

TMA4145 – Linear Methods

Franz Luef

2010 *Mathematics Subject Classification*. Primary

ABSTRACT. These notes are for the course TMA4145 – Linear Methods at NTNU and cover the following topics: Linear and metric spaces. Completeness, Banach spaces and Banach's fixed point theorem. Picard's theorem. Linear transformations. Inner product spaces, projections, and Hilbert spaces. Orthogonal sequences and approximations. Linear functionals, dual space, and Riesz' representation theorem. Spectral theorem, Jordan canonical form, and matrix decompositions.

Contents

| | |
|-----------------------------------------------------------------------------------------------|----|
| Chapter 1. Vector spaces and linear transformations | 1 |
| 1.1. Vector spaces and linear transformations | 1 |
| 1.1.1. Spanning sets and bases | 3 |
| 1.1.2. Linear transformations | 8 |
| Chapter 2. Real numbers and its topology | 11 |
| 2.1. Real Numbers | 11 |
| 2.1.1. Notation | 11 |
| 2.1.2. Real numbers | 11 |
| 2.1.3. Topology of \mathbb{R} | 20 |
| 2.1.4. Supplementary material | 23 |
| Chapter 3. Normed spaces and innerproduct spaces | 25 |
| 3.1. Normed spaces and innerproduct spaces | 25 |
| 3.1.1. Normed spaces | 25 |
| 3.1.2. Innerproduct spaces | 29 |
| 3.1.3. Bounded operators between normed spaces | 33 |
| Chapter 4. Banach spaces and Hilbert spaces | 37 |
| 4.1. Banach spaces and Hilbert spaces | 37 |
| 4.1.1. Completeness | 37 |
| 4.1.2. Equivalent norms | 42 |
| 4.1.3. Banach's Fixed Point Theorem aka Contraction Mapping Theorem | 44 |
| 4.1.4. Hilbert spaces | 49 |
| 4.1.5. Orthonormal bases for Hilbert spaces | 60 |
| Chapter 5. Topology of normed spaces and continuity | 63 |
| 5.1. Topology of normed spaces | 63 |
| Chapter 6. Linear mappings between finite dimensional vector spaces and matrix decompositions | 69 |
| 6.1. Linear mappings between finite dimensional vector spaces | 69 |
| 6.1.1. QR Decomposition | 78 |
| 6.1.2. Singular Value Decomposition | 79 |
| 6.1.3. Pseudoinverse and least squares method | 83 |
| 6.1.4. Nilpotent operators | 85 |
| 6.1.5. Jordan Normal Form | 87 |
| 6.1.6. Minimal polynomials | 93 |
| Chapter 7. Metric spaces | 95 |
| 7.1. Metric spaces | 95 |

| | |
|-----------------------------------------------------|-----|
| 7.1.1. Closed, open sets and complete metric spaces | 96 |
| Appendix A. Sets and functions | 99 |
| A.1. Sets and functions | 99 |
| Bibliography | 103 |

Vector spaces and linear transformations

1.1. Vector spaces and linear transformations

Vector spaces and linear mappings between them are a useful tool for engineers, scientists and mathematicians, aka Linear Algebra. In this chapter we review material from linear algebra.

We restrict our discussion to complex and real vector spaces, but many results in this section are true for general vector spaces.

Vector spaces formalize the notion of linear combinations of objects that might be vectors in the plane, polynomials, smooth functions, sequences. Many problems in engineering, mathematics and science are naturally formulated and solved in this setting due to their linear nature. Vector spaces are ubiquitous for several reasons, e.g. as linear approximation of a non-linear object, or as building blocks for more complicated notions, such as vector bundles over topological spaces.

A set V is a vector space if it is possible to build linear combinations out of the elements in V . More formally, on V we have the operations of addition of vectors and multiplication by scalars. The scalars will be taken from a field \mathbb{F} , which is either the real numbers \mathbb{R} or \mathbb{C} . In various situations \mathbb{F} might also be a finite field or a field different from \mathbb{R} and \mathbb{C} . If it is necessary we will refer to these vector spaces as real or complex vector spaces.

Developing an understanding of these vector spaces is one of the main objectives of this course. The axioms for a vector space specify the properties that addition of vectors and scalar multiplication.

DEFINITION 1.1.1. A *vector space* over a field \mathbb{F} is a set V together with the operations of addition $V \times V \rightarrow V$ and scalar multiplication $\mathbb{F} \times V \rightarrow V$ satisfying the following properties:

- (1) Commutativity: $u + v = v + u$ for all $u, v \in V$ and $(\lambda\mu)v = \lambda(\mu v)$ for all $\lambda, \mu \in \mathbb{F}$;
- (2) Associativity: $(u + v) + w = u + (v + w)$ for all $u, v, w \in V$;
- (3) Additive identity: There exists an element $0 \in V$ such that $0 + v = v$ for all $v \in V$;
- (4) Additive inverse: For every $v \in V$, there exists an element $w \in V$ such that $v + w = 0$;
- (5) Multiplicative identity: $1v = v$ for all $v \in V$;
- (6) Distributivity: $\lambda(u + v) = \lambda u + \lambda v$ and $(\lambda + \mu)u = \lambda u + \mu u$ for all $u, v \in V$ and $\lambda, \mu \in \mathbb{F}$.

The elements of a vector space are called vectors. Given v_1, \dots, v_n be in V and $\lambda_1, \dots, \lambda_n \in \mathbb{F}$ we call the vector

$$v = \lambda_1 v_1 + \dots + \lambda_n v_n$$

a *linear combination*.

Our focus will be on three classes of examples.

EXAMPLES 1.1.2. We define some useful vector spaces.

- **Spaces of n -tuples:** The set of tuples (x_1, \dots, x_n) of real and complex numbers are vector spaces \mathbb{R}^n and \mathbb{C}^n with respect to component-wise addition and scalar multiplication: $(x_1, \dots, x_n) + (y_1, \dots, y_n) = (x_1 + y_1, \dots, x_n + y_n)$ and $\lambda(x_1, \dots, x_n) = (\lambda x_1, \dots, \lambda x_n)$.
- The space of polynomials of degree at most n , denoted by \mathcal{P}_n , where we define the operations of multiplication and addition coefficient-wise: For $p(x) = a_0 + a_1 x + \dots + a_n x^n$ and $q(x) = b_0 + b_1 x + \dots + b_n x^n$ we define

$$(p+q)(x) = (a_0+b_0) + (a_1+b_1)x + \dots + (a_n+b_n)x^n \text{ and } (\lambda p)(x) = \lambda a_0 + \lambda a_1 x + \dots + \lambda a_n x^n$$

for $\lambda \in \mathbb{F}$.

The space of all polynomials \mathcal{P} is the vector space of polynomials of arbitrary degrees.

- **Sequence spaces:** s denotes the set of sequences, c the set of all convergent sequences, c_0 the set of all convergent sequences tending to 0, c_f the set of all sequences with finitely many non-zero elements.
- **Function spaces:** The set of continuous functions $C(I)$ on an interval of \mathbb{R} , popular choices for I are $[0, 1]$ and \mathbb{R} . We define addition and scalar multiplication as follows: For $f, g \in C(I)$ and $\lambda \in \mathbb{F}$

$$(f + g)(x) = f(x) + g(x) \quad \text{and} \quad (\lambda f)(x) = \lambda f(x).$$

We denote by $C^{(n)}(I)$ the space of n -times continuously differentiable functions on I and the space $C^\infty(I)$ of smooth functions on I is the space of functions with infinitely many continuous derivatives. More generally, the set $\mathcal{F}(X)$ of functions from a set X to \mathbb{F} is a vector space for the operations defined above. Note that $\mathcal{F}(\{1, 2, \dots, n\})$ is just \mathbb{F}^n and hence the first class of examples.

- **Spaces of matrices:** Denote by $\mathcal{M}_{m \times n}(\mathbb{C})$ the space of complex $m \times n$ matrices. The vector space $\mathcal{M}_{m \times n}(\mathbb{C})$ is isomorphic to \mathbb{C}^{mn} .

There are relations between the vector spaces in the aforementioned list. We start with clarifying their inclusion properties.

DEFINITION 1.1.3. A subset W of a vector space V is called a *subspace* if any linear combination of vectors of W is itself a vector in W .

If W is a subspace of V , then addition and scalar multiplication restricted to W , gives W the structure of a vector space.

Here are some examples of vector subspaces: $\mathcal{P}_n \subset \mathcal{P} \subset \mathcal{F}$, $C^\infty(I) \subset C^{(n)}(I) \subset C(I)$, $c_f \subset c_0 \subset c \subset s$. We define the linear span, $\text{span}W$, of a subset M of a vector space V to be the intersection of all subspaces of V containing M .

1.1.1. Spanning sets and bases. Let X be a complex vector space. Recall that a linear combination of vectors x_1, \dots, x_n in X is a vector $x \in X$ of the form

$$x = \alpha_1 x_1 + \alpha_2 x_2 + \cdots + \alpha_n x_n$$

for some scalars $\alpha_1, \dots, \alpha_n \in \mathbb{C}$.

The set of all possible linear combinations of the vectors x_1, \dots, x_n in X is called the *span* of x_1, \dots, x_n , denoted by $\text{span}\{x_1, \dots, x_n\}$.

Recall that a set of vectors $\{x_1, \dots, x_n\} \subset X$ is *linearly independent* if for all $\alpha_1, \dots, \alpha_n$ the equation

$$\alpha_1 x_1 + \cdots + \alpha_n x_n = 0$$

has only $\alpha_1 = \cdots = \alpha_n = 0$ as solution. If there exists a non-trivial linear combination of the x_i 's, then we call the $\{x_1, \dots, x_n\}$ *linearly dependent*.

We often will denote the set of vectors by S and call it linearly independent without explicitly specifying the vectors.

Here are a few elementary observations about linear independence.

LEMMA 1.1. $\{x_1, \dots, x_n\} \subset X$ is linearly dependent if and only if there exists a vector, e.g. x_j , that is a linear combination of the others, i.e.

$$\text{span}\{x_1, \dots, x_j, \dots, x_n\} = \text{span}\{x_1, \dots, x_{j-1}, x_{j+1}, \dots, x_n\}$$

EXAMPLE 1.1.4. $\{1, \cos x, \sin x\}$ is linearly independent in $C(\mathbb{R})$ and $\{1, \cos x, \sin x, \cos^2 x, \sin^2 x\}$ is linearly dependent in $C(\mathbb{R})$.

LEMMA 1.2. $\{x_1, \dots, x_n\} \subset X$ is linearly independent if and only if every $x \in \text{span}\{x_1, \dots, x_n\}$ can be written uniquely as a linear combination of elements of $\{x_1, \dots, x_n\}$.

PROOF. (\Rightarrow) Assume $\{x_1, \dots, x_n\}$ is linearly independent. Suppose there are two ways to express x :

$$x = \alpha_1 x_1 + \cdots + \alpha_n x_n$$

$$x = \alpha'_1 x_1 + \cdots + \alpha'_n x_n.$$

Then we have

$$0 = (\alpha_1 - \alpha'_1)x_1 + \cdots + (\alpha_n - \alpha'_n)x_n.$$

By linear independence all these scalars have to be zero, hence the representation is unique. Contradicting our assumption.

(\Leftarrow) Suppose every $x \in \text{span}\{x_1, \dots, x_n\}$ can be written uniquely as a linear combination of elements of $\{x_1, \dots, x_n\}$. Hence there exist unique scalars $\alpha_1, \dots, \alpha_n$ for every $x \in \text{span}\{x_1, \dots, x_n\}$ such that

$$x = \alpha_1 x_1 + \cdots + \alpha_n x_n.$$

In particular $x = 0$ is uniquely represented, hence the trivial decomposition $\alpha_1 = \cdots = \alpha_n = 0$ is the only way to represent the zero vector. Hence the set $\{x_1, \dots, x_n\}$ is linearly independent. \square

PROPOSITION 1.1.5 (Linear Dependence Lemma). *Suppose $\{x_1, \dots, x_n\}$ in X is linearly dependent and assume with out loss of generality that $x_1 \neq 0$. Then there exists a vector x_j for some $j \in \{2, \dots, n\}$ such that the following holds:*

- (1) $x_j \in \text{span}\{x_1, \dots, x_{j-1}\}$,
- (2) $\text{span}\{x_1, \dots, x_{j-1}, x_{j+1}, \dots, x_n\} = \text{span}\{x_1, \dots, x_n\}$.

There are two central notions in the theory of vector spaces:

DEFINITION 1.1.6. Let X be a vector space.

- (1) If there exists a set $S \subseteq X$ with $\text{span}(S) = X$, then we call S a *spanning set*. In case that S consists of finitely many elements $\{x_1, \dots, x_n\}$, then we say that X is *finite-dimensional*. Finally, if there exists no finite spanning set for X , then we call the vector space *infinite-dimensional*.
- (2) If there exists a linearly independent spanning set B for X , then we call B a *basis* for X .

EXAMPLE 1.1.7. (1) The space of polynomials of degree at most n is finite-dimensional, because the set of monomials $\{1, x, x^2, \dots, x^n\}$ is a spanning set and even a basis for \mathcal{P}_n .

- (2) The space of all polynomials \mathcal{P} is infinite dimensional.

Let us present the argument for this fact. We have to show that for any n there is only just the trivial linear combination of monomials $\{x_0(t) = 1, x_1(t) = t^2, \dots, x_n(t) = t^n\}$ that represents the zero function. We use induction: For $n = 0$ we have $\alpha_0 = 0$ if and only if $\alpha = 0$.

Suppose for n we know that

$$\alpha_0 x_0(t) + \dots + \alpha_n x_n(t) = 0 \quad \text{for all } t \in \mathbb{R}$$

only holds for $\alpha_0 = \alpha_1 = \dots = \alpha_n = 0$. Then we want to show that this is also true for $n + 1$. We reduce the latter case to the case n by differentiation. Suppose that

$$f(t) = \alpha_0 x_0(t) + \dots + \alpha_n x_n(t) + a_{n+1} x_{n+1}(t) = 0 \quad \text{for all } t \in \mathbb{R}.$$

Then

$$f'(t) = \alpha_1 t + \dots + n\alpha_n t^{n-1} + (n+1)a_{n+1} t^n = 0 \quad \text{for all } t \in \mathbb{R}.$$

Now the induction hypothesis implies that $\alpha_1 = \dots = \alpha_n = 0$ and by the induction base we get $a_0 = 0$. Hence $f(t)$ is identically zero. Hence the set of monomials is a linearly independent set of \mathcal{P} and it spans the space of polynomials by definition. Hence it is even a basis of infinite cardinality.

- (3) The space of continuous functions on the real-line, or the space of continuously differentiable function, or the space of infinitely often differentiable functions are infinite-dimensional vector spaces.

PROPOSITION 1.1.8 (Basis Reduction Theorem). *If $\{x_1, \dots, x_n\}$ is a spanning set for X , then either $\{x_1, \dots, x_n\}$ is a basis for X or some x_j 's can be removed from $\{x_1, \dots, x_n\}$ to obtain a basis.*

As a consequence we get that every finite-dimensional vector space has a basis.

PROPOSITION 1.1.9. *Every finite-dimensional vector space has a basis.*

An often used result is the following one:

PROPOSITION 1.1.10 (Basis Extension Theorem). *Let X be a finite-dimensional vector space. Then any linearly independent subset of X can be extended to a basis.*

PROPOSITION 1.1.11 (Exchange Lemma). *Suppose $\{x_1, \dots, x_m\}$ and $\{y_1, \dots, y_n\}$ are two bases for X . Then for each $i \in \{1, \dots, m\}$ there exists some $j \in \{1, \dots, n\}$ such that $\{x_1, \dots, x_{j-1}, y_j, x_{j+1}, \dots, x_m\}$ is a basis for X .*

COROLLARY 1.1.12. *Any two bases of a finite-dimensional vector space have the same number of elements.*

LEMMA 1.3. *Let X be a finite-dimensional vector space of dimension n . Then any set $\{x_1, \dots, x_n\}$ of n linearly independent vectors is a basis of X . In other words, any set of vectors $\{x_1, \dots, x_m\}$ with $m > n$ is linearly dependent.*

These observations motivate

DEFINITION 1.1.13. Suppose X has a basis $\{x_1, \dots, x_n\}$. Then we call the number of elements of this basis the *dimension* of X , denoted by $\dim(X)$. If X is infinite-dimensional, then we write $\dim(X) = \infty$.

EXAMPLE 1.1.14. $\dim(\mathbb{C}^n) = n$, $\dim(\mathcal{P}_n) = n + 1$ and $\dim(\mathcal{P}) = \infty$.

EXAMPLE 1.1.15. Consider the vector space \mathcal{P}_2 of polynomials of degree ≤ 2 . For each of the following sets, determine if it is linearly independent in \mathcal{P}_2 , if it spans \mathcal{P}_2 and if it is a basis in \mathcal{P}_2 :

- (1) $\{1 - x, 1 + x, x^2\}$.
- (2) $\{1 + x, 1 + x^2, x - x^2\}$.

(i) Consider the vector equation

$$\begin{aligned} c_1(1 - x) + c_2(1 + x) + c_3x^2 &= 0 \\ (c_1 + c_2) + (-c_1 + c_2)x + c_3x^2 &= 0. \end{aligned}$$

We then get

$$c_1 + c_2 = 0, -c_1 + c_2 = 0, c_3 = 0,$$

which implies that

$$c_1 = 0, c_2 = 0, c_3 = 0,$$

hence the set $\{1 - x, 1 + x, x^2\}$ is linearly independent in \mathcal{P}_2 .

Moreover, being a set with 3 vectors in a vector space of dimension 3, this set must be a basis, hence also span \mathcal{P}_2 .

(ii) Note that $(1 + x) - (1 + x^2) = x - x^2$, which shows that $\{1 + x, 1 + x^2, x - x^2\}$ is *not* linearly independent in \mathcal{P}_2 , hence it cannot be a basis either.

Moreover,

$$\text{span}\{1 + x, 1 + x^2, x - x^2\} = \text{span}\{1 + x, 1 + x^2\},$$

which is a subspace of dimension 2, so it cannot be the whole space \mathcal{P}_2 , which has dimension 3. This shows that $\{1 + x, 1 + x^2, x - x^2\}$ does not span \mathcal{P}_2 .

LEMMA 1.4. *Let \mathcal{P}_4 be the vector space of real polynomials of degree at most 4.*

a) *Show that the sets $U, V \subset \mathcal{P}_4$ defined by*

$$\begin{aligned} U &:= \{p \in \mathcal{P}_4 : p(-1) = p(1) = 0\}, \\ V &:= \{p \in \mathcal{P}_4 : p(1) = p(2) = p(3) = 0\} \end{aligned}$$

are subspaces of \mathcal{P}_4 .

- b) Determine the subspace $U \cap V$.
 c) Describe bases for U , V and $U \cap V$.

PROOF. a) We show that U is a subspace of \mathcal{P}_4 .

Let $p_1, \dots, p_n \in U$ and $\lambda_1, \dots, \lambda_n \in \mathbb{R}$.

Then $p_k(-1) = p_k(1) = 0$ for all $k = 1, \dots, n$. $k = 1, \dots, n$.

Consider the linear combination $p = \lambda_1 p_1 + \dots + \lambda_n p_n$. Then clearly

$$p(-1) = \lambda_1 p_1(-1) + \dots + \lambda_n p_n(-1) = \lambda_1 \cdot 0 + \dots + \lambda_n \cdot 0 = 0,$$

which shows that $p(-1) = 0$. Similarly, $p(1) = 0$.

Therefore, $p \in U$, so U is a subspace of \mathcal{P}_4 .

The same kind of argument shows that V is a subspace.

b) We clearly have

$$U \cap V = \{p \in \mathcal{P}_4 : p(-1) = p(1) = p(2) = p(3) = 0\}.$$

This is the set of all real polynomials of degree at most 4 with exactly 4 roots: $-1, 1, 2, 3$.

Let $p_0 := (x+1)(x-1)(x-2)(x-3)$.

Then $U \cap V = \{\lambda p_0 : \lambda \in \mathbb{R}\}$.

c) A basis in $U \cap V$ is clearly $\{p_0\}$, where p_0 is the polynomial defined above.

U consists of all real polynomials that have -1 and 1 as roots, so $p \in U$ if and only if

$$p = (x+1)(x-1)q,$$

where q is a polynomial of degree at most 2.

Therefore, a basis in U is given by

$$\{(x+1)(x-1), (x+1)(x-1)x, (x+1)(x-1)x^2\}.$$

Similarly, a basis for V is given by

$$\{(x-1)(x-2)(x-3), (x-1)(x-2)(x-3)x\}.$$

□

EXAMPLE 1.1.16 (Bernstein polynomials). Let \mathcal{P}_3 be the space of polynomials of degree at most 3.

- (1) Show that $\{B_0^3(x) = (1-x)^3, B_1^3(x) = 3x(1-x)^2, B_2^3(x) = 3x^2(1-x), B_3^3(x) = x^3\}$ is a basis for \mathcal{P}_3 , known as the Bernstein basis. Since $\{B_i^3(x) : i = 0, \dots, 3\}$ is a basis of \mathcal{P}_3 there exist unique coefficients $\alpha_0, \dots, \alpha_3$ for any $f \in \mathcal{P}_3$ such that

$$f(x) = \alpha_0 B_0^3(x) + \alpha_1 B_1^3(x) + \alpha_2 B_2^3(x) + \alpha_3 B_3^3(x).$$

- (2) On the other hand we have

$$f(x) = a_0 + a_1 x + a_2 x^2 + a_3 x^3.$$

Express α_i in terms of a_i for $i = 0, \dots, 3$. In other words, how does one convert a polynomial in monomial form to one in the Bernstein basis?

PROOF. a) We start by showing that the set is linearly independent, in other words that

$$\alpha_0 B_0^3 + \alpha_1 B_1^3 + \alpha_2 B_2^3 + \alpha_3 B_3^3 = 0 \quad \Rightarrow \quad \alpha_i = 0 \quad \forall i = 0, \dots, 3.$$

Assume that

$$\alpha_0 B_0^3 + \alpha_1 B_1^3 + \alpha_2 B_2^3 + \alpha_3 B_3^3 = 0.$$

Replacing the B_i^n 's with their definition, we get

$$\alpha_0(1-x)^3 + 3\alpha_1x(1-x)^2 + 3\alpha_2x^2(1-x) + \alpha_3x^3 = 0.$$

By multiplying out the brackets and rearranging we get

$$(\alpha_3 - 3\alpha_2 + 3\alpha_1 - \alpha_0)x^3 + (3\alpha_2 - 6\alpha_1 + 3\alpha_0)x^2 + (3\alpha_1 - 3\alpha_0)x + \alpha_0 = 0,$$

which means that the coefficients are all 0:

$$\begin{aligned}\alpha_3 - 3\alpha_2 + 3\alpha_1 - \alpha_0 &= 0 \\ 3\alpha_2 - 6\alpha_1 + 3\alpha_0 &= 0 \\ 3\alpha_1 - 3\alpha_0 &= 0 \\ \alpha_0 &= 0\end{aligned}$$

By repeatedly substituting from bottom and up, we get that $a_0 = a_1 = a_2 = a_3 = 0$. Hence the set is linearly independent. To show that it is a basis, it is enough to observe that the set has 4 elements, which is the same as the dimension of the vector space \mathcal{P}_3 .

The Bernstein basis for \mathcal{P}_n is given by $\{B_i^n\}_{i=1}^n$, where

$$B_i^n(x) = \binom{n}{i} x^i (1-x)^{n-i}.$$

b) We have

$$\alpha_0 B_0^3 + \alpha_1 B_1^3 + \alpha_2 B_2^3 + \alpha_3 B_3^3 = a_3 x^3 + a_2 x^2 + a_1 x + a_0.$$

By rearranging the left-hand side as we did in a) we get

$$(\alpha_3 - 3\alpha_2 + 3\alpha_1 - \alpha_0)x^3 + (3\alpha_2 - 6\alpha_1 + 3\alpha_0)x^2 + (3\alpha_1 - 3\alpha_0)x + \alpha_0 = a_3 x^3 + a_2 x^2 + a_1 x + a_0.$$

Since the coefficients have to match up, we get

$$\begin{aligned}a_0 &= \alpha_0 \\ a_1 &= 3\alpha_1 - 3\alpha_0 \\ a_2 &= 3\alpha_2 - 6\alpha_1 + 3\alpha_0 \\ a_3 &= \alpha_3 - 3\alpha_2 + 3\alpha_1 - \alpha_0,\end{aligned}$$

and after solving for the α_i 's we get

$$\begin{aligned}\alpha_0 &= a_0 \\ \alpha_1 &= \frac{a_1}{3} + a_0 \\ \alpha_2 &= \frac{a_2}{3} + \frac{2a_1}{3} + a_0 \\ \alpha_3 &= a_3 + a_2 + a_1 + a_0\end{aligned}$$

□

PROPOSITION 1.1.17. *Let M, N be subspaces of a finite-dimensional vector space X . Then*

$$\dim(M + N) + \dim(M \cap N) = \dim(M) + \dim(N).$$

1.1.2. Linear transformations. Let T be a linear transformation from X to Y . Then the kernel of T is

$$\ker(T) = \{x \in X : Tx = 0\}$$

and the range of T is

$$\operatorname{ran}(T) = \{y \in Y : y = Tx \text{ for some } x \in X\}.$$

The $\ker(T)$ is a subspace of X and the $\operatorname{ran}(T)$ is a subspace of Y . Suppose X and Y are finite dimensional vector spaces. Then one can construct bases for $\ker(T)$ and $\operatorname{ran}(T)$. We call the dimension of the $\ker(T)$ the *nullity* of T and the dimension of $\operatorname{ran}(T)$ the *rank* of T .

PROPOSITION 1.1.18. *Let X and Y be finite dimensional vector spaces. For a linear mapping $T : X \rightarrow Y$ we have*

$$\dim(X) = \dim(\ker(T)) + \dim(\operatorname{ran}(T)).$$

PROOF. Idea is to use the dimension formula for the sum of vector spaces. Let V be a n -dimensional vector space. Suppose $\{x_1, \dots, x_k\}$ is a basis for $\ker(T)$. Then there exist x_{k+1}, \dots, x_n in X such that $\{x_1, \dots, x_k, \dots, x_n\}$ is a basis for X . We denote by $S = \operatorname{span}\{x_{k+1}, \dots, x_n\}$. Then by construction we have

$$\ker(T) \cap S = \{0\}$$

and by the dimension formula for subspaces we have

$$\dim(X) = \dim(\ker(T)) + \dim(\operatorname{ran}(T)).$$

Note that $\operatorname{ran}(T) = T(S)$ and the restriction of T to S is injective. Hence $\dim(\operatorname{ran}(T(S))) = \dim(S) = \dim(\operatorname{ran}(T))$. Thus we have the desired assertion. \square

We associate two linear mappings to a basis $\mathcal{B} = \{x_1, \dots, x_n\}$ of a finite-dimensional vector space. Then each x can be uniquely expressed as

$$x = \alpha_1 x_1 + \cdots + \alpha_n x_n$$

and we define the *coefficient* map $C : X \rightarrow \mathbb{C}^n$ by

$$Cx = \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{pmatrix}$$

often denoted by $Cx = [x]_{\mathcal{B}}$, and the *synthesis* map $D : \mathbb{C}^n \rightarrow X$ by

$$x = \alpha_1 x_1 + \cdots + \alpha_n x_n.$$

Next we discuss the link between matrices and linear transformations. On the one hand a $m \times n$ matrix A defines a linear transformation from \mathbb{C}^n to \mathbb{C}^m by $Tx = Ax$.

On the other hand any linear transformation on finite-dimensional vector spaces can be represented in matrix form relative to a choice of bases.

We present the details for this assertion. Let $\mathcal{B} = \{x_1, \dots, x_n\}$ be a basis of X and

$\mathcal{C} = \{y_1, \dots, y_m\}$ be a basis of Y . Suppose T is a linear transformation $T : X \rightarrow Y$. Then

$$x = \sum_{i=1}^n \alpha_i x_i$$

yields

$$T(x) = \sum_{i=1}^n \alpha_i T(x_i)$$

and thus

$$[T(x)]_{\mathcal{C}} = \sum_{i=1}^n \alpha_i [T(x_i)]_{\mathcal{C}}.$$

We define a $m \times n$ matrix A which has as its j -th column $[[T(x_j)]_{\mathcal{C}}]$. Then we have

$$[Tx]_{\mathcal{C}} = A[x]_{\mathcal{B}}.$$

The matrix A represents T with respect to the bases \mathcal{B} and \mathcal{C} . Sometimes, we denote this A sometimes by $[T]_{\mathcal{B}}^{\mathcal{C}}$.

We address now the relation between the matrix representation of T depending on the change of bases. Suppose we have two bases $\mathcal{B} = \{x_1, \dots, x_n\}$ and $\mathcal{R} = \{y_1, \dots, y_n\}$ for X . Let $x = \sum_{j=1}^n \alpha_j x_j$. Then

$$[x]_{\mathcal{R}} = \sum_{j=1}^n \alpha_j \vec{x}_j_{\mathcal{R}}.$$

Define the $n \times n$ matrix P with j -th column $\vec{x}_j_{\mathcal{R}}$, and we call P the *change of bases matrix*:

$$[x]_{\mathcal{R}} = P[x]_{\mathcal{B}}$$

and by the invertibility of P we also have

$$[x]_{\mathcal{B}} = P^{-1}[x]_{\mathcal{R}}.$$

Let now \mathcal{C} and \mathcal{S} be two bases for Y . Then a linear transformation $T : X \rightarrow Y$ has two matrix representations:

$$A = [T]_{\mathcal{B}}^{\mathcal{C}} \text{ and } B = [T]_{\mathcal{R}}^{\mathcal{S}}.$$

In other words we have

$$[Tx]_{\mathcal{C}} = A[x]_{\mathcal{B}} \quad , \quad [Tx]_{\mathcal{S}} = B[x]_{\mathcal{R}}$$

for any $x \in X$. Let P be the change of bases matrix of size $n \times n$ such that $[x]_{\mathcal{R}} = P[x]_{\mathcal{B}}$ for any $x \in X$ and let Q be the invertible $m \times m$ matrix such that $[y]_{\mathcal{S}} = Q[y]_{\mathcal{C}}$.

Hence we get that

$$[Tx]_{\mathcal{S}} = BP[x]_{\mathcal{B}}$$

and

$$[y]_{\mathcal{S}} = [Tx]_{\mathcal{S}} = Q[Tx]_{\mathcal{C}} = QA[x]_{\mathcal{B}}$$

for any $x \in X$. Hence we get that

$$B = QAP^{-1} \text{ and } A = Q^{-1}BP.$$

In the case $X = Y$ we have $P = Q$ and we set $S = Q^{-1}$ to get $B = S^{-1}AS$. Then the matrices A and B represent the same linear transformation T on V with respect to different bases.

These observations motivate the definition.

DEFINITION 1.1.19. Two $m \times n$ matrices A and B are called *equivalent* if there exists an invertible matrix S such that $B = QAP^{-1}$. Furthermore, Two $n \times n$ matrices A and B are called *equivalent* if there exists an invertible matrix S such that $B = S^{-1}AS$.

Given a general $n \times n$ matrix A . Two similar matrices are “essentially the same”. The notion of similarity is of utmost importance for linear algebra. It allows one to classify matrices. We are going to show that it is possible that any matrix is similar to an upper triangular matrix, Schur’s theorem, and with more effort to get into a special upper triangular form, the Jordan normal form. Of special interest are matrices that are similar to diagonal matrices, which will turn out to be the normal matrices. The final statement is often referred to as “spectral theorem”.

For a matrix $A = (a_{ij})$ we define its *trace* to be the sum of its diagonal elements:

$$\operatorname{tr}(A) = a_{11} + \cdots + a_{nn}.$$

Real numbers and its topology

2.1. Real Numbers

2.1.1. Notation. We introduce some notation:

- (1) $\mathbb{N} = \{1, 2, 3, \dots\}$ the set of natural numbers,
- (2) $\mathbb{Q} = \{p/q : p, q \in \mathbb{Z}\}$ the set of rational numbers,
- (3) $\mathbb{Z} = \{\dots, -2, -1, 0, 1, 2, \dots\}$ the set of integers.
- (4) For real numbers a, b with $a < b$ we denote by $[a, b]$ the closed bounded interval, and by (a, b) the open bounded interval. The length of these bounded intervals is $b - a$.

2.1.2. Real numbers. The set \mathbb{Q} of rational numbers does not contain all the numbers one encounters in geometry or analysis, e.g. $x^2 - 5 = 0$ has no rational solution or Euler's number e is an irrational number.

PROPOSITION 2.1.1. *The equation*

$$x^2 - 3 = 0$$

has no solutions in \mathbb{Q} .

PROOF. We assume by contradiction that there is a rational number r such that $r^2 - 3 = 0$.

We represent r as a *reduced* fraction. That is, we write $r = \frac{p}{q}$ where p, q are integers, $q \neq 0$ and $\gcd(p, q) = 1$. We then have:

$$r^2 - 3 = 0 \implies r^2 = 3 \implies \frac{p^2}{q^2} = 3 \implies p^2 = 3q^2.$$

The last identity says that p^2 is a multiple of 3. Then p itself must be a multiple of 3 as well (why?), which means that $p = 3m$ for some integer m .

Substituting this into the identity $p^2 = 3q^2$ we get $9m^2 = 3q^2$, which implies $3m^2 = q^2$, and so q^2 must be a multiple of 3. But then q must also be a multiple of 3.

Let us step back and look at what we have: we started of with a completely reduced fraction $r = \frac{p}{q}$, assumed that $r^2 - 3 = 0$, which through a series of derivations led to the conclusion that both p and q must be multiples of 3. This contradicts the fraction $\frac{p}{q}$ being reduced.

Therefore, the equation $x^2 - 3 = 0$ cannot have any rational number as solution. \square

For the moment we do not introduce the set of real number \mathbb{R} in an informal manner. In the theory of metric spaces \mathbb{R} is constructed as the completion of \mathbb{Q} , as was originally done by A. L. Cauchy.

Real numbers may be realized as points on a line, the real line, where the irrational numbers correspond to the points that are not given by rational numbers $\mathbb{R} \setminus \mathbb{Q}$.

The real numbers have the Archimedean property:

LEMMA 2.1 (Archimedean property). *For any $x, y \in \mathbb{R}$ there exists a natural number n such that $nx > y$.*

As a consequence we deduce a close relation between \mathbb{Q} and \mathbb{R} .

PROPOSITION 2.1.2. *For $x, y \in \mathbb{R}$ with $x < y$ there exists a $r \in \mathbb{Q}$ such that $x < r < y$.*

PROOF. Goal: Find $m, n \in \mathbb{Z}$ such that

$$(2.1) \quad x < \frac{m}{n} < y.$$

First step: Choose the denominator of n large such that there exists an $m \in \mathbb{Z}$ such that $x \in (\frac{m-1}{n}, \frac{m}{n})$ are separating x and y . The Archimedean property of \mathbb{R} allows us to a $n \in \mathbb{N}$ with this property. More concretely, we pick $n \in \mathbb{N}$ large enough such that $1/n < y - x$ or equivalently

$$(2.2) \quad x < y - \frac{1}{n}$$

Second step: Inequality (2.1) is equivalent to $nx < m < ny$. From the first step we have n already chosen. Now we choose $m \in \mathbb{Z}$ to be the smallest integer greater than nx . In other words, we pick $m \in \mathbb{Z}$ such that $m - 1 \leq nx < m$. Thus we have $m - 1 \leq nx$, i.e. $m \leq nx + 1$. By inequality (2.2)

$$m \leq nx + 1 < n(y - \frac{1}{n}) + 1 = ny,$$

hence we have $m < ny$, i.e. $m/n < y$. Once more by (2.2) we have $x \leq m/n$. These two inequalities yield the desired assertion: $x < m/n < y$. \square

In an similar manner one may deduce the statement for irrational numbers.

PROPOSITION 2.1.3. *For $x, y \in \mathbb{R}$ with $x < y$ there exists a $r \in \mathbb{R} \setminus \mathbb{Q}$ such that $x < r < y$.*

PROOF. Pick your favorite irrational number, a popular choice is $\sqrt{2}$. Then by the density of the rational numbers there exists a rational number $r \in (x/\sqrt{2}, y/\sqrt{2})$. Hence $r\sqrt{2} \in (x, y)$. Note that $r\sqrt{2}$ is an irrational number in (x, y) that completes our argument. \square

The absolute value of $x \in \mathbb{R}$, denoted by $|x|$, is defined by

$$|x| = \begin{cases} -x & \text{if } x < 0, \\ 0 & \text{if } x = 0, \\ x & \text{if } x > 0. \end{cases}$$

Note that $|x| = \max\{x, -x\}$. We define the positive, x^+ and negative part, x^- of $x \in \mathbb{R}$:

$$x^+ = \max\{x, 0\}, \quad \text{and} \quad x^- = \max\{-x, 0\},$$

so we have $x = x^+ - x^-$ and $|x| = x^+ + x^-$.

