.
4. Let $C = F \oplus G = \text{diag}(F, G)$, $F \in \mathbb{F}^{l \times l}$, $G \in \mathbb{F}^{m \times m}$.
 - (a) Assume that ψ_F, ψ_G are the minimal polynomials of F, G respectively. Show that ψ_C , the minimal polynomial of C , is equal to $\frac{\psi_F \psi_G}{(\psi_F, \psi_G)}$.
 - (b) Show that $\text{nul}(C) = \text{nul}(F) + \text{nul}(G)$. In particular, if G is invertible, i.e. 0 is not eigenvalue of G then $\text{nul}(C) = \text{nul}(F)$.

3.5 Cyclic subspaces

Let $T : \mathbf{V} \rightarrow \mathbf{V}$ be a linear operator, and assume that \mathbf{V} is finite dimensional. Let $\mathbf{0} \neq \mathbf{u} \in \mathbf{V}$. Consider the sequence of vectors $\mathbf{u} = T^0\mathbf{u}, T\mathbf{u}, T^2\mathbf{u}, \dots$. Since $\dim \mathbf{V} < \infty$ there exists a positive integer $l \geq 1$ such that $\mathbf{u}, T\mathbf{u}, \dots, T^{l-1}\mathbf{u}$ linearly independent, and $\mathbf{u}, T\mathbf{u}, \dots, T^l\mathbf{u}$ are linearly dependent. Hence l is smallest integer such that

$$T^l\mathbf{u} = - \sum_{i=1}^l a_i T^{l-i}\mathbf{u}. \quad (3.17)$$

Clearly, $l \leq \dim \mathbf{V}$. The polynomial $\psi_{\mathbf{u}}(z) := z^l + \sum_{i=1}^l a_i z^{l-i}$ is called the *minimal* polynomial of \mathbf{u} , with respect to T . It is a monic polynomial of the minimal degree such that $\phi(T)\mathbf{u} = \mathbf{0}$. Its property is similar to the property of the minimal polynomial of T . Namely if a polynomial $\phi \in \mathbb{F}[z]$ annihilates \mathbf{u} , i.e. $\phi(T)\mathbf{u} = \mathbf{0}$ then $\psi_{\mathbf{u}}|\phi$. In particular, the minimal polynomial $\psi(z)$ of T is divisible by $\psi_{\mathbf{u}}$, since $\psi(T)\mathbf{u} = \mathbf{0}$. Clearly, every vector $\mathbf{w} \in \mathbf{U}$ can be uniquely represented as $\phi(T)\mathbf{u}$, where $\phi \in \mathbb{F}[z]$ and $\deg \phi \leq l-1$. Hence \mathbf{U} is T -invariant subspace. The subspace \mathbf{U} spanned by $\mathbf{u}, T\mathbf{u}, \dots$, is called a *cyclic* subspace, generated by \mathbf{u} . Note that in the basis $\mathbf{u}_1 = \mathbf{u}, \mathbf{u}_2 = T\mathbf{u}, \dots, \mathbf{u}_l = T^{l-1}\mathbf{u}$ the linear transformation $T|_{\mathbf{U}}$ is given by the matrix

$$\begin{bmatrix} 0 & 0 & 0 & \dots & 0 & -a_l \\ 1 & 0 & 0 & \dots & 0 & -a_{l-1} \\ 0 & 1 & 0 & \dots & 0 & -a_{l-2} \\ \vdots & \vdots & \vdots & \dots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 1 & -a_1 \end{bmatrix} \in \mathbb{F}^{l \times l}. \quad (3.18)$$

The above matrix is called the *companion* matrix, corresponding to the polynomial $\psi_{\mathbf{u}}(z)$. (Sometimes the transpose of the above matrix is called the companion matrix.)

Lemma 3.30 *Let $T : \mathbf{V} \rightarrow \mathbf{V}$, $\dim \mathbf{V} < \infty$, $\mathbf{0} \neq \mathbf{u}, \mathbf{w}$. Assume that $(\psi_{\mathbf{u}}, \psi_{\mathbf{w}}) = 1$. Then $\psi_{\mathbf{u}+\mathbf{w}}(z) = \psi_{\mathbf{u}}(z)\psi_{\mathbf{w}}(z)$.*

Proof. Let \mathbf{U}, \mathbf{W} be the cyclic invariant subspaces generated by \mathbf{u}, \mathbf{w} respectively. We claim that the T -invariant subspace $\mathbf{X} := \mathbf{U} \cap \mathbf{W}$ is the trivial subspace $\{\mathbf{0}\}$. Suppose to the contrary, there exists $\mathbf{0} \neq \mathbf{x} \in \mathbf{X}$. Let $\mathbf{X}_1 \subset \mathbf{X}$ be a nontrivial cyclic subspace generated by \mathbf{x} and $\psi_{\mathbf{x}}$ be the minimal polynomial corresponding to \mathbf{x} . Since $\mathbf{X}_1 \subset \mathbf{U}$ it follows that $\mathbf{x} = \phi(T)\mathbf{u}$. Hence $\psi_{\mathbf{u}}(T)\mathbf{x} = \mathbf{0}$. Thus $\psi_{\mathbf{x}}|\psi_{\mathbf{u}}$. Similarly $\psi_{\mathbf{x}}|\psi_{\mathbf{w}}$. This contradicts the assumption that $(\psi_{\mathbf{u}}, \psi_{\mathbf{w}}) = 1$. Hence $\mathbf{U} \cap \mathbf{W} = \{\mathbf{0}\}$. Let $\phi = \psi_{\mathbf{u}}\psi_{\mathbf{w}}$. Clearly

$$\phi(T)(\mathbf{u} + \mathbf{w}) = \psi_{\mathbf{w}}(T)(\psi_{\mathbf{u}}(T)\mathbf{u}) + \psi_{\mathbf{u}}(T)(\psi_{\mathbf{w}}(T)\mathbf{w}) = \mathbf{0} + \mathbf{0} = \mathbf{0}.$$

So ϕ is an annihilating polynomial of $\mathbf{u} + \mathbf{w}$. Let $\theta(z)$ be an annihilating polynomial of $\mathbf{u} + \mathbf{w}$. So $\mathbf{0} = (\theta(T)\mathbf{u}) + (\theta(T)\mathbf{w})$. Since \mathbf{U}, \mathbf{W} are T -invariant subspaces and $\mathbf{U} \cap \mathbf{W} = \{\mathbf{0}\}$, it follows that $\mathbf{0} = \theta(T)\mathbf{u} = \theta(T)\mathbf{w}$. Hence $\psi_{\mathbf{u}}|\theta, \psi_{\mathbf{w}}|\theta$. As $(\psi_{\mathbf{u}}, \psi_{\mathbf{w}}) = 1$ it follows that $\phi|\theta$. Hence $\psi_{\mathbf{u}+\mathbf{w}} = \phi$. \square

Theorem 3.31 *Let $T : \mathbf{V} \rightarrow \mathbf{V}$ be a linear operator and $1 \leq \dim V < \infty$. Let $\psi(z)$ be the minimal polynomial of T . Then there exists $\mathbf{0} \neq \mathbf{u} \in \mathbf{V}$ such that $\psi_{\mathbf{u}} = \psi$.*

Proof. Assume first that $\psi(z) = (\phi(z))^l$, where $\phi(z)$ is irreducible in $\mathbb{F}[z]$, i.e. $\phi(z)$ is not divisible by any polynomial θ such $1 \leq \deg \theta < \deg \phi$, and l is a positive integer. Let $\mathbf{0} \neq \mathbf{w} \in \mathbf{V}$. Recall that $\psi_{\mathbf{w}} | \psi$. Hence $\psi_{\mathbf{w}} = (\phi)^{l(\mathbf{w})}$, where $l(\mathbf{w}) \in [1, l]$ is a positive integer. Assume to the contrary that $1 \leq l(\mathbf{w}) \leq l - 1$ for each $\mathbf{0} \neq \mathbf{w} \in \mathbf{V}$. Then $(\phi(T))^{l-1}(\mathbf{w}) = \mathbf{0}$. As $(\phi(T))^{l-1}(\mathbf{0}) = \mathbf{0}$ we deduce that $(\phi(z))^{l-1}$ is annihilating polynomial of T . So $\psi | \phi^{l-1}$ which is impossible. Hence there exists $\mathbf{u} \neq \mathbf{0}$ such that $l(\mathbf{u}) = l$, i.e. $\psi_{\mathbf{u}} = \psi$.

Assume now the general case

$$\psi(z) = \prod_{i=1}^k (\phi_i(z))^{l_i}, \quad l_i \in \mathbb{N}, \quad \phi_i \text{ irreducible}, \quad (\phi_i, \phi_j) = 1 \text{ for } i \neq j, \quad (3.19)$$

where $k > 1$. Theorem 3.23 implies that $V = \bigoplus_{i=1}^k \mathbf{V}_i$, where $T : \mathbf{V}_i \rightarrow \mathbf{V}_i$ and $\phi_i^{l_i}$ is the minimal polynomial of $T|_{\mathbf{V}_i}$. The first case of the proof yields the existence $\mathbf{0} \neq \mathbf{u}_i \in \mathbf{V}_i$ such that $\psi_{\mathbf{u}_i} = \phi_i^{l_i}$ for $i = 1, \dots, k$. Let $\mathbf{w}_j = \mathbf{u}_1 + \dots + \mathbf{u}_j$. Use Lemma 3.30 to deduce that $\psi_{\mathbf{w}_j} = \prod_{i=1}^j \phi_i^{l_i}$. Hence $\psi_{\mathbf{u}} = \psi$ for $\mathbf{u} := \mathbf{w}_k$. \square

Let $\mathbf{U} \subset \mathbf{V}$ be a subspace of \mathbf{V} . Then \mathbf{V}/\mathbf{U} is the following object. We introduce the following relation \sim on the elements of V : $\mathbf{x} \sim \mathbf{y}$, which is also denoted as $\mathbf{x} \equiv \mathbf{y} \pmod{\mathbf{U}}$, if $\mathbf{x} - \mathbf{y} \in \mathbf{U}$. It is straightforward to show that \sim is an equivalence relation. Denote by $[\mathbf{x}]$ the equivalence class corresponding to \mathbf{x} . Then $\mathbf{V}/\mathbf{U} := \{[\mathbf{x}], \mathbf{x} \in \mathbf{V}\}$. It is straightforward to show that the following definitions are valid $a[\mathbf{x}] := [a\mathbf{x}]$, $[\mathbf{x}] + [\mathbf{y}] := [\mathbf{x} + \mathbf{y}]$, i.e. these definitions are independent of the choice of the representatives of the equivalence classes. Furthermore, under these two operations \mathbf{V}/\mathbf{U} is a vector space. Then V/\mathbf{U} is called the quotient vector space. Moreover, $\dim \mathbf{V}/\mathbf{U} = \dim \mathbf{V} - \dim \mathbf{U}$. Observe next that if T is a linear operator $T : \mathbf{V} \rightarrow \mathbf{V}$ and \mathbf{U} is an T -invariant subspace of \mathbf{V} , then T induces a linear transformation $\hat{T} : \mathbf{V}/\mathbf{U} \rightarrow \mathbf{V}/\mathbf{U}$. See Problem 4a.

