

# I - Linear Algebra Basics and Functional Analysis Snippets

Math 728 D - Machine Learning & Data Science - Spring 2019

# Contents

- 1 Basic Notions, Vector Spaces
- 2 Linear Mappings
- 3 Important Matrix Classes
- 4 How to Measure and Quantify
  - Norms, Normed Linear Spaces
  - Well-Posedness, Condition Numbers
  - Orthogonal Projections
- 5 Matrix Factorizations
  - Eigenvectors, Eigenvalues
  - Spectral Decompositions
  - **LR**- and **QR**-Factorization
  - Singular Value Decomposition (SVD)

# Examples

$\mathbb{R}, \mathbb{C}$  fields of real resp. complex numbers,  $\mathbb{K} \in \{\mathbb{R}, \mathbb{C}\}$

$\mathbb{K}^n := \{\mathbf{x} = (x_1, \dots, x_n)^\top : x_i \in \mathbb{K}, i = 1, \dots, n\}$

## Example 1

*Find the intersection of two lines in the plane*

$$L_1 = \{\mathbf{x} \in \mathbb{R}^2 : a_{1,1}x_1^* + a_{1,2}x_2^* = b_1\}, \quad L_2 = \{\mathbf{x} \in \mathbb{R}^2 : a_{2,1}x_1^* + a_{2,2}x_2^* = b_1\}$$

$$\mathbf{x}^* \in L_1 \cap L_2 \Leftrightarrow \begin{cases} a_{1,1}x_1^* + a_{1,2}x_2^* = b_1 \\ a_{2,1}x_1^* + a_{2,2}x_2^* = b_1 \end{cases} \Leftrightarrow \begin{pmatrix} a_{1,1} & a_{1,2} \\ a_{2,1} & a_{2,2} \end{pmatrix} \begin{pmatrix} x_1^* \\ x_2^* \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \end{pmatrix} \Leftrightarrow \mathbf{A}\mathbf{x}^* = \mathbf{b}$$

Matrix-vector multiplication:

$$\mathbf{A} = \begin{pmatrix} a_{1,1} & \dots & a_{1,n} \\ \vdots & & \vdots \\ a_{m,1} & \dots & a_{m,n} \end{pmatrix} = \underbrace{\left( \mathbf{a}^1, \dots, \mathbf{a}^n \right)}_{\text{columns of } \mathbf{A}} \in \mathbb{K}^{m \times n} \rightsquigarrow \mathbf{A}\mathbf{x} := x_1\mathbf{a}^1 + \dots + x_n\mathbf{a}^n \in \mathbb{K}^m$$

# The Core Questions

Given  $\mathbf{A} = (a_{i,j})_{i,j=1}^{m,n} \in \mathbb{K}^{m \times n}$ ,  $\mathbf{b} \in \mathbb{K}^m$ , find  $\mathbf{x} \in \mathbb{K}^n$  s.t.

$$\begin{array}{ccccccc}
 x_1 a_{1,1} & + \cdots + & x_n a_{1,n} & = & b_1 & & \\
 \vdots & & \vdots & \vdots & \vdots & \Leftrightarrow & \mathbf{Ax} = \mathbf{b} \Leftrightarrow \sum_{j=1}^n x_j \mathbf{a}^j = \mathbf{b} \\
 x_1 a_{m,1} & + \cdots + & x_n a_{m,n} & = & b_m & & 
 \end{array}$$

The driving issue in Linear Algebra is the “solvability” of such systems, more precisely:

under which conditions on  $\mathbf{A}$ ,  $\mathbf{b}$  does there

- exist exactly one solution  $\mathbf{x} \in \mathbb{K}^n$
- exist at least one solution  $\mathbf{x} \in \mathbb{K}^n$
- no solution?
- what is the structure of the set of solutions?

These questions can be best answered by viewing the matrix  $\mathbf{A}$  as a *mapping*

$$\mathbf{A} : \mathbb{K}^n \rightarrow \mathbb{K}^m, \quad \mathbf{A} : \mathbf{x} \mapsto \mathbf{Ax}$$

# What has been used?

$$x \in \mathbb{K}, \mathbf{a} \in \mathbb{K}^m \rightsquigarrow \mathbf{x}\mathbf{a} := \begin{pmatrix} xa_1 \\ \vdots \\ xa_m \end{pmatrix} \in \mathbb{K}^m, \quad \mathbf{a}, \mathbf{b} \in \mathbb{K}^m, \rightsquigarrow \mathbf{a} + \mathbf{b} := \begin{pmatrix} a_1 + b_1 \\ \vdots \\ a_m + b_m \end{pmatrix} \in \mathbb{K}^m$$

## Definition 2

A set  $\mathbb{V}$  is a *Vector Space* (linear space) over the field  $\mathbb{K}$  ( $\mathbb{K}$ -vector space) if the following holds:

- $\exists$  an “addition”  $+$  :  $\mathbb{V} \times \mathbb{V} \rightarrow \mathbb{V}$  which is commutative and associative
- $\exists$  a “multiplication”  $x \in \mathbb{K}, \mathbf{a} \in \mathbb{V} \mapsto \mathbf{x}\mathbf{a} \in \mathbb{V}$  which is distributive

## Example 3

- 1  $\mathbb{V} = \mathbb{K}^m$  is a vector space over  $\mathbb{K}$
- 2  $C^k([a, b]) := \{f : [a, b] \rightarrow \mathbb{R} : f^{(k)} \text{ is continuous in every } x \in [a, b]\}$  is a vector space over  $\mathbb{R}$  where
 
$$(af)(x) := af(x), \quad (f+g)(x) := f(x) + g(x), \quad x \in [a, b]$$

- 3
 
$$\mathbb{P}_n := \left\{ p(x) : p(x) = \sum_{j=0}^n a_j x^j, a_j \in \mathbb{R}, j = 0, \dots, n \right\}$$

is a vector space over  $\mathbb{R}$  with the same operations as in (2)

# Forms, Scalar (Inner) Products

Let  $\mathbb{X}, \mathbb{Y}$  be an  $\mathbb{K}$ -vector spaces

$\mathbb{K} = \mathbb{R}$ :  $b(\cdot, \cdot) : \mathbb{X} \times \mathbb{Y} \rightarrow \mathbb{R}$  is called **bilinear** if  $\forall x, z \in \mathbb{X}, y, v \in \mathbb{Y}, \alpha, \beta \in \mathbb{R}$

$$b(\alpha x + \beta z, y) = \alpha b(x, y) + \beta b(z, y), \quad b(x, \alpha y + \beta v) = \alpha b(x, y) + \beta b(x, v)$$

In other words  $b(\cdot, y) : \mathbb{X} \rightarrow \mathbb{R}, b(x, \cdot) : \mathbb{Y} \rightarrow \mathbb{R}$  are for fixed  $x \in \mathbb{X}, y \in \mathbb{Y}$  linear functionals

$b(\cdot, \cdot) : \mathbb{X} \times \mathbb{X} \rightarrow \mathbb{R}$  is called **symmetric** if  $b(x, y) = b(y, x), x, y \in \mathbb{X}$

$\mathbb{K} = \mathbb{C}$ :  $b(\cdot, \cdot) : \mathbb{X} \times \mathbb{Y} \rightarrow \mathbb{C}$  is called **sesquilinear** if for fixed  $y \in \mathbb{Y}, x \in \mathbb{X}, b(\cdot, y) : \mathbb{X} \rightarrow \mathbb{K}$  is linear and  $b(x, \cdot) : \mathbb{Y} \rightarrow \mathbb{K}$  is **semi-linear**, i.e.,  $b(x, \alpha y) = \bar{\alpha} b(x, y)$ . Hence  $\forall x, z \in \mathbb{X}, y, v \in \mathbb{Y}, \alpha, \beta \in \mathbb{R}$

$$b(\alpha x + \beta z, y) = \alpha b(x, y) + \beta b(z, y), \quad b(x, \alpha y + \beta v) = \bar{\alpha} b(x, y) + \bar{\beta} b(x, v)$$

A sesquilinear form  $b(\cdot, \cdot) : \mathbb{X} \times \mathbb{X} \rightarrow \mathbb{C}$  is called **hermitian** if  $b(x, y) = \overline{b(y, x)}$

## Definition 4

Let  $b(\cdot, \cdot) : \mathbb{X} \times \mathbb{X} \rightarrow \mathbb{K}$  be a symmetric, resp. hermitian form when  $\mathbb{K} = \mathbb{R}$  resp.  $\mathbb{K} = \mathbb{C}$ . Then  $b(\cdot, \cdot) : \mathbb{X} \times \mathbb{X} \rightarrow \mathbb{K}$  is called a **scalar product** (inner product, dot product) if it is **positive definite**, i.e.,

$$b(x, x) > 0 \quad \forall x \in \mathbb{X} \setminus \{0\}$$

# Examples

- $\mathbb{K} = \mathbb{R}, \mathbb{X} = \mathbb{R}^n, \mathbb{Y} = \mathbb{R}^m, \mathbf{A} \in \mathbb{R}^{n \times m}, b(\mathbf{x}, \mathbf{y}) := \mathbf{y}^\top \mathbf{A} \mathbf{x}$
- $\mathbb{X} = \mathbb{Y} = \mathbb{R}^n$

$$b(\mathbf{x}, \mathbf{y}) = \langle \mathbf{x}, \mathbf{y} \rangle_n := \mathbf{y}^\top \mathbf{x} = \sum_{k=1}^n x_k y_k =: \mathbf{x} \cdot \mathbf{y}$$

show that this is a scalar product on  $\mathbb{R}^n$

- $\mathbb{K} = \mathbb{C}, \mathbb{X} = \mathbb{C}^n$

$$b(\mathbf{x}, \mathbf{y}) = \langle \mathbf{x}, \mathbf{y} \rangle_n := \mathbf{y}^* \mathbf{x} := \sum_{k=1}^n x_k \bar{y}_k, \quad \mathbf{y}^* := \bar{\mathbf{y}}^\top$$

show that this is a scalar product on  $\mathbb{C}^n$

- $\mathbb{K} = \mathbb{R}, \mathbb{X} = C([a, b])$

$$\langle f, g \rangle_{[a,b]} := \int_a^b f(x)g(x)dx$$

is a scalar product on  $C([a, b])$

# Dimension

Recall:

- a vector space  $\mathbb{V}$  is generated by a set  $\{v_1, \dots, v_k\} \subset \mathbb{V}$  ( $\mathbb{V} = \text{span}\{v_1, \dots, v_k\}$ ) if every  $v \in \mathbb{V}$  can be written as linear combination of the  $v_j$ ,  $j = 1, \dots, k$ , i.e., for some  $c_1, \dots, c_k \in \mathbb{K}$

$$v = \sum_{j=1}^k c_j v_j$$

- there exists a minimal  $n \in \mathbb{N}$  such that  $\mathbb{V}$  can be generated by  $n$  elements of  $\mathbb{V}$ .  $n$  is called the *dimension* of  $\mathbb{V}$ ,  $n = \dim \mathbb{V}$
- a minimal generating set is called a *basis* of  $\mathbb{V}$ .

## Exercise 5

- Determine  $\dim \mathbb{K}^n$ ,  $\dim \mathbb{P}_n$ , exhibit bases
- What is the dimension of  $C^k([a, b])$ ?
- The elements of a basis  $\{v_1, \dots, v_n\}$  are linearly independent, i.e.,

$$\sum_{j=1}^n c_j v_j = 0 \Rightarrow c_j = 0, \quad j = 1, \dots, n$$



# Linear Mappings

## Definition 6

$\mathbb{X}, \mathbb{Y}$  vector spaces (over  $\mathbb{K}$ ),  $L : \mathbb{X} \rightarrow \mathbb{Y}$  is called a *linear mapping* (operator) if for any  $a, b \in \mathbb{K}$ ,  $v, w \in \mathbb{X}$

$$L(av + bw) = aLv + bLw \in \mathbb{Y}$$

- $L$  is injective (one-to-one) if  $L(v) = 0 \Rightarrow v = 0$
- $L$  is surjective (onto) if for  $y \in \mathbb{Y}$  there exists  $x \in \mathbb{X}$  such that  $L(x) = y$
- $L$  is bijective (one-to-one and onto, invertible) if  $L$  is injective and surjective

The set of all linear mappings from  $\mathbb{X}$  to  $\mathbb{Y}$  is denoted by  $\mathcal{L}(\mathbb{X}, \mathbb{Y})$

- Any matrix  $\mathbf{A} \in \mathbb{K}^{m \times n}$  induces a *linear mapping* (operator) from  $\mathbb{K}^n$  to  $\mathbb{K}^m$  by

$$\mathbf{x} \in \mathbb{K}^n \mapsto \mathbf{Ax} \in \mathbb{K}^m \quad (\text{recall } \mathbf{Ax} = x_1 \mathbf{a}^1 + \cdots + x_n \mathbf{a}^n) \rightsquigarrow \mathbf{A}(ax^1 + bx^2) = a\mathbf{Ax}^1 + b\mathbf{Ax}^2$$

- $\mathbb{V}$  a vector space over  $\mathbb{K}$ , a (linear) mapping  $\ell : \mathbb{V} \rightarrow \mathbb{K}$  is called a (linear) *functional*

### Examples:

$$- \mathbb{V} = C^0([a, b]), \ell(f) := \int_a^b f(x)\omega(x)dx$$

$$- \mathbb{V} \text{ as before, } \delta_{x_0} f := f(x_0)$$

$$- \mathbb{V} = \mathbb{R}^n, \mathbf{a} \in \mathbb{R}^n \text{ fixed, } \ell(\mathbf{x}) := \sum_{j=1}^n a_j x_j =: \mathbf{a} \cdot \mathbf{x} = \mathbf{a}^\top \mathbf{x}$$

# Linear Mappings

## some further notions and facts

- 1 If  $\mathbb{X}, \mathbb{Y}$  are  $\mathbb{K}$ -vector spaces then  $\mathcal{L}(\mathbb{X}, \mathbb{Y})$  is a  $\mathbb{K}$ -vector space; identify the “+” in  $\mathcal{L}(\mathbb{X}, \mathbb{Y})$
- 2  $\ker(L) := \{x \in \mathbb{X} : L(x) = 0\}$  is a linear subspace of  $\mathbb{X}$
- 3  $\text{ran}(L) := \{y \in \mathbb{Y} : \exists x \in \mathbb{X}, \text{ s.t. } y = L(x)\}$  is a linear subspace of  $\mathbb{Y}$
- 4 let  $\dim \mathbb{X} = n < \infty$ ,  $L \in \mathcal{L}(\mathbb{X}, \mathbb{Y})$  bijective implies  $\dim \mathbb{Y} = n$
- 5 compositions:  $L \in \mathcal{L}(\mathbb{X}, \mathbb{Y})$ ,  $B \in \mathcal{L}(\mathbb{V}, \mathbb{X})$ , then  $LB$ , defined by  $LB(v) := L(B(v))$ ,  $v \in \mathbb{V}$ , belongs to  $\mathcal{L}(\mathbb{V}, \mathbb{Y})$
- 6 let  $L \in \mathcal{L}(\mathbb{X}, \mathbb{Y})$  be bijective, then  $L^{-1} : \mathbb{Y} \rightarrow \mathbb{X}$ , defined by  $L(x) = y \Leftrightarrow L^{-1}(y) = x$ , belongs to  $\mathcal{L}(\mathbb{Y}, \mathbb{X})$ ;  $L^{-1}$  is called *inverse* of  $L$  and

$$LL^{-1} = I_{\mathbb{Y}}, \quad L^{-1}L = I_{\mathbb{X}}, \quad (I_{\mathbb{V}}v = v, v \in \mathbb{V}, \text{ identity on } \mathbb{V})$$

Thus, the set of bijective mappings in  $\mathcal{L}(\mathbb{X}, \mathbb{X})$  form a *group* with respect to composition as multiplication.

### Proposition 7

- every  $\mathbb{K}$ -vector space  $\mathbb{X}$  of dimension  $\dim \mathbb{X} = n$  is isomorphic to  $\mathbb{K}^n$ , in the sense that there exists a bijection  $B \in \mathcal{L}(\mathbb{X}, \mathbb{K}^n)$  such that  $\mathbb{K}^n = \{B(x) : x \in \mathbb{X}\}$ . Hence any two  $n$ -dimensional  $\mathbb{K}$ -vector spaces are isomorphic
- If  $\dim(\mathbb{X}) = n$ ,  $\dim(\mathbb{Y}) = m$ , then  $\mathcal{L}(\mathbb{X}, \mathbb{Y})$ ,  $\mathcal{L}(\mathbb{K}^n, \mathbb{K}^m)$ ,  $\mathbb{K}^{mn}$ ,  $\mathbb{K}^{m \times n}$  are all isomorphic to each other

# Proof of Proposition 7

Let  $\mathbb{X}$  be an  $n$ -dimensional  $\mathbb{K}$ -vector space and  $L \in \mathcal{L}(\mathbb{X}, \mathbb{X})$ . Then there exists a basis  $\Phi = \{\phi_1, \dots, \phi_n\} \subset \mathbb{X}$  of  $\mathbb{X}$ . Hence, every  $x \in \mathbb{X}$  has a unique representation

$$x = \sum_{k=1}^n x_k \phi_k \in \mathbb{X}. \quad (3.1)$$

Define  $B_\Phi : \mathbb{X} \rightarrow \mathbb{K}^n$  by  $B_\Phi(x) := \mathbf{x} = (x_1, \dots, x_n)^\top$  ( $B_\Phi$  depends on the choice of basis). Clearly  $B_\Phi$  is linear and injective (by definition of a basis), i.e.,  $B_\Phi \in \mathcal{L}(\mathbb{X}, \mathbb{K}^n)$ . Also  $M : \mathbb{K}^n \rightarrow \mathbb{X}$ , defined by  $M(\mathbf{x}) := \sum_{k=1}^n x_k \phi_k$  is linear, injective and obviously  $M(B_\Phi(x)) = x$ . Thus  $M = B_\Phi^{-1}$ .