For $x, y \in \mathbb{R}$ we measure the distance between x and y in \mathbb{R} by

$$(2.3) \quad d(x, y) = |x - y|,$$

the standard distance. By definition of d we have $d(x, y) = d(y, x)$.

LEMMA 2.2 (Triangle inequality). For x, y in \mathbb{R} we have $|x + y| \leq |x| + |y|$.

PROOF. For all $x \in \mathbb{R}$ we have $x \leq |x|$ and thus for $x, y \in \mathbb{R}$ we obtain $x + y \leq |x| + |y|$. By definition of $|\cdot|$ we also get that $-x - y \leq |x| + |y|$. Thus we have proved the desired assertion. \square

The triangle inequality has numerous consequences, such as

$$(2.4) \quad ||x| - |y|| \leq |x - y|.$$

The triangle inequality for $x = y + x - y$ yields $|x| - |y| \leq |x - y|$, and the interchange of x and y , i.e. $y = x + y - x$ gives $-(|x| - |y|) \leq |x - y|$. Hence we have the desired assertion.

We introduce two crucial notions: the infimum and supremum of a set. First we provide some preliminaries.

DEFINITION 2.1.4. Let A be a subset of \mathbb{R}

- If there exists $M \in \mathbb{R}$ such that $a \leq M$ for all $a \in A$, then M is an *upper bound* of A . We call A *bounded above*.
- If there exists $m \in \mathbb{R}$ such that $m \leq a$ for all $a \in A$, then m is a *lower bound* of A .
- If there exist lower and upper bounds, then we say that A is *bounded*. We call A *bounded below*.

DEFINITION 2.1.5 (Infimum and Supremum). Let A be a subset of \mathbb{R} .

- If m is a lower bound of A such that $m \geq m'$ for every lower bound m' , then m is called the *infimum* of A , denoted by $m = \inf A$. Furthermore, if $\inf A \in A$, then we call it the *minimum* of A , $\min A$.
- If M is an upper bound of A such that $m' \geq M$ for every upper bound M' , then M is called the *supremum* of A , denoted by $M = \sup A$. Furthermore, if $\sup A \in A$, then we call it the *maximum* of A , $\max A$.

Note that the infimum of a set A , as well as the supremum, are unique. The elementary argument is left as an exercise.

If $A \subset \mathbb{R}$ is not bounded above, then we define $\sup A = \infty$. Suppose that a subset A of \mathbb{R} is not bounded below, then we assign $-\infty$ as its infimum.

We state a different formulation of the notions $\inf A$ and $\sup A$ that is just a reformulation of the definition.

LEMMA 2.3. Let A be a subset of \mathbb{R} .

- Suppose A is bounded above. Then $M \in \mathbb{R}$ is the supremum of A if and only if the following two conditions are satisfied:
 - (1) For every $a \in A$ we have $a \leq M$.
 - (2) Given $\varepsilon > 0$, there exists $a \in A$ such that $M - \varepsilon < a$.
- Suppose A is bounded below. Then $m \in \mathbb{R}$ is the infimum of A if and only if the following two conditions are satisfied:
 - (1) For every $a \in A$ we have $m \leq a$.

- (2) Given $\varepsilon > 0$, there exists $a \in A$ such that $a < m + \varepsilon$.

LEMMA 2.4. Suppose A is a bounded subset of \mathbb{R} . Then $\inf A \leq \sup A$

For $c \in \mathbb{R}$ we define the *dilate* of a set A by $cA := \{b \in \mathbb{R} : b = ca \text{ for } a \in A\}$.

LEMMA 2.5 (Properties). Suppose A is a subset of \mathbb{R} .

- (1) For $c > 0$ we have $\sup cA = c \sup A$ and $\inf cA = c \inf A$.
- (2) For $c < 0$ we have $\sup cA = c \inf A$ and $\inf cA = c \sup A$.
- (3) Suppose A is contained in a subset B . If $\sup A$ and $\sup B$ exist, then $\sup A \leq \sup B$. In words, making a set larger, increases its supremum.
- (4) Suppose A is contained in a subset B . If $\inf A$ and $\inf B$ exist, then $\inf A \geq \inf B$. In words, making a set smaller increases its infimum.
- (5) Suppose $A \subset B$ are non-empty subsets of \mathbb{R} such that $x \leq y$ for all $x \in A$ and $y \in B$. Then $\sup A \leq \inf B$.
- (6) If A and B are non-empty subsets of \mathbb{R} , then $\sup(A + B) = \sup A + \sup B$ and $\inf(A + B) = \inf A + \inf B$

PROOF. (1) We prove that $\sup cA = c \sup A$ for positive c . Suppose $c > 0$. Then $cx \leq M \Leftrightarrow x \leq M/c$. Hence M is an upper bound of cA if and only if M/c is an upper bound of A . Consequently, we have the desired result.

- (2) Without loss of generality we set $c = -1$. Let $a \in A$ (we assume that the set A is non-empty, otherwise there is nothing interesting here). Then as a lower bound for A , $\inf A \leq a$. Moreover, as an upper bound for A , $a \leq \sup A$. Using transitivity, we conclude that $\inf A \leq \sup A$.

We now prove the second identity. Keep in mind that the supremum of a set is its **least upper bound**, while the infimum is its **greatest lower bound**.

For any $a \in A$, $\inf A \leq a$, so $-\inf A \geq -a$, showing that $-\inf A$ is an upper bound for $-A$. Therefore, $-\inf A \geq \sup(-A)$, which implies

$$\boxed{\inf A \leq -\sup(-A)}.$$

For any $a \in A$ we have $-a \in -A$, so $-a \leq \sup(-A)$, which implies $a \geq -\sup(-A)$. Therefore, $-\sup(-A)$ is a lower bound for A , so

$$\boxed{-\sup(-A) \leq \inf A}.$$

The two boxed inequalities prove the identity $\inf A = -\sup(-A)$.

- (3) Since $\sup B$ is an upper bound of B , it is also an upper bound of A , i.e. $\sup A \leq \sup B$.
- (4) Analogously to (iii).
- (5) Since $x \leq y$ for all $x \in A$ and $y \in B$, y is an upper bound of A . Hence $\sup A$ is a lower bound of B and we have $\sup A \leq \inf B$.
- (6) By definition $A + B = \{c : c = a + b \text{ for some } a \in A, b \in B\}$ and thus $A + B$ is bounded above if and only if A and B are bounded above. Hence $\sup(A + B) < \infty$ if and only if $\sup A$ and $\sup B$ are finite. Take $a \in A$ and $b \in B$, then $a + b \leq \sup A + \sup B$. Thus $\sup A + \sup B$ is an upper bound of $A + B$:

$$\sup(A + B) \leq \sup A + \sup B.$$

The reverse direction is a little bit more involved. Let $\varepsilon > 0$. Then there exists $a \in A$ and $b \in B$ such that

$$a > \sup A - \varepsilon/2, \quad b > \sup B - \varepsilon/2.$$

Thus we have $a + b > \sup A + \sup B - \varepsilon$ for every $\varepsilon > 0$, i.e. $\sup(A + B) \geq \sup A + \sup B$.

The other statements are assigned as exercises. \square

A property of utmost importance is the *completeness* of the real numbers.

THEOREM 2.6. *Let A be a non-empty subset of \mathbb{R} that is bounded above. Then there exists a supremum of A . Equivalently, if A is a non-empty subset of \mathbb{R} that is bounded below, then A has an infimum.*

We have noted above that the supremum of a bounded above set is unique. A different form to express the completeness property of \mathbb{R} is to consider the set of all upper bounds of a bounded above set A and the Theorem asserts that this set of upper bounds has a least element.

One reason for the relevance of the notions of supremum and infimum is in the formulation of properties of functions.

DEFINITION 2.1.6. Let f be a function with domain X and range $Y \subseteq \mathbb{R}$. Then

$$\sup_X f = \sup\{f(x) : x \in X\}, \quad \inf_X f = \inf\{f(x) : x \in X\}.$$

If $\sup_X f$ is finite, then f is bounded from above on A , and if $\inf_X f$ is finite we call f bounded from below. A function is bounded if both the supremum and infimum are finite.

LEMMA 2.7. *Suppose that $f, g : X \rightarrow \mathbb{R}$ and $f \leq g$, i.e. $f(x) \leq g(x)$ for all $x \in X$. If g is bounded from above, then $\sup_X f \leq \sup_X g$. Assume that f is bounded from below. Then $\inf_X f \leq \inf_X g$.*

PROOF. Follows from the definitions. \square

The supremum and infimum of functions do not preserve strict inequalities. Define $f, g : [0, 1] \rightarrow \mathbb{R}$ by $f(x) = x$ and $g(x) = x + 1$. Then we have $f < g$ and

$$\sup_{[0,1]} f = 1, \quad \inf_{[0,1]} f = 0, \quad \sup_{[0,1]} g = 2, \quad \inf_{[0,1]} g = 1.$$

Hence we have $\sup_{[0,1]} f > \inf_{[0,1]} g$.

LEMMA 2.8. *Suppose f, g are bounded functions from X to \mathbb{R} and c a positive constant. Then*

$$\sup_X (f + cg) \leq \sup_X f + c \sup_X g \quad \inf_X (f + cg) \geq \inf_X f + c \inf_X g.$$

The proof is left as an exercise. Try to convince yourself that the inequalities are in general strict, since the functions f and g may take values close to their suprema/infima at different points in X .

LEMMA 2.9. *Suppose f, g are bounded functions from X to \mathbb{R} . Then*

$$\left| \sup_X f - \sup_X g \right| \leq \sup_X |f - g|, \quad \left| \inf_X f - \inf_X g \right| \leq \sup_X |f - g|$$

LEMMA 2.10. Suppose f, g are bounded functions from X to \mathbb{R} such that

$$|f(x) - f(y)| \leq |g(x) - g(y)| \quad \text{for all } x, y \in X.$$

Then

$$\sup_X f - \inf_X f \leq \sup_X g - \inf_X g.$$

Recall that a sequence (x_n) of real numbers is an ordered list of numbers x_n , indexed by the natural numbers. In other words, (x_n) is a function f from \mathbb{N} to \mathbb{R} with $f(n) = x_n$. Hence we may define the if a sequence (x_n) is *bounded from above*, *bounded from below* and *bounded* as a special case of the above definitions, i.e. if there exists $M \in \mathbb{R}$ such that $x_n \leq M$ for all $n \in \mathbb{N}$, if there exists $m \in \mathbb{R}$ such that $x_n \geq m$ for all $n \in \mathbb{N}$ and if there exist m, M such that $m \leq x_n \leq M$.

We define the *lim sup* and *lim inf* of a sequence (x_n) . These notions reduce questions about the convergence of a sequence to ones about monotone sequences. We introduce two sequences associated to (x_n) by taking the supremum and infimum, respectively of the tails of $((x_k)_{k \geq n})_k$:

$$y_n = \sup\{x_k : k \geq n\}, \quad z_n = \inf\{x_k : k \geq n\}.$$

The sequences (y_n) and (z_n) are monotone sequences, because the supremum and infimum are taken over smaller sets for increasing n . Moreover, (y_n) is monotone decreasing and (z_n) is monotone increasing. Hence the limits of these sequences exist:

$$\begin{aligned} \limsup_{n \rightarrow \infty} x_n &:= \lim_{n \rightarrow \infty} y_n = \inf_{n \in \mathbb{N}} (\sup_{k \geq n} x_k), \\ \liminf_{n \rightarrow \infty} x_n &:= \lim_{n \rightarrow \infty} z_n = \sup_{n \in \mathbb{N}} (\inf_{k \geq n} x_k). \end{aligned}$$

We allow \limsup and \liminf to be $+\infty$ and $-\infty$. Note that we have $z_n \leq y_n$ and so by taking the limit as $n \rightarrow \infty$

$$\liminf_{n \rightarrow \infty} x_n \leq \limsup_{n \rightarrow \infty} x_n$$

. We illustrate these notions with some examples.

EXAMPLES 2.1.7. Consider the sequences.

- (1) $(x_n) = ((-1)^{n+1})$ has $\limsup x_n = 1$ and $\liminf x_n = -1$.
- (2) $(x_n) = (n^2)$ has $\limsup x_n = \infty$ and $\liminf x_n = \infty$.
- (3) $(x_n) = (2 - 1/n)$ has $\limsup x_n = 2$ and $\liminf x_n = 2$.

LEMMA 2.11. Let (x_n) and (y_n) be sequences in \mathbb{R} .

- (1) $\liminf(x_n + y_n) \geq \liminf x_n + \liminf y_n$,
- (2) $\limsup(x_n + y_n) \leq \limsup x_n + \limsup y_n$,
- (3) $\limsup(-x_n) = -\liminf x_n$ and $\liminf(-x_n) = -\limsup x_n$.

PROOF. (1) A sequence (or a subsequence) is a function from \mathbb{N} (or from a subset of \mathbb{N}) to \mathbb{R} , so the properties of the \inf and \sup for functions apply to sequences as well.

Then for every $n \in \mathbb{N}$ we have

$$\inf\{x_k + y_k : k \geq n\} \geq \inf\{x_k : k \geq n\} + \inf\{y_k : k \geq n\}.$$

Taking the limit on both sides of the inequality we have

$$\lim_{n \rightarrow \infty} \inf \{x_k + y_k : k \geq n\} \geq \lim_{n \rightarrow \infty} \inf \{x_k : k \geq n\} + \lim_{n \rightarrow \infty} \inf \{y_k : k \geq n\},$$

which proves that $\liminf(x_n + y_n) \geq \liminf x_n + \liminf y_n$.

(2)

$$\sup \{x_k + y_k : k \geq n\} \leq \sup \{x_k : k \geq n\} + \sup \{y_k : k \geq n\}.$$

Taking the limit on both sides of the inequality we have

$$\lim_{n \rightarrow \infty} \sup \{x_k + y_k : k \geq n\} \leq \lim_{n \rightarrow \infty} \sup \{x_k : k \geq n\} + \lim_{n \rightarrow \infty} \sup \{y_k : k \geq n\},$$

which proves that $\limsup(x_n + y_n) \leq \limsup x_n + \limsup y_n$.

(3) We established that for any bounded set A , we have $\inf A = -\sup(-A)$.

This is the same as saying

$$\sup(-A) = -\inf A.$$

When applied to the set $-A$ instead of A , this also shows that $\sup A = \sup(-(-A)) = -\inf(-A)$, so

$$\inf(-A) = -\sup A.$$

It is easy to see that these identities also hold when A is not bounded. We will apply these identities to sets related to the terms of a sequence, as follows.

$$\begin{aligned} \limsup(-x_n) &= \inf_{n \geq 1} \sup \{-x_k : k \geq n\} \\ &= \inf_{n \geq 1} (-\inf \{x_k : k \geq n\}) \\ &= -\sup_{n \geq 1} \inf \{x_k : k \geq n\} = -\liminf x_n. \end{aligned}$$

(4) Similarly,

$$\begin{aligned} \liminf(-x_n) &= \sup_{n \geq 1} \inf \{-x_k : k \geq n\} \\ &= \sup_{n \geq 1} (-\sup \{x_k : k \geq n\}) \\ &= -\inf_{n \geq 1} \sup \{x_k : k \geq n\} = -\limsup x_n. \end{aligned}$$

□

Note that for convergent sequences \limsup and \liminf are finite and equal. We recommend to prove this property.

PROPOSITION 2.1.8. *Let (x_n) be a sequence in \mathbb{R} . Then (x_n) converges if and only if $\liminf_{n \rightarrow \infty} x_n = \limsup_{n \rightarrow \infty} x_n$.*

Note that a sequence diverges to ∞ if and only if $\liminf_{n \rightarrow \infty} x_n = \limsup_{n \rightarrow \infty} x_n = \infty$ and that it diverges to $-\infty$ if and only if $\liminf_{n \rightarrow \infty} x_n = \limsup_{n \rightarrow \infty} x_n = -\infty$.

These considerations suggests that for non-convergent sequences the difference

$$\liminf_{n \rightarrow \infty} x_n - \limsup_{n \rightarrow \infty} x_n$$

measures the size of the oscillations in the sequene.

A central notion in analysis is the notion of a Cauchy sequence of objects, here we define it for real numbers.

DEFINITION 2.1.9. A sequence (x_n) in \mathbb{R} is a *Cauchy sequence* if for every $\varepsilon > 0$ there exists $N \in \mathbb{N}$ such that $|x_m - x_n| < \varepsilon$ for all $m, n \geq N$.

A theorem of utmost importance is that every Cauchy sequence converges to a real number.

THEOREM 2.12. *A sequence (x_n) converges in \mathbb{R} if and only if it is a Cauchy sequence.*

PROOF. One direction: Suppose (x_n) converges to a real number x . Then for every $\varepsilon > 0$ there exists $N \in \mathbb{N}$ such that $|x_n - x| < \varepsilon/2$ for all $n > N$. Hence by the triangle inequality we have

$$|x_n - x_m| \leq |x_n - x| + |x - x_m| \quad \text{for } m, n > N,$$

i.e. (x_n) is a Cauchy sequence.

Other direction: Suppose that (x_n) is a Cauchy sequence. Then there exists $N_1 \in \mathbb{N}$ such that $|x_m - x_n| < 1$ for all $m, n > N_1$, and that for $n > N_1$ we have

$$|x_n| \leq |x_n - x_{N_1}| + |x_{N_1}| \leq 1 + |x_{N_1}|.$$

Hence a Cauchy sequence is bounded with $|x_n| \leq \max\{|x_1|, \dots, |x_{N_1}|, 1 + |x_{N_1}|\}$ and \limsup, \liminf exist.

The aim is to show that $\limsup x_n = \liminf x_n$.

By the Cauchy property of (x_n) we have for a given $\varepsilon > 0$ a $N \in \mathbb{N}$ such that

$$x_n - \varepsilon < x_m < x_n + \varepsilon \quad \text{for all } m \geq n > N.$$

Consequently, we have for all $n > N$

$$x_n - \varepsilon \leq \inf\{x_m : m \geq n\} \quad \text{and} \quad \sup\{x_m : m \geq n\} \leq x_n + \varepsilon.$$

Thus we have

$$\sup\{x_m : m \geq n\} - \varepsilon \leq \inf\{x_m : m \geq n\} + \varepsilon$$

and for $n \rightarrow \infty$ we get that

$$\limsup x_n - \varepsilon \leq \liminf x_n + \varepsilon$$

for arbitrary $\varepsilon > 0$ and so

$$\limsup x_n \leq \liminf x_n.$$

□

In the proof we established that Cauchy sequences are bounded. Let us prove it in more detail.

LEMMA 2.13. *A Cauchy sequence (x_n) in \mathbb{R} is bounded.*

PROOF. The idea is that for a Cauchy sequence, all but finitely many of its terms are near each other, hence near (any) one of them. The remaining terms will have a maximum and a minimum, since they are finitely many. Let us formalize this idea.

Since (x_n) is a Cauchy sequence, for say $\varepsilon = 1$, there is $N \in \mathbb{N}$ such that for all $n, m \geq N$ we have $|x_n - x_m| \leq 1$.

In particular, choosing $m = N$, for all terms $n \geq N$ we have that $|x_n - x_N| \leq 1$. This implies that $-1 \leq x_n - x_N \leq 1$, from which we conclude

$$x_N - 1 \leq x_n \leq x_N + 1 \quad \text{for all } n \geq N.$$

Any *finite* set of numbers has a maximum and a minimum. Therefore, we can take

$$\begin{aligned} M &:= \max \{x_1, x_2, \dots, x_{N-1}, x_N + 1\} \text{ and} \\ m &:= \min \{x_1, x_2, \dots, x_{N-1}, x_N - 1\}. \end{aligned}$$

It is then clear that

$$m \leq x_n \leq M \quad \text{for all } n \geq 1,$$

proving that (x_n) is bounded. \square

We define the notion of a subsequence of a sequence (x_n) .

DEFINITION 2.1.10. Suppose (x_n) is a sequence in \mathbb{R} . Then a *subsequence* is a sequence of the form (x_{n_k}) where $n_1 < n_2 < \dots < x_{n_k} < \dots$.

An elementary observation is

LEMMA 2.14. *Every subsequence of a convergent sequence converges to the limit of the sequence.*

PROOF. Suppose that (x_n) is a convergent sequence with $\lim x_n = x$ and (x_{n_k}) is a subsequence. Given $\varepsilon > 0$. There exists $N \in \mathbb{N}$ such that $|x_n - x| < \varepsilon$ for all $n > N$. Since $n_k \rightarrow \infty$ as $k \rightarrow \infty$, there exists a $K \in \mathbb{N}$ such that $n_k > N$ for $k > K$, but then we have $|x_{n_k} - x| < \varepsilon$. Hence $\lim_{k \rightarrow \infty} x_{n_k} = x$. \square

COROLLARY 2.1.11. *If a sequence has subsequences that converge to different limits, then the sequence diverges.*

A well-known theorem due to Bolzano and Weierstraß deduces the convergence of a subsequence from its boundedness.

THEOREM 2.15 (Bolzano-Weierstraß). *Every bounded sequence (x_n) in \mathbb{R} has a convergent subsequence.*

PROOF. Suppose that (x_n) is a bounded sequence in \mathbb{R} . Hence there are m and M such that

$$m = \inf_n x_n \quad M = \sup_n x_n.$$

We define the closed interval $I_0 = [m, M]$ and divide it into two closed intervals L_0, R_0 :

$$L_0 = [m, (m + M)/2], \quad R_0 = [(m + M)/2, M].$$

Now, at least one of the intervals L_0, R_0 contains infinitely many terms of (x_n) . Choose I_1 to be the interval that contains infinitely many terms and pick $n_1 \in \mathbb{N}$ such that $x_{n_1} \in I_1$. Divide $I_1 = L_1 \cup R_1$, again one of these intervals contains infinitely many terms of (x_n) . Choose I_2 to be one of these intervals that contains infinitely many terms. We continue by dividing I_2 into two closed intervals, pick $n_2 > n_1$ such that $x_{n_2} \in I_2$. Continue in this manner we get a sequence of nested intervals (I_k) with $|I_k| = (M - m)/2^k$, and a sequence (x_{n_k}) such that $x_{n_k} \in I_{n_k}$. Given $\varepsilon > 0$. Since $|I_k| \rightarrow 0$ as $k \rightarrow \infty$, there exists a $K \in \mathbb{N}$ such that $|I_k| < \varepsilon$ for all $k > K$. Furthermore we have $|x_{n_j} - x_{n_k}| < \varepsilon$ for $j, k > K$, i.e. (x_{n_k}) is a Cauchy sequence and thus converges by Theorem 2.12. \square

The Bolzano-WeierstraßTheorem does not claim that the subsequence is unique, i.e. there might be convergent subsequences with different limits depending on the choice of L_k or R_k .

THEOREM 2.16. *If (x_n) is a bounded sequence in \mathbb{R} such that every convergent subsequence has the same limit x , then (x_n) converges to x .*

PROOF. We will show the contrapositive statement: Suppose a bounded sequence does not converge to x . Then (x_n) has a convergent subsequence with limit different from x .

If (x_n) does not converge to x , then there exists $\varepsilon_0 > 0$ such that $|x_n - x| \geq \varepsilon_0$ for infinitely many $n \in \mathbb{N}$. Hence there exists a subsequence (x_{n_k}) such that $|x_{n_k} - x| \geq \varepsilon_0$ for every $k \in \mathbb{N}$. Note that (x_{n_k}) is a bounded sequence and so by Bolzano-Weierstraßthere exists a convergent subsequence $(x_{n_{k_j}})$. If $\lim_j x_{n_{k_j}} = y$, then $|x - y| \geq \varepsilon_0$. In other words, x is not equal to y . \square

2.1.3. Topology of \mathbb{R} . In this section we treat some basic notions of topology for the real line. Generalizations of these notions and its manifestations in normed spaces and general metric spaces are going to be the pillars of this course.

We generalize the notion of open intervals (a, b) and closed intervals $[a, b]$.

DEFINITION 2.1.12 (Open sets). A subset O of \mathbb{R} is called *open* if for every $x \in O$ there exists an open interval I contained in O with $x \in I$.

DEFINITION 2.1.13 (Closed sets). A subset C of \mathbb{R} is called *closed* if the complement $C^c = \mathbb{R} \setminus C = \{x \in \mathbb{R} : x \notin C\}$ is open.

Note that the interval (a, b) is an open set and $[a, b]$ is closed. Observe further that by definition the empty set \emptyset and \mathbb{R} are open and closed.

PROPOSITION 2.1.14. *Suppose $\{I_j\}_{j \in J}$ is a collection of open intervals in \mathbb{R} with non-empty intersection $\bigcap_{j \in J} I_j \neq \emptyset$.*

- (1) *If J has finitely many elements, then $\bigcap_{j \in J} I_j$ is an open interval.*
- (2) *$\bigcup_{j \in J} I_j$ is an open interval for an arbitrary index set J .*

PROOF. We define open intervals $I_j = (a_j, b_j)$ for real numbers $a_j < b_j$, the interval bounds are also allowed to be $\pm\infty$, and set $I := \bigcup_{j \in J} I_j$.

- (1) We pick a point x in $\bigcup_{j=1}^n I_j$ and set $a := \max\{a_j : j = 1, \dots, n\}$ and $b := \min\{b_j : j = 1, \dots, n\}$. If all the a_j 's are $-\infty$, then $a = -\infty$, and if all the b_j 's are ∞ , then we have $b = \infty$.
Since $a_j < x < b_j$ for $j = 1, \dots, n$ we get that $x \in (a, b)$. Furthermore, we have that $\bigcap_{j \in J} (a_j, b_j) = (a, b)$.
- (2) We choose $x \in \bigcap_{j \in J} I_j$. Suppose $y \in \bigcup_{j \in J} I_j$. Then $y \in I_j$ for some $j \in J$. Since $x \in I_j$, the interval $(x, y) \subset I_j$ and thus in I . Hence I is the interval (a, b) , where $a = \inf\{a_j : j \in J\}$ or $-\infty$ and $b = \sup\{b_j : j \in J\}$ or ∞ . \square

The assumption in (i) cannot be weakened, e.g. $\bigcap_{n=1}^{\infty} (-1/n, 1/n) = \{0\}$. Hence an infinite intersection of open intervals is not necessarily an open interval. We show that the preceding statement is true for a more general class of sets, the open sets.

PROPOSITION 2.1.15. Let $\{O_j : j \in J\}$ be a family of open sets of \mathbb{R} .

- (1) $\bigcap_{j=1}^n O_j$ is an open set for any $n \in \mathbb{N}$.
- (2) $\bigcup_{j \in J} O_j$ is open for a general index set J .

PROOF. (1) We set $O = \bigcap_{j=1}^n O_j$. If $x \in O$, then $x \in O_j$ for $j = 1, \dots, n$. Since O_j 's are open, there are open intervals $I_j \subset O_j$ containing x . Hence, we have that $\bigcap_{j=1}^n I_j \subset \bigcap_{j=1}^n O_j$, the desired assertion.

- (2) Let x be in $\bigcup_{j \in J} O_j$. Then there exists some j such that $x \in O_j$ and thus an open interval I_j contained in O_j with $x \in I_j$ and consequently $I_j \subset O$. Hence O is an open set. □

We are in the position to introduce a notion of closedness between points, known as neighborhoods.

DEFINITION 2.1.16. Given $x \in \mathbb{R}$. Then a subset U of \mathbb{R} is called a *neighborhood* of x if there exists an open subset O of \mathbb{R} such that $x \in O \subset U$.

Due to the structure of \mathbb{R} we have that U is a neighborhood of x if and only if there exists a $\delta > 0$ such that $(x - \delta, x + \delta) \subset U$.

DEFINITION 2.1.17. For a subset A we introduce some notions.

- (1) The *closure* of a subset A of \mathbb{R} , denoted by \overline{A} , is the intersection of all closed sets containing A .
- (2) The *interior* of a subset A of \mathbb{R} , denoted by $\text{int}A$, is the union of all open subsets of \mathbb{R} contained in A .
- (3) The *boundary* of a subset A of \mathbb{R} , denoted by $\text{bd}A$, is the set $\overline{A} \setminus \text{int}A$.

Note that $\text{bd}A$ is a closed set and that the closure of a bounded subset of \mathbb{R} is bounded, too.

Here are some useful facts.

LEMMA 2.17. Suppose A is a subset of \mathbb{R} .

- (1) $\overline{A} = (\text{Int}(A^c))^c$ and $\text{int}(A) = (\overline{A^c})^c$
- (2) $\text{bd}A = \text{bd}(A^c) = \overline{A} \cap \overline{A^c}$
- (3) $\overline{A} = A \cup \text{bd}A = \text{int}A \cup \text{bd}A$

PROOF. (1) These identities are a consequence of the following general fact: B is a closed containing A if and only if B^c is open and $B^c \subset A^c$. The statement about the interior of A is the first statement for A^c instead of A .

- (2) $\text{bd}A = \overline{A} \setminus \text{int}A = \overline{A} \cap (\text{int}A)^c = \overline{A} \cap \overline{A^c}$, where we used (i) in the last step. Let us compute $\text{bd}A^c$: $\text{bd}A^c = \overline{A^c} \setminus \text{int}A^c = \overline{A^c} \cap (\text{int}A^c)^c = \overline{A^c} \cap \overline{A}$. Hence we have the desired assertions.

- (3) First note that $\text{int}A \cup \text{bd}A \subset A \cup \text{bd}A \subset \overline{A}$. Furthermore we have $\text{int}A \cup \text{bd}A = \text{int}A \cup (\overline{A} \setminus A) = \text{int}A \cup (\overline{A} \cap (\text{int}A)^c) = ((\text{int}A) \cup \overline{A}) \cap (\text{int}A \cup (\text{int}A)^c) = \overline{A}$. □

LEMMA 2.18. Suppose A is a subset of \mathbb{R} .

- (1) $\overline{A} = \{x \in \mathbb{R} : \text{every neighborhood of } x \text{ intersects } A\}$

- (2) $\text{int}(A) = \{x \in \mathbb{R} : \text{some neighborhood of } x \text{ is contained in } A\}$
 (3) $\text{bd}(A) = \{x \in \mathbb{R} : \text{every neighborhood of } x \text{ intersects } A \text{ and its complement}\}$

PROOF. (1) We choose an open neighborhood U of $x \in \mathbb{R}$ that does not intersect A , i.e. $A \subset U^c$. Since U^c is closed, we have that $\overline{A} \subset U^c$ and from $x \notin U^c$ we also have that $x \notin \overline{A}$. On the other hand, if $x \notin \overline{A}$, then $(\overline{A})^c$ is an open set containing x that is disjoint from A .

- (2) Follows from (i) and the preceding proposition.
 (3) Follows from (i), (ii) and the preceding proposition. □

DEFINITION 2.1.18. Let A be a subset of \mathbb{R} .

- (1) A point $x \in A$ is *isolated* in A if there exists a neighborhood U of x such that $U \cap A = \{x\}$.
 (2) A point $x \in \mathbb{R}$ is said to be an *accumulation point* of A if every neighborhood of x contains points in $A \setminus \{x\}$.

Note: Accumulation points of a set are not necessarily elements of the set. A well-known example is $A = \{1/n : n \in \mathbb{N}\}$ with 0 as accumulation point, which is clearly not in A .

The definition of an accumulation point makes only sense for sets with infinitely many elements.

Finally, an infinite closed set may not have accumulation points, e.g. $\mathbb{N} \subset \mathbb{R}$ has no accumulation points in \mathbb{R} .

LEMMA 2.19. *A point $x \in \mathbb{R}$ is an accumulation point of A if and only if every neighborhood of x contains infinitely many points of A .*

PROOF. One direction: Suppose every neighborhood of x contains infinitely many points of A , then x is an accumulation point of A .

Other direction: Suppose x is an accumulation point of A . For a neighborhood U of x , we choose $n_1 \in \mathbb{N}$ such that $(x - 1/n_1, x + 1/n_1) \subset U$. Take a point x_1 different from x in $A \cap (x - 1/n_1, x + 1/n_1)$. Now we repeat the procedure: Take $n_2 \geq n_1$ such that $x_1 \notin (x - 1/n_2, x + 1/n_2)$ and pick $x_2 \in A \cap (x - 1/n_2, x + 1/n_2)$ with $x_2 \neq x$. We continue in this way and get a sequence of points $(x_n) \subset A \cap U$. □

PROPOSITION 2.1.19. *Let A be a subset of \mathbb{R} . Then $\overline{A} = \{\text{isolated points of } A\} \cup \{\text{accumulation points of } A\}$.*

PROOF. Suppose $x \in \overline{A}$. Then if $x \in A$, then either x is isolated in A or every neighborhood of x contains points in A different from x . In the later case x is an accumulation point of A . Now assume $x \in \overline{A}$ and $x \notin A$. Then every neighborhood of x has a non-trivial intersection with A , and thus x is an accumulation point of A . In summary, we have that the closure of A is the union of the isolated points of A with the accumulation points of A .

For the converse we note: If x is isolated, then x is definitely in A . If x is an accumulation point of A , then $x \in \overline{A}$ □

DEFINITION 2.1.20. A subset A of \mathbb{R} is said to be *dense* in \mathbb{R} if its closure is equal to \mathbb{R} , i.e. $\overline{A} = \mathbb{R}$.

PROPOSITION 2.1.21. *The set of rational numbers, \mathbb{Q} , is dense in \mathbb{R} .*

PROOF. For an arbitrary $x \in \mathbb{R}$ we consider a neighborhood U of x . Then we know that U contains the interval $(x - \varepsilon, x + \varepsilon)$ for a sufficiently small $\varepsilon > 0$. By an earlier result we have that there exists a rational number in $(x - \varepsilon, x + \varepsilon)$. \square

We also have that the set of irrational numbers is dense in \mathbb{R} .

The property that \mathbb{Q} has only countably elements, but still is dense in \mathbb{R} is a very favorable property and occurs in various other situations. We say that \mathbb{R} is separable.

\mathbb{Q} is a dense subset of \mathbb{R} with empty interior and thus the boundary of Q is all of \mathbb{R} . The same is true for the set of irrational numbers.

2.1.4. Supplementary material.

THEOREM 2.20 (Nested Interval Theorem). *Let $\{I_j\}_{j=1}^{\infty}$ be a sequence of closed bounded intervals in \mathbb{R} , such that $I_j \subset I_{j+1}$ for all $j \in \mathbb{N}$. We assume in addition that the lengths of the intervals $|I_j|$ tends to zero. Then $I := \bigcap_{j \in \mathbb{N}} I_j = \{z\}$ for some $z \in \mathbb{R}$.*

PROOF. Without loss of generality we assume $I_j = [a_j, b_j]$. Then the assumptions yield that $a_1 \leq a_2 \leq \dots \leq b_2 \leq b_1$ and that for every $\varepsilon > 0$ there exist a $j \in \mathbb{N}$ such that $b_j - a_j \leq \varepsilon$.

We set $A := \{a_j : j \in \mathbb{N}\}$ and $B := \{b_j : j \in \mathbb{N}\}$, note that $a := \sup A < \infty$ and $b = \inf B < \infty$, and $a_j \leq a \leq b_j$ for $j \in \mathbb{N}$. Hence we have $[a, b] = \bigcap_{j=1}^{\infty} [a_j, b_j]$ and by the assumption on the shrinking of the interval lengths we get that $a = b = z$ for some $z \in \mathbb{R}$. \square

Normed spaces and innerproduct spaces

3.1. Normed spaces and innerproduct spaces

In this course vector spaces are equipped with additional structures in order to measure the distance between elements and formulate convergence of sequences of elements of vector spaces, or to provide quantitative and qualitative information on operators.

3.1.1. Normed spaces. The norm on a general vector space generalizes the notion of the length of a vector in \mathbb{R}^2 and \mathbb{R}^3 .

DEFINITION 3.1.1. A *normed space* $(X, \|\cdot\|)$ is a vector space X together with a function $\|\cdot\| : X \rightarrow \mathbb{R}$, the *norm* on X , such that for all $x, y \in X$ and $\lambda \in \mathbb{R}$:

- (1) *Positivity*: $0 \leq \|x\| < \infty$ and $\|x\| = 0$ if and only if $x = 0$;
- (2) *Homogeneity*: $\|\lambda x\| = |\lambda| \|x\|$;
- (3) *Triangle inequality*: $\|x + y\| \leq \|x\| + \|y\|$.

Normed spaces have a rich structure.

