# Part I:

# Linear Algebra

2024 / 25

( pages with a grey background were not treated this year and are not part of the exam; but the videos still exist if you are interested )

# Groups and fields

# Definition of a group

__Def__  A set $G$ of elements with an operation $+: G \times G \to G$ is called a <mark>__group__</mark> if the following properties hold:

(G1) Associativity: $\forall a, b, c \in G: \quad (a+b)+c = a + (b+c)$

(G2) Identity element: $\exists e \in G \; \forall g \in G: \quad e+g = g+e = g$

(G3) Inverse elements: $\forall a \in G \; \exists b \in G: \quad a+b = b+a = e$

The group is called a commutative group (Abelian group) if we have additionally that

(G4) $\quad \forall a, b \in G: \quad a+b = b+a$

# Examples of groups

- $(\mathbb{R}^n, +)$, $(\mathbb{R}^+, \cdot)$ are groups

- $(\mathbb{R}^-, \cdot)$ is <u>not</u> a group.

- $S_n := \{ \pi : \{1, \ldots, n\} \to \{1, \ldots, n\} \mid \pi \text{ is bijective} \}$

  $\circ : S_n \times S_n \longrightarrow S_n$ , $\quad \pi_1 \circ \pi_2 (i) = \pi_1 (\pi_2 (i))$

  $(S_n, \circ)$ is a group.

# Definition of a field

<u>Def</u>   A set $F$ with two operations $+, \cdot : F \times F \to F$ is
  called a ==field== if the following properties hold:

(F1)   $(F, +)$ is a commutative group, with identity element $0$.

(F2)   $(F \setminus \{0\}, \cdot)$ is a commutative group with id. el. $1$

(F3)   Distributivity: $\forall a, b, c \in F: \quad a \cdot (b+c) = a \cdot b + a \cdot c$

# Examples of fields

- $(\mathbb{R}, +, \cdot)$
- $(\mathbb{C}, +, \cdot)$

- $n \in \mathbb{Z}$, Consider $\mathbb{Z}_n := \{0, 1, \ldots, n-1\}$

  $a +_n b := (a+b) \bmod n$

  $a \cdot_n b := (a \cdot b) \bmod n$

  Then $(\mathbb{Z}_n, +_n, \cdot_n)$ is a field if and only if $n$ is prime.

Complex numbers

# Motivation

In machine learning, our data is often represented by real numbers: $\mathbb{R}$

In linear algebra, however, it often helps if we extend the real numbers to complex numbers: $\mathbb{C}$

We are not going to use a lot of maths of complex numbers, but at least used them to factorize polynomials.

Here are the very basics:

# Quadratic equations over $\mathbb{R}$

Quadratic equation: for given parameters $a, b, c \in \mathbb{R}$, want to find $x \in \mathbb{R}$ that satisfies the quadratic equation:

$$ax^2 + bx + c = 0$$

In school you learned the formula:

$$x_{1,2} = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

It has solutions in $\mathbb{R}$ if $b^2 - 4ac \geq 0$, otherwise it doesn't. Annoying!

## Imagine ...

Imagine that $\sqrt{-1}$ exists, give it a name: $i := \sqrt{-1}$

($i$ for 'imaginary number).

**Def**  A complex number is a number of the form $a + bi$ where $a, b \in \mathbb{R}$. We call $a$ the real part and $b$ the imaginary part of the number. Write $\mathbb{C}$ for the space of all such numbers: $\mathbb{C} = \{a + ib \mid a, b \in \mathbb{R}\}$.

**Observe:** $\mathbb{C}$ is a field.

# Quadratic equations over $\mathbb{C}$

Consider the quadratic equation again: $ax^2 + bx + c = 0$.

Observe that it now always has a solution in $\mathbb{C}$:

$$x_{1/2} = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

Case 1: $b^2 - 4ac \geq 0$. Then $x_{1,2} \in \mathbb{R}$ "as usual".

Case 2: $b^2 - 4ac < 0$. Then can write

$$\underbrace{\sqrt{b^2 - 4ac}}_{<0} = \sqrt{-1\underbrace{(4ac - b^2)}_{>0}} = \underbrace{\sqrt{-1}}_{i} \cdot \underbrace{\sqrt{4ac - b^2}}_{>0}$$

$$= i \cdot r \quad \text{where } r = \sqrt{4ac - b^2} \in \mathbb{R}.$$

# Fundamental theorem of algebra

## Theorem :

Consider a polynomial $a_0 + a_1 x + a_2 x^2 + \ldots + a_n x^n$.

(where the $a_i$ can be real or complex number).

If $a_n \neq 0$, this is a polynomial of degree $n$.

Any such polynomial has exactly roots $r_1, \ldots, r_n \in \mathbb{C}$

(not necessarily distinct) such that

$$a_0 + a_1 x + \ldots + a_n x^n = a_n (x - r_1)(x - r_2) \ldots (x - r_n)$$

# Outlook

Dealing with complex numbers is a rich field in mathematics, but we will not touch it ...

Vector space

# Definition of a vector space

<u>Def</u> Let $F$ be a field with id. elements $0$ and $1$.
A ==vector space== over the field $F$ is a set $V$ with a
mapping $+: V \times V \to V$ ("vector addition") and a mapping
$\cdot : F \times V \to V$ ("scalar multiplication") such that:

(V1) $(V, +)$ is a commutative group.

(V2) Multiplicative identity: $\forall v \in V : 1 \cdot v = v$

(V3) Distributive properties: $\forall a, b \in F \quad \forall u, v \in V$

$$a \cdot (u + v) = a \cdot u + a \cdot v$$
$$(a + b) u = a \cdot u + b \cdot u$$

Elements of $V$ are called vectors, elements of $F$ are
called scalars.

# Examples of vector spaces

- $\mathbb{R}^n$ with the standard operations.

- Function spaces:

  - $\mathbb{R}^X := \{ f : X \to \mathbb{R} \}$ the space of all real valued fcts

    on a set $X$. Define:

    $+ : \mathbb{R}^X \times \mathbb{R}^X \to \mathbb{R}^X$, $(f + g)(x) := f(x) + g(x)$

    $\cdot : \mathbb{R} \times \mathbb{R}^X \to \mathbb{R}^X$, $(\lambda \cdot f)(x) := \lambda \cdot (f(x))$

    Then $(\mathbb{R}^X, +, \cdot)$ is a real vector space.

  - $\mathcal{C}(X) := \{ f : X \to \mathbb{R} \mid f \text{ is continuous} \}$

  - $\mathcal{C}^r([a,b]) = \{ f : [a,b] \to \mathbb{R} \mid f \text{ is } r \text{ times cont.} \\ \text{differentiable} \}$

# Subspaces

<u>Def</u>   Let $V$ be a vector space, $U \subset V$ non-empty set.
We call $U$ a <mark>subspace</mark> of $V$ if it is closed under linear
combinations: $\forall \lambda, \mu \in F \; \forall u, v \in U : \lambda u + \mu \cdot v \in U$

<u>Examples</u>:  . $\ell(x)$ is a subspace of $\mathbb{R}^x$.

  . The set $S$ of symmetric matrices of size $n \times n$
    is a subspace of $\mathbb{R}^{n \times n}$.

# Basis and dimension

# Linear combinations

<u>Def</u>  $V$ vector space over $F$, $u_1, \ldots, u_n \in V$, $\lambda_1, \ldots, \lambda_n \in F$.

Then $\sum_{i=1}^{n} \lambda_i u_i$ is called a ==<u>linear combination</u>==. The set

of all lin. comb. of $u_1, \ldots, u_n$ is called the ==<u>span</u>==

( linear hull) of $u_1, \ldots, u_n$.


Notation:

$$\text{span}(u_1, \ldots, u_n) := \left\{ \sum_{i=1}^{n} \lambda_i u_i \;\middle|\; \lambda_i \in F \right\}.$$


The set $U := \{ u_1, \ldots, u_n \}$ is the ==<u>generator of</u>== span $(U)$.

# Linear independence

<u>Def</u>  A set of vectors $v_1, \ldots, v_n$ is called <mark><u>linearly independent</u></mark>
if the following holds:

$$\sum_{i=1}^{n} \lambda_i v_i = 0 \implies \lambda_1 = \ldots = \lambda_n = 0.$$

<u>Examples</u>:
- The vectors $\begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 2 \\ 1 \\ 0 \end{pmatrix} \begin{pmatrix} 3 \\ 4 \\ 1 \end{pmatrix} \in \mathbb{R}^3$ are lin. indep.

- The functions $\sin(x)$ and $\cos(x) \in \mathbb{R}^{\mathbb{R}}$ are lin. ind.

- Any set of $d+1$ vectors in $\mathbb{R}^d$ is lin. dependent.

- Any set that contains the $0$-vector is not lin. indep.

# Basis of a vector space

Def  A subset $B$ of a vectorspace $V$ is called a
   (Hamel) **basis** if

 (B1)  Span $(B) = V$

 (B2)  $B$ is lin. independent.

Example
 - The canonical basis of $\mathbb{R}^3$ is $\begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$

 - Another basis of $\mathbb{R}^3$ is given by
   $\begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$   or   $\begin{pmatrix} 0.5 \\ 0.8 \\ 0.4 \end{pmatrix}, \begin{pmatrix} 1.8 \\ 0.3 \\ 0.3 \end{pmatrix}, \begin{pmatrix} -2.2 \\ -1.3 \\ 3.5 \end{pmatrix}$

# Reducing a set to a basis

**Proposition:** If $U = \{u_1, \ldots, u_n\}$ spans a VS $V$, then the set $U$ can be reduced to a basis of $V$.

ML keyword: dictionary learning

# Proof sketch

- If $U$ is already lin. independent: done.

- If $U$ is lin. dependent:

$$\Rightarrow \exists \lambda_1, \ldots, \lambda_n \text{ not all } 0 \text{ such that } \sum \lambda_i u_i = 0.$$

Pick some $k$ with $\lambda_k \neq 0$. Then:

$$u_k = \sum_{i \neq k} \frac{\lambda_i}{\lambda_k} u_k$$

Denote $\tilde{U} := U \setminus \{u_k\}$. It is clear that $\text{span}(\tilde{U}) = \text{span}(U)$

If $\tilde{U}$ is not lin. independent, we repeat this process until the remaining set is lin. independent.

Note that this will eventually be the case, at least if the set only consists of one vector.

# Finite-dim vector space

**Def** A VS is called ==finite-dim== if it has a finite basis.

We need to do a bit more work to define the dim of a VS.

Can you come up with an example of an infinite space?

# Extending a set to a basis

**Prop** Let $U = \{u_1, ..., u_n\}$ be a set of lin. ind. vectors and let $\subset V$

$V$ be a finite-dim VS. Then $U$ can be extended to a basis of $V$.

**Proof** (Sketch) Let $w_1, ..., w_m$ be a basis of $V$. Consider the set

$\{u_1, ..., u_n, w_1, ..., w_m\}$. Remove vectors "from the end" until the remaining vectors are lin. independent.

- remaining set spans $V$
- remaining set is linearly ind. by construction
- remaining set contains $U$.

# Two finite bases have the same length

Corollary Let $V$ be a finite-dim VS. Then ==any two bases== ==of $V$ have the same length.==

Proof sketch : Let $\{b_1,..., b_n\}$ and $\{c_1,..., c_m\}$ be two bases, $n \leq m$. Consider the set $\{b_1,..., b_n, c_1\}$. By construction is lin. dependent. So by the same procedure as before, we can find a vector $b_i$ such that $\{b_1,.., b_{i-1}, b_{i+1}..., b_n, c_1\}$ is ind. Keep on applying this procedure: add vectors from $C$, remove vectors from $B$. At the end, this results in a set of $n$ vectors $\{c_{i_1},..., c_{i_n}\} \subset \{c_1,..., c_m\}$. By construction, they are lin. ind and span $V$. If now we had $m > n$, then the set $\{c_1,..., c_m\}$ cannot be lin. ind. any more.

# Dimension of a vector space

**Def** the length of a basis of a finite-dim VS is called the <mark>dimension</mark> of V.

# Linear Mappings

# Linear mapping

__Def__ Let $U, V$ VS over $F$. A mapping $f : U \to V$ is called
  __linear__ if $\forall u_1, u_2 \in U$, $\forall \lambda \in F$

$$f(u_1 + u_2) = f(u_1) + f(u_2)$$

$$f(\lambda u_1) = \lambda f(u_1)$$

The set of all linear mapping from $U \to V$ is denoted $\mathcal{L}(U, V)$.
If $U = V$, then we write $\mathcal{L}(U)$.

__Examples__ : $\bullet$ $T : \mathcal{C}[a,b] \to \mathbb{R}$, $f \mapsto \int_a^b f(x)\, dx$ (Integration)
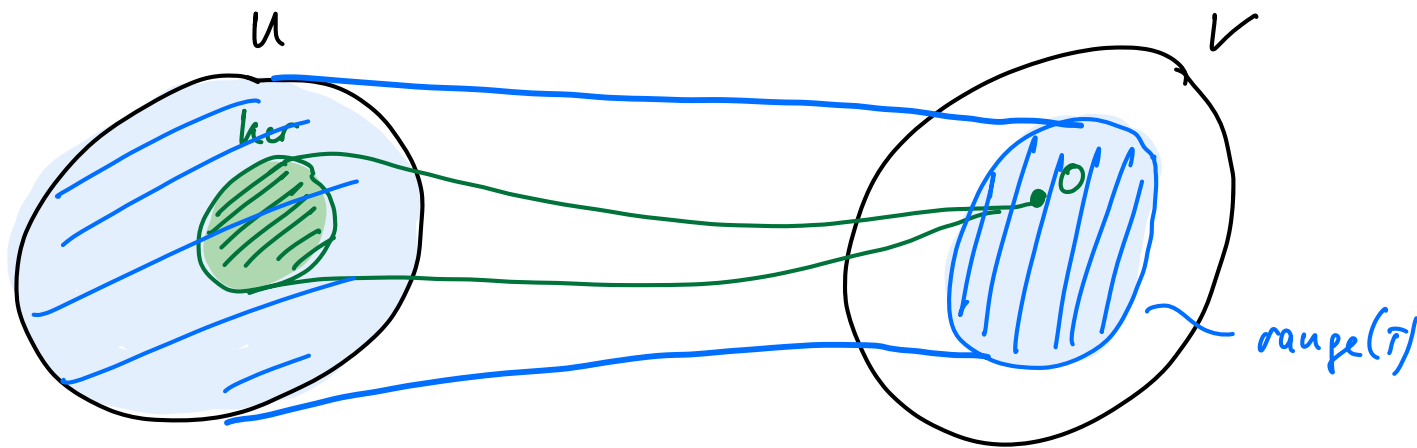
$\bullet$ $D : \mathcal{C}^\infty[a,b] \to \mathcal{C}^\infty[a,b]$, $f \mapsto f'$ (Differentiation)

**Def** $T \in \mathcal{L}(U, V)$. Then **kernel** of $T$ (**null space**) is defined as

$$\ker(T) := \text{null}(T) := \{ u \in U \mid Tu = 0 \} \subset U$$

The **range** of $T$ (**image** of $T$) is defined as

$$\text{range}(T) := \text{Im}(T) := \{ Tu \mid u \in U \} \subset V$$

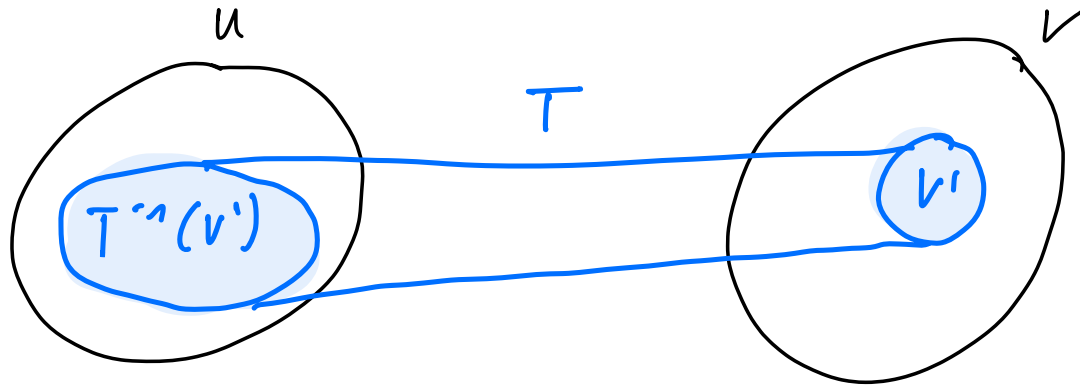# Properties of kernel and range

Proposition:

- ker $(T)$ and range $(T)$ are subspaces.

- $T$ is injective $\iff$ ker $T = \{0\}$

- $T$ is surjective $\iff$ range $T = V$

Proof: Exercise.

# Pre-image

Def : $V' \subset V$, $v'$ any set. The ==pre-image== of $V'$ is defined as

$$T^{-1}(v') := \{ u \in U \mid Tu \in V' \}.$$



Prop : If $V' \subset V$ is a subspace of $V$, then $T^{-1}(v')$ is a subspace of $U$.

Proof : exercise!

# Fundamental theorem for linear mappings

**Theorem:** Let $V$ be finite-dim, $W$ any VS, $T \in \mathcal{L}(V, W)$.
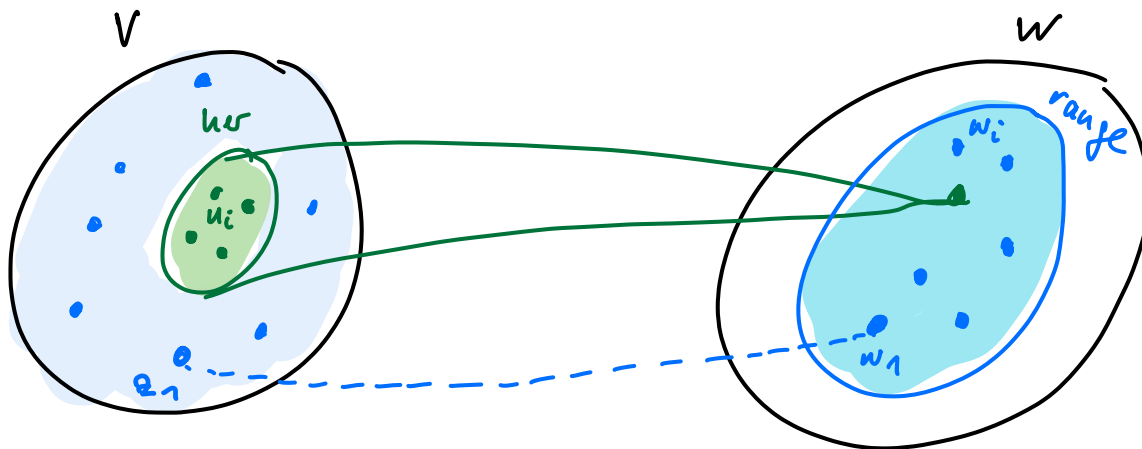
Let $u_1, \ldots, u_n$ be a basis of $\ker(T) \subset V$

Let $w_1, \ldots, w_m$ be a basis of $\text{range}(T) \subset W$.

Let $z_1 \in T^{-1}(w_1), \ldots, z_m \in T^{-1}(w_m)$. Then

Then $u_1, \ldots, u_n, z_1, \ldots, z_m \subset V$ form a basis of $V$.

In particular, $\dim(V) = \dim(\ker(T)) + \dim(\text{range}(T))$.

# Proof of the Theorem

Step 1: $V \subset \text{span}\{u_1, \ldots, u_n, z_1, \ldots, z_m\}$

Let $v \in V$, consider $Tv \in \text{range}(T)$.