Now suppose  $\mathbb{X}, \mathbb{Y}$  are  $K$ -vector spaces with bases  $\Phi = \{\phi_1, \dots, \phi_n\} \subset \mathbb{X}$ ,  $\Psi = \{\psi_1, \dots, \psi_m\} \subset \mathbb{Y}$ , respectively. In particular, for each  $k \in \{1, \dots, n\}$  one has  $B_\Psi(L(\phi_k)) = (c_{1,k}, \dots, c_{m,k})^\top =: \mathbf{c}^k \in \mathbb{K}^m$ . Thus, for  $\mathbf{C} = (\mathbf{c}^1, \dots, \mathbf{c}^n) \in \mathbb{K}^{m \times n}$

$$y := L(x) = L\left(\sum_{k=1}^n x_k \phi_k\right) = \sum_{k=1}^n x_k L(\phi_k) = \sum_{r=1}^m \left(\sum_{k=1}^n c_{r,k} x_k\right) \psi_r = B_\Psi^{-1}(\mathbf{C}\mathbf{x}) \in \mathbb{Y}.$$

In other words

$$L(x) = y \quad \Leftrightarrow \quad \mathbf{C}\mathbf{x} = \mathbf{y} \quad \text{for} \quad \mathbf{y} = B_\Psi(y), \quad \mathbf{x} = B_\Phi(x) \quad (3.2)$$

In that sense  $L \in \mathcal{L}(\mathbb{X}, \mathbb{Y})$  is **represented** by  $\mathbf{C} \in \mathcal{L}(\mathbb{K}^n, \mathbb{K}^m) \cong \mathbb{K}^{m \times n}$

# Linear Mappings, Matrices

some further notions and facts

- by (5) above  $\mathbf{A} \in \mathbb{K}^{m \times n}$ ,  $\mathbf{B} \in \mathbb{K}^{n \times r} \Rightarrow \mathbf{AB} = (\mathbf{Ab}^1, \dots, \mathbf{Ab}^r) \in \mathbb{K}^{r \times m}$
- change of bases: suppose  $\mathbf{a}^1, \dots, \mathbf{a}^n$ , and  $\mathbf{b}^1, \dots, \mathbf{b}^n$  are both bases for  $\mathbb{K}^n$ , then there exists a unique  $\mathbf{M} \in \mathbb{K}^{n \times n}$  such that  $\mathbf{b}^j = \mathbf{Ma}^j$ ,  $j = 1, \dots, n$   
In fact:  $\mathbf{A} := (\mathbf{a}^1, \dots, \mathbf{a}^n)$ ,  $\mathbf{B} := (\mathbf{b}^1, \dots, \mathbf{b}^n) \in \mathbb{K}^{n \times n}$  then  $\mathbf{B} = (\mathbf{BA}^{-1})\mathbf{A}$ , i.e.,  $\mathbf{M} = \mathbf{BA}^{-1}$
- $\mathbf{A} = (\mathbf{a}^1, \dots, \mathbf{a}^n) \rightsquigarrow \text{ran}(\mathbf{A}) = \text{span}\{\mathbf{a}^1, \dots, \mathbf{a}^n\}$
- $\text{rank}(\mathbf{A}) := \dim(\text{ran}(\mathbf{A})) = \#(\text{linearly independent columns of } \mathbf{A})$
- $\mathbf{A} \in \mathbb{K}^{m \times n}$ ,  $\ker(\mathbf{A}) = \{\mathbf{x} \in \mathbb{K}^n : \mathbf{Ax} = \mathbf{0}\} = \text{ran}(\mathbf{A}^\top)^\perp$ , where

$$\mathbf{A} = (a_{i,j})_{i,j=1}^{m,n} \in \mathbb{K}^{m \times n} \rightsquigarrow \mathbf{A}^\top := (a_{j,i})_{j,i=1}^{n,m} \in \mathbb{K}^{n \times m} \text{ (transpose of } \mathbf{A})$$

and

$$\mathbb{V} \subset \mathbb{K}^n \rightsquigarrow \mathbb{V}^\perp := \{\mathbf{y} \in \mathbb{K}^n : \mathbf{y} \cdot \mathbf{v} = 0, \mathbf{v} \in \mathbb{V}\}$$

## Exercise 8

- show that  $\#(\text{linearly independent columns of } \mathbf{A}) = \#(\text{linearly independent rows of } \mathbf{A})$
- show that for  $\mathbf{A} \in \mathbb{K}^{m \times n}$  one has  $m = \dim(\text{ran}(\mathbf{A})) + \dim(\ker(\mathbf{A}^\top))$

# Solvability of Linear Systems

Given  $\mathbf{A} \in \mathbb{K}^{m \times n}$ ,  $\mathbf{b} \in \mathbb{K}^m$ , solve  $\mathbf{Ax} = \mathbf{b}$ :

## Exercise 9

- *There is a unique solution if and only if  $\mathbf{b} \in \text{ran}(\mathbf{A})$ ,  $m \geq n$ ,  $\text{rank}(\mathbf{A}) = n$ .*
- *There exists a unique solution for **every**  $\mathbf{b} \in \mathbb{R}^m$  if and only if  $m = n$  and  $\text{rank}(\mathbf{A}) = n$ . In this case  $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$ .*
- *For  $n = m$  there exists a unique solution for **every**  $\mathbf{b} \in \mathbb{R}^m$  if and only if  $\ker(\mathbf{A}) = \{0\}$ .*
- *Recall:  $\det : \mathbb{K}^{n \times n} \rightarrow \mathbb{K}$  is a multilinear anti-symmetric mapping and  $\mathbf{A}^{-1}$  exists if and only if  $\det(\mathbf{A}) \neq 0$  in which case we call  $\mathbf{A}$  non-singular.*
- *More than one solution exists if and only if infinitely many solutions exist. The solution set is the affine space*

$$S(\mathbf{A}) = \mathbf{x}^0 + \ker(\mathbf{A}) = \{\mathbf{x} = \mathbf{x}^0 + \mathbf{y} : \mathbf{y} \in \ker(\mathbf{A})\}$$

Notice: searching for a “solution” of  $\mathbf{Ax} = \mathbf{b}$  makes sense only if  $\mathbf{A} \in \mathbb{K}^{n \times n}$  is non-singular, that is solvability depends **only** on  $\mathbf{A}$ , why? In all other cases the notion of “solution” has to be generalized.

# Hermitian and Positive Definite Matrices

For  $\mathbf{A} = (a_{k,j})_{k,j=1}^{m,n} \in \mathbb{K}^{m \times n}$  the **transpose**  $\mathbf{A}^\top$  is defined by  $\mathbf{A}^\top = (a_{j,k})_{k,j=1}^{n,m} \in \mathbb{K}^{n \times m}$  (reflect across the diagonal)

Recall:  $i$  imaginary unit,  $i^2 = -1$ ,  $z = x + iy \in \mathbb{C}$ ,  $\bar{z} := x - iy \rightsquigarrow z\bar{z} = |z|^2 = x^2 + y^2$

$\mathbb{K} = \mathbb{C}$ : **hermitian conjugate**:

$$\mathbf{A}^* := \overline{\mathbf{A}^\top} = (\overline{a_{j,k}})_{j=1,k}^{n,m} \in \mathbb{C}^{n \times m} \quad (\mathbf{A}^* = \mathbf{A}^\top \text{ when } \mathbb{K} = \mathbb{R})$$

## Definition 10

$\mathbf{A} \in \mathbb{K}^{n \times n}$  is called **symmetric** ( $\mathbb{K} = \mathbb{R}$ ), resp. **hermitian** ( $\mathbb{K} = \mathbb{C}$ ) if  $\mathbf{A} = \mathbf{A}^*$

A hermitian matrix  $\mathbf{A} \in \mathbb{K}^{n \times n}$  is called **positive (semi-)definite** if

$$\mathbf{x}^* \mathbf{A} \mathbf{x} > (\geq) 0 \quad \forall \mathbf{x} \in \mathbb{K}^n \setminus \{0\}$$

## Remark 11

For  $\mathbf{A} \in \mathbb{K}^{n \times n}$ , one has  $\langle \mathbf{A} \mathbf{x}, \mathbf{y} \rangle_n = \langle \mathbf{x}, \mathbf{A}^* \mathbf{y} \rangle_n$ , i.e.,  $\mathbf{A}^*$  is called the **adjoint** of  $\mathbf{A}$  for reasons to be explained later.

## Examples of Symmetric/Hermitian Positive Definite Matrices

- $\mathbf{I}$  is in each “club”
- For any  $\mathbf{B} \in \mathbb{K}^{m \times n}$  the matrices  $\mathbf{A} := \mathbf{B}^* \mathbf{B} \in \mathbb{K}^{n \times n}$ ,  $\mathbf{A} := \mathbf{B} \mathbf{B}^* \in \mathbb{K}^{m \times m}$  are hermitian positive semi-definite. Under which conditions on  $\mathbf{B}$  is which version positive definite?
- Suppose that  $\{\phi_1, \dots, \phi_n\} \subset \mathbb{X}$  where  $\mathbb{X}$  is a  $\mathbb{K}$ -vector space and  $\langle \cdot, \cdot \rangle_{\mathbb{X}}$  is a scalar product on  $\mathbb{X}$ . Then the **Gramian** matrix

$$\mathbf{G} := (\langle \phi_k, \phi_j \rangle_{\mathbb{X}})_{k,j=1}^{n,n} \in \mathbb{K}^{n \times n} \quad (4.1)$$

is hermitian positive semi-definite. Under which condition on the set  $\{\phi_1, \dots, \phi_n\} \subset \mathbb{X}$  is  $\mathbf{G}$  positive definite?

## Remark 12

*Any hermitian positive definite matrix is nonsingular and its inverse is also hermitian positive definite.*

**Proof:** Suppose  $\det(\mathbf{A}) = 0$ ,  $\Rightarrow \exists \mathbf{x} \in \mathbb{K}^n \setminus \{0\}$  s.t.  $\mathbf{A}\mathbf{x} = 0$ ,  $\Rightarrow \mathbf{x}^* \mathbf{A}\mathbf{x} = 0$  which is a contradiction for  $\mathbf{x} \neq 0$ . Now show that  $\mathbf{A}^{-1} = (\mathbf{A}^*)^{-1} \stackrel{!}{=} (\mathbf{A}^{-1})^* \Leftrightarrow \mathbf{I} = \mathbf{A}^* (\mathbf{A}^{-1})^*$ . Indeed, (since  $\mathbf{B}^* \mathbf{C}^* = (\mathbf{CB})^*$ ) one has  $\mathbf{A}^* (\mathbf{A}^{-1})^* = (\mathbf{A}^{-1} \mathbf{A})^* = \mathbf{I}^* = \mathbf{I}$ . □

## Further Properties of Hermitian Positive Definite Matrices

## Proposition 13

Let  $\mathbf{A} \in \mathbb{K}^{n \times n}$  be hermitian positive semi-definite, then:

- 1 All principal submatrices of  $\mathbf{A}$  are hermitian positive semi-definite. In particular, all diagonal entries are real and non-negative  $a_{k,k} \geq 0$ ,  $k = 1, \dots, n$ .
- 2 All eigenvalues of  $\mathbf{A}$  are real and nonnegative.
- 3 The maximum modulus entry of  $\mathbf{A}$  occurs on the diagonal of  $\mathbf{A}$ .

**Proof:** For any  $J \subset \{1, \dots, n\}$  let  $\mathbf{A}_J := (a_{j,k})_{(j,k) \in J \times J} \in \mathbb{K}^{\#(J) \times \#(J)}$ . For  $J \neq \emptyset$  and any  $\mathbf{x} \in \mathbb{K}^n \setminus \{0\}$  supported in  $J$  let  $\tilde{\mathbf{x}} := \mathbf{x}|_J \in \mathbb{K}^{\#(J)}$ . Since  $\tilde{\mathbf{x}} \neq \mathbf{0}$  we have  $0 < \mathbf{x}^* \mathbf{A} \mathbf{x} = \tilde{\mathbf{x}}^* \mathbf{A}_J \tilde{\mathbf{x}}$ . Since  $\tilde{\mathbf{x}} \neq \mathbf{0}$  is arbitrary (1) follows. As for (2), suppose  $\mathbf{A} \mathbf{x} = \lambda \mathbf{x}$ ,  $\lambda \neq 0$ . Then  $0 < \mathbf{x}^* \mathbf{A} \mathbf{x} = \lambda \mathbf{x}^* \mathbf{x} = \lambda \sum_{k=1}^n |x_k|^2$  which implies  $\lambda > 0$ . Concerning (3), we know from Theorem 30 (later below) that then all principal minors (determinants of principal submatrices) are nonnegative. Now suppose, that  $|a_{j,k}| = \max_{1 \leq r,s \leq n} |a_{r,s}|$  and  $j \neq k$ . Then (w.l.o.g.  $j < k$ )

$$0 \leq \det \begin{pmatrix} a_{j,j} & a_{j,k} \\ a_{j,k} & a_{k,k} \end{pmatrix} = a_{j,j} a_{k,k} - a_{k,j} a_{j,k} = a_{j,j} a_{k,k} - |a_{j,k}|^2$$

which means that  $|a_{j,k}|^2 \leq a_{j,j} a_{k,k}$  and hence  $\max\{a_{j,j}, a_{k,k}\} \geq |a_{j,k}|$ . □



# Unitary Matrices

$\mathbf{A} \in \mathbb{K}^{n \times n}$  is called **unitary** (orthogonal if  $\mathbb{K} = \mathbb{R}$ ):  $\mathbf{Q} \in \mathcal{O}_n = \mathcal{O}_n(\mathbb{K})$  if  $\mathbf{Q}^* \mathbf{Q} = \mathbf{I}$

Properties:

- $\mathbf{Q}^* \mathbf{Q} = \mathbf{I}$  means that the columns  $\mathbf{q}^j$ ,  $j = 1, \dots, n$ , are an orthonormal basis of  $\mathbb{K}^n$ , i.e.,

$$(\mathbf{q}^j)^* \mathbf{q}^k = \sum_{r=1}^n \overline{q_{r,j}} q_{r,k} = \delta_{j,k}, \quad j, k = 1, \dots, n$$

- $\mathbf{Q} \in \mathcal{O}_n \Leftrightarrow \mathbf{Q}^* \in \mathcal{O}_n$ , i.e., the rows of  $\mathbf{Q}$  also form an orthonormal basis
- $\mathbf{Q}, \tilde{\mathbf{Q}} \in \mathcal{O}_n \Rightarrow \mathbf{Q}\tilde{\mathbf{Q}} \in \mathcal{O}_n$ , i.e.,  $\mathcal{O}_n$  is a multiplicative group
- $|\det(\mathbf{Q})| = 1$  (since  $1 = \det(\mathbf{Q}^* \mathbf{Q}) = \det(\mathbf{Q}^*) \det(\mathbf{Q}) = \overline{\det(\mathbf{Q})} \det(\mathbf{Q}) = |\det(\mathbf{Q})|^2$ )

Permutation matrices belong to  $\mathcal{O}_n$ : let  $\pi : \{1, \dots, n\} \rightarrow \{1, \dots, n\}$  be a permutation,  $\mathbf{e}^j := (0, \dots, 0, 1, 0, \dots, 0)^\top = (\delta_{k,j})_{k=1}^j$  the  $j$ th coordinate vector. Then

$$\mathbf{P}_\pi := (\mathbf{e}^{\pi(1)}, \dots, \mathbf{e}^{\pi(n)}) \in \mathcal{O}_n \quad \text{and} \quad \mathbf{A}\mathbf{P}_\pi = (\mathbf{a}^{\pi(1)}, \dots, \mathbf{a}^{\pi(n)})$$

Hence  $\mathbf{P}_\pi^\top \mathbf{A}$  permutes the rows of  $\mathbf{A}$  according to  $\pi$

# Normed Linear Spaces

How to “measure” vectors, functions?

Let  $\mathbb{V}$  be a  $\mathbb{K}$ -vector space (finite or infinite dimensional). Any mapping  $\|\cdot\| : \mathbb{V} \rightarrow \mathbb{R}_+$  is called a **norm** on  $\mathbb{V}$  if

$$(N1) \quad \|v\| = 0 \quad \Rightarrow \quad v = 0$$

$$(N2) \quad \|\alpha v\| = |\alpha| \|v\|, \quad \alpha \in \mathbb{K}, v \in \mathbb{V}$$

$$(N3) \quad \|v + z\| \leq \|v\| + \|z\|, \quad v, z \in \mathbb{V} \text{ (triangle inequality)}$$

The pair  $(\mathbb{V}, \|\cdot\|)$  is called a **a normed linear space**. When the choice of the norm is clear we simply say that  $\mathbb{V}$  is a normed linear space.

## Remark 14

*The absolute value is a norm on  $\mathbb{V} = \mathbb{K}$ . Up to a scaling factor this is the only norm on  $\mathbb{K} = \mathbb{R}$ . For higher dimensional spaces many different norms exist. Their choice depends on the purpose or application. Every norm is a Lipschitz continuous mapping with Lipschitz constant equal to one, i.e.,*

$$|\|v\| - \|z\|| \leq \|v - z\|, \quad v, z \in \mathbb{V}. \quad (5.1)$$

To see (5.1), note that by (N3)  $\|v\| = \|v - z + z\| \leq \|v - z\| + \|z\| \Rightarrow \|v\| - \|z\| \leq \|v - z\|$ .