PROPOSITION 3.1.2. Let $(X, \|\cdot\|)$ be a normed space. Then $d : X \times X \rightarrow \mathbb{R}$ defined by $d(x, y) = \|x - y\|$ satisfies for all $x, y, z \in X$ (i) $d(x, y) \geq 0$ and $d(x, x) = 0$ if and only if $x = 0$ (positivity); (ii) $d(x, y) = d(y, x)$ (symmetry); (iii) $d(x, z) \leq d(x, y) + d(y, z)$ (triangle inequality).

The function $d(x, y) = \|x - y\|$ on the vector space X is an example of a distance function on X , aka as a metric. We will later discuss such distance functions on a general set.

PROOF. The properties (i)-(iii) are direct consequences of the axioms for a norm. In particular, (i) follows from property (1) of a norm, (ii) is derived from property (ii) of a norm for $\lambda = -1$ and (iii) is deduced from property (3) of a norm. \square

The metric d on X is also compatible with the linear structure of a vector space:

- *Translation invariance*: $d(x + z, y + z) = d(x, y)$ for all $x, y, z \in X$;
- *Homogeneity*: $d(\lambda x, \lambda y) = |\lambda| d(x, y)$ for all $x, y \in X$ and scalars $\lambda \in \mathbb{R}$.

The metric d on X gives us a way to generalize intervals in \mathbb{R} , called balls.

DEFINITION 3.1.3. For $r > 0$ and $x \in X$ we define the *open ball* $B_r(x)$ of radius r and center x as the set

$$B_r(x) = \{y \in X : \|x - y\| < r\},$$

and the *closed ball* $\overline{B}_r(x)$ of radius r and center x as

$$\overline{B}_r(x) = \{y \in X : \|x - y\| \leq r\}.$$

The translation invariance and the homogeneity imply that the ball $B_r(x)$ is the image of the unit ball $B_1(0)$ centered at the origin under the affine mapping $f(y) = ry + x$.

The balls $B_r(x)$ have another peculiar feature. Namely, these are convex subsets of X .

DEFINITION 3.1.4. Let X be a vector space.

- For two points $x, y \in X$ the *interval* $[x, y]$ is the set of points $\{z \mid z = \lambda x + (1 - \lambda)y, 0 \leq \lambda \leq 1\}$.
- A subset E of X is called *convex* if for any two points $x, y \in E$ the interval $[x, y]$ is also in E .

The notion of convexity is central to the theory of vector spaces and enters in an intricate manner in functional analysis, numerical analysis, optimization, etc. .

LEMMA 3.1. Let $(X, \|\cdot\|)$ be a normed vector space. Then the unit ball $B_1(0) = \{x \in X \mid \|x\| \leq 1\}$ is a convex set.

PROOF. For $x, y \in B_1(0)$ we have that $\|\lambda x + (1 - \lambda)y\| \leq |\lambda|\|x\| + |1 - \lambda|\|y\| = 1$, because $\|x\|, \|y\|$ are both less than or equal to 1. Thus $\lambda x + (1 - \lambda)y \in B_1(0)$. \square

The real numbers with the absolute value is a normed space $(\mathbb{R}, |\cdot|)$ and the open ball $B_r(x)$ is the open interval $(x - r, x + r)$ and $\overline{B}_r(x)$ is the closed interval $[x - r, x + r]$.

A fundamental class of metric spaces is \mathbb{R}^n with the ℓ^p -norms.

DEFINITION 3.1.5. For $p \in [1, \infty)$ we define the ℓ^p -norm $\|\cdot\|_p$ on \mathbb{R}^n by assigning to $x = (x_1, \dots, x_n) \in \mathbb{R}^n$ the number $\|x\|_p$:

$$\|x\|_p = (|x_1|^p + |x_2|^p + \dots + |x_n|^p)^{1/p}$$

. For $p = \infty$ we define the ℓ^∞ -norm $\|\cdot\|_\infty$ on \mathbb{R}^n by

$$\|x\|_\infty = \max |x_1|, \dots, |x_n|.$$

The notation for $\|\cdot\|_\infty$ is justified by the fact that it is the limit of the $\|\cdot\|_p$ -norms.

LEMMA 3.2. For $x \in \mathbb{R}^n$ we have that

$$\|x\|_\infty = \lim_{p \rightarrow \infty} \|x\|_p.$$

Some inequalities enter the stage: Hölder's inequality and Young's inequality. For $p \in (1, \infty)$ we define its *conjugate* q as the number such that

$$\frac{1}{p} + \frac{1}{q} = 1.$$

If $p = 1$, then we define its conjugate q to be ∞ and if $p = \infty$ then $q = 1$.

LEMMA 3.3 (Young's inequality). For $p \in (1, \infty)$ and q its conjugate we have

$$ab \leq \frac{a^p}{p} + \frac{b^q}{q},$$

for $a, b \geq 0$.

PROOF. Consider the function $f(x) = x^{p-1}$ and integrate this with respect to x from zero to a . Now take the inverse of f given by $f^{-1}(y) = y^{q-1}$ and integrate it from zero to V . Then the sum of these two integrals always exceeds the product ab , but the integrals are a^p/p and b^q/q . Hence we have established the desired inequality. \square

A consequence of Young's inequality is Hölder's inequality.

LEMMA 3.4. Suppose $p \in (1, \infty)$ and $x = (x_1, \dots, x_n)$ and $y = (y_1, \dots, y_n)$ are vectors in \mathbb{R}^n . Then

$$\left| \sum_{i=1}^n x_i y_i \right| \leq \left(\sum_{i=1}^n |x_i|^p \right)^{1/p} \left(\sum_{i=1}^n |y_i|^q \right)^{1/q}.$$

PROOF. Set $a_i = |x_i|/(\sum_{i=1}^n |x_i|^p)^{1/p}$ and $b_i = |y_i|/(\sum_{i=1}^n |y_i|^q)^{1/q}$. Then we have $\sum_i a_i^p = 1$ and $\sum_i b_i^q = 1$. By Young's inequality

$$\sum_{i=1}^n |x_i| |y_i| \leq \left(\sum_{i=1}^n |x_i|^p \right)^{1/p} \left(\sum_{i=1}^n |y_i|^q \right)^{1/q}.$$

\square

PROPOSITION 3.1.6. The space \mathbb{R}^n with the ℓ^p -norm $\|\cdot\|_p$ is a normed space for $p \in [1, \infty]$.

As an exercise I propose to draw the unit balls of $(\mathbb{R}^2, \|\cdot\|_1)$, $(\mathbb{R}^2, \|\cdot\|_2)$ and $(\mathbb{R}^2, \|\cdot\|_\infty)$.

PROOF. First we show that ℓ^p is a vector space for $p \in [1, \infty)$: For $\lambda \in \mathbb{F}$ and $x \in \ell^p$ we have $\lambda x \in \ell^p$. One has to work a little bit to see that for $x, y \in \ell^p$ also $x + y \in \ell^p$:

$$\begin{aligned} \|x + y\|_p^p &= \sum_{n=1}^{\infty} |x_n + y_n|^p \\ &\leq \sum_{n=1}^{\infty} |2 \max\{|x_n|, |y_n|\}|^p \\ &= 2^p \sum_{n=1}^{\infty} |\max\{|x_n|, |y_n|\}|^p \\ &\leq 2^p \left(\sum_{n=1}^{\infty} |x_n|^p + \sum_{n=1}^{\infty} |y_n|^p \right) = 2^p (\|x\|_p^p + \|y\|_p^p) < \infty. \end{aligned}$$

Positivity and homogeneity are consequences of the corresponding properties of the absolute value of a real number. The triangle inequality is the non-trivial assertion that we split up in three cases $p = 1$, $p = \infty$ and $p \in (1, \infty)$. Let $x = (x_1, \dots, x_n)$ and $y = (y_1, \dots, y_n)$ be points in \mathbb{R}^n .

(1) For $p = 1$ we have

$$\|x + y\|_1 = |x_1 + y_1| + \dots + |x_n + y_n| \leq |x_1| + |y_1| + \dots + |x_n| + |y_n| \leq \|x\|_1 + \|y\|_1$$

(2) For $p = \infty$ the argument is similar:

$$\begin{aligned}\|x + y\|_\infty &= \max\{|x_1 + y_1|, \dots, |x_n + y_n|\} \\ &= \max\{|x_1| + |y_1|, \dots, |x_n| + |y_n|\} \\ &= \max\{|x_1|, \dots, |x_n|\} + \max\{|y_1|, \dots, |y_n|\} = \|x\|_\infty + \|y\|_\infty.\end{aligned}$$

(3) The general case $p \in (1, \infty)$: The triangle inequality in this case is also known as *Minkowski's inequality*. We deduce it from Hölder's inequality

$$\begin{aligned}\|x + y\|_p^p &= \sum_{i=1}^n |x_i + y_i|^p \\ &\leq \sum_{i=1}^n |x_i + y_i|^{p-1} (|x_i| + |y_i|) \\ &\leq \sum_{i=1}^n |x_i + y_i|^{p-1} |x_i| + \sum_{i=1}^n |x_i + y_i|^{p-1} |y_i| \\ &\leq \left(\sum_{i=1}^n |x_i + y_i|^p \right)^{1/q} \left(\left(\sum_{i=1}^n |x_i|^p \right)^{1/p} + \left(\sum_{i=1}^n |y_i|^p \right)^{1/p} \right) \\ &= \|x + y\|_p^{1/q} (\|x\|_p + \|y\|_p)\end{aligned}$$

Dividing by $\|x + y\|_p^{1/q}$ and using $1 - 1/q = 1/p$ we arrive at Minkowski's inequality:

$$\|x + y\|_p \leq \|x\|_p + \|y\|_p.$$

□

EXAMPLE 3.1.7. Let \mathbb{C}^n be the vector space of complex n -tuples. There one also has $\|\cdot\|_p$ norms for $1 \leq p < \infty$:

$$\|x\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{1/p}, \quad x \in \mathbb{C}^n$$

where $x_i \in \mathbb{C}$ and $|x_i| = (x_i \bar{x}_i)^{1/2}$ denotes the modulus of x_i . The sup-norm of $x \in \mathbb{C}^n$ is defined by $\|x\|_\infty = \max |x_i| : i = 1, \dots, n$, where again $|\cdot|$ denotes the modulus of a complex number.

A natural generalization of the normed spaces $(\mathbb{R}^n, \|\cdot\|_p)$ is to replace tuples of finite length with ones of infinite length $x = (x_1, x_2, \dots)$ with $x_i \in \mathbb{R}$, i.e. $(\mathbb{R}^\infty, \|\cdot\|_p)$. The standard notation for these normed spaces is $(\ell^p, \|\cdot\|_p)$ because these are special classes of the Lebesgue spaces $L^p(\mathbb{N}, d\mu)$ for the counting measure. One often refers to these spaces as “little L^p ”-spaces.

EXAMPLE 3.1.8. For $1 \leq p < \infty$ the spaces $(\ell^p, \|\cdot\|_p)$ are normed spaces of convergent sequences $x = (x_i)_i$ such that

$$\|x\|_p = |x_1|^p + |x_2|^p + \dots < \infty,$$

and $(\ell^\infty, \|\cdot\|_\infty)$ is the space of bounded sequences $(x_i)_i$ with respect to the norm

$$\|x\|_\infty = \sup\{|x_i| : i = 1, 2, \dots\}.$$

We have the following inclusions:

$$\ell^1 \subset \ell^2 \subset \dots \ell^\infty.$$

For example $(1/n)_n$ is in ℓ^p for $p \geq 2$, but not in ℓ^1 .

EXERCISE 3.1.9. Suppose $p, q \in [1, \infty]$. Show that for $p < q$ the space ℓ^p is a proper subspace of ℓ^∞ .

Let us view these vectors of infinite length as real-valued sequences. Then the assumption $\|x\|_p$ imposes conditions on the structure of the sequences. For example, $\|x\|_\infty = \sup_i |x_i|$ is finite if and only if x is a bounded sequence, and $\|x\|_1 = \sum_{i=1}^\infty |x_i|$ is finite if the sequence (x_i) is absolutely summable. The norms $\|\cdot\|_p$ for $1 \leq p < \infty$ describe different notions of convergence, but $\|\cdot\|_\infty$ does not impose convergence but just boundedness.

PROPOSITION 3.1.10. For $1 \leq p \leq \infty$ the spaces $(\ell^p, \|\cdot\|_p)$ are normed spaces.

The proof of the finite-dimensional setting extends to the infinite-dimensional setting because Hölder's inequality is valid for ℓ^p -norms.

LEMMA 3.5 (Hölder's inequality). For $1 < p < \infty$ and q its conjugate index, $x \in \ell^p$ and $y \in \ell^q$ we have

$$\sum_{i=1}^\infty |x_i| |y_i| \leq \|x\|_p \|y\|_q.$$

EXAMPLE 3.1.11. Define a norm on $\mathcal{M}_{m \times n}(\mathbb{C})$ by picking a norm on \mathbb{C}^{mn} . For example, $\|A\|_{(p)} = (\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^p)^{1/p}$ or $\|A\|_{(\infty)} = \max |a_{ij}|$. The case $p = 2$ is of interest and is known as the Frobenius norm.

3.1.2. Innerproduct spaces. For vectors in \mathbb{R}^3 we have the 'dot product' aka 'scalar product' that assigns to a pair of vectors $x = (x_1, x_2, x_3)$ and $y = (y_1, y_2, y_3)$ the number

$$\langle x, y \rangle = x_1 y_1 + x_2 y_2 + x_3 y_3.$$

Pythagoras' theorem gives the length of $x = (x_1, x_2, x_3)$ as $\sqrt{x_1^2 + x_2^2 + x_3^2}$. Note that $\langle x, x \rangle = \sqrt{x_1^2 + x_2^2 + x_3^2}$. Innerproduct spaces are a generalization of these basic facts from Euclidean geometry to general vector spaces.

DEFINITION 3.1.12. Let X be a vector space. An *innerproduct* on X is a map $\langle \cdot, \cdot \rangle : X \times X \rightarrow \mathbb{F}$, which has the following properties:

- (1) (Linearity) For vectors $x_1, x_2, y \in X$ and scalars $\lambda_1, \lambda_2 \in \mathbb{F}$ we have $\langle \lambda_1 x_1 + \lambda_2 x_2, y \rangle = \lambda_1 \langle x_1, y \rangle + \lambda_2 \langle x_2, y \rangle$.
- (2) (Symmetry) For vectors $x, y \in X$ we have $\langle x, y \rangle = \langle y, x \rangle$ for $\mathbb{F} = \mathbb{R}$ and $\langle x, y \rangle = \overline{\langle y, x \rangle}$ for $\mathbb{F} = \mathbb{C}$.
- (3) (Positive definiteness) For any $x \in X$ we have $\langle x, x \rangle \geq 0$ and $\langle x, x \rangle = 0$ if and only if $x = 0$.

We call $(X, \langle \cdot, \cdot \rangle)$ an *innerproduct space* and denote by $\|x\| = \langle x, x \rangle^{1/2}$ the length of x .

We state a theorem of utmost importance about innerproduct spaces.

THEOREM 3.6 (Cauchy-Schwarz). *Suppose X is an innerproduct space. Then for all $x, y \in X$ we have*

$$|\langle x, y \rangle| \leq \|x\| \|y\|.$$

We have $|\langle x, y \rangle| = \|x\| \|y\|$ if and only if $x = \lambda y$ for some $\lambda \in \mathbb{F}$.

PROOF. Suppose $x \neq 0$ otherwise the inequality is trivial. Then we consider $z = \langle x, y \rangle x - \langle x, x \rangle y$. By the properties of an innerproduct we have

$$0 \leq \langle z, z \rangle = |\langle x, y \rangle|^2 \langle x, x \rangle - 2|\langle x, y \rangle|^2 \langle x, x \rangle + \langle x, x \rangle^2 \langle y, y \rangle,$$

hence we obtain

$$|\langle x, y \rangle|^2 \langle x, x \rangle \leq \langle x, x \rangle^2 \langle y, y \rangle$$

and after dividing through by the strictly positive number $\langle x, x \rangle$ we obtain the Cauchy-Schwarz inequality.

We have equality if and only if $z = 0$, which yields that $x = \lambda y$ for $\lambda = \langle x, x \rangle \langle x, y \rangle^{-1}$. \square

As first consequence we deduce that innerproduct spaces $(X, \langle \cdot, \cdot \rangle)$ are normed spaces for $\|x\| = \langle x, x \rangle^{1/2}$.

PROPOSITION 3.1.13. *For $(X, \langle \cdot, \cdot \rangle)$ the expression $\|x\| = \langle x, x \rangle^{1/2}$ defines a norm on X .*

PROOF. Homogeneity follows from the linearity of the innerproduct. The triangle inequality follows from Cauchy-Schwarz:

$$\|x + y\|^2 = \|x\|^2 + \|y\|^2 + 2\langle x, y \rangle \leq \|x\|^2 + \|y\|^2 + 2\|x\| \|y\|,$$

so the right side is $(\|x\| + \|y\|)^2$ and thus we have $\|x + y\| \leq \|x\| + \|y\|$. \square

The sequence space ℓ^2 was the first example of an innerproduct space, studied by D. Hilbert in 1901 in his work on Fredholm operators.

EXAMPLE 3.1.14. The sequence space ℓ^2 is an innerproduct space for real-valued sequences $(x_i), (y_i)$

$$\langle x, y \rangle = \sum_{i=1}^{\infty} x_i y_i$$

and

$$\langle x, y \rangle = \sum_{i=1}^{\infty} x_i \overline{y_i}$$

for complex-valued sequences.

The innerproduct $\langle \cdot, \cdot \rangle$ and its associated norm $\|\cdot\| = \langle \cdot, \cdot \rangle^{1/2}$ are related by the *polarization identity*.

LEMMA 3.7 (Polarization identity). *Let $(X, \langle \cdot, \cdot \rangle)$ be an innerproduct space with norm $\|\cdot\| = \langle \cdot, \cdot \rangle^{1/2}$.*

(1) *For a real innerproduct space we have $\langle x, y \rangle = \frac{1}{4}(\|x + y\|^2 - \|x - y\|^2)$ for all $x, y \in X$.*

(2) *For a complex innerproduct space we have $\langle x, y \rangle = \frac{1}{4} \sum_{k=1}^4 i^k \|x + i^k y\|^2$.*

PROOF. The arguments are based on the homogeneity properties of innerproducts.

- (1) $\|x + (-1)^k y\|^2 = \|x\|^2 + \|y\|^2 + (-1)^k \langle x, y \rangle$ for $k = 0, 1$. Adding these two identities yields the desired polarization identity.
- (2) Left as an exercise.

□

Jordan and von Neumann gave an elementary characterizations of norms that arise from innerproducts.

THEOREM 3.8 (Jordan-von Neumann). *Suppose $(X, \|\cdot\|)$ is a complex normed space. If the norm satisfies the parallelogram identity*

$$\|x - y\|^2 + \|x + y\|^2 = 2\|x\|^2 + 2\|y\|^2 \quad \text{for all } x, y \in X,$$

then X is an innerproduct space for the innerproduct

$$\langle x, y \rangle = \frac{1}{4} \sum_{k=1}^4 i^k \|x + i^k y\|^2.$$

The proof of this useful result is elementary and will be given in the supplement to the chapter.

Innerproduct spaces are the infinite-dimensional counterparts of $(\mathbb{R}^n, \|\cdot\|_2)$ and share many properties with these finite-dimensional spaces, in contrast to general normed spaces such as $C(I)$ with the sup-norm.

EXAMPLE 3.1.15. The supremums norm of $C[0, 1]$ does not come from an innerproduct. Use the polarization identity to show this fact.

A way to address this issue is to change the norm. Namely, if one instead equips $C(I)$ with the 2-norm $\|\cdot\|_2$ for functions, then one gets an innerproduct space.

LEMMA 3.9. *Let I be an interval of \mathbb{R} . Then the space of continuous complex-valued functions $C(I)$ is an innerproduct space for*

$$\langle f, g \rangle = \int_I f(x) \overline{g(x)} dx$$

for functions $f \in C(I)$ with finite norm

$$\|f\|_2 = \int_I |f(x)|^2 dx < \infty.$$

PROOF. We have $\langle \lambda f, g \rangle = \int_I \lambda f(x) \overline{g(x)} dx = \lambda \int_I f(x) \overline{g(x)} dx = \lambda \langle f, g \rangle$ for $\lambda \in \mathbb{C}$, and $\langle f, g \rangle = \int_I f(x) \overline{g(x)} dx = \int_I \overline{f(x)} g(x) dx$. Note that $|f(x)|^2$ is non-negative for $f \in C(I)$ and that it is zero for those $x \in I$ with $f(x) = 0$. By the properties of the integral we have shown the positivity of $\langle \cdot, \cdot \rangle$. □

Historical note: The Cauchy-Schwarz inequality for $(C(\mathbb{R}), \langle \cdot, \cdot \rangle)$ is due to Karl H. A. Schwarz in 1888 for continuous functions, and Cauchy for \mathbb{R}^n with the Euclidean innerproduct.

Innerproducts yield a generalization of the notion of *orthogonality* of elements.

DEFINITION 3.1.16. Two elements x, y in an innerproduct space $(V, \langle \cdot, \cdot \rangle)$ are *orthogonal* to each other if $\langle x, y \rangle = 0$

The theorem of Pythagoras is true for any innerproduct space $(X, \langle \cdot, \cdot \rangle)$.

PROPOSITION 3.1.17 (Pythagoras's Theorem). *Let $(X, \langle \cdot, \cdot \rangle)$ be an innerproduct space. For two orthogonal elements $x, y \in X$ we have*

$$\|x + y\|^2 = \|x\|^2 + \|y\|^2.$$

PROOF. The argument is based on the fact that $\langle x, x \rangle$ is a norm. By assumption we have $\langle x, y \rangle = 0$

$$\|x + y\|^2 = \|x\|^2 + 2\operatorname{Re} \langle x, y \rangle + \|y\|^2 = \|x\|^2 + \|y\|^2.$$

□

As an example we consider some orthogonal vectors in $(C([0, 1]), \langle \cdot, \cdot \rangle)$. For $m \neq n$ we define the exponentials $e_m(x) = e^{2\pi i m x}$ and $e_n(x) = e^{2\pi i n x}$. Then

$$\langle e_m, e_n \rangle = \int_0^1 e^{2\pi i(m-n)x} dx = (2\pi i(m-n))^{-2} (e^{2\pi i(m-n)} - 1) = 0.$$

Note that $\langle e_n, e_n \rangle = 1$ for any $n \in \mathbb{Z}$. With the help of Kronecker's delta function we may express this as $\langle e_m, e_n \rangle = \delta_{m,n}$.

The theorem of Pythagoras is now at our disposal in any innerproduct spaces such as ℓ^2 .

DEFINITION 3.1.18. A set of vectors $\{e_i\}_{i \in I}$ in an innerproduct space $(X, \langle \cdot, \cdot \rangle)$ is called an *orthogonal family* if $\langle e_i, e_j \rangle = 0$ for all $i \neq j$. In case that the orthogonal family $\{e_i\}_{i \in I}$ in V satisfies in addition $\|e_i\| = 1$ for any $i \in I$, then we refer to it as *orthonormal family*.

The set of vectors $\{e_i\}_{i \in I}$ is in general an infinite set. The exponentials $\{e^{2\pi i n x}\}_{n \in \mathbb{Z}}$ is an orthonormal family in $C[0, 1]$ with respect to $\langle \cdot, \cdot \rangle_2$ and is a system of utmost importance, e.g. it lies at the heart of Fourier analysis or more generally harmonic analysis.

Orthonormal families have an interesting property, known as *Bessel's inequality*.

PROPOSITION 3.1.19 (Bessel's inequality). *Suppose $\{e_i\}_{i \in I}$ is a countably infinite orthonormal family in an innerproduct space $(X, \langle \cdot, \cdot \rangle)$. Then for any $x \in X$ we have*

$$\sum_{i \in I} |\langle x, e_i \rangle|^2 \leq \|x\|^2.$$

Recall that a set I is *countably infinite* if there exists a bijection between I and the set of natural numbers \mathbb{N} , e.g. the set of integers \mathbb{Z} .

PROOF. It suffices to check the inequality for $I = \mathbb{N}$. Consider the vector $\tilde{x} = \sum_{i=1}^n \langle x, e_i \rangle e_i$ for each $n \in \mathbb{N}$. By the orthonormality of the set $\{e_i\}_{i \in I}$ we have

$$0 \leq \|x - \tilde{x}\|^2 = \|x\|^2 - \sum_{i=1}^n |\langle x, e_i \rangle|^2.$$

Thus the sequence of real numbers $(\sum_{i=1}^n |\langle x, e_i \rangle|^2)$ is bounded above and non-decreasing. Therefore it has a limit

$$\sum_{i=1}^{\infty} |\langle x, e_i \rangle|^2 \leq \|x\|^2.$$

□

The case of equality in Bessel's inequality characterizes an important properties of orthonormal systems and will be discussed in the chapter on Hilbert spaces.

For example Bessel's inequality for the set of exponentials $\{e^{2\pi inx}\}_{n \in \mathbb{Z}}$ in $(C[0, 1], \langle \cdot, \cdot \rangle_2)$ is a statement about the Fourier coefficients of f

$$\widehat{f}(n) = \int_0^1 f(x)e^{-2\pi nx} dx,$$

then we have

$$\sum_{n \in \mathbb{Z}} |\widehat{f}(n)|^2 \leq \|f\|_2^2.$$

Therefore we will refer to $(\langle x, e_i \rangle)_{i \in I}$ as the Fourier coefficients of $x \in X$ and of

$$\sum_{i \in I} \langle x, e_i \rangle e_i$$

as the *Fourier series* of x .

3.1.3. Bounded operators between normed spaces. Mappings between vector spaces are of interest in a wide range of applications. We restrict our focus to mappings that respect the vector space structure: linear mappings aka linear operators.

DEFINITION 3.1.20. Let X, Y be vector spaces over the same scalar field \mathbb{F} . Then a mapping $T : X \rightarrow Y$ is *linear* if

$$T(x + \lambda y) = Tx + \lambda Ty$$

for all $x, y \in X$ and $\lambda \in \mathbb{F}$. We denote by $\mathcal{L}(X, Y)$ the set of all linear operators between X and Y .

Linear mappings are a special class of functions between two sets. Hence it has the structure of a vector space. Here are some examples of linear mappings for the classes of vector spaces of our interest.

- (1) Linear mappings between \mathbb{F}^n and \mathbb{F}^m are given by $m \times n$ matrices A with entries in \mathbb{F} , $x \mapsto Ax$ for $x \in \mathbb{F}^n$.
- (2) On the space of polynomials \mathcal{P}_n of degree at most n we define the *differentiation operator* $Dp(x) = a_1x + \dots + ma_nx^{n-1}$, the operator $p \mapsto \int p(x)dx$ and the evaluation operator $Tp(x) = p(0)$.
- (3) Operators on sequence spaces: For an element of the vector space s , a sequence $x = (x_n)_n$, we define the *left shift* $Lx = (0, x_0, x_1, x_2, \dots)$, the *right shift* $Rx = (x_1, x_2, \dots)$ and the multiplication operator $T_a x = (a_0x_0, a_1x_1, \dots)$ for a sequence $a = (a_0, a_1, \dots) \in s$. On the vector space of convergent sequences c we define $Tx = \lim_n x_n$ for $x = (x_n) \in c$.
- (4) Operators on function spaces: The set of continuous functions $C(I)$ on an interval of \mathbb{R} , popular choices for I are $[0, 1]$ and \mathbb{R} . For $f \in C(I)$ we define the *integral operator* $f \mapsto \int k(x, y)f(y)dx$ for a function k defined on $I \times I$, the kernel of the operator, and the evaluation operator $Tf(x) = f(a)$ for $a \in I$. For a differentiable continuous function f we are able to study the *differentiation operator* $Df(x) = f'(x)$.

Norms on these spaces provide a tool to understand the properties of these mappings via the notion of *operator norm* that measures the size of the measure of distortion of x induced by T : For normed spaces $(X, \|\cdot\|_X)$, $(Y, \|\cdot\|_Y)$ and a linear mapping $T : X \rightarrow Y$ we are interested in operators such that there exists a constant c such that

$$\|Tx\|_Y \leq c\|x\|_X \quad \text{for all } x \in X.$$

Often we will omit the subscripts to ease the notation. The operators with a finite c are of particular relevance and are called *bounded operators*. We denote by $\mathcal{B}(X, Y)$ the set of all bounded linear operators from X to Y .

DEFINITION 3.1.21. Let T be a linear operator between the normed spaces $(X, \|\cdot\|_X)$ and $(Y, \|\cdot\|_Y)$. The *operator norm* of T is defined by

$$\|T\| = \sup\left\{\frac{\|Tx\|_Y}{\|x\|_X} : \|x\|_X \neq 0\right\}.$$

Sometimes we denote the operator norm of T by $\|T\|_{\text{op}}$.

LEMMA 3.10. For $T \in \mathcal{B}(X, Y)$ the following quantities are all equal to the operator norm $\|T\|$ of T :

- (1) $C_1 = \inf\{c \in \mathbb{R} : \|Tx\|_Y \leq c\|x\|_X\}$,
- (2) $C_2 = \sup\{\|Tx\|_Y : \|x\|_X \leq 1\}$,
- (3) $C_3 = \sup\{\|Tx\|_Y : \|x\|_X = 1\}$.

PROOF. The argument is based on some inequalities:

- (1) $C_2 \leq C_1$: By definition of C_1 we have $\|Tx\| \leq C_1\|x\|$. Hence for all $x \in \overline{B}_1(0)$ we have $\|Tx\| \leq C_1$ and thus we have $C_2 \leq C_1$.
- (2) $C_3 \leq C_2$: For all $x \in \overline{B}_1(0)$ we have $\|Tx\| \leq C_2$. Pick an x with $\|x\| = 1$ and define the sequence of vectors $(x_n = (1 - 1/n)v)_n$ which all have $\|x_n\| \leq 1$ and hence $\|Tx_n\| \leq C_2$ for all $n \in \mathbb{N}$. Taking the limit gives $\|Tx\| \leq C_2$ and thus $C_3 \leq C_2$.
- (3) $\|T\| \leq C_3$: By definition of C_3 we have $\|Tx\| \leq C_3$ for all x with $\|x\| = 1$. Take an arbitrary non-zero vector $x \in X$. Then $x/\|x\|$ has unit length and hence $\|T(\frac{x}{\|x\|})\| = \frac{\|Tx\|}{\|x\|} \leq C_3$, which establishes the desired inequality $\|T\| \leq C_3$.
- (4) We have $\|Tx\|\|x\| \leq \|T\|$ for all $x \in X$. Hence $\|Tx\| \leq \|T\|\|x\|$ for all $x \in X$. Hence we have $C_1 \leq \|T\|$. Hence we have $C_1 \leq C_2 \leq C_3 \leq \|T\| \leq C_1$ and so the assertion is established. □

These different expressions for the operator norm of a linear operator are elementary but nonetheless useful. Before we discuss some examples we note some properties of the operator norm.

PROPOSITION 3.1.22. For $S, T \in \mathcal{B}(X, Y)$ we have

- (1) $\|I\| = 1$ for the identity operator $I : X \rightarrow X$.
- (2) $\|\lambda S + \mu T\| \leq |\lambda|\|S\| + |\mu|\|T\|$ for $\lambda, \mu \in F$.
- (3) *Submultiplicativity*: $\|S \circ T\| \leq \|S\|\|T\|$.

PROOF. (1) By the definition of the operator norm we have $\|I\| = 1$.

- (2) The triangle inequality for norms yields the assertion.

(3) By definition we have

$$\|S \circ T\| = \sup\{\|STx\| : \|x\| = 1\} \leq \sup\{\|S\|\|Tx\| : \|x\| = 1\} = \|S\|\|T\|.$$

□

PROPOSITION 3.1.23. *The vector space $\mathcal{B}(X, Y)$ of bounded operators between two normed spaces is a normed spaces with respect to the operator norm.*

PROOF. The preceding proposition implies the homogeneity property and the triangle inequality. The operator norm is clearly positive definite, and we have $\|T\| = 0$ if and only if $T = 0$ because it is defined in terms of a norm on Y . fined in terms of a norm on Y . □

We treat some of the operators defined above.

- (1) The right shift $Rx = (0, x_0, x_1, x_2, \dots)$ has $\|R\| = 1$ and also the left shift $Lx = (x_2, x_3, \dots)$ $\|L\| = 1$ on ℓ^∞ . For the multiplication operator $T_a x = (a_0 x_0, a_1 x_1, \dots)$ for a sequence $a = (a_0, a_1, \dots) \in s$ we have $\|T_a\| = \|a\|_\infty$ on ℓ^∞ . Let us look at the right shift operator. The operator norm is given by $\|R\| = \sup\{\|Rx\|_\infty : \|x\|_\infty = 1\}$:

$$\|Rx\|_\infty = 0 + |x_0|^2 + |x_1|^2 + \dots = \|x\|_\infty = \|x\|_\infty,$$

for all $x \in \ell^\infty$, hence $\|R\| = 1$. In a similar way one gets the norms of the other operators.

- (2) The operator norm of the integral operator $T_k f(x) = \int_a^b k(x, y) f(y) dy$ on $C[a, b]$ with $\|\cdot\|_\infty$ for an interval of finite length with a kernel $k \in C([a, b] \times [a, b])$ is $(b - a) \|k\|_\infty$. Note that

$$\begin{aligned} \|T_k f\|_\infty &= \sup\left\{\left|\int_a^b k(x, y) f(y) dy\right| : x \in [a, b]\right\} \\ &\leq \sup\left\{\int_a^b |k(x, y)| |f(y)| dy : x \in [a, b]\right\} \\ &\leq \|k\|_\infty \|f\|_\infty (b - a), \end{aligned}$$

so we have $\|T_k f\|_\infty \leq \|k\|_\infty \|f\|_\infty (b - a)$ for all non-zero $f \in C[a, b]$, i.e. $\|T_k\| \leq \|k\|_\infty (b - a)$. For the constant function $f(x) = 1$ for all $x \in [a, b]$ we get $\|T_k\| = 1$.

Some classes of operators on a normed space X : (i) *isometries* on X are linear operators T with $\|Tx\| = \|x\|$ for all $x \in X$, (ii) projections are linear operators P on X satisfying $P^2 = P$. A different way is to specify norms $\|\cdot\|_a$ and $\|\cdot\|_b$ on \mathbb{C}^n and \mathbb{C}^m , respectively. Then these norms induce a norm on $\mathcal{M}_{m \times n}(\mathbb{C})$, known as the *induced* norm. From a general perspective that is the operator norm of the induced linear transformation.

EXAMPLE 3.1.24. Let $A : \mathbb{C}^n \rightarrow \mathbb{C}^n$ be a linear operator given by a matrix $A = (a_{ij})$ and we put on both spaces the 1-norm. Let $A = (a_1 | \dots | a_n)$. Then $\|A\|_{\text{op}} = \max_{1 \leq j \leq n} \|a_j\|_1$, i.e. it is the maximum column sum. We have $Ax = \sum_{j=1}^n a_{ij} x_j$ and thus

$$\|Ax\|_{\text{op}} = \|Ax\|_1 \leq \sum_{j=1}^n |a_{ij}| |x_j| \leq \|x\|_1 \max_j \|a_j\|_1.$$

Hence $\max_{\|x\|_1=1} \|Ax\|_1 \leq \max_j \|a_j\|_1$.

Let e_j be the j th standard basis vector for \mathbb{C}^n . Then $\|A\|_{\text{op}} = \max_j \|a_j\|_1$.

Banach spaces and Hilbert spaces

4.1. Banach spaces and Hilbert spaces

We extend the topological notions introduced for the real line to general normed spaces and we focus on completeness in this section. Complete normed spaces are nowadays called Banach spaces, after the numerous seminal contributions of the Polish mathematician Stefan Banach to these objects. The class of complete innerproduct spaces are named after David Hilbert, who introduced the sequence space ℓ^2 . His students made numerous contributions to the theory of innerproduct spaces, e.g. Erhard Schmidt, Hermann Weyl, Otto Toeplitz,...

4.1.1. Completeness. We start with the generalization of open and closed intervals in \mathbb{R} to general normed spaces.

DEFINITION 4.1.1. Let $(X, \|\cdot\|)$ be normed space.

- (1) $B_r(x) = \{y \in X : \|y - x\| < r\}$ denotes the *open ball* of radius r around a point $x \in X$.
- (2) $\overline{B}_r(x) = \{y \in X : \|y - x\| \leq r\}$ denotes the *closed ball* of radius r around a point $x \in X$.

For the sequence spaces ℓ^p open balls $B_r(x)$ around $x = (x_k)$ are all sequences $y = (y_k) \in \ell^p$ with $\|x - y\| < r$. In the setting of $(C(I), \|\cdot\|)$ the ball $B_\varepsilon(f)$ are all continuous functions g that are in an ε -strip of f .

Here are the the notions of a convergent sequence and Cauchy sequence in a normed space.