Theorem 3.32 *Let $T : \mathbf{V} \rightarrow \mathbf{V}$ be a linear operator on a finite dimensional space $\dim V > 0$. Then there exists a decomposition of V to a direct sum of r T -cyclic subspaces $V = \bigoplus_{i=1}^r \mathbf{U}_i$ with the following properties. Assume that ψ_i is the minimal polynomial of $T|_{\mathbf{U}_i}$ then ψ_1 is the minimal polynomial of T . Furthermore, $\psi_{i+1} | \psi_i$ for $i = 1, \dots, r - 1$.*

Proof. We prove the theorem by induction on $\dim \mathbf{V}$. For $\dim V = 1$, any $\mathbf{0} \neq \mathbf{u}$ is an eigenvector of T : $T\mathbf{u} = \lambda\mathbf{u}$, so \mathbf{V} is cyclic, and $\psi(z) = \psi_1(z) = z - \lambda$. Assume now that the theorem holds of all \mathbf{V} with $\dim V \leq n$. Let $\dim \mathbf{V} = n + 1$. Let $\psi(z)$ be the minimal polynomial of T . Assume that $m = \deg \psi$. Suppose first that T is nonderogatory, i.e. $m = n + 1$, which is the degree of the characteristic polynomial of T . Theorem 3.31 implies that \mathbf{V} is a cyclic subspace. So $r = 1$, and the theorem holds in this case.

Assume now that T is derogatory. So $m < n + 1$. Theorem 3.31 implies the existence of $\mathbf{0} \neq \mathbf{u}_1$ such that $\psi_{\mathbf{u}_1} = \psi$. Let \mathbf{U}_1 be the cyclic subspace generated by \mathbf{u}_1 . Let $\hat{V} := \mathbf{V}/\mathbf{U}_1$. So $1 \leq \dim \hat{V} = \dim \mathbf{V} - m \leq n$. We now apply the induction

hypothesis to $\hat{T} : \hat{\mathbf{V}} \rightarrow \hat{\mathbf{V}}$. So $\hat{V} = \bigoplus_{i=2}^r \hat{U}_i$, where \hat{U}_i is a cyclic subspace generated by $[\hat{\mathbf{u}}_i], i = 2, \dots, r$. The minimal polynomial of $[\hat{\mathbf{u}}_i]$ is $\psi_i, i = 2, \dots, r$. ψ_2 is the minimal polynomial of \hat{T} , and $\psi_{i+1}|\psi_i$ for $i = 2, \dots, r-1$. Observe first that for any polynomial $\phi(z)$ we have the identity $\phi(\hat{T})[\mathbf{x}] = [\phi(T)\mathbf{x}]$. Since $\psi(T)\mathbf{x} = \mathbf{0}$ it follows that $\psi_1 := \psi(z)$ is an annihilating polynomial of \hat{T} . Hence $\psi_2|\psi_1$.

Observe next that since $\psi_{i+1}|\psi_i$ for $i = 1, \dots, r-1$, it follows that $\psi_i|\psi_1 = \psi$. So $\psi = \theta_i\psi_i$, where θ_i is a monic polynomial and $i = 2, \dots, r$. Since $[\mathbf{0}] = \psi_i(\hat{T})[\hat{\mathbf{u}}_i] = [\psi_i(T)\hat{\mathbf{u}}_i]$ it follows that $\psi_i(T)\hat{\mathbf{u}}_i \in \mathbf{U}_1$. Hence $\psi_i(T)\hat{\mathbf{u}}_i = \omega_i(T)\mathbf{u}_1$ for some $\omega_i(z) \in \mathbb{F}[z]$. Hence $\mathbf{0} = \psi(T)\hat{\mathbf{u}}_i = \theta_i(T)\psi_i(T)\hat{\mathbf{u}}_i = \theta_i(T)\omega_i(T)\mathbf{u}_1$. Therefore, $\psi_{\mathbf{u}_1} = \psi$ divides $\theta_i\omega_i$. So $\psi_i|\omega_i$. Thus $\omega_i = \psi_i\alpha_i$. Define $\mathbf{u}_i = \hat{\mathbf{u}}_i - \alpha_i(T)\mathbf{u}_1$ for $i = 2, \dots, r$. The rest of the proof theorem follows from Problem 5. \square

Theorem 3.33 *Let $A \in \mathbb{F}^{n \times n}$. Then there exists r monic polynomials ψ_1, \dots, ψ_r of degree one at least, such that the following conditions. First $\psi_{i+1}|\psi_i$ for $i = 1, \dots, r$. Second, $\psi_1 = \psi$ is the minimal polynomial of A . Then A is similar to $\bigoplus_{i=1}^r C(\psi_i)$, where $C(\psi_i)$ is the companion matrix of the form (3.18) corresponding to the polynomial ψ_i .*

Decompose each $\psi_i = \psi_{i,1} \dots \psi_{i,t_i}$ to the product of its irreducible components, as in the decomposition of $\psi = \psi_1$ given in (3.19). Then A is similar to $\bigoplus_{i=1}^r \bigoplus_{l=1}^{t_i} C(\psi_{i,l})$.

Suppose finally that $\psi(z) = \prod_{j=1}^k (z - \lambda_j)^{m_j}$ where $\lambda_1, \dots, \lambda_k$ are the k distinct eigenvalues of A . Then A is similar to the Jordan canonical form given in Theorem 3.13.

Proof. Identify A with $T : \mathbb{F}^n \rightarrow \mathbb{F}^n$, where $T(\mathbf{x}) = A\mathbf{x}$. Use Theorem 3.32 to decompose $\mathbb{F}^n = \bigoplus_{i=1}^r \mathbf{U}_i$, where each \mathbf{U}_i is a cyclic invariant subspaces such that $T|\mathbf{U}_i$ has the minimal polynomial ψ_i . Since \mathbf{U}_i is cyclic then $T|\mathbf{U}_i$ is represented by $C(\psi_i)$ in the suitable basis of \mathbf{U}_i as shown in the beginning of this section. Hence A is similar to $\bigoplus_{i=1}^r C(\psi_i)$.

Consider next $T_i := T\mathbf{U}_i$. Use Theorem 3.23 to deduce that \mathbf{U}_i decomposes to a direct sum of T_i invariant subspaces $\bigoplus_{l=1}^{t_i} \mathbf{U}_{i,l}$, where the minimal polynomial of $T_{i,l} := T_i|\mathbf{U}_{i,l}$ is $\psi_{i,l}$. Since \mathbf{U}_i was cyclic, i.e., T_i was nonderogatory, it follows that each $T_{i,l}$ must be nonderogatory, i.e. $\mathbf{U}_{i,l}$ cyclic. (See Problem 3.) Recall that each $T_{i,l}$ is represented in a corresponding basis by $C(\psi_{i,l})$. Hence A is similar to $\bigoplus_{i=1}^r \bigoplus_{l=1}^{t_i} C(\psi_{i,l})$.

Suppose finally that $\psi(z)$ splits to linear factors. Hence $\psi_{i,l} = (z - \lambda_l)^{m_{i,l}}$. Hence $T_{i,l} - \lambda_l I$ is nilpotent of index $m_{i,l}$. So there exists a basis in $\mathbf{U}_{i,l}$ such that $T_{i,l}$ is represented by the Jordan block $J_{m_{i,l}}(\lambda_l)$. Therefore A is similar to a sum of corresponding Jordan blocks. \square

The polynomials ψ_1, \dots, ψ_k appearing in Theorem 3.32 are called *invariant polynomials* of T or its representation matrix A , in any basis. The polynomials $\psi_{i,1}, \dots, \psi_{i,t_i}, i = 1, \dots, k$ appearing in Theorem 3.33 are called the *elementary divisors* of A , or the corresponding linear transformation T represented by A . We now show that these polynomial are uniquely defined.

Lemma 3.34 *(Cauchy-Binet formula) For two positive integers $1 \leq p \leq m$ denote by $Q_{p,m}$ the set of all subsets α of $\{1, \dots, m\}$ of cardinality p . (See bottom of page 5 of these notes.) Let $A \in \mathbb{F}^{m \times n}, B \in \mathbb{F}^{n \times l}$ and denote $C = AB \in \mathbb{F}^{m \times l}$.*

Then for any integer $p \in [1, \min(m, n, p)]$, $\alpha \in \mathbb{Q}_{p,m}, \beta \in \mathbb{Q}_{p,l}$ the following identity holds.

$$\det C[\alpha, \beta] = \sum_{\gamma \in \mathbb{Q}_{p,n}} \det A[\alpha, \gamma] B[\gamma, \beta]. \quad (3.20)$$

Proof. It is enough to assume the case where $\alpha = \beta = \{1, \dots, p\}$. This is equivalent to the assumption that $p = m = l \leq n$. For $p = m = n = l$ the Cauchy-Binet formula reduces to $\det AB = (\det A)(\det B)$. So it is enough to consider the case $p = m = l < n$. Let $C = [c_{ij}]_{i,j=1}^p$. Then $c_{ij} = \sum_{t_j=1}^p a_{it_j} b_{t_j j}$ for $i, j = 1, \dots, p$. Use multilinearity of the determinant to deduce

$$\begin{aligned} \det C &= \det \left[\sum_{t_j=1}^p a_{it_j} b_{t_j j} \right]_{i,j=1}^p = \sum_{t_1, \dots, t_p=1}^p \det [a_{it_j} b_{t_j j}]_{i,j=1}^p = \\ &= \sum_{t_1, \dots, t_p=1}^p \det A[\{1, \dots, p\}, \{t_1, \dots, t_p\}] b_{t_1 1} b_{t_2 2} \dots b_{t_p p}. \end{aligned}$$