$\Rightarrow \exists \lambda_1, \ldots, \lambda_m$ s.t.

$$Tv = \lambda_1 w_1 + \cdots + \lambda_m w_m$$

$$= \lambda_1 T(z_1) + \cdots + \lambda_m T(z_m)$$

$$= T(\lambda_1 z_1 + \cdots + \lambda_m z_m)$$

$$\Rightarrow \quad Tv - T(\lambda_1 z_1 + \cdots + \lambda_m z_m) = 0$$

$$= T(v - (\lambda_1 z_1 + \cdots + \lambda_m z_m))$$

$$\underbrace{\qquad\qquad\qquad\qquad}_{\color{red}\in \ker(T)}$$

$$\Rightarrow \quad \exists \; \mu_1, \ldots, \mu_n \quad \text{s.t.} \quad v - (\lambda_1 z_1 + \cdots + \lambda_m z_m) = \mu_1 u_1 + \cdots + \mu_n u_n$$

$$\Rightarrow \quad v = \lambda_1 z_1 + \cdots + \lambda_m z_m + \mu_1 u_1 + \cdots + \mu_n u_n$$

Assume that $\mu_1 u_1 + \cdots + \mu_n u_n + \lambda_1 z_1 + \cdots + \lambda_m z_m = 0$ $\circledast$

$$\lambda_1 w_1 + \cdots + \lambda_m w_m = \lambda_1 T(z_1) + \cdots + \lambda_m T(z_m)$$

we add something
that is 0
(because in kernel)
$$= \lambda_1 T(z_1) + \cdots + \lambda_m T(z_m) + \underbrace{\mu_1 T(u_1) + \cdots + \mu_n T(u_n)}_{= 0}$$

exploit
linearity
$$= T(\underbrace{\lambda_1 z_1 + \cdots + \lambda_m z_m + \mu_1 u_1 + \cdots + \mu_n u_n}_{= 0 \text{ by } \circledast}) = 0$$

$\Rightarrow \lambda_1 w_1 + \cdots + \lambda_m w_m = 0$  $\Rightarrow \lambda_1 = \cdots = \lambda_m = 0$
$w_1, \ldots, w_m$
basis

$\Rightarrow \mu_1 u_1 + \cdots + \mu_n u_n = 0$  by $\circledast$

$\Rightarrow \mu_1 = \cdots = \mu_n = 0$  because $u_1, \ldots, u_n$ basis.

# Injective <=> Surjective <=> bijective

**Prop**  $T \in \mathcal{L}(V, V)$, $V$ finite-dim. Then the following three statements are equivalent:

(i)  $T$ injective.

(ii)  $T$ surjective.

(iii)  $T$ bijective.

**Proof** Direct consequence of theorem. (can you see why?. Exercise!)

⚠️ Does not hold in $\infty$-dim spaces!

# Matrices and linear maps

# Matrices

A **matrix** is the following object:

$$A = \underbrace{\left\{ \overset{\overbrace{n \; col.}}{\begin{pmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & & \\ a_{m1} & \cdots & a_{mn} \end{pmatrix}} \right.}_{m \; rows} = \left( a_{ij} \right)_{\substack{i = 1 \cdots m \\ j = 1 \cdots n}}$$

# Matrices represent linear maps

Consider $T \in \mathcal{L}(V, W)$, $V, W$ finite-dim,

let $v_1, \ldots, v_n$ be a basis of $V$

$\qquad w_1, \ldots, w_m \qquad$ basis of $W$

- If we know the results of $T$ applied to the basis vectors $v_i$, then we can express $T(v)$ for arbitrary $v$:

$v = \lambda_1 v_1 + \ldots + \lambda_n v_n \quad$ arbitrary vector

$T(v) = T(\lambda_1 v_1 + \ldots + \lambda_n v_n)$

$\qquad = \lambda_1 T(v_1) + \ldots + \lambda_n T(v_n)$

- For basis vector $v_j$, we can express the image $T(v_j)$ in basis $w_1, \ldots, w_m$:

  There exist coefficients $a_{1j}, \ldots, a_{mj}$ s.t.

$$T(v_j) = a_{1j} w_1 + \cdots + a_{mj} w_m$$

- We now stack these coefficients in a matrix called $M(T)$:

m rows, one for each basis vector of $W$

$$M := \begin{pmatrix} a_{11} & \cdots & a_{1j} & \cdots & a_{1n} \\ \vdots & & \vdots & & \vdots \\ a_{m1} & \cdots & a_{mj} & \cdots & a_{mn} \end{pmatrix}$$

col $j$

$\underbrace{\hspace{4cm}}$ n cols, one for each basis vectors of $V$

$=$ matrix of mapping $T$ with respect to the bases $v_1, \ldots, v_n$ of $V$ $w_1 \ldots w_m$ of $W$.

- The result of $T(v)$ can now be expressed by a matrix-vector multiplication:

$$T(v) = \sum_{j=1}^{n} \lambda_j \, T(v_j)$$

$$= \sum_{j=1}^{n} \lambda_j \sum_{i=1}^{m} a_{ij} \, w_i$$

$$= \sum_{i=1}^{m} \underbrace{\left( \sum_{j=1}^{n} a_{ij} \lambda_j \right)}_{\left( M \cdot \begin{pmatrix} \lambda_1 \\ \vdots \\ \lambda_n \end{pmatrix} \right)_i} w_i$$

$i$-th entry of product of matrix $M$ with vector $\lambda$

# Notation for matrices of linear maps

<u>Notation</u>  Let $T : V \to W$ be linear, let $\mathcal{B}$ a basis of $V$, $\mathcal{C}$ basis of $W$. We denote by

$$M(T, \mathcal{B}, \mathcal{C})$$

the matrix corresponding to $T$ wrt bases $\mathcal{B}$ and $\mathcal{C}$.

# Properties of matrices

$V, W$ vector spaces, consider the basis fixed. Let $S, T \in \mathcal{L}(V, W)$.

Then:

- $M(S + T) = M(S) + M(T)$

- $M(\lambda S) = \lambda M(S)$

- $T: U \to V, \; S: V \to W$ linear, then
$$M(S \circ T) = M(S) \cdot M(T)$$

# Matrix transpose

**Def** Given a matrix $A = (a_{ij})_{ij} \in F^{m \times n}$, the ==transpose matrix== is given as

$$(A^t)_{kj} = A_{jk}$$

Notation: $A^t$, $A'$

If $F = \mathbb{C}$, then the ==conjugate transpose== matrix is defined as

$$(A^*)_{ij} = \overline{a_{ji}}$$

# Sum of vector spaces

**Def**  Assume that we have $U_1, U_2$ subspaces of $V$.

The ==sum== of the two spaces is defined as

$$U_1 + U_2 := \{ u_1 + u_2 \mid u_1 \in U_1, u_2 \in U_2 \}$$

The sum is called a ==direct sum==, if each element in the sum can be written in exactly one way.

Notation:  $U_1 \oplus U_2$

# Complement of a subspace

<u>Prop</u>  Suppose $V$ is finite-dim, and $U \subset V$ is a subspace. Then there exists a subspace $W \subset V$ such that $U \oplus W = V$.

<u>Proof</u>  (sketch)  Let the set $\{u_1, \dots, u_n\}$ basis of of $U$. Extend it to a basis of $V$, say the resulting set is

$$\underbrace{\{u_1, \dots, u_n}_{\leadsto U}, \underbrace{v_1, \dots, v_m\}}_{\leadsto W} . \quad \text{Define}$$

$$W = \text{span}\{v_1, \dots, v_m\}.$$

# Invertible maps and matrices

# Inverse of a linear map

<u>Def</u> $T \in \mathcal{L}(V, W)$ is called <mark>invertible</mark> if there exists
a linear map $S \in \mathcal{L}(W, V)$ such that

$$S \circ T = Id_V \quad \text{and} \quad T \circ S = Id_W$$

The map $S$ is called the <mark>inverse</mark> of $T$, denoted by $T^{-1}$.


<u>Rem</u> • Inverse maps are unique.

• not every lin. map is invertible

# Characterizing invertability

**Prop**   A linear map is invertible iff it is injective and surjective.

**Proof**   "⟹" Invertible ⟹ injective:

suppose $T(u) = T(v)$. Then $\underline{u} = T^{-1}(T(u))$

$= T^{-1}(Tv) = \underline{v}$     ⟹ injective

Invertible ⟹ surjective:

Consider $w \in W$. Then

$w = T(T^{-1}(w))$  ⟹  $w \in$ range of $T$

⟹ surjective.

"$\Leftarrow$" inj & surj $\Rightarrow$ invertible.

- Let $w \in W$. There exists unique $v \in V$ s.th. $T(v) = w$. (because $T$ surj.)

  Define the mapping: $S(w) = v$. Clearly here $\underline{T \circ S = Id}$.

- Let $v \in V$. Then $T\left( \underline{(S \circ T)(v)} \right) =$

  $= \underbrace{(T \circ S)}_{Id}(Tv) = Id \circ Tv = T\underline{v}$. Because $T$ is injective, then

  $(S \circ T)v = v \quad \Rightarrow \quad \underline{S \circ T = Id}$

- Linearity: $T(Sw_1 + Sw_2) = \underbrace{TS}_{Id}w_1 + \underbrace{TS}_{Id}w_2 = w_1 + w_2$

inj. $\curvearrowright \Rightarrow \quad Sw_1 + Sw_2$ is the unique el. in $V$ that $T$ maps to $w_1 + w_2$

  $\Rightarrow$ By def. of $S$, we thus get $S(w_1 + w_2) = Sw_1 + Sw_2$

  Similarly for scalar multiplication.

# Inverse matrix

<u>Def</u>    A square <mark>matrix</mark> $A \in F^{n \times n}$ is <mark>invertible</mark> if there exists

a square matrix $B \in F^{n \times n}$ such that

$$A \cdot B = B \cdot A = Id \quad \leftarrow \begin{pmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & 1 \end{pmatrix}.$$

The matrix $B$ is called the inverse matrix, and is denoted

by $A^{-1}$.

# Inverting maps $\hat{=}$ inverting matrices

**Prop**  The inverse matrix represents the inverse of the corr. lin.
map, that is:   $T : V \to V$

$$M(T^{-1}) = \left(M(T)\right)^{-1}$$

matrix of ( inverse map )       inverse matrix (of the original matrix)

In particular, a matrix is invertible iff the corr. map
is invertible.

**Proof** : Exercise.

# Properties of inverse matrices

- The inverse matrix does __not__ always exist.

- $\left(A^{-1}\right)^{-1} = A$ , $\left(A \cdot B\right)^{-1} = B^{-1} \cdot A^{-1}$

- $A^t$ invertible $\iff$ $A$ invertible,

  $$\left(A^t\right)^{-1} = \left(A^{-1}\right)^t$$

- $A \in F^{n \times n}$ invertible $\iff$ rank$(A) = n$

- The set of all invertible matrices is called general linear group:

  $$GL(n, F) = \left\{ A \in F^{n \times n} \mid A \text{ invertible} \right\}$$

# Change of basis

# Representing the identity

Consider the identity mapping $J: V \to V, \ x \mapsto x$.
Assume we fix a basis of $V$ (both in source and target space), then the corr. matrix looks as follows:

$$M(J, B, B) = \begin{pmatrix} 1 & \cdots & 0 \\ 0 & & 1 \end{pmatrix}$$

Now consider $A = \{a_1, ..., a_n\}$ and $B = \{b_1, ..., b_n\}$ both bases on $V$. How does the matrix of the id. mapping

$$J: (V, A) \to (V, B) \qquad \text{look like?}$$

Because $B$ is basis, we can write each of the vectors in $A$
as lin. comb. of vectors in $B$:

$$a_1 = \boxed{t_{11} b_1 + t_{21} b_2 + \dots + t_{n1} b_n}$$

$$a_2 = \dots$$

Now we form the corr. matrix $T$

$$T = \begin{pmatrix} t_{11} & \dots & t_{1n} \\ \vdots & & \vdots \\ t_{n1} & \dots & b_{nn} \end{pmatrix}$$

This matrix represents the identity mapping:

- In the basis $\mathcal{A}$, the first basis vector $a_1$ has the representation $\begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$.

$a_1 = 1 \cdot a_1 + 0 \cdot a_2 + 0 \cdot a_3 \cdots + 0 \cdot a_n$

- $T \cdot \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} = \begin{pmatrix} t_{11} \\ t_{21} \\ \vdots \\ t_{n1} \end{pmatrix}$

This vector gives us $Ta_1$ expressed in basis $\mathcal{B}$

by def of coefficients $t_{ij}$

- $t_{11} b_1 + \cdots + t_{n1} b_n = a_1$

- $T a_1 = a_1 .$

# Change of basis is invertible

**Prop**   Let $\mathcal{A}, \mathcal{B}$ be two bases of $V$. Then the matrices $M(\mathrm{Id}, \mathcal{A}, \mathcal{B})$ and $M(\mathrm{Id}, \mathcal{B}, \mathcal{A})$ are invertible and each is the inverse of each other.

**Proof**   Exercise / skipped.
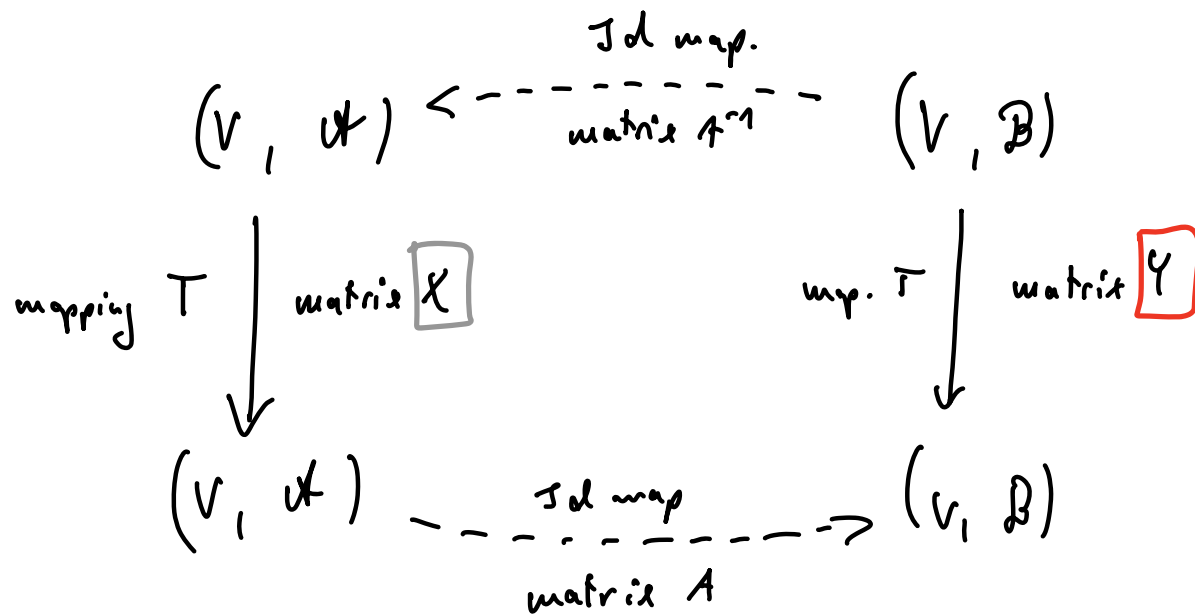
# Change of basis for an arbitrary mapping

**Prop** Let $\mathcal{A}, \mathcal{B}$ be two bases of $V$. Consider the transformation matrix

$A = M(\text{Id}, \underline{\mathcal{A}}, \underline{\mathcal{B}})$, and $A^{-1} = M(\text{Id}, \underline{\mathcal{B}}, \underline{\mathcal{A}})$.

Let $T: V \to V$ linear, and $\boxed{X} = M(T, \underline{\mathcal{A}}, \underline{\mathcal{A}})$. Then

$\boxed{Y := A \cdot X \cdot A^{-1}}$ represents $T$ in basis $\mathcal{B}$, that is

$Y = M(T, \underline{\mathcal{B}}, \underline{\mathcal{B}})$.

# Rank of a matrix

# Rank of a matrix

__Def__ $A \in F^{m \times n}$. The ==column rank== of $A$ is

$$\dim \left( \mathrm{span} \left( \text{column vectors of } A \right) \right)$$

The ==row rank== is defined accordingly.

__Prop__ For a matrix, the row and column rank always coincide. We now call it the ==rank== of the matrix.

__Prop__ $T \in \mathcal{L}(V, W)$. Then $\mathrm{rank}\,(M(T)) = \dim(\mathrm{range}\,(T))$.

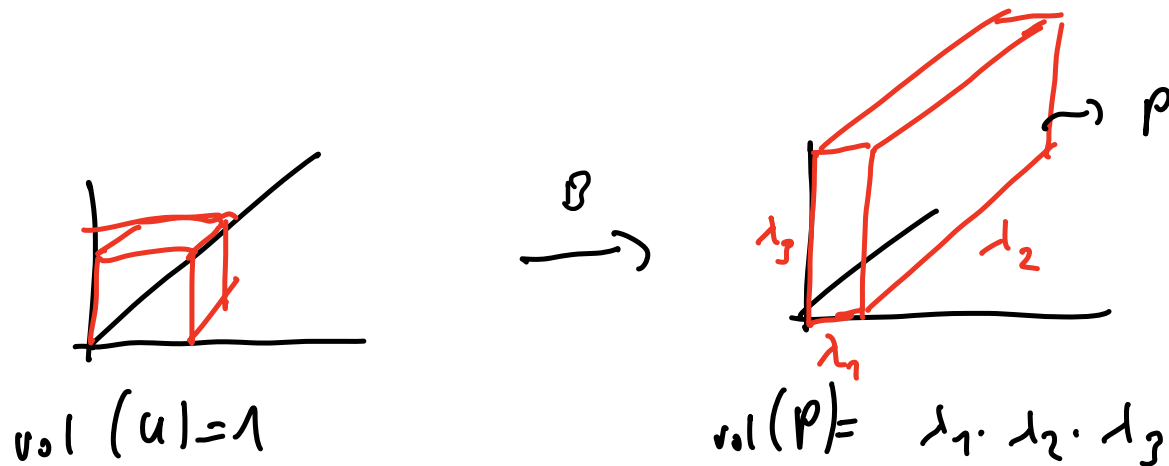__Proofs:__ skipped

# The determinant

# Motivation to study the determinant: geometry!

Consider the standard basis of $\mathbb{R}^3$: $e_1 = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$, $e_2 = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$, $e_3 = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$.

Consider a linear mapping that just streches these vectors: $T = \begin{pmatrix} \lambda_1 & & 0 \\ & \lambda_2 & \\ 0 & & \lambda_3 \end{pmatrix}$

Let $U$ be the unit cube and $P = T(U)$ its mapping.

The volume of $P$ is then $\lambda_1 \cdot \lambda_2 \cdot \lambda_3$



$\text{vol}(U) = 1$

$\text{vol}(P) = \lambda_1 \cdot \lambda_2 \cdot \lambda_3$

Want to define a quantity "det" that tells us how volumes change under _arbitrary_ linear mappings.

Which properties would such a mapping $d$ have to satisfy?

- $\text{vol}(C) = 1$, so we would like that $d(J) = 1$. $\leadsto$ (D3)

- $T = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix}$ with respect to standard basis $\Rightarrow d(T) = \lambda_1 \cdots \lambda_n$

  $\leadsto$ (D2)

- If the linear mapping does not have full rank
  ( the image $\text{vol}(u)$ is not "full-dimensional", but for example the cube in $\mathbb{R}^3$ is just mapped to a "plane" in 2 dim), then volume is $0$, hence we would like

  $d(u) = 0.$   $\leadsto$ (D2)

Now let's look at the formal definition:

# Definition of the determinant

**Def** Consider a linear mapping $d: F^{n \times n} \to F$. Then $d$ is called a $\boxed{\text{determinant}}$ if:

(D1) $d$ is linear in each column of the matrix:

Let $A$ be a matrix with columns $a_1, \dots, a_n$.

Consider column $a_i$, assume $a_i = a_i' + a_i''$ for some $a_i', a_i'' \in F^{n \times 1}$. Then it holds that

$$\cdot \det\left( (a_1, \dots, a_i, \dots, a_n) \right) =$$

$$\det\left( (a_1, \dots, a_i', \dots, a_n) \right) + \det\left( (a_1, \dots, a_i'', \dots, a_n) \right)$$

$$\cdot \det\left( (a_1, \dots, \lambda a_i, \dots, a_n) \right) = \lambda \cdot \det\left( (a_1, \dots, a_i, \dots, a_n) \right)$$

$$\begin{pmatrix} | & | \\ a_1 & a_2 & \cdots \\ | & | \end{pmatrix}$$

(D2) $d$ is alternating: if $A$ has two identical columns, then $\det A = 0$.

(D3) $d$ is normed: $\det \begin{pmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & 1 \end{pmatrix} = 1$.

# Existence and uniqueness

Theorem : The mapping $\alpha$ exists and is unique.

Proof: Skipped

# Properties of the determinant

Based on (D1), (D2), (D3) we can now prove many important properties of the determinant:

- The determinant of an linear mapping does **not** depend on the basis.

- $\det(c \cdot A) = c^n \det(A)$

- $\det(A \cdot B) = (\det A) \cdot (\det B)$

- $\det(A^t) = \det(A)$

- $\det(A^{-1}) = 1/\det(A)$  (if $A$ is invertible)

- A invertible $\Longleftrightarrow$ $\det(A) \neq 0$

- $\det(A+B) \neq \det(A) + \det(B)$

- If A is upper triangular, that is

$$A = \begin{pmatrix} \lambda_1 & & * \\ & \ddots & \\ 0 & & \lambda_n \end{pmatrix}$$

then $\det A = \lambda_1 \cdot \ldots \cdot \lambda_n$.