DEFINITION 4.1.2. Let $(X, \|\cdot\|)$ be a normed space.

- (1) A sequence $(x_k)_{k \in \mathbb{N}}$ converges to $x \in X$ if for a given $\varepsilon > 0$ there exists a N such that $\|x - x_k\| < \varepsilon$ for $k \geq N$.
- (2) A sequence $(x_k)_{k \in \mathbb{N}}$ is a *Cauchy sequence* if for any $\varepsilon > 0$ there exists a N such that $\|x_m - x_n\| < \varepsilon$ for all $m, n > N$.

This notion of sequences is a natural generalization of the one for real and complex numbers. Note that the elements of the sequences are vectors in a normed space. For example, a sequence in ℓ^2 is a sequence where the elements themselves are also sequences. The difference between the the normed space \mathbb{Q} and the real numbers \mathbb{R} viewed as normed space is that not all Cauchy sequences in \mathbb{Q} converge to a rational number but that is the case for \mathbb{R} .

DEFINITION 4.1.3. A normed space $(X, \|\cdot\|)$ is called *complete* if every Cauchy sequence (x_k) in X has a limit x belonging to X . Moreover, a complete normed space is referred to as *Banach space* and a complete innerproduct space is known as *Hilbert space*.

Let us start with an elementary observation that is a straightforward consequence of the definitions.

LEMMA 4.1. *A subspace M of a Banach space is complete if and only if M is closed.*

THEOREM 4.2. *For $p \in [1, \infty]$ the normed space $(\mathbb{R}^n, \|\cdot\|_p)$ is complete.*

The infinite-dimensional counterpart of the previous is also true, but its proof is more intricate.

THEOREM 4.3. *For $p \in [1, \infty]$ the normed spaces $(\ell^p, \|\cdot\|_p)$ are complete.*

PROOF. We show the completeness of ℓ^1 and that of ℓ^∞ , since the arguments for $1 < p < \infty$ are analogous to the ones for ℓ^1 and the case of ℓ^∞ requires a slightly different reasoning. We discuss the case of real-valued sequences.

(1) *Completeness of ℓ^1* : The argument is split into three steps.

Step 1: Find a candidate for the limit. Let $(x_n)_n$ be a Cauchy sequence in ℓ^1 . We denote the n -th element of the sequence by $x_n = (x_1^{(n)}, x_2^{(n)}, \dots)$. Note that $|x_1^{(m)} - x_1^{(n)}| \leq \|x_m - x_n\|_1$, so the first coordinates $(x_1^{(n)})_n$ are a Cauchy sequence of real numbers and hence converge to some real number z_1 . Similarly, the other coordinates converge: $z_j = \lim_{n \rightarrow \infty} x_j^{(n)}$. Hence our candidate for the limit of (x_n) is the sequence $z = (z_1, z_2, \dots)$.

Step 2: Show that z is in ℓ^1 . We have that

$$\sum_{j=1}^N |z_j| = \sum_{j=1}^N \lim_n |x_j^{(n)}| = \lim_n \sum_{j=1}^N |x_j^{(n)}|,$$

where the interchange of the limit with the sum of a finite number of real numbers is no problem. Since Cauchy sequences are bounded, there is a constant $C > 0$ such that $\|x_n\|_1 < C$ for all n . Thus for any N

$$\sum_{j=1}^N |x_j^{(n)}| \leq \sum_{j=1}^{\infty} |x_j^{(n)}| = \|x_n\|_1 < C.$$

Letting $n \rightarrow \infty$ we find that

$$\sum_{j=1}^N |z_j| \leq \|x_n\|_1 < C$$

for arbitrary N . Hence we have $z \in \ell^1$.

Step 3: Show the convergence. We want to prove that $\|x_n - z\|_1 \rightarrow 0$ for $n \rightarrow \infty$.

Given $\varepsilon > 0$, pick N_1 so that if $m, n > N_1$ then $\|x_m - x_n\|_1 < \varepsilon$. Hence for any fixed N and $m, n > N_1$, we find

$$\sum_{j=1}^N |x_j^{(m)} - x_j^{(n)}| \leq \sum_{j=1}^{\infty} |x_j^{(m)} - x_j^{(n)}| = \|x_m - x_n\|_1 < \varepsilon.$$

Fix $n > N_1$ and N , let $m \rightarrow \infty$ to obtain

$$\sum_{j=1}^N |x_j^{(n)} - z_j| = \lim_{m \rightarrow \infty} \sum_{j=1}^N |x_j^{(n)} - x_j^{(m)}| \leq \varepsilon.$$

Since this is true for all N we have demonstrated that

$$\|x_n - z\|_1 < \varepsilon.$$

That is our desired conclusion.

(2) *Completeness of ℓ^∞* : The argument is split into three steps.

Step 1: Find a candidate for the limit. Let $(x_n)_n$ be a Cauchy sequence in ℓ^∞ . We denote the n -th element of the sequence by $x_n = (x_1^{(n)}, x_2^{(n)}, \dots)$.

Note that $|x_k^{(m)} - x_k^{(n)}| \leq \|x_m - x_n\|_\infty$ for all k and all $m, n > N$, so the k -th coordinates $(x_k^{(n)})_n$ are a Cauchy sequence of real numbers and hence converge to some real number z_k . Similarly, the other coordinates converge: $z_k = \lim_{m \rightarrow \infty} x_k^{(m)}$.

Hence our candidate for the limit of (x_n) is the sequence $z = (z_1, z_2, \dots)$.

Step 2: Show that z is in ℓ^∞ . We have that

$$\sup\{|z_j| : j = 1, \dots, N\} = \sup\{\lim_n |x_j^{(n)}| : j = 1, \dots, N\} = \lim_n \{\sup |x_j^{(n)}| : j = 1, \dots, N\},$$

where the interchange of the limit with the sum of a finite number of real numbers is no problem. Since Cauchy sequences are bounded, there is a constant $C > 0$ such that $\|x_n\|_\infty < C$ for all n . Thus for any N

$$\lim_n \{\sup |x_j^{(n)}| : j = 1, \dots, N\} \leq \|x_n\|_\infty < C.$$

Thus we find that $\|x_n\|_\infty < C$, i.e. we have $z \in \ell^\infty$.

Step 3: Show the convergence. We want to prove that $\|x_n - z\|_\infty \rightarrow 0$ for $n \rightarrow \infty$.

Given $\varepsilon > 0$, pick N_1 so that if $m, n > N_1$ then

$$|x_m^{(k)} - x_n^{(k)}| \leq \|x_m - x_n\|_\infty < \varepsilon$$

for all k . Taking limits as $m \rightarrow \infty$ we have

$$|z_k - x_n^{(k)}| \leq \varepsilon$$

Taking supremum in k , we obtain

$$\sup_k |z_k - x_n^{(k)}| \leq \varepsilon$$

for all $n > N_1$, i.e. $\|x_n - z\|_\infty \leq \varepsilon$ for all $n > N$. Consequently we have that x_n converges to z in $(\ell^\infty, \|\cdot\|_\infty)$. □

The completeness of the space of function spaces for a closed and bounded interval is of utmost importance in many arguments.

THEOREM 4.4. *For a finite interval $[a, b]$ the normed space $C[a, b]$ with respect to the sup-norm $\|\cdot\|_\infty$ is complete.*

For the proof we have to discuss notions of convergence for sequences of functions. Observe that the $\|f - g\|_\infty$ -norm measures the distance between two functions by looking at the point they are the furthest apart.

LEMMA 4.5. *For $f, g \in C[a, b]$ we have that $\sup\{|f(x) - g(x)| : x \in [a, b]\}$ is finite, and there is a $y \in [a, b]$ such that $d_\infty(f, g) = \sup\{|f(x) - g(x)| : x \in [a, b]\}$.*

PROOF. We show that $d(x) = |f(x) - g(x)|$ is continuous on $[a, b]$ and thus by the Extreme Value Theorem the assertion follows. The continuity of d is deduced from

$$|d(x) - d(y)| \leq ||f(x) - g(x)| - |f(y) - g(y)|| \leq |f(x) - f(y)| + |g(y) - g(x)|.$$

Since f and g are continuous at x there is for any given $\varepsilon > 0$ a $\delta > 0$ such that $|f(x) - f(y)| < \varepsilon/2$ and $|g(x) - g(y)| < \varepsilon/2$ for $|x - y| < \delta$. Hence

$$|d(x) - d(y)| \leq |f(x) - f(y)| + |g(y) - g(x)| < \varepsilon/2 + \varepsilon/2 = \varepsilon$$

for all $y \in [a, b]$ with $|x - y| < \delta$. Consequently d is continuous. \square

DEFINITION 4.1.4. Let (f_n) be a sequence of functions on a set X .

- We say that (f_n) *converges pointwise* to a limit function f if for a given $\varepsilon > 0$ and $x \in X$ there exists an N so that

$$|f_n(x) - f(x)| < \varepsilon \quad \text{for all } n \geq N.$$

- We say that (f_n) *converges uniformly* to a limit function f if for a given $\varepsilon > 0$ there exists an N so that

$$|f_n(x) - f(x)| < \varepsilon \quad \text{for all } n \geq N$$

holds for all $x \in X$.

There is a substantial difference between these two definitions. In pointwise convergence, one might have to choose a different N for each point $x \in X$. In the case of uniform convergence there is an N that holds for all $x \in X$. Note that uniform convergence implies pointwise convergence. If one draws the graphs of a uniformly convergent sequence, then one realizes that the definition amounts for a given $\varepsilon > 0$ to have a N so that the graphs of all the f_n for $n \geq N$, lie in an ε -band about the graph of f . In other words, the f_n 's get uniformly close to f . Hence uniform convergence means that the maximal distance between f and f_n goes to zero. We prove this assertion in the next proposition.

PROPOSITION 4.1.5. *Let (f_n) be a sequence of continuous functions on $[a, b]$. Then the following are equivalent:*

- (1) (f_n) *converges uniformly to f .*
- (2) $\sup\{|f_n(x) - f(y)| : x \in [a, b]\} \rightarrow 0$ *as $n \rightarrow \infty$.*

PROOF. Assertion (i) \Rightarrow (ii): Assume that (f_n) converges uniformly to f . Then for any $\varepsilon > 0$ there exists a N such that $|f_n(x) - f(x)| < \varepsilon$ for all $x \in [a, b]$ and all $n > N$. Hence $\sup\{|f_n(x) - f(y)| : x \in [a, b]\} \leq \varepsilon$ for all $n > N$. Since this holds for all $\varepsilon > 0$, we have demonstrated that $\sup\{|f_n(x) - f(y)| : x \in [a, b]\} \rightarrow 0$ for $n \rightarrow \infty$.

Assertion (ii) \Rightarrow (i): Assume that $\sup\{|f_n(x) - f(y)| : x \in [a, b]\} \rightarrow 0$ for $n \rightarrow \infty$. Given an $\varepsilon > 0$, there is a N such that $\sup\{|f_n(x) - f(y)| : x \in [a, b]\} < \varepsilon$ for all $n > N$. Thus we have $|f_n(x) - f(y)| < \varepsilon$ for all $x \in [a, b]$ and all $n > N$, i.e. (f_n) converges uniformly to f . \square

A reformulation of this result is that a sequence converges in $(C[a, b], \|\cdot\|_\infty)$ to f is equivalent to the uniform convergence of (f_n) to f .

PROPOSITION 4.1.6. *A sequence (f_n) converges to f in $(C[a, b], \|\cdot\|_\infty)$ if and only if (f_n) converges uniformly to f .*

Uniform convergence has an important property.

THEOREM 4.6. *Let (f_n) be a uniformly convergent sequence in $C(I)$ with limit f . Then the limit function f is continuous on I .*

PROOF. Let $y \in I$ and $\varepsilon > 0$ be given. By the uniform convergence of $f_n \rightarrow f$, there exists an N such that $n \geq N$ implies that

$$|f_n(x) - f(x)| \leq \varepsilon/3 \quad \text{for all } x \in I.$$

The continuity of f_N implies that there exists a $\delta > 0$ such that

$$|f_N(x) - f_N(y)| \leq \varepsilon/3 \quad \text{for } |x - y| \leq \delta.$$

We want to show that f is continuous. For all x such that $|x - y| < \delta$ we have that

$$\begin{aligned} |f(x) - f(y)| &\leq |f(x) - f_N(x)| + |f_N(x) - f_N(y)| + |f_N(y) - f(y)| \\ &< \varepsilon/3 + \varepsilon/3 + \varepsilon/3 = \varepsilon. \end{aligned}$$

□

Convergence of a sequence in $(C[a, b], \|\cdot\|_\infty)$ to $f \in C[a, b]$ is equivalent to uniform convergence of the sequence to f .

Finally we are in the position to prove our main theorem on continuous functions: Completeness of $(C[a, b], \|\cdot\|_\infty)$.

PROOF. Assume that (f_n) is a Cauchy sequence in $(C[a, b], \|\cdot\|_\infty)$. Then we have to show that there exists a function $f \in C[a, b]$ that has (f_n) as its limit.

Fix $x \in [a, b]$ and note that $|f_n(x) - f_m(x)| \leq \|f_n - f_m\|_\infty$. Since (f_n) is a Cauchy sequence $(f_n(x))$ is a Cauchy sequence in \mathbb{R} . Since \mathbb{R} is complete, $(f_n(x))$ converges to a point $f(x)$ in \mathbb{R} . In other words, $f_n \rightarrow f$ pointwise.

Next we show that $f \in C[a, b]$. Since (f_n) is a Cauchy sequence, we have for any $\varepsilon > 0$ a N such that $\|f_n - f_m\| < \varepsilon/2$ for all $m, n > N$. Hence we have $|f_n(x) - f_m(x)| < \varepsilon/2$ for all $x \in [a, b]$ and for all $m, n > N$. Letting $m \rightarrow \infty$ yields for all $x \in [a, b]$ and all $n > N$:

$$|f_n(x) - f(x)| = \lim_{m \rightarrow \infty} |f_n(x) - f_m(x)| \leq \varepsilon/2 < \varepsilon.$$

Consequently, $f_n \rightarrow f$ converges uniformly. Now by the preceding proposition f is a continuous function on $[a, b]$. In other words, we have established that $(C[a, b], \|\cdot\|_\infty)$ is a Banach space. □

THEOREM 4.7. *The normed space of bounded operators $(B(X, Y), \|\cdot\|_{\text{op}})$ is complete if and only if Y is a Banach space.*

The Banach space $(B(X, \mathbb{C}), \|\cdot\|_{\text{op}})$ is known as the *dual space* of X , denoted by X' , and its elements are referred to as *functionals* on X .

PROOF. Let (T_n) be a Cauchy sequence in $B(X, Y)$, so for any $\varepsilon > 0$ there exists a $N \in \mathbb{N}$ such that for all $m, n \geq N$ we have $\|T_m - T_n\|_{\text{op}} < \varepsilon$. Hence for any $x \in X$ we have

$$\|(T_m - T_n)x\|_Y \leq \|T_m - T_n\|_{\text{op}} \|x\|_X < \varepsilon \|x\|_X.$$

Hence for all $x \in X$ the sequence $(T_n x)$ is a Cauchy sequence in Y . Since Y is a Banach space, it has a limit denoted by Tx , and thus we define $Tx = \lim_{n \rightarrow \infty} T_n x$. The limit operator T is linear and bounded.

$$\|Tx\|_Y \leq \sup_n \|T_n x\|_Y \leq \|x\|_X \sup_n \|T_n\|_{\text{op}},$$

and thus we have $\|T\|_{\text{op}} \leq \sup_n \|T_n\|_{\text{op}}$, i.e. $T \in \mathcal{B}(X, Y)$.

We show that $\|T_n - T\|_{\text{op}} \rightarrow 0$. We assume otherwise that $\|T_n - T\|_{\text{op}}$ does not converge to 0. Then there exists an $\varepsilon > 0$ and a subsequence $(T_{n_k})_k$ of (T_n) such that

$$\|T_{n_k} - T\|_{\text{op}} \geq \varepsilon \quad \text{for all } k.$$

Consequently, for every k there exists a $x_k \in X$ with $\|x_k\| = 1$ and

$$\|T_{n_k}(x_k) - T_m(x_k)\| \geq \varepsilon.$$

By assumption (T_n) is a Cauchy sequence, so one can choose a N_0 such that for all $m, n_k \geq N_0$ we have

$$\|T_{n_k}(x_k) - T_m(x_k)\| \leq \varepsilon/2$$

and this gives

$$\varepsilon \leq \|T_{n_k}(x_k) - T(x_k)\|_Y \leq \|T_{n_k}(x_k) - T_m(x_k)\|_Y + \|T_m(x_k) - T(x_k)\|_Y.$$

Hence for all $m \geq N_0$ we have

$$\|T_m(x_k) - T(x_k)\|_Y \geq \varepsilon/2.$$

That is a contradiction to the definition of T , thus we have $T_m(x_k) - T(x_k) \rightarrow 0$ in $(\mathcal{B}(X, Y), \|\cdot\|_{\text{op}})$. \square

4.1.2. Equivalent norms. On a vector space X one may define different norms. We describe a way to compare these norms that respects basic properties, e.g. convergent sequences.

DEFINITION 4.1.7. Given a vector space X . Two metrics $\|\cdot\|_a$ and $\|\cdot\|_b$ are called *equivalent* if there exist (positive) constants C_1, C_2 such that

$$C_1 \|x\|_a \leq \|x\|_b \leq C_2 \|x\|_a \quad \text{for all } x \in X.$$

LEMMA 4.8. *Given three norms $\|\cdot\|_a, \|\cdot\|_b$ and $\|\cdot\|_c$ on X . Suppose $\|\cdot\|_a$ and $\|\cdot\|_c$ are equivalent and $\|\cdot\|_b$ and $\|\cdot\|_c$. Then $\|\cdot\|_a$ and $\|\cdot\|_b$ are equivalent norms.*

PROOF. We have constants C_1, C_2, C'_1, C'_2 such that

$$C_1 \|x\|_c \leq \|x\|_a \leq C_2 \|x\|_c$$

and

$$C'_1 \|x\|_c \leq \|x\|_b \leq C'_2 \|x\|_c.$$

Hence $\|x\|_c \leq C_1^{-1} \|x\|_a$ and thus

$$\|x\|_b \leq C'_2 C_1^{-1} \|x\|_a,$$

which by the second set of inequalities gives

$$\|x\|_b \leq C'_2 C_1^{-1} \|x\|_a.$$

In a similar way, we use

$$C'_1 C_2^{-1} \|x\|_a \leq C'_1 \|x\|_c$$

to obtain

$$C'_1 C_2^{-1} \|x\|_b \leq \|x\|_a$$

and thus $\|x\|_b$ and $\|x\|_a$ are equivalent

$$C_1' C_2^{-1} \|x\|_a \leq \|x\|_b \leq C_2' C_1^{-1} \|x\|_a.$$

□

PROPOSITION 4.1.8. *Suppose $\|\cdot\|_a$ and $\|\cdot\|_b$ are equivalent metrics on X . If a sequence (x_n) converges with respect to the $\|\cdot\|_a$, then it converges with respect to $\|\cdot\|_b$.*

PROOF. Suppose $\lim \|x_n - x\|_a = 0$. Then we have

$$C_1 \|x_n - x\|_a \leq \|x_n - x\|_b \leq C_2 \|x_n - x\|_a$$

and hence $\lim \|x_n - x\|_b = 0$. □

LEMMA 4.9. *Suppose $\|\cdot\|_a$ and $\|\cdot\|_b$ are equivalent metrics on X . Then $(X, \|\cdot\|_a)$ and $(X, \|\cdot\|_b)$ are both Banach spaces or the two normed spaces are incomplete.*

On a finite-dimensional vector space X all norms are equivalent. Since any real finite-dimensional vector space X is isomorphic to \mathbb{R}^n (after a choice of basis and using the coefficient mapping and synthesis mapping as isomorphism) we just have to show this statement for \mathbb{R}^n .

THEOREM 4.10. *All norms on \mathbb{R}^n are equivalent.*

PROOF. By Lemma 4.8 it suffices to show a norm $\|\cdot\|$ on \mathbb{R}^n is equivalent to a fixed norm. We fix the $\|\cdot\|_1$ on \mathbb{R}^n . Suppose e_1, \dots, e_n is a basis for \mathbb{R}^n . Then any $x \in \mathbb{R}^n$ has a unique expansion

$$x = \sum_{i=1}^n a_i e_i$$

and its 1-norm is defined by

$$\|x\|_1 = \sum_{i=1}^n |a_i|.$$

The proof may be broken up into four steps. Step 1 is the reduction of the general case to the situation that we have to show that $\|\cdot\|$ is equivalent to $\|\cdot\|_1$. Step 2 is the elementary observation that it suffices to check the desired assertion

$$C_1 \|x\|_1 \leq \|x\| \leq C_2 \|x\|_1$$

not for all $x \in X$ but just for elements in the unit ball of $\|\cdot\|_1$. Namely, the preceding inequalities are true for $x = 0$. Let us assume $x \neq 0$. Then we can divide the inequalities by $\|x\|_1$:

$$C_1 \leq \|x/\|x\|_1\| \leq C_2.$$

Since the elements we have to check our inequalities are now in $B_1(0)$ defined by the $\|\cdot\|_1$.

The next step paves the way to make the problem accessible to methods from analysis. Step 4: $\|\cdot\|$ is continuous under $\|\cdot\|_1$. Explicitly, we have to show that for a given $\varepsilon > 0$ there exists a $\delta > 0$ such that $\|x - x'\|_1 < \delta$ implies that $|\|x\| - \|x'\|| < \varepsilon$. We know that

$$|\|x\| - \|x'\|| \leq \|x - x'\|.$$

Let us relate the $\|\cdot\|$ with $\|\cdot\|_1$. We represent x and x' with respect to the basis $\{e_1, \dots, e_n\}$:

$$x = \sum_{i=1}^n a_i e_i \quad \text{and} \quad x' = \sum_{i=1}^n a'_i e_i.$$

The triangle inequality implies

$$\|x - x'\| \leq \sum_{i=1}^n |a_i - a'_i| \|e_i\| \leq (\max_i \|e_i\|) \|x - x'\|_1.$$

Choose $\delta = \varepsilon / \max_i \|e_i\|$. Then we get the desired statement: If $\|x - x'\|_1 < \delta$, then

$$\| \|x\| - \|x'\| \| \leq \|x - x'\| \varepsilon.$$

The final step is to use the the Extreme Value Theorem for the continuous function $\|\cdot\|$ on \mathbb{R}^n and note that the set $\{x \in X : \|x\|_1 = 1\}$ is closed and bounded. Then $\|\cdot\|$ has to achieve its minimum and maximum on the unit ball for the 1-norm:

$$C_1 := \max\{\|x\| : \|x\|_1 = 1\} \quad \text{and} \quad C_2 := \min\{\|x\| : \|x\|_1 = 1\}.$$

By definition we have $C_2 \geq C_1$ and hence

$$C_1 \leq \|x\| \leq C_2$$

for $x \in X$ with $\|x\|_1 = 1$. □

In the infinite-dimensional setting one has norms on vector spaces that are not equivalent. Let us take the space of continuous functions $C[0, 1]$ and complete it with respect to $\|\cdot\|_2$ and $\|\cdot\|_\infty$. Then you have seen in the exercises that $(C[0, 1], \|\cdot\|_2)$ is not complete, but $(C[0, 1], \|\cdot\|_\infty)$ is a Banach space.

A consequence of the equivalence of norms on \mathbb{R}^n is that a sequence in \mathbb{R}^n converges in norm if and only if converges coordinate-wise.

PROPOSITION 4.1.9. *Let $\|\cdot\|$ be a norm on \mathbb{R}^n , and (x_j) a sequence in \mathbb{R}^n . Then $\|x_j - x\| \rightarrow 0$ if and only if $x_j^{(i)} \rightarrow x^{(i)}$ for $i = 1, \dots, n$.*

PROOF. (\Leftarrow) Since all norms are equivalent we are allowed to pick a norm most appropriate for our problem. We pick the sup-norm.

Suppose $\lim_j \|x_j - x\| = 0$. Denote the components of x_j by $x_j = (x_j^{(1)}, \dots, x_j^{(n)})$. Then $x_j^{(i)}$ converges to $x^{(i)}$ for $i = 1, \dots, n$.

(\Rightarrow) For this direction we use the 1-norm. Suppose $x_j^{(i)} \rightarrow x^{(i)}$ for $i = 1, \dots, n$. Then $\|x_j - x\|_1 = \sum_{j=1}^n |x_j^{(i)} - x^{(i)}| \rightarrow 0$. □

4.1.3. Banach's Fixed Point Theorem aka Contraction Mapping Theorem. In 1922 Banach established a theorem on the convergence of iterations of contractions that has become a powerful tool in applied and pure mathematics. Suppose we have a bounded operator T acting on a normed space X . Take a point x_0 in X and build the sequence of iterates $x_0, x_1 = Tx_0, x_2 = Tx_1 = T^2x_0, \dots, x_{n+1} = Tx_n$. The basic question is about the existence of the limit of this sequence $x = \lim_n x_n = \lim_n T^n x_0$. The limit x of the iterates (x_n) is a fixed point of the continuous map T :

$$T(x) = T(\lim_n x_n) = \lim_n T(x_n) = \lim_n x_{n+1} = \lim_n x_n = x.$$

A mapping on a normed space X is called a *contraction* if there exists a $0 < K < 1$ such that

$$\|Tx - Ty\| \leq K\|x - y\| \quad x, y \in X.$$

EXAMPLE 4.1.10. Let T be a bounded linear operator on a normed space X . If $\|T\| < 1$, then T is a contraction on X . By assumption we have $\|Tx - Tx'\| = \|T(x - x')\| \leq \|T\|(\|x - x'\|) < \|x - x'\|$ for all $x, x' \in X$.

THEOREM 4.1.11 (Banach Fixed Point). *Let M be a closed subspace of a Banach space X . Any contraction f on M has a unique fixed point \tilde{x} and the fixed point is the limit of every sequence generated from an arbitrary nonzero point $x_0 \in M$ by iteration $(x_n)_n$, where $x_{n+1} = f(x_n)$ for $n \geq 1$.*

REMARK 4.1.11. Open and closed sets are defined in the following section.

PROOF. Let $x_0 \in M$ be arbitrary. Define $x_{n+1} = f(x_n)$ for $n = 1, 2, \dots$. By the contractivity of T we have

$$\|x_n - x_{n-1}\| = \|f(x_{n-1}) - f(x_{n-2})\| \leq K\|x_{n-1} - x_{n-2}\|$$

and iterations yields

$$\|x_n - x_{n-1}\| \leq K^{n-1}\|x_{n-1} - x_{n-2}\|.$$

The existence of a fixed point is based on the completeness of M . Hence we proceed to show that $(x_n)_n$ is a Cauchy sequence. Let m, n be greater than N and we choose $m \geq n$. Then by the preceding inequality and the triangle inequality we have

$$\begin{aligned} \|x_m - x_n\| &\leq \|x_m - x_{m-1}\| + \|x_{m-1} - x_{m-2}\| + \dots + \|x_{n+1} - x_n\| \\ &\leq (K^{m-1} + K^{m-2} + \dots + K^n)\|x_1 - x_0\| \\ &\leq (K^N + K^{N+1} + \dots)\|x_1 - x_0\| \\ &= K^N(1 - K)^{-1}\|x_1 - x_0\|. \end{aligned}$$

Since $0 < K < 1$, $\lim_N K^N = 0$ and thus (x_n) is a Cauchy sequence. Consequently, (x_n) converges to a point \tilde{x} by the completeness of X . Furthermore \tilde{x} is a fixed point by the contractivity of T .

Uniqueness: Suppose there is another fixed point \tilde{y} of T . Then $\|\tilde{x} - \tilde{y}\| = \|T\tilde{x} - T\tilde{y}\| \leq K\|\tilde{x} - \tilde{y}\|$ and $\|\tilde{x} - \tilde{y}\| > 0$. Thus we deduce that $K \geq 1$ which is a contradiction to the contractivity of T . \square

Two well-known applications are Newton's method for finding roots of general equations and the theorem of Picard-Lindelöf on the existence of solutions of ordinary differential equations.

Newton's method:

How does one compute $\sqrt{3}$ up to a certain precision, i.e. we are interested in error estimates? Idea: Formulate it in the form $x^2 - 3 = 0$ and try to use a method that allows to compute zeros of general equations.

Newton came up with a method to solve $g(x) = 0$ for a differentiable function

$g : I \rightarrow \mathbb{R}$.

Suppose x_0 is an approximate solution or starting point. Define recursively

$$x_{n+1} = x_n - \frac{g(x_n)}{g'(x_n)} \quad \text{for } n \geq 0.$$

Then (x_n) converges to a solution \tilde{x} , provided certain assumptions on g hold.

If $x_n \rightarrow \tilde{x}$, then by continuity of g we get $g(\tilde{x}) = 0$.

When does Newton's method lead to a convergent sequence of iterates? Idea: Apply Banach's Fixed Point Theorem.

Set $f(x) := x - \frac{g(x)}{g'(x)}$. Then given $x_0 \in I$ and $x_{n+1} = x_n - \frac{g(x_n)}{g'(x_n)} = f(x_n)$. Moreover, $f(\tilde{x}) = \tilde{x}$ if and only if $g(\tilde{x}) = 0$.

Let us restrict our discussion to the computation of $\sqrt{3}$. The Banach space X is the space of real numbers \mathbb{R} and $g(x) = x^2 - 3$, so

$$f(x) = x - \frac{x^2 - 3}{2x} = \frac{1}{2}\left(x + \frac{3}{x}\right)$$

on $[\sqrt{3}, \infty) \rightarrow [\sqrt{3}, \infty)$. Note that $[\sqrt{3}, \infty)$ is a closed set of \mathbb{R} containing $\sqrt{3}$. For $x \geq 0$ we have $\frac{1}{2}(x + 3/x) \geq \sqrt{3x/x} = \sqrt{3}$. Compute f' and note that a differentiable function $f : I \rightarrow \mathbb{R}$ with a bounded derivative is Lipschitz continuous with constant L (Homework):

$$f'(x) = \frac{1}{2}\left(1 - \frac{3}{x^2}\right)$$

and note that its range is contained in $[0, 1/2]$ for $x \geq \sqrt{3}$. Hence we have $L = 1/2$ and by Banach's Fixed Point Theorem $\frac{1}{2}(x_n + \frac{3}{x_n}) \rightarrow \sqrt{3}$.

Let's pick $x_0 = 2$ and thus $x_1 = 7/4$ and so $|x_1 - x_0| = 1/4$. Furthermore, we have

$$|x_n - \sqrt{3}| \leq \frac{(1/2)^n}{1 - 1/2} |x_1 - x_0| = \frac{1}{2^n} \cdot 2 \cdot \frac{1}{4} = \frac{1}{2^{n+1}}.$$

Hence

$$|x_n - \sqrt{3}| \leq \frac{1}{2^{n+1}}.$$

For $n = 4$, we have $|x_n - \sqrt{3}| \leq 1/1024 < 0.001$.

Existence and uniqueness of solutions of an ordinary differential equation (ODE) – Picard-Lindelöf Theorem.

Consider the following general initial value problem:

$$(4.1) \quad x'(t) = \frac{dx}{dt} = f(t, x) \quad \text{and } x(t_0) = x_0$$

for a function $f : A \subset \mathbb{R}^2 \rightarrow \mathbb{R}$ with $t_0 \in I$.

DEFINITION 4.1.12. Let I be an interval and $t_0 \in I$. A differentiable function $x : I \rightarrow \mathbb{R}$ is a *solution* of the IVP (4.1) if for all $t \in \mathbb{R}$ we have $x'(t) = f(t, x(t))$ and $x(t_0) = x_0$.

We say that the IVP has a *local solution* if there exists a $\delta > 0$ such that (4.1) has a solution x on $(x_0 - \delta, x_0 + \delta)$

EXAMPLE 4.1.13. The IVP $x'(t) = rx$, $x(0) = A$ has as solution $x(t) = Ae^{rt}$ on \mathbb{R} .

Now we can state the theorem of Picard-Lindelöf and in its proof we will also show how to construct approximately a solution to IVPs.

THEOREM 4.12 (Picard-Lindelöf). *Consider the initial value problem:*

$$(4.2) \quad x'(t) = \frac{dx}{dt} = f(t, x) \quad \text{and } x(t_0) = x_0,$$

where $f : U \times V \rightarrow \mathbb{R}$ is a function, U, V are intervals with t_0 in the interior of U and x_0 in the interior of V .

Assume that f is continuous and uniformly Lipschitz in x :

$$|f(t, x) - f(t, x')| \leq L|x - x'| \quad \text{for all } t \in U, x, x' \in V.$$

Then the IVP has a unique local solution.

PROOF. We start with a more precise formulation of the assumptions on f .

We have that f is a continuous function defined $f : U \times V \rightarrow \mathbb{R}$ on the intervals $U = [t_0 - a, t_0 + a]$, $V = [x_0 - b, x_0 + b]$ for $a, b > 0$, such that

$$|f(t, x) - f(t, x')| \leq L|x - x'| \quad \text{for all } t \in U, x, x' \in V.$$

The assumptions on f imply that it is bounded, i.e. there exists a $M > 0$ such that $|f(t, x)| \leq M$ for all $(t, x) \in U \times V$. Hence, the theorem of Picard-Lindelöf asserts that for $\delta < \min a, 1/L, b/M$ the IVP has a solution on $[t_0 - \delta, t_0 + \delta]$.

A key step in the proof is the reformulation of the theorem in terms of an integral equation.

LEMMA 4.13. *The IVP has a solution if and only if*

$$x(t) = x_0 + \int_{t_0}^t f(s, x(s))ds.$$

PROOF. We define φ on U by $\varphi(t) = f(t, x(t))$. By the Fundamental Theorem of Analysis $x_0 + \int_{t_0}^t \varphi(s)ds$ is the anti-derivative of f whose value at t_0 is x_0 . \square

The next step is an iterative procedure to solve the integral equation, also known as *Picard iteration*.

We define an operator Φ by

$$\Phi(x)(t) = x_0 + \int_{t_0}^t f(s, x(s))ds.$$

Then x solves the integral equation

$$x(t) = x_0 + \int_{t_0}^t f(s, x(s))ds.$$

if and only if $\Phi(x) = x$. We are going to specify the space of functions on which Φ acts later.

Consequently, we have reduced the IVP to finding a fixed point for Φ . The latter will be done with the help of an iteration scheme, the Picard iterations.

$$x_0(t) := x_0, \quad x_{n+1} := x_n + \int_{t_0}^t f(s, x_n(s)) ds \quad , n \geq 1,$$

or equivalently

$$x_0(t) := x_0 \quad x_{n+1} = \Phi(x_n).$$

Choose a δ such that $\delta < \min a, 1/L, b/M$ and consider the Banach space $X = (C[t_0 - \delta, t_0 + \delta], \|\cdot\|_\infty)$. As closed subset of X we pick

$$A = \{x \in C[t_0 - \delta, t_0 + \delta] : x(t) \in [x_0 - b, x_0 + b] \text{ for all } t\}.$$

Let us show that A is closed in X .

Suppose $(x_n) \subset A$ converges to $x \in X$ wrt $\|\cdot\|_\infty$. Then $x_n(t) \rightarrow x(t)$ for all t . For a fixed t we have $x_n(t) \in [x_0 - b, x_0 + b]$ which converges to $x(t)$ with values in $[x_0 - b, x_0 + b]$.

Now we show that for $x \in A$ also $\Phi(x) \in A$. Since $x \in A$ we have

$$x(t) \in [x_0 - b, x_0 + b] \text{ for all } t \in [t_0 - \delta, t_0 + \delta],$$

so we have $|x(t) - x_0| \leq b$ for all $t \in [t_0 - \delta, t_0 + \delta]$.

Consider

$$|\Phi(x)(t) - x_0| = \left| \int_{t_0}^t f(s, x(s)) ds \right| \leq \int_{t_0}^t |f(s, x(s))| ds \leq M|t - t_0|,$$

which yields that

$$|\Phi(x)(t) - x_0| \leq M\delta \text{ for } \delta < b/M.$$

Finally, we demonstrate that Φ is a contraction on A . Concretely, there exists a constant $q < 1$ such that

$$\|\Phi(x) - \Phi(y)\|_\infty \leq q\|x - y\|_\infty$$

for $x, y \in A$. Hence we have to get some control of the term $|\Phi(x)(t) - \Phi(y)(t)|$:

$$\begin{aligned} |\Phi(x)(t) - \Phi(y)(t)| &\leq \int_{t_0}^t |f(s, x(s)) - f(s, y(s))| ds \\ &\leq \int_{t_0}^t L|x(s) - y(s)| ds \\ &\leq \int_{t_0}^t L\|x - y\|_\infty ds \\ &\leq L|t - t_0|\|x - y\|_\infty \\ &\leq \delta L\|x - y\|_\infty. \end{aligned}$$

Hence we have

$$\|\Phi(x) - \Phi(y)\|_\infty \leq \delta L\|x - y\|_\infty,$$

so $q = \delta L < 1$.