Observe next that $\det A[\{1, \dots, p\}, \{t_1, \dots, t_p\}] = 0$ if $t_i = t_j$ for some $1 \leq i < j \leq p$, since the the columns t_i and t_j in $A[\{1, \dots, p\}, \{t_1, \dots, t_p\}]$ are equal. Consider the sum

$$\sum_{\{t_1, t_2, \dots, t_p\} = \gamma \in \mathbb{Q}_{p,n}} A[\{1, \dots, p\}, \{t_1, \dots, t_p\}] b_{t_1 1} b_{t_2 2} \dots b_{t_p p}.$$

The above arguments yield that this sum is $\det(A[\langle p \rangle, \gamma] C[\gamma, \langle p \rangle]) = (\det A[\langle p \rangle, \gamma]) (\det C[\gamma, \langle p \rangle])$. This establishes (3.20). \square

Proposition 3.35 *Let $A(z) \in \mathbb{F}^{m \times n}[z]$. Denote by the rank of $A(z)$, denoted by $r = \text{rank } A(z)$ the size of the biggest minor of $A(z)$ which is not a zero polynomial. For an integer $k \in [1, r]$ let $\delta_k(z)$ be the greatest common divisor of all $k \times k$ minors of $A(z)$, which is assumed to be a monic polynomial. Assume that $\delta_0 = 1$. Then $\delta_i | \delta_{i+1}$ for $i = 0, \dots, r-1$.*

Proof. Expand a nonzero $(k+1) \times (k+1)$ minor of $A(z)$ to deduce that $\delta_k(z) | \delta_{k+1}(z)$. In particular, $\delta_i | \delta_j$ for $1 \leq i < j \leq r$. \square

Proposition 3.36 *Let $A(z) \in \mathbb{F}[z]^{m \times n}$. Assume that $P \in \text{GL}(m, \mathbb{F}), Q \in \text{GL}(n, \mathbb{F})$. Let $B(z) = PA(z)Q$. Then*

1. $\text{rank } A(z) = \text{rank } B(z) = r$.
2. $\delta_k(A(z)) = \delta_k(B(z))$ for $k = 0, \dots, r$.

Proof. Use Cauchy-Binet to deduce that $\text{rank } B(z) \leq \text{rank } A(z)$. Since $A(z) = P^{-1}B(z)Q^{-1}$ it follows that $\text{rank } A(z) \leq \text{rank } B(z)$. Hence $\text{rank } A(z) = \text{rank } B(z) = r$. Use the Cauchy-Binet to deduce that $\delta_k(A(z)) | \delta_k(B(z))$. As $A(z) = P^{-1}B(z)Q^{-1}$ we deduce also that $\delta_k(B(z)) | \delta_k(A(z))$. \square

Definition 3.37 *Let $A \in \mathbb{F}^{n \times n}$ and define $A(z) := zI_n - A$. Then the polynomials the polynomials of degree 1 at least in the sequence $\phi_i = \frac{\delta_{n-i+1}}{\delta_{n-i}(z)}, i = 1, \dots, n$ are called the invariant polynomials of A .*

Note that the product of all invariant polynomials is $\det(zI_n - A)$. In view of Proposition 3.36 we obtain that similar matrices have the same invariant polynomials. Hence for linear transformation $T : \mathbf{V} \rightarrow \mathbf{V}$ we can define its invariant polynomials of T by any representation matrix of T in a basis of \mathbf{V} .

Theorem 3.38 *Let $T : \mathbf{V} \rightarrow \mathbf{V}$. Then the invariant polynomials of T are the polynomials ψ_1, \dots, ψ_r appearing in Theorem 3.32.*

Proof. From the proof of Theorem 3.32 it follows that \mathbf{V} has a basis in which T is represented by $C := \bigoplus_{i=1}^r C(\psi_i)$. It is easy to see that for any monic polynomial θ if $\text{degree } m \geq 1$ $\delta_{m-1}(zI_m - C(\theta))$ is 1. Hence $\delta_{n-r}(C) = 1$, and $\delta_{n-i} = \psi_r \dots \psi_i$ for $1 \leq i < r$. Hence $\psi_i = \frac{\delta_{n-i+1}}{\delta_{n-i}}$. \square

Note that the above theorem implies that $\psi_{i+1} | \psi_i$ for $i = 1, \dots, r-1$.

The irreducible factors $\phi_{i,1}, \dots, \phi_{i,t_i}$ of ψ_i given in Theorem 3.33 for $i = 1, \dots, r$ are called the *elementary divisors* of A . The matrices $\bigoplus_{i=1}^k C(\psi_i), \bigoplus_{i=l=1}^{k,t_i} C(\psi_{i,l})$ are called the *rational canonical forms* of A . Those are the canonical forms in the case that the characteristic polynomial of A does not split to linear factors over \mathbb{F} .

Problems

1. Let $\phi(z) = z^l + \sum_{i=1}^l a_i z^{l-i}$. Denote by $C(\phi)$ the matrix (3.18). Show
 - (a) $\phi(z) = \det(zI - C(\phi))$.
 - (b) Show that the minimal polynomial of $C(\phi)$ is ϕ , i.e. $C(\phi)$ is nonderogatory.
 - (c) Assume that $A \in \mathbb{F}^{n \times n}$ is nonderogatory. Show that A is similar to $C(\phi)$, where ϕ is a characteristic polynomial of A .
2. Let $T : \mathbf{V} \rightarrow \mathbf{V}, \dim \mathbf{V} < \infty, \mathbf{0} \neq \mathbf{u}, \mathbf{w}$. Show that $\psi_{\mathbf{u}+\mathbf{w}} = \frac{\psi_{\mathbf{u}}\psi_{\mathbf{w}}}{(\psi_{\mathbf{u}}, \psi_{\mathbf{w}})}$.
3. Let $T : V \rightarrow V$ and assume that V is a cyclic space, i.e. T is nondegenerate. Let ψ be the minimal and the characteristic polynomial. Assume that $\psi = \psi_1\psi_2$, where $\text{deg } \psi_1, \text{deg } \psi_2 \geq 1$, and $(\psi_1, \psi_2) = 1$. Show that there exists $\mathbf{0} \neq \mathbf{u}_1, \mathbf{u}_2$ such that $\psi_i = \psi_{\mathbf{u}_i}, i = 1, 2$. Furthermore $\mathbf{V} = \mathbf{U}_1 \oplus \mathbf{U}_2$, where \mathbf{U}_i is the cyclic subspace generated by \mathbf{u}_i .
4. Let T is a linear operator $T : \mathbf{V} \rightarrow \mathbf{V}$ and \mathbf{U} is a subspace of \mathbf{V} . Show
 - (a) Assume that $T\mathbf{U} \subseteq \mathbf{U}$. Show that T induces a linear transformation $\hat{T} : \mathbf{V}/\mathbf{U} \rightarrow \mathbf{V}/\mathbf{U}$, i.e. $\hat{T}[\mathbf{x}] := [T\mathbf{x}]$.
 - (b) Assume that $T\mathbf{U} \not\subseteq \mathbf{U}$. Show that $\hat{T}[\mathbf{x}] := [T\mathbf{x}]$ does not make sense.
5. In this problem we finish the proof of Theorem 3.32.
 - (a) Show that ψ_i is an annihilating polynomial of \mathbf{u}_i for $i \geq 2$.
 - (b) By considering the vector $[\hat{\mathbf{u}}_i] = [\mathbf{u}_i]$ show that $\psi_{\mathbf{u}_i} = \psi_i$ for $i \geq 2$.
 - (c) Show that the vectors $\mathbf{u}_i, T\mathbf{u}_i, \dots, T^{\text{deg } \psi_i - 1}\mathbf{u}_i, i = 1, \dots, r$ are linearly independent. (**Hint:** Assume linear dependence of all vectors, and then quotient this dependence by \mathbf{U}_1 . Then use the induction hypothesis on \hat{V} .)

- (d) Let \mathbf{U}_i be the cyclic subspace generated by \mathbf{u}_i for $i = 2, \dots, r$. Conclude the proof of Theorem 3.32.
6. Let the assumptions of Theorem 3.32 hold. Show that the characteristic polynomial of T is equal to $\prod_{i=1}^r \psi_i$.

4 Applications of Jordan Canonical form

4.1 Functions of Matrices

Let $A \in \mathbb{C}^{n \times n}$. Consider the iterations

$$\mathbf{x}_l = A\mathbf{x}_{l-1}, \quad \mathbf{x}_{l-1} \in \mathbb{C}^n, \quad l = 1, \dots \quad (4.1)$$

Clearly $\mathbf{x}_l = A^l \mathbf{x}_0$. To compute \mathbf{x}_l from \mathbf{x}_{l-1} one need to perform $n(2n - 1)$ flops, (operations: n^2 multiplications and $n(n - 1)$ additions). If we want to compute \mathbf{x}_{10^8} we need to $10^8 n(2n - 1)$ operations, if we simply program the iterations (4.1). If $n = 10$ it will take us some time to do these iterations, and we will probably run to the roundoff error, which will render our computations meaningless. Is there any better way to find \mathbf{x}_{10^8} ? The answer is *yes*, and this is the purpose of this section. To do that we need to give the correct way to find directly A^{10^8} , or for that matter any $f(A)$, where $f(z)$ is either polynomial, or more complex functions as $e^z, \cos z, \sin z$, an entire function $f(z)$, or even more special functions.

Theorem 4.1 *Let $A \in \mathbb{C}^{n \times n}$ and*

$$\det(zI_n - A) = \prod_{i=1}^k (z - \lambda_i)^{n_i}, \quad \psi(z) = \prod_{i=1}^k (z - \lambda_i)^{m_i}, \quad (4.2)$$

$$1 \leq m := \deg \psi = \sum_{i=1}^k m_i \leq n = \sum_{i=1}^k n_i, \quad 1 \leq m_i \leq n_i, \quad \lambda_i \neq \lambda_j \text{ for } i \neq j, \quad i, j = 1, \dots, k,$$

where $\psi(z)$ is the minimal polynomial of A . Then there exists unique m linearly independent matrices $Z_{ij} \in \mathbb{C}^{n \times n}$ for $i = 1, \dots, k$ and $j = 0, \dots, m_i - 1$, which depend on A , such that for any polynomial $f(z)$ the following identity holds

$$f(A) = \sum_{i=1}^k \sum_{j=0}^{m_i-1} \frac{f^{(j)}(\lambda_i)}{j!} Z_{ij}. \quad (4.3)$$

($Z_{ij}, i = 1, \dots, k, j = 0, \dots, m_i - 1$ are called the A -components.)