Same for lower triangular matrices.

# Computing the determinant (in theory)

Special cases:

$n = 1$    $\det(a) = a$

$n = 2$    $\det \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} = a_{11} \cdot a_{22} - a_{12} \cdot a_{21}$

$n = 3$    $\det \begin{pmatrix} a & b & c \\ d & e & f \\ g & h & i \end{pmatrix} = a \cdot \det \begin{pmatrix} e & f \\ h & i \end{pmatrix} - b \cdot \det \begin{pmatrix} d & f \\ g & i \end{pmatrix}$

$$+ c \cdot \det \begin{pmatrix} d & e \\ g & h \end{pmatrix}$$

In general, there exists the <u>formula of Laplace</u> that expresses the determinant of an $n \times n$ matrix as a lin. comb. of det. of many $(n-1) \times (n-1)$ - submatrices.

# Alternative definition of determinant

There exists a more straight-forward definition of det. However, starting from this definition, proving the geometric properties is more cumbersome.

## Def (alternative):

- If $V$ is a vector space over $\mathbb{C}$ and $T: V \to V$, then $\det(T)$ is the product of all eigenvalues (repeated according to multiplicity).

- If $V$ is a vector space over $\mathbb{R}$ and $T: V \to V$, then $\det(T)$ is the product of its eigenvalues <u>over $\mathbb{C}$</u>, repeated according to multiplicity.

We skip the proof that this is the same object as "our" determinant.

# Computing the determinant (in practice)

- LU decomposition:

  Any matrix $A$ can be written as a product

$$A = L \cdot U$$

where $L$ is a lower triangular matrix and $U$ upper triangular:

$$L = \begin{pmatrix} \ell_{11} & & O \\ & \ell_{22} & \\ * & & \ddots & \ell_{nn} \end{pmatrix}, \quad U = \begin{pmatrix} u_{11} & & * \\ & \ddots & \\ O & & u_{nn} \end{pmatrix}$$

- $\det(A) = \det(L \cdot U) = \det(L) \cdot \det(U) =$

$$= \left( \prod_{i=1}^{n} \ell_{ii} \right) \cdot \left( \prod_{i=1}^{n} u_{ii} \right)$$

# Geometric intuition again

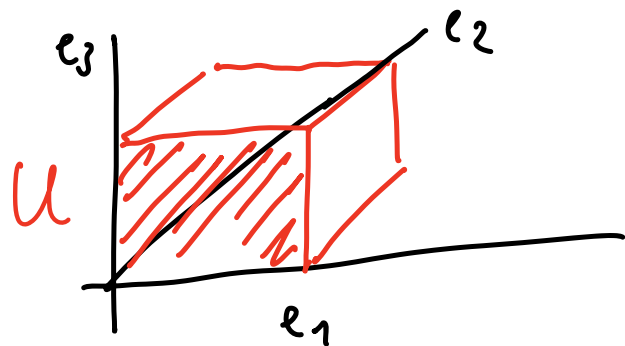**Theorem**: Consider an $n \times n$ matrix $A$ with columns $(a_1 | a_2 | \cdots | a_n) = A$.

Consider the unit cube $U = \{ c_1 e_1 + \ldots + c_n e_n \mid 0 \leq c_i \leq 1 \}$ and its image $P$ under the mapping $A$:
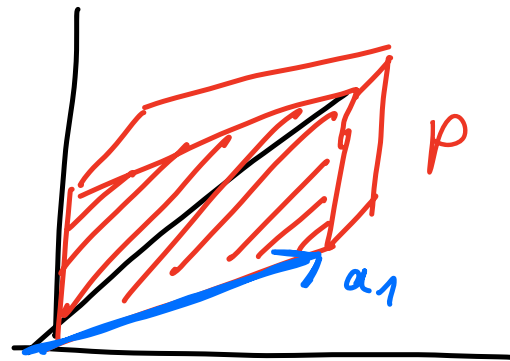
$$U \longmapsto P := \{ c_1 a_1 + \ldots + c_n a_n \mid 0 \leq c_i \leq 1 \} \quad \text{parallelotope.}$$

Then $\det(A)$ gives us the (signed) volume of $P$

**Proof** (for general matrices $A$!): skipped.

# Applications to integrals

**Proposition:** $\Omega \subset \mathbb{R}^n$ open subset, $\sigma : \Omega \to \mathbb{R}^n$ differentiable, $f : \sigma(\Omega) \to \mathbb{R}$. Then:

$$\int_{\sigma(\Omega)} f(y)\, dy = \int_{\Omega} f(\sigma(x)) \, \underbrace{\left| \det(\sigma'(x)) \right| dx}_{\text{volume element}}$$

(volume element under $dy$ on left side)

derivative, linear (label pointing to $\sigma'(x)$)

**Intuition:** $\sigma$ differentiable, that is we can locally (on a small ball $B$ around $x$) approximate $\sigma$ by a linear function. And volumes thus get stretched by the factor given by the determinant.

# Another occurrence of the det in ML

density of the multivariate Gaussian:

$$p(x) = \frac{1}{(2\pi)^{d/2}} \frac{1}{\det(\Sigma)^{1/2}} \exp\left(-\frac{1}{2}(x-\mu)^t \Sigma (x-\mu)\right)$$

normalization term related to the volume, so that
density integrates to 1 in the end

Eigenvalues, eigenvectors, eigenspaces

# Eigenvalues

**Def** Let $T: V \to V$. A scalar $\lambda \in F$ is called an
==eigenvalue== if there exists a $v \in V$, $v \neq 0$,
such that $Tv = \lambda \cdot v$. A vector $v \neq 0$ with this
property is called an ==eigenvector== corresponding to
eigenvalue $\lambda$. The set of all eigenvectors of $\lambda$ is
called the ==eigenspace== $E(\lambda, T) = \ker(T - \lambda I)$.

$$
\begin{cases}
Tv = \lambda v \\
Tv - \lambda v = 0 \\
Tv - \lambda I v = 0 \\
(T - \lambda I) v = 0
\end{cases}
$$

# Geometric intuition

- Eigenvalue/eigenvector realizes a "stretching"
$$v \longmapsto \lambda v$$

- Many mappings do not have eigenvectors for example, a rotation over $\mathbb{R}^2$.

stretching

$v$ $\longrightarrow$ $2 \cdot v$

rotation $R$

$v$ $\longrightarrow$ $Rv$

# Eigenvectors are not unique

- If $\lambda$ is an eigenvalue, it has many eigenvectors!

  For example, if $v$ is eigenvector, then also

  $a \cdot v$ $(a \in K)$ is an eigenvector!

$$T(a \cdot v) = a \cdot T(v) = a \cdot \lambda \cdot v = \lambda(a \cdot v)$$

# Linear (in)dependence of eigenvectors?

- Eigenvectors corr. to <u>distinct</u> eigenvalues are linearly independent. <span style="color:red">( easy exercise; assume they are dependent and drive a contradiction )</span>

- Eigenvectors that corr. to the same eigenvalue do not need to be independent

  <span style="color:red">Simple example: $v$ eig. $\Rightarrow$ $c \cdot v$ eig. but $v$ and $c \cdot v$ are not lin. ind.</span>

  They can be lin. independent:

  <span style="color:red">Easy example: $A = I$, then every vector $v$ is an eigenvector of eigenvalue 1.</span>

- The eigenspace $E(\lambda, T)$ is always a lin. subspace of $V$.

**Theorem :** Every operator $T: V \to V$ on a finite-dim, <u>complex</u> (!) vs V has at least one eigenvalue.

<u>Proof</u>

Step 1:
Getting started

Let $n = \dim V$. Choose a vector $v \in V$, $v \neq 0$. Then the set

$$v, Tv, T^2 v, \ldots, T^n v$$

has to be linearly dependent (it consists of $n+1$ vectors in an $n$-dim space). Find coefficients $a_0, a_1, \ldots, a_n$ such that

$$p(T) := a_0 v + a_1 Tv + \ldots + a_n T^n v = 0.$$

Now we want to show that we can factorize this "polynomial of operators":

$$a_0 + a_1 T + a_2 T^2 + \ldots + a_n T^n$$

$$\stackrel{!}{=} c(T - \lambda_1 I)(T - \lambda_2 I) \ldots (T - \lambda_n I).$$

Excursion to complex-valued polynomials:

Consider a polynomial on $\mathbb{C}$ with these coefficients:

$$p(z) := a_0 + a_1 \cdot z + \ldots + a_n z^n$$

*n and m can be different, eg if $a_n = 0$.*

Over $\mathbb{C}$, we can factorise it:

$$p(z) = c \cdot (z - \lambda_1)(z - \lambda_2) \ldots (z - \lambda_m)$$

Not difficult to see that we can then also factorise $p(T)$:

*\* below*

$$p(T) = a_0 + a_1 T + \ldots + a_n T^n$$

$$= c \cdot (T - \lambda_1 I)(T - \lambda_2 I) \ldots (T - \lambda_m I) \quad ⊛$$

Hence, $0 = a_0 v + a_1 Tv + \ldots + a_n T^n v =$

$$= c(T - \lambda_1 I)(T - \lambda_2 I) \ldots (T - \lambda_n I) \cdot v$$

$\Rightarrow v \in \ker(\text{big operator } \circledast)$

$\Rightarrow$ there must exist $i \in \{0, \ldots, m\}$ such that

$$(T - \lambda_i I) \quad \text{not injective}$$

$\Rightarrow \lambda_i$ is an eigenvalue of $T$ !

$\triangle$ ! The corresponding eigenvector is not necessarily $v$ !

Reason:

$$c(T - \lambda_1 I) \ldots \underbrace{(T - \lambda_i I)}_{\text{not injective}} \underbrace{(T - \lambda_{i+1}) \ldots (T - \lambda_n I) v}_{=: w} = 0$$

$w$ is the eigenvector !

(we started with any vector $v$)

# ✗ Polynomials of operators

Let $p(x) = \sum_{i=0}^{n} a_i x^i$ be a polynomial.

For a linear operator $T$, define $p(T) := \sum_{i=0}^{n} a_i T^i$

Then the following properties ensure that "factorization" works:

Let $p, q$ be two polynomials. Then:

- If the polynomial factorizes, so does the related operator:

$$(p \cdot q)(T) = p(T) \cdot q(T)$$

- Order of factors does not matter:

$$p(T) \cdot q(T) = q(T) \cdot p(T).$$

Moreover, if $p(T) \cdot q(T) v = 0 \implies$

# Diagonalization and Triangularization

# Diagonalizable matrices

<u>Def</u>   An operator $T \in \mathcal{L}(V)$ is <mark><u>diagonalizable</u></mark> if there exists a basis $\mathcal{B}$ of $V$ such that the corr. matrix is diagonal:

$$M(T, \mathcal{B}, \mathcal{B}) = \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_n \end{pmatrix}.$$

⚠ Not always the case, neither over $\mathbb{R}$ nor over $\mathbb{C}$!

Will later see some special cases where it always works (symmetric matrices).

# When is a matrix diagonalizable?

**Proposition:** Let $T$ be an operator on a finite-dim $=n$

vector space $V$. Let $\lambda_1, \ldots, \lambda_m$ denote the distinct eigenvalues of $T$. Then the following statements are equivalent:

(a) $T$ is diagonalizable.

(b) $V$ has a basis consisting of eigenvectors of $T$

(c) $\dim V = \dim(\text{eigenspace}(\lambda_1)) + \ldots + \dim(\text{eigenspace}(\lambda_m))$.

**Remark:** Later will see that this is also equivalent to saying:

(d) $V$ has $n$ eigenvalues if counted with multiplicity, and for all of them the algebraic and geometric multiplicities are the same.

# Triangular matrices

A matrix is called ==upper triangular,== if it has the form

$$\begin{pmatrix} \lambda_1 & & * \\ & \ddots & \\ O & & \lambda_n \end{pmatrix}$$

# Geometric intuition

**Prop** $T \in \mathcal{L}(V)$, $\mathcal{B} = \{v_1, v_2 \cdots, v_n\}$ a basis.

Then equivalent:

(a) $M(T, \mathcal{B})$ is upper triangular.

(b) $T v_j \in \text{span}\{v_1, \cdots, v_j\}$ $\qquad \forall j = 1, \cdots, n$

**Proof idea:**

$$T v_1 = \begin{pmatrix} \lambda_1 & a_{12} & a_{13} \\ & \lambda_2 & a_{23} \\ 0 & & \lambda_3 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} \lambda_1 \\ 0 \\ 0 \end{pmatrix} = \lambda_1 \cdot v_1$$

$$T v_2 = \begin{pmatrix} \lambda_1 & a_{12} & a_{13} \\ & \lambda_2 & a_{23} \\ 0 & & \lambda_3 \end{pmatrix} \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} = \begin{pmatrix} a_{12} \\ \lambda_2 \\ 0 \end{pmatrix} = a_{12} \overbrace{\begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}}^{v_1} + \lambda_2 \overbrace{\begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}}^{v_2}$$

$$\in \text{span}(v_1, v_2)$$

# When are triangular matrices invertible?

**Proposition:** Let $T: V \to V$ have an upper-triangular matrix. Then $T$ is invertible if and only if all entries $\lambda_1, \ldots, \lambda_n$ in the diagonal are non-zero.

$$\begin{pmatrix} \lambda_1 & & * \\ & \ddots & \\ 0 & & \lambda_n \end{pmatrix}$$

# Proof

Proof "$\Leftarrow$" : Let $v_1, \ldots, v_n$ the basis for which $T$ is upper-triangular.

Assume all $\lambda_i \neq 0$.

- By triangular form, $T v_1 = \lambda_1 v_1 \overset{ass.}{\neq} 0$, thus $v_1 \in \text{range}(T)$. $\neq 0$

- By triangular form, $T v_2 = \underbrace{a_1 v_1}_{\in \text{range } T} + \underbrace{\lambda_2 v_2}_{\in \text{range } T}$ for some scalar $a$.

  $\underbrace{\lambda_2 v_2}_{\text{also has to be in range } T}$

  Thus $\lambda_2 v_2 \in \text{range}(T)$, and because $\lambda_2 \neq 0$, then also $v_2 \in \text{range}(T)$.     etc

- In this way, we can see that $v_1, \ldots, v_n$ are all in range $(T)$.
- Thus $T$ surjective, thus injective, thus invertible.

Proof "$\Rightarrow$"   Assume $T$ invertible.

- Clearly $t_1 \neq 0$ ( otherwise $Tv_1 = 0$ thus not invertible).

- Suppose $t_j = 0$.   Then  $T$ maps  $\text{span}(v_1, \ldots, v_j)$ into

  $\text{span}(v_1, \ldots, v_{j-1})$

- Thus $T$ is not injective on the subspace spanned by $v_1, \ldots, v_j$,

  so there exists $v$ in this subspace s.t.  $Tv = 0$.

  But then $T$ is not invertible.

# Entries of diagonal are eigenvalues

Prop  Suppose $T \in \mathcal{L}(V)$, $V$ any finite-dim VS,
T has an upper triangular form. Then the entries
on the diagonal are precisely the eigenvalues
of T.

Holds both over $\mathbb{R}$ and $\mathbb{C}$

# Proof

Fix any $\lambda \in F$ and consider $T - \lambda J$:

$$T = \begin{pmatrix} \lambda_1 & & * \\ & \ddots & \\ 0 & & \lambda_n \end{pmatrix}, \quad T - \lambda J = \begin{pmatrix} \lambda_1 - \lambda & & * \\ & \ddots & \\ 0 & & \lambda_n - \lambda \end{pmatrix}$$

$\lambda$ eigenvalue of $T$ $\iff$ $T - \lambda J$ not invertible

$\iff$ One of the diagonal entries of $T - \lambda J$ is zero

<span style="color:red">previous proposition</span>

$\iff$ $\lambda = \lambda_i$ for some $i$.

# Algebraic multiplicity of eigenvalues

**Definition:** The number of times each eigenvalue occurs on the diagonal is called the ==algebraic multiplicity of the eigenvalue.==

Remark: More often, the algebraic mult. is defined using the characteristic polynomial, which we skipped.
The two definitions are equivalent.

**Definition:** The ==geometric multiplicity== of an eigenvalue is the dimension of the corresponding eigenspace.

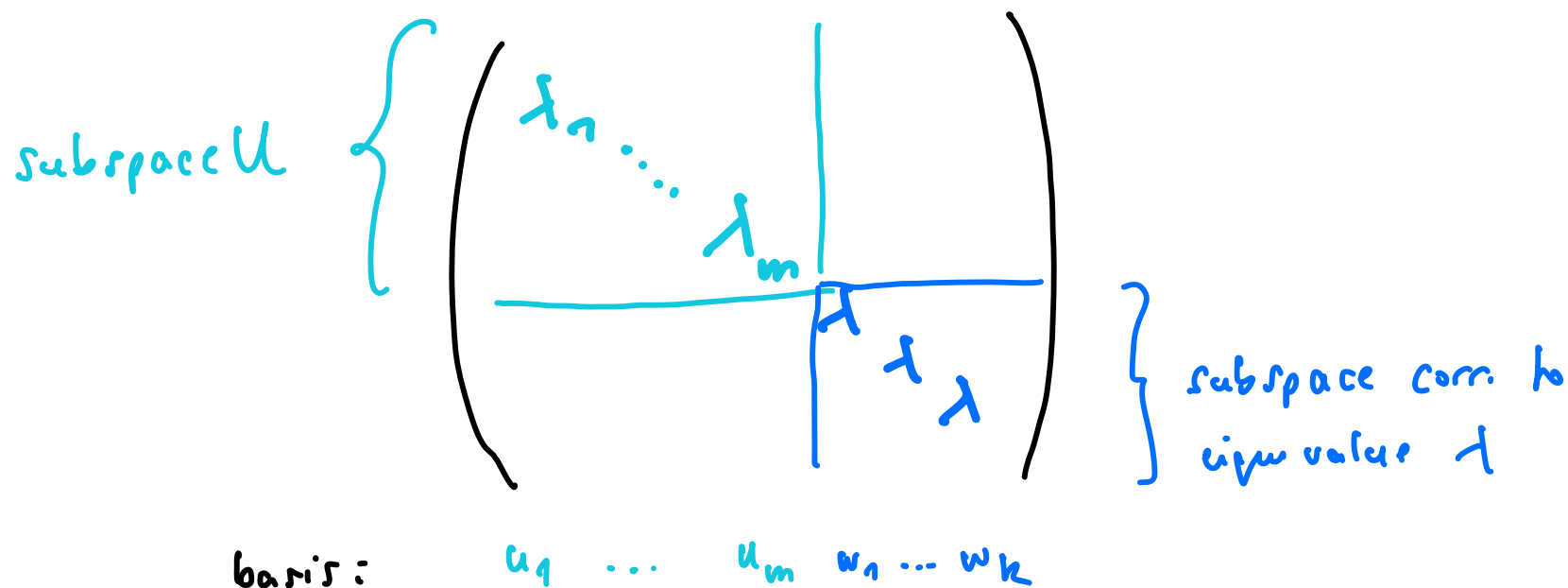⚠ In general, the two do not agree!

# Over $\mathbb{C}$ each matrix can be triangularized

**Prop**    $V$ __complex__ finite-dim VS, $T \in \mathcal{L}(V)$. Then $M(T)$ has an upper triangular form for some basis.

⚠ Does not hold over $\mathbb{R}$!

Proof idea as an image:

- split off the part of the space that belongs to eigenvectors for $\lambda$
- On remainder $U$ apply induction hypothesis
- Then add any basis for $eig(\lambda)$ and show that it does not destroy upper-triang. form.

subspace $U$ $\left\{\begin{pmatrix} \lambda_1 & \cdots & & \\ & & \lambda_m & \\ & & & \lambda & \\ & & & & \lambda \end{pmatrix}\right\}$ subspace corr. to eigenvalue $\lambda$

basis: $u_1 \ \cdots \ u_m \ w_1 \cdots w_k$

**Proof:** induction on n

- know already that one complex eigenvalue exists: $\lambda$
- Consider subspace $U := \text{range} (T - \lambda I)$    "complement of eig $(\lambda)$"
- Easy to see that $U$ is invariant under $T$, that is $TU \subset U$ and that $\dim (\text{range} (u)) < n$.
- We can now apply the induction hypothesis to the operator $T : U \to U$: exists basis $u_1, \ldots, u_m$ $(m < n)$ such that $T_u$ has upper triangular form.
- Now extend $u_1, \ldots, u_m$ by $w_1, \ldots, w_k$ to a basis of $V$ again.
- $T w_1 = T w_1 - \lambda w_1 + \lambda w_1 = (T - \lambda I) w_1 + \lambda w_1 \in \text{span} \{u_1 \ldots u_m, w_1\}$

  Similarly: $T w_i \in \text{span} \{u_1, \ldots, u_m, w_1, \ldots, w_i\}$.
- Thus $T$ is upper-triangular wrt basis $\{u_1 \ldots u_m, w_1, \ldots, w_k\}$.