Application of Banach's Fixed Point Theorem yields that there exists a unique $\tilde{x} \in A$ such that

$$\tilde{x}(t) = \tilde{x}_0 + \int_{t_0}^t f(s, \tilde{x}(s)) ds.$$

□

EXAMPLE 4.1.14. Consider the following IVP:

$$x'(t) = \sin(tx), \quad x(0) = 1.$$

Thus $|f(x, t)| = |\sin(tx)| \leq 1$, i.e. $M = 1$.

$\frac{\partial}{\partial x} f(t, x) = |t \cos(tx)| \leq |t| \leq a$, so $L = a$ and $\delta < \min\{a, 1/\delta, b\}$. For $a = b = 1$ we get $\delta < 1$. We have $t_0 = 1$ and $x_0 = 1$.

Choose $x_0(t) = 1$ and so $x_1(t) = 1 + \int_0^t \sin(s) ds = 1 - \cos t$, $x_2(t) = 1 + \int_0^t \sin(1 - \cos(s)) ds$. Note that x_2 is hard to compute analytically, but there are methods based on numerical integration.

In the next example we show that the assumption of continuity of f cannot be weakened.

EXAMPLE 4.1.15. Consider the IVP $x'(t) = f(x, t)$ for

$$f(t, x) = \begin{cases} 1 & \text{if } t \geq 0 \\ 0 & \text{if } t < 0 \end{cases}$$

and $x(0) = 0$. Then we have

$$x(t) = \begin{cases} t + c & \text{if } t \geq 0 \\ c_2 & \text{if } t < 0 \end{cases}$$

and thus

$$x(t) = \begin{cases} t & \text{if } t \geq 0 \\ 0 & \text{if } t < 0 \end{cases}$$

Hence x is not differentiable at 0. Consequently, the IVP has no solution.

4.1.4. Hilbert spaces. Banach spaces arising from innerproduct spaces are known as Hilbert spaces. These are easier to handle than general Banach spaces.

DEFINITION 4.1.16. A Hilbert space is an innerproduct space $(X, \langle \cdot, \cdot \rangle)$ such that the induced norm $\|\cdot\| = \langle \cdot, \cdot \rangle$ is complete.

Let M be a subspace of X . Denote by M^\perp , its *orthogonal complement*, the set of all $x \in X$ that are orthogonal to all the elements of M . Formally we have

$$M^\perp = \{x \in X : \langle x, y \rangle = 0 \text{ for all } y \in M\}.$$

The linearity of an innerproduct implies that M is a vector space.

LEMMA 4.14. Let M be a subspace of $(X, \langle \cdot, \cdot \rangle)$. Then M^\perp is a closed subspace of X .

PROOF. Let (x_n) be a sequence in M^\perp converging to $x \in X$. We have to show that $x \in M^\perp$. Since $\langle x_n, y \rangle = 0$ for all $y \in M$ we note that

$$|\langle x_n - x, y \rangle| \leq \|x_n - x\| \|y\| \rightarrow 0.$$

Hence we have

$$\langle x_n, y \rangle \rightarrow \langle x, y \rangle,$$

but $\langle x_n, y \rangle = 0$ for all n . Consequently, $\langle x, y \rangle = 0$ and so $x \in M^\perp$. □

By definition of M^\perp we have that M and M^\perp are disjoint subspaces of X . For any proper closed subspace M of X its orthogonal complement M^\perp is non-empty and there are sufficiently many elements in M^\perp that allows one to decompose elements in X with respect to M and M^\perp . The precise formulations of these facts and their proofs are the main parts of our treatment of Hilbert spaces.

The best approximation property holds for proper closed subspaces of Hilbert spaces.

THEOREM 4.15 (Best Approximation Theorem). *Suppose M is a proper closed subspace of a Hilbert space X . Then for any $x \in X$ there exists a unique element $z \in M$ such that*

$$\|x - z\| = \inf_{m \in M} \|x - m\|.$$

The quantity $\inf_{m \in M} \|x - m\|$ measures the distance of x from M . In the chapter on metric spaces we show that it defines an honest metric on X .

REMARK 4.1.17. In general the theorem is not true in Banach spaces. Take ℓ^∞ and as closed subspace c_0 , the space of sequences converging to zero. For $x = (1, 1, 1, \dots)$ there exists no sequence in c_0 attaining the minimal distance 1.

PROOF. Denote by $d = \inf_{m \in M} \|x - m\|^2$. Note that d is finite, since the real numbers $\|x - m\|$ for $m \in M$ are all nonnegative and bounded below by 0. Since d is the greatest lower bound of this set, there exists a sequence $(m_k) \subset M$ such that for each $\varepsilon > 0$ there exists an N such that $\|x - m_k\|^2 \leq d + \varepsilon$ for all $k \geq N$.

Claim: The sequence (m_k) is a Cauchy sequence. Applying the parallelogram identity to $x - m_k$ and $x - m_l$ we get

$$\|2x - m_k - m_l\|^2 + \|m_k - m_l\|^2 = 2(\|x - m_k\|^2 + \|x - m_l\|^2),$$

which yields to

$$\|x - \frac{m_k + m_l}{2}\|^2 + \|m_k - m_l\|^2/2 = (\|x - m_k\|^2 + \|x - m_l\|^2)/2.$$

Since $\frac{m_k + m_l}{2} \in M$ we have $\|x - \frac{m_k + m_l}{2}\|^2 \geq d$ and so we have

$$\|m_k - m_l\|^2 \leq 2(\|x - m_k\|^2 + \|x - m_l\|^2) - 4d.$$

For any $\varepsilon > 0$ there exists a N such that $\|x - m_k\|^2 \leq d + \varepsilon/4$ for all $k \geq N$. Then we have for all $m, m \geq N$ that

$$\|m_k - m_l\|^2 \leq 2(\|x - m_k\|^2 + \|x - m_l\|^2) - 4d \leq \varepsilon.$$

Hence we have demonstrated that (m_k) is a Cauchy sequence. Since M is closed, (m_k) converges to some element $z \in M$ and we have that $\|x - z\|^2 = d$ and so z is the vector in M closest to x . We have established the existence of a closest vector. The uniqueness goes as follows: Suppose there is another element $y \in M$ such that $\|x - y\|^2 = d$. Consider the sequence (y, z, y, z, \dots) , and note that it is a Cauchy sequence by the same argument as for (m_k) . Hence $y = z$ and so z is the unique solution to our approximation problem. \square

There is a characterization of best approximations in Hilbert spaces in terms of the orthogonal complement.

THEOREM 4.16 (Characterization of Best Approximation). *Suppose M is a proper closed subspace of a Hilbert space X . Then for any $x \in X$ there exists a best approximation $\tilde{x} \in M$ if and only if $x - \tilde{x} \in M^\perp$.*

PROOF. *First step:* Suppose $x - \tilde{x} \in M^\perp$. Then for any $y \in M$ with $y \neq \tilde{x}$ we have $\|y - x\|^2 = \|y - \tilde{x} + \tilde{x} - x\|^2$. Note that $y - \tilde{x} \in M$ and $\tilde{x} - x \in M^\perp$ so we have $\langle y - \tilde{x}, \tilde{x} - x \rangle = 0$. Hence Pythagoras yields $\|y - x\|^2 = \|y - \tilde{x}\|^2 + \|\tilde{x} - x\|^2$. By assumption $y - \tilde{x} \neq 0$ so we arrive at the desired assertion $\|y - x\|^2 > \|\tilde{x} - x\|^2$. *Second step:* Suppose \tilde{x} minimizes $\|x - \tilde{x}\|$. We assume that there exists a $y \in M$ of unit length such that $\langle x - \tilde{x}, y \rangle = \delta \neq 0$. Consider the element $z = \tilde{x} + \delta y$.

$$\begin{aligned} \|x - z\|^2 &= \|x - \tilde{x} - \delta y\|^2 \\ &= \langle x - \tilde{x}, x - \tilde{x} \rangle + \langle x - \tilde{x}, \delta y \rangle - \langle \delta y, x - \tilde{x} \rangle + \langle \delta y, \delta y \rangle \\ &= \|x - \tilde{x}\|^2 - |\delta|^2 - |\delta|^2 + |\delta|^2 \\ &= \|x - \tilde{x}\|^2 - |\delta|^2. \end{aligned}$$

Thus we have $\|x - z\|^2 \leq \|x - \tilde{x}\|^2$. Contradiction to the assumption that \tilde{x} minimizes $\|x - \tilde{x}\|$. \square

THEOREM 4.17 (Projection Theorem). *Let M be a closed subspace of a Hilbert space X . Then every $x \in X$ can be uniquely written as $x = y + z$ where $y \in M$ and $z \in M^\perp$.*

PROOF. For $x \in X$ there exists a best approximation $y \in M$. Note that $x = y + x - y$ with $y \in M$ and $x - y \in M^\perp$. Furthermore we have $M \cap M^\perp = \{0\}$ (if $x \in M \cap M^\perp$, then $\langle x, x \rangle = 0 = \|x\|^2$ and thus $x = 0$.) which completes the proof. \square

COROLLARY 4.1.18. *Let M be proper closed subspace of a Hilbert space X . Then $M^\perp \neq \{0\}$.*

PROOF. If $x \neq 0$, then the decomposition $x = y + z$ has a $z \neq 0$. Since $z \in M^\perp$ we have $M^\perp \neq \{0\}$. \square

Recall that a *projection* on a normed space X is a linear mapping $P : X \rightarrow X$ satisfying $P^2 = P$.

Here is a reformulation of the preceding theorem in terms of projections, justifying the name.

PROPOSITION 4.1.19. *For any closed subspace M of a Hilbert space X , there is a unique projection P on X satisfying:*

- (1) $\text{ran}(P) = M$ and $\text{ran}(I - P) = M^\perp$.
- (2) $\|Px\| \leq \|x\|$ for all $x \in X$. Moreover, $\|P\| = 1$.

PROOF. (1) The decomposition of $x \in X$ into $x = y + z$ for $y \in M, z \in M^\perp$ allows one to define $Px := y$. By definition $\text{ran}(P) \subseteq M$ and if $x \in M$, then $Px = x$. Thus $P^2 = P$ and $M \subseteq \text{ran}(P)$.

Once more, by $x = y + z$ we have $(I - P)x = z \in M^\perp$ and as above we deduce that $\text{ran}(I - P) = M^\perp$.

- (2) By Pythagoras we have $\|x\|^2 = \|Px\|^2 + \|z\|^2$ and thus we have $\|Px\| \leq \|x\|$. Hence $\|P\| \leq 1$. On the other hand, there exists $x \in X$ with $Px \neq 0$ and $\|P(Px)\| = \|Px\|$, so that $\|P\| \geq 1$. Hence we conclude that $\|P\| = 1$. \square

EXAMPLE 4.1.20. Let M be the line $\{t\xi : t \in \mathbb{R}\}$ given by a unit vector $\xi \in X$. Then

$$P_\xi x = \langle \xi, x \rangle \xi$$

projects a vector orthogonally onto its component in direction ξ

We state some consequences of the projection theorem. In the mathematics literature the tensor product notation $\xi \otimes \xi$ is used to refer to P_ξ .

PROPOSITION 4.1.21. *Let X be a Hilbert space.*

- (1) *For any closed subspace M of X we have $M^{\perp\perp} = M$.*
(2) *For any set A in X we have $A^{\perp\perp} = \overline{\text{span}(A)}$.*

PROOF. (1) For any $x \in M$ we have $\langle x, y \rangle = 0$ for every $y \in M^\perp$. In other words, x is orthogonal to M^\perp , so $x \in (M^\perp)^\perp$.

Conversely, suppose that $x \in M^{\perp\perp}$. Since M is closed, we can decompose $x = y + z$ with $y \in M$ and $z \in M^\perp$. Since $x \in M^{\perp\perp}$ we have $\langle x, z \rangle = 0$. Furthermore, we have $x \in M \subseteq M^{\perp\perp}$, so we also have $\langle x, y \rangle = 0$. Consequently, $\|z\|^2 = \langle z, z \rangle = \langle x - y, z \rangle = \langle x, z \rangle - \langle y, z \rangle = 0$. Hence $z = 0$ and we have deduced that $x \in M$.

- (2) For a general set A in X we note that $\overline{\text{span}(A)}$ is the smallest closed subspace containing A . We set $M = \text{span}(A)$. Then we have $M \subset \overline{M}$ and thus $\overline{M}^\perp \subseteq M^\perp$. Consequently, $M^{\perp\perp} \subseteq \overline{M}^{\perp\perp}$. But \overline{M} is closed in X so $\overline{M}^{\perp\perp} = M^{\perp\perp}$. Since $\overline{M}^{\perp\perp} = M^{\perp\perp}$ we get that $M^{\perp\perp} \subseteq \overline{M}^{\perp\perp}$. Finally, $M \subseteq M^{\perp\perp}$ and $M^{\perp\perp}$ closed implies $\overline{M} \subseteq M^{\perp\perp}$, which completes the argument. \square

COROLLARY 4.1.22. *A subset A in a Hilbert space X is dense if and only if $A^\perp = \{0\}$. Moreover, $A^\perp = \{0\}$ is equivalent to x orthogonal to A and hence $x = 0$. In words, $\overline{\text{span}(A)} = X$ if and only if the only element orthogonal to every element in A is the zero vector.*

PROOF. Suppose $\overline{\text{span}(A)} = X$. Then A is a closed linear subspace and hence $A^\perp = A^{\perp\perp\perp} = X^\perp = 0$.

Conversely, $\overline{\text{span}(A)} = A^{\perp\perp} = 0^\perp = X$. \square

Many interesting theorems in analysis are about the identification of the dual spaces of normed spaces. A topic one is at the heart of functional analysis. Here we restrict our focus to the Hilbert space setting since its proof relies on the projection theorem.

Recall that the dual space X' of a normed space X is the space of bounded operators from X to \mathbb{C} .

LEMMA 4.18. *For $\varphi \in X'$ we have that $\ker(\varphi)$ is a closed subspace of X .*

PROOF. Let (x_n) be a sequence in $\ker(\varphi)$ converging to $x \in X$. Then $\varphi(x_n) = 0$ for all n and so $|\varphi(x_n) - \varphi(x)| \leq \|\varphi\| \|x - x_n\|$. Thus we have $\varphi(x) = 0$. \square

THEOREM 4.19 (Riesz representation theorem). *Let X be a Hilbert space. For each $\xi \in X$ define $\varphi_\xi(x) = \langle x, \xi \rangle$. Then $\varphi_\xi \in X'$ is a bounded linear functional on X .*

Furthermore, every $\varphi \in X'$ is of the form φ_ξ for some $\xi \in X$.

The final assertion of the theorem is the subtle part and is due to F. Riesz.

PROOF. The Cauchy-Schwarz inequality gives $|\varphi_\xi(x)| \leq \|x\| \|\xi\|$ and thus $\varphi_\xi \in X'$.

Converse statement: For any $x, z \in X$ and a non-zero $\varphi \in X'$. Then $\varphi(x)z - \varphi(z)x \in \ker(\varphi)$.

Let us pick z in $\ker(\varphi)^\perp$, which we can do by the projection theorem, to get

$$0 = \langle z, \varphi(x)z - \varphi(z)x \rangle = \varphi(x)\|z\|^2 - \varphi(z)\langle x, z \rangle.$$

Hence,

$$\varphi(x) = \frac{\varphi(z)}{\|z\|^2} \langle x, z \rangle.$$

We set $\xi = \frac{\overline{\varphi(z)}}{\|z\|^2} z$. Then we have $\varphi(x) = \langle x, \xi \rangle$.

Since $\xi \rightarrow \varphi_\xi$ preserves sums and differences we have that $\|\varphi\|$ obeys the parallelogram law. Hence the theorem of Jordan-von Neumann implies that X' is a Hilbert space.

Uniqueness: Suppose $\tilde{\xi}$ is another representation of φ of the form $\varphi_{\tilde{x}}$. Then $\langle x, \xi - \tilde{\xi} \rangle = \langle x, \xi \rangle - \langle \tilde{x}, \xi \rangle = 0$ and $x = \tilde{x}$. \square

The theorem yields that any bounded linear functional φ on ℓ^2 is of the form

$$\varphi(x) = \sum_{n=1}^{\infty} x_n \xi_n \quad \text{for a unique } \xi \in \ell^2.$$

A different description of operators is one consequence of Riesz' theorem, because it implies the existence of the adjoint of an operator.

LEMMA 4.20. *Suppose $T \in B(X)$, X a Hilbert space, and $x, x' \in X$.*

- (1) *If $\langle x, y \rangle = \langle x', y \rangle$ for all $y \in X$, then we have $x = x'$.*
- (2) $\|T\| = \sup\{\|Tx\| = \sup\{|\langle Tx, y \rangle| : x, y \in X \text{ with } \|x\|, \|y\| \leq 1\}$.

For motivation of the general result we indicate the main idea for linear operators T on \mathbb{C}^2 . We represent T with respect to the standard basis of \mathbb{C}^2 , so $T = Ax$ for a matrix $A = (a_{ij})$. We look for a matrix $B = (b_{ij})$ such that

$$\langle Ax, y \rangle = \langle x, By \rangle$$

for all $x, y \in \mathbb{C}^2$. Concretely, we have

$$\left\langle \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}, y \right\rangle = \left\langle x, \begin{pmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} \right\rangle$$

and so

$$\left\langle \begin{pmatrix} a_{11}x_1 + a_{12}x_2 \\ a_{21}x_1 + a_{22}x_2 \end{pmatrix}, y \right\rangle = \left\langle x, \begin{pmatrix} b_{11}y_1 + b_{12}y_2 \\ b_{21}y_1 + b_{22}y_2 \end{pmatrix} \right\rangle$$

The equation is equivalent to

$$\begin{aligned} a_{11}x_1\overline{y_1} + a_{12}x_2\overline{y_1} + a_{21}x_1\overline{y_2} + a_{22}x_2\overline{y_2} &= \\ &= x_2\overline{b_{11}y_1} + x_1\overline{b_{12}y_2} + x_2\overline{b_{21}y_1} + x_2\overline{b_{22}y_2} \end{aligned}$$

to hold for all $x_1, x_2, y_1, y_2 \in \mathbb{C}$. Hence we deduce that

$$a_1 1 = \overline{b_{11}}, a_1 2 = \overline{b_{21}}, a_2 1 = \overline{b_{12}}, a_2 2 = \overline{b_{22}}.$$

Thus

$$B = \begin{pmatrix} \overline{a_{11}} & \overline{a_{21}} \\ \overline{a_{12}} & \overline{a_{22}} \end{pmatrix}$$

is the *conjugate-transpose* of A . The adjoint of T , denoted by T^* , is in this way linked to the original transform.

THEOREM 4.21 (Adjoint). *Let T be a bounded operator on a Hilbert space X . Then there exists a unique operator $T^* \in \mathcal{B}(X)$ such that*

$$\langle Tx, y \rangle = \langle x, T^*y \rangle \quad \text{for all } x, y \in X.$$

The operator T^ is called the adjoint of T .*

PROOF. Fix $y \in X$ and let $\varphi : X \rightarrow \mathbb{C}$ be defined by $\varphi(x) = \langle Tx, y \rangle$. Then φ is linear and by Cauchy-Schwarz it is bounded:

$$|\varphi(x)| \leq |\langle Tx, y \rangle| \leq \|Tx\| \|y\| \leq \|T\| \|x\| \|y\|.$$

Hence φ is a bounded linear functional on X and so by the Riesz representation theorem there exists a unique $\xi \in X$ such that $\varphi(x) = \langle x, \xi \rangle$ for all $x \in X$.

The vector ξ depends on the vector $y \in X$. In order to keep track of this fact we set $T^*y := \xi$. Hence we have defined an operator T^* from X to X based on the structure of bounded linear functionals on X . In summary, we have demonstrated the existence of an operator T^* on X such that

$$\langle Tx, y \rangle = \langle x, T^*y \rangle \quad \text{for all } x, y \in X.$$

(1) T^* is linear.

$$\begin{aligned} \langle x, T^*(\lambda y_1 + \mu y_2) \rangle &= \langle Tx, \lambda y_1 + \mu y_2 \rangle \\ &= \overline{\lambda} \langle Tx, y_1 \rangle + \overline{\mu} \langle Tx, y_2 \rangle \\ &= \overline{\lambda} \langle x, T^*y_1 \rangle + \overline{\mu} \langle x, T^*y_2 \rangle \\ &= \langle x, \lambda T^*y_1 + \mu T^*y_2 \rangle. \end{aligned}$$

(2) T^* is bounded. We use the Cauchy-Schwarz inequality:

$$\begin{aligned} \|T^*y\|^2 &= \langle T^*y, T^*y \rangle = \langle TT^*y, y \rangle \\ &\leq \|TT^*y\| \|y\| \\ &\leq \|T\| \|T^*y\| \|y\|. \end{aligned}$$

Hence we have shown

$$\|T^*y\|^2 \leq \|T\| \|T^*y\| \|y\|$$

If $\|T^*y\| > 0$, then we can through and obtain the desired result: $\|T^*y\| \leq \|T\| \|y\|$. Suppose $\|T^*y\| = 0$. Then the desired inequality holds, too. Consequently, we have proved that

$$\|T^*\| \leq \|T\|.$$

- (3) T^* is unique. Suppose there exists another $S \in B(X)$ such that $\langle Tx, y \rangle = \langle x, Sy \rangle$ for all $x, y \in X$. Then we have

$$\langle x, Sy \rangle = \langle x, T^*y \rangle \quad y \in Y$$

and by a well-known fact about innerproducts we deduce that $T^*y = Sy$ for all $y \in Y$. Hence T^* is unique. \square

We collect a few properties of the adjoint.

LEMMA 4.22. *Let S, T be in $B(X)$ and $\lambda, \mu \in \mathbb{C}$.*

- (1) $(\lambda S + \mu T)^* = \bar{\lambda}S^* + \bar{\mu}T^*$;
- (2) $(ST)^* = T^*S^*$.
- (3) *If T is invertible, then T^* is also invertible and $(T^*)^{-1} = (T^{-1})^*$.*

PROOF. The proofs of (i) and (iii) are left as an exercise. Here we show the second assertion:

$$\langle x, (ST)^*y \rangle = \langle STx, y \rangle = \langle Tx, S^*y \rangle = \langle x, T^*S^*y \rangle$$

holds for all $x \in X$ and so we have $(ST)^* = T^*S^*$. \square

We continue with some useful facts about T^* .

LEMMA 4.23. *Let T be a bounded operator on a Hilbert space X .*

- (1) $(T^*)^* = T$;
- (2) $\|T^*\| = \|T\|$;
- (3) $\|T^*T\| = \|T\|^2$ (C^* -algebra identity)

PROOF. (1) For $x, y \in X$ we have

$$\begin{aligned} \langle y, (T^*)^*x \rangle &= \langle T^*y, x \rangle \\ &= \overline{\langle x, T^*y \rangle} \\ &= \overline{\langle Tx, y \rangle} \\ &= \langle y, Tx \rangle, \end{aligned}$$

so $(T^*)^*x = Tx$ for all $x \in X$.

- (2) In the proof of the existence of the adjoint we established that $\|T^*\| \leq \|T\|$. Applying this result to T^{**} and using (i) yields $\|T\| \leq \|T^*\|$. Hence we have $\|T^*\| = \|T\|$.
- (3) By (ii) we have $\|T^*\| = \|T\|$ that implies

$$\|T^*T\| \leq \|T^*\| \|T\| = \|T\|^2.$$

For the reverse inequality we use

$$\begin{aligned} \|Tx\|^2 &= \langle Tx, Tx \rangle \\ &= \langle T^*Tx, x \rangle \\ &\leq \|T^*Tx\| \|x\| \\ &\leq \|T^*T\| \|x\|^2 \end{aligned}$$

to deduce $\|T\|^2 \leq \|T^*T\|$. \square

Some examples should help to build up some intuition on adjoint operators.

EXAMPLE 4.1.23 (Operators on ℓ^2). (1) The adjoint of $Lx = (0, x_1, x_2, \dots)$ on ℓ^2 is the right shift operator $Rx = (x_2, x_3, \dots)$.

By definition

$$\langle (0, x_1, x_2, \dots), (y_1, y_2, \dots) \rangle = \langle x, L^*y \rangle$$

for all $x, y \in \ell^2$. We denote L^*y by $z = (z_n)$. Therefore we have

$$x_1\overline{y_2} + x_2\overline{y_3} + \dots = x_1\overline{z_1} + x_2\overline{z_2} + \dots.$$

This equation is true for all x_i if $z_1 = y_2, z_2 = y_3, \dots$. Hence by the uniqueness of the adjoint

$$L^*y = (y_2, y_3, \dots),$$

i.e. $L^* = R$.

(2) The adjoint of the multiplication operator T_a for $a \in \ell^\infty$ is the multiplication operator for the sequence \overline{a} .

$$\langle T_ax, y \rangle = \langle x, T_a^*y \rangle$$

Hence

$$a_1x_1\overline{y_1} + a_2x_2\overline{y_2} + \dots = x_1\overline{\overline{a_1}y_1} + x_2\overline{\overline{a_2}y_2} + \dots,$$

which by the uniqueness of the adjoint gives that $T_{\overline{a}}$ is the adjoint of T_a .

A useful class of operators are acting on spaces of continuous functions $C[a, b]$. In order to determine their adjoints we have to define an innerproduct on $C[a, b]$. We use a continuous analog of the ℓ^2 -innerproduct. For $f, g \in C[a, b]$ we define

$$\langle f, g \rangle = \int_a^b f(t)\overline{g(t)}dt.$$

LEMMA 4.24. *The space $(C[a, b], \langle \cdot, \cdot \rangle)$ is an innerproduct space with associated norm*

$$\|f\|_2 = \left(\int_a^b |f(t)|^2 dt \right)^{1/2},$$

which is not complete.

The proof is one of the homework problems.

Define the space $L^2[a, b]$ to be the completion of $C[a, b]$ with respect to $\|\cdot\|_2$, i.e. we add all the limits of Cauchy sequences in $C[a, b]$ to it. The notation has a deeper reason, because this space is an example of a Lebesgue space. More generally, one could define $L^p[a, b]$ for $p \geq 1$ as the completions of $C[a, b]$ for the norm $\|f\|_p = \left(\int_a^b |f(t)|^p dt \right)^{1/p}$. These spaces are of utmost importance for analysis. Due to the lack of measure theory we are not in the position to exploit these spaces further.

EXAMPLE 4.1.24.

The multiplication operator T_a on $L^2[0, 1]$ defined by $a \in C[0, 1]$ has $T_{\overline{a}}$ as its adjoint.

$$\langle T_af, g \rangle = \int_0^1 a(t)f(t)\overline{g(t)}dt = \overline{\int_0^1 f(t)\overline{a(t)}g(t)dt} = \langle f, T_{\overline{a}}g \rangle.$$

We introduce some classes of operators defined in terms of the adjoint.

DEFINITION 4.1.25. Let T be a bounded operator on a Hilbert space X .

- (1) T is called *normal* if $T^*T = T^*T$.
- (2) T is called *unitary* if $T^*T = T^*T = I$.
- (3) T is called *selfadjoint* if $T = T^*$.

EXAMPLES 4.1.26 (Operators on ℓ^2). (1) The multiplication operator T_a for $a \in \ell^\infty$ is normal, since $T_a^*T_a = T_a^*T_a = T_{|a|^2}$. Hence it is unitary if $|a| = 1$ as in the example $(1, i, -1, -i, \dots) = (-i^k)_{k=0}^\infty$. T_a is selfadjoint if and only if a is real-valued.

- (2) The shift operator is not normal: $L^*L = I$ and $LL^*y = (y_2, y_3, \dots) \neq I$. Hence L is not unitary.

We state a few properties of unitary operators. We denote the set of all unitary operators on X by \mathcal{U}

LEMMA 4.25. For S, T in \mathcal{U} we have that ST and TS are also in \mathcal{U} . The identity operator is a unitary operator. Unitary operators are invertible and $T^{-1} = T^*$.

PROOF. Since $(ST)^*(ST) = T^*S^*ST^*$ we get from $S^*S = I$ and $T^*T = I$ that ST is also unitary. The invertibility follows from the definition of unitary operators. \square

In some problems it is of interest to have control over linear operators that preserve the norm, known as isometries.

DEFINITION 4.1.27. Let X be a normed space. Then $T \in B(X)$ is called an *isometry* if $\|Tx\| = \|x\|$ for all $x \in X$.

We settle the structure of isometries for Hilbert spaces.

PROPOSITION 4.1.28. Let T be a bounded operator on a Hilbert space X .

- (1) T is an isometry of X if and only if $T^*T = I$.
- (2) T is unitary then T is an isometry of X .

PROOF. (1) Suppose that $T^*T = I$. Then

$$\|Tx\|^2 = \langle Tx, Tx \rangle = \langle T^*Tx, x \rangle = \langle Ix, x \rangle = \|x\|^2,$$

so T is an isometry.

Conversely, suppose that T is an isometry. Then

$$\langle T^*Tx, x \rangle = \langle Tx, Tx \rangle = \|Tx\|^2 = \|x\|^2 = \langle Ix, x \rangle.$$

Hence $T^*T = I$.

- (2) Suppose that T is unitary. By (i) T is an isometry. \square

EXAMPLE 4.1.29. The shift operator $Rx = (0, x_1, x_2, \dots)$ is an isometry on ℓ^2 , but it is not a unitary operator.

EXAMPLE 4.1.30. Let U be a linear transformation on a finite-dimensional innerproduct space X . Consider U as a matrix relative to an orthonormal basis on X . Show that the following statements are equivalent.

- (1) U is unitary, i.e. $U^*U = I = UU^*$.
- (2) The columns of U are an orthonormal basis of X .

(3) The rows of U are an orthonormal basis of X .

PROOF. We will show that 1 is equivalent to 2. A similar argument can be applied to show that 1 is equivalent to 3.

Let $\mathcal{B} = (x_1, x_2, \dots, x_n)$ be an orthonormal basis on X . We consider U as the matrix

$$U = (U_1|U_2|\dots|U_n)$$

relative to \mathcal{B} with columns U_1, U_2, \dots, U_n . Observe that

$$\begin{aligned} \langle U_i, U_j \rangle &= \left\langle \sum_{k=1}^n U_{i,k} x_k, \sum_{l=1}^n U_{j,l} x_l \right\rangle \\ &= \sum_{k=1}^n \sum_{l=1}^n U_{i,k} \overline{U_{j,l}} \langle x_k, x_l \rangle \\ &= \sum_{k=1}^n U_{i,k} \overline{U_{j,k}} \\ &= U_i U_j^* \end{aligned}$$

1 \Rightarrow 2

Assume that $UU^* = I$, in other words, $U_i U_j^* = \delta_{i,j}$ for $i, j = 1, \dots, n$. Then we have

$$\langle U_i, U_j \rangle = U_i U_j^* = \delta_{i,j},$$

hence (U_1, U_2, \dots, U_n) is an orthonormal system of vectors in X . To show that it is a basis for X it is enough to note that X has dimension n , and the system consists of n vectors.

2 \Rightarrow 1

Assume that the columns U_1, U_2, \dots, U_n of U are an orthonormal basis of X , i.e.

$$\langle U_i, U_j \rangle = \delta_{i,j},$$

for $i, j = 1, \dots, n$. Then we have

$$U_i U_j^* = \langle U_i, U_j \rangle = \delta_{i,j},$$

hence we have $UU^* = I$. Since the columns of U form an orthonormal basis of X , the matrix U is invertible, and we get

$$\begin{aligned} UU^* &= I \\ \Rightarrow (UU^*)^{-1} &= I \\ \Rightarrow (U^*)^{-1}U^{-1} &= I \\ \Rightarrow U^*(U^*)^{-1}U^{-1}U &= U^*IU \\ \Rightarrow IU^{-1}U &= U^*U \\ \Rightarrow II &= U^*U \\ \Rightarrow I &= U^*U. \end{aligned}$$

□

We close our discussion of the adjoint, a notion of utmost importance.

PROPOSITION 4.1.31. *Let T be a bounded operator on a Hilbert space X .*

- (1) $\ker(T) = (\operatorname{ran}(T^*))^\perp$;
- (2) $\ker(T^*) = (\operatorname{ran}(T))^\perp$.

Equivalent formulation:

$$\overline{\operatorname{ran}(T)} = (\ker(T^*))^\perp, \quad \ker(T) = (\operatorname{ran}(T^*))^\perp$$

and consequently:

$$X = \ker(T) \oplus \overline{\operatorname{ran}(T)}.$$

PROOF. (1) $\ker(T) \subseteq (\operatorname{ran}(T^*))^\perp$: Let $x \in \ker(T)$ and let $z \in \operatorname{ran}(T^*)$,
i.e. there exists a $y \in X$ such that $z = T^*y$. Hence

$$\langle x, z \rangle = \langle x, T^*y \rangle = \langle Tx, y \rangle = 0$$

and we have shown that $z \in (\operatorname{ran}(T^*))^\perp$.

$(\operatorname{ran}(T^*))^\perp \subseteq \ker(T)$: Let $x \in (\operatorname{ran}(T^*))^\perp$. As $T^*Tx \in \operatorname{ran}(T^*)$ we have

$$\langle Tx, Tx \rangle = \langle x, T^*Tx \rangle = 0,$$

hence $Tx = 0$ and so $x \in \ker(T)$.

(2) By part (i) we have

$$\ker(T^*) = (\operatorname{ran}(T^{**}))^\perp = \operatorname{ran}(T) = \{0\}.$$

For the equivalent formulation note, that we have as above $\operatorname{ran}(T) = (\ker(T^*))^\perp$, but since $(\ker(T^*))^\perp$ is closed we also get $\overline{\operatorname{ran}(T)} \subseteq (\ker(T^*))^\perp$. The rest of the argument follows similar lines as before. \square

COROLLARY 4.1.32. *Let T be a bounded operator on a Hilbert space X . Then $\ker(T^*) = \{0\}$ if and only if $\operatorname{ran}(T)$ is dense in X*

PROOF. Assume that $\ker(T^*) = \{0\}$. Then

$$\ker(T^*)^\perp = \{0\}^\perp = X$$

and the assertion (ii) of the proposition implies that

$$\ker(T^*)^\perp = (\operatorname{ran}(T))^{\perp\perp} = \operatorname{ran}(T).$$

Thus we have $\operatorname{ran}(T)$ is dense in X .

Suppose $\operatorname{ran}(T)$ is dense in X . Then by $(\operatorname{ran}(T))^{\perp\perp} = \overline{\operatorname{ran}(T)} = X$ and

$$\ker(T^*) = \operatorname{ran}(T)^\perp = ((\operatorname{ran}(T))^{\perp\perp})^\perp = X^\perp = \{0\}.$$

\square

The corollary allows one to check if the range of an operator is dense in a Hilbert space by determining its adjoint and the computation of the kernel of the adjoint. In general, this is a good strategy, because it is very difficult to compute the range of an operator. Another important application of the preceding theorem is the Fredholm alternative.

THEOREM 4.26 (Fredholm alternative). *Suppose T is a bounded linear operator on a Hilbert space X with closed range. Then the equation*

$$Tx = b, \quad b \in X$$

has a solution x in X for every $b \in X$ if and only if

$$b \in (\ker(T^*))^\perp.$$

Hence operators with a closed range have a general criterion of existence. For example if $T \in \mathcal{B}(X)$ satisfies for all $x \in X$ and estimate of the form

$$\|Tx\| \geq c\|x\| \quad \text{for some } c > 0.$$

EXAMPLE 4.1.33. The range of the right shift operator R on ℓ^2 is closed since it consists of $\{(0, x_2, x_3, \dots) : x_i \in \mathbb{C}\}$. The left shift is L not invertible since its kernel is one-dimensional and spanned by $(1, 0, 0, \dots)$.

The equation

$$Rx = b \Leftrightarrow (0, x_1, x_2, \dots) = (b_1, b_2, \dots)$$

is solvable if and only if $b_1 = 0$, or $b \in (\ker(L))^\perp$.

On the other hand

$$Lx = b$$

is solvable for all $b \in \ell^2$ despite of L not being injective.

4.1.5. Orthonormal bases for Hilbert spaces. Hilbert spaces have one more property distinguishing them from Banach spaces: the existence of orthonormal bases.

DEFINITION 4.1.34. An *orthonormal basis* of a Hilbert space X is a set of vectors $\{e_j\}_{j \in J}$ such that $\text{span}\{e_j\}$ is dense in X and $\langle e_i, e_j \rangle = 0$ for $i \neq j$ and $\|e_i\| = 1$ for $i \in J$.