The same argument shows that  $\|z\| - \|v\| \leq \|v - z\|$ .

# Examples

## Vector and Sequence Norms

$p$ -norms on  $\mathbb{K}^n$ ,  $1 \leq p \leq \infty$ :

$$\|\mathbf{x}\|_p := \left( \sum_{j=1}^n |x_j|^p \right)^{1/p}, \quad 1 \leq p < \infty, \quad \|\mathbf{x}\|_\infty := \max_{j=1, \dots, n} |x_j|$$

For any fixed positive “weights”  $\omega_j, j = 1, \dots, n$ , the “weighted” counterparts

$$\|\mathbf{x}\|_{p, \omega} := \left( \sum_{j=1}^n \omega_j |x_j|^p \right)^{1/p} \quad (5.2)$$

define norms as well.

$p$ -norms on  $\mathbb{K}^\infty$  (sequences),  $1 \leq p \leq \infty$ :

$$\|\mathbf{x}\|_p := \left( \sum_{j=1}^{\infty} |x_j|^p \right)^{1/p}, \quad 1 \leq p < \infty, \quad \|\mathbf{x}\|_\infty := \sup_{j=1, \dots, n} |x_j|$$

In this case one often uses the notation  $\|\cdot\|_{\ell_p}$  and  $\ell_p$  denotes the space of all sequences for which  $\|\mathbf{x}\|_{\ell_p}$  is finite. When using specific index sets  $\mathcal{I}$  different from  $\mathbb{N}$ , one writes  $\ell_p(\mathcal{I})$  viewing a sequence as a **function** mapping a discrete domain  $\mathcal{I}$  to  $\mathbb{K}$ . To stress the dependence on a finite dimension  $n$  we also write briefly  $\ell_p^n = \ell_p(\{1, \dots, n\})$

# Examples

## Vector and Sequence Norms

$\|\cdot\|_p, \|\cdot\|_{\ell_p}$  are indeed norms:

- For  $p = 1, \infty$  this follows directly corresponding properties of the absolute value.
- For  $1 < p < \infty$  properties (N1), (N2) for  $\|\cdot\|_p, \|\cdot\|_{\ell_p}$  follow directly as well. The triangle inequality (N3) is in this case less obvious. The main tool for verifying (N3) as well is the following important inequality:

**Hölder's Inequality:** for  $\frac{1}{p} + \frac{1}{p^*} = 1$  one has

$$\left| \sum_{j \in \mathcal{I}} x_j \bar{y}_j \right| \leq \|\mathbf{x}\|_{\ell_p(\mathcal{I})} \|\mathbf{y}\|_{\ell_{p^*}(\mathcal{I})}, \quad \mathbf{x} \in \ell_p(\mathcal{I}), \mathbf{y} \in \ell_{p^*}(\mathcal{I}). \quad (5.3)$$

We'll show later for the special case  $p = 2$  how this implies the triangle inequality (N3).

Defining  $\mathbf{x}\bar{\mathbf{y}} := (x_j \bar{y}_j)_{j \in \mathcal{I}}$ , (5.3) implies

$$\|\mathbf{x}\bar{\mathbf{y}}\|_{\ell_1} \leq \|\mathbf{x}\|_{\ell_p(\mathcal{I})} \|\mathbf{y}\|_{\ell_{p^*}(\mathcal{I})}.$$

# Operator Norms

One can “measure” the mapping properties of a linear operator by monitoring how it distorts the unit ball:

$$\|L\|_{\mathcal{L}(\mathbb{X}, \mathbb{Y})} := \sup_{x \in \mathbb{X}; \|x\|_{\mathbb{X}}=1} \|L(x)\|_{\mathbb{Y}} = \sup_{x \in \mathbb{X} \setminus \{0\}} \frac{\|L(x)\|_{\mathbb{Y}}}{\|x\|_{\mathbb{X}}}$$

Verify:  $\|L\|_{\mathcal{L}(\mathbb{X}, \mathbb{Y})}$  is a norm on  $\mathcal{L}(\mathbb{X}, \mathbb{Y})$

The operator norm  $\|L\|_{\mathcal{L}(\mathbb{X}, \mathbb{Y})}$  is induced by and thus depends on the norms  $\|\cdot\|_{\mathbb{X}}$ ,  $\|\cdot\|_{\mathbb{Y}}$ .

- $L$  is called **bounded** if  $\|L\|_{\mathcal{L}(\mathbb{X}, \mathbb{Y})} < \infty$
- $\rightsquigarrow \|L(x)\|_{\mathbb{Y}} = \frac{\|L(x)\|_{\mathbb{Y}}}{\|x\|_{\mathbb{X}}} \|x\|_{\mathbb{X}} \leq \sup_{y \in \mathbb{X} \setminus \{0\}} \frac{\|L(y)\|_{\mathbb{Y}}}{\|y\|_{\mathbb{X}}} \|x\|_{\mathbb{X}} \rightsquigarrow$

$$\|L(x)\|_{\mathbb{Y}} \leq \|L\|_{\mathcal{L}(\mathbb{X}, \mathbb{Y})} \|x\|_{\mathbb{X}}, \quad x \in \mathbb{X}. \quad (5.4)$$

In other words, if  $L \in \mathcal{L}(\mathbb{X}, \mathbb{Y})$  is bounded there exists a constant  $C < \infty$  such that  $\|L(x)\|_{\mathbb{Y}} \leq C\|x\|_{\mathbb{X}}$ ,  $x \in \mathbb{X}$ , and  $C = \|L\|_{\mathcal{L}(\mathbb{X}, \mathbb{Y})}$  is the smallest such constant.

- A bounded linear operator is Lipschitz continuous with Lipschitz constant  $\|L\|_{\mathcal{L}(\mathbb{X}, \mathbb{Y})}$ :

$$\|L(x) - L(z)\|_{\mathbb{Y}} = \|L(x - z)\|_{\mathbb{Y}} \leq \|L\|_{\mathcal{L}(\mathbb{X}, \mathbb{Y})} \|x - z\|_{\mathbb{X}};$$

- $L \in \mathcal{L}(\mathbb{X}, \mathbb{Y})$ ,  $R \in \mathcal{L}(\mathbb{Y}, \mathbb{W}) \rightsquigarrow \|R \circ L\|_{\mathcal{L}(\mathbb{X}, \mathbb{W})} \leq \|R\|_{\mathcal{L}(\mathbb{Y}, \mathbb{W})} \|L\|_{\mathcal{L}(\mathbb{X}, \mathbb{Y})}$ .

# Matrix Norms

When viewing  $\mathbf{A} \in \mathbb{K}^{m \times n}$  as a “vector” in  $\mathbb{K}^{mn}$  one could, in principle, use  $p$ -norms for matrices as well. However, these norms would not reflect the properties of  $\mathbf{A}$  when viewed as a linear operator from  $\mathbb{K}^n$  to  $\mathbb{K}^m$ . Here are some operator norms for matrices: let  $\mathbf{A} \in \mathbb{K}^{m \times n} \equiv \mathcal{L}(\mathbb{K}^n, \mathbb{K}^m)$

$$\|\mathbf{A}\|_\infty = \|\mathbf{A}\|_{\mathcal{L}(\ell_\infty, \ell_\infty)} = \sup_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{A}\mathbf{x}\|_\infty}{\|\mathbf{x}\|_\infty} = \max_{j=1, \dots, m} \sum_{k=1}^n |a_{j,k}| \quad (\text{maximal row sum})$$

$$\|\mathbf{A}\|_1 = \|\mathbf{A}\|_{\mathcal{L}(\ell_1, \ell_1)} = \sup_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{A}\mathbf{x}\|_1}{\|\mathbf{x}\|_1} = \max_{k=1, \dots, n} \sum_{j=1}^m |a_{j,k}| \quad (\text{maximal column sum})$$

The so called **spectral norm**, induced by the Euclidean norm  $\|\cdot\|_2$  requires the notion of **eigenvalue** of a matrix. The set of eigenvalues of  $\mathbf{A}$  is often called “spectrum” of  $\mathbf{A}$ . This will be discussed in more detail a little later. Here it suffices to know that  $\lambda \in \mathbb{K}$ ,  $\mathbf{x} \in \mathbb{K}^n$  are called eigenvalue and eigenvector of  $\mathbf{A}$  if  $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$ , i.e.,  $\mathbf{A}$  just stretches or shortens an eigenvector while preserving its direction. Moreover (see Proposition 13 later below) the matrix  $\mathbf{A}^* \mathbf{A}$  has only nonnegative eigenvalues. We denote by  $\lambda_{\max}(\mathbf{A}^* \mathbf{A})$  the largest eigenvalue of  $\mathbf{A}^* \mathbf{A}$ . Then one can show that for  $\mathbf{A} \in \mathbb{K}^{m \times n}$

$$\|\mathbf{A}\|_2 = \|\mathbf{A}\|_{\mathcal{L}(\ell_2^n, \ell_2^m)} = \sqrt{\lambda_{\max}(\mathbf{A}^* \mathbf{A})} \quad (5.5)$$

# Unitary Matrices “Love” the Euclidean Norm

## Proposition 15

Let  $\mathbf{Q} \in \mathcal{O}_n$

- 1 Unitary matrices map Euclidean spheres into themselves:  $\|\mathbf{Q}\mathbf{x}\|_2 = \|\mathbf{x}\|_2, \mathbf{x} \in \mathbb{K}^n$ ;
- 2 they have minimal condition numbers:  $\kappa_2(\mathbf{Q}) = 1$ ;
- 3 multiplying a matrix with a unitary matrix does not change the spectral norm of the matrix:  $\|\mathbf{A}\|_2 = \|\mathbf{QA}\|_2 = \|\mathbf{AQ}\|_2$ ;
- 4 multiplying a matrix with a unitary matrix does not change the spectral norm of the matrix:  $\kappa_2(\mathbf{A}) = \kappa_2(\mathbf{QA}) = \kappa_2(\mathbf{AQ})$

**Proof:** (1):  $\|\mathbf{Q}\mathbf{x}\|_2^2 = (\mathbf{Q}\mathbf{x})^* \mathbf{Q}\mathbf{x} = \mathbf{x}^* \mathbf{Q}^* \mathbf{Q}\mathbf{x} = \mathbf{x}^* \mathbf{x} = \|\mathbf{x}\|_2^2$ ;

(2): (1)  $\Rightarrow \|\mathbf{Q}\|_2 = 1 \Rightarrow \|\mathbf{Q}^*\|_2 = \|\mathbf{Q}^{-1}\|_2 = 1 \Rightarrow \kappa_2(\mathbf{Q}) = \|\mathbf{Q}\|_2 \|\mathbf{Q}^{-1}\|_2 = 1$ ;

# Well-Posedness

A problem is called **well posed** if

- 1 There **exists** a solution;
- 2 the solution is **unique**;
- 3 the solution depends **continuously** on the data

The problem is called **ill-posed** if at least one of these requirements fails to hold.

- To deal with ill-posed data one has to **regularize** the problem, i.e., one seeks a nearby well-posed problem and solves this one
- In finite dimensions “continuous dependence” is independent of a specific norm. In the infinite dimensional case the choice of the norm may be essential.



# Condition of a Problem, Condition Numbers

The notion of **Condition** of a problem attempts to quantify how strongly the output depends on perturbations of the input data.

The **relative condition** of a linear operator  $L \in \mathcal{L}(\mathbb{X}, \mathbb{Y})$  can be estimated by the (relative) **condition number**  $\kappa_{\mathbb{X}, \mathbb{Y}}(L)$  of  $L$  which is the smallest constant  $\kappa$  satisfying

$$\frac{\|L(x) - L(\tilde{x})\|_{\mathbb{Y}}}{\|L(x)\|_{\mathbb{Y}}} \leq \kappa \frac{\|x - \tilde{x}\|_{\mathbb{X}}}{\|x\|_{\mathbb{X}}}, \quad x, \tilde{x} \in \mathbb{X}, \quad (5.6)$$

hence quantifying how the relative output accuracy is controlled by the relative input error.

## Proposition 16

One has

$$\kappa_{\mathbb{X}, \mathbb{Y}}(L) = \frac{\sup_{x \in \mathbb{X}; \|x\|_{\mathbb{X}}=1} \|L(x)\|_{\mathbb{Y}}}{\inf_{x; \|x\|_{\mathbb{X}}=1} \|L(x)\|_{\mathbb{Y}}} = \frac{\|L\|_{\mathcal{L}(\mathbb{X}, \mathbb{Y})}}{\inf_{x; \|x\|_{\mathbb{X}}=1} \|L(x)\|_{\mathbb{Y}}} \quad (5.7)$$

If  $L$  is invertible one has

$$\kappa_{\mathbb{X}, \mathbb{Y}}(L) = \|L\|_{\mathcal{L}(\mathbb{X}, \mathbb{Y})} \|L^{-1}\|_{\mathcal{L}(\mathbb{Y}, \mathbb{X})}. \quad (5.8)$$

Thus,  $L$  and  $L^{-1}$  have the same condition number.

$\kappa_{\mathbb{X}, \mathbb{Y}}(L)$  is the ratio of maximal expansion and maximal compression that can be caused by  $L$ .

# Proof of Proposition 16

$$\|L(x) - L(\tilde{x})\|_Y = \|L(x - \tilde{x})\|_Y \stackrel{(5.4)}{\leq} \|L\|_{\mathcal{L}(X,Y)} \|x - \tilde{x}\|_X, \quad (5.9)$$

and

$$\|L(x)\|_X = \frac{\|L(x)\|_Y}{\|x\|_X} \|x\|_X \geq \inf_{x' \neq 0} \frac{\|L(x')\|_Y}{\|x'\|_X} \|x\|_X = \inf_{\|x'\|_X=1} \|L(x')\|_Y \|x\|_X \quad (5.10)$$

provide

$$\frac{\|L(x) - L(\tilde{x})\|_Y}{\|L(x)\|_Y} \leq \frac{\|L\|_{\mathcal{L}(X,Y)}}{\inf_{\|x'\|_X=1} \|L(x')\|_Y} \frac{\|x - \tilde{x}\|_X}{\|x\|_X}$$

which shows that  $\kappa_{X,Y}(L) \leq \frac{\|L\|_{\mathcal{L}(X,Y)}}{\inf_{\|x'\|_X=1} \|L(x')\|_Y}$ . Show that one cannot do better.

Regarding (5.8), notice

$$\frac{1}{\inf_{\|x'\|_X=1} \|L(x')\|_Y} = \sup_{x' \neq 0} \frac{\|x'\|_X}{\|L(x')\|_Y} \stackrel{(y=L(x'))}{=} \sup_{y \in Y} \frac{\|L^{-1}(y)\|_X}{\|y\|_Y} = \|L^{-1}\|_{\mathcal{L}(Y,X)}$$

which completes the proof □

# Condition Numbers and Residuals

Let  $\mathbf{A} \in \mathbb{K}^{n \times n} \equiv \mathcal{L}(\mathbb{K}^n, \mathbb{K}^n)$  be nonsingular, let  $\|\cdot\|$  be any norm on  $\mathbb{K}^n$ , and  $\kappa(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^{-1}\|$  the corresponding condition number.

**Question:** suppose that  $\tilde{\mathbf{x}}$  is an approximate solution of  $\mathbf{A}\mathbf{x} = \mathbf{b}$ . The **residual**  $\mathbf{r} := \mathbf{b} - \mathbf{A}\tilde{\mathbf{x}}$  is a computable and hence known quantity. What does the residual tell us about the unknown error  $\tilde{\mathbf{x}} - \mathbf{x}$ ?

$$\frac{\|\mathbf{x} - \tilde{\mathbf{x}}\|}{\|\mathbf{x}\|} \leq \kappa(\mathbf{A}) \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|} \quad (5.11)$$

Hence, the smaller  $\kappa(\mathbf{A})$  the more accurate is the information provided by the (known) residual about the (unknown) error. To see (5.11), note

$$\|\mathbf{x} - \tilde{\mathbf{x}}\| = \|\mathbf{A}^{-1} \mathbf{A}(\mathbf{x} - \tilde{\mathbf{x}})\| = \|\mathbf{A}^{-1}(\mathbf{b} - \mathbf{A}\tilde{\mathbf{x}})\| \leq \|\mathbf{A}^{-1}\| \|\mathbf{b} - \mathbf{A}\tilde{\mathbf{x}}\| = \|\mathbf{A}^{-1}\| \|\mathbf{r}\|, \quad \|\mathbf{b}\| \leq \|\mathbf{A}\| \|\mathbf{x}\|$$

which gives (5.11).