# Normed spaces

ML keywords: loss functions, regularizers, sparsity,

# Metric space

Definition: Let $X$ be a set. A function $d: X \times X \longrightarrow \mathbb{R}$ is called a ==metric== if the following conditions hold:
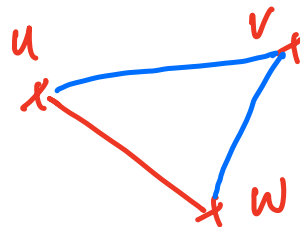
$\forall u, v, w \in X$

(1) $d(x, y) > 0$ if $x \neq y$ and

$\quad\quad d(x, x) = 0$

(2) $d(x, y) = d(y, x)$ $\quad$ (symmetry)

(3) $d(u, v) + d(v, w) \geq d(u, w)$



ML keywords:
inductive bias
k-nearest neighbor

# Norm on a vector space

<u>Def</u>  Let $V$ be a vector space. A <mark>norm</mark> on $V$ is
a function $\|\cdot\| : V \to \mathbb{R}$ such that $\forall x, y \in V$, $\lambda \in F$
the conditions are true:

(N1)  $\| \lambda \cdot x \| = |\lambda| \cdot \|x\|$  ( homogeneous)

(N2)  $\| x + y \| \leq \|x\| + \|y\|$  (triangle inequality)

(N3)  $\|x\| = 0 \Leftarrow x = 0$  $\left.\rule{0pt}{4em}\right\}$ (definiteness)

(N4)  $\|x\| = 0 \Rightarrow x = 0$

$\|\cdot\|$ is a <mark>semi-norm</mark> if (N1) - (N3) are satisfied.

# Euclidean norm on $\mathbb{R}^d$

Euclidean norm on $\mathbb{R}^d$: $\|x\| = \left( \sum_{i=1}^{d} x_i^2 \right)^{1/2}$

Intuition    norm $(x)$ = "length of $x$"
                        = distance $(x, 0)$

Every norm induces a metric: $d(x,y) := \|x - y\|$

But not vice versa (try to find a counter-example!)

# p - Norms on $\mathbb{R}^d$

**Def**   Consider $V = \mathbb{R}^d$.   Define $\|\cdot\|_p : \mathbb{R}^d \to \mathbb{R}$,

$$\|x\|_p := \left( \sum_{i=1}^{d} |x_i|^p \right)^{1/p} \quad \text{for } 0 < p < \infty$$

not always
a norm

The case $p = 2$ coincides with the Euclidean norm.

$$x = \begin{pmatrix} x_1 \\ \vdots \\ x_d \end{pmatrix} \in \mathbb{R}^d$$

Furthermore we define:

$$\|x\|_\infty := \max |x_i| \quad (\text{is a norm})$$

$$\|x\|_0 := \text{number of non-zero coordinates}$$

is not a proper
norm, see later.

$$= \sum_{i=1}^{d} \mathbb{1}\{x_i \neq 0\}$$

# Unit ball of a norm

Def

The ==unit ball== of a norm is the set of points such that norm $\leq 1$:
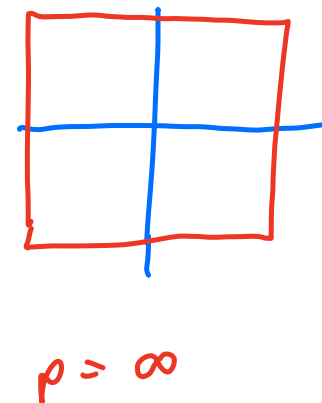
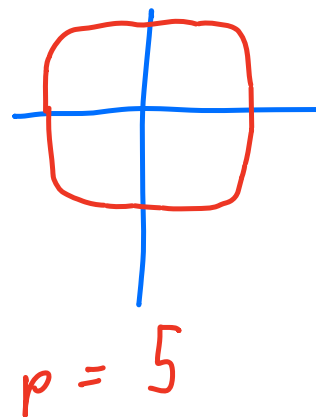$$B_p := \left\{ x \in \mathbb{R}^d \mid \|x\|_p \leq 1 \right\}$$

The ==unit sphere== is the set of points such that norm $= 1$:

$$S_p := \left\{ x \in \mathbb{R}^d \mid \|x\|_p = 1 \right\}$$

# Illustration: unit balls on $\mathbb{R}^2$

**$p \geq 1$:**

(convex balls)

$p = 1$

$p = 2$

$p = 5$

$p = \infty$

**$p < 1$:**

(balls not convex)

$p = 0.5$

$p = 0.1$

$p = 0$

# When is a "p-norm" a norm?

**Prop** $\|\cdot\|_p$ is a norm on $\mathbb{R}^d$ iff $p \geq 1$.

**Proof**
Sketch

- Definiteness (N3, N4) hold for any $p > 0$

- Homogeneity holds for any $p > 0$

- Triangle inequality: this is the critical point, it only holds for $p \geq 1$. This is due to the famous Minkowski inequality, which holds iff $p \geq 1$:

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

■

All norms on $\mathbb{R}^n$ are equivalent

# Equivalent norms (definition)

Definition: Let $V$ be a vector space and $\|\cdot\|_a$ and $\|\cdot\|_b$ two norms on $V$. Then the two norms are called (topologically) equivalent if there exist constants $\alpha, \beta > 0$ such that

$$\forall x \in V: \quad \alpha \|x\|_a \leq \|x\|_b \leq \beta \cdot \|x\|_a$$

**Theorem** All norms on $\mathbb{R}^n$ are (topologically) equivalent.

**Proof:** W.l.o.g. we prove that if $\|\cdot\|$ is any norm on $\mathbb{R}^d$, then it is equivalent to $\|\cdot\|_\infty$.

$\rightarrow$

Let $x = \sum x_i e_i$ the representation of $x$ in the standard basis of $\mathbb{R}^d$.

$$\|x\| = \left\| \sum_{i=1}^{d} x_i e_i \right\|$$

$$\leq \sum_i \| x_i e_i \|$$

$$= \sum_i \underbrace{|x_i|}_{\leq \|x\|_\infty} \|e_i\|$$

$$\leq \sum_i \|x\|_\infty \cdot \|e_i\|$$

$$= \|x\|_\infty \cdot \underbrace{\sum_i \|e_i\|}_{=: c_1}$$

**Second inequality:** $\exists c_2 > 0 \; \forall x \; \|x\|_\infty \leq c_2 \cdot \|x\|$

Let $S := \{x \in \mathbb{R}^d \mid \|x\|_\infty = 1\}$ be the unit sphere wrt $\|\cdot\|_\infty$

Consider $f: S \to \mathbb{R}$, $x \longmapsto \|x\|$.

- The mapping $f$ is continuous wrt $\|\cdot\|_\infty$:

  ( this follows directly from the fact that

  $$|f(x) - f(y)| = |\, \|x\| - \|y\| \,|$$
  $$\leq \|x - y\| \leq c_1 \cdot \|x - y\|_\infty \; ).$$

- The $S$ is closed and bounded, thus by Theorem of Heine-Borel, $S$ compact. Any continuous mapping on a compact set takes its min and max.

$$\tilde{c}_2 := \min \{f(x) \mid x \in S\}$$

- Because $0 \notin S$ (sphere, not ball), we can conclude from the definiteness that $\tilde{c}_2 \neq 0$.

- Now compute: For $x \in S$ we have

$$\tilde{c}_2 \leq \|x\| = \left\| \frac{x}{1} \right\| \overset{\downarrow}{=} \left\| \frac{x}{\|x\|_\infty} \right\| = \frac{\|x\|}{\|x\|_\infty}$$

choice of $\tilde{c}_2$

$$\Rightarrow \quad \|x\|_\infty \leq \frac{1}{\tilde{c}_2} \|x\|$$

$$=: c_2$$

Scalar product

# Scalar product, definition

__Def.__ Consider vector space $V$. A mapping $\langle \cdot , \cdot \rangle : V \times V \to \mathbb{R}$

is called a scalar product if

linearity
$$
\begin{cases}
(S1) & \langle x_1 + x_2 , y \rangle = \langle x_1 , y \rangle + \langle x_2 , y \rangle \\
(S2) & \langle \lambda x , y \rangle = \lambda \langle x , y \rangle
\end{cases}
$$

symmetry
$$
\begin{cases}
(S3) & \langle x , y \rangle = \langle y , x \rangle \qquad (\text{if on } \mathbb{R}) \\
& \langle x , y \rangle = \overline{\langle y , x \rangle} \qquad (\text{if on } \mathbb{C})
\end{cases}
$$

complex conjugate

positive
definite
$$
\begin{cases}
(S4) & \langle x , x \rangle \geq 0 \\
(S5) & \langle x , x \rangle = 0 \iff x = 0
\end{cases}
$$

# Examples

- Euclidean scalar product on $\mathbb{R}^n$: $\quad x = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}, \; y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$

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

- On $\mathbb{C}^n$, $\quad \langle x, y \rangle = \sum_{i=1}^{n} x_i \overline{y_i}$

- $\mathcal{C}([a,b]):$ $\quad \langle f, g \rangle = \int_{a}^{b} f(t) \, g(t) \, dt$

  is a scalar product (but space would not be complete)

# Scalar product and angles

**Prop**   Consider the standard scalar product on $\mathbb{R}^n$. Then

$$\langle v, w \rangle = \|v\| \|w\| \cos(\alpha)$$

where $\alpha$ is the angle enclosed by $v$ and $w$.

**Proof**  • In a general triangle we have
$$\|v-w\|^2 = \|v\|^2 + \|w\|^2 - 2\|v\|\|w\| \cos\alpha$$

• $\|v-w\|^2 = \langle v-w, v-w \rangle = \|v\|^2 + \|w\|^2 - 2\langle v, w \rangle$

$$\Rightarrow \langle v, w \rangle = \|v\| \|w\| \cos(\alpha).$$

ML keyword:
cosine similarity

# Banach and Hilbert spaces

Def    A vector space with a norm is called a
<mark>normed space</mark>. If a normed space is complete
(each Cauchy sequence converges), then V is called
a <mark>Banach space</mark>. A VS with a scalar product is called
a <mark>pre-Hilbert-space</mark>. If it is additionally complete,
then it is called <mark>Hilbert space</mark>.

ML keyword:
RKHS: reproducing
kernel Hilbert
space

# Relationship between norm and scalar product

Scalar product $\Rightarrow$ norm

$\not\Leftarrow$

Consider a VS with a scalar product $\langle \cdot, \cdot \rangle$. Define $\| \cdot \| : V \to \mathbb{R}$ as $\boxed{\| x \| := \sqrt{\langle x, x \rangle}}$. Then $\| \cdot \|$ is a norm on $V$, the norm induced by $\langle \cdot, \cdot \rangle$.

The other way round does not work in general!

(can you find a counterexample?)

# Relationship between norm and metric

norm $\Rightarrow$ metric

$\not\Leftarrow$

Consider a VS $V$ with norm $\|\cdot\|$. Then

$$d: V \times V \to \mathbb{R} \quad , \quad d(x,y) := \|x - y\|$$

is a metric on $V$, the metric induced by the norm.

The other direction does not work in general.
(Can you find a counterexample?)

# Important inequalities

Consider $u, v \in \mathbb{R}^d$, and denote by $\|\cdot\|_p$ the $p$-norm.

## Cauchy - Schwarz inequality:

$$|\langle u, v \rangle| \leq \|u\|_2 \|v\|_2$$

## Hölder inequality:

Let $p, q \geq 1$ with $1/p + 1/q = 1$. Then:

$$|\langle u, v \rangle| \leq \sum_{i=1}^{d} |u_i v_i| \leq \|u\|_p \|v\|_q$$

# Orthonormal basis
# and
# orthogonal projections

# Orthogonal vectors and sets

**Def** Consider a pre-Hilbert space $V$.

- Two **vectors** $v_1, v_2 \in V$ are called **orthogonal** if $\langle v_1, v_2 \rangle = 0$.

  Notation: $v_1 \perp v_2$

- Two **sets** $V_1, V_2 \subset V$ are called orthogonal if

  $$\forall v_1 \in V_1 \ \forall v_2 \in V_2 : \quad \langle v_1, v_2 \rangle = 0$$

- For a set $S \subset V$ we define its **orthogonal complement** $S^\perp$ as follows:

  $$S^\perp := \{ v \in V \mid v \perp s \ \forall s \in S \}$$

# Orthonormal vectors and sets

Def:
___

- Two ==vectors== $v_1, v_2$ are called ==orthonormal== if they are orthogonal and additionally the two vectors have norm 1:

$$\cdot \langle v_1, v_2 \rangle = 0$$

$$\cdot \|v_1\| = 1, \quad \|v_2\| = 1$$

- A ==set== of vectors $v_1, v_2, \ldots, v_n$ is called orthonormal if any two vectors are orthonormal.

# Orthogonal/orthonormal basis

We are particularly interested in orthogonal/orthonormal bases of a space:

Proposition: In an orthonormal basis $u_1, ..., u_n$, the representation of a vector $v$ is given as

$v = \sum a_i u_i$

$$v = \sum_{i=1}^{n} <v, u_i> u_i$$

Proof exercise

⚠️ We still don't know whether an orthonormal basis always exists...

# Projection (general case)

Def    $A \in \mathcal{L}(V)$ is called a <mark>projection</mark> if $A^2 = A$.

blue vector gets
projected on red
vector (not
orthogonal)

# Orthogonal projection

**Theorem & Def** : Let $U$ be a finite-dim subspace of a pre-Hilbert space $H$. Then there exists a linear projection $P_U : H \to U$ with $\ker(P_U) = U^\perp$. $P_U$ is then called the <mark>orthogonal projection</mark> of $H$ on $U$.

**Proof idea:** Let $u_1, ..., u_k$ be an orthogonal basis of $U$.

Define $\quad p_u : V \to U \;\text{by}\quad p_u(w) = \sum_{i=1}^{k} \frac{\langle w, u_i \rangle}{\|u_i\|} u_i$

Obviously:

- $p_u$ linear

- $w \in U^+$

  $\Rightarrow p_u(w) = 0.$

# How to make a given basis orthogonal?

Proposition: Let $v_1, ..., v_n$ be any basis of a pre-Hilbert space. Then we can transform it into an orthonormal basis $u_1, ..., v_n$.

Proof:

## Gram-Schmidt orthogonalization

Is a procedure that takes any basis $v_1, ..., v_n$ of a finite-dim VS and transforms it into another basis $u_1, ..., u_n$ that is orthogonal:

Intuition: iterative procedure

Step 1:  $u_1 := \dfrac{v_1}{\|v_1\|}$ ,  $U_1 := \text{span}\{u_1\}$

Step k: Assume that we already identified $u_1, \dots, u_{k-1}$.

- Project $v_k$ on $U_{k-1}$, and keep "the rest":

$$\tilde{u}_k := v_k - P_{U_{k-1}}(v_k)$$

- Renormalize:

$$u_k = \tilde{u}_k / \|\tilde{u}_k\|$$

Works in theory      (would need to prove that, skipped)

In practice, it quickly results in large numerical errors.

$\leadsto$ QR factorization, see later.

Orthogonal matrices

# Orthogonal matrices, definition

<u>Def</u> Let $Q \in \mathbb{R}^{n \times n}$ be a matrix with orthonormal (!) column vectors (wrt Euclidean scalar product). Then $Q$ is called an ==orthogonal (!) matrix.==

If $Q \in \mathbb{C}^{n \times n}$ and the columns are orthonormal (wrt the standard scalar product on $\mathbb{C}^n$), then it is called a ==<u>unitary matrix.</u>==

# Orthogonal matrices, examples

- Identity: $\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$ , Reflection: $\begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}$

- Permutation of coordinates: $\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$

- Rotation in $\mathbb{R}^2$: $\begin{pmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{pmatrix}$

- Rotation in $\mathbb{R}^3$:

  - Rotation about one of the axes:

  $$R_{\theta,1} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & \cos\theta & -\sin\theta \\ 0 & \sin\theta & \cos\theta \end{pmatrix}$$

  - General rotation: can be written as a product of "elementary" rotations

# Properties of orthonormal matrices

Let $Q$ be orthogonal. Then:

- Columns are orthogonal $\iff$ rows are orthogonal     <span style="color:red">* see also slide below</span>

- $Q$ is always invertible, and $Q^{-1} = Q^t$

- $Q$ realizes an <u>isometry</u>: $\forall v \in V: \|Qv\| = \|v\|$

- $Q$ preserves angles: $\langle Qu, Qv \rangle = \langle u, v \rangle \quad \forall u, v \in V$

- $|\det Q| = 1$

The respective properties also hold for unitary matrices $U$:

- $U^{-1} = \overline{U}^t$

<u>Proofs</u>: assignment!

# Orthogonal rows and columns

⚠️

Consider the projection matrix $A = \begin{pmatrix} 0 & 0 \\ 2 & 1 \end{pmatrix}$. The columns are obviously not orthogonal. The rows formally satisfy that $\langle \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 2 \\ 1 \end{pmatrix} \rangle = 0$. The property that "rows orthogonal $\Leftrightarrow$ cols orthogonal" does not hold here. But note that $A$ is __not__ an orthogonal matrix because the later requires all rows/cols to have norm $1$ (in particular, also full rank).

The statement "rows orthogonal $\Leftrightarrow$ columns orthogonal" clearly does not hold for arbitrary matrices.

# Representation of isometric mappings

**Theorem** Let $S \in \mathcal{L}(V)$ for a real VS $V$. Then equivalent:

(a) $S$ is an isometry: $\|Sv\| = \|v\|$   $\forall v \in V$

(b) There exists an orthonormal basis of $V$ such that the matrix of $S$ has the following form:

$$M = \begin{pmatrix} \boxed{} & & & & & & 0 \\ & \boxed{} & & & & & \\ & & \boxed{} & & & & \\ & & & \boxed{} & & & \\ & & & & \boxed{} & & \\ & & & & & \ddots & \\ 0 & & & & & & \boxed{} \end{pmatrix}$$

where each of the little blocks

- either is a $1 \times 1$ matrix ($\hat{=}$ a real number) being $1$ or $-1$

- or is a $2 \times 2$ rotation matrix: $\begin{pmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{pmatrix}$

**Proof** skipped.

# Symmetric matrices

In ML, symmetric matrices are all over the place:
covariance matrix, kernel matrix, similarity matrix, ...

# Symmetric matrices

Def  A matrix $A \in \mathbb{R}^{n \times n}$ is called symmetric if $A = A^t$.

A matrix $A \in \mathbb{C}^{n \times n}$ is called hermitean if $A = \bar{A}^t$.

# Hermitian matrices have real-valued eigenvalues and orthogonal eigenvectors

**Prop**    Let $A \in \mathbb{C}^{n \times n}$ be hermitian. Then __all__ eigenvalues of $A$ are real-valued. Eigenvectors that correspond to different eigenvalues are orthogonal.

**Proof** •   $\lambda$ eig. value of $A$ with eigenvector $x$. Then

$$\lambda \langle x, x \rangle = \langle \lambda x, x \rangle = \langle Ax, x \rangle = \qquad \text{hermitian}$$

$$= \langle x, Ax \rangle = \langle x, \lambda x \rangle = \bar{\lambda} \langle x, x \rangle$$

property of scalar prod. on $\mathbb{C}$

$$\Rightarrow \lambda = \bar{\lambda} \in \mathbb{R}$$

- $(\lambda_1, x_1)$ , $(\lambda_2, x_2)$ eigs of $A$, $\lambda_1 \neq \lambda_2$. Then:

$$\underline{\lambda_1 \langle x_1, x_2 \rangle} = \dots \text{as above} \dots = \underline{\lambda_2 \langle x_1, x_2 \rangle}$$

$$\Rightarrow \quad 0 = \lambda_1 \langle x_1, x_2 \rangle - \lambda_2 \langle x_1, x_2 \rangle =$$

$$= (\lambda_1 - \lambda_2) \langle x_1, x_2 \rangle$$

$\Rightarrow$ either $\lambda_1 = \lambda_2$

or if $\lambda_1 \neq \lambda_2$, then $\langle x_1, x_2 \rangle = 0$

$\Rightarrow x_1 \perp x_2$ .