We know that $\overline{\text{span}\{e_j\}} = X$ if and only if $\langle e_j, x \rangle = 0$ for all $j \in J$ implies that $x = 0$.

In general an orthonormal basis may have uncountably many elements, e.g. the space of almost periodic functions. In the case that $\{e_j\}_{j \in J}$ is a countable set, then the Hilbert space X is separable.

THEOREM 4.27. *Any Hilbert space has an orthonormal basis.*

The proof relies on the axiom of choice and is a well-known application of Zorn's lemma.

From now on we will assume that the orthonormal basis of a Hilbert space is countable. An important example is the exponential basis $\{e^{2\pi i n x} : n \in \mathbb{Z}\}$ of the Hilbert space $L^2[0, 1]$. The theory of Fourier series has been of great influence in the development of the theory of Hilbert spaces.

PROPOSITION 4.1.35. *Let M be a closed subspace of a Hilbert space X such that M has a Hilbert basis $\{e_n\}_{n \in \mathbb{N}}$. Then the following are equivalent:*

- (1) $\sum_{n=1}^{\infty} a_n e_n$ converges in M .
- (2) (a_n) lies in ℓ^2 .

PROOF. Denote the partial sums of (e_n) by $s_N = \sum_{n=1}^N a_n e_n$. We assume $N > M$ without loss of generality. Then

$$\begin{aligned} \|s_N - s_M\|^2 &= \langle s_N - s_M, s_N - s_M \rangle \\ &= \left\langle \sum_{n=M+1}^N a_n e_n, \sum_{m=M+1}^N a_m e_m \right\rangle \\ &= \sum_{n=M+1}^N a_n \overline{a_n} \langle e_n, e_n \rangle \\ &= \sum_{n=M+1}^N |a_n|^2. \end{aligned}$$

Suppose that $(a_n) \in \ell^2$. Then the preceding computation yields that (s_n) is a Cauchy sequence in M . Since M is closed, (s_n) converges to a s in M . Conversely, suppose that (s_n) converges. Then $\|s_N - s_M\|$ converges to zero. Thus $(\sum_{n=1}^N |a_n|^2)$ is a Cauchy sequence in \mathbb{C} and hence must converge as $N \rightarrow \infty$. \square

In the discussion of innerproduct spaces we established the Bessel inequality for finitely many orthonormal vectors. Hence we obtain the result for countable bases.

PROPOSITION 4.1.36 (Bessel's inequality). *Suppose a closed subspace M of a Hilbert space X has a countable orthonormal basis (e_n) . Then we have*

$$\sum_{n=1}^{\infty} |\langle x, e_n \rangle|^2 \leq \|x\|^2.$$

The preceding two propositions yields that the general Fourier series $\sum_n \langle x, e_n \rangle e_n$. Moreover, we are able to use it to express the projection onto M .

THEOREM 4.28. *Suppose a closed subspace M of a Hilbert space X has a countable orthonormal basis (e_n) . Then the projection of x onto M is given by*

$$Px = \sum_{n=1}^{\infty} \langle x, e_n \rangle e_n.$$

PROOF. We have that $\sum_{n=1}^{\infty} \langle x, e_n \rangle e_n$ converges to a vector y in M and from the orthonormal basis property we have

$$\langle e_m, x - y \rangle = \langle e_m, x \rangle - \sum_{n=1}^{\infty} \langle e_n, x \rangle \langle e_m, e_n \rangle = 0$$

for all $m \in \mathbb{N}$. Thus $\langle e_m, x - y \rangle = 0$, i.e. $x - y \in (\text{span}\{e_m\})^\perp = M^\perp$. Consequently, y is the closest point to x . \square

The case M equal to X is of special interest and is known as Parseval's identity.

THEOREM 4.29 (Parseval's identity). *If $\{e_n\}$ is a countable basis for the Hilbert space X , then any $x \in X$ can be decomposed as*

$$x = \sum_{n=1}^{\infty} \langle x, e_n \rangle e_n.$$

If $x = \sum_{n=1}^{\infty} \langle x, e_n \rangle e_n$ and $y = \sum_{n=1}^{\infty} \langle y, e_n \rangle e_n$, then

$$\langle x, y \rangle = \sum_{n=1}^{\infty} \langle x, e_n \rangle \overline{\langle y, e_n \rangle}.$$

In particular,

$$\|x\|^2 = \sum_{n=1}^{\infty} |\langle x, e_n \rangle|^2.$$

PROOF. The statement about the decomposition of x follows from $Px = x$ for all $x \in X$ for $M = X$. The remaining assertions are elementary computations. \square

Two Hilbert spaces X and Y are called *isomorphic* if there exists a unitary operator T from X to Y with $\text{ran}(T) = Y$.

THEOREM 4.30 (Riesz-Fischer theorem). *Any separable Hilbert space X is isomorphic to ℓ^2 . Suppose (e_n) is an orthonormal basis of X . Then the isomorphism $T : X \rightarrow \ell^2$ is given by $x \mapsto (\langle x, e_n \rangle)_{n \in \mathbb{N}}$.*

PROOF. Bessel's inequality yields that the Fourier coefficients $(\langle x, e_n \rangle)$ are in ℓ^2 . T is linear and by Parseval's identity J preserves innerproducts: $\langle x, y \rangle = \langle Tx, Ty \rangle$. T is surjective: It maps $\sum_n a_n e_n$ to (a_n) which lies in ℓ^2 . Hence T is an isometry between X and ℓ^2 . \square

Topology of normed spaces and continuity

5.1. Topology of normed spaces

Normed spaces are a generalizations of \mathbb{R} with the absolute value. Definitions and theorems based on the absolute value of a real number generalize verbatim to general normed spaces.

- DEFINITION 5.1.1. (1) A set $U \subset X$ is a *neighborhood* of $x \in X$ if $B_r(x) \subset U$ for some $r > 0$.
- (2) A set $O \subset X$ is *open* if every $x \in O$ has a neighborhood U contained in O .
- (3) A set $C \subset X$ is *closed* if its complement $C^c = X \setminus C$ is open.

Note that the definition of open sets depends on the norm. In other words, open sets with respect to one norm need not be open with respect to another norm.

LEMMA 5.1. *Let $(X, \|\cdot\|)$ be normed space. Then $B_r(x)$ is open and $\overline{B_r(x)}$ is closed for $x \in X$ and $r > 0$.*

PROOF. The proof goes along the same lines as in the case of the real line. Suppose that $y \in B_r(x)$ and choose ε as $\varepsilon = r - d(x, y) > 0$. The triangle inequality yields that $B_\varepsilon(y) \subset B_r(x)$, i.e. $B_r(x)$ is open. We show that $X \setminus \overline{B_r(x)}$ is open. For $y \in X \setminus \overline{B_r(x)}$ we set $\varepsilon = d(x, y) - r > 0$ and once more by the triangle inequality we deduce that $B_\varepsilon(y) \subset X \setminus \overline{B_r(x)}$. Hence $X \setminus \overline{B_r(x)}$ is open and $\overline{B_r(x)}$ is closed. \square

DEFINITION 5.1.2. For a subset A of $(X, \|\cdot\|)$ we introduce some notions.

- (1) The *closure* of a subset A of X , denoted by \overline{A} , is the intersection of all closed sets containing A .
- (2) The *interior* of a subset A of X , denoted by $\text{int}A$, is the union of all open subsets of X contained in A .
- (3) The *boundary* of a subset A of X , denoted by $\text{bd}A$, is the set $\overline{A} \setminus \text{int}A$.

We continue with some definitions

DEFINITION 5.1.3. Let A be a subset of $(X, \|\cdot\|)$.

- (1) A point $x \in A$ is *isolated* in A if there exists a neighborhood U of x such that $U \cap A = \{x\}$.
- (2) A point $x \in \mathbb{R}$ is said to be an *accumulation point* of A if every neighborhood of x contains points in $A \setminus \{x\}$.

DEFINITION 5.1.4. A subset A of $(X, \|\cdot\|)$ is said to be *dense* in \mathbb{R} if its closure is equal to X , i.e. $\overline{A} = X$. If the dense subset A is countable, then X is called *separable*.

In other words, a subset A of a normed space X is dense in X if for each $x \in X$ and each $\varepsilon > 0$ there exists a vector $y \in A$ such that

$$\|x - y\| < \varepsilon.$$

The relevance of a dense subset of a normed space is that it provides a way to approximate elements of the normed space by ones from the dense subset up to any given precision.

LEMMA 5.2. *Suppose A is a dense subspace of a normed space X . For any $x \in X$ there exists a sequence of elements $x_k \in A$ such that $\|x_k - x\| \rightarrow 0$ as $k \rightarrow \infty$.*

PROOF. For $x \in X$ there exists an x_k such that $\|x_k - x\| < 1/k$ for $k = 1, 2, \dots$. By construction x_k converges to x . \square

The next results have been proved in the section on real numbers and these are also true for normed spaces. The proofs of these results are along the same lines as the ones for the real line.

LEMMA 5.3. *Let $\{O_j : j \in J\}$ be a family of open sets of $(X, \|\cdot\|)$.*

- (1) $\bigcap_{j=1}^n O_j$ is an open set for any $n \in \mathbb{N}$.
- (2) $\bigcup_{j \in J} O_j$ is open for a general index set J .

Note that open and closed subset of a normed space also applies to subspaces, since these are sets with some extra properties. For the most part we are going to discuss closed subspaces of a normed space.

LEMMA 5.4. *Suppose A is a subset of $(X, \|\cdot\|)$.*

- (1) $\overline{A} = (\text{Int}(A^c))^c$ and $\text{int}(A) = \overline{(A^c)^c}$
- (2) $\text{bd}A = \text{bd}(A^c) = \overline{A} \cap \overline{A^c}$
- (3) $\overline{A} = A \cup \text{br}A = \text{int}A \cup \text{bd}A$

LEMMA 5.5. *Suppose A is a subset of $(X, \|\cdot\|)$.*

- (1) $\overline{A} = \{x \in X : \text{every neighborhood of } x \text{ intersects } A\}$
- (2) $\text{int}(A) = \{x \in X : \text{some neighborhood of } x \text{ is contained in } A\}$
- (3) $\text{bd}(A) = \{x \in X : \text{every neighborhood of } x \text{ intersects } A \text{ and its complement}\}$

LEMMA 5.6. *A point x in a normed space $(X, \|\cdot\|)$ is an accumulation point of A if and only if every neighborhood of x contains infinitely many points of A .*

DEFINITION 5.1.5. A set K in a metric space X is called *compact* if every sequence in K contains a subsequence converging to a point in K .

The Bolzano-Weierstrass theorem implies that a bounded and closed subset of R^n is compact.

We collect all notions of continuity required in this course.

DEFINITION 5.1.6 (Different types of continuity). Let $(X, \|\cdot\|)$ and $(Y, \|\cdot\|)$ be two normed spaces, let $A \subset X$ and let $f : A \rightarrow Y$ be a function.

- (1) We say that f is *continuous* at a point $a \in A$ if for all $\varepsilon > 0$ there is $\delta > 0$ such that for all $x \in A$ with $\|x - a\| < \delta$ we have $\|f(x) - f(a)\| < \varepsilon$.
- (2) We say that f is *continuous* on A if it is continuous at each point of A .
- (3) We say that f is *uniformly continuous* on A if for all $\varepsilon > 0$ there is $\delta > 0$ such that for all $x, y \in A$ with $\|x - y\| < \delta$ we have $\|f(x) - f(y)\| < \varepsilon$.

(4) We say that f is Lipschitz (with Lipschitz constant $L \in \mathbb{R}$) if

$$\|f(x) - f(x')\| \leq L \|x - x'\| \quad \text{for all } x, x' \in A.$$

LEMMA 5.7. *If $f: A \rightarrow Y$ is a Lipschitz function, where $A \subset X$ and X, Y are normed spaces, then f is continuous at every point $a \in A$. Moreover, f is uniformly continuous.*

PROOF. Let $a \in A$. We assume that f is Lipschitz with Lipschitz constant $L > 0$ and we show that f is continuous at a .

Let $\varepsilon > 0$. Put $\delta := \frac{\varepsilon}{L}$, so if $\|x - a\| < \delta$, then

$$\|f(x) - f(a)\| \leq L \|x - a\| < L \delta = L \frac{\varepsilon}{L} = \varepsilon,$$

so $\|f(x) - f(a)\| < \varepsilon$.

Since $\varepsilon > 0$ was arbitrary, this proves the continuity of f at a . Since $a \in A$ was arbitrary, this proves the continuity of f everywhere on A . Since the δ is independent of the choice of a we deduce that f is uniformly continuous. \square

Here is a useful criterion for continuity of a function.

PROPOSITION 5.1.7. *Let $f: A \rightarrow Y$ be a function, where $A \subset X$ and X, Y are normed spaces. Let $a \in A$. Then the following two statements are equivalent.*

(i) f is continuous at a .

(ii) For every sequence $(x_n) \subset A$, if $x_n \rightarrow a$ then $f(x_n) \rightarrow f(a)$.

PROOF. i) \Rightarrow (ii): We assume that f is continuous at a .

Let $(x_n) \subset A$ be a sequence such that $x_n \rightarrow a$. We prove that $f(x_n) \rightarrow f(a)$.

Let $\varepsilon > 0$. Since f is continuous at a , there is $\delta > 0$ such that if $\|x - a\| < \delta$ then $\|f(x) - f(a)\| < \varepsilon$.

Since $x_n \rightarrow a$, there is $N \in \mathbb{N}$ such that for all $n \geq N$ we have $\|x_n - a\| < \delta$. From the above, if $n \geq N$ we must then have $\|f(x_n) - f(a)\| < \varepsilon$.

As ε was arbitrary, this proves that $f(x_n) \rightarrow f(a)$.

(i) \Leftarrow (ii): We assume by contradiction that f is *not* continuous at a . Let us write down carefully what that means.

Firstly, we recall the definition of continuity. f is continuous at the point $a \in A$ means:

for all $\varepsilon > 0$ there is $\delta > 0$ such that for all $x \in A$ with $\|x - a\| < \delta$ we have $\|f(x) - f(a)\| < \varepsilon$.

Next, we formulate the *negation* of this statement.

The function f is *not* continuous the point $a \in A$ means:

there is $\varepsilon_0 > 0$ such that for all $\delta > 0$ there is an element of A , which we denote by x_δ , such that $\|x_\delta - a\| < \delta$ but $\|f(x_\delta) - f(a)\| \geq \varepsilon_0$.

For every $n \geq 1$, we may choose $\delta = \frac{1}{n}$. Then for some element of A , which we denote by x_n , we have that $\|x_n - a\| < \frac{1}{n}$ but $\|f(x_n) - f(a)\| \geq \varepsilon_0$.

We have thus obtained a sequence $(x_n) \subset A$ such that $\|x_n - a\| < \frac{1}{n} \rightarrow 0$, so $x_n \rightarrow a$. However, since $\|f(x_n) - f(a)\| \geq \varepsilon_0$, the sequence $f(x_n) \not\rightarrow f(a)$, which is a contradiction.

Hence f must be continuous at a . \square

LEMMA 5.8. *et $I \subset \mathbb{R}$ be an interval and let $f: I \rightarrow \mathbb{R}$ be a differentiable function. Assume that for some $L \in \mathbb{R}$ we have*

$$(5.1) \quad |f'(x)| \leq L \quad \text{for all } x \in I.$$

Then f is Lipschitz with Lipschitz constant L .

PROOF. We use the mean value theorem (also called Rolle's theorem). Since f is differentiable everywhere throughout the interval I , for any two points $a, b \in I$ with $a < b$, there is $c \in (a, b)$ such that

$$f'(c) = \frac{f(a) - f(b)}{a - b}.$$

From here we get, using (5.1), that

$$|f(a) - f(b)| = |f'(c)| |a - b| \leq L |a - b|,$$

which proves that f is Lipschitz with Lipschitz constant L . \square

The norm and the innerproduct are continuous mappings.

LEMMA 5.9. *Let X be a normed space. Then $x \rightarrow \|x\|$ is continuous and moreover Lipschitz continuous with constant 1.*

PROOF. By the triangle inequality we have

$$\|x\| - \|y\| = \|x - y + y\| - \|y\| \leq \|x - y\| + \|y\| - \|y\| = \|x - y\|,$$

and if $\|y\| > \|x\|$ we get

$$\| \|x\| - \|y\| \| \leq \|x - y\|.$$

Hence $\|\cdot\|$ is a Lipschitz continuous and in particular continuous. \square

LEMMA 5.10. *Let X be an innerproduct space. Then the innerproduct is continuous in each component.*

PROOF. We have to show that $x \rightarrow \langle x, y \rangle$ is continuous for a fixed $y \in X$. By the symmetry of innerproducts this also yields the continuity with respect to the second component.

By Cauchy-Schwarz

$$|\langle x - x', y \rangle| \leq \|x - x'\| \|y\|$$

for a fixed y . Hence for $\varepsilon > 0$ we take $\delta \|y\|$ in the definition of continuity or by noticing that we have a bounded map. \square

EXAMPLE 5.1.8. For $a = (a_n) \in \ell^\infty$ we define $\varphi(x) = \sum_n a_n x_n$ for $(x_n) \in \ell^1$. Then φ is continuous, i.e. a bounded linear functional on ℓ^1 . First we show that φ is well-defined.

$$|\varphi(x)| \leq \sum_n |a_n| |x_n| \leq \|a\|_\infty \sum_n |x_n| = \|a\|_\infty \|x\|_1.$$

Furthermore this yields that φ is a bounded linear mapping from ℓ^1 to \mathbb{C} and hence continuous.

Linear mapping between normed spaces are an important class of continuous functions.

PROPOSITION 5.1.9. *Let X and Y be normed spaces. For a linear transformation $T : X \rightarrow Y$ the following conditions are equivalent:*

- (1) T is uniformly continuous.
- (2) T is continuous on X .
- (3) T is continuous at 0.
- (4) T is a bounded operator.

PROOF. We will show the following implications to demonstrate the assertions.

From the definitions we have (i) implies (ii) and (ii) implies (iii).

(iii) \Rightarrow (iv) By the continuity of T at 0 there exists a $\delta > 0$ for $\varepsilon = 1$ such that $\|Tx\| < \varepsilon = 1$ for $\|x\| \leq \delta$. We want to show that there exists a constant $C > 0$ such that

$$\|Tx\| \leq C\|x\| \quad \text{for all } x \text{ with } \|x\| \leq 1$$

Note that for $x \in \overline{B_1(0)}$ we have $\frac{\delta x}{2} \in B_\delta(0)$:

$$\|\frac{\delta x}{2}\| = \delta\|x\|/2 \leq \delta/2 < \delta.$$

Hence $\|T(\frac{\delta x}{2})\| < 1$. Since T is linear transformation this condition is equivalent to $\|T(\frac{\delta x}{2})\| = \delta\|T(x)\|/2 < 1$ and thus $\|Tx\| \leq 2/\delta$ for $x \in B_1(0)$. In other words, T is a bounded operator.

(iv) \Rightarrow (i) Since T is linear we have

$$\|Tx - Ty\| = \|T(x - y)\| \leq C\|x - y\|$$

for all $x, y \in X$. Let $\varepsilon > 0$ and $\delta = \varepsilon/C$. Then for all $x, y \in X$ with $\|x - y\| < \delta$

$$\|Tx - Ty\| = \|T(x - y)\| \leq C\|x - y\| \leq C\varepsilon/C = \varepsilon.$$

Hence T is uniformly continuous. □

We just state the equivalence between continuity and the boundedness of a linear mapping as a separate statement due to its relevance.

PROPOSITION 5.1.10 (Boundedness \Leftrightarrow Continuity). *A linear operator T between two normed spaces X and Y is continuous if and only if it is bounded.*

Linear mappings between finite dimensional vector spaces and matrix decompositions

6.1. Linear mappings between finite dimensional vector spaces

Finite-dimensional vector spaces and linear mappings between them are a useful tool for engineers, scientists and mathematicians, aka Linear Algebra. In this chapter we present some basic results from Linear Algebra.

We restrict our discussion to complex vector spaces, i.e. the scalars in our linear combinations are complex numbers.

DEFINITION 6.1.1. A complex number λ is called an *eigenvalue* of a linear transformation $T : V \rightarrow W$ if there exists a non-zero $x \in X$ such that $Tx = \lambda x$. In other words, $x \in \ker T - \lambda I$. The subspace $E_\lambda = \ker T - \lambda I$ is called the eigenspace of T for the eigenvalue λ . The dimension of E_λ is called the *geometric multiplicity* of λ . The set $\sigma(T)$ of \mathbb{C}

$$\sigma(T) = \{z \in \mathbb{C} : T - zI \text{ is not invertible}\}$$

is known as the spectrum of T .

Note that E_λ consists of the eigenvectors of T and the zero vector 0 . For finite-dimensional vector spaces $\sigma(T)$ is the set of all eigenvalues counting multiplicities of T .

THEOREM 6.1. *Suppose T is a linear transformation on a finite-dimensional complex vector space. Then there exists an eigenvalue $\lambda \in \mathbb{C}$ for an eigenvector x of T .*

PROOF. We assume that $\dim(X) = n$ and choose any non-zero vector x in X . Consider the following set of $n + 1$ vectors in X :

$$\{x, Tx, T^2x, \dots, T^n x\}.$$

Since $n + 1$ vectors in an n -dimensional vector space X are linearly independent, there exists a non-trivial linear combination:

$$a_0x + a_1Tx + \dots + a_nT^n x = 0.$$

Let us denote by $p(t) = a_0 + a_1t + \dots + a_nt^n$ the polynomial that after replacing t by the linear transformation T and powers of T by the corresponding iterates of T . Then the non-trivial linear combination among the vectors turns into a polynomial equation in T :

$$p(T) = 0.$$

By the Fundamental Theorem of Algebra any polynomial can be written as a product of linear factors:

$$p(t) = c(t - \lambda_1)(t - \lambda_2) \cdots (t - \lambda_n), \quad \lambda_i \in \mathbb{C}, c \neq 0.$$

Hence $p(T)$ has a factorization of the form:

$$p(T) = c(T - \lambda_1 I)(T - \lambda_2 I) \cdots (T - \lambda_n I).$$

Hence $p(T)$ is a product of linear mappings $T - \lambda_j I$ for $j = 1, \dots, n$. We know that $p(T)x = 0$ for a non-zero $x \neq 0$, which implies that at least one of these linear mappings is not invertible. Thus it has to have a non-trivial kernel, let's say $y \in \ker(T - \lambda_i I)$, which yields that y is an eigenvector for the eigenvalue λ_i . Consequently, we have shown the desired assertion. \square

PROPOSITION 6.1.2 (Gersgorin's disks theorem). *For any $n \times n$ matrix the spectrum is contained in the following union of disks*

$$\cup_{i=1}^n \{z \in \mathbb{C} : |z - a_{ii}| \leq \sum_{j=1, j \neq i}^n |a_{ij}|\}.$$

The disks $B_i = \{z \in \mathbb{C} : |z - a_{ii}| \leq \sum_{j=1, j \neq i}^n |a_{ij}|\}$ centered at a_{ii} and if radius $r_i \sum_{j=1, j \neq i}^n |a_{ij}|$ are called Gersgorin disks.

PROOF. Let λ be an eigenvalue of A and eigenvector x . In components the eigenvalue equation $Ax = \lambda x$ is the set of equations:

$$\sum_{j=1}^n a_{ij}x_j = \lambda x_i \quad \text{for } i = 1, \dots, n.$$

Hence

$$(\lambda - a_{ii})x_i = \sum_{j=1, j \neq i}^n a_{ij}x_j$$

and by the triangle inequality

$$|\lambda - a_{ii}||x_i| \leq \sum_{j=1, j \neq i}^n |a_{ij}||x_j| \leq \sum_{j=1, j \neq i}^n |a_{ij}|\|x\|_\infty.$$

Choose $i \in \{1, \dots, n\}$ to be the largest component of x , i.e. $|x_i| = \|x\|_\infty$ we obtain the conclusion after dividing through by $\|x\|_\infty$. \square

PROPOSITION 6.1.3. *Eigenvalues of a matrix A corresponding to distinct eigenvalues are linearly independent.*

PROOF. Suppose $\lambda_i \neq \lambda_k$ for $i \neq k$ and $Ax_i = \lambda_i x_i$ for $x_i \neq 0$. We assume that $\{x_1, \dots, x_n\}$ is linearly dependent. Hence there exists a linear dependence relation with the fewest number of elements, say m . Thus there exist a_1, \dots, a_m such that

$$\sum_{j=1}^m a_j x_j = 0.$$

Application of A to this linear dependence relation yields

$$\sum_{j=1}^m a_j Ax_j = \sum_{j=1}^m a_j \lambda_j x_j = 0.$$

Multiplication of the last equation by λ_m and subtracting from the linear dependence relation gives

$$\sum_{j=1}^m (a_j \lambda_j - a_j \lambda_m) x_j = 0.$$

Hence the coefficient for x_m is zero. Therefore we have found a linear combination with $m-1$ vectors, contrary to our assumption of m being the smallest such linear combination. \square

DEFINITION 6.1.4. A $n \times n$ matrix A is called diagonalizable if it has n linearly independent eigenvectors.

Note that the set of eigenvectors of a diagonalizable matrix is consequently a basis for \mathbb{C}^n .

By definition a diagonalizable $n \times n$ matrix A has eigenvalues $\lambda_1, \dots, \lambda_n$ and associated eigenvectors u_1, \dots, u_n satisfying:

$$\begin{aligned} Au_1 &= \lambda u_1 \\ &\vdots \\ Au_n &= \lambda u_n. \end{aligned}$$

Collect the eigenvectors of A into one matrix: $U = (u_1|u_2|\dots|u_n)$; and the eigenvalues of A into the diagonal matrix

$$D = \begin{pmatrix} \lambda_1 & 0 & \dots & \dots & 0 \\ \vdots & \lambda_2 & & 0 & \dots & 0 \\ \vdots & 0 & \ddots & \ddots & & \lambda_n \end{pmatrix}.$$

Then the eigenvalue equations turn into a matrix equation:

$$AU = UD.$$

Since A is diagonalizable, the eigenvectors are a basis for \mathbb{C}^n . Hence U is invertible and we have

$$A = UDU^{-1}.$$

Sometimes U is an unitary matrix, i.e. the eigenvectors yield an orthonormal basis for \mathbb{C}^n . Then we have $A = UDU^*$.

A well-known criterion for the non-invertibility of a matrix is the vanishing of its determinant. Hence eigenvalues are the zeros of the polynomial $p_A(z) = \det(zI - A)$, known as the *characteristic polynomial*.

LEMMA 6.2. *Similar matrices have the same characteristic equation.*

PROOF. Let A and B be similar matrices. Thus there exists an invertible matrix S such that $B = S^{-1}AS$.

$$p_B(z) = \det(zI - S^{-1}AS) = \det(zS^{-1}S - S^{-1}AS) = \det(S^{-1}(zI - A)S) = p_A(z).$$

\square

As an important consequence of the existence of an eigenvector for linear mappings between complex finite-dimensional vector spaces we prove Schur's triangularization theorem, our first classification theorem. Before we introduce a refined version of similarity. Namely, if the matrix S in the definition of similar matrices may be chosen as a unitary matrix, then we call the matrices A and B *unitarily equivalent*.

THEOREM 6.3 (Triangularization Theorem). *Given a $n \times n$ matrix with eigenvalues $\lambda_1, \dots, \lambda_n$, counting multiplicities. There exists a unitary $n \times n$ matrix U such that*

$$A = UTU^*$$

for an upper triangular matrix T with the eigenvalues on the diagonal. Hence any matrix is similar to an upper triangular matrix.

We refer to the decomposition of the theorem as *Schur form*.

PROOF. We proceed by induction on n . For $n = 1$, there is nothing to show. Suppose that the result is true up to matrices of size $n - 1$.

Let A be a $n \times n$ matrix with eigenvalues $\lambda_1, \dots, \lambda_n$ counting multiplicities. Choose a normalized eigenvector u_1 for the eigenvalue λ_1 . Then we extend u_1 to a basis $\{u_1, \dots, u_n\}$ of \mathbb{C}^n and we choose this basis to be orthonormal. Relative to this basis the matrix is of the form

$$A = U \begin{pmatrix} \lambda_1 & x & \cdots & x \\ 0 & & & \\ \vdots & A_{n-1} & & \\ 0 & & & \end{pmatrix} U^{-1},$$

where U is the matrix of the system $\{u_1, \dots, u_n\}$ relative to the canonical basis. Since this is a unitary matrix, the similarity, is actually a unitary equivalence. By the induction hypothesis there exists a $(n - 1) \times (n - 1)$ -matrix V such that VAV^* is upper triangular. Set \tilde{V} to be the $n \times n$ matrix where $v_1 1 = 1$ and the other entries of the first column and row are zero. Then \tilde{V} is a unitary matrix and $U\tilde{V}$ is the desired unitary matrix. \square

EXAMPLE 6.1.5. Find the Schur form of $A = \begin{pmatrix} 5 & 7 \\ -2 & -4 \end{pmatrix}$.

First step: Find an eigenvalue of A and associated eigenvector. The characteristic polynomial is $\lambda^2 - \lambda - 6 = 0$ and so $\lambda_1 = -2$ and $\lambda_2 = 3$. An eigenvector for $\lambda_1 = -2$ is $x_1 = \begin{pmatrix} 1 \\ -1 \end{pmatrix}$.

The second step is to complete it to a basis of \mathbb{C}^2 . In our case we take the eigenvector to the second eigenvalue and note that the corresponding set of vectors is linearly independent: $x_2 = \begin{pmatrix} 7 \\ -2 \end{pmatrix}$.

Third step: Use a orthonormalization procedure, e.g. Gram-Schmidt, to turn the system $\{x_1, x_2\}$ into a basis $\{u_1 = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ -1 \end{pmatrix}, u_2 = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ 1 \end{pmatrix}\}$.

Final step: Form the matrix $U = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ -1 & 1 \end{pmatrix}$. Computation of $U^*AU = \begin{pmatrix} 2 & 9 \\ 0 & 3 \end{pmatrix}$, which has the eigenvalues of A on its diagonal and is upper triangular.

Schur's triangularization theorem has a number of important consequences.

THEOREM 6.4 (Cayley-Hamilton). *Given a $n \times n$ matrix. Then*

$$p_A(A) = 0,$$

where $p_A(A)$ is the characteristic polynomial of A .

We state a refined version of Schur's triangularization theorem

THEOREM 6.5 (Schur normal form). *Given a $n \times n$ matrix A with distinct eigenvalues $\lambda_1, \dots, \lambda_k$ with $k \leq n$. Then A is unitarily equivalent to*

$$\begin{pmatrix} T_1 & 0 & \cdots & 0 \\ 0 & T_2 & \ddots & 0 \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \cdots & 0 & T_k \end{pmatrix}$$

where T_i has the form

$$\begin{pmatrix} \lambda_i & x & \cdots & x \\ 0 & \lambda_i & \ddots & x \\ \vdots & \ddots & \ddots & x \\ 0 & \cdots & 0 & \lambda_i \end{pmatrix}$$

We present an interplay on the structure of diagonalizable matrices and the notions from our discussion of normed spaces. Let $\mathcal{M}_n(\mathbb{C})$ denote the vector space of complex $n \times n$ matrices, and by \mathcal{D} the set of diagonalizable $n \times n$ matrices.

LEMMA 6.6. *$\mathcal{M}_n(\mathbb{C})$ is a normed vector space with respect to the Frobenius norm*

$$\|A\|_F = \text{tr}(A^*A)^{1/2}$$

and this norm comes from an innerproduct on $\mathcal{M}_n(\mathbb{C})$:

$$\langle A, B \rangle = \text{tr}(B^*A).$$

Furthermore $\|A\|_F$ is unitarily equivalent $\|UAV\|_F = \|A\|_F$ for unitary matrices U, V .

We leave the proof as an exercise. Use the identification between $\mathcal{M}_n(\mathbb{C})$ and \mathbb{C}^{n^2} and note that then the Frobenius norm is the Euclidean norm on the latter space. A computation yields the following useful fact:

LEMMA 6.7. *Let U be a unitary $n \times n$ matrix. Then $\text{tr}(A) = \text{tr}(UA)$. Furthermore, we have $\text{tr}(AB) = \text{tr}(BA)$ for any $n \times n$ matrices A and B .*

Note that

$$\text{tr}(A^*A) = \sum_{i,j=1}^n |a_{ij}|^2.$$

LEMMA 6.8. *If A and B are unitarily equivalent, then*

$$\sum_{i,j=1}^n |a_{ij}|^2 = \sum_{i,j=1}^n |b_{ij}|^2.$$

PROOF. From $\sum_{i,j=1}^n |a_{ij}|^2 = \text{tr}(A^*A)$ we want to show that this equals $\text{tr}(B^*B)$:

$$\text{tr}(B^*B) = \text{tr}((UAU^*)^*UAU^*) = \text{tr}(UA^*AU^*) = \text{tr}(U^*UA^*A) = \text{tr}(A^*A).$$

□

PROPOSITION 6.1.6. *The set of diagonalizable matrices \mathcal{D} is dense in $\mathcal{M}_n(\mathbb{C})$ with respect to the Frobenius norm. More explicitly, given $A \in \mathcal{M}_n(\mathbb{C})$ and $\varepsilon > 0$. There exists a diagonalizable matrix $\tilde{A} \in \mathcal{M}_n(\mathbb{C})$ such that*

$$\sum_{i,j=1}^n |a_{ij} - \tilde{a}_{ij}|^2 < \varepsilon.$$

We have the Schur form for A

PROOF.

$$A = U \begin{pmatrix} \lambda_1 & x & \cdots & x \\ 0 & \lambda_2 & \ddots & x \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \cdots & \cdots & \lambda_n \end{pmatrix} U^*,$$

for a unitary matrix and eigenvalues $\lambda_1, \dots, \lambda_n$ counting multiplicities. Define small perturbations of these eigenvalues λ_j such that these new numbers $\tilde{\lambda}_1, \dots, \tilde{\lambda}_n$ are all distinct. We add multiples of a number η to the λ_j 's:

$$\tilde{\lambda}_j = \lambda_j + j\eta, \quad \eta > 0$$

and fixed at the end of the proof. Set \tilde{A}

$$U \begin{pmatrix} \tilde{\lambda}_1 & x & \cdots & x \\ 0 & \tilde{\lambda}_2 & \ddots & x \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \cdots & \cdots & \tilde{\lambda}_n \end{pmatrix} U^*,$$

where we only change the diagonal entries of the upper triangular matrix. Now \tilde{A} is diagonalizable and we have

$$\text{tr}((A - \tilde{A})^*(A - \tilde{A})) = \sum_{i,j=1}^n |a_{ij} - \tilde{a}_{ij}|^2$$

Since the diagonal matrix with entries $\lambda_1 - \tilde{\lambda}_1, \dots, \lambda_n - \tilde{\lambda}_n$ is unitarily equivalent to $A - \tilde{A}$ we deduce that

$$\text{tr}((A - \tilde{A})^*(A - \tilde{A})) = \sum_{j=1}^n |\lambda_j - \tilde{\lambda}_j|^2.$$

By the definition of $\tilde{\lambda}_j$ this gives

$$\sum_{j=1}^n |\lambda_j - \tilde{\lambda}_j|^2 = \eta^2 \sum_{j=1}^n j^2 = \eta^2 n(n+1)/2.$$

Consequently,

$$\sum_{j=1}^n |\lambda_j - \tilde{\lambda}_j|^2 \leq \varepsilon$$

for $\eta \leq 2\varepsilon/(n(n+1))$. □

THEOREM 6.9. *Given a $n \times n$ matrix A . Let p_A be the characteristic polynomial of A . Then A annihilates p_A , in other words $p_A(A) = 0$.*

PROOF. Schur's triangularization theorem gives that A is unitarily equivalent to an upper triangular matrix T , $A = UTU^*$ for a unitary matrix U . The powers of A are also similar to powers of T via the same matrix U :

$$A^j = UT^jU^*,$$

e.g. $A^2 = UTU^*UTU^* = UT^2U^*$ since $U^*U = I$. Hence the characteristic polynomials of A and T are also unitarily equivalent:

$$p_A(A) = Up_T(T)U^*.$$

Consequently, $p_A(A) = 0$ if and only if $p_T(T) = 0$. The case $p_T(T) = 0$ is definitely more accessible than the general one, and one can show by a matrix decomposition argument that the latter is true. \square

EXAMPLE 6.1.7. We check the statement for a general 2×2 upper triangular matrix

$$T = \begin{pmatrix} a & b \\ 0 & c \end{pmatrix}.$$

We have to compute T^2

$$T^2 = \begin{pmatrix} a^2 & ab + bc \\ 0 & c^2 \end{pmatrix}.$$

The characteristic polynomial of T is $p_T(z) = z^2 - (a + c)z + ac$. For $z^i \mapsto T^i$ we get

$$p_T(T) = T^2 - (a + c)T + acI^0 = T^2 - (a + c)T + acI,$$

which is equal to

$$\begin{pmatrix} a^2 & ab + bc \\ 0 & c^2 \end{pmatrix} - (a + c) \begin{pmatrix} a & b \\ 0 & c \end{pmatrix} + \begin{pmatrix} ac & 0 \\ 0 & ac \end{pmatrix} = 0.$$

THEOREM 6.10 (Spectral theorem). *Given $A \in \mathcal{M}_n(\mathbb{C})$. Then the following statements are equivalent:*

- (1) A is normal.
- (2) A is unitarily diagonalizable. Hence there exists a unitary matrix U such that $A = UDU^*$, where D is a diagonal matrix with the eigenvalues of A as entries of the diagonal, the columns of U are the corresponding eigenvectors of A .
- (3) $\sum_{i,j=1}^n |a_{ij}|^2 = \sum_{i,j=1}^n |\lambda_i|^2$, where $\lambda_1, \dots, \lambda_n$ are the eigenvalues of A counting multiplicities.

