

# Foundation of Data Analysis

Note: The dates on the lecture notes correspond with the dates on the videos, which were recorded one year prior to SS21.

---

---

---

---



## Chapter 1 : Preliminaries on Linear Algebra

### 1.1 Matrix Notations

#### Assumptions / Notations / Definitions

(i)  $I, J, L$  are finite index sets.

(ii) We mainly consider fields  $\mathbb{K} \in \{\mathbb{R}, \mathbb{C}\}$ ,

for  $\alpha \in \mathbb{C}$ ,  $\alpha = a + i \cdot b$ ,  $a, b \in \mathbb{R}$ , we consider the complex conjugate

$$\bar{\alpha} := a - i \cdot b \in \mathbb{C}.$$

(iii) Define  $\mathbb{K}^{I \times J}$  the vector space of matrices  $A \in \mathbb{K}^{I \times J}$ .

The entries of  $A$  are denoted by  $A_{ij}$ ,  $i \in I$  row,  $j \in J$  column.

For  $a_{ij} \in \mathbb{K}$  we can form  $A := (a_{ij})_{i \in I, j \in J} \in \mathbb{K}^{I \times J}$ .

(iv) Transposed matrix  $A^T \in \mathbb{K}^{J \times I}$  for  $A \in \mathbb{K}^{I \times J}$ .

(v)  $A \in \mathbb{K}^{I \times I}$  is symmetric if  $A^T = A$ .

(vi)  $A \in \mathbb{K}^{I \times I}$  is hermitian if  $A^H = A$  for  $A^H = \overline{A^T}$ .

For  $\mathbb{K} = \mathbb{R}$  : symmetric  $\Leftrightarrow$  hermitian.

(vii) For  $A \in \mathbb{K}^{I \times J}$  denote the  $i$ th row of  $A$  by  $A^{(i)} = (a_{ij})_{j \in J}$  and the  $j$ th column of  $A$  by  $A_{(j)} = (a_{ij})_{i \in I}$ .

(viii) Matrix-vector-multiplication :  $A \in \mathbb{K}^{I \times J}$ ,  $x \in \mathbb{K}^J$ , then

$$(Ax)_i = \sum_{j \in J} a_{ij} x_j.$$

Bsp:

$$\begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} a_{11} \cdot x_1 + a_{12} \cdot x_2 \\ a_{21} \cdot x_1 + a_{22} \cdot x_2 \end{pmatrix}$$

(ix) Matrix-Matrix-multiplication :  $A \in \mathbb{K}^{I \times J}$ ,  $B \in \mathbb{K}^{J \times L}$ , then

$$(AB)_{il} = \sum_{j \in J} A_{ij} \cdot B_{jl}$$

Bsp:

$$\begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} \cdot \begin{pmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{pmatrix} = \begin{pmatrix} a_{11}b_{11} + a_{12}b_{21} & * \\ * & * \end{pmatrix}$$

(x) Define the Kronecker Symbol:  $\delta: \mathbb{I} \times \mathbb{I} \rightarrow \{0, 1\}$

$$\delta_{ij} := \begin{cases} 1, & \text{if } i = j \\ 0, & \text{if } i \neq j \end{cases}$$

(xi) Denote the unit vector  $e^{(i)} \in \mathbb{K}^{\mathbb{I}}$  by  $e^{(i)} := (\delta_{ij})_{j \in \mathbb{I}}$ ,  
the identity matrix by  $I := (\delta_{ij})_{i \in \mathbb{I}, j \in \mathbb{I}}$ .

(xii) Denote the range of  $A \in \mathbb{K}^{\mathbb{I} \times \mathbb{J}}$  by

$$\text{range}(A) := \{Ax : x \in \mathbb{K}^{\mathbb{J}}\} = \text{span}\{A_{(j)} : j \in \mathbb{J}\}.$$

$$\begin{aligned} \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} &= \begin{pmatrix} x_1 \cdot a_{11} + x_2 \cdot a_{12} \\ x_1 \cdot a_{21} + x_2 \cdot a_{22} \end{pmatrix} \\ &= x_1 \cdot \begin{pmatrix} a_{11} \\ a_{21} \end{pmatrix} + x_2 \begin{pmatrix} a_{12} \\ a_{22} \end{pmatrix} \\ &= x_1 \cdot A_{(1)} + x_2 \cdot A_{(2)} \end{aligned}$$

(xiii) For  $x, y \in \mathbb{K}^{\mathbb{I}}$ , we define the Euclidean scalar product by

$$\langle x, y \rangle = y^T x = \sum_{i \in \mathbb{I}} x_i \overline{y_i}.$$

We can rewrite for  $A \in \mathbb{K}^{\mathbb{I} \times \mathbb{J}}$ ,  $x \in \mathbb{K}^{\mathbb{J}}$

$$(Ax)_i = \langle A^{(i)}, \bar{x} \rangle \quad \text{and} \quad (AB)_{il} = \langle A^{(l)}, \bar{B}_{(i)} \rangle.$$

"undo the conjugate of definition"

(xiv)  $x, y \in \mathbb{K}^{\mathbb{I}}$  are orthogonal, if  $\langle x, y \rangle = 0$ , short  $x \perp y$ .

The sets  $X, Y \subset \mathbb{K}^{\mathbb{I}}$  are mutually orthogonal, short  $X \perp Y$ , if  
 $\forall x \in X$  and  $y \in Y$ :  $\langle x, y \rangle = 0$ , i.e.  $x \perp y$ .

For two subspaces  $X, Y \subset \mathbb{K}^{\mathbb{I}}$ , it is sufficient to check  
if their bases are mutually orthogonal.

$\left( \begin{array}{l} \text{scalar product is linear, i.e. } \langle ax, y \rangle = a \langle x, y \rangle \\ \quad \langle x+y, z \rangle = \langle x, z \rangle + \langle y, z \rangle \\ \text{and each element } x \in X, y \in Y \text{ is linear combination} \\ \text{of the basis.} \end{array} \right)$

(xv) A family of vectors  $X = \{x_v\}_{v \in \mathbb{F}} \subset \mathbb{K}^{\mathbb{I}}$  is called orthogonal  
if the vectors  $x_v$  are pair-wise orthogonal, i.e.

$$\langle x_v, x_{v'} \rangle = 0 \quad \forall v, v' \in \mathbb{F}, v \neq v'.$$

The family is called orthonormal, if it is orthogonal and

$$\langle x_v, x_v \rangle = 1 \quad \forall v \in \mathbb{F}.$$



## 1.2. Matrix Rank

$$A = \begin{pmatrix} * & * & * \\ * & * & * \\ * & * & * \end{pmatrix} \quad \begin{matrix} \uparrow \text{rows } I \\ \downarrow \text{Columns } J \end{matrix}$$

Proposition 1.1 let  $A \in \mathbb{K}^{I \times J}$ . The following statements are equivalent  
and define rank  $r = \text{rank}(A) \in \mathbb{N}$ .

- (i)  $r = \dim(\text{range}(A))$
- (ii)  $r = \dim(\text{range}(A^H))$
- (iii)  $r$  is the max. number of linear independent rows of  $A$ .
- (iv)
- (v)  $r$  is minimal with the property

$$A = \sum_{i=1}^r \underbrace{a_i b_i^H}_{\text{matrix}}, \text{ where } a_i \in \mathbb{K}^I, b_i \in \mathbb{K}^J.$$

$$ab^H = \underbrace{\begin{pmatrix} a_1 \\ \vdots \\ a_r \end{pmatrix}}_{\mathbb{K}^{I \times r}} \cdot \underbrace{\begin{pmatrix} b_1 & \cdots & b_r \end{pmatrix}}_{r \times J} = \begin{pmatrix} \uparrow \leftarrow \rightarrow \\ \downarrow \leftarrow \rightarrow \end{pmatrix} = A \in \mathbb{K}^{I \times J}$$

$$\Rightarrow A = \underbrace{\begin{pmatrix} \uparrow \leftarrow \rightarrow \\ \downarrow \leftarrow \rightarrow \end{pmatrix}}_{r \text{ times (minimum)}} + \dots + \underbrace{\begin{pmatrix} \uparrow \leftarrow \rightarrow \\ \downarrow \leftarrow \rightarrow \end{pmatrix}}_{r \text{ times (minimum)}} \quad \begin{matrix} \text{is sum of tensor} \\ \text{products} \\ a \otimes b := ab^H. \end{matrix}$$

- (vi)  $r$  is maximal with the property that there exists an invertible  $r \times r$  submatrix of  $A$ .
- (vii)  $r$  is the number of positive singular values (proven later).