Proof. We start first with $A = J_n(\lambda)$. So $J_n(\lambda) = \lambda I_n + H_n$, where $H_n := J_n(0)$. Thus H_n is a nilpotent matrix, with $H_n^n = \mathbf{0}$ and H_n^j has 1's on the j -th subdiagonal and all other elements are equal 0 for $j = 0, 1, \dots, n - 1$. Hence $I_n = H_n^0, H_n, \dots, H_n^{n-1}$ are linearly independent.

Let $f(z) = z^l$. Then

$$A^l = (\lambda I_n + H_n)^l = \sum_{j=0}^l \binom{l}{j} \lambda^{l-j} H_n^j = \sum_{j=0}^{\min(l, n-1)} \binom{l}{j} \lambda^{l-j} H_n^j.$$

The last equality follows from the equality $H^j = \mathbf{0}$ for $j \geq n$. Note that $\psi(z) = \det(zI_n - J_n(\lambda)) = (z - \lambda)^n$, i.e. $k = 1$ and $m = m_1 = n$. From the above equality we conclude that $Z_{1j} = H_n^j$ for $j = 0, \dots$ if $f(z) = z^l$ and $l = 0, 1, \dots$. With this definition of Z_{1j} (4.3) holds for $K_l z^l$, where $K_l \in \mathbb{C}$ and $l = 0, 1, \dots$. Hence (4.3) holds for any polynomial $f(z)$ for this choice of A .

Assume now that A is a direct sum of Jordan blocks as in (3.6): $A = \bigoplus_{i=1}^{k, l_i} J_{m_{ij}}(\lambda_i)$. Here $m_i = m_{i1} \geq \dots \geq m_{il_i} \geq 1$ for $i = 1, \dots, k$, and $\lambda_i \neq \lambda_j$ for $i \neq j$. Thus (4.2) holds with $n_i = \sum_{j=1}^{l_i} m_{ij}$ for $i = 1, \dots, k$. Let $f(z)$ be a polynomial. Then $f(A) = \bigoplus_{i=1}^{k, l_i} f(J_{m_{ij}}(\lambda_i))$. Use the results for $J_n(\lambda)$ to deduce

$$f(A) = \bigoplus_{i=1}^{k, l_i} \sum_{r=0}^{m_{ij}-1} \frac{f^{(r)}(\lambda_i)}{r!} H_{m_{ij}}^r.$$

Let $Z_{ij} \in \mathbb{C}^{n \times n}$ be a block diagonal matrix of the following form. For each integer $l \in [1, k]$ with $l \neq i$ all the corresponding blocks to $J_{l_r}(\lambda_l)$ are equal to zero. In the block corresponding to $J_{m_{ir}}(\lambda_i)$ Z_{ij} has the block matrix $H_{m_{ir}}^j$ for $j = 0, \dots, m_i - 1$. Note that each Z_{ij} is a nonzero matrix with 0–1 entries. Furthermore, two different Z_{ij} and $Z_{i'j'}$ do not have a common 1 entry. Hence $Z_{ij}, i = 1, \dots, k, j = 0, \dots, m_i - 1$ are linearly independent. It is straightforward to deduce (4.3) from the above identity.

Let $B \in \mathbb{C}^{n \times n}$. Then $B = UAU^{-1}$ where A is the Jordan canonical form of B . Recall that A and B have the same characteristic polynomial. Let $f(z) \in \mathbb{C}[z]$. Then (4.3) holds. Clearly

$$f(B) = Uf(A)U^{-1} = \sum_{i=1}^k \sum_{j=0}^{m_i-1} \frac{f^{(j)}(\lambda_i)}{j!} UZ_{ij}U^{-1}.$$

Hence (4.3) holds for B , where $UZ_{ij}U^{-1}, i = 1, \dots, k, j = 0, \dots, m_{ij}-1$ are the B -components.

The uniqueness of the A -components follows from the existence and uniqueness of the Lagrange-Sylvester interpolation polynomial as explained below.

□

Theorem 4.2 (The Lagrange-Sylvester interpolation polynomial). *Let $\lambda_1, \dots, \lambda_k \in \mathbb{C}$ be k -distinct numbers. Let m_1, \dots, m_k be k positive integers and let $m = m_1 + \dots + m_k$. Let $s_{ij}, i = 1, \dots, k, j = 0, \dots, m_i - 1$ be any m complex numbers. Then there exists a unique polynomial $\phi(z)$ of degree at most $m - 1$ satisfying the conditions $\phi^{(j)}(\lambda_i) = s_{ij}$ for $i = 1, \dots, k, j = 0, \dots, m_i - 1$ satisfying the conditions. (For $m_i = 1, i = 1, \dots, k$ ϕ is the Lagrange interpolating polynomial.)*

Proof. The Lagrange interpolating polynomial is given by the formula

$$\phi(z) = \sum_{i=1}^k \frac{(z - \lambda_1) \dots (z - \lambda_{i-1})(z - \lambda_{i+1}) \dots (z - \lambda_k)}{(\lambda_i - \lambda_1) \dots (\lambda_i - \lambda_{i-1})(\lambda_i - \lambda_{i+1}) \dots (\lambda_i - \lambda_k)} s_{i0}.$$

In the general case one determines $\phi(z)$ as follows. Let $\psi(z) := \prod_{i=1}^k (z - \lambda_i)^{m_i}$. Then

$$\phi(z) = \psi(z) \sum_{i=1}^k \sum_{j=0}^{m_i-1} \frac{t_{ij}}{(z - \lambda_i)^{m_i-j}} = \sum_{i=1}^k \sum_{j=1}^{m_i-1} t_{ij} (z - \lambda_i)^j \psi_i(z), \quad (4.4)$$

$$\psi_i = \frac{\psi(z)}{(z - \lambda_i)^{m_i}} = \prod_{j \neq i} (z - \lambda_j)^{m_j}, \quad \frac{d^l \psi_i}{dz^l}(\lambda_r) = 0 \text{ for } l = 0, \dots, m_r - 1 \text{ and } r \neq i.$$

Now start to determine $t_{i0}, t_{i1}, \dots, t_{i(m_i-1)}$ recursively for each fixed value of i . This is done by using the values of $\phi(\lambda_i), \phi'(\lambda_i), \dots, \phi^{(m_i-1)}(\lambda_i)$ in the above formula for $\phi(z)$. Note that $\deg \phi \leq m - 1$. It is straightforward to show that

$$t_{i0} = \frac{\phi(\lambda_i)}{\psi_i(\lambda_i)}, \quad t_{i1} = \frac{\phi'(\lambda_i) - t_{i0} \psi_i'(\lambda_i)}{\psi_i(\lambda_i)}, \quad t_{i2} = \frac{\phi''(\lambda_i) - t_{i0} \psi_i''(\lambda_i) - 2t_{i1} \psi_i'(\lambda_i)}{2\psi_i(\lambda_i)}. \quad (4.5)$$

The uniqueness ϕ is shown as follows. Assume that $\theta(z)$ is another Lagrange-Sylvester polynomial of degree less than m . Then $\omega(z) := \phi(z) - \theta(z)$ must be divisible by $(z - \lambda_i)^{m_i}$, since $\omega^{(j)}(\lambda_i) = 0$ for $j = 0, \dots, m_i - 1$, for each $i = 1, \dots, k$. Hence $\psi(z) | \omega(z)$. As $\deg \omega(z) \leq m - 1$ it follows that $\omega(z)$ is the zero polynomial, i.e. $\phi(z) = \theta(z)$. \square

Proof of the uniqueness of A-components. Let $\phi_{ij}(z)$ be the Lagrange-Sylvester polynomial given by the data $s_{i'j'}, i' = 1, \dots, k, j' = 1, \dots, m_{i'} - 1$. Assume $s_{ij} = j!$ and all other $s_{i'j'} = 0$. Then (4.3) yields that $Z_{ij} = \phi_{ij}(A)$. \square

Example: Find the polynomials ϕ_{ij} for $\psi(z) = z^3(z - 1)^2(z + 1)$.

Solution: $\lambda_1 = 0, m_1 = 3, \lambda_2 = 1, m_2 = 2, \lambda_3 = -1, m_3 = 1$.

$$\psi_1 = (z - 1)^2 z, \quad \psi_2 = z^3 (z + 1), \quad \psi_3 = z^3 (z - 1)^2.$$

Proposition 4.3 *Let $A \in \mathbb{C}^{n \times n}$. Assume that the minimal polynomial $\psi(z)$ be given by (4.2) and denote $m = \deg \psi$. Then for each integers $u, v \in [1, n]$ denote by $a_{uv}^{(l)}$ and $(Z_{ij})_{uv}$ the (u, v) entries of A^l and of the A-component Z_{ij} respectively. Then $(Z_{ij})_{uv}, i = 1, \dots, k, j = 0, \dots, m_i - 1$ are the unique solutions of the following system with m unknowns*

$$\sum_{i=1}^k \sum_{j=0}^{m_i-1} \binom{l}{j} \lambda_i^{\max(l-j, 0)} (Z_{ij})_{uv} = a_{uv}^{(l)}, \quad l = 0, \dots, m - 1. \quad (4.6)$$

(Note that $\binom{l}{j} = 0$ for $j > l$.)