# What can we conclude from here?

- Each matrix $\in \mathbb{C}^{n \times n}$ has $n$ eigenvalues $\in \mathbb{C}$
- For hermitian matrices $\in \mathbb{C}^{n \times n}$ all $n$ eigenvalues now live in $\mathbb{R}$
- The eigenvectors still might live in $\mathbb{C}^{n \times n}$ (and not in $\mathbb{R}^{n \times n}$).

Symmetric matrices $\in \mathbb{R}^{n \times n}$ are a special case of hermitian, (because $A = \underset{sym}{A^t} = \underset{\mathbb{R}}{\bar{A}^t}$), so it also has $n$ real-valued eigenvalues. But the eigenvectors might still live in $\mathbb{C}^n$.

# Self-adjoint operators

**Def**   An operator $T \in \mathcal{L}(V)$ on a pre-Hilbert space $V$
is called <mark>self-adjoint</mark> if

$$\langle Tv, w \rangle = \langle v, Tw \rangle.$$

Sometimes it is called a <mark>Hermitean operator</mark> (on $\mathbb{C}^n$)
                                        <mark>symmetric operator</mark> (on $\mathbb{R}^n$).

**Remark**   Over $\mathbb{C}^n$, self-adjoint operators are represented by
hermitean matrices.   On $\mathbb{R}^n$, self-adjoint op. are represented
by symmetric matrices.

# Self-adjoint operator has real-valued eigenvalue and eigenvector

Prop: $V$ vector space over $\mathbb{C}$ or $\mathbb{R}$, $T \in \mathcal{L}(V)$, $T \neq 0$, self-adjoint.

Then $T$ has at least one eigenvalue and it is <u>real-valued</u>.

In particular in case of $\mathbb{R}$, also the eigenvector lives in $\mathbb{R}^n$.

⚠ Already know: over $\mathbb{C}$, every matrix has an eigenvalue. But it could be a complex number. Now: if $T$ is self-adjoint, have a <u>real</u> eigenvalue.

Proof

$n := \dim V$. Chose $v \neq 0$ ($\in \mathbb{R}^n$), and consider

$$v, \, Tv, \, T^2 v, \, \ldots, \, T^n v. \, \in \mathbb{R}^n$$

These vectors have to be lin. dependent ($n+1$ vectors, $\dim = n$)

There exist $a_0, \ldots, a_n \in \mathbb{R}$ (not all $0$) such that

$$a_0 v + a_1 Tv + \ldots + a_n T^n v = 0.$$

Consider the polynomial with these coefficients and decompose it over $\mathbb{R}$ into linear and quadratic factors:

$$a_0 + a_1 x + a_2 x^2 + \cdots + a_n x^n =$$

decompose over $\mathbb{R}$!
$c \neq 0, \ b_i, c_i, \lambda_i \in \mathbb{R}$
$M + m \geq 1$

$$= \underset{\neq 0}{c} \underbrace{(x^2 + b_1 x + c_1) \cdots (x^2 + b_M x + c_M)}_{\substack{\text{quadratic terms that cannot} \\ \text{be decomposed into linear} \\ \text{terms;} \\ \text{in particular they satisfy} \\ b_i^2 < 4 c_i \quad (\text{otherwise one} \\ \text{could factorize them by the} \\ \text{quadratic formula})} \cdot \underbrace{(x - \lambda_1) \cdots (x - \lambda_m)}_{\text{linear terms}}$$

Replace the $x$ by $T$:

$$0 = (a_0 + a_1 T + ... + a_n T^n) v = \left( c (\underbrace{...}_{quadr.}) \underbrace{(----)}_{lin. \, terms} \right) \cdot v$$

Now can prove: the "quadratic operators" are invertible. (see (*) below).

So we multiply the equation with their inverses and obtain

$$0 = (T - \lambda_1 I) \cdot ... (T - \lambda_n I) v$$

Because we had chosen $v \neq 0$, at least one of the terms $(T - \lambda_i I)$ has to exist (i.e. $m \geq 1$)

But then also the product $(T - \lambda_1 I) \cdot ... \cdot (T - \lambda_m I)$ is not injective, thus at least one of the factors just not injective, thus $\lambda_j$ eigenvalue. over $\mathbb{R}$. In particular, the eigenvector lives in $\mathbb{R}^n$ as well.

(\*)

**Proposition:** Suppose $T$ self-adjoint and $b, c \in \mathbb{R}$ satisfy

$b^2 < 4c$. Then $T^2 + bT + cI$ is invertible.

**Proof** Intuition:

$$x^2 + bx + c = \left(x + \frac{b}{2}\right)^2 + \left(c - \frac{b^2}{4}\right) > 0$$

$\underbrace{\phantom{\left(x + \frac{b}{2}\right)^2}}_{\geq 0}$   $\underbrace{\phantom{\left(c - \frac{b^2}{4}\right)}}_{> 0 \text{ by ass.}}$

So $x^2 + bx + c$ is an invertible real number.

Now do it: Consider any $v \neq 0$.

$$\langle (T^2 + bT + cI)v, v \rangle =$$

$$= \underbrace{\langle T^2 v, v \rangle}_{\substack{= \langle Tv, Tv \rangle \\ = \|Tv\|^2}} + \underbrace{b \langle Tv, v \rangle}_{\leq \|Tv\| \|v\| \text{ Cauchy-Schwartz}} + c \|v\|^2$$

$$\geq \|Tv\|^2 - |b| \|Tv\| \|v\| + c \|v\|^2$$

$$= \left( \|Tv\| - \frac{|b| \|v\|}{2} \right)^2 + \left( c - \frac{b^2}{4} \right) \|v\|^2 \quad > 0$$

$$\underbrace{\phantom{\left( \|Tv\| - \frac{|b| \|v\|}{2} \right)^2}}_{\geq 0} \quad \underbrace{\phantom{\left( c - \frac{b^2}{4} \right)}}_{> 0 \text{ by ass.}}$$

$$\Rightarrow \quad (T^2 + bT + cI)v \neq 0 \quad \text{for all} \quad v \neq 0$$

$$\Rightarrow \quad T^2 + bT + cI \quad \text{invertible.} \qquad \blacksquare$$

# Spectral theorem for symmetric / hermitian matrices

**Theorem:** A symmetric matrix $A \in \mathbb{R}^{n \times n}$ is orthogonally diagonalizable: there exists an orthogonal matrix $Q \in \mathbb{R}^{n \times n}$ and a diagonal matrix $D \in \mathbb{R}^{n \times n}$ s.t.

$$A = Q D Q^t = \sum_{i=1}^{n} \lambda_i \, q_i \, q_i^t \; .$$

$$D = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix}, \quad Q = \begin{pmatrix} | & | & \\ q_1 & q_2 & \cdots \\ | & | & \end{pmatrix}$$

Proof (sketch)     By induction on $n := \dim V$

Base case $n = 1$:     clear

Inductive step $n - 1 \rightsquigarrow n$

previous theorem

$\lambda \in \mathbb{R}$

- A symmetric $\Rightarrow$ A has at least one eigenvector $u \in \mathbb{R}^n$

- $U := \mathrm{span}\{u\}$.  U is invariant under A.     $\forall v \in U: Av \in U$

- Consider $U^\perp$ and the restriction of A to $U^\perp$.

  On $U^\perp$, A is again a symmetric operator

  and $\dim(U^\perp) = n - 1$.

- Apply the induction hypothesis on this space of dim $n-1$.

  Does the job!

# Complex version of this theorem

(often not as relevant to ML as the real version above)

**Theorem** A hermitian matrix $A \in \mathbb{C}^{n \times n}$ is unitarily diagonalizable: there exists a unitary matrix $U$ and a diagonal matrix $D$ s.th.

$$A = U D \bar{U}^t$$

In particular, the entries of $D$ are real-valued.

# Positive definite matrices

# Positive definite matrices

Def    A matrix $A \in \mathbb{R}^{n \times n}$ is called

positive definite (pd) [semi-definite] (psd) if $\forall x \in \mathbb{R}^n$, $x \neq 0$:

$$x^t A x > 0.$$
$$\geq$$

Def    A matrix $A$ is called a Gram matrix

if there exists a set of vectors $v_1, \ldots, v_n$ s.th.

$$a_{ij} = \langle x_i, x_j \rangle.$$

Observe:  On $\mathbb{C}^{n \times n}$, Gram matrices are hermitian

On $\mathbb{R}^{n \times n}$, Gram matrices are symmetric.

# Characterization of pd matrices over $\mathbb{C}$

**Theorem** : $A \in \mathbb{C}^{n \times n}$ hermitean. Then equivalent:

(i)   $A$ is psd (pd)

(ii)   All eigenvalues of $A$ are $\geq 0$  ($> 0$)

(iii)   The mapping $\langle \cdot, \cdot \rangle_A : \mathbb{C}^n \times \mathbb{C}^n \longrightarrow \mathbb{C}$ with

$$\langle x, y \rangle_A := \bar{y}^t A x$$

satisfies all properties of a scalar product

except one: if $\langle x, x \rangle_A = 0$ this does not

imply $x = 0$.

(iv)   $A$ is a Gram matrix of $n$ vectors $a_{ij} = \langle x_i, x_j \rangle$

which are not necessarily lin. independent

which are lin. independent

⚠️ Observe that implicitly this theorem implies that over $\mathbb{C}$ we have pd $\Rightarrow$ self-adjoint.

Over $\mathbb{R}$, this is **not** true: pd $\not\Rightarrow$ symmetric

Thus to get a similar characterization of pd matrices over $\mathbb{R}$, we need to add a symmetry condition.

Example:     $A = \begin{pmatrix} 1 & 1 \\ -1 & 1 \end{pmatrix}$

on $\mathbb{R}$, if $x \neq 0$

$x^t A x = x_1^2 + x_2^2 > 0$

• So $A$ is pd but not symmetric.

• Over $\mathbb{C}$, the same matrix is not pd because $x_1^2 + x_2^2$ can be negative !

# Characterization of symmetric, pd matrices over $\mathbb{R}$

$$\begin{bmatrix} 0 & & D \end{bmatrix}$$

Add

# Roots of psd matrices

**Theorem:** Let $A \in \mathbb{R}^{n \times n}$ be symmetric, psd. Then there exists a matrix $B \in \mathbb{R}^{n \times n}$, $B$ psd such that $A = B^2$. Sometimes $B$ is called the square root of $A$, sometimes denoted as $B = (A)^{1/2}$.

# Proof

- Spectral theorem $\Rightarrow$ $A = U D U^t$, $\quad D = \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_n \end{pmatrix}$

- psd $\Rightarrow$ eigenvalues $\lambda_i \geq 0$

- Define $\sqrt{D} := \begin{pmatrix} \sqrt{\lambda_1} & & \\ & \ddots & \\ & & \sqrt{\lambda_n} \end{pmatrix}$ and set

$$B := U \sqrt{D} U^t.$$ Does the job. $\blacksquare$

Remark: more generally, one can define $A^{1/k}$ for any $k \in \mathbb{N}$. And because we can also define $A^{-1}$ (in case pd) and $A^k$, one can define $A^{p/q}$ for $p, q \in \mathbb{Z}$.

# Important psd matrices for ML

- Consider a **data matrix** $X \in \mathbb{R}^{n \times d}$ : $n$ data pts, $d$ dims

- Then the matrix $C := X^t X \in \mathbb{R}^{d \times d}$ is called the **covariance matrix**.

- The matrix $K := X X^t \in \mathbb{R}^{n \times n}$ is called a **kernel matrix** (for the linear kernel), and in this special case it is also the **Gram matrix**.

- All these matrices are symmetric and positive semi-definite (prove it!)

# Variational characterization of eigenvalues

ML keyword:
spectral clustering

Literature: Bhatia : Matrix Analysis.

# Rayleigh coefficient

<u>Def</u>   Let $A \in \mathbb{R}^{n \times n}$ be a symmetric matrix.

$$R_A : \mathbb{R}^n \setminus \{0\} \to \mathbb{R} , \quad x \mapsto \frac{x^t A x}{x^t x} =: \frac{x^t A x}{\|x\|^2}$$

is called the Rayleigh coefficient.

**Rayleigh coeff.** $\rightsquigarrow$ **first eigenvalue**

**Prop**

Let $A$ be symmetric, let $\lambda_1 \leq \lambda_2 \leq \ldots \leq \lambda_n$ be the eigenvalues and $v_1, \ldots, v_n$ the eigenvectors of $A$.

Then:

$$\min_{x \in \mathbb{R}^n} R_A(x) = \min_{\|x\|=1} x^t A x = \lambda_1 \quad, \quad \text{attained at } x = v_1$$

$$\max_{x \in \mathbb{R}^n} R_A(x) = \max_{\|x\|=1} x^t A x = \lambda_n \quad, \quad \text{attained at } v_n.$$

# Intuition for the proposition

Assume $A$ is expressed in terms of the basis $v_1, \ldots, v_n$ by

$$A = \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_n \end{pmatrix}.$$

Let $y$ be a vector, also represented in this basis:

$$y = y_1 v_1 + y_2 v_2 + \cdots + y_n v_n \qquad y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

$\circledR$  $y^t A y = \underline{\lambda_1 y_1^2} + \cdots + \lambda_n y_n^2$

Among the vectors $\begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \ldots, \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{pmatrix}$    $\begin{pmatrix} 0.5 \\ 0.5 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$

the smallest result of $y^t A y$ would be given by the vector $\begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$, and the value would be $\lambda_1$

$$y_1 = v_1$$

# Formal proof (sketch)

Assume we start with the standard basis.

Let $Q = \begin{pmatrix} \mid & & \mid \\ v_1 & \cdots & v_n \\ \mid & & \mid \end{pmatrix}$ be the basis transformation that brings $A$

in diagonal form: $Q$ orthogonal such that

$$A = Q^t \underline{\Lambda} Q \quad \text{with} \quad \underline{\Lambda} \text{ diagonal}.$$

For a vector $x = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$ in the original basis, we now consider

the transformed vector $y := Q^t x$ and compute

its Rayleigh coefficient:

$$R_A(\gamma) = \frac{\overset{Y^t}{\overbrace{(Q^b x)^t}} \overset{A}{\overbrace{(Q^t \Lambda Q)}} \overset{Y}{\overbrace{(Q^b x)}}}{\underset{Y^b}{\underbrace{(Q^t x)^t}} \underset{Y}{\underbrace{(Q^b x)}}}$$

$$(Q^t x)^t = x^t Q$$

$$= \frac{x^t \cancel{Q} \; Q^t \Lambda Q \; \cancel{Q^b} x}{x^t \cancel{Q} \; \cancel{Q^t} x}$$

$$= \frac{x^t \Lambda x}{x^t x} = \frac{\lambda_1 x_1^2 + \cdots + \lambda_u x_u^2}{\|x\|}$$

Now we look at the minimum of $R_A(y)$:

$$\min_{\|y\|=1} R_A(y) = \min_{\|x\|=1} \lambda_1 x_1^2 + \dots + \lambda_n x_n^2$$

This min. is attained for $x = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$, that is

$y = Q^t x = v_1$, with value $R(y) = \lambda_1$.

# Rayleigh coeff ↝ second eigenvalue

**Prop**

Consider a symmetric matrix $A$ with eigenvalues

$$\lambda_1 \leq \lambda_2 \leq \ldots \leq \lambda_n$$

problem

Consider the optimization

$$\min_{\substack{\|x\| = 1 \\ x \perp v_1}} R(x)$$

This problem is solved by $x = v_2$, $R(v_2) = \lambda_2$.

# Proof intuition

Consider operator $A$ restricted to the space

$$V_1^{\perp} := \left( \text{span} \{ v_1 \} \right)^{\perp}.$$ We know that on this

space, $A$ is invariant and symmetric, so we can

apply Rayleigh to this "smaller" space.

$$V_1^{\perp} = \text{span} \{ v_2 , \cdots , v_n \}$$

If we apply Rayleigh to $V_1^{\perp}$, then we get

the solutions $\lambda_2, v_2$.

# Min-max-Theorem / Courant-Fischer-Weyl

**Theorem**   $A \in \mathbb{R}^{n \times n}$ symmetric; eigenvalues $\lambda_1 \leq \ldots \leq \lambda_n$. Then:

$$\lambda_k = \min_{\substack{U \text{ subspace} \\ \dim U = k}} \quad \max_{x \in U \setminus \{0\}} \quad R_A(x) \tag{1}$$

$$= \max_{\substack{U \text{ subspace,} \\ \dim U = n-k+1}} \quad \min_{x \in U \setminus \{0\}} \quad R_A(x) \tag{2}$$

# Proof intuition

- For case $k = 1$, case (2) is pretty much what we have proved already:

$$\max_{\substack{U \text{ subspace} \\ \dim U = n-k+1}} \min_{x \in U \setminus \{0\}} R_A(x) =$$

$\underbrace{n-k+1}_{n}$ can drop it

previous result

$$= \min_{x \in V \setminus \{0\}} R_A(x) \overset{\downarrow}{=} \lambda_1.$$

Case (1) follows similar principles.

- case $k = 2$ .... similar to the previous statement.

- General case: induction.

# Matrix norms

ML keywords:
- low rank approximation
- Perturbation analysis

# Motivation

Want to quantify the "similarity"
of various matrices.

So we define "norms" on the space
of all matrices.

⚠️ Not all of them are proper norms.

# Definitions of matrix norms

Given a matrix $A \in \mathbb{R}^{m \times n}$. Define the following norms:

$$\|A\|_{max} = \|A\|_{\infty} = \max_{ij} |a_{ij}|$$

$$\|A\|_F = \sqrt{\sum_{ij} a_{ij}^2} = \sqrt{\operatorname{tr}(A^t A)}$$

trace, see later

$$= \sqrt{\sum \sigma_i^2} \quad \text{where } \sigma_i \text{ are the singular values of } A.$$

Frobenius norm

see next slide

$$\|A\|_2 = \sigma_{max}(A) \quad \text{where} \quad \sigma_{max} \text{ is the largest singular value}$$

$$= \max_{x \neq 0} \frac{\|Ax\|}{\|x\|} \quad \text{Euclidean norm on vectors in } \mathbb{R}^m$$

"Operator norm", "spectral norm"

$$\|A\|_* = tr\left(\sqrt{A^t A}\right) \quad \text{nuclear norm}$$

trace see later

Many more matrix norms exist ...

# Simple inequalities

Let $A \in \mathbb{R}^{m \times n}$. Then:

$$\|A\|_2 \leq \|A\|_F \leq \sqrt{n} \, \|A\|_2$$

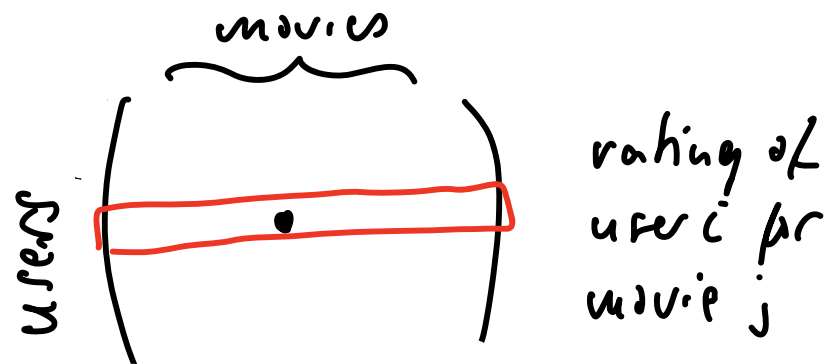$$\frac{1}{\sqrt{n}} \|A\|_\infty \leq \|A\|_2 \leq \sqrt{m} \, \|A\|_\infty$$

$$\|A\|_2 \leq \sqrt{\|A\|_1 \cdot \|A\|_\infty}$$

Many, many more... see eg the book on Matrix Analysis by Bhatia.

# Singular value decomposition

# UL motivation: recommender systems

- Netflix ratings: huge matrix!



movies

users

rating of user $i$ for movie $j$

- ratings of a particular user are "not random" but have structure
- "compress" the matrix into something much smaller that also better represents the "structure" of this matrix.