Remark: •  $r_{\max} = \min \{ \# I, \# J \}$

- rank of real valued matrix does not change when considered in  $\mathbb{C}$ .
- Denote the set of matrices of bound rank  $r \leq k$  by  $\mathcal{R}_k = \{ A \in \mathbb{K}^{I \times J}, \text{rank}(A) \leq k \}$ .

$\mathcal{R}_k$  is not a vector space.

$$\text{e.g. } \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} + \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

# 1.3. Matrix Norms

## 1.3.1. Norms

Let  $\mathbb{K}$  be a field,  $V$  a vector space over  $\mathbb{K}$ .

We call a function  $\|\cdot\|: V \rightarrow [0, \infty)$  a norm, if

- (i)  $\|v\| = 0 \iff v = 0, v \in V.$
- (ii)  $\|\lambda v\| = |\lambda| \|v\|, v \in V, \lambda \in \mathbb{K}.$

(iii)  $\|v+w\| \leq \|v\| + \|w\|, v, w \in V$ , called triangle inequality.

Rem: Norms are always continuous as a consequence of the (inverse) triangle inequality.

$$|\|v\| - \|w\|| \leq \|v-w\|, \text{ for all } v, w \in V.$$

We call the pair  $(V, \|\cdot\|)$  normed vector space.

## 1.3.2. Scalar Products, Hilbert spaces, projections

A normed vector space  $(V, \|\cdot\|)$  is a pre-Hilbert space if its norm is defined by

$$\|v\| = \sqrt{\langle v, v \rangle}, \quad v \in V,$$

where  $\langle \cdot, \cdot \rangle: V \times V \rightarrow \mathbb{K}$  is a scalar product on  $V$ , i.e. it fulfills the following properties:

"How does the vector  $v$  look like, seen from  $y$ ?"

(i) Positive definiteness

$$\langle v, v \rangle > 0 \text{ for } v \neq 0$$

(ii) Conjugate symmetry

$$\langle v, w \rangle = \overline{\langle w, v \rangle} \text{ for } v, w \in V.$$

(iii) Linearity in first argument

$$\langle u + \lambda v, w \rangle = \langle u, w \rangle + \lambda \langle v, w \rangle \text{ for } u, v, w \in V, \lambda \in \mathbb{K}.$$

$$(ii) + (iii) \text{ gives } \langle w, u + \lambda v \rangle = \langle w, u \rangle + \bar{\lambda} \langle w, v \rangle \text{ for } u, v, w \in V, \lambda \in \mathbb{K}.$$

We refer to this pre-Hilbert space by  $(V, \langle \cdot, \cdot \rangle)$ .

Remark: In a pre-Hilbert space  $(V, \|\cdot\|)$  over  $\mathbb{K}$ , the triangle inequality follows from the Schwarz inequality

$$|\langle v, w \rangle| \leq \|v\| \|w\|, \quad v, w \in V.$$

Note:  
An inner product always gives a norm, by taking the square root. But a norm is only induced by an inner product if the parallelogram law holds:

$$\|x+y\|^2 + \|x-y\|^2 = 2(\|x\|^2 + \|y\|^2)$$

Example: Euclidean norm on  $\mathbb{K}^I$  (generated by  $\langle v, w \rangle = v w^\top$ ):

$$\|v\|_2 = \sqrt{\sum_{i \in I} (v_i)^2}$$

We call a pre-Hilbert space  $(V, \langle \cdot, \cdot \rangle)$  which is complete, i.e. each Cauchy-sequence is convergent wrt. its induced norm  $\|\cdot\|$ , a Hilbert space.

Remark: Each finite dimensional pre-Hilbert space is a Hilbert space.

A set  $C \subseteq V$ ,  $V$  vector space, is called convex, if

$$\forall v, w \in C, t \in [0,1] : tv + (1-t)w \in C.$$

Given a closed, convex set  $C$  in a Hilbert space  $(V, \langle \cdot, \cdot \rangle)$ , we define a projection on  $C$  of any vector  $v \in V$  as

$$P_C(v) = \underset{w \in C}{\operatorname{argmin}} \|v - w\|.$$

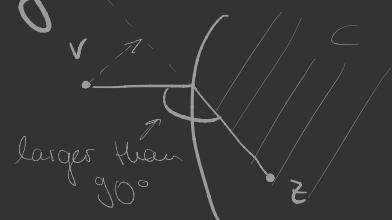
vector of minimal distance in  $C$



Such projections fulfill the following inequality:

$$\operatorname{Re}(\langle z - P_C(v), v - P_C(v) \rangle) \leq 0.$$

for all  $z \in C$ .



Note: Convexity of  $C$  ensures uniqueness of projection.

This projection is also called orthogonal projection.

Calculation of  $P_C(v)$ : In case  $C = W \subset V$  is a closed linear subspace of  $V$  and let  $\{w_v\}_{v \in F}$  be an orthonormal basis of  $W$ , then

$$P_W(v) = \sum_{v \in F} \langle v, w_v \rangle w_v \quad \text{for all } v \in V.$$

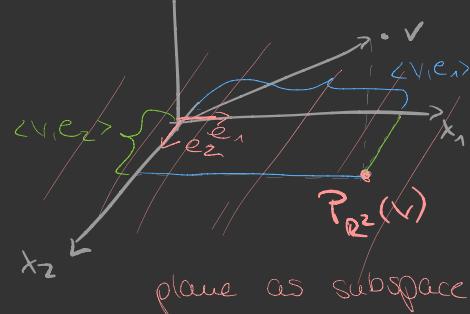
$\langle v, e_i \rangle =$  how long does  $v$  look like from  $x_i$  axis?

Theorem: Pythagoras-Fourier-Theorem



$$\|P_W(v)\|^2 = \sum_{v \in F} |\langle v, w_v \rangle|^2$$

for all  $v \in V$ .



Remark: Applied to the special case  $V = W$ , we have that  $P_W = \text{id}_W$  and hence the orthonormal expansion

$$v = \sum_{\varphi \in S} \langle v, w_\varphi \rangle w_\varphi$$

and the norm equivalence

$$\|v\|^2 = \sum_{\varphi \in S} |\langle v, w_\varphi \rangle|^2.$$

### 1.3.3. Trace and matrix norms

For  $A \in \mathbb{K}^{I \times I}$ , we define the trace of A by

$$\text{tr}(A) = \sum_{i \in I} A_{ii}.$$

Proposition:

a) Let  $A \in \mathbb{K}^{I \times J}$  and  $B \in \mathbb{K}^{J \times I}$ . Then  $\text{tr}(AB) = \text{tr}(BA)$ .

b) Circularity property

For  $A, B, C$ :  $\text{tr}(ABC) = \text{tr}(BCA)$ .

c) Invariance under unitary transformations, i.e.

$$\text{tr}(A) = \text{tr}(UAU^*)$$

for  $A \in \mathbb{K}^{I \times I}$  and unitary matrix  $U \in \mathbb{K}^{I \times I}$ .

d)  $\text{tr}(A) = \sum_{i \in I} \lambda_i$  where  $\{\lambda_i \mid i \in I\}$  set of eigenvalues of A.