Proof. Consider the equality (4.3) for $f(z) = z^l$ where $l = 0, \dots, m - 1$. Restricting these equalities to (u, v) entries we deduce that $(Z_{ij})_{uv}$ satisfy the system (4.6). Thus the systems (4.6) are solvable for each pair $(u, v), u, v = 1, \dots, n$. Let $X_{ij} \in \mathbb{C}^{n \times n}, i = 1, \dots, k, j = 1, \dots, m_i - 1$ such that $((X_{ij})_{uv})$ satisfy the system (4.6) for each $u, v \in [1, n]$. Hence $f(A) = \sum_{i=1}^k \sum_{j=0}^{m_i-1} \frac{f^{(j)}(\lambda_i)}{j!} T_{ij}$ for $f(z) = z^l$ and $l = 0, \dots, m - 1$. Hence the above equality holds for any polynomial $f(z)$ of degree

less than m . Apply the above formula to the Lagrange-Sylvester polynomial ϕ_{ij} as given in the proof of the uniqueness of the A -components. Then $\phi_{ij}(A) = X_{ij}$. So $X_{ij} = Z_{ij}$. Thus each system (4.6) has a unique solution. \square

The algorithm for finding the A -components and its complexity.

1. (a) Set $i = 1$.
 (b) Compute and store A^i . Check if I_n, A, \dots, A^i are linearly independent. If independent, set $i = i + 1$ and go to (b).
 (c) $m = i$ and express $A^m = \sum_{i=1}^m a_i A^{m-i}$. Then $\psi(z) = z^m - \sum_{i=1}^m a_i z^{m-i}$ is the minimal polynomial.
 (d) Find the k roots of $\psi(z)$ and their multiplicities: $\psi(z) = \prod_{i=1}^k (z - \lambda_i)^{m_i}$.
 (e) Find the A -components by solving n^2 systems (4.6).
2. The maximum complexity to find $\psi(z)$ happens when $m = n$. Then we need to compute and store I_n, A, A^2, \dots, A^n . So we need n^3 storage space. Viewing I_n, A, \dots, A^i as row vectors arranged as $i \times n^2$ matrix $B_i \in \mathbb{C}^{i \times n^2}$, we bring B_i to a row echelon form: $C_i = U_i B_i, U_i \in \mathbb{C}^{i \times i}$. Note that C_i is essentially upper triangular. Then we add $i + 1$ -th row: A^{i+1} to the B_i to obtain $C_{i+1} = U_{i+1} B_{i+1}$. (C_i is $i \times i$ submatrix of C_{i+1} .) To get C_{i+1} from C_i we need $2in^2$ flops. In the case $m = n$ C_{n^2+1} has last row zero. So to find $\psi(z)$ we need at most Kn^4 flops. ($K \leq 2?$). The total storage space is around $2n^3$.

Now to find the roots of $\psi(z)$ with certain precision will take a polynomial time, depending on the precision.

To solve n^2 systems with n variables, given in (4.6), use Gauss-Jordan for the augmented matrix $[S \ T]$. Here $S \in \mathbb{C}^{n \times n}$ stands for the coefficient of the system (4.6), depending on $\lambda_1, \dots, \lambda_k$. $T \in \mathbb{C}^{n \times n^2}$ given the "left-hand side" of n^2 systems of (4.6). One needs around n^3 storage space. Bring $[S \ T]$ to $[I_n \ Q]$ using Gauss-Jordan to find A -components. To do that we need about n^4 flops.

In summary, we need storage of $2n^3$ and around $4n^4$ flops. (This would suffice to find the roots of $\psi(z)$ with good enough precision.)

Problems

1. Let $A \in \mathbb{C}^{4 \times 4}$ be given as in Problem 3 of Section 3.3. Assume that the characteristic polynomial of A is $z^2(z - 1)^2$.
 (a) Use Problem 4 of Section 3.4 to find the Jordan canonical form of A .
 (b) Assume that the minimal polynomial of A is $z(z - 1)^2$. Find all the A -components.
 (c) Give the explicit formula for any A^l .
2. Let $A \in \mathbb{C}^{n \times n}$ and assume that $\det(zI_n - A) = \prod_{i=1}^k (z - \lambda_i)^{n_i}$, and the minimal polynomial $\psi(z) = \prod_{i=1}^k (z - \lambda_i)^{m_i}$ where $\lambda_1, \dots, \lambda_k$ are k distinct eigenvalues of A . Let $Z_{ij}, j = 0, \dots, m_i - 1, i = 1, \dots, k$ are the A -components.
 (a) Show that $Z_{ij}Z_{pq} = \mathbf{0}$ for $i \neq p$.
 (b) What is the exact formula for $Z_{ij}Z_{ip}$?

4.2 Power stability, convergence and boundedness of matrices

Corollary 4.4 *Let $A \in \mathbb{C}^{n \times n}$. Assume that the minimal polynomial $\psi(z)$ be given by (4.2) and denote by $Z_{ij}, i = 1, \dots, k, j = 0, \dots, m_j - 1$ the A -components. Then for each positive integer l*

$$A^l = \sum_{i=1}^k \sum_{j=0}^{m_i-1} \binom{l}{j} \lambda_i^{\max(l-j, 0)} Z_{ij}. \quad (4.7)$$

If we know the A -components then to compute A^l we need only around $2mn^2 \leq 2n^3$ flops! Thus we need at most $4n^4$ flops to compute A^l , including the computations of A -components, without dependence on l ! (Note that $\lambda_i^j = e^{j \log \lambda_i}$.) So to find $\mathbf{x}_{10^8} = A^{10^8} \mathbf{x}_0$ discussed in the beginning of the previous section we need about 10^4 flops. So to compute \mathbf{x}_{10^8} we need about $10^4 10^2$ flops compared with $10^8 10^2$ flops using the simple minded algorithm explained in the beginning of the previous section. There are much simpler algorithms to compute A^l which are roughly of the order $(\log_2 l)^2 n^3$ of computations and $(\log_2 l)^2 n^2$ ($4n^2$?) storage. See Problem ? However roundoff error remains a problem for large l .

Definition 4.5 *Let $A \in \mathbb{C}^{n \times n}$. A is called power stable if $\lim_{l \rightarrow \infty} A^l = \mathbf{0}$. A is called power convergent if $\lim_{l \rightarrow \infty} A^l = B$ for some $B \in \mathbb{C}^{n \times n}$. A is called power bounded if there exists $K > 0$ such that the absolute value of every entry of every $A^l, l = 1, \dots$ is bounded above by K .*

Theorem 4.6 *Let $A \in \mathbb{C}^{n \times n}$. Then*

1. *A is power stable if and only if each eigenvalue of A is in the interior of the unit disk: $|z| < 1$.*
2. *A is power convergent if and only if each eigenvalue λ of A satisfies one of the following conditions*
 - (a) $|\lambda| < 1$;
 - (b) $\lambda = 1$ and each Jordan block of the JCF of A with an eigenvalue 1 is of order 1, i.e. 1 is a simple zero of the minimal polynomial of A .
3. *A is power bounded if and only if each eigenvalue λ of A satisfies one of the following conditions*
 - (a) $|\lambda| < 1$;
 - (b) $|\lambda| = 1$ and each Jordan block of the JCF of A with an eigenvalue λ is of order 1, i.e. λ is a simple zero of the minimal polynomial of A .

Proof. Consider the formula (4.4). Since the A -components $Z_{ij}, i = 1, \dots, k, j = 0, \dots, m_i - 1$ are linearly independent we need to satisfy the conditions of the theorem for each term in (4.4), which is $\binom{l}{j} \lambda_i^{l-j} Z_{ij}$ for $l \gg 1$. Note that for a fixed j $\lim_{l \rightarrow \infty} \binom{l}{j} \lambda_i^{l-j} = 0$ if and only if $|\lambda_i| < 1$. Hence we deduce the condition 1 of the theorem.

Note that the sequence $\binom{l}{j} \lambda_i^{l-j}, l = j, j+1, \dots$, converges if and only if either $|\lambda_i| < 1$ or $\lambda_i = 1$ and $j = 0$. Hence we deduce the condition 2 of the theorem.

Note that the sequence $\binom{l}{j}\lambda_i^{l-j}, l = j, j+1, \dots$, is bounded if and only if either $|\lambda_i| < 1$ or $|\lambda_i| = 1$ and $j = 0$. Hence we deduce the condition \mathcal{B} of the theorem. \square

Corollary 4.7 *Let $A \in \mathbb{C}^{n \times n}$ and consider the iterations $\mathbf{x}_l = A\mathbf{x}_{l-1}$ for $l = 1, \dots$. Then for any \mathbf{x}_0*

1. $\lim_{l \rightarrow \infty} \mathbf{x}_l = \mathbf{0}$ if and only if A is power stable.
2. $\mathbf{x}_l, l = 0, 1, \dots$ converges if and only if A is power convergent.
3. $\mathbf{x}_l, l = 0, 1, \dots$ is bounded if and only if A is power bounded.

Proof. If A satisfies the conditions of an item i Theorem 4.6 then the corresponding condition i of the corollary clearly holds. Assume that the conditions of an item i of the corollary holds. Choose $\mathbf{x}_0 = \mathbf{e}_j = (\delta_{1j}, \dots, \delta_{nj})^\top$ for $j = 1, \dots, n$ to deduce the corresponding condition i of Theorem 4.6. \square

Theorem 4.8 *Let $A \in \mathbb{C}^{n \times n}$ and consider the nonhomogeneous iterations*

$$\mathbf{x}_l = A\mathbf{x}_{l-1} + \mathbf{b}_l, \quad l = 0, \dots \quad (4.8)$$

Then

1. $\lim_{l \rightarrow \infty} \mathbf{x}_l = \mathbf{0}$ for any $\mathbf{x}_0 \in \mathbb{C}^n$ and any sequence $\mathbf{b}_0, \mathbf{b}_1, \dots$ satisfying the condition $\lim_{l \rightarrow \infty} \mathbf{b}_l = \mathbf{0}$ if and only if A is power stable.
2. The sequence $\mathbf{x}_l, l = 0, 1, \dots$ converges for any \mathbf{x}_0 and any sequence $\mathbf{b}_0, \mathbf{b}_1, \dots$ satisfying the condition $\sum_{l=0}^l \mathbf{b}_l$ converges if and only if A is power convergent.
3. The sequence $\mathbf{x}_l, l = 0, 1, \dots$ is bounded for any \mathbf{x}_0 and any sequence $\mathbf{b}_0, \mathbf{b}_1, \dots$ satisfying the condition $\sum_{l=0}^l \|\mathbf{b}_l\|_\infty$ converges if and only if A is power bounded. (Here $\|(x_1, \dots, x_n)\|_\infty = \max_{i \in [1, n]} |x_i|$.)