Why does this imply immediately also

$$\frac{\|\mathbf{b} - \mathbf{A}\tilde{\mathbf{x}}\|}{\|\mathbf{b}\|} = \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|} \leq \kappa(\mathbf{A}) \frac{\|\mathbf{x} - \tilde{\mathbf{x}}\|}{\|\mathbf{x}\|} ? \quad (5.12)$$

# Further Comments

- Applying an operator  $L \in \mathcal{L}(\mathbb{X}, \mathbb{Y})$  to some input  $x \in \mathbb{X}$ , i.e.,  $L : x \mapsto L(x)$  (e.g. a matrix vector multiplication  $\mathbf{A} : \mathbf{x} \mapsto \mathbf{Ax}$ ) and solving an “operator equation”  $L(x) = y$ , i.e.,  $y \mapsto L^{-1}(y)$  ( $\mathbf{x} \mapsto \mathbf{A}^{-1}\mathbf{b}$ ) are mathematical operations with the same condition numbers.
- A **linear** operator on a **finite dimensional** space is always **bounded** (hence continuous). In fact, on account of Proposition 7, this has to be verified only for matrices and some fixed norm. For instance, the maximal row-sum norm is trivially bounded. Whenever  $\mathbf{A}$  is non-singular  $\mathbf{A}^{-1}$  has for the same reason a finite norm as well. Therefore,

$$\mathbf{Ax} = \mathbf{b} \quad \text{is a well-posed problem} \quad \Leftrightarrow \mathbf{A}^{-1} \text{ exists} \quad \Leftrightarrow \kappa(\mathbf{A}) < \infty \quad (\text{any norm})$$

# Stability of a Basis

Suppose that  $\mathbb{X}$  is a finite-dimensional normed linear space and  $\Phi = \{\phi_1, \dots, \phi_n\} \subset \mathbb{X}$  is a basis for  $\mathbb{X}$ . We wish to quantify the quality of the basis in the following sense: suppose we know only perturbed coefficients  $\tilde{x}_j$ ,  $j = 1, \dots, n$  of the representation  $x = \sum_{j=1}^n x_j \phi_j$ , how much does  $\tilde{x} = \sum_{j=1}^n \tilde{x}_j \phi_j$  differ from  $x$ ? If we decide to measure the perturbation of the coefficient vector in the norm  $\|\cdot\|_p$  on  $\mathbb{K}^n$ , say, we ask to estimate the relative error  $\frac{\|\tilde{x} - x\|_{\mathbb{X}}}{\|x\|_{\mathbb{X}}}$  in terms of the relative error

$$\frac{\|\tilde{x} - x\|_p}{\|x\|_p} = \frac{\|B_{\Phi}(\tilde{x}) - B_{\Phi}(x)\|_p}{\|B_{\Phi}(x)\|_p} \stackrel{???}{\leftrightarrow} \frac{\|\tilde{x} - x\|_{\mathbb{X}}}{\|x\|_{\mathbb{X}}} \quad (5.13)$$

Thus, the relation between these relative errors is precisely described by the condition number  $\kappa_{\mathbb{X}, \ell_p^n}(B_{\Phi})$  of the mapping  $B_{\Phi} : \mathbb{X} \rightarrow \mathbb{K}^n$ , induced by (3.1). By the previous comments,  $B_{\Phi}$  is a linear and bijective mapping between finite dimensional spaces and thus has a bounded condition number. Hence, the basis  $\Phi$  is “the more stable” the smaller  $\kappa_{\mathbb{X}, \ell_p^n}(B_{\Phi})$  is. Stability of a basis  $\Phi$  (in a quantitative sense) is therefore synonymous to a small condition number of the coordinate mapping  $B_{\Phi}$ .

# Finite vs. Infinite Dimensional Spaces

$$\ell_p(\mathcal{I}) := \{\mathbf{x} \in \mathbf{K}^{\mathcal{I}} : \|\mathbf{x}\|_{\ell_p(\mathcal{I})} < \infty\}, \quad B_r(\ell_p(\mathcal{I})) := \{\mathbf{x} \in \mathbf{K}^{\mathcal{I}} : \|\mathbf{x}\|_{\ell_p(\mathcal{I})} \leq r\} \text{ ball of radius } r.$$

There are essential differences between  $\#\mathcal{I} = n < \infty$  and  $\#\mathcal{I} = \infty$ :  $(\mathbb{K}^{\mathcal{I}} = (\mathbb{K} \setminus \{\infty\})^{\mathcal{I}})$

- $\|\mathbf{x}\|_p < \infty$  for all  $\mathbf{x} \in \mathbb{K}^n$ . For  $\#\mathcal{I} = \infty$  there exist  $\mathbf{x} \in \mathbb{K}^{\mathcal{I}}$  for which  $\|\mathbf{x}\|_{\ell_p(\mathcal{I})} = \infty$ , i.e.,  $\ell_p(\mathcal{I})$  is a strict subset of  $\mathbb{K}^{\mathcal{I}}$  and  $\ell_p(\mathcal{I})$  differs from  $\ell_q(\mathcal{I})$  when  $p \neq q$ .
- For  $\#\mathcal{I} = \infty$  closed balls  $B_r(\ell_p(\mathcal{I}))$  are not compact (they are for  $\#\mathcal{I} < \infty$ ).
- Visualize  $B_1(\ell_p^2)$  for different  $p$ , including  $p = 1, 2, \infty$
- Determine the volume of  $B_1(\ell_p^n)$  for  $p = 1, 2, \infty$ .
- If you draw  $N$  points in  $[-1, 1]^n = B_1(\ell_\infty^n)$  randomly according to the uniform distribution. How many of those points do you expect to find on average in the Euclidean ball  $B_1(\ell_2^n)$  and in the  $\ell_1$ -ball  $B_1(\ell_1^n)$ ?
- Linear operators on infinite dimensional spaces are no longer automatically bounded.

# Finite vs. Infinitely Dimensional Spaces

How different can different norms be?

## Proposition 17

Let  $\mathbb{V}$  be  $\mathbb{K}$ -vector space. If  $\dim(\mathbb{V}) = n < \infty$  then *all* norms on  $\mathbb{V}$  are *equivalent*, i.e., for any two norms  $\|\cdot\|_*$ ,  $\|\cdot\|_{**}$  on  $\mathbb{V}$  there exists constants  $0 < c, C < \infty$  such that

$$c\|v\|_* \leq \|v\|_{**} \leq C\|v\|_*, \quad \forall v \in \mathbb{V}. \quad (5.14)$$

This is in essence a consequence of the fact that any closed bounded subsets of a finite-dimensional space is *compact* and that continuous functions on compact sets attain their extrema in those sets.

In particular, recall:

## Theorem (Heine-Borel):

All bounded closed sets in  $\mathbb{K}^n$  are compact.

**Sketch of Proof of Prop. 17:** Note that (5.14) is equivalent to

$$c \leq \|v\|_{**} \leq C \quad \forall v \in B_1((\mathbb{V}, \|\cdot\|_*)). \quad (5.15)$$

Now, fix a basis  $\{\phi_1, \dots, \phi_n\}$  for  $\mathbb{V}$ . As shown earlier every  $v \in \mathbb{V}$  has a unique coefficient vector  $\mathbf{v} = \Phi(v) \in \mathbb{K}^n$ . Define  $\|v\|_n := \|\mathbf{v}\|_\infty$ . Then it suffices to show that (5.15) holds for any norm  $\|\cdot\|_{**}$  on  $\mathbb{V}$  and  $\|\cdot\|_* = \|\cdot\|_n$ . To that end, recall

Moreover, we know that every continuous function attains its minimum and maximum on a compact set. The function  $F(\mathbf{v}) := \left\| \sum_{j=1}^n v_j \phi_j \right\|_{**} = \|B_\Phi^{-1}(\mathbf{v})\|_{**}$  is clearly continuous because  $B_\Phi^{-1} : \mathbf{v} \rightarrow \sum_{j=1}^n v_j \phi_j$  is easily seen to be continuous and  $v \rightarrow \|v\|_{**}$  is by (5.1) even Lipschitz continuous. Since by Heine-Borel,  $B_1(\ell_\infty^n)$  is compact,  $F(\mathbf{v})$  attains its minimum and maximum on  $B_1(\ell_\infty^n)$ . Since by linear independence,  $B_\Phi^{-1}(\mathbf{v}) \neq 0$  iff  $\mathbf{v} \neq \mathbf{0}$  (N1) implies that  $c := \min_{\mathbf{v} \in B_1(\ell_\infty^n)} \|B_\Phi^{-1}(\mathbf{v})\|_{**} > 0$ . The upper bound is handled analogously.  $\square$



# Finite vs. Infinitely Dimensional Spaces

## Remark 18

The statement in Proposition 17 needs to be read with caution. The equivalence constants  $c, C$  in (5.14) *depend* on the dimension  $n$ . For instance,

$$\|\mathbf{x}\|_\infty = \max_{j=1,\dots,n} |x_j| \leq \left( \sum_{j=1}^n |x_j|^2 \right)^{1/2} = \|\mathbf{x}\|_2 \leq \sqrt{n} \|\mathbf{x}\|_\infty.$$

Hence, the larger the dimension the more the choice of a norm matters.

In infinite dimensions there could be issues with *convergence*. In that regard, the  $\ell_p$ -spaces are not so bad in the following sense

## Definition 19

A normed linear space  $(\mathbb{V}, \|\cdot\|)$  is called *complete* if every *Cauchy sequence* in  $\mathbb{V}$  converges to an element in  $\mathbb{V}$ , i.e.,

$$\|v_k - v_n\| \rightarrow 0, \quad k, n \rightarrow \infty, \quad \Rightarrow \quad \exists v \in \mathbb{V}, \text{ s.t. } \|v_k - v\| \rightarrow 0, \quad k \rightarrow \infty.$$

Complete normed linear spaces are called *Banach spaces*.

# Norms on Function Spaces

## Remark 20

Obviously, completeness is useful to be sure that, if one constructs a sequence, the Cauchy property (which one may be able to verify) already ensures that this sequence has a limit in the space one is working in.

- All finite dimensional spaces, endowed with any norm, are always complete.
- $\ell_p(\mathcal{I})$  are Banach spaces.

Things get more subtle when dealing with function spaces. One can formally define again for  $1 \leq p \leq \infty$  and a given function  $f : \Omega \subset \mathbb{R}^n \rightarrow \mathbb{K}$ , say:

$$\|f\|_{L_p(\Omega)} := \left( \int_{\Omega} |f(x)|^p dx \right)^{1/p}, \quad 1 \leq p < \infty, \quad \|f\|_{L_{\infty}(\Omega)} := \sup_{x \in \Omega} |f(x)|. \quad (5.16)$$

This is an obvious analogy to sequence norms replacing discrete arguments  $j \in \mathcal{I}$  and summation by continuous arguments  $x$  and integration.

Again, the verification of properties (N1), (N2) (and (N3) when  $p = \infty$ ) is straight forward. For  $1 < p < \infty$  the triangle inequality (N3) - also referred to as Minkowski's inequality, rests again the continuous version of [Hölder's lequality](#)

$$\left| \int_{\Omega} f(x) \overline{g(x)} dx \right| \leq \|f\|_{L_p(\Omega)} \|g\|_{L_{p^*}(\Omega)}, \quad \frac{1}{p} + \frac{1}{p^*}. \quad (5.17)$$

# $L_p$ -Spaces Continued

Can be skipped, just for interested readers

When  $\Omega$  is a bounded domain and  $f$  is continuous, the expressions (5.16) are well defined and  $C(\Omega; \mathbb{K}) := \{f : \Omega \rightarrow \mathbb{K} : f \text{ continuous in } \Omega\}$ , endowed with any of the norms  $\|\cdot\|_{L_p(\Omega)}$  becomes a normed linear space. However, completeness is an issue.

## Remark 21

*Consider  $\Omega = (-1, 1) \subset \mathbb{R}$ ,  $f_n(x) := 1$ ,  $x \in (-1, 0]$ ,  $f_n(x) := \max\{0, 1 - xn\}$ . One easily checks that this is a Cauchy sequence for every  $1 \leq p < \infty$  but not for  $p = \infty$ . All  $f_n$  are continuous, but the  $f_n$  don't have a limit in  $C((-1, 1))$ . In fact, the pointwise limit is  $f(x) = 1$ ,  $x \in (-1, 0]$ ,  $f(x) = 0$ ,  $x \in (0, 1)$  which is also the limit in  $L_p((-1, 1))$  for  $1 \leq p < \infty$ . This hints at the facts:*

- $(C(\Omega), L_p(\Omega))$  is not a complete normed linear space for  $p < \infty$ ;
- $(C(\Omega), \|\cdot\|_{L_\infty(\Omega)})$  is complete and hence a Banach space.

Regarding a suitable notion of normed linear spaces for  $1 \leq p < \infty$ , completeness is an issue. First, one should employ the right notion of integration is used. The conceptually simple Riemann integration will not lead to complete spaces. Instead, integration is always understood in the [Lebesgue](#) sense based on Measure Theory. Starting from measurable sets and measurable functions, the space  $L_p(\Omega)$  is defined as the closure of step functions with respect to the above norms. So these spaces are complete by definition. Then the following has to be kept in mind:

- The elements of the spaces  $(L_p(\Omega), \|\cdot\|_{L_p(\Omega)})$  are, strictly speaking only [equivalence classes](#) of functions, where elements of one class differ only on sets of measure zero. For instance, points have measure zero in  $\mathbb{R}$ , points, lines, curves have measure zero in  $\mathbb{R}^2$ , etc. Therefore, it does not make sense to ask for point values of elements in  $L_p(\Omega)$ . Keeping this in mind, we simply speak of "functions" in  $L_p(\Omega)$ .
- In the definition of  $L_\infty(\Omega)$  which is a strictly larger Banach space than  $C(\Omega)$  with the same norm.

# Pre-Hilbert Spaces

Norms induced by scalar products play an important role:

## Remark 22

Suppose  $\mathbb{V}$  is  $\mathbb{K}$ -vector space with a scalar product  $\langle \cdot, \cdot \rangle_{\mathbb{V}}$ . Then

$$\|v\|_{\mathbb{V}} := \langle v, v \rangle_{\mathbb{V}}^{1/2} \quad (5.18)$$

is a norm on  $\mathbb{V}$ . A linear space with a scalar product and associated norm is called a *Pre-Hilbert space*. If  $\mathbb{V}$  is complete under this norm it is called a *Hilbert space*. Hilbert spaces are in some sense closest to finite dimensional Euclidean spaces.

Properties (N1), (N2) of a norm follow directly from the properties of a scalar product. The triangle inequality (N3) follows from the *Cauchy-Schwarz Inequality*

$$|\langle v, w \rangle_{\mathbb{V}}| \leq \|v\|_{\mathbb{V}} \|w\|_{\mathbb{V}}, \quad v, w \in \mathbb{V}. \quad (5.19)$$

In fact,

$$\begin{aligned} \|v + w\|_{\mathbb{V}}^2 &= \langle v + w, v + w \rangle_{\mathbb{V}} = \langle v, v \rangle_{\mathbb{V}} + \langle w, w \rangle_{\mathbb{V}} + \langle v, w \rangle_{\mathbb{V}} + \langle w, v \rangle_{\mathbb{V}} \\ &= \|v\|_{\mathbb{V}}^2 + \|w\|_{\mathbb{V}}^2 + \langle v, w \rangle_{\mathbb{V}} + \langle w, v \rangle_{\mathbb{V}} \leq \|v\|_{\mathbb{V}}^2 + \|w\|_{\mathbb{V}}^2 + 2\|v\|_{\mathbb{V}}\|w\|_{\mathbb{V}} \\ &= (\|v\|_{\mathbb{V}} + \|w\|_{\mathbb{V}})^2 \end{aligned}$$

# Proof of the Cauchy Schwarz Inequality

Assume that  $v, w \neq 0$  (otherwise the inequality is trivial). Then for any  $\lambda \in \mathbb{K}$

$$0 \leq \langle v - \lambda w, v - \lambda w \rangle_V = \|v\|_V^2 + |\lambda|^2 \|w\|_V^2 - \lambda \langle w, v \rangle_V - \bar{\lambda} \langle v, w \rangle_V.$$

Choose  $\lambda = \frac{\langle v, w \rangle_V}{\|w\|_V^2}$  to obtain

$$0 \leq \|v\|_V^2 + \frac{|\langle w, v \rangle_V|^2}{\|w\|_V^4} \|w\|_V^2 - 2 \frac{|\langle w, v \rangle_V|^2}{\|w\|_V^2} = \|v\|_V^2 - \frac{|\langle w, v \rangle_V|^2}{\|w\|_V^2},$$

which implies (5.19). □

# Examples

- $(\mathbb{V}, \|\cdot\|) = (\mathbb{K}^n, \|\cdot\|_2)$ ,

$$\langle \mathbf{x}, \mathbf{y} \rangle_n = \mathbf{y}^* \mathbf{x}, \quad \|\mathbf{x}\|_2^2 = \mathbf{x}^* \mathbf{x} = \sum_{j=1}^n |x_j|^2$$

- $(\ell_2(\mathbb{N}), \|\cdot\|_{\ell_2(\mathbb{N})})$

$$\langle \mathbf{x}, \mathbf{y} \rangle_{\mathbb{N}} = \sum_{j \in \mathbb{N}} x_j \bar{y}_j = \mathbf{y}^* \mathbf{x}$$

- $(L_2(\Omega), \|\cdot\|_{L_2(\Omega)})$ ,

$$\langle f, g \rangle_{L_2(\Omega)} = \int_{\Omega} f(x) \overline{g(x)} dx, \quad \|f\|_{L_2(\Omega)}^2 = \int_{\Omega} |f(x)|^2 dx$$

- $(L_{2,\pi}, \|\cdot\|_{L_2(-\pi,\pi)})$ ,  $2\pi$ -periodic square integrable functions

$$\langle f, g \rangle_{(-\pi,\pi)} := \frac{1}{2\pi} \int_{-\pi}^{\pi} f(x) \overline{g(x)} dx,$$

$$e_k(x) := e^{ikx}, \quad \hat{f}(k) := \langle f, e_k \rangle_{(-\pi,\pi)} = \frac{1}{2\pi} \int_{-\pi}^{\pi} f(x) e^{-ikx} dx \quad (\text{Fourier coefficients})$$

# Dual Spaces

Let  $(\mathbb{X}, \|\cdot\|_{\mathbb{X}})$  be a Banach space. The collection of all **bounded linear functionals**  $\mathbb{X}^* := \mathcal{L}(\mathbb{X}, \mathbb{K})$  is also a Banach space under the norm

$$\|g\|_{\mathbb{X}^*} := \sup_{x \in \mathbb{X} \setminus \{0\}} \frac{g(x)}{\|x\|_{\mathbb{X}}}. \quad (5.20)$$

Examples:

- $C^k(\Omega) :=$  collection of all continuously differentiable functions on  $\Omega$ ,  
 $\|f\|_{C^k(\Omega)} := \max_{0 \leq j \leq k} \|f^{(j)}\|_{L^\infty(\Omega)}$ ; then  $\delta_{x_0}^{(k)} : f \mapsto f^{(k)}(x_0)$  belongs to  $(C^k(\Omega))^*$  because

$$\|\delta_{x_0}^{(k)}\|_{(C^k(\Omega))^*} = \sup_{f \in C^k(\Omega)} \frac{\delta_{x_0}^{(k)}(f)}{\|f\|_{C^k(\Omega)}} = \sup_{f \in C^k(\Omega)} \frac{f^{(k)}(x_0)}{\|f\|_{C^k(\Omega)}} \leq 1$$

- $\mathbf{A} \in \mathbb{K}^{m \times n}$ ,  $\mathbf{y} \in \mathbb{K}^n$  fixed, then  $g(\mathbf{x}) := (\mathbf{A}\mathbf{y})^* \mathbf{x}$  belongs to  $(\ell_2^m)^* = \mathcal{L}((\mathbb{K}^m, \|\cdot\|_2), \mathbb{K})$  and  $\|g\|_{(\ell_2^m)^*} = \|\mathbf{A}\mathbf{y}\|_2$

## Theorem 23

**Riesz Representation Theorem:** let  $(\mathbb{V}, \|\cdot\|_{\mathbb{V}})$  be a Hilbert space. There exists a linear mapping  $R = R_{\mathbb{V}^* \rightarrow \mathbb{V}} : \mathbb{V}^* \rightarrow \mathbb{V}$  such that for any  $g \in \mathbb{V}^*$ ,  $v \in \mathbb{V}$ , one has  $g(v) = \langle v, Rg \rangle_{\mathbb{V}}$  and  $\|Rg\|_{\mathbb{V}} = \|g\|_{\mathbb{V}^*}$ , i.e.,  $R$  is an isometry  $\|R\|_{\mathcal{L}(\mathbb{V}^*, \mathbb{V})} = 1$ .