Definition: Frobenius-Norm (aka Schur norm, Hilbert-Schmidt norm)

For  $A \in \mathbb{K}^{I \times J}$  we define the Frobenius norm  $\|\cdot\|_F$  as

$$\|A\|_F := \sqrt{\sum_{i \in I, j \in J} |A_{ij}|^2}.$$

It is generated by scalar product  $\langle \cdot, \cdot \rangle_F$  defined by

$$\begin{aligned} \langle A, B \rangle_F &:= \sum_{i \in I} \sum_{j \in J} \overline{A_{ij}} B_{ij} \\ &= \text{tr}(AB^*) = \text{tr}(B^*A). \end{aligned}$$

$$AB^* = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} \begin{pmatrix} \bar{b}_{11} & \bar{b}_{21} \\ \bar{b}_{12} & \bar{b}_{22} \end{pmatrix} = \begin{pmatrix} a_{11} \cdot \bar{b}_{11} + a_{12} \cdot \bar{b}_{21} & * \\ * & a_{21} \cdot \bar{b}_{12} + a_{22} \cdot \bar{b}_{22} \end{pmatrix}.$$

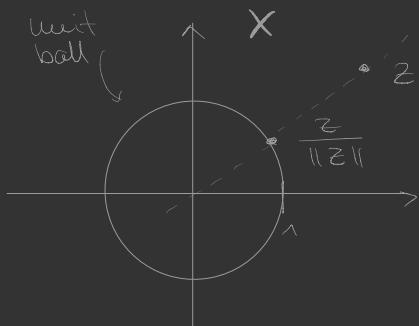
In particular  $\|A\|_F^2 = \text{tr}(AA^*) = \text{tr}(A^*A) = \langle A, A \rangle_F$ .

### Definition:

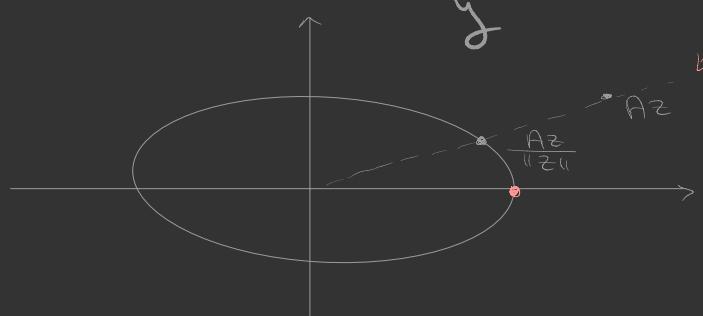
Let  $X = \mathbb{K}^I$ ,  $Y = \mathbb{K}^J$  and  $(X, \|\cdot\|_X)$ ,  $(Y, \|\cdot\|_Y)$  normed vector spaces. Then the associated matrix norm is

$$\|A\| := \|A\|_{X \rightarrow Y} := \sup_{z \neq 0} \frac{\|Az\|_Y}{\|z\|_X}, \quad A \in \mathbb{K}^{I \times J}$$

↑ search for the "worst possible"  
transformation.



$A \cdot z$



Example: For  $\|\cdot\|_X = \|\cdot\|_Y = \|\cdot\|_2$ , the associated matrix norm  $\|\cdot\|_\infty := \|\cdot\|_{X \rightarrow Y}$  is called spectral norm, often denoted by  $\|A\|$ .

Definition: A matrix norm  $\|\cdot\|$  is said to be

(i) unitary invariant, if

$$\|A\| = \|UAV^*\|$$

for unitary matrices  $U, V$ .

(ii) submultiplicative, if

$$\|AB\| \leq \|A\| \cdot \|B\|.$$

Proof uses the property of  
the trace.

! unitary matrices  
represent rotations

→ no change of the  
magnitude.

Remark: Both  $\|\cdot\|_F$  and  $\|\cdot\|_\infty$  are unitary invariant and submultiplicative.

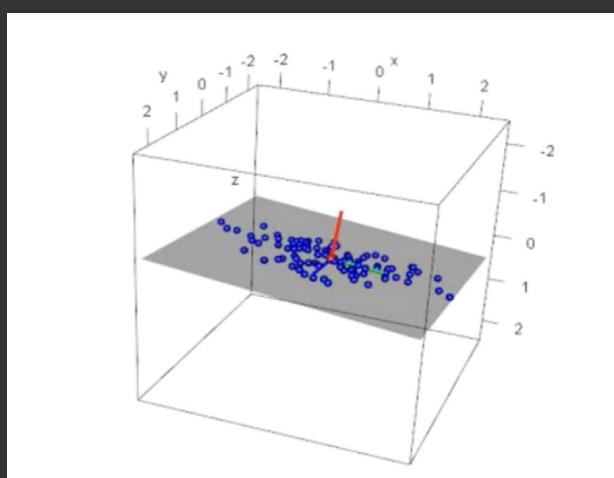
$$\text{Also, } \|AB\|_F \leq \|A\|_\infty \cdot \|B\|_F \leq \|A\|_F \cdot \|B\|_F.$$

# Chapter 2 : Singular Value Decomposition

- Example :
- We can understand data as a stream of numbers,  
e.g. representation of color in an image in a matrix/table  
or spatial movement represented by geological coordinates.
  - If we want to extract some information by simultaneously looking at all the data together
    - ~> want to find correlation between the independent streams (images, paths)...
  - For this, combine all the data in a single matrix and study its algebraic properties.  
These properties encode fundamental properties of the data.
  - What we don't do: multi-dimensional tensors.  
~> often transformed into "flat" matrix due to theoretical methods

Idea :  $A \in \mathbb{R}^{I \times J}$  with  $n = |I|$  and  $d = |J|$ , where each row represents a data point in  $\mathbb{R}^d$ .  
Want to now find the best  $k$ -dim. Subspace to represent the data points

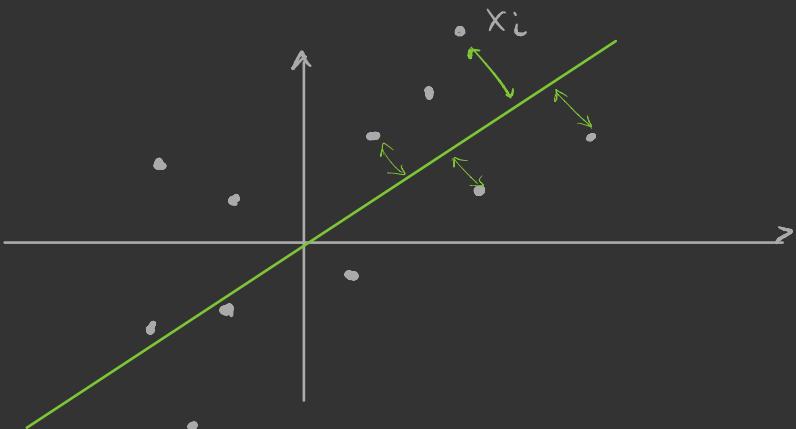
Aim : dimension reduction.



~> cut nicely through the "data cloud"

Do this by minimizing sum of Squares between subspace and data points.

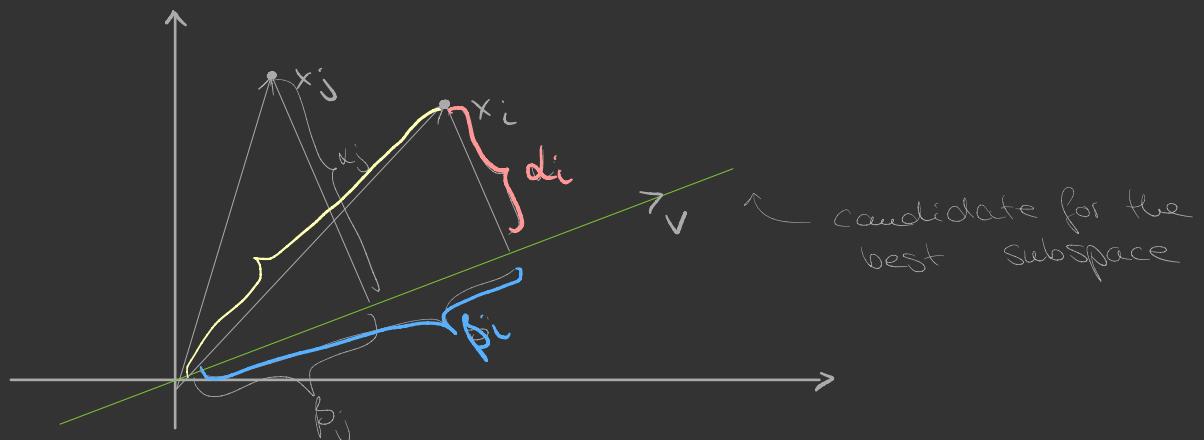
=> find the direction of maximal spread of the points in the data cloud.