Proof. Assume that $\mathbf{b}_l = \mathbf{0}$. Since \mathbf{x}_0 is arbitrary we deduce the necessity of all the conditions from Theorem 4.6. The sufficiency of the above conditions follow from the Jordan Canonical Form of A as follows.

Let $J = U^{-1}AU$ where U is an invertible matrix and J is the Jordan canonical form of A . By letting $\mathbf{y}_l := U^{-1}\mathbf{x}_l$ and $\mathbf{c}_l = U^{-1}\mathbf{b}_l$ it is enough to prove the sufficiency part of the theorem for the case where A is sum of Jordan blocks. In this case system (4.8) reduces to independent systems of equations for each Jordan block. Thus it is left to prove the theorem when $A = J_n(\lambda)$.

1. We show that if $A = J_n(\lambda)$ and $|\lambda| < 1$, then $\lim_{l \rightarrow \infty} \mathbf{x}_l = \mathbf{0}$ for any \mathbf{x}_0 and $\mathbf{b}_l, l = 1, \dots$ if $\lim_{l \rightarrow \infty} \mathbf{b}_l = \mathbf{0}$. We prove this claim by the induction on n . For $n = 1$ (4.8) reduces to

$$x_l = \lambda x_{l-1} + b_l, \quad x_0, x_l, b_l \in \mathbb{C} \text{ for } l = 1, \dots \quad (4.9)$$

It is straightforward to show, e.g. use induction that

$$x_l = \sum_{i=0}^l \lambda^i b_{l-i} = b_l + \lambda b_{l-1} + \dots + \lambda^l b_0 \quad l = 1, \dots, \text{ where } b_0 := x_0. \quad (4.10)$$

Let $\beta_m = \sup_{i \geq m} |b_i|$. Since $\lim_{l \rightarrow \infty} b_l = 0$, it follows that each β_m is finite, the sequence $\beta_m, m = 0, 1, \dots$ decreasing and $\lim_{m \rightarrow \infty} \beta_m = 0$. Fix m . Then for $l > m$

$$\begin{aligned} |x_l| &\leq \sum_{i=0}^l |\lambda|^i |b_{l-i}| = \sum_{i=0}^{l-m} |\lambda|^i |b_{l-i}| + |\lambda|^{l-m} \sum_{j=1}^m |\lambda|^j |b_{m-j}| \leq \\ &\beta_m \sum_{i=0}^{l-m} |\lambda|^i + |\lambda|^{l-m} \sum_{j=1}^m |\lambda|^j |b_{m-j}| \leq \beta_m \sum_{i=0}^{\infty} |\lambda|^i + |\lambda|^{l-m} \sum_{j=1}^m |\lambda|^j |b_{m-j}| = \\ &\frac{\beta_m}{1-|\lambda|} + |\lambda|^{l-m} \sum_{j=1}^m |\lambda|^j |b_{m-j}| \rightarrow \frac{\beta_m}{1-|\lambda|} \text{ as } l \rightarrow \infty. \end{aligned}$$

That is $\limsup_{l \rightarrow \infty} |x_l| \leq \frac{\beta_m}{1-|\lambda|}$. As $\lim_{m \rightarrow \infty} \beta_m = 0$ it follows that $\limsup_{l \rightarrow \infty} |x_l| = 0$, which is equivalent to the statement $\lim_{l \rightarrow \infty} x_l = 0$. This proves the case $n = 1$.

Assume that the theorem holds for $n = k$. Let $n = k + 1$. View $\mathbf{x}_l^\top := (x_{1,l}, \mathbf{y}_l^\top)^\top$, $\mathbf{b}_l = (b_{1,l}, \mathbf{c}_l^\top)^\top$, where $\mathbf{y}_l = (x_{2,l}, \dots, x_{k+1,l})^\top$, $\mathbf{c}_l \in \mathbb{C}^k$ are the vectors composed of the last k coordinates of \mathbf{x}_l and \mathbf{b}_l respectively. Then (4.8) for $A = J_{k+1}(\lambda)$ for the last k coordinates of \mathbf{x}_l is given by the system $\mathbf{y}_l = J_k(\lambda) \mathbf{y}_{l-1} + \mathbf{c}_l$ for $l = 1, 2, \dots$. Since $\lim_{l \rightarrow \infty} \mathbf{c}_l = \mathbf{0}$ the induction hypothesis yields that $\lim_{l \rightarrow \infty} \mathbf{y}_l = \mathbf{0}$. The system (4.8) for $A = J_{k+1}(\lambda)$ for the first coordinate is $x_{1,l} = \lambda x_{1,l-1} + (x_{2,l-1} + b_{1,l})$ for $l = 1, \dots$. From induction hypothesis and the assumption that $\lim_{l \rightarrow \infty} \mathbf{b}_l = \mathbf{0}$ we deduce that $\lim_{l \rightarrow \infty} x_{2,l-1} + b_{1,l} = 0$. Hence from the case $k = 1$ we deduce that $\lim_{l \rightarrow \infty} x_{1,l} = 0$. Hence $\lim_{l \rightarrow \infty} \mathbf{x}_l = \mathbf{0}$. The proof of this case is concluded.

2. Assume that A satisfies the each eigenvalue λ of A satisfies the following conditions: either $|\lambda| < 1$, or $\lambda = 1$ and each Jordan block corresponding to 1 is of order 1. As we pointed out we assume that A is a direct sum of its Jordan form. So first we consider $A = J_k(\lambda)$ with $|\lambda| < 1$. Since we assumed that $\sum_{l=1}^{\infty} \mathbf{b}_l$ converges we deduce that $\lim_{l \rightarrow \infty} \mathbf{b}_l = \mathbf{0}$. Thus, by part 1 we get that $\lim_{l \rightarrow \infty} \mathbf{x}_l = \mathbf{0}$.

Assume now that $A = (1) \in \mathbb{C}^{1 \times 1}$. Thus we consider (4.9) with $\lambda = 1$. (4.10) yields that $x_l = \sum_{i=0}^l b_i$. By the assumption of the theorem $\sum_{i=1}^{\infty} \mathbf{b}_i$ converges, hence the sequence $x_l, l = 1, \dots$ converges.

3. As in the part 2 it is enough to consider the case $J_1(\lambda)$ with $|\lambda| = 1$. Note that (4.10) yields that $|x_l| \leq \sum_{i=0}^l |b_i|$. The assumption that $\sum_{i=1}^{\infty} |\mathbf{b}_i|$ converges imply that $|x_l| \leq \sum_{i=0}^{\infty} |b_i| < \infty$.

□

Remark 4.9 *The stability, convergence and boundedness of the nonhomogeneous systems:*

$$\begin{aligned}\mathbf{x}_l &= A_l \mathbf{x}_{l-1}, \quad A_l \in \mathbb{C}^{n \times n}, \quad l = 1, \dots, \\ \mathbf{x}_l &= A_l \mathbf{x}_{l-1} + \mathbf{b}_l, \quad A_l \in \mathbb{C}^{n \times n}, \quad \mathbf{b}_l \in \mathbb{C}^n \quad l = 1, \dots,\end{aligned}$$

are much harder to analyze. (If time permits we revisit these problems later on in the course.)

Problems

1. Consider the nonhomogeneous system $\mathbf{x}_l = A_l \mathbf{x}_{l-1}$, $A_l \in \mathbb{C}^{n \times n}$, $l = 1, \dots$. Assume that the sequence $A_l, l = 1, \dots$, is periodic, i.e. $A_{l+p} = A_l$ for all $l = 1, \dots$, and a fixed positive integer p .
 - (a) Show that for each $\mathbf{x}_0 \in \mathbb{C}^n$ $\lim_{l \rightarrow \infty} \mathbf{x}_l = \mathbf{0}$ if and only if $B := A_p A_{p-1} \dots A_1$ is power stable.
 - (b) Show that for each $\mathbf{x}_0 \in \mathbb{C}^n$ the sequence $\mathbf{x}_l, l = 1, \dots$, converges if and only if the following conditions satisfied. First, B is power convergent, i.e. $\lim_{l \rightarrow \infty} B^l = C$. Second, $A_i C = C$ for $i = 1, \dots, p$.
 - (c) Find a necessary and sufficient conditions such that for each $\mathbf{x}_0 \in \mathbb{C}^n$ the sequence $\mathbf{x}_l, l = 1, \dots$, is bounded.

4.3 e^{At} and stability of certain systems of ODE

Recall that the exponential function e^z has the MacLaurin expansion

$$e^z = 1 + z + \frac{z^2}{2} + \frac{z^3}{6} + \dots = \sum_{l=0}^{\infty} \frac{z^l}{l!}.$$

Hence for each $A \in \mathbb{C}^{n \times n}$ one defines

$$e^A := I_n + A + \frac{A^2}{2} + \frac{A^3}{6} + \dots = \sum_{l=0}^{\infty} \frac{A^l}{l!}.$$

More generally, if $t \in \mathbb{C}$ then

$$e^{At} := I_n + At + \frac{A^2 t^2}{2} + \frac{A^3 t^3}{6} + \dots = \sum_{l=0}^{\infty} \frac{A^l t^l}{l!}.$$

Hence e^{At} satisfies the matrix differential equation

$$\frac{d e^{At}}{dt} = A e^{At} = e^{At} A. \quad (4.11)$$

Also one has the standard identity $e^{At} e^{Au} = e^{A(t+u)}$ for any complex numbers t, u .