# Singular value decomposition

Proposition   Consider $A \in \mathbb{R}^{m \times n}$ of rank $r$. Then we can write $A$ in the form

$$A = U \cdot \Sigma \cdot V^t$$

where $U \in \mathbb{R}^{m \times m}$, $V \in \mathbb{R}^{n \times n}$ are orthogonal matrices and $\Sigma \in \mathbb{R}^{m \times n}$ is "diagonal".
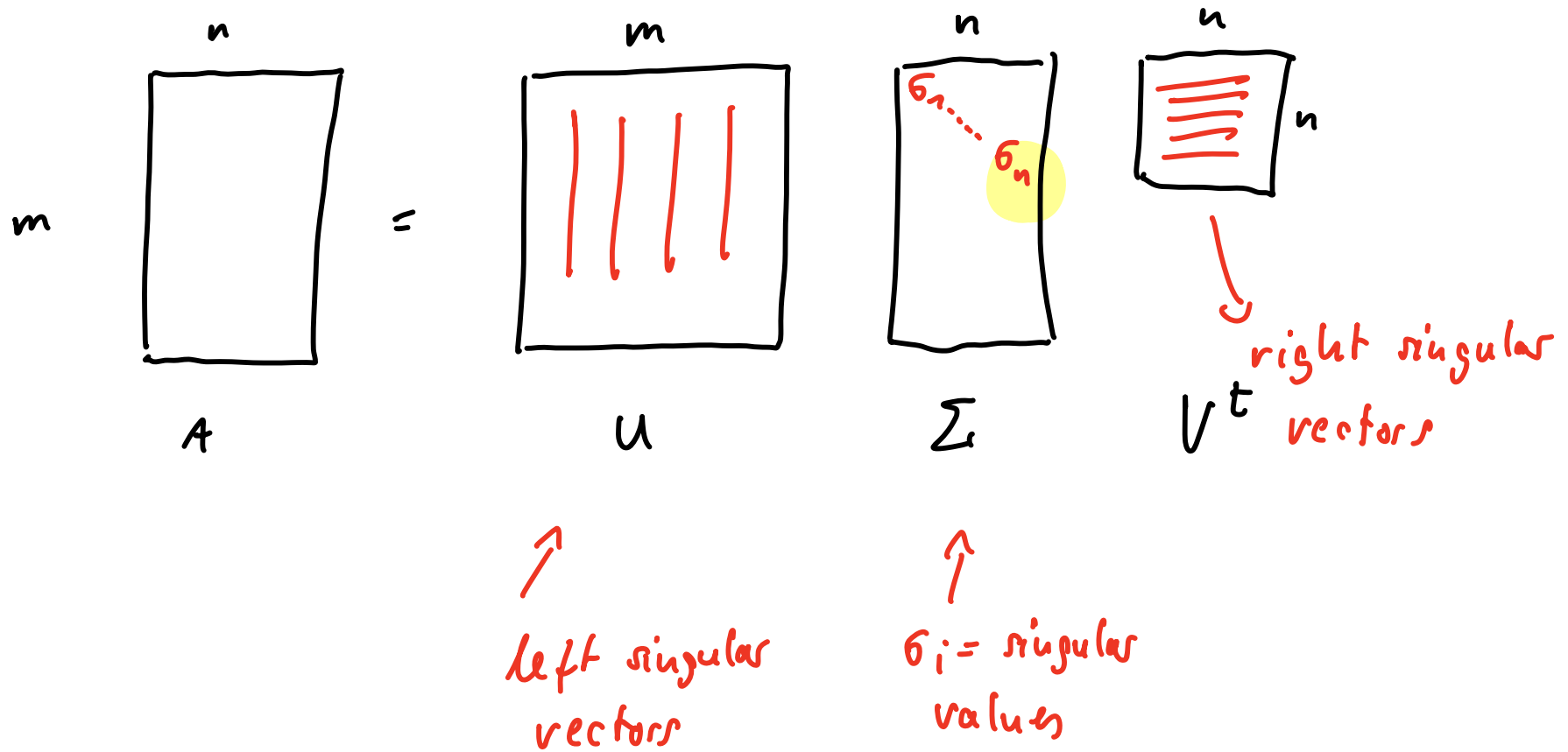
$$m \begin{pmatrix} \begin{matrix} \sigma_1 & & 0 \\ & \ddots & \\ 0 & & \sigma_n \end{matrix} \\ 0 \end{pmatrix} \overset{n}{} \qquad m \begin{pmatrix} \begin{matrix} \sigma_1 & & 0 \\ & \ddots & \\ 0 & & \sigma_m \end{matrix} & 0 \end{pmatrix} \overset{n}{}$$

Exactly $r$ of the diagonal values $\sigma_1, \sigma_2, \ldots$ are non-zero.

# Illustration

$$A = U \Sigma V^t$$



$n$

$m$

$A$

$m$

$U$

$n$

$\sigma_1 \cdots$

$\sigma_n$

$\Sigma$

$n$

$n$

$V^t$

right singular vectors

↑ left singular vectors

↑ $\sigma_i =$ singular values

# Proof sketch

Given $A \in \mathbb{R}^{m \times n}$, we consider $\underline{B} := \underbrace{A^t}_{n \times m} \underbrace{A}_{m \times n} \in \mathbb{R}^{n \times n}$.

observe:
- $B$ is symmetric:
$$(A^t A)^t = A^t (A^t)^t = A^t A$$

- $B$ is positive semi-definite:

$$x^t B x = \langle x, Bx \rangle = \langle x, A^t A x \rangle$$

$$\langle (A^t)^t x, Ax \rangle$$

$$= \langle Ax, Ax \rangle$$

$$= \| Ax \|^2 \geq 0$$

$\langle x, Cy \rangle$

$= \langle C^t x, Y \rangle$

So there exists an orthonormal basis of eigenvectors $v_1, \ldots, v_n$ with eigenvalues $\lambda_1, \ldots, \lambda_n \geq 0$.

Define:

- $\Sigma = \begin{pmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_n \end{pmatrix} \in \mathbb{R}^{m \times n}$ where $\sigma_i = \sqrt{\lambda_i}$

- $U = \begin{pmatrix} | \\ u_i \\ | \end{pmatrix}$ matrix with columns $u_i := \dfrac{A v_i}{\sigma_i}$

- $V = \begin{pmatrix} | \\ v_i \\ | \end{pmatrix}$ matrix with $v_i$ as columns

Now we need to show that with these definitions we have $A = U \cdot \Sigma \cdot V^t$:

- Columns of $U \cdot \Sigma$ are given as

$$\sigma_i \cdot u_i = \sigma_i \frac{A v_i}{\sigma_i} = A v_i$$

- Now multiply with $V^t$:

  - rows of $V^t$ are the $v_i$,

  - exploit that if $i \neq j$ then $v_i \perp v_j$ and $\| v_i \| = 1$.

  - The terms consisting of $i, j$ with $i \neq j$ cancel, the terms with $i = j$ will result in a factor of $1$.

  So we will be left with matrix $A$.

Finally, it is easy to see that $U, V$ are orthogonal matrices, and that the number of non-zero entries in $\Sigma$ coincides with the rank.

# Basic properties of SVD

- The rank of a matrix coincides with the number of non-zero singular values

- If the matrix $A$ has rank $r$, then

$$\ker(A) = \text{span}\left\{v_{r+1}, \ldots, v_n\right\}$$

$$\text{range}(A) = \text{span}\left\{u_1, \ldots, u_r\right\}$$

Proof: Exercise

# Key differences between SVD (A) and eig (A)

- SVD always exists, no matter how A looks like!
  can be rectangular, does not need to be symmetric,...

- $U, V$ are orthonormal! (not true for eigenvectors in general).

- singular values are always real and non-negative.

# Key differences between SVD $(A)$ and eig $(A)$

If $A \in \mathbb{R}^{n \times n}$ is <u>symmetric</u>, then the SVD is "nearly the same" as the eigenvalue decomposition: $(\lambda_i, v_i)$ are the eigenvalues/vectors of $A$, then

$(|\lambda_i|, v_i)$ are the singular values / vectors of $A$.

In particular, left- and right singular vectors are the same (up to signs).

# Relationship SVD($A$) and eig($\underline{AA^t}$)

symmetric!

- For general (not nec. square) matrices $A$:

  Left-singular vectors of $A$ are the eigenvectors of $AA^t$.
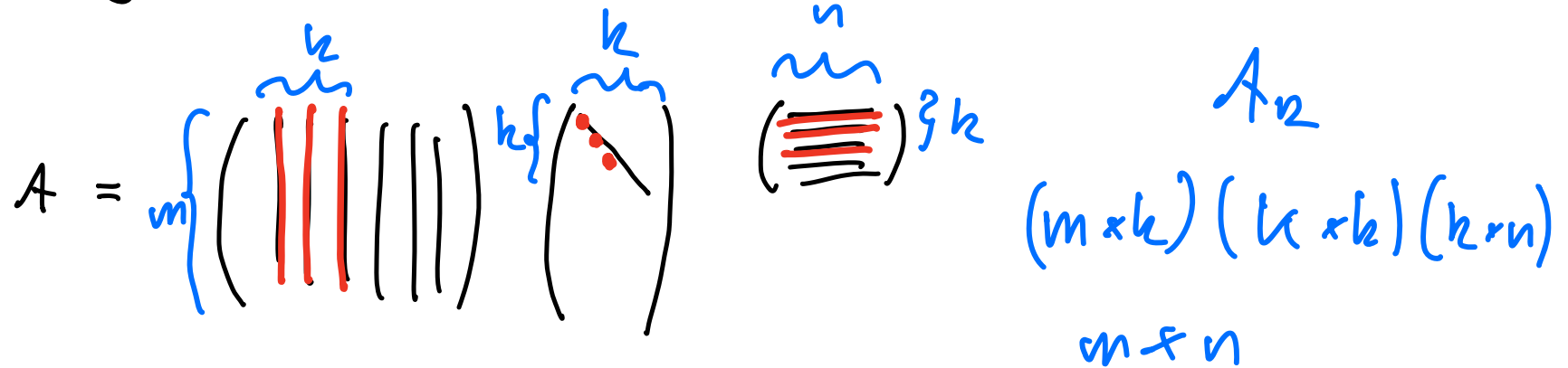
  Right- $\hspace{10cm} A^tA$.

- $\lambda \neq 0$ is an eigenvalue of $AA^t$ $\iff$

  $\sqrt{|\lambda|} \neq 0$ is singular value of $A$

# Rank-k-approximation

Given SVD $A = U \Sigma V^t$, entries $\sigma_1, \sigma_2, \ldots$ sorted in decreasing order, $k \in \mathbb{N}$. Now we are going to define a new matrix $A_k$ by the following procedure:

$$A = m\left\{ \overbrace{\begin{pmatrix} | | | & | | \\ | | | & | | \end{pmatrix}}^{k} {\scriptstyle k}\left\{ \overbrace{\begin{pmatrix} \ddots \\ & \ddots \end{pmatrix}}^{k} \overbrace{\begin{pmatrix} \equiv \\ \equiv \end{pmatrix}}^{n} \right\}k \right.$$

$$A_k$$
$$(m \times k)(k \times k)(k \times n)$$
$$m \neq n$$

- take first $k$ cols of $U$,
  first $k$ entries of $\Sigma$
  first $k$ rows of $V^t$
$\left. \right\}$ $A_k$

More formally:

$$A_k = \sum_{i=1}^{k} \sigma_i u_i v_i^t$$

Observe: $\mathrm{rank}(A_k) = k$.

# Best rank-k approximation

**Theorem** (Eckart - Young - Mirsky)

Let $\|\cdot\|$ be either the Frobenius-norm $\|\cdot\|_F$ or the two-norm $\|\cdot\|_2$.

Consider a matrix $A \in \mathbb{R}^{m \times n}$, $A_k$ the low-rank-matrix constructed above, and $B \in \mathbb{R}^{m \times n}$ any other rank-k matrix.

Then $\|A - A_k\| \leq \|A - B\|$.

In particular,

$$\|A - A_k\| = \begin{cases} \sigma_{k+1} & \text{in case of } \|\cdot\|_2 \\ \left( \sum_{i=k+1}^{\min(m,n)} \sigma_i^2 \right)^{1/2} & \text{in case } \|\cdot\|_F \end{cases}$$

Proof: skipped.

# SVD and matrix norms

Consider $A \in \mathbb{R}^{m \times n}$ with singular values $\sigma_1 \geq \sigma_2 \geq \ldots \geq \sigma_p$.

$(p = \min\{n, m\})$.

Then:

- $\|A\|_F = \sigma_1^2 + \ldots + \sigma_p^2$

- $\|A\|_2 = \sigma_1$

Proof: skipped

# Relation to machine learning

Computational reasons for a rank-$k$-approximation:

The running time of many ML algorithms scales heavily in the (implicit) dim. Often they can be implemented efficiently if matrices are (sparse or) low rank.

Statistical reason: ML only works if the data "is simple".

Typical assumptions:
- Lives on a low-dim manifold ($\rightsquigarrow$ locally low rank)

( • is sparse )

Replacing a matrix with its SVD is supposed to get rid of noise.

$$\boxed{\text{Pseudo-inverse}}$$

# Pseudo-inverse

**Definition**   for $A \in \mathbb{R}^{m \times n}$, a pseudo-inverse of $A$ is defined as the matrix $A^{\#} \in \mathbb{R}^{n \times m}$ which satisfies the following conditions:

<span style="color:red">If $A$ would be invertible</span>
<span style="color:red">$A A^{-1} = Id \Rightarrow A A^{-1} A = A$</span>

<span style="color:red">$\neq Id$ in general</span>

$(1) \quad A A^{\#} A = A$

<span style="color:red">$\neq Id$ in general</span>

$(2) \quad A^{\#} A A^{\#} = A^{\#}$

$\left. \begin{array}{} \end{array} \right\}$ "nearly inverse"

$(3) \quad \left( A A^{\#} \right)^{t} = A A^{\#}$

$(4) \quad \left( A^{\#} A \right)^{t} = A^{\#} A$

$\left. \begin{array}{} \end{array} \right\}$ symmetry

# Intuition:

- Consider a projection from $\mathbb{R}^3 \to \mathbb{R}^2$,

$$A \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$

- Cannot invert, obviously (inverting would mean to reconstruct the original point).

- But I could "invent" a reconstruction, for example: $R : \mathbb{R}^2 \to \mathbb{R}^3$,

$$R \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} x_1 \\ x_2 \\ 17 \end{pmatrix}$$

- Now we have: $A \underset{\shortparallel}{R} A = A$

$$A \overset{\#}{A} A = A$$

**Intuition: how could we define a pseudo-inverse?**

Consider a diagonal matrix $\quad A = \begin{pmatrix} \lambda_1 & & & \\ & \ddots & & \\ & & \lambda_k & 0 \dots \\ & & & & 0 \end{pmatrix}$

with $\lambda_1, \dots, \lambda_k \neq 0$.

Then set $\quad A^\# := \begin{pmatrix} 1/\lambda_1 & & & \\ & \ddots & & \\ & & 1/\lambda_k & 0 \dots \\ & & & & 0 \end{pmatrix}$

Does the job.

How to do it for general matrices? $\leadsto$ SVD!

# Moore-Penrose pseudoinverse

Proposition: Let $A \in \mathbb{R}^{m \times n}$, $A = U \Sigma V^t$ its SVD. Then the following

matrix is a pseudo-inverse:

$$A^{\#} := V \Sigma^{\#} U^t \quad \text{with} \quad \Sigma^{\#} \in \mathbb{R}^{m \times n}$$

$$\Sigma^{\#}_{ii} = \begin{cases} 1/\Sigma_{ii} & \text{if } \Sigma_{ii} \neq 0 \\ 0 & \text{otherwise} \end{cases}$$

Proof: easy, just do it.

**Proposition:** If $A$ is invertible, then $A^{-1} = A^{\#}$.

**Proof:** eory, just do it.

# Trace of a matrix

# Trace

<u>Def</u> The <mark>trace</mark> of a square matrix $A \in F^{n \times n}$ is the sum of its diagonal elements:

$$\text{tr}(A) = \sum_{i=1}^{n} a_{ii} .$$

# Properties

- $\text{tr}: \mathbb{R}^{n \times n} \to \mathbb{R}$ is a linear operator

  In particular, $\quad \text{tr}(A+B) = \text{tr}(A) + \text{tr}(B)$.

- $\text{tr}(A \cdot B) = \text{tr}(B \cdot A)$

  $\triangle! \quad \text{tr}(A \cdot B) \neq \text{tr}(A) \cdot \text{tr}(B)$

- trace does <u>not</u> depend on the basis:

  Let $T \in \mathcal{L}(V)$, and $\mathcal{B}_1$ and $\mathcal{B}$ two bases of $V$. Then

  $\text{tr}(M(T, \mathcal{B}_1)) = \text{tr}(M(T, \mathcal{B}_2))$.

- The trace of an operator equals the sum of its complex eigenvalues, summed according to multiplicity:

$$\tilde{A} = \begin{pmatrix} \lambda_1 & & * \\ & \ddots & \\ 0 & & \lambda_n \end{pmatrix} \quad \text{wrt some basis } v_1, \ldots, v_n$$

$$\Rightarrow \text{tr}(\tilde{A}) = \sum_{i=1}^{n} \lambda_i$$

- trace equals the negative of the coefficient in front of $t^{n-1}$ in the char. polynomial

$$p_A(t) = t^n + \boxed{a_{n-1}} t^{n-1} + \ldots + f$$

# Trace vs determinant

$\text{tr}(A) = \underline{\text{sum}}$ of eigenvalues

$\det(A) = \underline{\text{product}}$ of eigenvalues

# Riddle

Consider a real-valued matrix $A \in \mathbb{R}^{n \times n}$.

Over $\mathbb{C}$, we can always find eigenvalues $\lambda_1, \ldots, \lambda_n \in \mathbb{C}$ and bring the matrix in triangular form, $\tilde{A}$. Then:

$$\text{tr}(A) = \sum_{i=1}^{n} a_{ii} \in \mathbb{R}$$

$$\text{tr}(\tilde{A}) = \underbrace{\sum \lambda_i}_{\in \mathbb{C}} \in \mathbb{C}$$

$\rangle$ but the two are identical because trace is independent of basis, hence $\text{tr}(\tilde{A}) \in \mathbb{R}$ ! ?!?

So even over $\mathbb{C}$, the trace always is a real number. Seems confusing...

Let's look at an example:

**Example :** Consider a rotation matrix

$$A = \begin{pmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{pmatrix}$$

- $A$ does not have any real eigenvalues.

- The trace is given as $2 \cdot \cos\theta$.

- The char. poly. of $A$ is

$$p(t) = \det(A - tI) = \det\begin{pmatrix} (\cos\theta) - t & -\sin\theta \\ \sin\theta & (\cos\theta) - t \end{pmatrix}$$

$$= (\cos\theta - t)^2 + \sin^2\theta$$

$$= t^2 - 2\cos\theta \cdot t + \underbrace{\cos^2\theta + \sin^2\theta}_{=1}$$

$$= t^2 - (2\cos\theta) \cdot t + 1$$

- The roots of the char. pol.

$$\lambda_1, \lambda_2 = \frac{2\cos\theta \pm \sqrt{(2\cos\theta)^2 - 4}}{2}$$

$4\cos^2\theta - 4$

$= 4(\cos^2\theta - 1)$

$= 4(-\sin\theta)$

eigenvalues $\longrightarrow = \cos\theta \pm i\cdot\sin\theta$

$\in \mathbb{C}$

- The matrix has a diagonal representation

$$\begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix}$$

$$\operatorname{tr}\begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix} = \cos\theta + i\sin\theta + \cos\theta - i\sin\theta$$

$$\longrightarrow = 2\cos\theta$$

Sum in $\mathbb{R}$

# What is going on?

For any matrix $A \in \mathbb{C}^{n \times n}$, if $\lambda \in \mathbb{C}$ is an eigenvalue, then also $\bar{\lambda}$ is an eigenvalue.

Reason: the complex eigenvalues have their roots in solving quadratic equations, as in the last example:

$$\lambda_{1,2} = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

- If $b^2 - 4ac > 0$, eigenvalues $\in \mathbb{R}$

- If $b^2 - 4ac < 0$, then eigenvalues are

$$\lambda_{1,2} = \frac{-b \pm i\sqrt{4ac - b^2}}{2a} = \frac{-b}{2a} \pm i\sqrt{4ac - b^2}$$

so they are complex conjugates, and $\lambda_1 + \lambda_2 \in \mathbb{R}$.

Matrix series

# Spectral radius

<u>Def</u>   The ==spectral radius== of a matrix $A \in \mathbb{R}^{n \times n}$ or $A \in \mathbb{C}^{n \times n}$ is

defined as

$$\rho(A) := \max \{ |\lambda| \; ; \; \lambda \in \mathbb{C} \text{ eigenvalue of } A \}$$

( Note that even for real matrices, this definition looks at all complex

eig' of $A$ ).

**Proposition**

$$\lim_{k \to \infty} A^k = 0 \quad \Longleftrightarrow \quad \rho(A) < 1$$

in Frobenius norm: $\|A^k - 0\|_F \to 0$

**Proof** "$\Rightarrow$" Let $v$ the eigenvector of the eig $\lambda$ that defines $\rho(A)$.

$$0 \overset{ass.}{=} \lim_{k \to \infty} A^k v \overset{eig}{=} \lim_{k \to \infty} \lambda^k v = \left( \lim_{k \to \infty} \lambda^k \right) v$$

This implies $\lim \lambda^k = 0$, thus $|\lambda| < 1$.