① First step: Try to find a 1-dim subspace, i.e. line through the origin.

→ later: finding best fitting k-dim subspace is found by applying ① k-times.

Therefore:  $x_i = A^{(i)}$ , i.e. the  $i$ -th row of matrix A.



Pythagoras Theorem gives:

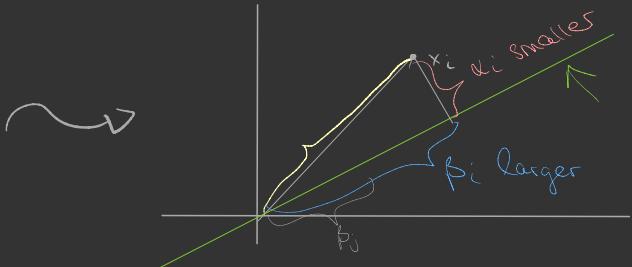
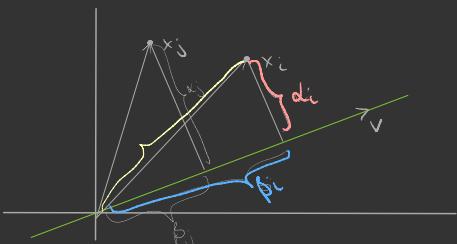
$$a^2 + b^2 = c^2$$

$$\|x\|^2 = (\text{length of projection})^2 + (\text{distance of point to line})^2$$

Want to find  $v$  to

minimize the distance of points to the line

$\Leftrightarrow$  maximize the length of the projection for each point



Let now  $v$  be the unit vector along the line.

Then the length of the projection is given by

$$|\langle A^{(i)}, v \rangle|$$

To do this for all data points simultaneously,

we want to maximize  $\|Av\|_2^2 = \sum_{i \in I} |\langle A^{(i)}, v \rangle|^2$ ,

meaning find the vector  $v \in \mathbb{R}^d$  which gets maximally distorted by matrix  $A$  (since we start with  $\|v\|_2^2 = 1$ ) or want to find the direction in which the  $d$ -dimensional sphere gets maximally deformed.

We call this vector the first singular vector  $v_1 \in \mathbb{R}^d$  or first principle component, thus

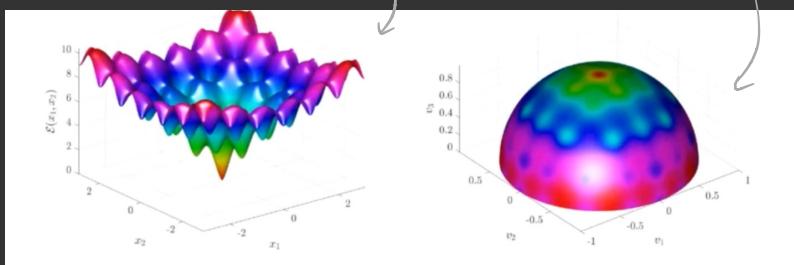
$$v_1 = \underset{\substack{v \in \mathbb{R}^d \\ \|v\|_2=1}}{\operatorname{argmax}} \|Av\|_2$$

compare  
spectral  
norm

and value  $\sigma_1(A) = \|Av_1\|_2$  is called the first singular value of  $A$ .

## Examples

1) function we want to find the minimizes for



set of potential minimizes

2)



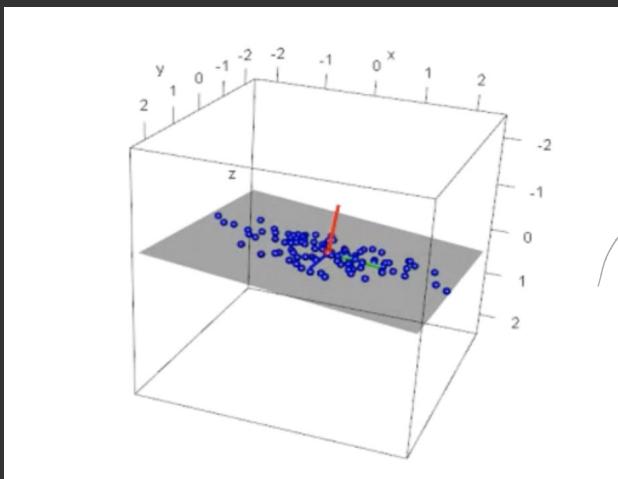
Figure 12: Samples from the 10K US Adult Faces Database [7] and one instance of outlier.

Want to compute the "most typical face"



no "Eigenface" computed with SVD.

② Go from 1-dim. subspace to 2-dim. subspace.



greedy approach: we keep  $v_1$  and then try to find a new vector  $v$  which is perpendicular to  $v_1$ , and again maximizes the distortion:

$$v_2 = \underset{\|v\|=1}{\operatorname{argmax}} \|Av\|_2 .$$

$$\langle v_1, v \rangle = 0$$

maximally "different" to  $v_1$

Then  $\Theta_2(A) = \|Av_2\|_2$  is second singular value of A

→ continue like this for  $k$ -th vector  $v_k$ , s.t.

$$v_k = \underset{\|v\|=1}{\operatorname{argmax}} \|Av\|_2 .$$

$$\langle v_1, v \rangle = 0, \dots, \langle v_{k-1}, v \rangle = 0$$

Note: Process stops when we have found  $v_1, \dots, v_r$  singular vectors and

$$0 = \underset{\|v\|=1}{\operatorname{max}} \|Av\|_2 .$$

$$\langle v_1, v \rangle = \dots = \langle v_r, v \rangle = 0$$

Reason:  $r = \text{rank}(A)$  i.e. dimension of  $\text{range}(A)$ .

Problem: Hard optimization algorithm, but can be made much easier: "Power Method" ( $\leadsto$  later)

And: Does one-by-one optimization actually give the optimal subspace?

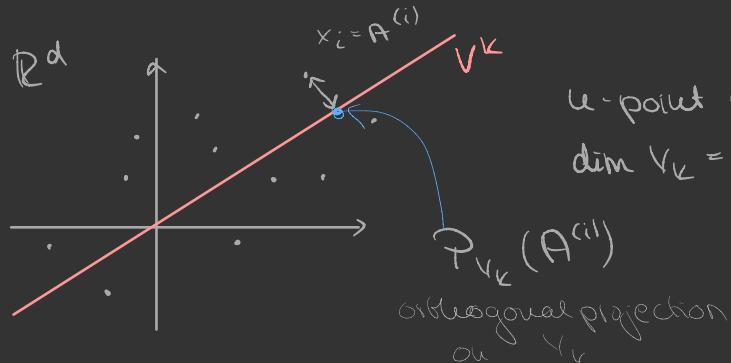
Yes, due to Pythagoras theorem and properties of the Euclidean distance.  $\Rightarrow$  Proposition 2.1.

## Worlesung 4

| 29.04.20

Proposition 2.1: Let  $A \in \mathbb{K}^{I \times J}$  and  $v_1, v_2, \dots, v_r$  be singular values defined above. For  $1 \leq k \leq r$ , let  $V_k$  be the subspace spanned by  $v_1, v_2, \dots, v_k$ . Then for each  $k$ ,  $V_k$  is the best fit subspace for  $A$ .