Proposition 4.10 *Let $A \in \mathbb{C}^{n \times n}$ and consider the system of linear system of n ordinary differential equations with constant coefficients $\frac{d\mathbf{x}(t)}{dt} = A\mathbf{x}(t)$, where $\mathbf{x}(t) \in \mathbb{C}^n$, satisfying the initial conditions $\mathbf{x}(t_0) = \mathbf{x}_0$. Then $\mathbf{x}(t) = e^{A(t-t_0)} \mathbf{x}_0$*

is the unique solution to the above system. More generally, let $\mathbf{b}(t) \in \mathbb{C}^n$ be any continuous vector function on \mathbb{R} and consider the nonhomogeneous system of n ordinary differential equations with the initial condition:

$$\frac{d\mathbf{x}(t)}{dt} = A\mathbf{x}(t) + \mathbf{b}(t), \quad \mathbf{x}(t_0) = \mathbf{x}_0. \quad (4.12)$$

Then this system has a unique solution of the form

$$\mathbf{x}(t) = e^{A(t-t_0)}\mathbf{x}_0 + \int_{t_0}^t e^{A(t-u)}\mathbf{b}(u)du. \quad (4.13)$$

Proof. The uniqueness of the solution of (4.12) follows from the uniqueness of solutions to system of ODE (Ordinary Differential Equations). The first part of the proposition follows from (4.11). To deduce the second part one does the *variations of parameters*. Namely one tries a solution $\mathbf{x}(t) = e^{A(t-t_0)}\mathbf{y}(t)$ where $\mathbf{y}(t) \in \mathbb{C}^n$ is unknown vector function. Hence

$$\mathbf{x}' = (e^{A(t-t_0)})'\mathbf{y}(t) + e^{A(t-t_0)}\mathbf{y}'(t) = Ae^{A(t-t_0)}\mathbf{y}(t) + e^{A(t-t_0)}\mathbf{y}'(t) = A\mathbf{x}(t) + e^{A(t-t_0)}\mathbf{y}'(t).$$

Substitute this expression of $\mathbf{x}(t)$ to (4.12) to deduce the differential equation $\mathbf{y}' = e^{-A(t-t_0)}\mathbf{b}(t)$. Since $\mathbf{y}(t_0) = \mathbf{x}_0$ this simple equation have a unique solution $\mathbf{y}(t) = \mathbf{x}_0 + \int_{t_0}^t e^{A(u-t_0)}\mathbf{b}(u)du$. Now multiply by $e^{A(t-t_0)}$ and use the fact that $e^{At}e^{Au} = e^{A(u+v)}$ to deduce (4.13). \square

Note: The second term in the formula (4.13) can be considered as a perturbation term to the solution $\frac{d\mathbf{x}(t)}{dt} = A\mathbf{x}(t), \mathbf{x}(t_0) = \mathbf{x}_0$, i.e. to the system (4.12) with $\mathbf{b}(t) \equiv \mathbf{0}$.

Use (4.3) for e^{zt} and the observation that $\frac{d^j e^{zt}}{dz^j} = t^j e^{zt}, j = 0, 1, \dots$ to deduce:

$$e^{At} = \sum_{j=1}^k \sum_{i=0}^{m_i-1} \frac{t^j e^{\lambda_i t}}{j!} Z_{ij}. \quad (4.14)$$

We can substitute this expression for e^{At} in (4.13) to get a simple expression of the solution $\mathbf{x}(t)$ of (4.12).

Definition 4.11 Let $A \in \mathbb{C}^{n \times n}$. A is called *exponentially stable*, or *simple stable*, if $\lim_{t \rightarrow \infty} e^{At} = \mathbf{0}$. A is called *exponentially convergent* if $\lim_{t \rightarrow \infty} e^{At} = B$ for some $B \in \mathbb{C}^{n \times n}$. A is called *exponentially bounded* if there exists $K > 0$ such that the absolute value of every entry of every $e^{At}, t \in [0, \infty)$ is bounded above by K .

Theorem 4.12 Let $A \in \mathbb{C}^{n \times n}$. Then

1. A is stable if and only if each eigenvalue of A is in the left half of the complex plane: $\Re z < 0$.
2. A is exponentially convergent if and only if each eigenvalue λ of A satisfies one of the following conditions
 - (a) $\Re \lambda < 0$;

(b) $\lambda = 2\pi l\sqrt{-1}$ for some integer l , and each Jordan block of the JCF of A with an eigenvalue λ is of order 1, i.e. λ is a simple zero of the minimal polynomial of A .

3. A is exponentially bounded if and only if each eigenvalue λ of A satisfies one of the following conditions

(a) $\Re\lambda < 0$;

(b) $\Re\lambda = 0$ and each Jordan block of the JCF of A with an eigenvalue λ is of order 1, i.e. λ is a simple zero of the minimal polynomial of A .

Proof. Consider the formula (4.14). Since the A -components $Z_{ij}, i = 1, \dots, k, j = 0, \dots, m_i - 1$ are linearly independent we need to satisfy the conditions of the theorem for each term in (4.14), which is $\frac{t^j}{j!} e^{\lambda_i t} Z_{ij}$. Note that for a fixed j $\lim_{t \rightarrow \infty} \frac{t^j}{j!} e^{\lambda_i t} = 0$ if and only if $\Re\lambda_i < 0$. Hence we deduce the condition 1 of the theorem.

Note that the function $\frac{t^j}{j!} e^{\lambda_i t}$ converges as $t \rightarrow \infty$ if and only if either $\Re\lambda_i < 0$ or $e^{\lambda_i} = 1$ and $j = 0$. Hence we deduce the condition 2 of the theorem.

Note that the function $\frac{t^j}{j!} e^{\lambda_i t}$ is bounded for $t \geq 0$ if and only if either $\Re\lambda_i < 0$ or $|e^{\lambda_i}| = 1$ and $j = 0$. Hence we deduce the condition 3 of the theorem. \square

Corollary 4.13 Let $A \in \mathbb{C}^{n \times n}$ and consider the system of differential equations $\frac{d\mathbf{x}(t)}{dt} = A\mathbf{x}(t), \mathbf{x}(t_0) = \mathbf{x}_0$. Then for any \mathbf{x}_0

1. $\lim_{t \rightarrow \infty} \mathbf{x}(t) = \mathbf{0}$ if and only if A is stable.
2. $\mathbf{x}(t)$ converges as $t \rightarrow \infty$ if and only if A is exponentially convergent.
3. $\mathbf{x}(t), t \in [0, \infty)$ is bounded if and only if A is exponentially bounded.

Theorem 4.14 Let $A \in \mathbb{C}^{n \times n}$ and consider the system of differential equations (4.12). Then for any $\mathbf{x}_0 \in \mathbb{C}^n$

1. $\lim_{t \rightarrow \infty} \mathbf{x}(t) = \mathbf{0}$ for any continuous function $\mathbf{b}(t)$, such that $\lim_{t \rightarrow \infty} \mathbf{b}(t) = \mathbf{0}$, if and only if A is stable.
2. $\mathbf{x}(t)$ converges as $t \rightarrow \infty$ for any continuous function $\mathbf{b}(t)$, such that $\int_{t_0}^{\infty} \mathbf{b}(u) du$ converges, if and only if A is exponentially convergent.
3. $\mathbf{x}(t), t \in [0, \infty)$ is bounded for any continuous function $\mathbf{b}(t)$, such that $\int_{t_0}^{\infty} |\mathbf{b}(u)| du$ converges if and only if A is exponentially bounded.

Proof. The necessity of the conditions of the theorem follow from Corollary 4.13 by choosing $\mathbf{b}(t) \equiv \mathbf{0}$.

1. Suppose that A is stable. Then Corollary 4.13 yields that $\lim_{t \rightarrow \infty} e^{At} \mathbf{x}_0 = \mathbf{0}$. Thus show that $\lim_{t \rightarrow \infty} \mathbf{x}(t) = \mathbf{0}$, it is enough to show that the second term in (4.13) tends to $\mathbf{0}$. Use (4.14) to show that it is enough to demonstrate that

$$\lim_{t \rightarrow \infty} \int_{t_0}^t (t-u)^j e^{\lambda(t-u)} g(u) du = 0, \text{ where } \Re\lambda < 0,$$

for any continuous $g(t) \in [t_0, \infty)$, such that $\lim_{t \rightarrow \infty} g(t) = 0$. For $\epsilon > 0$ there exists $T = T(\epsilon)$ such that $|g(t)| \leq \epsilon$ for $t \geq T(\epsilon)$. Let $t > T(\epsilon)$. Then

$$\begin{aligned} \left| \int_{t_0}^t (t-u)^j e^{\lambda(t-u)} g(u) du \right| &= \left| \int_{t_0}^{T(\epsilon)} (t-u)^j e^{\lambda(t-u)} g(u) du + \int_{T(\epsilon)}^t (t-u)^j e^{\lambda(t-u)} g(u) du \right| \\ &\leq \left| \int_{t_0}^{T(\epsilon)} (t-u)^j e^{\lambda(t-u)} g(u) du \right| + \left| \int_{T(\epsilon)}^t (t-u)^j e^{\lambda(t-u)} g(u) du \right| \\ &\leq \left| \int_{t_0}^{T(\epsilon)} (t-u)^j e^{\lambda(t-u)} g(u) du \right| + \epsilon \int_{T(\epsilon)}^t (t-u)^j e^{\Re \lambda(t-u)} du. \end{aligned}$$

Consider the first term in the last inequality. Since $\lim_{t \rightarrow \infty} t^j e^{\lambda t} = 0$ it follows that the first term converges to zero. The second term bounded by ϵK for $K := \int_0^\infty t^j e^{\Re \lambda t} dt$. Hence as $\epsilon \rightarrow 0$ we deduce that $\lim_{t \rightarrow \infty} \int_{t_0}^t (t-u)^j e^{\lambda(t-u)} g(u) du = 0$.

2. Using part 1 we deduce the result for any eigenvalue λ with $\Re \lambda < 0$. It is left to discuss the case $\lambda = 0$. We assume that the Jordan blocks of A corresponding to $\lambda = 0$ are of length one. So the A -component corresponding to $\lambda = 0$ is Z_{10} . The corresponding term is

...

5 Perron-Frobenius theorem

Denote by $\mathbb{R}_+^{m \times n}$, the set of nonnegative matrices $A = [a_{ij}]_{i,j=1}^{m,n}$, where each entry $a_{ij} \geq 0$. We denote a nonnegative matrix by $A \geq 0$. We say that A is a nonzero nonnegative if $A \geq 0$ and $A \neq 0$. We denote that by $A \succeq 0$. We say that A is a positive matrix if all the entries of A are positive, i.e. $a_{ij} > 0$ for all i, j . We denote that by $A > 0$. Similar notation is for vectors, i.e. $n = 1$. From now on we assume that all matrices are square matrices, unless stated otherwise. A nonnegative square matrix $A \in \mathbb{R}_+^{n \times n}$ is called *irreducible*, if $(I + A)^{n-1} > 0$. A nonnegative matrix is called *primitive* if $A^p > 0$ for some positive integer p .