In the proof we make use of two useful statements. An elementary computation yields the following fact.

LEMMA 6.11. *Suppose A and B are unitarily equivalent. Then A is normal if and only if B is normal, i.e. A is normal if and only if UAU^* is normal for some unitary matrix U .*

LEMMA 6.12. *An upper triangular matrix is normal if and only if it is diagonal.*

PROOF. (\Rightarrow) Suppose T is an upper triangular matrix. Then the n, n -th entry of TT^* is $|t_{nn}|^2$ while the n, n -th entry of T^*T is $|t_{nn}|^2 + \sum_{i=1}^{n-1} |t_{in}|^2$. If T is normal, then these two entries have to be the same. Hence $t_{in} = 0$ for $i = 1, \dots, n - 1$. Repeating this argument for the entries $n - 1, n - 1, \dots, 1$ gives that T is diagonal. (\Leftarrow) If T is diagonal, then T is certainly normal. \square

SPECTRAL THEOREM. (i) \Leftrightarrow (ii) By Schur's theorem A is unitarily equivalent to an upper triangular matrix T . Then we know that A is normal if and only if T is normal, which is normal if and only if T is diagonal. In other words, A is unitarily equivalent to a diagonal matrix.

(ii) \Leftrightarrow (iii) Suppose A is unitarily equivalent to a diagonal matrix D where the diagonal entries of D are the eigenvalues $\lambda_1, \dots, \lambda_n$ of A . Then

$$\sum_{i,j=1}^n |a_{ij}|^2 = \operatorname{tr}(A^*A) = \operatorname{tr}(D^*D) = \sum_{i=1}^n |\lambda_i|^2.$$

(ii) \Leftrightarrow (ii) By Schur's theorem A is unitarily equivalent to a triangular matrix T :

$$\sum_{i=1}^n |\lambda_i|^2 = \sum_{i,j=1}^n |a_{ij}|^2 = \operatorname{tr}(A^*A) = \operatorname{tr}(T^*T) = \sum_{i=1}^n |t_{ii}|^2 + \sum_{i,j=1, i \neq j}^n |t_{ij}|^2.$$

Since the diagonal entries of T are the eigenvalues of A we have that

$$\sum_{i=1}^n |\lambda_i|^2 = \sum_{i=1}^n |t_{ii}|^2.$$

Hence $t_{ij} = 0$ for $i \neq j$, i.e. T is diagonal and A is unitarily equivalent to a diagonal matrix. \square

The matrix $\begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}$ is not normal. This matrix and its higher-dimensional analogs are going to play a crucial role in the Jordan Normal Form.

LEMMA 6.13. Let X be a finite-dimensional Hilbert space and $T : X \rightarrow X$ a linear mapping.

- Show that T is a normal if and only if $\|Tx\| = \|T^*x\|$ for all $x \in X$.
- Suppose that T is normal and let x be an eigenvector of T with eigenvalue λ . Show that x is also an eigenvector of T^* with eigenvalue $\bar{\lambda}$.

PROOF. a) We will first prove the *only if* part. Suppose T is normal. We have

$$\|Tx\|^2 = \langle Tx, Tx \rangle = \langle x, T^*Tx \rangle = \langle x, TT^*x \rangle = \langle T^*x, T^*x \rangle = \|T^*x\|^2.$$

We will now prove the *if* part. Suppose $\|Tx\| = \|T^*x\|$ for all $x \in X$. We have

$$\begin{aligned} \|Tx\| &= \|T^*x\|, \quad \forall x \in X \\ \Rightarrow \langle Tx, Tx \rangle &= \langle T^*x, T^*x \rangle, \quad \forall x \in X \\ \Rightarrow \langle T^*Tx, x \rangle &= \langle TT^*x, x \rangle, \quad \forall x \in X \\ \Rightarrow \langle T^*Tx, x \rangle - \langle TT^*x, x \rangle &= 0, \quad \forall x \in X \\ \Rightarrow \langle T^*Tx - TT^*x, x \rangle &= 0, \quad \forall x \in X \\ \Rightarrow \langle (T^*T - TT^*)x, x \rangle &= 0, \quad \forall x \in X \end{aligned}$$

Now, let $A = T^*T - TT^*$, so we have $\langle Ax, x \rangle = 0$ for every $x \in X$. We will now show that every eigenvalue of A is zero. Let λ be an eigenvalue of A , i.e. $Ax = \lambda x$ for some non-zero $x \in X$. We get

$$0 = \langle Ax, x \rangle = \langle \lambda x, x \rangle = \lambda \langle x, x \rangle,$$

and since $x \neq 0$ we have $\langle x, x \rangle \neq 0$ and thus $\lambda = 0$. Also note that A is normal, since it is hermitian. It follows from the Spectral Theorem that we can write $A = UDU^*$, where D is a diagonal matrix with the eigenvalues of A as entries on the diagonal. Since all eigenvalues of A are zero, U is the zero matrix, and thus A is the zero matrix. It follows that $T^*T = TT^*$, which was what we needed to prove.

b) We use the result in a) and get

$$\begin{aligned} (T - \lambda I)x &= 0 \\ \Rightarrow \|(T - \lambda I)x\| &= 0 \\ \Rightarrow \|(T - \lambda I)^*x\| &= 0 \quad (\text{The matrix } T - \lambda I \text{ is normal since it is a sum of normal matrices}) \\ \Rightarrow \|(T^* - \bar{\lambda}I)x\| &= 0 \\ \Rightarrow (T^* - \bar{\lambda}I)x &= 0, \end{aligned}$$

hence x is also an eigenvector of T^* with eigenvalue $\bar{\lambda}$. □

Recall that selfadjoint matrices, $A = A^*$, are normal. Consequently our spectral theorem for normal matrices implies the spectral theorem for selfadjoint matrices.

THEOREM 6.14. *Suppose A is a selfadjoint $n \times n$ matrix. Then A is unitarily equivalent to a diagonal matrix, and the eigenvalues of A are real.*

PROOF. The fact about the diagonalizability follows from the Spectral Theorem for unitary matrices. Now let U be the unitary matrix implementing this similarity: $A = UDU^*$. Then we have $A^* = U\bar{D}U^*$. Hence A is selfadjoint if and only if the diagonal entries of D are real. Since these entries are the eigenvalues of A , we have proved that eigenvalues of a selfadjoint matrix are real numbers. □

In the case of unitary matrices we can also use the spectral theorem to deduce some information about the eigenvalues.

PROPOSITION 6.1.8. *A matrix A is unitary if and only if all of the eigenvalues of A have modulus one.*

LEMMA 6.15. *Let A be a $n \times n$ matrix. Then A^*A and AA^* are selfadjoint matrices.*

DEFINITION 6.1.9. A complex selfadjoint matrix A on an n -dimensional innerproduct space $(X, \langle \cdot, \cdot \rangle)$ is said to be *positive definite* if $\langle Ax, x \rangle > 0$ for all non-zero vectors $x \in X$. If A satisfies the weaker condition $\langle Ax, x \rangle \geq 0$ for all non-zero vectors $x \in X$, then we call A *semi-positive definite*.

The notion of positivity is also of interest in the infinite dimensional setting, where it lies at the heart of the theory of operator algebras. We restrict our discussion to mappings between finite-dimensional vector spaces.

REMARK 6.1.10. The matrix $\begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}$ is not positive definite. Hence one cannot deduce from the positivity of the matrix entries its positive definiteness.

For 2×2 matrices there is a way to state some explicit conditions on the matrix entries by just examining the quadratic form $\langle Ax, x \rangle$. Completion the squares yields that A is positive definite if and only if its pivots are positive. A good way to think about positive definite matrices is to understand its relation with the spectrum.

LEMMA 6.16. *A complex selfadjoint $n \times n$ matrix A is positive if and only if all its eigenvalues $\lambda_1, \dots, \lambda_n$ are positive.*

PROOF. (\Leftarrow) Suppose A is positive definite. Then $\langle Ax, x \rangle$ is positive for all non-zero vectors. In particular, also for eigenvectors. Let x be an eigenvector of A . Then $\langle Ax, x \rangle = \langle \lambda x, x \rangle = \lambda \|x\|^2 > 0$ and thus $\lambda > 0$.

(\Rightarrow) By the spectral theorem A is unitarily equivalent to a diagonal matrix given by its eigenvalues. Hence $\langle Ax, x \rangle$ is positive for all non-zero vectors. \square

For the singular value decomposition we have to know that a $m \times n$ matrix A of rank r give rise to positive semi-definite matrices AA^* and A^*A .

LEMMA 6.17. *Let A be a $m \times n$ matrix A of rank r . Then AA^* is a positive semi-definite $m \times m$ matrix with r positive eigenvalues and $m - r$ zero eigenvalues. Furthermore, A^*A is a positive semi-definite $n \times n$ matrix with r positive eigenvalues and $n - r$ zero eigenvalues.*

PROOF. Note that $(AA^*)^* = AA^*$, i.e. AA^* is selfadjoint. By the spectral theorem AA^* is unitarily equivalent to a diagonal matrix D with the eigenvalues as its entries. Hence we have $UAA^*U^* = D$, but $\langle UAA^*U^*x, x \rangle = \|UAx\|^2 \geq 0$. Hence $\lambda_1, \dots, \lambda_n \geq 0$. By assumption A has rank r , so has AA^* and thus AA^* has r non-zero eigenvectors. The argument for A^*A goes along similar lines. \square

A complex number z may be written as $z = \operatorname{Re}(z) + i\operatorname{Im}(z)$. From the perspective of matrix theory a number is a 1×1 matrix and so one might wonder if there is an analogous way to decompose general matrices. Indeed that is the case: The decomposition in real and imaginary part is based on replacing complex conjugation of a number by its multi-variate analogue, complex conjugation+transpose. Any $m \times n$ matrix A has a *Cartesian decomposition*

$$A = \operatorname{Re}(T) + i \operatorname{Im}(T),$$

where $\operatorname{Re}(T)$ denotes the *real part* of A and $\operatorname{Im}(T)$ denotes the *imaginary part* of A

$$\operatorname{Re}(T) = \frac{A + A^*}{2}, \quad \operatorname{Im}(T) = \frac{A - A^*}{2},$$

and $\operatorname{Re}(T)$, $\operatorname{Im}(T)$ are selfadjoint matrices. (Note that this holds true for arbitrary bounded operators.)

6.1.1. QR Decomposition. The Gram-Schmidt orthonormalization procedure may be expressed in terms of a matrix decomposition, the *QR-decomposition*.

Given n linearly independent vectors u_1, \dots, u_n in \mathbb{C}^n . We wish to construct vectors v_1, \dots, v_n such that $\{v_1, \dots, v_k\}$ is an orthonormal basis for $V_k := \operatorname{span}\{u_1, \dots, u_k\}$ for all $k = 1, \dots, n$.

Suppose v_1, \dots, v_{k-1} have been constructed and are an orthonormal set. Try to find a vector \tilde{v}_k in V_k

$$\tilde{v}_k = u_k + \sum_{i=1}^{k-1} c_{k,i} v_i$$

such that \tilde{v}_k is orthogonal to V_{k-1} . Hence we have the conditions

$$0 = \langle \tilde{v}_k, v_i \rangle = \langle u_k, v_i \rangle + c_{k,i}$$

for $i = 1, \dots, k-1$, which yields $c_{k,i} = -\langle u_k, v_i \rangle$. Finally, we normalize it to obtain $v_k = \tilde{v}_k / \|\tilde{v}_k\|$.

PROPOSITION 6.1.11 (QR decomposition). *Given an invertible $n \times n$ matrix. Then there exists a $n \times n$ matrix Q with orthonormal columns and an upper triangular $n \times n$ matrix R such that $A = QR$. The decomposition is unique.*

PROOF. Let u_1, \dots, u_n be the columns of the invertible matrix A . Then the set $\{u_1, \dots, u_n\}$ is linearly independent. The Gram-Schmidt procedure yields orthonormal vectors $\{v_1, \dots, v_n\}$ such that for each $j = 1, \dots, n$ we have

$$u_j = \sum_{k=1}^j r_{k,j} v_k = \sum_{k=1}^n r_{k,j} v_k$$

with $r_{k,j} = 0$ for $k > j$. Hence $R = (r_{i,j})$ is an $n \times n$ upper triangular matrix. In other words, this n equations reduce in matrix form $A = QR$, where Q has as columns $\{v_1, \dots, v_n\}$.

Uniqueness: Suppose that $A = Q_1 R_1 = Q_2 R_2$ where Q_1, Q_2 are unitary and R_1, R_2 are upper triangular with positive diagonal entries. Then

$$M := R_1 R_2^{-1} = Q_1^* Q_2.$$

Since M is unitary matrix which is also upper triangular, it must be diagonal with positive diagonal entries, and furthermore of modulus one. Hence $M = I$ and thus we have established the uniqueness of the QR decomposition. \square

$$\begin{pmatrix} 6 & 6 & 1 \\ 3 & 6 & 1 \\ 2 & 1 & 1 \end{pmatrix} = \begin{pmatrix} 6 & -2 & -3 \\ 3 & 6 & 2 \\ 2 & -3 & 6 \end{pmatrix} \begin{pmatrix} 7 & 8 & 11/7 \\ 0 & 3 & 1/7 \\ 0 & 0 & 5/7 \end{pmatrix}$$

6.1.2. Singular Value Decomposition. We present a way to factorize an arbitrary complex matrix, namely the *singular value decomposition* (SVD). The SVD is a standard tool in computational and numerical linear algebra.

DEFINITION 6.1.12. Given an $m \times n$ matrix A of rank r . Let $\sigma_1^2 \geq \dots \geq \sigma_r^2$ be the positive eigenvalues of A^*A . The numbers $\sigma_1, \dots, \sigma_r$ the *singular values* of A .

Since the matrix A^*A is of size $n \times n$, it has n eigenvalues and so we define the singular values to the $n - r$ zero eigenvalues to be 0, i.e. $\sigma_j := 0$ for $j = r + 1, \dots, n$.

THEOREM 6.18 (SVD). *Given an $m \times n$ matrix A of rank r . Let $\sigma_1 \geq \dots \geq \sigma_r$ be the positive singular values of A . Let Σ be the $m \times n$ diagonal matrix with $\sigma_1, \dots, \sigma_r$ in the first r diagonal entries and zeros elsewhere. Then there exist unitary matrices U and V , of sizes $m \times m$ and $n \times n$, respectively, such that*

$$A = U \Sigma V^*.$$

The decomposition in the theorem is often called the *full SVD*.

PROOF. Note that $D = \Sigma^* \Sigma$ is a real $n \times n$ diagonal matrix with $\sigma_1^2 \geq \sigma_2^2 \geq \dots \geq \sigma_r^2$ and zeros everywhere. The matrix A^*A is a selfadjoint matrix with r positive eigenvalues $\sigma_1^2 \geq \sigma_2^2 \geq \dots \geq \sigma_r^2$ and $n - r$ eigenvalues equal to zero. The spectral theorem yields that there exists a unitary matrix V such that

$$V^* A^* A V = D.$$

The ij th entry of $V^* A^* A V$ is the innerproduct of columns j and i of AV . Hence the preceding equation yields that the columns of AV are pairwise orthogonal. Furthermore, when $1 \leq i \leq r$ then the length of column j is σ_j . Let U_r denote

the $m \times r$ matrix with $\frac{1}{\sigma_j}$ (column j of AV) as its j th column. The r columns of U_r are then an orthonormal set. Now complete U_r to an $m \times m$ matrix by using an orthonormal basis for the orthogonal complement of the column space of U_r for the remaining $m - r$ columns. Hence

$$AV = U\Sigma$$

and hence $AV = U\Sigma V^*$. \square

There is other ways to write the SVD. Since only the first r diagonal entries of Σ are non-zero, the last $m - r$ columns of U and the last $n - r$ columns of V are superfluous. Let $\tilde{\Sigma}$ be the $r \times r$ matrix $\text{diag}(\sigma_1, \dots, \sigma_r)$. Replace the $n \times n$ matrix U and the $m \times m$ matrix V by the $(m - r) \times (m - r)$ matrix U_r and by the $r \times n$ matrix V_r consisting of the first r rows, respectively. Hence,

$$A = U_r \tilde{\Sigma} V_r.$$

Summary: Any matrix A has an SVD with a unique diagonal matrix Σ , but the unitary matrices U and V are not uniquely determined by the matrix A . It is just the way these unitaries are used that is specified: Namely, $A(\text{column } j \text{ of } V) = \sigma_j(\text{column } j \text{ of } U)$, or in matrix form:

$$AV = U\Sigma V^*.$$

DEFINITION 6.1.13. The vectors u_1, u_2, \dots, u_m and v_1, \dots, v_n are called the *left* and *right singular vectors*. Based on our results implying the Fredholm alternative the property of singular vectors is not surprising:

PROPOSITION 6.1.14. *Let A be a $m \times n$ matrix of rank r . Then*

$$\begin{aligned} \text{ran}(A) &= \text{span}\{u_1, \dots, u_r\}, \ker(A^*) = \text{span}\{u_{r+1}, \dots, u_m\} \\ \text{ran}(A^*) &= \text{span}\{v_1, \dots, v_r\}, \ker(A) = \text{span}\{v_{r+1}, \dots, v_n\}. \end{aligned}$$

Hence we have

$$\text{ran}(A) \oplus \ker(A^*) = \mathbb{C}^m$$

and

$$\text{ran}(A^*) \oplus \ker(A) = \mathbb{C}^n.$$

Or in terms of basis: The columns of V^* are an orthonormal basis for \mathbb{C}^n and the columns of U are an orthonormal basis for \mathbb{C}^m . Then A maps the j th basis vector of \mathbb{C}^n to a multiple of the j th basis vector of \mathbb{C}^m , where the multiplier is given by the singular value σ_j . If we order the singular values decreasingly, then σ_1 is the largest factor by which the length of a basis vector is multiplied. We now show that this is the largest factor by which the length of any vector is multiplied. In other words, the operator norm of the linear transformation induced by A is equal to the largest singular value. The operator norm of a matrix is often known as the *spectral norm*.

PROPOSITION 6.1.15. *Let A be a $m \times n$ matrix with singular values $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$. Then the operator norm of A equals σ_1 :*

$$\|A\| = \sigma_1.$$

PROOF. The equation $AV = U\Sigma$ gives for the first column vector v_1 of V that $\|Av_1\| = \sigma_1$. Let x be a vector of length one in \mathbb{C}^n . Then the SVD gives $Ax = U\Sigma V^*x$. Since V is unitary, also V^* is unitary and hence an isometry. Let us denote $V^*x = y$. Then $\|y\| = 1$ and the vector Σy is the vector where the j th component gets multiplied by σ_j . Hence $\|\Sigma y\| \leq \sigma_1\|y\|$. Since U is unitary

$$\|Ax\| = \|U\Sigma y\| \leq \sigma_1.$$

□

A complex number may be written in polar form $z = |z|e^{2\pi i\varphi}$. The polar decomposition of a matrix A decomposes it as a product of a unitary matrix and a positive definite matrix. If one looks at the eigenvalues of these matrices, then the first one has only eigenvalues of modulus one and the other has only positive eigenvalues. Hence in terms of the spectrum of the matrices the polar decomposition is a natural generalization of the one for complex numbers.

PROPOSITION 6.1.16 (Polar decomposition). *Given a $n \times n$ matrix A . There exist a unitary matrix U and a positive definite matrix R such that*

$$A = UR.$$

PROOF. The SVD decomposition gives us unitary $n \times n$ matrices U and V such that

$$A = U\Sigma V^* = UV^*V\Sigma V^*.$$

Note that UV^* is unitary as a product of two unitary matrices and $V\Sigma V^*$ is positive definite, since Σ is positive definite. Hence $V\Sigma V^*$ is the replacement of the length of a complex number and UV^* the one for the phase factor. □

Consequently, the SVD gives in a straightforward manner the polar decomposition. There is also a version of this result for general bounded operators on a Hilbert space.

EXAMPLE 6.1.17. Determine the singular value decomposition of $\begin{pmatrix} 3 & 2 & 2 \\ 2 & 3 & -2 \end{pmatrix}$.

Write

$$A = \begin{pmatrix} 3 & 2 & 2 \\ 2 & 3 & -2 \end{pmatrix}$$

We follow the procedure for singular value decomposition. We have

$$A^*A = \begin{pmatrix} 3 & 2 \\ 2 & 3 \\ 2 & -2 \end{pmatrix} \begin{pmatrix} 3 & 2 & 2 \\ 2 & 3 & -2 \end{pmatrix} = \begin{pmatrix} 13 & 12 & 2 \\ 12 & 13 & -2 \\ 2 & -2 & 8 \end{pmatrix}$$

We find the eigenvalues of A^*A by solving

$$\det \begin{pmatrix} 13 - \lambda & 12 & 2 \\ 12 & 13 - \lambda & -2 \\ 2 & -2 & 8 - \lambda \end{pmatrix} = 0$$

We get the characteristic equation

$$\begin{aligned} & (13 - \lambda)(13 - \lambda)(8 - \lambda) + 12 \cdot (-2) \cdot 2 + 2 \cdot 12 \cdot (-2) \\ & - (13 - \lambda)(-2)(-2) - 12 \cdot 12 \cdot (8 - \lambda) - 2 \cdot (13 - \lambda) \cdot 2 \end{aligned}$$

which simplifies to $-\lambda(\lambda - 9)(\lambda - 25)$ and has the solutions $\lambda_1 = 25$, $\lambda_2 = 9$ and $\lambda_3 = 0$. We now find normalized eigenvectors for each eigenvalue.

$$\underline{\lambda_1 = 25}$$

$$\begin{pmatrix} 13-25 & 12 & 2 \\ 12 & 13-25 & -2 \\ 2 & -2 & 8-25 \end{pmatrix} \sim \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 1 \\ 1 & -1 & -\frac{17}{2} \end{pmatrix}$$

$$v_1 = \begin{pmatrix} \frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} \\ 0 \end{pmatrix} \text{ is a normalized eigenvector for } \lambda_1 = 25.$$

$$\underline{\lambda_2 = 9}$$

$$\begin{pmatrix} 13-9 & 12 & 2 \\ 12 & 13-9 & -2 \\ 2 & -2 & 8-9 \end{pmatrix} \sim \begin{pmatrix} 0 & 0 & 0 \\ 0 & 1 & \frac{1}{4} \\ 1 & 0 & -\frac{1}{4} \end{pmatrix}$$

$$v_2 = \begin{pmatrix} \frac{\sqrt{2}}{6} \\ -\frac{\sqrt{2}}{6} \\ \frac{2\sqrt{2}}{3} \end{pmatrix} \text{ is a normalized eigenvector for } \lambda_2 = 9.$$

$$\underline{\lambda_3 = 0}$$

$$\begin{pmatrix} 13 & 12 & 2 \\ 12 & 13 & -2 \\ 2 & -2 & 8 \end{pmatrix} \sim \begin{pmatrix} 0 & 0 & 0 \\ 0 & 1 & -2 \\ 1 & 0 & 2 \end{pmatrix}$$

$$v_3 = \begin{pmatrix} \frac{2}{3} \\ -\frac{2}{3} \\ -\frac{1}{3} \end{pmatrix} \text{ is a normalized eigenvector for } \lambda_3 = 0.$$

We get the singular value decomposition $A = U\Sigma V^*$, where

$$V = (v_1|v_2|v_3) = \begin{pmatrix} \frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{6} & \frac{2}{3} \\ \frac{\sqrt{2}}{2} & -\frac{\sqrt{2}}{6} & -\frac{2}{3} \\ 0 & \frac{2\sqrt{2}}{3} & -\frac{1}{3} \end{pmatrix},$$

$$\Sigma = \begin{pmatrix} \sigma_1 & 0 & 0 \\ 0 & \sigma_2 & 0 \end{pmatrix} = \begin{pmatrix} \sqrt{\lambda_1} & 0 & 0 \\ 0 & \sqrt{\lambda_2} & 0 \end{pmatrix} = \begin{pmatrix} 5 & 0 & 0 \\ 0 & 3 & 0 \end{pmatrix},$$

and

$$U = (U_1|U_2) = \left(\frac{Av_1}{\|Av_1\|} \mid \frac{Av_2}{\|Av_2\|} \right) = \begin{pmatrix} \frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} & -\frac{\sqrt{2}}{2} \end{pmatrix}.$$

Explicitly, we have

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

6.1.3. Pseudoinverse and least squares method. Given a $m \times n$ matrix A and a vector $b \in \mathbb{C}^m$. Then we are interested in solutions of

$$Ax = b.$$

There will be a solution, if b lies in the column space of A or in other words in the range of A . If that is not the case, there exists no solution, but we still can project b onto the range of A . Given our knowledge about projections in Hilbert spaces, this vector will be the best approximation of all vectors out of the range of A . This projection may also be expressed in terms of the matrix and the resulting object if the pseudoinverse of A . Since in the Euclidean case this amounts to expressions with squares. This method is known as *least squares*. We will approach this circle of ideas from the SVD.

More formally, we are interested in the solutions of the following problem:

$$\min_{x \in \mathbb{C}^n} \|Ax - b\|_2.$$

We denote the set of these optimal solutions to $Ax = b$ by X_{opt} .

PROPOSITION 6.1.18.

$$X_{\text{opt}} = \ker(A) + A^\dagger b,$$

where A^\dagger is the pseudoinverse of A . In terms of the reduced SVD of A its pseudoinverse equals

$$A^\dagger = V_r \tilde{\Sigma}^{-1} U_r^*,$$

and $x_{MN} = A^\dagger b$ is the minimal norm solution.

If A has full column rank, meaning $m > n$, then $A^\dagger = (A^*A)^{-1}A^*$.

PROOF. For the $m \times n$ matrix A we use its SVD $A = U\Sigma V^*$ to rewrite our minimization problem

$$\|Ax - b\|_2^2 = \|U\Sigma V^*x - UU^*b\|_2^2 = \|\Sigma V^*x - U^*b\|_2^2,$$

since the Euclidean norm is invariant under unitary transformations. Now we make changes of variables:

$$\tilde{x} = V^*x \quad \text{and} \quad \tilde{b} = U^*b.$$

Then our optimal solutions are minimizing

$$\min \|\Sigma \tilde{x} - \tilde{b}\|_2^2.$$

The solution depends on the r non-zero eigenvalues. Hence we split our vectors into blocks of length r and $n - r$, respectively. Given $x = (x_1, \dots, x_r, x_{r+1}, \dots, x_n)$. We denote by

$$\underline{x}_r = (x_1, \dots, x_r) \quad \text{and} \quad \underline{x}_{n-r} = (x_{r+1}, \dots, x_n)$$

and analogously for b . Then our optimal solutions X^* are of the form

$$x^* = \min_{\underline{x}_r} \|\tilde{\Sigma} \tilde{x}_r - \tilde{b}_r\|_2^2 + \|\tilde{b}_{m-r}\|_2^2.$$

Since $\tilde{\Sigma}$ is invertible, the optimal choice is $\tilde{x}_r = \tilde{\Sigma}^{-1} \tilde{b}_r$. Hence the minimum is determined by \tilde{b}_{m-r} , which is going to be optimal for all \tilde{b} with $\tilde{b}_{r+1} = \dots = \tilde{b}_m$. Consequently, the optimal solutions are determined by the first r components of b and so the solutions have $n - r$ free variables x_{r+1}^*, \dots, x_n^* . Since $x = V\tilde{x}$ we have

that the first r components of the optimal solutions x are $V_r \tilde{x}_r$ and we also have $\tilde{b}_r = \tilde{U}_r^* b$. Thus an optimal solution is of the form

$$x^* = V_r \tilde{\Sigma}^{-1} U_r^* b + V_{n-r} z, \quad z \in \mathbb{C}^{n-r}.$$

Note that $V_r \tilde{\Sigma}^{-1} U_r^* b = A^\dagger b$. The mapping $b \rightarrow A^\dagger b$ is the orthogonal projection of b onto the range of A . The last $n - r$ columns of U are an orthonormal basis for $\ker(A)$. Hence we have established the desired assertion. Furthermore, the Projection Theorem implies that this solution is the unique solution of minimal norm.

The condition of full column rank, $r = n \leq m$, is equivalent to $\ker(A) = \{0\}$ and in this case $(A^* A)^{-1}$ exists. Furthermore, this yields that the unique optimal solution is

$$\tilde{x} = A^\dagger b = (A^* A)^{-1} A^* b.$$

□

Another way to arrive at the statement for matrices of full column rank is to recall that minimal norm solutions have the property that $b - Ax$ is orthogonal to the range of A . Since the orthogonal complement of the range space of A is the row space of A^* we have

$$A^*(b - Ax) = 0$$

or

$$A^* Ax = A^* b,$$

which are the *normal equations* for your linear system. If A has full column rank, we can invert $A^* A$ and hence our optimal solution is given by $(A^* A)^{-1} A^* b$. Another way of putting it, is that the pseudoinverse A^\dagger of a matrix with *full column rank* is given by $(A^* A)^{-1} A^*$.

REMARK 6.1.19. The name pseudoinverse has its origins in the fact that A^\dagger is a left inverse for A with full column rank but not a right inverse: $A^\dagger A = I$ but $AA^\dagger \neq I$, the latter actually describes the orthogonal projection onto the range of A . In the case of matrices of full column rank one may compute it explicitly: $(A^* A)^{-1} A^* A = I$

EXAMPLE 6.1.20. Solve the equation

$$-x_1 + 2x_2 + 2x_3 = b, \quad \text{for } b \in \mathbb{R},$$

and explain in which sense your result has to be interpreted.

We let $A = \begin{pmatrix} -1 & 2 & 2 \end{pmatrix}$ and rewrite the equation as $Ax = b$. The Singular Value Decomposition gives that

$$A = U \Sigma V^*,$$

where

$$U = (1), \quad \Sigma = \begin{pmatrix} 3 & 0 & 0 \end{pmatrix}, \quad V = \begin{pmatrix} -\frac{1}{3} & \frac{2}{\sqrt{5}} & \frac{2}{3\sqrt{5}} \\ \frac{2}{3} & 0 & \frac{\sqrt{5}}{3} \\ \frac{2}{3} & \frac{1}{\sqrt{5}} & \frac{4}{3\sqrt{5}} \end{pmatrix}.$$

The pseudoinverse of A is

$$A^\dagger = V\Sigma^+U^* = \begin{pmatrix} -\frac{1}{3} & \frac{2}{\sqrt{5}} & \frac{2}{3\sqrt{5}} \\ \frac{2}{3} & 0 & \frac{\sqrt{5}}{3} \\ \frac{2}{3} & \frac{1}{\sqrt{5}} & \frac{4}{3\sqrt{5}} \end{pmatrix} \begin{pmatrix} \frac{1}{3} \\ 0 \\ 0 \end{pmatrix} (1) = \begin{pmatrix} -\frac{1}{9} \\ \frac{2}{9} \\ \frac{2}{9} \end{pmatrix}$$

The solutions of the equation $Ax = b$ are given by

$$x = A^\dagger b + \ker A = \begin{pmatrix} -\frac{1}{9} \\ \frac{2}{9} \\ \frac{2}{9} \end{pmatrix} b + \ker A.$$

6.1.4. Nilpotent operators. An ingredient in understanding the Jordan normal form (JNF) of a linear transformation are nilpotent operators.

DEFINITION 6.1.21. Let T be a linear transformation on a finite-dimensional vector space X . We say that T is *nilpotent* if there exists a power of the matrix such that $T^k = 0$. The minimal exponent, e , such that $T^e = 0$ and $T^{e-1} \neq 0$, is the *index of nilpotency* of T .

The matrix N_p defined by

$$N_p = \begin{pmatrix} 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 1 & \cdots & 0 \\ \vdots & \ddots & \ddots & 1 & \vdots \\ 0 & 0 & 0 & \vdots & 1 \\ 0 & 0 & 0 & \vdots & 0 \end{pmatrix}$$

is a nilpotent matrix of index $p - 1$.

PROPOSITION 6.1.22. Let T be a linear transformation on a finite-dimensional vector space X . Then T is nilpotent if and only if the spectrum $\sigma(T) = \{0\}$. In other words T is nilpotent if and only if 0 is the only eigenvalues of T .

PROOF. (\Leftarrow) Suppose T is nilpotent and λ is an eigenvalue of T . Then there exists a non-zero $x \in X$ such that $Tx = \lambda x$. Then there exists a p such that

$$0 = T^p x = \lambda^p x$$

and hence $T^p = 0$ implies that $\lambda = 0$.

(\Rightarrow) Suppose $\sigma(T) = \{0\}$. Then T is similar to a triangular matrix with all zeros on the diagonal. The powers of an upper-triangular matrix become eventually the zero matrix. Hence T is nilpotent. \square

LEMMA 6.19. Let N be matrix such that $N^{k-1} \neq 0$ and $N^k = 0$, i.e. N is a nilpotent matrix. Then $I - N$ is invertible and its inverse is given by

$$(I - N)^{-1} = I + N + N^2 + \cdots + N^{k-1}.$$

PROOF. a) It suffices to show that

$$(I - N)(I + N + N^2 + \cdots + N^{k-1}) = I$$

and

$$(I + N + N^2 + \cdots + N^{k-1})(I - N) = I.$$

Indeed, we have

$$\begin{aligned} & (I - N)(I + N + N^2 + \cdots + N^{k-1}) \\ &= I + N + N^2 + \cdots + N^{k-1} - (N + N^2 + N^3 + \cdots + N^k) \\ &= I - N^k \\ &= I \end{aligned}$$

and

$$\begin{aligned} & (I + N + N^2 + \cdots + N^{k-1})(I - N) \\ &= I + N + N^2 + \cdots + N^{k-1} - (N + N^2 + N^3 + \cdots + N^k) \\ &= I - N^k \\ &= I. \end{aligned}$$

□

EXAMPLE 6.1.23. Find the singular value decomposition of

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

Describe the kernel and range of A and A^* in terms of the left and right singular vectors.

The procedure for singular value decomposition gives us $A = U\Sigma V^*$, where

$$U = \begin{pmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix}, \quad \Sigma = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix}, \quad V = \begin{pmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{pmatrix}.$$

The range of A is spanned by the left singular vectors corresponding to the non-zero singular values, i.e.

$$\text{ran}(A) = \text{span} \left\{ \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \right\}.$$

The kernel of A is spanned by the right singular vectors corresponding to the vanishing singular values, i.e.

$$\text{ker}(A) = \text{span} \left\{ \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \right\}.$$

The singular value decomposition of A^* is $A^* = V\Sigma U^*$, so we get

$$\text{ran}(A^*) = \text{span} \left\{ \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \right\}.$$

and

$$\text{ker}(A^*) = \text{span} \left\{ \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} \right\}.$$

6.1.5. Jordan Normal Form. The key objective is to describe a refinement of Schur's form, the *Jordan normal form* of a linear operator on a finite-dimensional vector space X . Suppose A is the matrix representation of T with respect to a basis in X .

Given a $n \times n$ matrix A with distinct eigenvalues $\lambda_1, \dots, \lambda_k$ with $k \leq n$. Then A is unitarily equivalent to

$$\begin{pmatrix} J_1(\lambda_1) & 0 & \cdots & 0 \\ 0 & J_2(\lambda_2) & \ddots & 0 \\ \vdots & \ddots & \vdots & \\ 0 & \cdots & 0 & J_k(\lambda_k) \end{pmatrix}$$

where a general upper triangular matrix T_i is replaced by a *Jordan block*:

$$J_i(\lambda_i) = \begin{pmatrix} \lambda_i & 1 & 0 & \cdots & 0 \\ 0 & \lambda_i & 1 & \ddots & 0 \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & \cdots & \ddots & \ddots & 1 \\ 0 & \cdots & \cdots & 0 & \lambda_i \end{pmatrix}.$$

The basis of X that allows us to express T in this particular “almost”-diagonal form, is the main focus of this section. For this purpose we have to extend the definition of eigenspace to *generalized* eigenspace.