Proof:



$u$ -point data cloud,  
 $\dim V_k = k \leq r = \text{rank}(A)$

$$\sum_{i=1}^u \|A^{(i)} - P_V(A^{(i)})\|_2^2 = \delta(V) \quad (\text{distance of } x_i = A^{(i)} \text{ to } V)$$

$$V_k = \underset{\substack{V \\ \dim(V) \leq k}}{\operatorname{argmin}} \delta(V)$$

$$\underset{\substack{V \\ \dim(V) \leq k}}{\operatorname{argmax}} \sum_{i=1}^u \|P_V(A^{(i)})\|_2^2$$

$\hat{V}$  optimal Subspace ;  $V_k = \text{span}\{v_1, \dots, v_k\}$ .

to be proved

by induction:

$$k=1: V_1 = \text{span}\{v_1\} = \hat{V}, \text{ true by construction of } v_1.$$







We claim that the first singular vector  $z$  of  $B$  is orthogonal to  $v_1$ .

If this were not the case, then  $0 \neq z_1 = \langle z, v_1 \rangle v_1$ .

$$\left\| B \frac{z - z_1}{\|z - z_1\|_2} \right\|_2 \stackrel{\text{def.}}{=} \left\| \frac{Bz - Bz_1}{\|z - z_1\|_2} \right\|_2$$

$$Bz_1 = \underbrace{\langle z_1, v_1 \rangle}_{=0} Bv_1 = 0$$

$\Rightarrow$

$$\left\| \frac{Bz}{\|z - z_1\|_2} \right\|_2 > \|Bz\|_2$$

$$\|z\|_2 < 1$$

$z$  maximizes  
 $\|Bz\|_2$ .

$$\Rightarrow z \perp v_1.$$

Let  $v$  be orthogonal to  $v_1$ , then

$$\underline{Bv} = \underline{Av} - \theta_1 u_1 \underbrace{\langle v_1, v \rangle}_{=0} = \underline{Av}$$

Hence, the singular vectors of  $B$  are those of  $A$  orthogonal to  $v_1$ , i.e.  $\{v_2, \dots, v_r\}$

$$\Rightarrow \text{rank}(B) = r-1.$$

$$\Rightarrow \{u_2, \dots, u_r\} \text{ ONB}.$$

Remains to show:  $u_j \perp u_1$ ,  $j = 2, \dots, r$ .

Assume  $\exists j$  s.t.  $\langle u_j, u_1 \rangle \neq 0$ . For simplicity  $\mathbb{K} = \mathbb{R}$ . Wlog we assume  $\langle u_j, u_1 \rangle > 0$ .

Then for  $\epsilon > 0$

$$A \left( \frac{v_1 + \epsilon v_j}{\|v_1 + \epsilon v_j\|_2} \right) = \frac{Av_1 + \epsilon \cdot Av_j}{\underbrace{\|v_1\|_2^2 + \epsilon^2 \|v_j\|_2^2}_{=1} + 2\epsilon \underbrace{\langle v_1, v_j \rangle}_{=0}}$$

*construction of unit vector*

$$= \frac{\theta_1 u_1 + \epsilon \theta_j u_j}{\sqrt{1 + \epsilon^2}} =: w$$

Pythagoras Thm.

$$\|w\|_2 = |\langle w, u_1 \rangle| = \langle u_1, \frac{\theta_1 u_1 + \epsilon \theta_j u_j}{\sqrt{1 + \epsilon^2}} \rangle$$

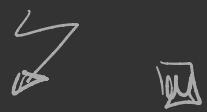
Taylor expansion

$$= (\theta_1 + \epsilon \theta_j \underbrace{\langle u_1, u_j \rangle}_{>0}) \left( 1 - \frac{\epsilon^2}{2} + \Theta(\epsilon^4) \right)$$

$\left\{ \begin{array}{l} \text{following the} \\ \text{induction iteration} \end{array} \right.$

$$= \sigma_1 + \underbrace{\varepsilon \sigma_j \langle u_1, u_j \rangle}_{>0} - O(\varepsilon^2) > \sigma_1.$$

with larger norm  
than  $\arg \max$  value

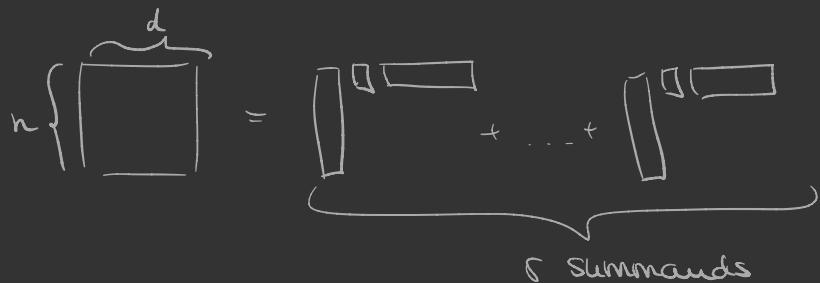


Theorem 2.3: Let  $A \in \mathbb{C}^{n \times d}$  with right singular vectors  $v_1, \dots, v_r$ , left singular vectors  $u_1, \dots, u_r$  and  $\sigma_1, \dots, \sigma_r$  the corresponding singular values. Then

$$A = \sum_{k=1}^r \sigma_k u_k v_k^H.$$

Interpretation:

$$A = \sigma_1 u_1 v_1^H + \sigma_2 u_2 v_2^H + \dots + \sigma_r u_r v_r^H = 0.$$



To spell out the matrix for  $A \in \mathbb{R}^{n \times d}$ , one needs  $\boxed{n \times d}$  scalar entries.

With this decomposition we need  $\boxed{r \cdot (n+d+1)}$  scalar entries.

$\Rightarrow (n+d+1)r \ll n \times d$  whenever  $r \ll \min\{n, d\}$ ,  
i.e. we save lots of space for sparse, low rank matrices,  
we are compressing the information of  $A$ .

$\Rightarrow$  enables us to extract relevant information from our data.

Example: image as a matrix containing pixel-wise color values.



neglect summands with  
smaller singular values

even though very  
pixelated, but  
still contain  
sufficient information

# Vorlesung 5

5.5.20

What can we use it for?

$$\text{Given } A = \sum_{i=1}^r \sigma_i u_i v_i^T \text{ for } A \in \mathbb{K}^{I \times J}.$$

- ▷ Best rank- $k$  approximation remember: Singular values are ordered by magnitude

Idea: discard information with "negligible" singular values  
 ↗ compress matrix

$$A_k = \sum_{l=1}^k \sigma_l u_l v_l^T \quad \text{"truncated SVD"}$$

Note:  $A_k$  has rank  $k$  by construction

Lemma 2.6: Rows of  $A_k$  are orthogonal projections of rows of  $A$  onto subspace  $V_k = \text{Span}\{v_1, \dots, v_k\}$

Proof: Assume  $\mathbb{K} = \mathbb{R}$ .

Consider  $A^{(i)}$  arbitrary row of  $A$  and  $V_k = \text{Span}\{v_1, \dots, v_k\}$  with  $V_k$  optimal  $k$ -dim best fitting space, i.e.

$$V_k = \arg \min_{\dim V \leq k} \sum_{i=1}^n \|A^{(i)} - P_V(A^{(i)})\|_2^2.$$

We know that

$$P_{V_k}(A^{(i)})^T = \underbrace{\sum_{l=1}^k \langle A^{(i)}, v_l \rangle v_l^T}_{\text{PF-Theorem for orthogonal projection}}$$

and

$$\langle A^{(i)}, v_l \rangle = (Av_l)_i.$$

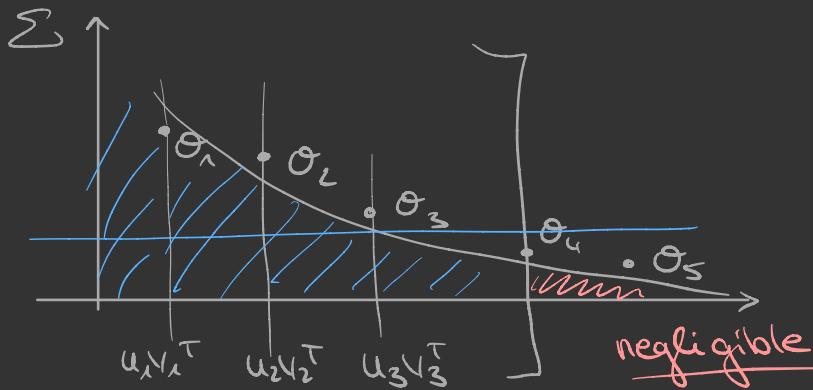
With this we get that

$$\sum_{l=1}^k A v_l \cdot v_l^T = \sum_{l=1}^k \underbrace{\sigma_l u_l v_l^T}_{} = A_k,$$

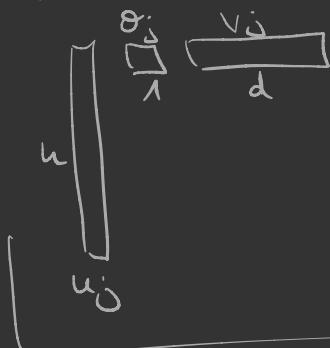
i.e. the rows of  $A_k$  are the orthogonal projections of the rows of  $A$  onto  $V_k$ .