# Orthogonal Projections and Best Approximations

Given a (finite dimensional) subspace  $\mathbb{U}$  of a normed linear space  $\mathbb{V}$  and given any  $v \in \mathbb{V}$ , the problem of **best approximation** to  $v$  from  $\mathbb{U}$  is to find an element  $u(v) \in \mathbb{U}$  that is closest to  $v$  with respect to a given norm. In general, this is a difficult problem but when  $\mathbb{V}$  is a Hilbert space, it amounts to a **linear projection**.

## Theorem 24

Let  $(\mathbb{V}, \|\cdot\|_{\mathbb{V}})$  be a Hilbert space and let  $\mathbb{U} \subset \mathbb{V}$  be a finite dimensional subspace (this holds in greater generality). Given any  $v \in \mathbb{V}$ , then some  $u(v) \in \mathbb{U}$  satisfies

$$\|v - u(v)\|_{\mathbb{V}} = \min_{u \in \mathbb{U}} \|v - u\|_{\mathbb{V}} \quad (5.21)$$

if and only if

$$\langle v - u(v), u \rangle_{\mathbb{V}} = 0 \quad \forall u \in \mathbb{U} \quad \text{i.e.,} \quad v - u(v) \perp \mathbb{U} \quad (5.22)$$

## Remark 25

- 1 If  $u \in \mathbb{U}$  and  $\langle u, w \rangle_{\mathbb{V}} = 0, \forall w \in \mathbb{U}$  then  $u = 0$ , i.e., the only element in a linear space that is orthogonal to all elements in the same space is the zero element. In fact, for  $w = u$  one has  $0 = \langle u, u \rangle_{\mathbb{V}} = \|u\|_{\mathbb{V}}^2$  which, by (N1) implies  $u = 0$ .
- 2 Given  $v \in \mathbb{V}$  there exists at most one  $\bar{u} \in \mathbb{U}$  such that  $\langle v - \bar{u}, u \rangle_{\mathbb{V}} = 0, \forall u \in \mathbb{U}$ . In fact, suppose  $u_1, u_2 \in \mathbb{U}$  have that property. Then  $0 = \langle v - u_1, w \rangle_{\mathbb{V}} - \langle v - u_2, w \rangle_{\mathbb{V}} = \langle u_2 - u_1, w \rangle_{\mathbb{V}}, \forall w \in \mathbb{U}$ , so, by (1),  $u_1 = u_2$ .



# Proof of Theorem 24

The proof uses a standard variational argument. Consider a candidate  $u \in \mathbb{U}$  for (5.21) and any  $t \in \mathbb{R}$ ,  $w \in \mathbb{U}$ . We wish to see under which circumstances a perturbation  $u + tw \in \mathbb{U}$  could do better.

$$\begin{aligned} \|v - (u + tw)\|_{\mathbb{V}}^2 &= \langle (v - u) - tw, (v - u) - tw \rangle_{\mathbb{V}} = \|v - u\|_{\mathbb{V}}^2 + t^2 \|w\|_{\mathbb{V}}^2 \\ &\quad - 2t \operatorname{Re}(\langle v - u, w \rangle_{\mathbb{V}}) \end{aligned} \quad (5.23)$$

(5.21)  $\Rightarrow$  (5.22): Suppose  $u = u(v)$  is the minimizer but there exists a  $w \in \mathbb{U}$ ,  $w \neq 0$  such that  $\langle v - u, w \rangle_{\mathbb{V}} = \alpha \neq 0$ , and hence  $\langle v - u, iw \rangle_{\mathbb{V}} \neq 0$ . Without loss of generality we can assume  $\beta := \operatorname{Re}(\alpha) > 0$ ,  $\|w\|_{\mathbb{V}} = 1$ . Then, by optimality of  $u = u(v)$  we must have

$$0 < \|v - (u + tw)\|_{\mathbb{V}}^2 - \|v - u\|_{\mathbb{V}}^2 = t^2 - 2t\beta \quad \forall t \in \mathbb{R}.$$

But choosing  $t = \beta$ , yields  $\beta^2 - 2\beta^2 < 0$ , which is a contradiction.

(5.22)  $\Rightarrow$  (5.21): By (5.23) we have in this case  $\|v - (u + w)\|_{\mathbb{V}}^2 - \|v - u\|_{\mathbb{V}}^2 = \|w\|_{\mathbb{V}}^2$  for all  $w \in \mathbb{U}$ , which implies (5.21). □

# Orthogonal Projections and Best Approximations

## Remark 26

*Theorem 24 doesn't explicitly state existence of a minimizer only the equivalence of (5.21) and (5.22). Existence of a minimizer in (5.21) can be argued as follows: the search can be restricted to a ball in  $\mathbb{U}$ . Since  $\mathbb{U}$  is finite dimensional this ball is compact. Since  $F(u) = \|v - u\|_{\mathbb{V}}$  is a continuous function it attains its minimum. The uniqueness of the minimizer, by the equivalence of (5.21) and (5.22), follows from Remark 25, (1).*

## Theorem 27

*Let  $(\mathbb{V}, \|\cdot\|_{\mathbb{V}})$  be a Hilbert space and  $\mathbb{U} \subset \mathbb{V}$  a finite dimensional subspace. For every  $v \in \mathbb{V}$  there exists a unique  $u = u(v) \in \mathbb{U}$  such that*

$$\langle v - u(v), w \rangle_{\mathbb{V}} = 0, \quad \forall w \in \mathbb{U}, \quad (5.24)$$

*i.e., the difference  $v - u(v)$  is perpendicular to the subspace  $\mathbb{U}$ . Hence,  $P_{\mathbb{U}} : \mathbb{V} \rightarrow \mathbb{U}$  given by  $P_{\mathbb{U}}v = u(v)$  is well-defined and has the following properties:*

- 1  $P_{\mathbb{U}}$  is linear and idempotent  $P_{\mathbb{U}} \circ P_{\mathbb{U}} = P_{\mathbb{U}}$ , i.e.,  $P_{\mathbb{U}}$  is a **projector**.
- 2  $P_{\mathbb{U}}$  is **self-adjoint**, i.e.,  $\langle P_{\mathbb{U}}v, z \rangle_{\mathbb{V}} = \langle v, P_{\mathbb{U}}z \rangle_{\mathbb{V}}$ ,  $v, z \in \mathbb{V}$ .
- 3  $\|P_{\mathbb{U}}\|_{\mathcal{L}(\mathbb{V}, \mathbb{V})} = 1$ .

# Proof of Theorem 27

$P_U$  is well-defined: This follows from Remark 25, (2) and Remark 26.

(1): It is also linear because for all  $w \in U$  one has for all  $w \in U$

$$\begin{aligned} \langle P_U(v_1 + v_2) - (P_U v_1 + P_U v_2), w \rangle_V &= \langle P_U(v_1 + v_2) - (v_1 + v_1), w \rangle_V + \langle (v_1 + v_1) \\ &\quad - (P_U v_1 + P_U v_2), w \rangle_V \\ &= 0 + \langle v_1 - P_U v_1, w \rangle_V + \langle v_2 - P_U v_2, w \rangle_V = 0 \end{aligned}$$

Since both  $P_U(v_1 + v_2)$  and  $P_U v_1 + P_U v_2$  belong to  $U$ , Remark 25, (1), implies

$P_U(v_1 + v_2) = P_U v_1 + P_U v_2$ . In the same way one verifies  $P_U(\alpha v) = \alpha P_U v$  which shows the first part of (1). Moreover,  $0 = \langle P_U v - P_U(P_U v), w \rangle_V$  which, by Remark 25, (1), again shows that  $P_U v = P_U(P_U v)$  which is (1).

(2):  $0 = \langle P_U v, z - P_U z \rangle_V \rightsquigarrow \langle P_U v, z \rangle_V = \langle P_U v, P_U z \rangle_V$ . Likewise,  $0 = \langle v - P_U v, P_U z \rangle_V \rightsquigarrow \langle v, P_U z \rangle_V = \langle P_U v, P_U z \rangle_V$ . Thus  $\langle P_U v, z \rangle_V = \langle v, P_U z \rangle_V$ .

(3):  $\|P_U v\|_V^2 = \langle P_U v, P_U v \rangle_V \stackrel{(1),(2)}{=} \langle v, P_U v \rangle_V \stackrel{CS}{\leq} \|P_U v\|_V \|v\|_V \Rightarrow \|P_U\|_{\mathcal{L}(V,V)} \leq 1$ ,  
 $\|P_U u\|_U \stackrel{(1)}{=} \|u\|_U \Rightarrow (3)$ . □

**Pythagoras' Theorem:**

$$\begin{aligned} \|v - P_U v\|_V^2 &= \langle v - P_U v, v - P_U v \rangle_V = \langle v, v \rangle_V + \langle P_U v, P_U v \rangle_V - \langle P_U v, v \rangle_V - \langle v, P_U v \rangle_V \\ &= \|v\|_V^2 + \|P_U v\|_V^2 - 2\|P_U v\|_V^2 = \|v\|_V^2 - \|P_U v\|_V^2 \end{aligned}$$

# Computing Orthogonal Projections/Best Approximations

Suppose  $\Phi = \{\phi_1, \dots, \phi_n\}$  is a basis for  $\mathbb{U}$ . Since

$$\langle v - P_{\mathbb{U}}v, u \rangle_{\mathbb{V}} = 0, \quad \forall u \in \mathbb{U} \quad \Leftrightarrow \quad \langle v - P_{\mathbb{U}}v, \phi_k \rangle_{\mathbb{V}} = 0, \quad \text{for } k = 1, \dots, n,$$

substituting  $P_{\mathbb{U}}v = \sum_{j=1}^n u_j \phi_j$ , yields (see (4.1))

$$\sum_{j=1}^n u_j \langle \phi_j, \phi_k \rangle_{\mathbb{V}} = \langle v, \phi_k \rangle_{\mathbb{V}} =: b_k, \quad k = 1, \dots, n, \quad \Leftrightarrow \quad \mathbf{G}_{\Phi} \mathbf{u} = \mathbf{b} \quad (5.25)$$

where  $\mathbf{G}_{\Phi} = (\langle \phi_k, \phi_j \rangle_{\mathbb{V}})_{j,k=1}^n$  is the **Gramian matrix** associated with  $\Phi$ . We already know that  $\mathbf{G}_{\Phi}$  is **hermitian positive definite** and hence non-singular, so that the computing  $P_{\mathbb{U}}v$  amounts to solving a linear system of equations for the unknown expansion coefficient vector  $\mathbf{u} = (u_1, \dots, u_n)^{\top}$ .

**Orthonormal Bases:** The orthogonal projection is very easy to compute if the basis  $\Phi$  is **orthonormal**, i.e.,  $\langle \phi_j, \phi_k \rangle_{\mathbb{V}} = \delta_{j,k}$ ,  $j, k = 1, \dots, n$ . In fact, then  $\mathbf{G}_{\Phi} = \mathbf{I}$  and

$$u_j = \langle v, \phi_j \rangle_{\mathbb{V}}, \quad j = 1, \dots, n. \quad (5.26)$$

**Example:**  $\mathbb{V} = L_{2,\pi}$ , show that the  $e_k(x) = e^{ikx}$ ,  $k \in \mathbb{Z}$ , form an orthonormal system respect to  $\langle f, g \rangle = \frac{1}{2\pi} \int_{-\pi}^{\pi} f(x) \overline{g(x)} dx$ , so that orthogonal projection to the space spanned by the  $2n + 1$  first harmonics  $e_k$ ,  $-n \leq k \leq n$ , is the **Fourier partial sum**

$$S_n(f; x) = \sum_{|k| \leq n} \hat{f}(k) e^{ikx}, \quad \hat{f}(k) = \langle f, e_k \rangle, \quad |k| \leq n.$$

# Gram-Schmidt Orthogonalization

Let  $\mathbb{U}$  be an  $n$ -dimensional linear space with scalar product  $\langle \cdot, \cdot \rangle$  and norm  $\| \cdot \| = \langle \cdot, \cdot \rangle^{1/2}$ . Given a basis  $\{\phi_1, \dots, \phi_n\}$  of  $\mathbb{U}$  one can always generate an **orthonormal** basis  $\{\psi_1, \dots, \psi_n\}$  as follows:

- normalize  $\psi_1 := \frac{\phi_1}{\|\phi_1\|}$ ;
- while  $k < n$ , given an orthonormal basis  $\{\psi_1, \dots, \psi_k\}$  for  $\text{span}\{\phi_1, \dots, \phi_k\}$ , let

$$\psi'_{k+1} := \phi_{k+1} - \sum_{j=1}^k \langle \phi_{k+1}, \psi_j \rangle \psi_j$$

and normalize

$$\psi_{k+1} := \frac{\psi'_{k+1}}{\|\psi'_{k+1}\|}.$$

Show that  $\{\psi_1, \dots, \psi_n\}$ , generated in this fashion, is indeed an orthonormal basis of  $\mathbb{U}$ .

# Orthonormal Bases Continued

## Remark 28

An orthonormal basis is “optimally stable” in the sense of (5.13), i.e.,

$$\kappa_{\mathbb{U}, \ell_2^n}(B_{\Psi}) = 1. \quad (5.27)$$

In fact, for any  $u = \sum_{j=1}^n u_j \psi_j \in \mathbb{U}$ ,  $\Psi = \{\psi_1, \dots, \psi_n\}$  orthonormal, one has

$$\|u\|^2 = \langle u, u \rangle = \left\langle \sum_{j=1}^n u_j \psi_j, \sum_{k=1}^n u_k \psi_k \right\rangle = \mathbf{u}^* \mathbf{G}_{\Psi} \mathbf{u} = \mathbf{u}^* \mathbf{u} = \|\mathbf{u}\|_2^2$$

which implies  $\|B_{\Psi}^{-1}(\mathbf{u})\| = \|B_{\Psi}(u)\|_2 = 1$ .

Consider the monomial basis  $\phi_j(x) := x^j$ ,  $j = 0, \dots, n$  of  $\mathbb{U} := \mathbb{P}_n$  equipped with the scalar product  $\langle f, g \rangle := \int_0^1 f(x)g(x)dx$  ( $\mathbb{K} = \mathbb{R}$ ). What can you say about the condition number of  $B_{\Phi}$  in this case. Intuitively, it is very large for large  $n$ . What does this mean about the Gram-Schmidt process turning  $\Phi$  into an orthonormal basis  $\Psi$  for  $\mathbb{P}_n$ ? Since this is a change of bases it can be represented by a matrix whose condition describes the stability of the change of bases.