"$\Leftarrow$" Proof for symmetric matrices:

$$A = U D U^t, \quad \text{then} \quad A^k = U D^k U^t.$$

$$\|A\|_F = \|D\|_F \to 0.$$

For general matrices, need to exploit another normal form, the Jordan normal form. Skipped.

# Neumann series

**Prop** . The series $\sum_{i=0}^{\infty} A^k$ converges if and only if $\rho(A) < 1$.

- If $\rho(A) < 1$, then $(I - A)$ is invertible and

$$(I - A)^{-1} = \sum_{i=0}^{\infty} A^k$$

over $\mathbb{R}$:

$$\frac{1}{1-a} = \sum_{k=0}^{\infty} a^k$$

**Proof intuition:** for symmetric matrices easy, can apply matrix powers to diagonal matrix $D$ as above.

In general a bit more work, skipped.

# Matrix exponential

For any matrix $A$ the series

$$\exp(A) = \sum_{n=0}^{\infty} \frac{A^n}{n!} = I + A + \frac{A^2}{2} + \ldots$$

converges. It is called the matrix exponential.
The matrix $\exp(A)$ is always invertible and

$$\left(\exp(A)\right)^{-1} = \exp(-A)$$

# Implementing matrix operations

Literature: Golub/van Loan: Matrix computations

# Triangular matrices are great

If we have a matrix in triangular form, then many standard quantities can be easily computed:

- Solving a linear system : $Ax = b$

$$\begin{pmatrix} a_{11} & a_{12} & a_{13} \\ 0 & a_{22} & a_{23} \\ 0 & 0 & a_{33} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ b_3 \end{pmatrix}$$

Starting from bottom: $x_3 = b_3 / a_{33}$ !

Then plug in the second-to-last row, $x_2 = (b_2 - a_{23} x_3) / a_{22}$

...

- eigenvalues are the entries on the diagonal

- det $(A)$ is the product of the diagonal entries

So many numerical algorithms are based on triangular matrices.

# Linear system: theory

Solving linear systems is about the most basic task and occurs as a substep of more complex algorithms all over the place.

$$Ax = b$$

**Prop** $A \in \mathbb{R}^{m \times n}$, $b \in \mathbb{R}^m$.

(i) The set of solutions of $Ax = 0$ is given by her$(A)$ and is a subspace of $\mathbb{R}^m$.

(ii) The system $Ax = b$ has at least one solution iff $b \in \text{Range}(A)$

(iii) If $w$ is a solution of $Ax = b$, then the full set of solutions is given by $w + \text{her}(A) := \{w + v \mid v \in \text{her}(A)\}$

# Gauss algorithm to solve linear systems

Intuition:
- Take the first equation and use it to eliminate the first variable in all other equations.
- Then use the second equation to eliminate the second variable from all the remaining equations.
- and so on.
- At the end, we are left with an upper triangular matrix
- We then solve the system starting from the bottom...

Of course one can do clever tricks such as selecting the best row for replacing others (pivoting, rearranging, ...).

Computational complexity in general case: $O(n^3)$

# LU decomposition

Idea: One can see that the Gaussian elimination algorithm implicitly does something more general.

Given a square matrix $A$, it decomposes $A$ into a product

$$A = L \cdot U$$

where $L$ is a lower triangular matrix and $U$ an upper triangular matrix.

LU decomposition exists under certain conditions. In particular for symmetric pd matrices it always exists.

Computational complexity: $O(n^3)$

# LU decomposition to solve a linear system

$O(n^3)$    Once we have a LU decomposition,

the solution to the problem $Ax = b$ can then be found

by solving two linear systems (trivial due to triang. form):

$O(n^2)$        (1)      $Ly = b$

$O(n^2)$        (2)      $Ux = y$

$\Rightarrow Ax = \underbrace{LUx}_{y} \overset{(2)}{=} Ly \overset{(1)}{=} b.$

# LU decomposition to compute the inverse

To compute the inverse of a matrix, we have to find a matrix $X$ that solves $A \cdot X = I$. This corresponds to solving $n$ linear systems:
$$A x_i = e_i$$

Once we have the $LU$-decomposition of $A$, we can simply solve them.

$$O(n^3)$$

# Cholesky factorization

For the special case of symmetric, psd matrices, one can simplify the LU decomposition:

$$A = L L^t$$

Is numerically very stable, uses less memory than LU, and is a bit faster than LU ( still $O(n^3)$, but with better constants)

# QR decomposition

Any matrix $A \in \mathbb{R}^{n \times n}$ can be decomposed as

$$A = QR$$

where $Q$ is an orthogonal matrix and $R$ is upper triangular.

Several algorithms exist, with different advantages / disadvantages.

Computational complexity: $O(n^3)$

# QR decomposition to find orthogonal basis

In particular, if $A$ has full rank (i.e., its columns form a basis of $\mathbb{R}^n$), then $Q$ contains an orthonormal basis of $\mathbb{R}^n$.

More generally, the first $k$ columns of $Q$ form an orthonormal basis of the subspace spanned by the first columns of $A$.

# QR decomposition to compute _all_ eigenvalues of _dense_ matrices A

Iterative procedure:

- $A^{(0)} =: A$

- For $k = 1, 2, \ldots$

    - Compute QR factorization of $A^{(k-1)}$:

    $$A^{(k-1)} = Q^{(k)} R^{(k)}$$

    - Recombine in reversed order:

    $$A^{(k)} := R^{(k)} Q^{(k)}$$

Note: $A_{k+1} = R_k Q_k = (Q_k^{-1} Q_k) R_k Q_k = Q_k^{-1} (Q_k R_k) Q_k = Q_k^{-1} A_k Q_k$

Under certain assumptions, $A^k$ converges to a triangular matrix or even to a diagonal matrix (eg if $A$ is symmetric).

# Largest eigenvalue and eigenvector

Let $A \in \mathbb{R}^{n \times n}$ with eigenvalues $\lambda_1, \ldots, \lambda_m$ such that
$$|\lambda_1| > |\lambda_2| \quad \ldots \quad > |\lambda_m|.$$

## Power method (vanilla version)

- start with a random vector $x_0$

- iteratively compute
$$v_{k+1} = \frac{A \cdot v_k}{\|A v_k\|}$$

- If $|\lambda_1| > |\lambda_2|$, then $v_k$ converges to the eigenvector $v_1$.
  The speed of convergence depends on the spectral gap $|\lambda_2| / |\lambda_1|$.

- Problematic if eigenspace has dim $> 1$, or if it is unknown whether $A$ has an eigenvalue in the first place.

# Iterative methods for sparse matrices

Many of the algorithms we have seen (LU, QR, ...) cannot really exploit ==sparsity== of a matrix. Bad, in ML many matrices are very sparse.

Alternatively, one uses ==iterative methods== that are based on matrix-vector-multiplications (example: power method)

# Conjugate gradient method for linear systems of sparse symmetric pd matrices

- want to solve $Ax = b$
- Consider the minimization problem $\min \Phi(x)$ with

$$\Phi(x) = \frac{1}{2} x^t A x - x^t b.$$

Minimum is achieved by setting $x := A^{-1} b$.

- So we can find the solution $x$ to our system $Ax = b$ by minimizing $\Phi$.

- The gradient is $\nabla \Phi(x) = Ax - b$ and can be computed just with a matrix-vector product (good for sparsity). ...

- Now apply optimization methods (conjugate gradient descent).

# Skipped material

Not treated this year, but the videos still exist if you are interested.

# Quotient spaces

<u>Def</u>  Consider a set $S$. A subset $R \subset S \times S$ is called an equivalence relation on $S$ if $\forall x, y, z \in S$:

(E1)    $(x, x) \in R$           (reflexivity)

(E2)    $(x, y) \in R \Rightarrow (y, x) \in R$        (symmetry)

(E3)     $(x, y) \in R$, $(y, z) \in R \Rightarrow (x, z) \in R$  (transitivity)

Notation:   $(x, y) \in R \iff x \sim y$

**Example** $V$ $VS_1$ $W \subseteq V$ subspace.

$$v \sim u \quad :\Longleftrightarrow \quad v - u \in W$$

**Example:** Consider the space $\mathcal{L}(\mathbb{R})$ of all functions $f: \mathbb{R} \to \mathbb{R}$ that are Lebesgue integrable. Define

$$f \sim g \quad :\Longleftrightarrow \quad f = g \text{ almost everywhere}$$

<u>Def</u>    The equivalence class of an element $a \in S$
under equivalence relation $\sim$ is defined as

$$[a] := \{ b \in S \mid b \sim a \}$$

<u>Prop</u>    Two equivalence classes $[a]$ and $[b]$ are either
identical or disjoint.


<u>Consequence</u> :    An equ. relation on $S$ results in a disjoint
partition of equivalence classes.

# Constructing quotient spaces:

$V$ $VS$, $W \subset V$ subspace, equivalence relation

$$v \sim u \iff v - u \in W$$

Denote the equivalence classes as $[v]$.

Observe: the equ. classes have the form

$$[v] = v + W = \{u \in S \mid \exists w \in W : u = v + w\}$$

$$= \{v + w \mid w \in W\} \subset V$$

Define the ==quotient "space"== as

$$V/W := \{ [v] \mid v \in V \}$$

$$[v], [u] \in V/W$$

$$[v] + [u] \quad :\Longleftrightarrow \quad [v+u]$$

$$\lambda [v] \quad :\Longleftrightarrow \quad [\lambda v]$$

These operations are <u>well-defined</u>:

- suppose $v' \sim v$ $\quad$ (i.e. $v' \in [v]$, $[v] = [v']$)

$$u' \sim u$$

$$[v] + [u] \overset{?}{=} [v'] + [u']$$

$$v \sim v' \iff \exists\, w \in W \quad v - v' = w$$

$$u \sim u' \iff \exists\, \tilde{w} \in W: \quad u - u' = \tilde{w}$$

$$[v] + [u] = [v+u]$$
$$[v'] + [u'] = [v'+u']$$
$\big)?$  $\quad$ $(v+u) \sim (v'+u')$

$$(v+u) - (v'+u')$$

$$= \underbrace{(v-v')}_{w} + \underbrace{(u-u')}_{\tilde{w}} \in W$$

- similarly, for scalar mult.

$$\left(V/W, +, \cdot\right) \text{ is a vector space: exercise.}$$

**Prop :** Consider $g: V \to V/W$, $v \mapsto [v]$. Then:

- $g$ is linear

- $\ker(g) = W$

- $\text{range}(g) = V/W$

- If $V$ has finite dim, then $\dim V/W = \dim V - \dim W$.

# Characteristic polynomial

Motivation :

$$Av = \lambda v$$

$A$ $n \times n$ - matrix
$v \neq 0$

$$\Rightarrow \quad (A - \lambda I) v = 0$$

$$\Leftrightarrow \quad v \in \ker (A - \lambda I)$$

$$\Leftrightarrow \quad \text{rank} (A - \lambda I) < n$$

$$\Leftrightarrow \quad \det (A - \lambda I) = 0$$

**Def** The ==characteristic polynomial== of an $n \times n$-matrix $A$ is defined as

$$P_A(t) := \det(A - t \cdot I)$$

Example:

$$A = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix}$$

$$\det(A - t \cdot I) = \det\left( \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} - t \cdot \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \right)$$

$$= \det \begin{pmatrix} a_{11} - t & a_{12} \\ a_{21} & a_{22} - t \end{pmatrix}$$

$$= (a_{11} - t)(a_{22} - t) - a_{12} \cdot a_{21}$$

$$= t^2 + t(-a_{11} - a_{22}) - a_{12} \cdot a_{21} + a_{11} \cdot a_{22}$$

# Observations

- $p_A(t)$ is a polynomial with degree $n$

- Char. pol. does **not** depend on the basis:

  **Proof** Consider $A$, basis transformation matrix $U$.
  Want to look at char. pol. of $UAU^{-1}$.

  $$\det\left(UAU^{-1} - t\cdot I\right)$$

  $$= \det\left(UAU^{-1} - t\cdot U\cdot U^{-1}\right)$$

  $$= \det\left(U\,(A - t\cdot I)\,U^{-1}\right)$$

  $$= \det(U)\cdot \det(A - t\cdot I)\cdot \det(U^{-1})$$

  $$= \det(A - tI)$$

- The roots of the characteristic poly. correspond exactly to the eigenvalues of $A$.

- Over $\mathbb{C}$, the char. poly. always has $n$ roots, so the matrix has "$n$ eigenvalues" (not nec. distinct).

- $A$ is invertible $\Leftrightarrow$ $0$ is <u>not</u> an eigenvalue.

$$\text{If } 0 \text{ is an eigenvalue, } u.v \text{ with}$$
$$A v = 0 \cdot v = 0$$
$$\Leftrightarrow \ker(A) \text{ non-trivial} \Leftrightarrow A \text{ not invertible}$$

- Let $A \in \mathcal{L}(V)$, $\lambda$ eig. of $A$. Then $\lambda^k$ is an eig. of $A^k$.

- Let $A$ be invertible; $\lambda$ eig of $A$. Then $1/\lambda$ is an eig. of $A^{-1}$.

<u>Def</u>    For an operator $A$ with eigenvalue $\lambda$, we define its
**geometric multiplicity** as the dimension of the
corr. eigenspace $E(\lambda, A)$.

The **algebraic multiplicity** is the multiplicity of the
root $\lambda$ in the char. poly.


In general, the two notions <u>do not</u> coincide.

# Computing eigs in theory

- Write down the char. pol., find the roots.
  $\rightsquigarrow$ eigenvalues

- To compute the eigenvectors, solve the lin. system
  $$A x = \lambda x$$

In practice ... see later (numerics)
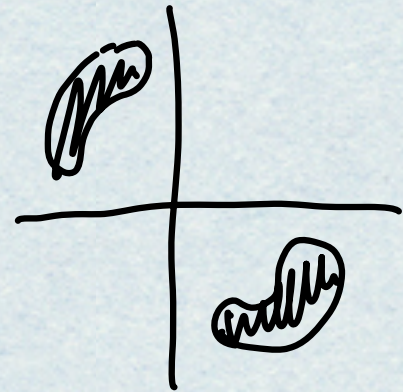
Norms can be characterized
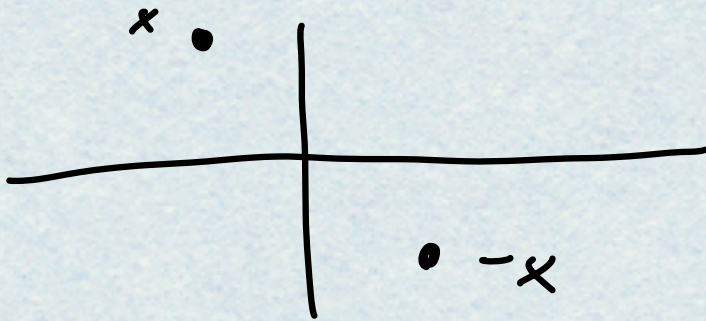by convex sets

# Convex set

**Def** Consider a real VS $V$, $S \subset V$. $S$ is called
<mark>convex</mark> if $\forall t$, $0 \le t \le 1$ and $\forall x, y \in S$,
$$t \cdot x + (1-t) y \in S \leftarrow$$

## Intuition

# Symmetric set

**Def**   A set $C \subset V$ is called ==symmetric== if

$$x \in C \implies -x \in C$$

# Convex sets induce norms

**Theorem:**

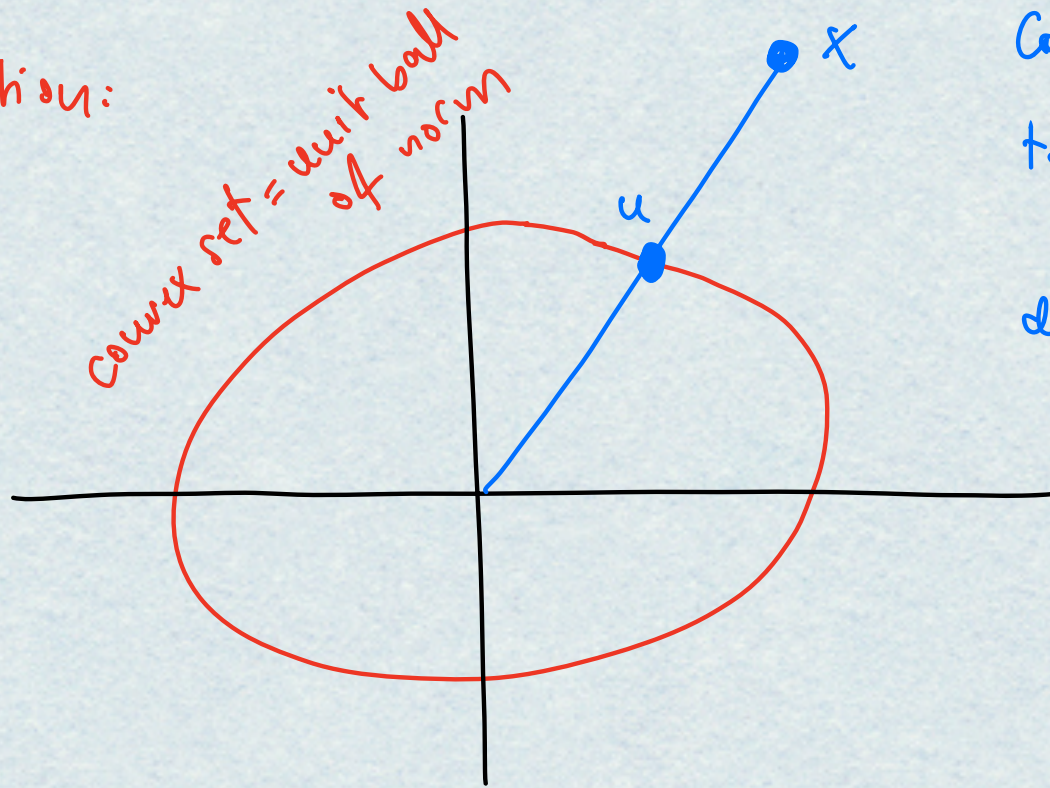(1) Let $C \subset \mathbb{R}^d$ closed, convex, symmetric and has non-empty interior. Define

$$p(x) := \inf\left\{ t > 0 \;\Big|\; \frac{x}{t} \in C \right\}. \quad \text{Then}$$

$$= \inf\{ t > 0 \mid x \in t \cdot C \}, \quad \text{this might be more intuitive.}$$

$p$ is a semi-norm. If $C$ is bounded, then $p$ is a norm, and its unit ball coincides with $C$:

$$C = \left\{ x \in \mathbb{R}^d \;\Big|\; p(x) \leq 1 \right\}$$

(2) For any norm $\|\cdot\|$ on $\mathbb{R}^d$, the set $C := \{x \in \mathbb{R}^d \mid \|x\| \leq 1\}$ is bounded, symmetric, closed, convex, and has non-empty interior.

Illustration: part (1)

convex set = unit ball of norm

Consider factor $a$ by which I need to multiply $x$ to end up on the unit sphere. Then define $\|x\| := 1/a$.
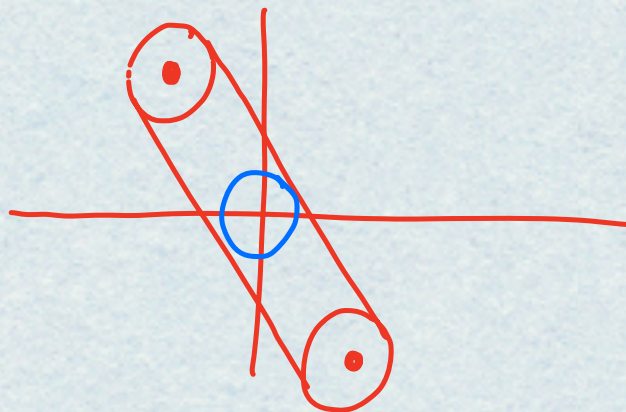
# Proof :

The next couple of pages treat the proof of this theorem.