Interpretation:



$$(d+n+1)(r-k) =$$



⇒ compress a matrix which has in principle  $n \times d$  components to a matrix defined by only  $(n+d+1)k$  parameters.

⇒ For  $k \ll r \ll \min\{n, d\}$ ,  $A_k$  is truly low complexity i.e.

$$(n+d+1)k \ll (n+d+1)r \ll n \times d.$$

Approximating  $A$  in terms of  $\|\cdot\|_F$  is done optimally by  $A_k$ .

Remark:  $A_k$  is the  $k$ -rank best approximation to  $A$  in any Schatten- $p$ -norm

$$\|A\|_p = \sum_{l=1}^n \sigma_l^{1/p} = \|\Sigma\|_p$$

↑  
 p-norm of  
 singular values

$$\text{with } \Sigma = (\sigma_1, \dots, \sigma_r)^T$$

This means for all  $p$ :

$$A_k = \underset{B \in R_k}{\operatorname{argmin}} \|A - B\|_p$$

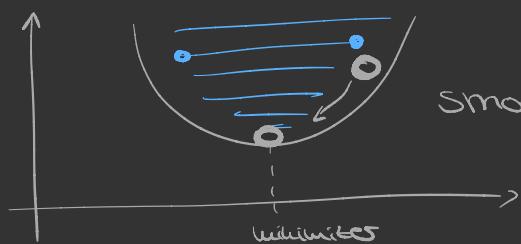
## Optimization problems

### ▷ "Easy" optimizations:

↳ typically unconstrained

↳ objective function to be minimized is smooth  
(differentiable with differential which is at least Lipschitz-continuous)

↳ objective function is convex.



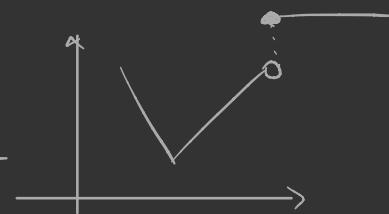
Smooth + epigraph is convex

use gradient descent for minimization

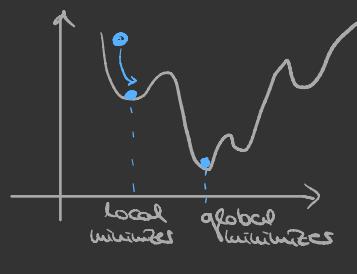
### ▷ "Hard" optimizations:

↳ constrained domain which is very complicated  
+ not convex.

↳ non-smooth objective function  
gradient information is not available



↳ non-convex functions  
with multiple local minimizers



$$\text{in our case: } A_k = \underset{B \in \mathbb{R}^{k \times n}}{\operatorname{argmin}} \|A - B\|_p$$

- norms are convex
- depending on  $p$ , the norm might not be differentiable ( $p=1$ )
- set of  $\mathbb{R}^k$  is no subspace i.e.  $A + B = C$  for  $\operatorname{rank}(C) > k$

but Solvable in closed form with SVD.





We have proven that

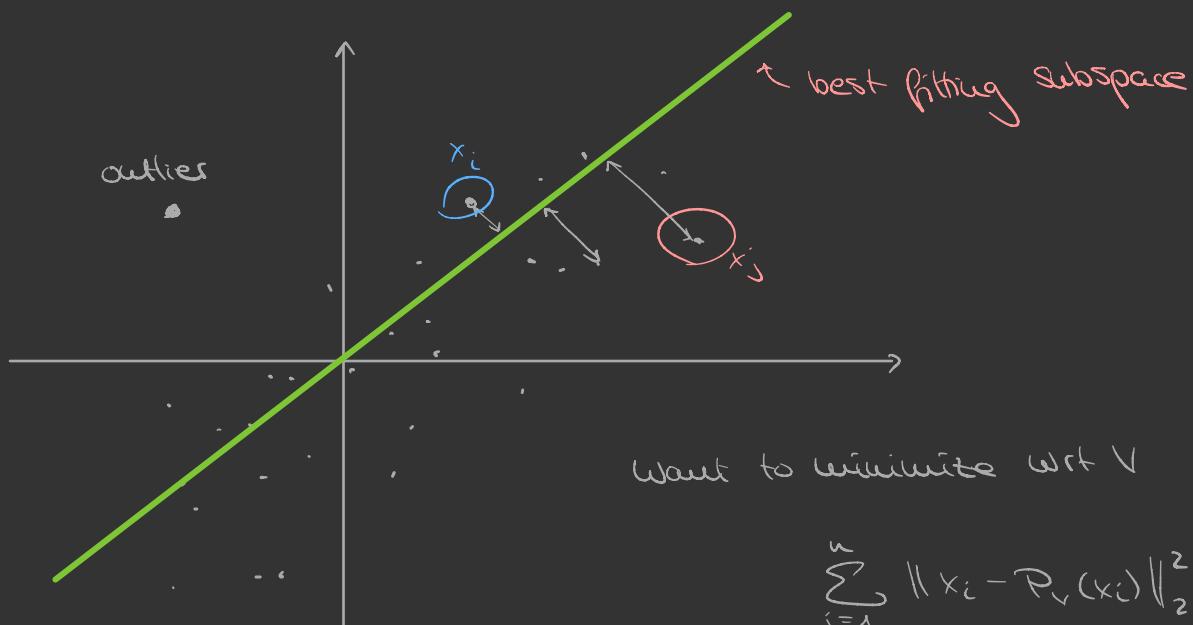
$$\theta_{k+1}^2 \geq \|A - B\|^2 \geq \|(A - B)z\|_2^2 = \|Az\|_2^2 \geq \theta_k^2$$

↑  
assumption

$$\Rightarrow \|A - B\| \geq \|A - A_k\|.$$

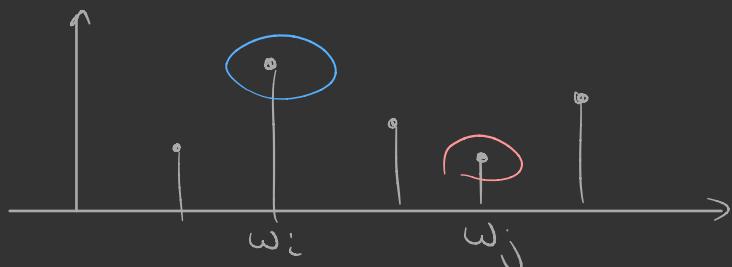
□

### Repetition on the intuition



But: there might be some points that are more important than others

→ introduce a sequence of weights  $w_j \geq 0$ , i.e.



then  $x_i$  is better approximated than  $x_j$ .

→ good way to handle outliers and counterweights that

e.g.  $\|\cdot\|_2$  strongly considers outliers.

→ alternative: use norms with  $0 < p < 1$ .

but this gives us a non-smooth, non-convex and constraint optimization problem.

Also: we can use other distances than the perpendicular distance, but it heavily depends on the application (would give more weight to a specific feature).