# Eigenvectors, Eigenvalues

Let  $\mathbb{X}$  be a  $\mathbb{K}$ -vector space. Then  $\lambda \in \mathbb{K}$ ,  $x \in \mathbb{X}$  are called eigenvalue, resp. (right-)eigenvector of  $L \in \mathcal{L}(\mathbb{X}, \mathbb{X})$  if

$$Lx = \lambda x \quad \text{We also say } (\lambda, x) \text{ is an } \text{eigenpair} \text{ of } L. \text{ Also } (\lambda^{-1}, x) \text{ is an eigenpair of } L^{-1} \text{ if it exists} \quad (6.1)$$

## Proposition 29

*When  $\mathbb{X}$  is finite dimensional there always exists an eigenpair  $(\lambda, x) \in \mathbb{C} \times \mathbb{C}^n$ .  $x$  is determined only up to normalization.*

**Proof:** Fix bases  $\Phi, \Psi$  for  $\mathbb{X}, \mathbb{Y}$  as before. Then, by Proposition 7 and (3.2) we have

$$L(x) = \lambda x \Leftrightarrow \mathbf{C}x = \lambda x \Leftrightarrow (\mathbf{C} - \lambda \mathbf{I})x = \mathbf{0}, \Leftrightarrow \det(\mathbf{C} - \lambda \mathbf{I}) = 0. \quad (6.2)$$

By Laplace's expansion rule of determinants, it follows that is a polynomial in  $\mathbb{P}_n(\mathbb{K})$  of degree at most  $n$ . This is called the **characteristic polynomial**. By the Fundamental Theorem of Algebra, any polynomial in  $\mathbb{P}_n(\mathbb{K})$  has exactly  $n$  (possibly complex) roots (counting multiplicities).  $\square$

It therefore suffices (in the finite dimensional case) to understand eigenproblems for matrices  $\mathbf{A} \in \mathbb{K}^{n \times n}$ ,  $\mathbb{K} \in \{\mathbb{R}, \mathbb{C}\}$ :

$$\mathbf{A}x = \lambda x \quad (6.3)$$

## Theorem 30

*For  $\mathbf{A} \in \mathbb{C}^{n \times n}$  one has  $\det(\mathbf{A}) = \prod_{k=1}^n \lambda_k$ , where  $\lambda_1, \dots, \lambda_n \in \mathbb{C}$  are the  $n$  eigenvalues of  $\mathbf{A}$ .*

# Spectral Theorem for Hermitian Matrices

## Theorem 31

Let  $\mathbf{A} \in \mathbb{K}^{n \times n}$  be hermitian, then:

- 1 all eigenvalues of  $\mathbf{A}$  are real;
- 2 eigenvectors corresponding to distinct eigenvalues are pairwise orthogonal;
- 3 there exists an orthonormal basis of  $\mathbb{K}^n$ , consisting of eigenvectors.

Hence,  $\mathbf{A}$  is diagonalizable, i.e., there exists a unitary matrix  $\mathbf{U} \in \mathcal{O}_n(\mathbb{K})$  such that

$$\mathbf{U}^* \mathbf{A} \mathbf{U} = \Lambda = \text{diag}(\lambda_1, \dots, \lambda_n), \quad \lambda_j \in \mathbb{R}, j = 1, \dots, n.$$

**Proof:** (1): suppose  $\mathbf{A}\mathbf{x} = \lambda\mathbf{x} \Rightarrow \mathbf{x}^* \mathbf{A}\mathbf{x} = \lambda\mathbf{x}^* \mathbf{x}$ . Since  $\mathbf{x}^* \mathbf{x} > 0$  we also have

$$\overline{\lambda\mathbf{x}^* \mathbf{x}} = \overline{\mathbf{x}^* \mathbf{A}\mathbf{x}} = (\mathbf{x}^\top \overline{\mathbf{A}\mathbf{x}})^\top = \overline{\mathbf{x}}^\top \overline{\mathbf{A}}^\top \mathbf{x} = \mathbf{x}^* \overline{\mathbf{A}}^* \mathbf{x} = \mathbf{x}^* \mathbf{A}\mathbf{x} = \lambda\mathbf{x}^* \mathbf{x} \Rightarrow \overline{\lambda} = \lambda.$$

(2): Assume that  $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$ ,  $\mathbf{A}\mathbf{y} = \mu\mathbf{y}$  but  $\mathbf{y}^* \mathbf{x} \neq 0$ . Since  $(\mathbf{A}^*)^{-1} = \mathbf{A}^{-1}$  one has  $\mathbf{y}^* \mathbf{x} = \mathbf{y}^* \mathbf{A}^* \mathbf{A}^{-1} \mathbf{x} = \overline{\mu} \lambda^{-1} \mathbf{y}^* \mathbf{x} = \mu \lambda^{-1} \mathbf{y}^* \mathbf{x}$ , where we have used (1). Since  $\mathbf{y}^* \mathbf{x} \neq 0$ , this can only hold when  $\lambda = \mu$  which is a contradiction.

(3): By Proposition 29, there exists an eigenpair  $(\lambda_1, \mathbf{u}^1)$  ( $\|\mathbf{u}^1\|_2 = 1$ ) for  $\mathbf{A}$ . Let

$$\mathbf{V}_1 := \{\mathbf{z} \in \mathbb{K}^n : \mathbf{z}^* \mathbf{u}^1 = \langle \mathbf{u}^1, \mathbf{z} \rangle_n = 0\}$$

be the orthogonal complement of  $\text{span}\{\mathbf{u}^1\}$ . Observe next that  $\mathbf{A}$  maps  $\mathbf{V}_1$  into itself.



**Proof of (3) continued:**

In fact, when  $\mathbf{z} \in \mathbb{V}_1$

$$(\mathbf{Az})^* \mathbf{u}^1 = \mathbf{z}^* \mathbf{A}^* \mathbf{u}^1 = \mathbf{z}^* \mathbf{A} \mathbf{u}^1 = \mathbf{z}^* \lambda \mathbf{u}^1 = \lambda \mathbf{z}^* \mathbf{u}^1 = 0,$$

i.e.,  $\mathbf{z} \in \mathbb{V}_1$  implies  $\mathbf{Az} \in \mathbb{V}_1$ . Therefore, the operator  $L(\mathbf{x}) := \mathbf{Ax}$  belongs to  $\mathcal{L}(\mathbb{V}_1, \mathbb{V}_1)$ . By Proposition 29 there exists an eigenpair  $(\lambda_2, \mathbf{u}^2) \in \mathbb{C} \times \mathbb{V}_1$  and  $\langle \mathbf{u}^1, \mathbf{u}^2 \rangle_n = 0$ . Clearly  $\dim \mathbb{V}_1 = n - 1$ . Now, by the same reasoning, we consider the orthogonal complement  $\mathbb{V}_2$  to  $\mathbf{u}^2$  in  $\mathbb{V}_1$  (which is therefore also orthogonal to  $\mathbf{u}^1$ ) and find the next eigenpair  $(\lambda_3, \mathbf{u}^3) \in \mathbb{V}_2$ . Since each time the dimension of the subsequent orthogonal complement decreases by one, this process terminates after  $n - 1$  steps.  $\square$

**Remark 32**

- *Not every matrix  $\mathbf{A} \in \mathbb{K}^{n \times n}$  can be diagonalized (is similar to a diagonal matrix, meaning there exists a matrix  $\mathbf{C}$  such that  $\mathbf{C}^{-1} \mathbf{A} \mathbf{C} = \text{diag}(\lambda_1, \dots, \lambda_n)$ ). It is diagonalizable iff there exists a basis of eigenvectors.*
- *For  $\mathbb{K} = \mathbb{R}$  the eigenvectors are in  $\mathbb{R}^n$ , i.e.,  $\mathbf{U} \in \mathcal{O}_n(\mathbb{R})$  is an orthogonal matrix.*

# Spectral Theorem for Unitary Matrices

## Theorem 33

$\mathbf{Q} \in \mathcal{O}_n(\mathbb{K})$  unitary, then:

- 1 all eigenvalues of  $\mathbf{Q}$  have absolute value equal to one;
- 2 eigenvectors corresponding to distinct eigenvalues are pairwise orthogonal;
- 3 there exists an orthonormal basis of  $\mathbb{K}^n$ , consisting of eigenvectors.

Hence,  $\mathbf{Q}$  is diagonalizable, i.e., there exists a unitary matrix  $\mathbf{U} \in \mathcal{O}_n(\mathbb{K})$  such that

$$\mathbf{U}^* \mathbf{Q} \mathbf{U} = \Lambda = \text{diag}(\lambda_1, \dots, \lambda_n), \quad |\lambda_j| = 1, \quad j = 1, \dots, n.$$

**Proof:** (1): assume that  $\mathbf{Q}\mathbf{x} = \lambda\mathbf{x}$ ,  $\mathbf{x} \neq \mathbf{0}$ ,  $\rightsquigarrow \bar{\lambda}\lambda\mathbf{x}^* \mathbf{x} = (\mathbf{Q}\mathbf{x})^* \mathbf{Q}\mathbf{x} = \mathbf{x}^* \mathbf{Q}^* \mathbf{Q}\mathbf{x} = \mathbf{x}^* \mathbf{x} \Rightarrow |\lambda|^2 = 1$ .

(2): Assume that  $\mathbf{Q}\mathbf{x} = \lambda\mathbf{x}$ ,  $\mathbf{Q}\mathbf{y} = \mu\mathbf{y}$  but  $\mathbf{y}^* \mathbf{x} \neq 0 \rightsquigarrow \mathbf{y}^* \mathbf{x} = \mathbf{y}^* \mathbf{Q}^* \mathbf{Q}\mathbf{x} = \bar{\mu}\lambda\mathbf{y}^* \mathbf{x} \Rightarrow \bar{\mu}\lambda = 1 \Rightarrow \lambda/\mu = 1$  which is a contradiction to  $\mu \neq \lambda$ .

(3): Let  $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$  and  $\mathbb{V}_1 := \{\mathbf{z} \in \mathbb{K}^n : \mathbf{z}^* \mathbf{x} = 0\}$ . Now for any  $\mathbf{z} \in \mathbb{V}_1$  one has

$(\mathbf{Q}\mathbf{z})^* \mathbf{x} = \mathbf{z}^* \mathbf{Q}^* \mathbf{x} = \mathbf{z}^* \mathbf{Q}^{-1} \mathbf{x} = \lambda^{-1} \mathbf{z}^* \mathbf{x} = 0$ . Hence, the orthogonal complement of any eigenvector is an invariant subspace of  $\mathbf{Q}$ , i.e.,  $\mathbf{Q}\mathbb{V}_1 \subseteq \mathbb{V}_1$ . As in the proof of Theorem 31, one can therefore successively peel off invariant subspaces of decreasing dimension which are orthogonal to previously found eigenvectors. □

# Application: Direct Solvers for Systems of Linear Equations

**Common principle:** given  $\mathbf{A} \in \mathbb{K}^{n \times n}$  find a **factorization**  $\mathbf{A} = \mathbf{CB}$ , where both  $\mathbf{B}, \mathbf{C} \in \mathbb{K}^{n \times n}$  are “easy to invert”

then, since  $\mathbf{Ax} = \mathbf{C}(\mathbf{Bx}) = \mathbf{b}$

$$\underbrace{\mathbf{C}(\mathbf{Bx})}_{=: \mathbf{y}} = \mathbf{b}$$

first solve  $\mathbf{Cy} = \mathbf{b} \rightarrow \mathbf{y}$ , then solve  $\mathbf{Bx} = \mathbf{y} \rightarrow \mathbf{x}$

- One “difficult” solve is traded against two “easy solves”
- Such factorizations are typically generated in a step-wise fashion using, for instance Gauß elimination, in which case  $\mathbf{C}$  is a product of “simple” lower triangular matrices, or rotations, in which case  $\mathbf{C}$  is a unitary matrix successively built from products of “simple” unitary matrices.
- The factors belong to special matrix classes to be discussed below

# Easy to invert matrices ...

- $\mathbf{I} = (\delta_{i,j})_{i,j=1}^{n,n}$  trivial  $\mathbf{I}^{-1} = \mathbf{I}$

- Upper/lower triangular matrices:

$$\mathbf{R} = (r_{i,j})_{i,j=1}^{n,n}, \quad j < i \Rightarrow r_{i,j} = 0, \quad \mathbf{L} = (\ell_{i,j})_{i,j=1}^{n,n}, \quad j > i \Rightarrow \ell_{i,j} = 0,$$

$\mathbf{R}$  non-singular if and only if  $\det(\mathbf{R}) = \prod_{j=1}^n r_{j,j} \neq 0$ .

Backsubstitution:  $\mathbf{R}\mathbf{x} = \mathbf{b} \Rightarrow x_n = b_n/r_{n,n}$ , knowing  $x_n, x_{n-1}, \dots, x_{n-j} \Rightarrow$

$$x_{n-j-1} = \left( b_{n-j-1} - \sum_{k=0}^j r_{n-j-1, n-k} x_{n-k} \right) / r_{n-j-1, n-j-1}$$

Complexity:  $\sim n^2$  flops

- same for  $\mathbf{L}\mathbf{x} = \mathbf{b}$
- Orthogonal ( $\mathbb{K} = \mathbb{R}$ )/unitary ( $\mathbb{K} = \mathbb{C}$ ) matrices:  $\mathbf{Q} \in \mathcal{O}_n$  iff  $\mathbf{Q}^* \mathbf{Q} = \mathbf{I}$ . i.e.

$$\mathbf{Q}\mathbf{x} = \mathbf{b} \Leftrightarrow \mathbf{x} = \mathbf{Q}^* \mathbf{b}$$

# Properties of Triangular Matrices

- An upper triangular matrix  $\mathbf{R} \in \mathbb{K}^{n \times n}$  is non-singular iff  $r_{j,j} \neq 0, j = 1, \dots, n$ , and

$$\det(\mathbf{R}) = \prod_{j=1}^n r_{j,j}.$$

- The inverse of an upper (lower) triangular non-singular matrix is upper (lower) triangular.
- Products of upper (lower) triangular matrices (of the same dimension) are again upper (lower) triangular. Hence, such matrices form a multiplicative group with  $\mathbf{I}$  as the neutral element.

## Exercise 34

*Prove the above statements*

# LR Factorization

## Theorem 35

For every nonsingular  $\mathbf{A} \in \mathbb{R}^{n \times n}$  there exists a permutation matrix  $\mathbf{P}$  a *normalized* lower triangular matrix  $\mathbf{L}$ , i.e.,  $\ell_{j,j} = 1, j = 1, \dots, n$ , and an upper triangular matrix  $\mathbf{R}$ , such that

$$\mathbf{PA} = \mathbf{LR}$$

- The construction of this factorization is a by-product of Gauß-elimination with pivoting, i.e.,  $\mathbf{A}$  is reduced to an upper triangular matrix  $\mathbf{R}$  by columnwise elimination of all non-zero entries below the diagonal, starting with the first column. More precisely, at the  $j$ th step rows are exchanged so as to move the largest entry in column  $j$  to the diagonal position (multiplication by a permutation matrix  $\mathbf{P}_j$ , followed by an elimination step, realized by multiplication with a specific lower triangular normalized matrix  $\mathbf{L}_j$  (Frobenius matrix). Thus

$$\mathbf{R} = \mathbf{L}_{n-1} \mathbf{P}_{n-1} \mathbf{L}_{n-2} \cdots \mathbf{L}_1 \mathbf{P}_1 \mathbf{A} = \underbrace{\tilde{\mathbf{L}}_{n-1} \cdots \tilde{\mathbf{L}}_1}_{=:\mathbf{L}^{-1}} \underbrace{\mathbf{P}_{n-1} \cdots \mathbf{P}_1}_{=:\mathbf{P}} \mathbf{A}$$

- The computational cost is  $\sim n^3/3$  multiplications

## Exercise 36

What is the complexity of computing  $\mathbf{A}^{-1}$ ?

What is the complexity of computing  $\det(\mathbf{A})$ ?

# QR Factorization

## Theorem 37

For every  $\mathbf{A} \in \mathbb{K}^{m \times n}$  there exist  $\mathbf{Q} \in \mathcal{O}_m$  and an upper triangular  $\mathbf{R} \in \mathbb{K}^{m \times n}$  (with the natural interpretation when  $m \neq n$ ) such that

$$\mathbf{A} = \mathbf{QR}$$

## Remark 38

There is neither any restriction on the dimension  $m \times n$  nor on ranks.