Proof $\boxed{p(x) \text{ is well defined}}$

Want to prove: given $x \in \mathbb{R}^d$, the set $\left\{ t > 0 \mid \frac{x}{t} \in C \right\} \neq \emptyset$

We are going to prove: $\exists \, \varepsilon > 0$ such that
$$B_\varepsilon(0) = \left\{ e \in \mathbb{R}^a \mid \|e\|_2 < \varepsilon \right\} \quad \subset \quad C.$$

Intuition:

- By ass, $C$ has at least one interior point

  $v \in C^o \implies \exists \varepsilon$ such that

  $B_\varepsilon (v) \subset C$

  "$v + B_\varepsilon (0) = \{ v + e \mid e \in B_\varepsilon (0) \}$

- By symmetry, $v + e \in C \implies -(v + e) \in C$

- By convexity, $\frac{1}{2}(v+e) + \frac{1}{2}(-v-e) = e \in C$

  So $B_\varepsilon (0) \subset C$, so the set $\{ t > 0 \mid \frac{x}{t} \in C \}$

  is non-empty.

The infimum of $\overbrace{\inf \{ t > 0 \mid \frac{x}{t} \in C \}}^{S_x}$ exists

because $S_x \subset \mathbb{R}$, $0$ is a lower bound.

Now we need to prove all axioms of a norm:

$$p(0) = 0$$

- Have seen: $0 \in C$

- $\forall t > 0: \quad \dfrac{0}{t} = 0 \in C$

- $\inf \left\{ t \mid \dfrac{0}{t} \in C \right\} = 0.$

$$\Rightarrow p(0) = 0$$

$$\boxed{p(\alpha x) = |\alpha|\, p(x)}$$

- For all $\alpha > 0$ we have

$$p(\alpha \cdot x) = \inf\left\{ t > 0 \;\middle|\; \frac{\alpha x}{t} \in C \right\} \xleftarrow{s := \frac{t}{\alpha}}$$

$$= \inf\left\{ \alpha \cdot s > 0 \;\middle|\; \frac{x}{s} \in C \right\}$$

$$= \alpha \cdot \underbrace{\inf\left\{ s > 0 \;\middle|\; \frac{x}{s} \in C \right\}}_{p(x)}$$

$$\Rightarrow p(\alpha x) = \alpha \cdot p(x)$$

- By symmetry we also get

$$p(-x) = \inf\left\{ t > 0 \mid \frac{-x}{t} \in C \right\} = \quad \checkmark \quad -\frac{x}{t} \in C \Rightarrow \frac{x}{t} \in C$$

$$= \inf\left\{ t > 0 \mid \frac{x}{t} \in C \right\} = p(x)$$

- Combining the two statements gives homogeneity.

$\boxed{\triangle - \text{inequality}}$ Consider $x, y \in \mathbb{R}^d$, $s, t > 0$ such that

$$\frac{x}{s} \in C, \quad \frac{y}{t} \in C.$$

Observe: $\dfrac{s}{s+t} + \dfrac{t}{s+t} = 1$. Thus, by convexity,

$$\underbrace{\frac{s}{s+t}}_{} \cdot \underbrace{\frac{x}{s}}_{\in C} + \underbrace{\frac{t}{s+t}}_{} \cdot \underbrace{\frac{y}{t}}_{C} \quad \in \quad C \quad \cancel{\otimes}$$

two scalars that sum up to 1

$$p(x+y) = \inf \left\{ u > 0 \ \bigg| \ \frac{x+y}{u} \in C \right\}$$

$$\leq \underbrace{\inf \left\{ s > 0 \ \bigg| \ \frac{x}{s} \in C \right\}}_{p(x)} + \underbrace{\inf \left\{ t > 0 \ \bigg| \ \frac{y}{t} \in C \right\}}_{p(y)}$$

$$\frac{s}{s+t} \ \frac{x}{s} + \frac{t}{s+t} \ \frac{y}{t} \ \in \ C$$

$$\frac{x+y}{s+t} \ \in \ C$$

$$\boxed{\frac{x+y}{s+t}} = u_0$$

$$p(x+y) = \inf \left\{ u > 0 \,\middle|\, \frac{x+y}{u} \in C \right\} \overset{!}{\leq} u_0$$

$$= \underset{p(x)}{\underline{s}} + \underset{p(y)}{\underline{t}}$$

$s$ was chosen such that   $\dfrac{x}{s} \in C$

$t$   $\dfrac{y}{t} \in C$

Consider a sequence $(s_i)_{i \in \mathbb{N}}$ such that

$$\frac{x}{s_i} \in C \quad \text{and} \quad s_i \to p(x)$$

Similarly $(t_i)_{i \in \mathbb{N}}$ such that $\frac{y}{t_i} \in C$ and $t_i \to p(y)$.

By the argument above, we know that

$$\forall i: \quad p(x+y) \leq s_i + t_i$$

$$\underbrace{s_i}_{\downarrow} \quad \underbrace{t_i}_{\downarrow}$$
$$p(x) \qquad p(y)$$

$$\Rightarrow \quad p(x+y) \leq p(x) + p(y).$$

$$\boxed{p(x) = 0 \implies x = 0}$$

$$p(x) = 0 \iff \inf\left\{ t \geq 0 \;\middle|\; \frac{x}{t} \in C \right\} = 0$$

$\implies$ There exists a sequence $(t_k)_{k \in \mathbb{N}}$ such that

$$t_k \to 0 \quad \text{and} \quad \frac{x}{t_k} \in C \quad \forall k.$$

Now assume that $x \neq 0$. Then the sequence

$$\left( \frac{x}{t_k} \right)_{k \in \mathbb{N}} \text{ is unbounded.} \quad \underset{\text{Contradiction because}}{\text{$\curvearrowright$}} \quad C \text{ is bounded.}$$

Spaces of functions:

(continuous, differentiable, integrable )
$\mathcal{C}^0$                  $\mathcal{C}^1$                    $L_p$

# Space of continuous functions

**Def** Let $T$ be a metric space,
$$C^b(T) := \{ f : T \to \mathbb{R} \mid f \text{ is continuous and bounded} \}$$

As norm on $C^b(T)$ we now use

$$\| f \|_\infty := \sup_{t \in T} | f(t) |$$

$\exists c \in \mathbb{R} :$
$\forall t \in T : | f(t) | < c$

**Prop** Then the space $C^b(T)$ with norm $\| \cdot \|_\infty$ is a Banach space: a complete, normed vector space.
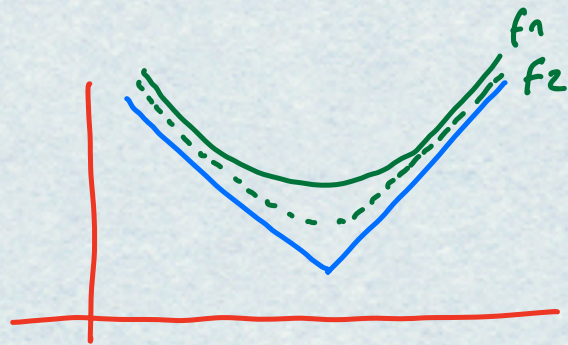
## Proof outline:

- need to check vector space axioms
- norm axioms
- completeness: follows from the fact that $\|\cdot\|_\infty$ induces uniform convergence

# Space of differentiable functions

Let $[a, b] \subset \mathbb{R}$, $\mathcal{C}^1([a, b]) = \{f : [a, b] \to \mathbb{R} \mid f$ is cont. differentiable$\}$

Which norm?

- Consider $\|\cdot\|_\infty$. With this norm, $\mathcal{C}^1$ is not complete!



limit function, not differentiable!

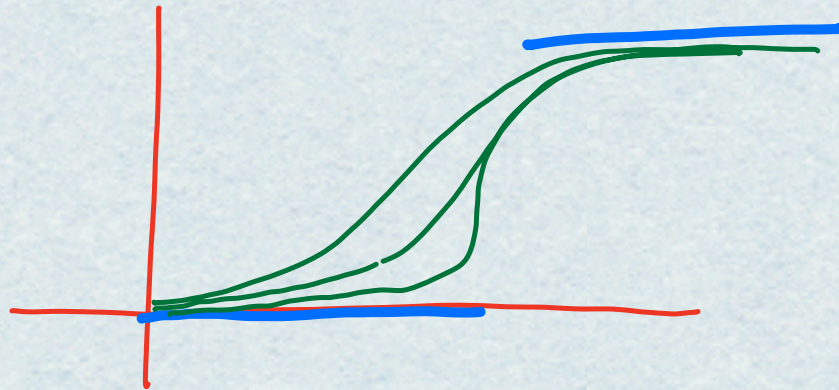- Consider $\|f\| := \sup_{t \in [a,b]} \max\{|f(t)|, |f'(t)|\}$

$$\|\|f\|\| := \|f\|_\infty + \|f'\|_\infty$$

$\mathcal{C}^1([a,b])$ with any of these two norms is a Banach space.

# Spaces of integrable functions?

- Consider $\mathcal{C}^b([a,b])$ with the norm $\|f\|_1 := \int_a^b |f(t)| \, dt$
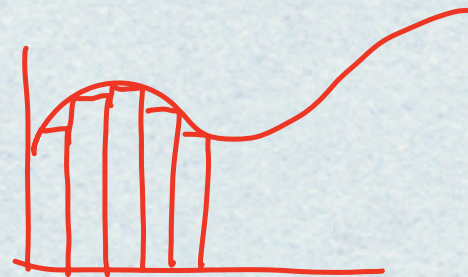
Can see: $\|\cdot\|_1$ is a norm, but the space is not complete.

- Consider $R([a,b])$ of all Riemann-integrable functions on $[a,b] \subset \mathbb{R}$, together with $\|\cdot\|_1$.

However, on $R([a,b])$, $\|\cdot\|_1$ is not a norm: it is not true that

$$\|f\| = 0 \implies f = 0$$

$\int f \, dt = 0$
but $f \neq 0$

# The space $\mathcal{L}_p$

<u>Def</u> For $1 \leq p < \infty$, we define

$$\mathcal{L}_p([a,b]) := \left\{ f : [a,b] \to \mathbb{R} \mid f \text{ measurable wrt the Lebesgue measure} \right.$$
$$\left. \text{and } \int |f(t)|^p \, dt < \infty \right\}$$

$$\|f\|_p := \left( \int |f|^p \, d\lambda \right)^{1/p}$$

Proposition 1 : $\|f\|_p$ is a semi-norm on $\mathcal{L}_p$.

Proof : • Vector space (clear)

• semi-norm : observe: $\|f\|_p = 0 \Rightarrow f = 0$ almost everywhere

so we do **not** have that $\|f\|_p = 0 \Rightarrow f = 0$.          a.e.

⚠ $\|f\|_p$ is not a norm! For example, the function

$$f(x) = \begin{cases} 1 & \text{if } x = 17 \\ 0 & \text{otherwise} \end{cases}$$

has integral 0, but is not

the 0-function.

**Proposition 2** $\mathcal{L}_p$ is complete under $\|\cdot\|_p$.

**Proof:** If $(f_i)_{i \in \mathbb{N}}$ is a Cauchy-sequence in $\mathcal{L}_p$, then want to prove that $\lim_{i \to \infty} f_i \in \mathcal{L}_p$.

This is equivalent to proving the following:

Let $(f_i)_i$ be a sequence such that

$$a := \sum_{i=1}^{\infty} \|f_i\|_p < \infty \qquad \text{⊘🏃}$$

then there exists $f \in \mathcal{L}_p$ such that $f_i \to f$ (in $\|\cdot\|_p$).

Define

$$\hat{g} := \sum_{i=1}^{\infty} |f_i|$$

N.ok: this might not yet be a well-defined fct from $[a,b]$ to $\mathbb{R}$, might be $\infty$ at certain points.

$$\hat{g}_n := \sum_{i=1}^{n} |f_i| \in \mathcal{L}_p$$

By Minkowski,

$$\|\hat{g}_n\|_p \overset{\text{def}}{=} \|\sum_{i=1}^{n} |f_i| \|_p \overset{\text{Mink.}}{\leq} \sum_{i=1}^{n} \||f_i|\|_p \overset{\text{ass. }*}{<} a$$

$$\hat{g}_n \to g \qquad \text{monotonously}$$

By theorem of monotonic convergence, $g$ is measurable

and we have

$$\lim_{n \to \infty} \int \hat{g}_n^p \, d\lambda = \int \lim \hat{g}_n^p \, d\lambda$$

$$\overset{\text{by def}}{=} \int \hat{g}^p \, d\lambda$$

$$\leq a^p.$$

$\Rightarrow$ $\hat{g} < \infty$ a.e., that is there exists a set $N$

of measure $0$ such that on $[a,b] \setminus N$, $\hat{g}$ is finite.

Now we can define

$$g(t) = \begin{cases} \hat{g}(t) & t \in [a,b] \setminus N \\ 0 & t \in N \end{cases} \qquad \in \mathcal{L}_p$$

From this it now follows that $f(t) = \sum_{i=1}^{\infty} f_i(t)$, $t \notin N$

exists. For $t \in N$, we put $f(t) = 0$.

Now $f$ is measurable, and in $\mathcal{L}_p$

because $\int |f|^p \, d\lambda \leq \int \hat{g}^p \, d\lambda < \infty$

Finally, $\sum_{n=1}^{\infty} f_n$ converges to $f$ in $\|\cdot\|_p$

because of the theorem of dominated convergence.

# From $\mathcal{L}_p$ to $L_p$

We constructed a space $\mathcal{L}_p$ with the Lebesgue integral as a semi-norm. This means, given $f \in \mathcal{L}_p$, we can change the $p$ values of $f$ in a set of measure $0$, resulting in $\tilde{f}$, but the norm "does not see a difference":

$$\| f - \tilde{f} \| = 0$$

To fix this, we want to consider functions to be "equivalent" if they only differ by a set of measure $0$.

The formal construction goes as follows:

Define $N := \ker(\|\cdot\|_p) := \{ f \in \mathcal{L}_p \mid \|f\|_p = 0 \}$

is a subspace of $\mathcal{L}_p$.

Now consider the quotient space of $\mathcal{L}_p$ wrt this subspace:

$$ L_p([a,b]) := \mathcal{L}_p([a,b]) / N $$

Define a norm on $L_p$ by $\|\underbrace{[f]}_{\text{new}}\|_p = \underbrace{\|f\|_p}_{\text{old}}$

This norm is well-defined: if $f, \tilde{f} \in [f]$,

then $\|f\|_p = \|\tilde{f}\|_p$.

This "norm" is a norm, because

$$\|[f]\|_p = 0 \implies [f] = [0].$$

Conclusion: $L_p$ with $\|\cdot\|_p$ is a Banach space!

For simplicity, in future we write $\|f\|_p$ for $\|[f]\|_p$.

Elements ("functions") in $L_p$ are equivalence classes $[f]$ consisting of all functions that coincide a.e.

⚠ • It does not make sense to evaluate $f(0)$ because $\{0\}$ has Lebesgue measure 0.

• quite annoying for machine learning, where we always want to evaluate functions on input points.

• often use alternative spaces instead, for example reproducing kernel Hilbert spaces.

Operator norm

# Continuous = bounded

**Theorem** $X, Y$ normed spaces, $T: X \to Y$ linear. Then the following statements are equivalent:

(i) $T$ is continuous.  $\longrightarrow$ $\begin{cases} \forall x \in X \ \forall \varepsilon > 0 \ \exists \delta > 0 \ \forall y \in X: \\ \|x-y\| < \delta \Rightarrow \|Tx - Ty\| < \varepsilon \end{cases}$
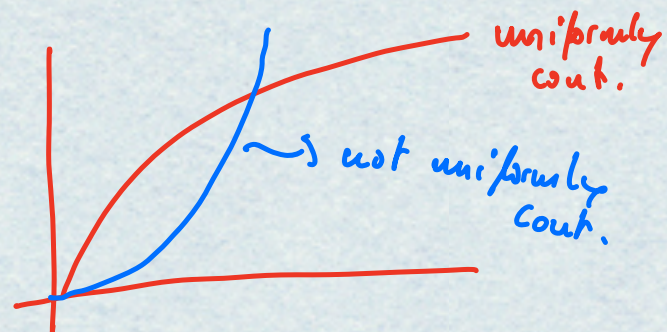
(ii) $T$ is continuous at $0$.

(iii) $T$ is bounded:

$$\exists M > 0 \ \forall x \in X: \|Tx\| \leq M \cdot \|x\|$$

(iv) $T$ is uniformly continuous.

$$\forall \varepsilon > 0 \ \exists \delta > 0 \ \forall x \in X \ \forall y \in X:$$
$$\|x-y\| < \delta \Rightarrow \|Tx - Ty\| < \varepsilon$$



uniformly cont.

$\longrightarrow$ not uniformly cont.

# Definition

**Def** $X, Y$ normed spaces, $T : X \to Y$ linear and continuous.

$$\|T\| := \sup_{x \in X} \frac{\|Tx\|}{\|x\|} = \sup_{\substack{x \in X \\ \|x\| \leq 1}} \|Tx\| = \sup_{\substack{x \in X \\ \|x\| = 1}} \|Tx\|$$

is called the ==operator norm== of $T$.

*Observe: coincides with the matrix norm*
*$\|\cdot\|_2$ as we had defined it earlier.*

# Examples

- **Evaluation operator:** $T: \mathcal{C}[0,1] \longrightarrow \mathbb{R}$ , $Tf = f(0)$.

  As norms consider $\|\cdot\|_\infty$ on $\mathcal{C}[0,1]$, $|\cdot|$ on $\mathbb{R}$. Then $\|T\| = 1$.

$$\sup_{f \in \mathcal{C}[0,1]} \frac{|Tf|}{\|f\|_\infty} = \sup \frac{|f(0)|}{\|f\|_\infty} = \text{exercise} \ldots \ldots = 1$$

- **Integral operator:** $T: \mathcal{C}[0,1] \longrightarrow \mathbb{R}$ , $Tf = \int_0^1 f(t)\, dt$

  With the same norms as above, $T$ is cont. and has $\|T\| = 1$.

# Examples

- <u>Differential operator</u>: $\mathcal{D}: \mathcal{C}^1[0,1] \longrightarrow \mathcal{C}[0,1], \; f \mapsto f'$.

  - Consider $\|\cdot\|_\infty$ on $\mathcal{C}^1$ and $\mathcal{C}$. Then $\mathcal{D}$ is linear, but not continuous!

  - Consider $\||f|\| := \|f\|_\infty + \|f'\|_\infty$ on $\mathcal{C}^1$. With this norm, $\mathcal{D}$ is continuous and bounded.

# Dual space

# Dual space

**Definition**   $V$ VS, $T : V \to F$ is called a $\overset{\text{linear}}{\underset{\smile}{}}$ functional.

Given a vector space $V$, the algebraic ==dual space $V^*$== consists

of all linear functionals on $V$:

$$V^* := \mathcal{L}(V, F).$$

If $V$ is a normed VS, then the space of all linear,

continuous functionals from $V$ to $F$ is called the

(topological) ==dual space $V'$== of $V$.

**Remark** If $V$ is finite dim, then $V^* = V'$ because then linear mappings are always continuous. In general, this is not true.

We endow the dual space with the operator norm
$$\|T\| := \sup_{x \in X} \frac{\|Tx\|}{\|x\|}.$$

**Proposition:** $V'$ is a vector space, and the operator norm is indeed a norm on $V'$.

**Prop:** If $V$ is a normed VS (but not necessarily complete), then $V'$ with the operator norm is a Banach space.

# Examples

- $K \subset \mathbb{R}$ compact set, $\mathcal{C}(K)$ space of cont. fcts with $\|\cdot\|_\infty$. Then $(\mathcal{C}(K))'$ is equivalent to the space $M(K)$, the space of all (Radon) measures over $K$.

- $S \subset \mathbb{R}$ measurable set, $1 \leq p < \infty$, $q$ such that $\frac{1}{p} + \frac{1}{q} = 1$. Then: the dual of $L_p(S)$ is given as $L^q(S)$.

# Riesz representation theorem

Theorem: $H$ Hilbert space, $H'$ its dual. Then the

mapping $\Phi : H \to H'$, $y \mapsto \langle \cdot , y \rangle$

is bijective, isometric, and satisfies $\Phi(\lambda x) = \bar{\lambda}\, \Phi(y)$.

Stated differently: for any mapping $x' \in H'$ there exists

a unique $y \in H$ such that $x'(x) = \langle x, y \rangle$.

# Adjoint operator

# Definition

**Def** Let $T \in \mathcal{L}(H_1, H_2)$, $H_1, H_2$ Hilbert spaces. Then there exists an operator $T^*: H_2 \to H_1$ such that

$$\langle Tx, y \rangle_{H_2} = \langle x, T^* y \rangle_{H_1}.$$

for all $x \in H_1$, $y \in H_2$. $T^*$ is called the adjoint of $T$.

**Remark** The existence of this operator is a consequence of the Riesz representation theorem.

**Def** The operator $T: H_1 \to H_1$ is called self-adjoint

$$\text{if} \quad \langle Tx, y \rangle = \langle x, Ty \rangle$$