## Vorlesung 6

| 06.05.20

Everything so far boils down to the calculation of the SVD.

However: calculation of right singular vectors is itself an optimization problem  
 $\Rightarrow$  maximize over convex function

$$v_k = \underset{\|v\|_2=1}{\operatorname{argmax}} \|Av\|_2$$

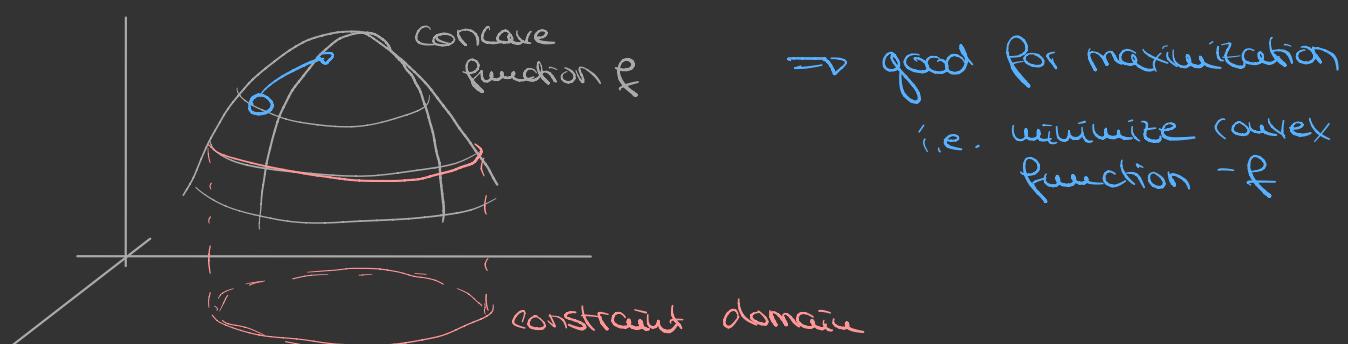
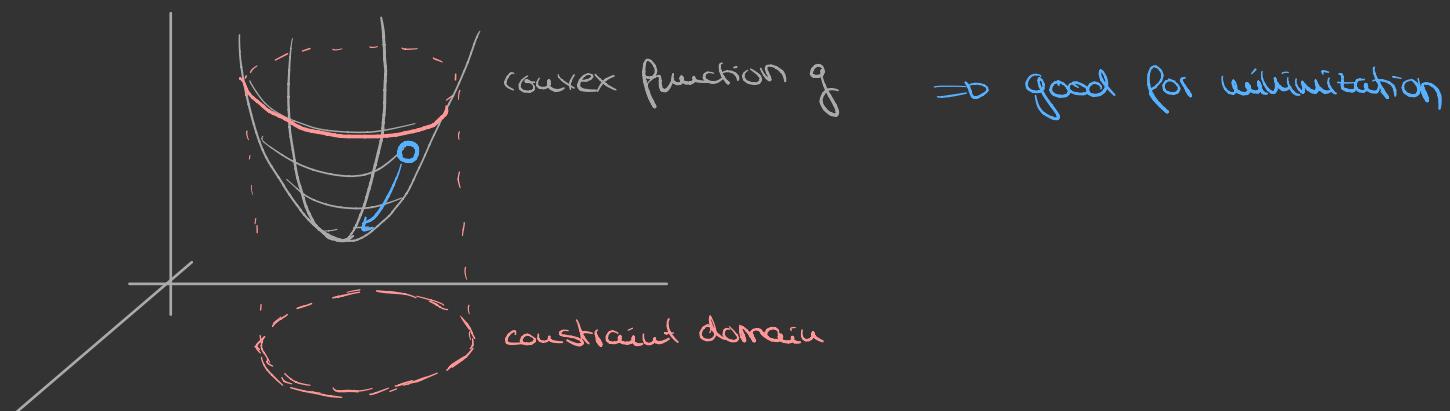
Ball is convex,  
but we need to stay on sphere  
 $\rightarrow$  surface of ball

$\left. \begin{matrix} & \\ & \end{matrix} \right\}$

$\left. \begin{matrix} & \\ & \end{matrix} \right\} \text{linear constraints } \checkmark$

$\checkmark$  we could optimize over convex ball if minimizing, since it's scaleable, but here we need to maximize.

More details on this assertion:



$\Rightarrow$  Since scaling any solution that already is close to a maximum to the border of the ball will only further decrease the function value (function is convex), we can simply consider the relaxation

$$v_k = \underset{\|v\|_2 \leq 1}{\operatorname{argmax}} \|Av\|_2$$

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

Additional challenges: "curse of dimensionality"

- > We usually work with Big Data, i.e. several billions of data points in millions of dimensions with pictures for example
- > the more dimensions, the more directions  $v$  we can choose from, so the complexity of our algorithm will scale depending on our dimension  $d$ .

Goal: Find an algorithm with complexity linear in  $d$ .

Method: "Power Method"

1) Easy case first:

Assume  $A \in \mathbb{R}^{n \times n}$  and symmetric. Then the left singular vectors equal the right singular vectors. Then

$$A = \sum_{k=1}^r \sigma_k v_k v_k^\top$$

For SVD we have

$$\begin{aligned} A^2 &= \left( \sum_{k=1}^r \sigma_k v_k v_k^\top \right) \left( \sum_{l=1}^r \sigma_l v_l v_l^\top \right) \\ &= \sum_{k,l} \sigma_k \sigma_l v_k v_k^\top \underbrace{v_l v_l^\top}_{= \delta_{kl} \text{ due to orthogonality}} \\ &= \sum_{k=1}^r \sigma_k^2 v_k v_k^\top. \end{aligned}$$

Induction shows:  $A^m = \sum_{k=1}^r \sigma_k^m v_k v_k^\top$ .

In case  $\sigma_1 \gg \sigma_2$ , we would get that

$$\lim_{m \rightarrow \infty} \frac{A^m}{\sigma_1^m} = v_1 v_1^\top.$$

Indeed:

$$\frac{A^m}{\sigma_1^m} = \sum_{k=1}^r \frac{\sigma_k^m}{\sigma_1^m} v_k v_k^\top = v_1 v_1^\top + \underbrace{\sum_{k=2}^r \left( \frac{\sigma_k}{\sigma_1} \right)^m v_k v_k^\top}_{\ll 1 \xrightarrow{m \rightarrow \infty} 0}.$$

III

We can easily extract  $v_1$  from  $v_1 v_1^\top$  by taking the square root of the diagonal of the matrix:

$$\begin{pmatrix} v_1^{(1)} \\ v_1^{(2)} \\ \vdots \\ v_1^{(d)} \end{pmatrix} \cdot \begin{pmatrix} v_1^{(1)} & v_1^{(2)} & \dots & v_1^{(d)} \\ v_1^{(2)} v_1^{(1)} & (\sqrt{v_1^{(2)}})^2 & & \\ \vdots & \ddots & \ddots & \\ v_1^{(d)} v_1^{(1)} & \dots & \ddots & (\sqrt{v_1^{(d)}})^2 \end{pmatrix}.$$

While calculation of  $\sigma_1$  is an difficult optimisation problem corresponding with the spectral norm  $\|A\|$ , it is easy to compute the Frobenius norm  $\|A\|_F$ .

Spectral norm:  $\sigma_1 = \max_{\|v\|_2=1} \|Av\|_2 = \max_{v \neq 0} \frac{\|Av\|_2}{\|v\|_2}$

Frobenius norm:  $\|A\|_F = \left( \sum_i \sum_j |A_{ij}|^2 \right)^{1/2}$

With this

$$\|A^m\|_F = \left( \sum_{i,j} \underbrace{(A \times \dots \times A)}_{m\text{-times}} )_{ij} \right)^{1/2}$$

we get

$$\lim_{m \rightarrow \infty} \frac{A^m}{\|A^m\|_F} = v_1 v_1^\top \quad \text{as soon as } \sigma_1 \gg \sigma_2.$$

Proof: Exercise.