The factorization is computed by repeated multiplication by special “elementary” unitary matrices so as to successively eliminate sub-triangular entries:

$$\mathbf{Q}_{n-1} \mathbf{Q}_{n-2} \cdots \mathbf{Q}_1 \mathbf{A} = \mathbf{R} \quad \Leftrightarrow \quad \underbrace{\mathbf{Q}_1^* \mathbf{Q}_2^* \cdots \mathbf{Q}_{n-2}^* \mathbf{Q}_{n-1}^*}_{=:\mathbf{Q}} \mathbf{R} = \mathbf{A}$$

An ingredient: dyades - rank-one matrices:  $\mathbf{x} \in \mathbb{K}^m, \mathbf{y} \in \mathbb{K}^n \rightsquigarrow$

$$\mathbf{xy}^* = (\bar{y}_1 \mathbf{x}, \dots, \bar{y}_n \mathbf{x}) \in \mathbb{K}^{m \times n}$$

Hence, for  $\mathbf{z} \in \mathbb{K}^n$  one has  $(\mathbf{xy}^*)\mathbf{z} = \mathbf{x}(\mathbf{y}^*\mathbf{z}) = \langle \mathbf{z}, \mathbf{y} \rangle_n \mathbf{x}$ , i.e.,  $\text{rank}(\mathbf{xy}^*) = 1$

# Householder-Reflections

Given  $\mathbf{a} \in \mathbb{K}^n$  let

$$\mathbf{Q}_v := \mathbf{I} - \frac{2}{\mathbf{v}^* \mathbf{v}} \mathbf{v} \mathbf{v}^*$$

The following facts can be verified by corresponding calculations:

- 1  $\mathbf{Q}_v^* = \mathbf{Q}_v$
- 2  $\mathbf{Q}_v^2 = \mathbf{I}$ , i.e.,  $\mathbf{Q}_v \in \mathcal{O}_n$
- 3  $\mathbf{Q}_v = \mathbf{Q}_{\alpha \mathbf{v}}$  for any  $\alpha \in \mathbb{K} \setminus \{0\}$
- 4  $\mathbf{Q}_v \mathbf{v} = -\mathbf{v}$ ; interpret this geometrically in terms of the hyperplane

$$H_v := \{\mathbf{x} \in \mathbb{K}^n : \mathbf{v}^* \mathbf{x} = 0\}$$

Claim:

Given  $\mathbf{a} \in \mathbb{K}^n$ , take  $\mathbf{v}(\mathbf{a}) = \mathbf{v} := \mathbf{a} + \operatorname{sgn}(a_1) \|\mathbf{a}\|_2 \mathbf{e}^1$ ,  $\mathbf{e}^j := (0, \dots, 0, 1, 0, \dots, 0)^\top$ . Then

$$\mathbf{Q}_v \mathbf{a} = -\operatorname{sgn}(a_1) (\|\mathbf{a}\|_2 \mathbf{e}^1,$$

i.e.,  $\mathbf{Q}_v$  rotates  $\mathbf{a}$  into a multiple of the first coordinate vector. The particular sign is chosen so as to avoid numerical cancellation in the first component of  $\mathbf{a} + \operatorname{sgn}(a_1) \|\mathbf{a}\|_2 \mathbf{e}^1$



# QR-Factorization

Recall:

$$\mathbf{v}(\mathbf{a}) = \mathbf{v} := \mathbf{a} + \operatorname{sgn}(a_1) \|\mathbf{a}\|_2 \mathbf{e}^1 \quad (6.4)$$

- Given  $\mathbf{A} = (\mathbf{a}^1, \mathbf{a}^2, \dots, \mathbf{a}^n)$ ,  $\mathbf{a}^j \in \mathbb{K}^m$ , choose  $\mathbf{v}^1 := \mathbf{v}(\mathbf{a}^1) \rightsquigarrow$

$$\mathbf{Q}_1 := \mathbf{Q}_{\mathbf{v}^1} \mathbf{A} = \begin{pmatrix} \tilde{a}_{1,1} & 0 \\ 0 & \tilde{\mathbf{A}}_1 \end{pmatrix} =: \mathbf{A}^1, \quad \tilde{a}_{1,1} = a_{1,1} + \operatorname{sgn}(a_{1,1}) \|\mathbf{a}^1\|_2, \quad \tilde{\mathbf{A}}_1 \in \mathbb{K}^{(m-1) \times (n-1)}$$

- Given  $\mathbf{A}_j$ , let  $\tilde{\mathbf{Q}}_{j+1} = \mathbf{Q}_{\mathbf{v}((\tilde{\mathbf{A}}_j)^1)} \in \mathcal{O}_{n-j}$  where  $(\tilde{\mathbf{A}}_j)^1$  is the first column in  $\tilde{\mathbf{A}}_j \in \mathbb{K}^{(m-j) \times (n-j)}$  and set

$$\mathbf{Q}_{j+1} := \begin{pmatrix} \mathbf{I}_j & 0 \\ 0 & \tilde{\mathbf{Q}}_{j+1} \end{pmatrix} \in \mathcal{O}_m \rightsquigarrow \mathbf{Q}_{j+1} \mathbf{A}_j = \mathbf{A}_{j+1} \quad (\mathbf{I}_j \text{ is the } j \times j \text{ identity matrix})$$

Then  $\mathbf{A}_n =: \mathbf{R}$  is upper triangular and

$$\mathbf{Q}_{n-1} \cdots \mathbf{Q}_1 \mathbf{A} = \mathbf{R}, \quad \rightsquigarrow \quad \mathbf{A} = \mathbf{Q}_1^* \cdots \mathbf{Q}_{n-1}^* \mathbf{R} = \mathbf{QR}$$

# Orthonormalization in $\mathbb{K}^n$

Suppose that  $\{\mathbf{b}^1, \dots, \mathbf{b}^n\}$  is a basis for  $\mathbb{K}^n$ . One could generate an orthonormal basis for  $\mathbb{K}^n$  by applying Gram-Schmidt orthogonalization. Alternatively, let  $\mathbf{B}$  be the matrix with columns  $\mathbf{b}^j$  and compute a  $QR$  factorization

$$\mathbf{B} = \mathbf{Q}\mathbf{R} \quad \Leftrightarrow \quad \mathbf{Q} = \mathbf{B}\mathbf{R}^{-1},$$

i.e., the columns  $\mathbf{q}^j$  of  $\mathbf{Q}$  form an orthonormal basis of  $\mathbb{K}^n$  as well. Recall that  $\kappa_2(\mathbf{B}) = \kappa_2(\mathbf{R})$  (since  $\mathbf{Q}$  is unitary). But  $\mathbf{Q}$  can be computed, **without inverting  $\mathbf{R}$** , by a successive application of Householder reflections, each being unitary and numerically stable. This process being stable the interrelation between the basis  $\{\mathbf{q}^1, \dots, \mathbf{q}^n\}$  and the basis  $\{\mathbf{b}^1, \dots, \mathbf{b}^n\}$  is given by the transformation  $\mathbf{R}^{-1}$  whose condition number equals the one of  $\mathbf{B}$ . Thus if  $\{\mathbf{b}^1, \dots, \mathbf{b}^n\}$  has poor stability properties  $\mathbf{R}^{-1}$  is poorly conditioned. This is, however, not used in the computation of  $\mathbf{Q}$ !

Compare this with the Gram-Schmidt process!

# Definition and Existence

The following factorization has a multitude of applications in [data science and machine learning](#). Unlike the spectral decompositions or *LR*-factorization it holds for matrices of [any](#) dimension (as the *QR*-factorization)

## Theorem 39

For every matrix  $\mathbf{A} \in \mathbb{K}^{m \times n}$  there exist unitary matrices  $\mathbf{U} \in \mathcal{O}_m$ ,  $\mathbf{V} \in \mathcal{O}_n$  and a diagonal matrix  $\mathbf{S} \in \mathbb{R}^{m \times n}$

$$\mathbf{S} = \text{diag}(\sigma_1, \dots, \sigma_p, \mathbf{0}, \dots, \mathbf{0}), \quad p := \min\{m, n\},$$

with

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p \geq 0, \quad (6.5)$$

so that

$$\mathbf{U}^* \mathbf{A} \mathbf{V} = \mathbf{S} \quad \Leftrightarrow \quad \mathbf{A} = \mathbf{U} \mathbf{S} \mathbf{V}^* \quad (6.6)$$

The second relation in (6.6) can be written as

$$\mathbf{A} = \sum_{k=1}^p \sigma_k \mathbf{u}^k (\mathbf{v}^k)^* \quad (6.7)$$

where  $\mathbf{u}^k, \mathbf{v}^k$  are the columns of  $\mathbf{U}, \mathbf{V}$ , respectively.

**Proof of Theorem 39:** Assume  $\mathbf{A} \neq \mathbf{0}$ , let

$$\sigma_1 := \|\mathbf{A}\|_2 = \max_{\|\mathbf{x}\|_2=1} \|\mathbf{A}\mathbf{x}\|_2 > 0.$$

Let  $\mathbf{v} \in \mathbb{K}^n$ , with  $\|\mathbf{v}\|_2 = 1$  a vector satisfying  $\|\mathbf{A}\|_2 = \|\mathbf{A}\mathbf{v}\|_2$  and  $\mathbf{u} := \frac{1}{\sigma_1}\mathbf{A}\mathbf{v} \in \mathbb{K}^m$ . Then one has for  $\mathbf{u}$  that  $\|\mathbf{u}\|_2 = \|\mathbf{A}\mathbf{v}\|_2/\sigma_1 = 1$ . We can extend the vectors  $\mathbf{v}$  and  $\mathbf{u}$  to orthonormal bases  $\{\mathbf{v}, \tilde{\mathbf{v}}_2, \dots, \tilde{\mathbf{v}}_n\}$  resp.  $\{\mathbf{u}, \tilde{\mathbf{u}}_2, \dots, \tilde{\mathbf{u}}_m\}$  of  $\mathbb{K}^n$  resp.  $\mathbb{K}^m$ . We view the elements of these bases as columns of corresponding unitary matrices  $\mathbf{V}_1 \in \mathcal{O}_n(\mathbb{K})$ ,  $\mathbf{U}_1 \in \mathcal{O}_m(\mathbb{K})$ :

$$\begin{aligned} \mathbf{V}_1 &= \begin{pmatrix} \mathbf{v} & \tilde{\mathbf{V}}_1 \end{pmatrix} \in \mathbb{K}^{n \times n}, & \text{unitary,} \\ \mathbf{U}_1 &= \begin{pmatrix} \mathbf{u} & \tilde{\mathbf{U}}_1 \end{pmatrix} \in \mathbb{K}^{m \times m}, & \text{unitary.} \end{aligned}$$

Since  $\tilde{\mathbf{u}}_i^* \mathbf{A}\mathbf{v} = \sigma_1 \tilde{\mathbf{u}}_i^* \mathbf{u} = 0$ ,  $i = 2, \dots, m$  the matrix  $\mathbf{U}_1^* \mathbf{A}\mathbf{V}_1$  has the form

$$\mathbf{A}_1 := \mathbf{U}_1^* \mathbf{A}\mathbf{V}_1 = \begin{pmatrix} \sigma_1 & \mathbf{w}^* \\ \mathbf{0} & \mathbf{B} \end{pmatrix} \in \mathbb{K}^{m \times n},$$

with  $\mathbf{w} \in \mathbb{K}^{n-1}$ . From

$$\left\| \mathbf{A}_1 \begin{pmatrix} \sigma_1 \\ \mathbf{w} \end{pmatrix} \right\|_2 = \left\| \begin{pmatrix} \sigma_1^2 + \mathbf{w}^* \mathbf{w} \\ \mathbf{B}\mathbf{w} \end{pmatrix} \right\|_2 \geq \sigma_1^2 + \mathbf{w}^* \mathbf{w} = \left\| \begin{pmatrix} \sigma_1 \\ \mathbf{w} \end{pmatrix} \right\|_2^2$$

and  $\|\mathbf{A}\|_2 = \|\mathbf{A}_1\|_2$  it follows that

$$\sigma_1 = \|\mathbf{A}_1\|_2 \geq \frac{\|\mathbf{A}_1 \begin{pmatrix} \sigma_1 \\ \mathbf{w} \end{pmatrix}\|_2}{\left\| \begin{pmatrix} \sigma_1 \\ \mathbf{w} \end{pmatrix} \right\|_2} \geq \sqrt{\sigma_1^2 + \mathbf{w}^* \mathbf{w}} \Rightarrow \mathbf{w} = \mathbf{0} \Rightarrow \mathbf{U}_1^* \mathbf{A}\mathbf{V}_1 = \begin{pmatrix} \sigma_1 & \mathbf{0} \\ \mathbf{0}^* & \mathbf{B} \end{pmatrix} \in \mathbb{K}^{m \times n}.$$

This proves the assertion for  $m = 1$  or  $n = 1$ .

For  $m, n > 1$  one can use induction. Suppose that  $\mathbf{U}_2^* \mathbf{B} \mathbf{V}_2 = \mathbf{S}_2$  with  $\mathbf{U}_2 \in \mathcal{O}_{m-1}(\mathbb{K})$ ,  $\mathbf{V}_2 \in \mathcal{O}_{n-1}(\mathbb{K})$  and  $\mathbf{S}_2 \in \mathbb{R}_+^{(m-1) \times (n-1)}$  is diagonal. For the largest diagonal entry  $\sigma_2$  of  $\mathbf{S}_2$  we have again  $\sigma_2 := \|\mathbf{B}\|_2 \leq \|\mathbf{U}_1^* \mathbf{A} \mathbf{V}_1\|_2 = \|\mathbf{A}\|_2 = \sigma_1$ . Furthermore, with the unitary matrices

$$\mathbf{U} = \mathbf{U}_1 \begin{pmatrix} 1 & \mathbf{0}^* \\ \mathbf{0} & \mathbf{U}_2 \end{pmatrix}, \quad \mathbf{V} = \mathbf{V}_1 \begin{pmatrix} 1 & \mathbf{0}^* \\ \mathbf{0} & \mathbf{V}_2 \end{pmatrix}$$

we obtain the decomposition

$$\mathbf{U}^* \mathbf{A} \mathbf{V} = \begin{pmatrix} \sigma_1 & \mathbf{0}^* \\ \mathbf{0} & \mathbf{S}_2 \end{pmatrix}$$

from which the assertion follows by induction. □

# Basic Properties

## Proposition 40

For  $\mathbf{A} \in \mathbb{K}^{m \times n}$ ,  $\mathbf{U}^* \mathbf{A} \mathbf{v} = \mathbf{S}$  as before with

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > \sigma_{r+1} = \dots = \sigma_p = 0, \quad p = \min\{m, n\},$$

one has:

- 1  $\mathbf{A} \mathbf{v}^j = \sigma_j \mathbf{u}^j, \quad \mathbf{A}^* \mathbf{u}^j = \sigma_j \mathbf{v}^j, j = 1, \dots, p.$
- 2  $\text{rank}(\mathbf{A}) = r,$
- 3  $\text{range}(\mathbf{A}) = \text{span}\{\mathbf{u}^1, \dots, \mathbf{u}^r\}, \text{ker}(\mathbf{A}) = \text{span}\{\mathbf{v}^{r+1}, \dots, \mathbf{v}^n\}.$
- 4  $\|\mathbf{A}\|_2 = \sigma_1.$
- 5 *The strictly positive singular values  $\sigma_k, k \leq r$  are the square roots of the (strictly positive) eigenvalues of  $\mathbf{A}^* \mathbf{A}$ :*

$$\sigma_j = \sqrt{\lambda_j(\mathbf{A}^* \mathbf{A})}, \quad j = 1, \dots, r.$$

(1) - (4) follow immediately from  $\mathbf{A} = \mathbf{U} \mathbf{S} \mathbf{V}^*$  (and its hermitian conjugate) and  $\|\mathbf{A}\|_2 = \|\mathbf{S}\|_2 = \sigma_1.$

Since  $\mathbf{A}^* \mathbf{A} = \mathbf{V} \mathbf{S}^* \mathbf{U}^* \mathbf{U} \mathbf{S} \mathbf{V}^* = \mathbf{V} \text{diag}(\sigma_1^2, \dots, \sigma_r^2, \mathbf{0}, \dots, \mathbf{0}) \mathbf{V}^*$ , we have determined the spectral decomposition of  $\mathbf{A}^* \mathbf{A}$  which confirms (5).

# More Properties and the Pseudo-Inverse of $\mathbf{A}$

- Likewise one could have argued  $\mathbf{AA}^* = \mathbf{USS}^*\mathbf{U}^*$  which reveals the roles of  $\mathbf{U}, \mathbf{V}$  regarding the eigenspaces of  $\mathbf{AA}^*, \mathbf{A}^*\mathbf{A}$ , respectively.
- Denoting by  $\mathbf{U}_r, \mathbf{V}_r$  ( $r = \text{rank}(\mathbf{A})$  as above) the matrices formed by the first  $r$  columns of  $\mathbf{U}, \mathbf{V}$ , respectively, which refines (6.7)

$$\mathbf{A} = \sum_{k=1}^r \sigma_k \mathbf{u}^k (\mathbf{v}^k)^*. \quad (6.8)$$

If  $r \ll p = \min\{m, n\}$ , computing first  $c_j := (\mathbf{v}^j)^* \mathbf{x}$ ,  $j = 1, \dots, r$  and then summing  $\sum_{k=1}^r \sigma_k c_k \mathbf{u}^k$ , requires roughly  $r(n+m)$  operations as opposed to the order of  $mn$  operations when applying  $\mathbf{A}$  to  $\mathbf{x}$  directly.

- For  $\mathbf{S}$  as above define  $\mathbf{S}^\dagger := \text{diag}(\sigma_1^{-1}, \dots, \sigma_r^{-1}, \mathbf{0}, \dots, \mathbf{0}) \in \mathbb{R}^{n \times m}$  (the dimension of the transpose with inverses of the positive singular values on the diagonal). One then **defines the pseudo-inverse** of  $\mathbf{A}$  as:

$$\mathbf{A}^\dagger := \mathbf{VS}^\dagger \mathbf{U}^* \in \mathbb{K}^{n \times m}. \quad (6.9)$$

One checks that

$$\mathbf{A}^\dagger = \begin{cases} \mathbf{A}^{-1} & \text{if } m = n \text{ and } \mathbf{A} \text{ is non-singular,} \\ (\mathbf{A}^* \mathbf{A})^{-1} \mathbf{A}^* & \text{if } m > n \text{ and } \text{rank}(\mathbf{A}) = n. \end{cases}$$

# An Application of the Pseudo-Inverse Least Squares

- $\mathbf{A} \in \mathbb{K}^{n \times n}$  non-singular  $\rightsquigarrow \mathbf{Ax} = \mathbf{b}$  has a unique solution  $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$   
 $\mathbf{b} \rightarrow \mathbf{x}$  is a linear process;
- $\mathbf{A} \in \mathbb{K}^{m \times n}$ ,  $m > n$ , it makes no longer sense to pose  $\mathbf{Ax} = \mathbf{b}$ . Instead: find

$$\hat{\mathbf{x}} \in \underset{\mathbf{x} \in \mathbb{K}^n}{\operatorname{argmin}} \|\mathbf{Ax} - \mathbf{b}\|_2 \quad (6.10)$$

**Comment:** one could take as well any other norm but, as shown later, the Euclidean norm has significant advantages and comes with a favorable statistical interpretation

- $\mathbf{A}$  has full rank  $n$   $\rightsquigarrow$  (6.10) has the unique solution  $\hat{\mathbf{x}} = (\mathbf{A}^* \mathbf{A})^{-1} \mathbf{A}^* \mathbf{b}$   
 again  $\mathbf{b} \rightarrow \hat{\mathbf{x}}$  is a linear process;
- $\operatorname{rank}(\mathbf{A}) < \min\{n, m\}$   $\rightsquigarrow$  (6.10) has infinitely many solutions, but  $\rightsquigarrow \hat{\mathbf{x}} = \mathbf{A}^\dagger \mathbf{b}$  is  
 the unique minimizer of **minimal Euclidean norm**

See Theorem 45 later below for more details.