Theorem 5.1 *Let $A \in \mathbb{R}_+^{n \times n}$. Denote by $\rho(A)$, the spectral radius of A , i.e. the maximum of the absolute values of the eigenvalues of A . The the following conditions hold.*

1. $\rho(A)$ is an eigenvalues of A .
2. There exists $\mathbf{x} \succeq \mathbf{0}$ such that $A\mathbf{x} = \rho(A)\mathbf{x}$.
3. Assume that λ is an eigenvalue of A , $\lambda \neq \rho(A)$ and $|\lambda| = \rho(A)$. Then $\text{index}(\lambda) \leq \text{index}(\rho(A))$. Furthermore $\zeta = \frac{\lambda}{\rho(A)}$ is a root of unity of order n at most, i.e. $\zeta^m = 1$ and $m \leq n$.
4. Suppose that A is primitive. Then $\rho(A) > 0$, $\rho(A)$ is a simple root of the characteristic polynomial, and there exists $\mathbf{x} > \mathbf{0}$ such that $A\mathbf{x} = \rho(A)\mathbf{x}$. Furthermore if λ is an eigenvalue of A different from $\rho(A)$ then $|\lambda| < \rho(A)$.

5. Suppose that A is irreducible. Then $\rho(A) > 0$, $\rho(A)$ is a simple root of the characteristic polynomial, and there exists $\mathbf{x} > \mathbf{0}$ such that $A\mathbf{x} = \rho(A)\mathbf{x}$. Let $\lambda_0 := \rho(A), \lambda_1, \dots, \lambda_{m-1}$ be all distinct eigenvalues of A satisfying $|\lambda| = \rho(A)$. Then all these eigenvalues are simple roots of the characteristic polynomial of A . Furthermore $\lambda_j = \rho(A)\zeta^j$, for $j = 0, \dots, m-1$, where $\zeta = e^{\frac{2\pi\sqrt{-1}}{m}}$ is a primitive m -th root of 1. Finally, the characteristic polynomials of ζA and A are equal.

We first prove this result for symmetric nonnegative matrix.

Proposition 5.2 *Let $A = A^\top \geq 0$ be a symmetric matrix with nonnegative entries. Then the conditions of Theorem 5.1 hold.*

Proof. Recall that A is diagonalizable by an orthogonal matrix. Hence the index of each eigenvalue of A is 1. Next recall that each eigenvalue is given by the Rayleigh quotient $\frac{\mathbf{x}^\top A \mathbf{x}}{\mathbf{x}^\top \mathbf{x}}$. For $\mathbf{x} = (x_1, \dots, x_n)^\top$ let $|\mathbf{x}| := (|x_1|, \dots, |x_n|)^\top$. Note that $\mathbf{x}^\top \mathbf{x} = |\mathbf{x}|^\top |\mathbf{x}|$. Also $|\mathbf{x}^\top A \mathbf{x}| \leq |\mathbf{x}|^\top A |\mathbf{x}|$. Hence $|\frac{\mathbf{x}^\top A \mathbf{x}}{\mathbf{x}^\top \mathbf{x}}| \leq \frac{|\mathbf{x}|^\top A |\mathbf{x}|}{|\mathbf{x}|^\top |\mathbf{x}|}$. Thus, the maximum of the Rayleigh quotient is achieved for some $\mathbf{x} \succeq \mathbf{0}$. This shows that $\rho(A)$ is an eigenvalue of A , and there exists an eigenvector $\mathbf{x} \succeq \mathbf{0}$ corresponding to A .

Suppose that $A > 0$. Then $A\mathbf{x} = \rho(A)\mathbf{x}$, $\mathbf{x} \succeq \mathbf{0}$. So $A\mathbf{x} > \mathbf{0}$. Hence $\rho(A) > 0$ and $\mathbf{x} > \mathbf{0}$. Assume that there exists a linear independent eigenvector \mathbf{y} corresponding to $\rho(A)$. It can be chosen to be orthogonal to the positive eigenvector \mathbf{x} corresponding to $\rho(A)$, i.e $\mathbf{x}^\top \mathbf{y} = 0$. Hence \mathbf{y} has positive and negative coordinates. Therefore

$$\left| \frac{\mathbf{y}^\top A \mathbf{y}}{\mathbf{y}^\top \mathbf{y}} \right| < \frac{|\mathbf{y}|^\top A |\mathbf{y}|}{|\mathbf{y}|^\top |\mathbf{y}|} \leq \rho(A).$$

This gives the contradiction $\rho(A) < \rho(A)$. Therefore, $\rho(A)$ is a simple root of the characteristic polynomial of A . In a similar way we deduce that each eigenvalue λ of A , different from $\rho(A)$ satisfy $|\lambda| < \rho(A)$.

Since the eigenvectors of A^p are the same of A we deduce the same results if A is a nonnegative symmetric matrix which is primitive. Suppose that A is irreducible. Then $(tI + A)^{n-1} > 0$ for any $t > 0$. Clearly $\rho(tI + A) = \rho(A) + t$. Hence the eigenvector corresponding to $\rho(A)$ is positive. Therefore $\rho(A) > 0$. Each eigenvalue of $tI + A$ is of the form $\lambda + t$ for some eigenvalue λ of A . Since $(tI + A)$ is primitive for any $t > 0$ it follows that $|\lambda + t| < \rho(A) + t$. Let $t \rightarrow 0^+$ to deduce that $|\lambda| \leq \rho(A)$. Since all the eigenvalues of A are real, we can have an equality for some eigenvalue $\lambda \neq \rho(A)$ if and only if $\lambda = -\rho(A)$. This will be the case if A has the following form $C = \begin{bmatrix} 0 & B \\ B^\top & 0 \end{bmatrix}$ for some $B \in \mathbb{R}_+^{m \times l}$. It can be shown that if A is a symmetric irreducible matrix such that $-\rho(A)$ then $A = PCP^\top$ for some permutation matrix P . We already showed before that C and $-C$ has the same eigenvalues, counted with their multiplicities. Hence $-A$ and A have the same characteristic polynomial. \square

The next result is basic in proof of Theorem 5.1, and is due to Wielandt.

Lemma 5.3 *Let $0 < A \in \mathbb{R}_+^{n \times n}$. For any positive vector $\mathbf{x} = (x_1, \dots, x_n)^\top$ let $r(A, \mathbf{x}) = \max_{i=1, \dots, n} \frac{A\mathbf{x}}{x_i}$. The $\rho(A) \leq r(A, \mathbf{x})$. Equality holds if and only if*

$A\mathbf{x} = \rho(A)\mathbf{x}$. There is a unique positive eigenvector $\mathbf{x} > \mathbf{0}$ corresponding to $\rho(A)$, up to a multiple by a constant. Any other eigenvalue λ of A satisfies $|\lambda| \leq \rho(A)$.

Proof. Since for any $a > 0, \mathbf{x} > \mathbf{0}$ we have that $r(A, a\mathbf{x}) = r(A, \mathbf{x})$ it is enough to assume that \mathbf{x} is a probability vector, i.e. $\|\mathbf{x}\|_1 = \sum_{i=1}^n x_i = 1$. Since each entry a_{ij} of A is positive, it follows that $a_{ij} \geq t > 0$ for some t . Let $\alpha(\mathbf{x}) = \min_{i=1, \dots, n} x_i$. Then $r(A, \mathbf{x}) \geq \frac{t}{\alpha(\mathbf{x})}$. Consider $\inf_{\mathbf{x} > \mathbf{0}, \|\mathbf{x}\|_1=1} r(A, \mathbf{x})$. In view of $r(A, \mathbf{x}) \geq \frac{t}{\alpha(\mathbf{x})}$ it follows that $r(A, \mathbf{x}) \rightarrow \infty$ when \mathbf{x} tends to the boundary of the set $\mathbf{x} > \mathbf{0}, \|\mathbf{x}\|_1 = 1$. Hence $\inf_{\mathbf{x} > \mathbf{0}, \|\mathbf{x}\|_1=1} r(A, \mathbf{x})$ is achieved at some vector $\mathbf{y} > \mathbf{0}, \|\mathbf{y}\|_1 = 1$. So $r(A, \mathbf{x}) \geq r(A, \mathbf{y})$ for any $\mathbf{x} > \mathbf{0}$. From the definition of $r(A, \mathbf{y})$ it follows that $A\mathbf{y} \leq r(A, \mathbf{y})\mathbf{y}$. We claim that $A\mathbf{y} = r(A, \mathbf{y})\mathbf{y}$. Suppose not. So $(r(A, \mathbf{y})\mathbf{y} - A\mathbf{y}) \succeq \mathbf{0}$. Hence $A(r(A, \mathbf{y})\mathbf{y} - A\mathbf{y}) > \mathbf{0}$. Let $\mathbf{z} := \frac{1}{\|A\mathbf{y}\|_1} A\mathbf{y}$. Clearly $\mathbf{z} > \mathbf{0}, \|\mathbf{z}\|_1 = 1$. The above inequality yields that $r(A, \mathbf{y})\mathbf{z} > A\mathbf{z}$. Hence $r(A, \mathbf{y}) > r(A, \mathbf{z})$. This contradicts the minimality of $r(A, \mathbf{y})$. So $\mathbf{y} = (y_1, \dots, y_n)$ is an eigenvector of A corresponding to the eigenvalue $r := r(A, \mathbf{y})$. $D = \text{diag}(y_1, \dots, y_n)$. Consider the positive matrix $B = D^{-1}AD$. It is straightforward to see that $B\mathbf{1} = r\mathbf{1}$. Hence $\|B\|_\infty = r$. Therefore any eigenvalue λ of B satisfies $|\lambda| \leq r$. Since A and B are similar the eigenvalues of A and B are the same.

It is left to show that $\text{rank}(rI - A) = n - 1$. Assume not. Then there exists another linearly independent eigenvector \mathbf{z} corresponding to r . By considering $\mathbf{y} + t\mathbf{z}$ we can choose $t \in \mathbb{R}$ such that $\mathbf{y} + t\mathbf{z} \succeq \mathbf{0}$ and has at least one coordinate zero. Now $\mathbf{0} < A(\mathbf{y} + t\mathbf{z}) = r(\mathbf{y} + t\mathbf{z})$. So $\mathbf{y} + t\mathbf{z} > \mathbf{0}$ contrary to our assumptions. \square

References

- [1] G.H. Golub and C.F. Van Loan. Matrix Computation, *John Hopkins Univ. Press, 3rd Ed.*, Baltimore, 1996.