DEFINITION 6.1.24. Let T be a linear transformation on a vector space X .

- (1) A non-zero vector $x \in X$ is called a *generalized eigenvector* of T corresponding to a scalar λ if

$$(T - \lambda I)^p x = 0$$

for some positive integer p .

- (2) Suppose X is n dimensional and A is the matrix representation of T for a basis of X . A non-zero vector $x \in \mathbb{C}^n$ is called a *generalized eigenvector* of a $n \times n$ matrix A corresponding to the scalar λ if $(T - \lambda I)^p x = 0$ for some positive integer p .
- (3) The *generalized eigenspace* \tilde{E}_λ corresponding to λ is

$$\tilde{E}_\lambda = \{x \in X : (A - \lambda I)^p x = 0 \text{ for some positive integer } p\}.$$

Note that \tilde{E}_λ consists of the zero vector and all generalized eigenvectors corresponding to λ , since $\tilde{E}_\lambda = \ker((T - \lambda I)^p)$. Furthermore, let p be the smallest positive integer such that $(T - \lambda I)^p x = 0$, then $(T - \lambda I)^{p-1} x \neq 0$ and is an eigenvector of T corresponding to λ (since $0 = (T - \lambda I)^p x = (T - \lambda I)(T - \lambda I)^{p-1} x$ and hence $y = (T - \lambda I)^{p-1} x$ satisfies $Ty = \lambda y$). Hence, the scalars in the definition of generalized eigenvectors and generalized eigenspaces are eigenvalues of T , as the name suggested. Consequently, $T - \lambda I$ is a nilpotent operator of exponent p with eigenvalue λ .

DEFINITION 6.1.25. A subspace M of X is called *T -invariant* for a linear operator T if $T(M) \subseteq M$.

The range of a linear operator is such an invariant subspace. We collect a few properties of generalized eigenspaces without proof.

PROPOSITION 6.1.26. *Let T be a linear operator on a vector space X . Suppose λ is an eigenvalue of T .*

- \tilde{E}_λ is a T -invariant subspace of X containing the eigenspace E_λ corresponding to λ .
- For any scalar μ different from the eigenvalue λ , the restriction of $(T - \mu I)$ to \tilde{E}_λ is injective.

DEFINITION 6.1.27. Suppose T is a linear operator on a finite-dimensional vector space X and let λ be an eigenvalue of T .

- (1) If the characteristic polynomial contains a factor of the form $(x - \lambda)^a$. Then we call a the *algebraic multiplicity* of the eigenvalue λ .
- (2) The *geometric multiplicity* g of the eigenvalue λ equals the dimension of the eigenspace associated with λ : $g := \dim(\ker(T - \lambda I))$.

Observe that the algebraic multiplicity of an eigenvalue equals the number of times λ appears on the diagonal of the upper-triangular matrix in the Schur form. Note that the geometric multiplicity of an eigenvalue is always less than or equal to the algebraic multiplicity. In case the sum of the geometric multiplicities is less than the sum of the algebraic multiplicities, then T has not enough eigenvectors to form a basis for X and the T is not invertible.

PROPOSITION 6.1.28. *Suppose T is a linear operator on a complex finite-dimensional vector space X . Given an eigenvalues of T with algebraic multiplicity a . Then $\tilde{E}_\lambda = \ker((T - \lambda I)^a)$.*

The proof is omitted, since it is not essential for understanding the construction of Jordan blocks. The next statement is a crucial observation towards the Jordan normal form.

EXAMPLE 6.1.29. Find the generalized eigenspaces of

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

We start by finding the eigenvectors of A . The characteristic polynomial is

$$\begin{vmatrix} \lambda - 1 & -2 & -3 \\ 0 & \lambda - 1 & -1 \\ 0 & 1 & \lambda + 1 \end{vmatrix} = (\lambda - 1)((\lambda - 1)(\lambda + 1) + 1) = \lambda^2(\lambda - 1)$$

We have the eigenvalues $\lambda = 0$ with algebraic multiplicity 2 and $\lambda = 1$ with algebraic multiplicity 1. We find the generalized eigenspace for each eigenvalue.

$\lambda = 0$

The generalized eigenspace of $\lambda = 0$ is $\ker(A^2)$. We have

$$A^2 = \begin{pmatrix} 1 & 1 & 2 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix},$$

and the kernel of A^2 , i.e. the solutions to $A^2x = 0$ are

$$x = \begin{pmatrix} 1 \\ -1 \\ 0 \end{pmatrix} r + \begin{pmatrix} 2 \\ 0 \\ -1 \end{pmatrix} s, \quad r, s \in \mathbb{C}.$$

Hence, the generalized eigenspace of $\lambda = 0$ is

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

$\lambda = 1$

The generalized eigenspace of $\lambda = 1$ is $\ker(A - I)$. We have

$$A - I = \begin{pmatrix} 0 & 2 & 3 \\ 0 & 0 & 1 \\ 0 & -1 & -2 \end{pmatrix},$$

and the kernel of $A - I$, i.e. the solutions to $(A - I)x = 0$ are

$$x = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} r, \quad r \in \mathbb{C}.$$

Hence, the generalized eigenspace of $\lambda = 1$ is

$$\text{span} \left\{ \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \right\}.$$

THEOREM 6.20. *Suppose T is a linear operator on a complex finite-dimensional vector space X . Let $\lambda_1, \dots, \lambda_k$ be the distinct eigenvalues of T with corresponding algebraic multiplicities a_i . If \mathcal{B}_i is a basis for \tilde{E}_{λ_i} for $1 \leq i \leq k$, then we have*

- (1) *The bases $\mathcal{B}_1, \dots, \mathcal{B}_k$ are pair-wise disjoint; $\mathcal{B}_i \cap \mathcal{B}_j = \emptyset$ for $i \neq j$.*
- (2) *$\mathcal{B} = \mathcal{B}_1 \cup \dots \cup \mathcal{B}_k$ is a basis for V .*
- (3) *$\dim(\tilde{E}_{\lambda_i}) = a_i$ for $i = 1, \dots, k$.*

A proof is in the textbook by Friedberg et al., Linear Algebra. As a consequence we state a criterion for diagonalizable operators.

COROLLARY 6.1.30. *Suppose T is a linear operator on a complex finite-dimensional vector space X . Then T is diagonalizable if and only if $\tilde{E}_{\lambda} = E_{\lambda}$ for all eigenvalues λ of T .*

The proof amounts to the fact that T is diagonalizable if and only if $\dim(E_{\lambda}) = \dim(\tilde{E}_{\lambda})$ for all eigenvalues λ of T . Now we have by definition $E_{\lambda} \subseteq \tilde{E}_{\lambda}$. Recall that two subspaces of a finite-dimensional vector space have the same dimension if and only if they are equal.

The problem is now reduced to the quest of finding bases for the generalized eigenspaces of a linear operator. We have observed that $T - \lambda I$ is a nilpotent operator of index equal to the algebraic multiplicity of λ . Recall that we have discussed a canonical construction of a basis associated to a nilpotent operator. Following Friedberg et al we define some notions related to these bases.

DEFINITION 6.1.31. Suppose T is a linear operator on a finite-dimensional vector space X and let x be a generalized eigenvector associated to an eigenvalue λ of T . Suppose p is the smallest positive integer such that $(T - \lambda I)^p x = 0$. Then the set

$$\Gamma = \{(T - \lambda I)^{p-1}x, (T - \lambda I)^{p-2}x, \dots, (T - \lambda I)x, x\}$$

is called a *cycle of generalized eigenvectors* of T corresponding to λ . The vector $(T - \lambda I)^{p-1}x$ is called the *initial vector* and x is known as the *end vector* of the cycle. The cardinality, p , of this set is called *length* of the cycle Γ .

Note that the initial vector $(T - \lambda I)^{p-1}x$ of a cycle of generalized eigenvectors of T is the only eigenvector of T in this cycle. If x is an eigenvector of T corresponding to the eigenvalue λ (see the remark after the definition of generalized eigenspaces), then we consider the eigenvector x as a cycle $\{x\}$ of length 1. Our discussion of nilpotent operators yields that Γ is linearly independent.

THEOREM 6.21. *Suppose T is a linear operator on a finite-dimensional vector space X .*

- (1) *Then each cycle Γ of generalized eigenvectors the subspace W spanned by Γ is T -invariant. Furthermore, the restriction of T to W has a matrix representation with respect to Γ that has the form of a Jordan block.*
- (2) *The generalized eigenspace \tilde{E}_λ corresponding to an eigenvalue λ of T has a basis consisting of a union of disjoint cycles of generalized eigenvectors.*
- (3) *There exists a basis of V constructed out of disjoint cycles of generalized eigenvectors corresponding to all the distinct eigenvalues of T with respect to which the matrix representation of T has Jordan canonical form.*

Let us discuss some examples.

EXAMPLE 6.1.32. Given the matrix

$$A = \begin{pmatrix} 3 & 1 & -2 \\ -1 & 0 & 5 \\ -1 & -1 & 4 \end{pmatrix}$$

. We find the characteristic polynomial to be

$$p_A(x) = \det(A - xI) = -(x - 3)(x - 2)^2,$$

so the eigenvalues of A are $\lambda_1 = 3$ with algebraic multiplicity 1 and $\lambda_2 = 2$ with algebraic multiplicity 2. By our general result on the dimension of generalized eigenspaces we have $\dim(\tilde{E}_{\lambda_1}) = 1$ and $\dim(\tilde{E}_{\lambda_2}) = 2$ and $\tilde{E}_{\lambda_1} = \ker(A - 3I)$, $\tilde{E}_{\lambda_2} = \ker(A - 2I)^2$. Since the generalized eigenspace for λ_1 is one-dimensional,

it is equal to the eigenspace which is spanned by the eigenvector $\begin{pmatrix} -1 \\ 2 \\ 1 \end{pmatrix}$. Hence

$\mathcal{B}_1 = \text{span}\left\{\begin{pmatrix} -1 \\ 2 \\ 1 \end{pmatrix}\right\}$ is a basis for \tilde{E}_{λ_1} . The basis of \tilde{E}_{λ_2} is either a union of

two cycles of length 1 or a single cycle of length 2. Since cycles of length 1 are eigenvectors for A , this is not possible. Hence we have a single cycle of length 2. Now, a vector $v \in \mathbb{C}^3$ is the end vector of such a cycle if and only if

$$(A - 2I)v \neq 0 \text{ and } (A - 2I)^2v = 0.$$

In other words, what are the solutions of $(A - 2I)^2x = 0$? These are spanned by the vectors $\begin{pmatrix} 1 \\ -3 \\ -1 \end{pmatrix}$ and $\begin{pmatrix} -1 \\ 2 \\ 0 \end{pmatrix}$. Observe that $v = \begin{pmatrix} -1 \\ 2 \\ 0 \end{pmatrix}$ satisfies $(A - 2I)v = \begin{pmatrix} 1 \\ -3 \\ -1 \end{pmatrix}$ and so our cycle of generalized eigenvectors is

$$\mathcal{B}_2 = \left\{ \begin{pmatrix} -1 \\ 2 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ -3 \\ -1 \end{pmatrix} \right\}$$

a basis for \tilde{E}_{λ_2} . Then union of $\mathcal{B}_1 \cup \mathcal{B}_2$ is a basis

$$\mathcal{B} = \left\{ \begin{pmatrix} -1 \\ 2 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ -3 \\ -1 \end{pmatrix}, \begin{pmatrix} -1 \\ 2 \\ 0 \end{pmatrix} \right\}$$

with respect to which A has Jordan canonical form:

$$[T]_{\mathcal{B}} = \begin{pmatrix} 3 & 0 & 0 \\ 0 & 2 & 1 \\ 0 & 0 & 2 \end{pmatrix}.$$

The matrix A is similar to $[T]_{\mathcal{B}}$:

$$[T]_{\mathcal{B}} = Q^{-1}AQ,$$

where Q is the matrix $Q = \begin{pmatrix} -1 & 1 & -1 \\ 2 & -3 & 2 \\ 1 & -1 & 0 \end{pmatrix}$, whose columns are the vectors of the basis \mathcal{B} and we have $Q^{-1} = \begin{pmatrix} -2 & -1 & 1 \\ -2 & -1 & 0 \\ -1 & 0 & -1 \end{pmatrix}$.

EXAMPLE 6.1.33. Let T be the linear operator on \mathcal{P}_2 defined by

$$Tf(x) = -f(x) - f'(x).$$

In the monomial basis $\mathcal{M} = \{1, x, x^2\}$ for \mathcal{P}_2 we have

$$A = [T]_{\mathcal{M}} = \begin{pmatrix} -1 & -1 & 0 \\ 0 & -1 & -2 \\ 0 & 0 & -1 \end{pmatrix}$$

and the characteristic polynomial is $p_T(x) = -(x + 1)^3$. Hence $\lambda = -1$ is an eigenvalue of algebraic multiplicity 3 and so we have $\mathcal{P}_2 = \tilde{E}_{\lambda=-1}$. Consequently, \mathcal{M} is a basis for $\tilde{E}_{\lambda=-1}$. Observe that

$$\dim(E_{\lambda=2}) = 3 - \text{rank}(A + I) = 3 - 2 = 1.$$

A basis for $\tilde{E}_{\lambda=-1}$ cannot be the union of two or three cycles because the initial vector of each cycle is an eigenvector. Since there are no two linearly independent eigenvectors, we must have a single cycle Γ of length 3. Thus Γ determines a single Jordan block of size 3, which in our case has the form

$$[T]_{\Gamma} = \begin{pmatrix} -1 & 1 & 0 \\ 0 & -1 & 1 \\ 0 & 0 & -1 \end{pmatrix}.$$

How does a basis Γ of $\tilde{E}_{\lambda=-1}$ look like for which T has Jordan normal form? Recall the discussion of canonical bases associated to nilpotent operators. In our example we take $f(x) = x^2$. Then

$$\Gamma = \{(T+1)^2 f(x), (T+1)f(x), f(x)\} = \{2, -2x, x^2\}.$$

EXAMPLE 6.1.34. Determine the Jordan normal form of $\begin{pmatrix} 1 & 1 & 1 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$. We start

by finding the eigenvalues of the matrix. Since the matrix is triangular, its eigenvalues are the diagonal entries, so $\lambda_1 = 1$, $\lambda_2 = 1$ and $\lambda_3 = 1$. We find the eigenvectors corresponding to $\lambda = 1$:

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

We see that the solutions are given by

$$x = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} r + \begin{pmatrix} 0 \\ 1 \\ -1 \end{pmatrix} s, \quad r, s \in \mathbb{C}.$$

Thus the geometric multiplicity of the eigenvalue $\lambda = 1$ is two. This means that there are two blocks in the Jordan normal form, and it follows that the Jordan normal form is

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

EXAMPLE 6.1.35 (ODEs and Jordan normal form). Solving ordinary differential equations is a well-known applications of the Jordan normal form (JNF). Here we treat the 2×2 case. Given the system

$$\begin{pmatrix} x_1' \\ x_2' \end{pmatrix} = \begin{pmatrix} \lambda & 1 \\ 0 & \lambda \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$

with initial values $x_1(0)$ and $x_2(0)$ determining the solutions $x_1(t)$ and $x_2(t)$. Explicitly, we want to solve

$$\begin{aligned} x_1' &= \lambda x_1 + x_2 \\ x_2' &= \lambda x_2 \end{aligned}$$

by backward substitution. We have

$$x_2(t) = x_2(0)e^{\lambda t}$$

and

$$x_1'(t) = \lambda x_1(t) + x_2(t) = \lambda x_1(t) + x_2(0)e^{\lambda t}.$$

Hence

$$x_1'(t) - \lambda x_1(t) = x_2(0)e^{\lambda t}$$

which becomes

$$x_2(0) = e^{-\lambda t}(x_1'(t) - \lambda x_1(t)).$$

Note that

$$(e^{-\lambda t} x_1(t))' = e^{-\lambda t}((x_1'(t) - \lambda x_1(t))).$$

Thus we have

$$x_2(0) = (e^{-\lambda t} x_1(t))',$$

which after integration becomes

$$x_2(0)t = e^{-\lambda t}x_1(t) + e^{-\lambda t}x_1(0)$$

and after solving for $x_1(t)$:

$$x_1(t) = (x_1(0) + x_2(0)t)e^{\lambda t}.$$

6.1.6. Minimal polynomials. Let T be a linear operator on a finite-dimensional vector space. Then we say that a polynomial p *annihilates* T if $p(T) = 0$, the zero matrix. The set of all *annihilating polynomials* is non-empty, since the theorem of Cayley-Hamilton shows that the characteristic polynomial p_T is annihilating T . Let us now consider the monic polynomial of least degree that annihilates T , the *minimal polynomial* m_T . A polynomial is called *monic* if it is of the form $x^n + a_{n-1}x^{n-1} + \cdots + a_0$.

THEOREM 6.22. *Let m_T be the minimal polynomial of a linear operator T on a finite-dimensional vector space.*

- (1) *Suppose p annihilates T . Then m_T divides p . Hence m_T is a divisor of the characteristic polynomial p_T .*
- (2) *The minimal polynomial m_T is unique.*

PROOF. (1) By the division algorithm for polynomials there exist polynomials q and r such that

$$p(x) = q(x)m_T(x) + r(x),$$

where the degree of r is less than the degree of m_T . Now we have

$$p(T) = q(T)m_T(T) + r(T),$$

which by assumption equals

$$0 = p(T) = q(T)0 + r(T).$$

Consequently, $r(T) = 0$, which is impossible because m_T is the polynomial of least degree that annihilates T , so $r = 0$.

- (2) Suppose m_1 and m_2 are two minimal polynomials of T . Then m_1 divides m_2 , but both polynomials have the same degree. Hence $m_2(t) = cm_1(t)$ for some non-zero scalar c . By definition minimal polynomials are monic, so $c = 1$ and we have $m_1 = m_2$. □

There is a relation between the characteristic polynomial, the minimal polynomial and the eigenvalues of a linear operator.

PROPOSITION 6.1.36. *Let T be a linear operator on a finite-dimensional vector space. Then a scalar λ is an eigenvalue of T if and only if $m_T(\lambda) = 0$. In other words, the characteristic polynomial and the minimal polynomial have the same zeros.*

PROOF. (\Leftarrow) Let p_T be the characteristic polynomial of T . We know that the minimal polynomial m_T divides p_T , so there exists a polynomial q such that $p_T(x) = q(x)m_T(x)$. If λ is a zero of the minimal polynomial, then

$$p_T(\lambda) = q(\lambda)m_T(\lambda) = 0$$

and λ is an eigenvalue of T .

(\Rightarrow) Suppose that λ is an eigenvalue of T and x an associated eigenvector. Then

$$0 = m_T(x) = m_T(\lambda)x$$

for a non-zero vector x . Thus $m_T(\lambda) = 0$, i.e. λ is a zero of the minimal polynomial. \square

COROLLARY 6.1.37. *Let T be a linear operator on a finite-dimensional vector space with distinct eigenvalues $\lambda_1, \dots, \lambda_k$. Suppose that the characteristic polynomial is of the form*

$$p_T(x) = (x - \lambda_1)^{n_1} \cdots (x - \lambda_k)^{n_k}.$$

Then there exist integers m_1, \dots, m_k such that $1 \leq m_i \leq n_i$ for $i = 1, \dots, k$ and the minimal polynomial is

$$m_T(x) = (x - \lambda_1)^{m_1} \cdots (x - \lambda_k)^{m_k}.$$

The integers n_i are the algebraic multiplicities of λ_i and m_i are equal to the geometric multiplicities of λ_i for $i = 1, \dots, k$.

COROLLARY 6.1.38. *Let T be a linear operator on a finite-dimensional vector space. Then T is diagonalizable if and only if the minimal polynomial is of the form*

$$m_T(x) = (x - \lambda_1) \cdots (x - \lambda_k)$$

EXAMPLE 6.1.39. Let D be the differentiation operator on \mathcal{P}_2 with the monomial basis \mathcal{M} . Then

$$[D]_{\mathcal{M}} = \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 2 \\ 0 & 0 & 0 \end{pmatrix}.$$

The characteristic polynomial $p_D(x) = -x^3$. Now $D(x^2) \neq 0$, hence $D^2 \neq 0$ and $m_D(x) = x^3$.

EXAMPLE 6.1.40. $\begin{pmatrix} 1 & 1 & 1 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$. Since its Jordan normal form is

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

the minimal polynomial is $(x - 1)^2$.

Metric spaces

7.1. Metric spaces

Metric spaces are generalizations of the real line, or more generally of normed spaces.

DEFINITION 7.1.1. A set X is called a *metric space* if there exists a function $d : X \times X \rightarrow [0, \infty)$ satisfying

- (1) $d(x, y) = 0$ if and only if $x = y$;
- (2) $d(x, y) = d(y, x)$
- (3) $d(x, z) \leq d(x, y) + d(y, z)$.

The function d is called a *metric* on X . The metric space is often denoted by (X, d) .

The class of metrics is much richer as metrics arising from norms on vector spaces. Here are two examples.

EXAMPLES 7.1.2. (1) *Hamming distance*: Let X be a set of n -tuples (x_1, \dots, x_n) , where x_i is either 0 or 1 for $i = 1, \dots, n$. Given two elements $x, y \in X$. We define the metric

$$d_H(x, y) = \text{the number of components } x_i \text{ such that } x_i \neq y_i,$$

aka the Hamming metric and of relevance in coding theory where it serves as a measure of how much a message gets distorted during transmission.

(2) *Discrete metric*: Given a set X . The *discrete metric* is defined by

$$d(x, y) = \begin{cases} 1 & \text{if } x \neq y \\ 0 & \text{if } x = y. \end{cases}$$

The discrete metric has properties different from the ones we are used to from \mathbb{R}^n .

There are ways to produce new metrics from old ones.

LEMMA 7.1. *If (Y, d) is a metric space and $h : X \rightarrow Y$ is an injective map, then $d_*(x, y) = d(h(x), h(y))$ defines a metric on X .*

PROOF. Since h is injective, the metric properties of d transfer to the ones for d_* . For example, $d_*(x, y) = 0$ if and only if $h(x) = h(y)$, which holds only for $x = y$ by the injectivity of h . \square

An example of interest is the $X = (-\pi/2, \pi/2)$, $Y = \mathbb{R}$ and $h = \tan$.

DEFINITION 7.1.3. Let (X, d) be a metric space. Then $B_r(x) = \{y \in X : d(x, y) < r\}$ is the *open ball* centered at x and of radius r .

An important class of normed spaces are the Banach spaces and complete metric spaces are also of interest.

DEFINITION 7.1.4. Let (X, d) be a metric space. A sequence (x_n) is a *Cauchy sequence* if for any $\varepsilon > 0$ there exists an index N such that $d(x_n, x_m) < \varepsilon$ for all $m, n \geq N$. If every Cauchy sequence in X has a limit in X , then (X, d) is called a *complete metric space*.

The completeness of metric spaces depends on the distance.

EXAMPLE 7.1.5. The metric space $(-\pi/2, \pi/2)$ with the standard distance $d(x, y) = |x - y|$ is not complete. In contrast $(-\pi/2, \pi/2)$ with the metric $d_*(x, y) = |\tan x - \tan y|$ is complete. The endpoints in this metric are no longer detected as missing, since the metric stretches distances near the endpoints.

LEMMA 7.2. *Given a metric space (X, d) . Then (X, d') is a metric space where $d'(x, y) = d(x, y)/(1 + d(x, y))$.*

PROOF. Mads: Please prove this statement. \square

7.1.1. Closed, open sets and complete metric spaces. Definitions and properties of open and closed sets, sequences and other notions for normed spaces have natural counterparts in the setting of metric spaces.

DEFINITION 7.1.6. (1) A set $U \subset X$ is a *neighborhood* of $x \in X$ if $B_r(x) \subset U$ for some $r > 0$.

(2) A set $O \subset X$ is *open* if every $x \in O$ has a neighborhood U contained in O .

(3) A set $C \subset X$ is *closed* if its complement $C^c = X \setminus C$ is open.

Note that the definition of open sets depends on the norm. In other words, open sets with respect to one norm need not be open with respect to another norm.

LEMMA 7.3. *Let (X, d) be a metric space. Then $B_r(x)$ is open and $\overline{B_r}(x)$ is closed for $x \in X$ and $r > 0$.*

PROOF. The proof goes along the same lines as in the case of normed spaces. Suppose that $y \in B_r(x)$ and choose ε as $\varepsilon = r - d(x, y) > 0$. The triangle inequality yields that $B_\varepsilon(y) \subset B_r(x)$, i.e. $B_r(x)$ is open.

We show that $X \setminus \overline{B_r}(x)$ is open. For $y \in X \setminus \overline{B_r}(x)$ we set $\varepsilon = d(x, y) - r > 0$ and once more by the triangle inequality we deduce that $B_\varepsilon(y) \subset X \setminus \overline{B_r}(x)$. Hence $X \setminus \overline{B_r}(x)$ is open and $\overline{B_r}(x)$ is closed. \square

DEFINITION 7.1.7. For a subset A of (X, d) we introduce some notions.

(1) The *closure* of a subset A of X , denoted by \overline{A} , is the intersection of all closed sets containing A .

(2) The *interior* of a subset of A of X , denoted by $\text{int}A$, is the union of all open subsets of X contained in A .

(3) The *boundary* of a subset A of X , denoted by $\text{bd}A$, is the set $\overline{A} \setminus \text{int}A$.

We continue with some definitions

DEFINITION 7.1.8. Let A be a subset of (X, d) .

(1) A point $x \in A$ is *isolated* in A if there exists a neighborhood U of x such that $U \cap A = \{x\}$.

- (2) A point $x \in \mathbb{R}$ is said to be an *accumulation point* of A if every neighborhood of x contains points in $A \setminus \{x\}$.

DEFINITION 7.1.9. A subset A of (X, d) is said to be *dense* in \mathbb{R} if its closure is equal to X , i.e. $\overline{A} = X$. If the dense subset A is countable, then X is called *separable*.

In other words, a subset A of a metric space X is dense in X if for each $x \in X$ and each $\varepsilon > 0$ there exists a vector $y \in A$ such that

$$d(x, y) < \varepsilon.$$

The next results have been proved in the section on real numbers and these are also true for metric spaces. The proofs of these results are along the same lines as the ones for the real line.

LEMMA 7.4. Let $\{O_j : j \in J\}$ be a family of open sets of (X, d) .

- (1) $\bigcap_{j=1}^n O_j$ is an open set for any $n \in \mathbb{N}$.
- (2) $\bigcup_{j \in J} O_j$ is open for a general index set J .

Note that open and closed subset of a normed space also applies to subspaces, since these are sets with some extra properties. For the most part we are going to discuss closed sets of a metric space.

LEMMA 7.5. Suppose A is a subset of (X, d) .

- (1) $\overline{A} = (\text{Int}(A^c))^c$ and $\text{int}(\overline{A}) = \overline{(\text{Int} A)}$
- (2) $\text{bd} A = \text{bd}(A^c) = \overline{A} \cap \overline{A^c}$
- (3) $\overline{A} = A \cup \text{br} A = \text{int} A \cup \text{bd} A$

LEMMA 7.6. Suppose A is a subset of $(X, \|\cdot\|)$.

- (1) $\overline{A} = \{x \in X : \text{every neighborhood of } x \text{ intersects } A\}$
- (2) $\text{int}(A) = \{x \in X : \text{some neighborhood of } x \text{ is contained in } A\}$
- (3) $\text{bd}(A) = \{x \in X : \text{every neighborhood of } x \text{ intersects } A \text{ and its complement}\}$

LEMMA 7.7. A point x in a metric space $(X, \|\cdot\|)$ is an accumulation point of A if and only if every neighborhood of x contains infinitely many points of A .

LEMMA 7.8. Let M be a subset of (X, d) . Then (M, d) with d as restriction to M of the distance on X is complete if and only if closed.

We collect all notions of continuity required in this course.

DEFINITION 7.1.10 (Different types of continuity). Let (X, d_X) and (Y, d_Y) be two metric spaces, let $A \subset X$ and let $f: A \rightarrow Y$ be a function.

- (1) We say that f is *continuous* at a point $a \in A$ if for all $\varepsilon > 0$ there is $\delta > 0$ such that for all $x \in A$ with $d_X(x, a) < \delta$ we have $d_Y(f(x), f(a)) < \varepsilon$.
- (2) We say that f is *continuous* on A if it is continuous at each point of A .
- (3) We say that f is *uniformly continuous* on A if for all $\varepsilon > 0$ there is $\delta > 0$ such that for all $x, y \in A$ with $d_X(x, y) < \delta$ we have $d_Y(f(x), f(y)) < \varepsilon$.
- (4) We say that f is Lipschitz (with Lipschitz constant $L \in \mathbb{R}$) if

$$d_Y(f(x), f(x')) \leq L d_X(x, x') \quad \text{for all } x, x' \in A.$$

LEMMA 7.9. If $f: A \rightarrow Y$ is a Lipschitz function, where $A \subset X$ and X, Y are metric spaces, then f is continuous at every point $a \in A$. Moreover, f is uniformly continuous.

Here is a useful criterion for continuity of a function.

PROPOSITION 7.1.11. *Let $f: A \rightarrow Y$ be a function, where $A \subset X$ and X, Y are metric spaces. Let $a \in A$. Then following two statements are equivalent.*

(i) *f is continuous at a .*

(ii) *For every sequence $(x_n) \subset A$, if $x_n \rightarrow a$ then $f(x_n) \rightarrow f(a)$.*

The proof is the same as the one given in the setting of normed spaces.

Banach's fixed point theorem holds for general metric spaces with the same proof as for normed spaces.

THEOREM 7.10 (Banach Fixed Point). *Let M be a closed subset of a metric space X . Any contraction f on M has a unique fixed point \tilde{x} and the fixed point is the limit of every sequence generated from an arbitrary nonzero point $x_0 \in M$ by iteration $(x_n)_n$, where $x_{n+1} = f(x_n)$ for $n \geq 1$.*

Sets and functions

A.1. Sets and functions

In order to formalize our intuition about collections of objects we use the framework of set theory. The relation between sets and their elements will be described by functions.

DEFINITION A.1.1. A *set* is a collection of distinct objects, its *elements*. If an object x is an element of a set X , we denote it by $x \in X$. If x is not an element of X , then we write $x \notin X$.

A set is uniquely determined by its elements. Suppose X and Y are sets. Then they are identical, $X = Y$, if they have the same elements. More formalized, $X = Y$ if and only if for all $x \in X$ we have $x \in Y$, and for all $y \in Y$ we have $y \in X$.

The *empty set* is the set with no elements, denoted by \emptyset .

DEFINITION A.1.2. Suppose X and Y are sets. Then Y is a subset of X , denoted by $Y \subseteq X$, if for all $y \in Y$ we have $y \in X$.

If $Y \subseteq X$, one says that Y is contained in X . If $Y \subseteq X$ and $X \neq Y$, then Y is a proper subset of X and we use the notation $Y \subset X$.

The most direct way to prove that two sets E and F are equal is to show that

$$x \in E \iff x \in F$$

for any element x .

(Another way is to prove a double inclusion: if $x \in E$ then $x \in F$, establishing that $E \subseteq F$ and if $x \in F$, then $x \in E$, establishing that $F \subseteq E$. You may, of course, do it this way.)

Here are a few constructions of sets.

DEFINITION A.1.3. Let X and Y be sets.

- The *union* of X and Y , denoted by $X \cup Y$, is defined by

$$X \cup Y = \{z \mid z \in X \text{ or } z \in Y\}.$$

- The *intersection* of X and Y , denoted by $X \cap Y$, is defined by

$$X \cap Y = \{z \mid z \in X \text{ and } z \in Y\}.$$

- The *difference set* of X from Y , denoted by $X \setminus Y$, is defined by

$$X \setminus Y = \{z \in X \mid z \notin Y\}.$$

If all sets are contained in one set X , then the difference set $X \setminus Y$ is called the *complement* of Y .

- The *Cartesian product* of X and Y , denoted by $X \times Y$, is the set

$$X \times Y = \{(x, y) | x \in X, y \in Y\},$$

i.e. the set of all ordered pairs (x, y) , with $x \in X$ and $y \in Y$.

Here are some basic properties of sets.

LEMMA A.1. *Let A, B and C be sets.*

- (1) $A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$ and $A \cup (B \cap C) = (A \cup B) \cap (A \cup C)$
(*distributative law*)
- (2) $(A \cup B)^c = A^c \cap B^c$ and $(A \cap B)^c = A^c \cup B^c$ (*De Morgan's laws*)
- (3) $A \setminus (B \cup C) = (A \setminus B) \cap (A \setminus C)$ and $X \setminus (B \cap C) = (A \setminus B) \cup (A \setminus C)$
- (4) $(A^c)^c = A$.

PROOF. (ii) Let us prove one of de Morgan's relations. Let us use the most direct approach. Keep in mind that $x \in E^c \iff x \notin E$. We then have:

$$\begin{aligned} x \in (A \cup B)^c &\iff x \notin A \cup B \iff x \notin A \text{ and } x \notin B \\ &\iff x \in A^c \text{ and } x \in B^c \iff x \in A^c \cap B^c. \end{aligned}$$

This proves the identity.

(iv)

$$x \in (A^c)^c \iff x \notin A^c \iff x \in A.$$

□

Let X and Y be sets. A function with *domain* X and *codomain* Y , denoted by $f : X \rightarrow Y$, is a relation between the elements of X and Y satisfying the properties: for all $x \in X$, there is a unique $y \in Y$ such that $(x, y) \in f$, we denote it by: $f(x) = y$.

By definition, for each $x \in X$ there is exactly one $y \in Y$ such that $f(x) = y$. We say that y the *image* of x under f . The *graph* $G(f)$ of a function f is the subset of $X \times Y$ defined by

$$G(f) = \{(x, f(x)) | x \in X\}.$$

The *range* of a function $f : X \rightarrow Y$, denoted by $\text{range}(f)$, or $f(A)$, is the set of all $y \in Y$ that are the image of some $x \in X$:

$$\text{range}(f) = \{y \in Y | \text{there exists } x \in X \text{ such that } f(x) = y\}.$$

The *pre-image* of $y \in Y$ is the subset of all $x \in X$ that have y as their image. This subset is often denoted by $f^{?1}(y)$:

$$f^{?1}(y) = \{x \in X | f(x) = y\}.$$

Note that $f^{?1}(y) = \emptyset$ if and only if $y \in Y \setminus \text{range}(f)$.

The following notions are central for the theory of functions.

DEFINITION A.1.4. Let $f : X \rightarrow Y$ be a function.

- (1) Then we call f *injective* or *one-to-one* if $f(x_1) = f(x_2)$ implies $x_1 = x_2$, i.e. no two elements of the domain have the same image.
- (2) Then we call f *surjective* or *onto* if $\text{range}(f) = Y$, i.e. each $y \in Y$ is the image of at least one $x \in X$.
- (3) Then we call f *bijective* if f is both injective and surjective.

Let $f : X \rightarrow Y$ and $g : Y \rightarrow Z$ be two functions so that the codomain of f coincides with the domain of g . Then we define the *composition*, denoted by $g \circ f$, as the function $g \circ f : X \rightarrow Z$, defined by $x \mapsto g(f(x))$.

For every set X , we define the *identity map*, denoted by id_X or id for short: $\text{id}_X : X \rightarrow X$ is defined by $\text{id}_X(x) = x$ for all $x \in X$. The identity map is a bijection.

If f is a bijection, then it is invertible. Hence, the inverse relation is also a function, denoted by f^{-1} . It is the unique bijection $Y \rightarrow X$ such that $f^{-1} \circ f = \text{id}_X$ and $f \circ f^{-1} = \text{id}_Y$.

LEMMA A.2. *Let $f : X \rightarrow Y$ and $g : Y \rightarrow Z$ be bijections. Then $g \circ f$ is also a bijection and $(g \circ f)^{-1} = f^{-1} \circ g^{-1}$.*

LEMMA A.3. *Let $f : X \rightarrow Y$ be a function and let $C, D \subset Y$. Prove that*

$$f^{-1}(C \cup D) = f^{-1}(C) \cup f^{-1}(D)$$

where

$$f^{-1}(E) := \{x \in X : f(x) \in E\}$$

is the pre-image of a set $E \subset Y$.

PROOF.

$$\begin{aligned} x \in f^{-1}(C \cup D) &\iff f(x) \in C \cup D \iff f(x) \in C \text{ or } f(x) \in D \\ &\iff x \in f^{-1}(C) \text{ or } x \in f^{-1}(D) \iff x \in f^{-1}(C) \cup f^{-1}(D). \end{aligned}$$

□

Bibliography