## Remark 41

*Finding  $\hat{\mathbf{x}}$  in (6.10) is called a **Least Squares** problem (minimizing the squares of residual components). It plays an eminent role for the design of **estimators** in machine learning, especially regression. It can be seen as generalizing the solution of linear systems which is well-posed only if one has as many equations (conditions) as unknowns. Least squares methods will therefore be discussed in more detail later.*



# Best Low-Rank Approximation

Many applications of the SVD, used later, are based on the following fact:

## Theorem 42

Let  $\mathbf{A} = \mathbf{USV}^*$  and  $\text{rank}(\mathbf{A}) = r \leq p := \min\{m, n\}$ . Defining the truncated matrices  $\mathbf{U}_k, \mathbf{V}_k$  as before, one has for  $k \leq p$

$$\|\mathbf{A} - \mathbf{U}_k \mathbf{S}_k \mathbf{V}_k^*\|_2 = \min\{\|\mathbf{A} - \mathbf{B}\|_2 : \mathbf{B} \in \mathbb{K}^{m \times n}, \text{rank}(\mathbf{B}) \leq k\} = \sigma_{k+1}. \quad (6.11)$$

Moreover, defining the Frobenius norm  $\|\mathbf{A}\|_F := \left(\sum_{j,k=1}^{m,n} |a_{j,k}|^2\right)^{1/2} = (\text{trace}(\mathbf{A}^* \mathbf{A}))^{1/2}$ , one has

$$\|\mathbf{A} - \mathbf{U}_k \mathbf{S}_k \mathbf{V}_k^*\|_F^2 = \sum_{j>k} \sigma_j^2. \quad (6.12)$$

**Proof:** Since the spectral norm of a diagonal matrix is the maximal diagonal entry in absolute value, one has  $\|\mathbf{A} - \mathbf{U}_k \mathbf{S}_k \mathbf{V}_k^*\|_2 = \|\mathbf{U}(\mathbf{S} - \mathbf{S}_k)\mathbf{V}^*\|_2 = \|\mathbf{S} - \mathbf{S}_k\|_2 = \sigma_{k+1}$ . It remains to show that for every matrix  $\mathbf{B}$  of rank  $k$  one has  $\|\mathbf{B} - \mathbf{A}\|_2 \geq \sigma_{k+1}$ . Now suppose  $\mathbf{B} \in \mathbb{K}^{m \times n}$  has rank  $\leq k$ . Then  $\dim(\ker(\mathbf{B})) \geq n - k$ . Hence  $\mathbb{U} := \ker(\mathbf{B}) \cap \text{span}\{\mathbf{v}^1, \dots, \mathbf{v}^{k+1}\} \neq \{\mathbf{0}\}$ . Let  $\mathbf{z} = \sum_{j=1}^{k+1} \alpha_j \mathbf{v}^j \in \mathbb{U}$ ,  $\sum_{j=1}^{k+1} |\alpha_j|^2 = 1$ . Then

$$\|\mathbf{A} - \mathbf{B}\|_2^2 \geq \|(\mathbf{A} - \mathbf{B})\mathbf{z}\|_2^2 = \|\mathbf{Az}\|_2^2 \stackrel{\text{Prop. 40, (1)}}{=} \sum_{j=1}^{k+1} \sigma_j^2 |\alpha_j|^2 \geq \sigma_{k+1}^2 \sum_{j=1}^{k+1} |\alpha_j|^2 = \sigma_{k+1}^2. \quad \square$$

# A Greedy Characterization Principal Component Analysis (PCA)

Given a **point cloud**  $\{\mathbf{a}^j : j = 1, \dots, n\} \subset \mathbb{R}^m$ , find its **best simultaneous approximation by a line**  $L_1 := \{\mathbf{x} = \mathbf{t}\mathbf{u} : \mathbf{u} \in \mathbb{R}^m, \|\mathbf{u}\|_2 = 1, \mathbf{t} \in \mathbb{R}\}$ , i.e.,

$$\mathbf{u}^1 = \operatorname{argmin}_{\substack{\mathbf{u} \in \mathbb{R}^m \\ \|\mathbf{u}\|_2 = 1}} \sum_{i=1}^n \|\mathbf{a}^i - (\mathbf{u}^* \mathbf{a}^i) \mathbf{u}\|_2^2. \quad (6.13)$$

## Lemma 43

The unit vector  $\mathbf{u}^1$ , defined by (6.14) is characterized as follows

$$\mathbf{u}^1 = \operatorname{argmax}_{\substack{\mathbf{u} \in \mathbb{R}^m \\ \|\mathbf{u}\|_2 = 1}} \mathbf{u}^* \mathbf{A} \mathbf{A}^* \mathbf{u}, = \max_{\substack{\mathbf{u} \in \mathbb{R}^m \\ \|\mathbf{u}\|_2 = 1}} \mathbf{u}^* \mathbf{A} \mathbf{A}^* \mathbf{u} \lambda_{\max}(\mathbf{A} \mathbf{A}^*) = \sigma_1^2. \quad (6.14)$$

## Remark 44

Thus,  $\mathbf{u}^1$  is the normalized eigenvector for the largest eigenvalue of  $\mathbf{A} \mathbf{A}^*$ . This matrix will later be seen as a “covariance” matrix in a statistical context, and the direction  $\mathbf{u}^1$  maximizes the “variance”.

**Proof of Lemma 43:** Recall that

$$\text{trace}(\mathbf{A}^* \mathbf{A}) = \sum_{j=1}^n (\mathbf{A}^* \mathbf{A})_{j,j} = \sum_{k=1}^n \sum_{j=1}^m a_{k,j}^2 = \sum_{j=1}^n \|\mathbf{a}^j\|_2^2.$$

Therefore

$$\begin{aligned} \sum_{j=1}^n \|\mathbf{a}^j - (\mathbf{u}^* \mathbf{a}^j) \mathbf{u}\|_2^2 &= \sum_{j=1}^n \|\mathbf{a}^j\|_2^2 - 2(\mathbf{a}^j)^* (\mathbf{u}^* \mathbf{a}^j) \mathbf{u} + (\mathbf{u}^* \mathbf{a}^j)^2 = \sum_{j=1}^n \|\mathbf{a}^j\|_2^2 - (\mathbf{u}^* \mathbf{a}^j)^2 \\ &= \text{trace}(\mathbf{A}^* \mathbf{A}) - \|\mathbf{A}^* \mathbf{u}\|_2^2 = \text{trace}(\mathbf{A}^* \mathbf{A}) - \mathbf{u}^* \mathbf{A} \mathbf{A}^* \mathbf{u}. \end{aligned} \quad (6.15)$$

Maximizing  $\mathbf{u}^* \mathbf{A} \mathbf{A}^* \mathbf{u}$  over  $\|\mathbf{u}\|_2 = 1$ , minimizes the left hand side of (6.15) as claimed. As hinted at in Remark 44, this yields the first eigenvector and maximal eigenvalue of  $\mathbf{A} \mathbf{A}^*$  and hence, by Proposition 40 (4), the first column of the matrix  $\mathbf{U}$  in the SVD and the square of the first singular value  $\sigma_1$ . □

# A Greedy Characterization

## Principal Component Analysis (PCA)

- $\mathbf{u}^1$  given by (6.14);
- given  $\mathbf{u}^1, \dots, \mathbf{u}^k$ , determine  $\mathbf{u}^{k+1}$  by

$$\mathbf{u}^{k+1} = \underset{\substack{\mathbf{u} \in \mathbb{R}^m, \|\mathbf{u}\|_2=1 \\ \mathbf{u} \perp \mathbf{u}^1, \dots, \mathbf{u}^k}}{\operatorname{argmax}} \mathbf{u}^* \mathbf{A} \mathbf{A}^* \mathbf{u}, \quad (6.16)$$

$\rightsquigarrow$  this is the successive construction of the columns  $\mathbf{u}^j$  of the  $\mathbf{U}$  in SVD  $\mathbf{A} = \mathbf{U} \mathbf{S} \mathbf{V}^*$ .

**Statistical Interpretation:** If  $\mathbf{a}^j$  are random samples of an  $m$ -variate random variable the  $\mathbf{A} \mathbf{A}^*$  is the corresponding **covariance matrix** which is diagonalized by the above process PCA (i.e., by the SVD).  $\mathbf{u}^1$  is the direction that maximizes the variance - largest contribution to the total variance of the underlying distribution - providing most information.

If the distribution is normal with mean zero -  $\mathcal{N}(0, \sigma)$  - then the  $\mathbf{u}^j$  are uncorrelated and statistically independent  $\rightsquigarrow$  cluster analysis, pattern recognition, feature discrimination, ... more on this later ...

**Example:** consumer behavior; the vector  $\mathbf{a}^j$  contains characteristic consumer traits such as age, gender, income, debts, location, etc.  $\rightsquigarrow$  PCA helps identifying those possibly few characteristics with largest impact on buying patterns

# Least Squares Method

## Theorem 45

Let  $\mathbf{A} \in \mathbb{K}^{m \times n}$ ,  $\mathbf{b} \in \mathbb{K}^m$  and consider

$$\|\mathbf{Ax} - \mathbf{b}\|_2 \rightarrow_{\mathbf{x} \in \mathbb{K}^n} \min \quad (6.17)$$

Then

- 1  $\mathbf{x} \in \mathbb{K}^n$  is a solution of (6.17) if and only if

$$\mathbf{A}^* \mathbf{Ax} = \mathbf{A}^* \mathbf{b} \quad \text{called "normal equations"}. \quad (6.18)$$

- 2 (6.17) has a unique solution  $\mathbf{x}$  if and only if  $n \leq m$  and  $\text{rank}(\mathbf{A}) = n$ .
- 3 For any  $n, m \in \mathbb{N}$  and  $\text{rank}(\mathbf{A}) \leq \min\{m, n\}$  arbitrary, there exists a unique  $\tilde{\mathbf{x}} \in \mathbb{K}^n$  satisfying

$$\|\mathbf{A}\tilde{\mathbf{x}} - \mathbf{b}\|_2 = \min_{\mathbf{x} \in \mathbb{K}^n} \|\mathbf{Ax} - \mathbf{b}\|_2 \text{ and } \|\tilde{\mathbf{u}}\|_2 \leq \|\mathbf{x}'\|_2 \text{ for which } \|\mathbf{Ax}' - \mathbf{b}\|_2 = \min. \quad (6.19)$$

Moreover,  $\tilde{\mathbf{x}}$  is given by  $\tilde{\mathbf{x}} = \mathbf{A}^\dagger \mathbf{b}$ . (*minimal norm solution*)

**Proof of Theorem 45:** Let  $\mathbb{V} = \mathbb{K}^m$ ,  $\mathbb{U} := \text{ran}(\mathbf{A}) = \{\mathbf{u} = \mathbf{Ax} : \mathbf{x} \in \mathbb{K}^n\} \subset \mathbb{V}$ . **ad (1):** Then, by Theorem 24,

$$\begin{aligned} \|\mathbf{Ax} - \mathbf{b}\|_2 &= \min_{\mathbf{x}' \in \mathbb{K}^n} \|\mathbf{Ax}' - \mathbf{b}\|_2 && \Leftrightarrow && \mathbf{u}^*(\mathbf{Ax} - \mathbf{b}) = 0, \quad \forall \mathbf{u} \in \mathbb{U} = \text{ran}(\mathbf{A}) \\ &&& \Leftrightarrow && (\mathbf{Az})^*(\mathbf{Ax} - \mathbf{b}) = 0, \quad \forall \mathbf{z} \in \mathbb{K}^n \\ &&& \Leftrightarrow && \mathbf{z}^*(\mathbf{A}^*\mathbf{Ax} - \mathbf{A}^*\mathbf{b}) = 0, \quad \forall \mathbf{z} \in \mathbb{K}^n \\ &&& \stackrel{\text{Remark 25,(1)}}{\Leftrightarrow} && \mathbf{A}^*\mathbf{Ax} - \mathbf{A}^*\mathbf{b} = 0, \quad \rightsquigarrow (6.18). \end{aligned}$$

**ad (2):** If  $m \geq n = \text{rank}(\mathbf{A})$  then  $\mathbf{A}^*\mathbf{A}$  is hermitian positive definite and hence non-singular (see Remark 12), which by (1) confirms (2).

By (2), it suffices to consider the case  $\text{rank}(\mathbf{A}) < \min\{m, n\}$ . In this case the solution set

$$S(\mathbf{A}, \mathbf{b}) := \{\mathbf{x} \in \mathbb{K}^n : \|\mathbf{Ax} - \mathbf{b}\|_2 = \min\}$$

is easily seen to be given by

$$S(\mathbf{A}, \mathbf{b}) = \mathbf{x}^0 + \ker(\mathbf{A}) = \{\mathbf{x}^0 + \mathbf{z} : \mathbf{z} \in \mathbb{K}^n, \mathbf{Az} = 0\}, \quad (6.20)$$

where  $\mathbf{x}^0$  is any fixed minimizer of  $\|\mathbf{Ax} - \mathbf{b}\|_2$ . Taking  $\mathbb{U} := \ker(\mathbf{A}) \subset \mathbb{V} := \mathbb{K}^n$ , we know again from Theorem 24 that there exists a unique  $\bar{\mathbf{z}} \in \ker(\mathbf{A}) = \mathbb{U}$  such that

$\|\mathbf{x}^0 - \bar{\mathbf{z}}\|_2 = \min_{\mathbf{z} \in \ker(\mathbf{A})} \|\mathbf{x}^0 - \mathbf{z}\|_2$  and  $\bar{\mathbf{z}}$  is characterized by  $\mathbf{z}^*(\mathbf{x}^0 - \bar{\mathbf{z}}) = 0$ ,  $\mathbf{z} \in \ker(\mathbf{A})$ . Hence  $\tilde{\mathbf{x}} := \mathbf{x}^0 - \bar{\mathbf{z}} \in S(\mathbf{A}, \mathbf{b})$  is the unique minimal norm minimizer.

**Proof of Theorem 45 continued:** It remains to show that  $\tilde{\mathbf{x}}$  is given by  $\mathbf{A}^\dagger \mathbf{b}$ .

To that end, let  $\mathbf{y} := \mathbf{A}^\dagger \mathbf{b}$  and recall that  $\mathbf{A} = \mathbf{USV}^*$ ,  $\mathbf{A}^\dagger = \mathbf{VS}^\dagger \mathbf{U}^*$ .

Show that  $\mathbf{y} \in S(\mathbf{A}, \mathbf{b})$  by verifying that  $\mathbf{A}^* \mathbf{A} \mathbf{y} = \mathbf{A}^* \mathbf{b}$  (see (1)):

$$\begin{aligned} \mathbf{A}^* \mathbf{A} \mathbf{y} &= (\mathbf{USV}^*)^* (\mathbf{USV}^*) \mathbf{A}^\dagger \mathbf{b} = \mathbf{VS}^* \mathbf{U}^* \mathbf{USV}^* \mathbf{VS}^\dagger \mathbf{U}^* \mathbf{b} \\ &= \mathbf{VS}^* \mathbf{SS}^\dagger \mathbf{U}^* \mathbf{b} = \mathbf{VS}^* \mathbf{U}^* \mathbf{b} = \mathbf{A}^* \mathbf{b}, \end{aligned}$$

i.e.,  $\mathbf{y}$  is a solution of the normal equations and hence a minimizer.

Show that  $\mathbf{y}$  has minimal Euclidean norm:

As shown before, this is equivalent to showing that  $\mathbf{y} \perp \ker(\mathbf{A})$ . By Proposition 40, it suffices to show that  $(\mathbf{v}^k)^* \mathbf{y} = 0$ ,  $k = r + 1, \dots, n$ , where  $r = \text{rank}(\mathbf{A})$ . To that end,

$$(\mathbf{v}^k)^* \mathbf{y} = (\mathbf{v}^k)^* \mathbf{A}^\dagger \mathbf{b} = (\mathbf{v}^k)^* \mathbf{VS}^\dagger \mathbf{U}^* \mathbf{b} = (\mathbf{v}^k)^* \mathbf{V}_r \mathbf{S}_r^\dagger \mathbf{U}_r^* \mathbf{b} = \mathbf{0}$$

since  $k > r$  and the columns of  $\mathbf{V}$  are pairwise orthogonal. Hence  $\mathbf{y}$  solves Problem (6.19) and, by uniqueness, agrees with  $\tilde{\mathbf{x}}$ .  $\square$

# References I

- [1] G. H. Golub, C. F. van Loan.  
*Matrix Computations*.  
3. Aufl., Oxford University Press, 1996.
- [2] Y. Saad.  
*Iterative Methods for Sparse Linear Systems*.  
2. Aufl. SIAM Publications, 2003.
- [3] R. DeVore and G.G. Lorentz,  
*Constructive Approximation*, vol. 303,  
Springer Grundlehren, Springer, Berlin-Heidelberg, 1993.