Problem with this calculation

▷ need to calculate  $A^m$  for  $m \rightarrow \infty$ .

$$A_{m+1} = A \times A^m \quad \text{▷ polynomial complexity wrt. dim. of } A \text{ i.e. } C \geq n.d.$$

▷ need to calculate the Frobenius norm

$$\|A_{m+1}\|_F = \frac{\|A_{m+1}\|}{\|A_{m+1}\|_F}$$

2) If  $A \in \mathbb{R}^{n \times d}$  is least symmetric,  $n \neq d$ , we use

$$A^T A = \dots = \sum_{k=1}^r \sigma_k^2 v_k v_k^T$$

which itself is again squared and symmetric.

Then

$$\lim_{m \rightarrow \infty} \frac{(A^T A)^m}{\|A^T A\|_F} \stackrel{\sigma_1 \gg \sigma_2}{=} v_1 v_1^T.$$

Now: To reduce complexity we fix any vector  $x \in \mathbb{R}^d$ ,  $\|x\|_2 = 1$  and calculate

$$(A^T A)^m \cdot x = A^T A \left( A^T A \cdot \underbrace{\dots (A^T A x)}_{\in \mathbb{R}^d} \right)$$

which only requires a single matrix-matrix-multiplication, followed by significantly less complex matrix-vector-multip.

Then

$$\lim_{m \rightarrow \infty} \frac{(A^T A)^m x}{\|A^T A\|_F} = v_1 v_1^T \cdot x = \underbrace{\langle v_1, x \rangle}_{\in \mathbb{R}} \cdot v_1.$$

Two issues at this point

- we need  $\sigma_1 \gg \sigma_2$  for it to work
- What if we pick  $x$ , s.t.  $\langle v_1, x \rangle = 0$  ?  
numerically

Then our limit will approx. 0, and we have no useful result.

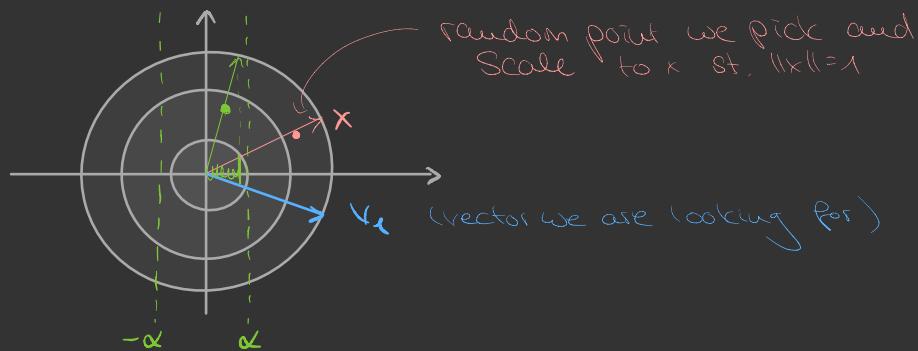
$\Rightarrow$  To solve the "curse of dimension", we will start using  
Probability (randomness) (adding to LA + optimization)

From this point on, we enter the space of randomized algorithms.



Example: Playing darts:

Terrible dart player will have same probability for any point on the target.



$\Rightarrow$  Probability to pick a point  $x$  which would be scaled to be  $v_1$  exactly is equal to zero

$\Rightarrow$  Probability to pick a point such that  $|<v_1, x>| \leq \alpha$  for  $\alpha \in \mathbb{R}$  small is volume of the area as part of  $B^d \subseteq \mathbb{R}^d$ .

Remember:  $V(d) = \text{volume of } B^d$   
 $\approx \frac{1}{d! \pi} \left( \frac{\pi e}{d} \right)^{d/2}$

Considering the probability to choose  $x \in S$ ,  $S \subseteq B^d$ , then

$$P(x \in S) = \frac{V_d(S)}{V(d)} \in [0, 1].$$

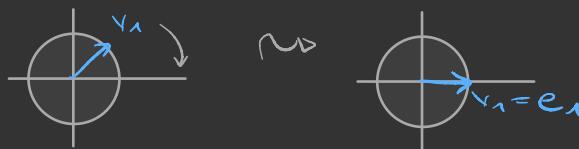
With this, we get that probability to pick  $x \in \mathbb{R}^d$  at random, such that  $|<v_1, x>| \leq \alpha$ ,  $\alpha$  small, is given by

Lemma 28: Let  $x \in \mathbb{R}^d$  be a unit  $d$ -dimensional vector of components  $x_1, \dots, x_d \in \mathbb{R}$ . We pick this vector at random. \*  
 Then, for  $\alpha > 0$

$$P(|x_1| \geq \alpha) \leq 1 - 2\alpha \sqrt{d-1}. \quad (1)$$

Note: If  $\alpha$  is small enough, the probability in (1) is very close to 1.

Remark: This corresponds with the assumption that  $v_1 = e_1$ . We can use this assumption wlog, because we can always rotate the coordinate system accordingly.



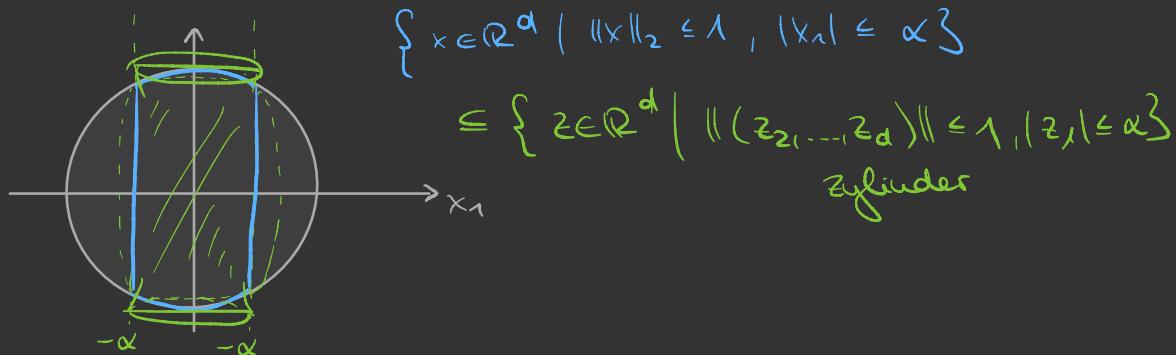
$$\text{Then } \langle x, v_1 \rangle = \langle x, e_1 \rangle = x_1.$$

Proof: Firstly, note that  $P(|x_1| \geq \alpha) = 1 - P(|x_1| < \alpha)$ .

Hence, we will prove that

$$P(|x_1| < \alpha) \leq 2\alpha \sqrt{d-1}.$$

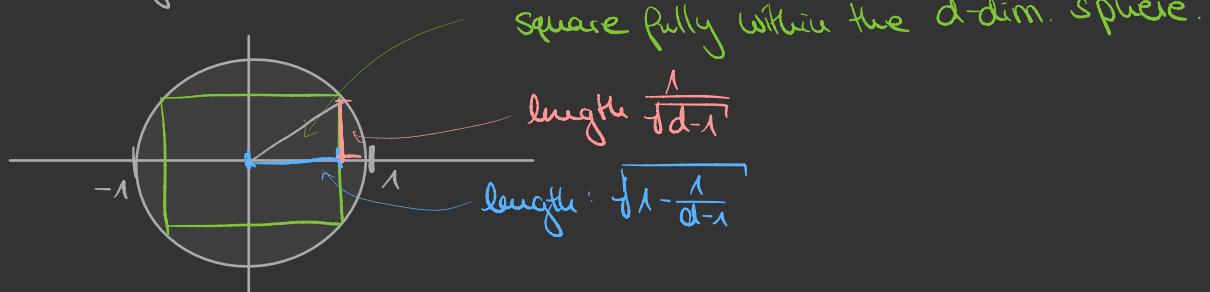
For this, let's picture the situation first



$$\Rightarrow \text{Vol}(\{\}) = V(d-1)2\alpha \quad \left( \text{Volumen Zylinder mit } (d-1) \text{ dim. Basis} \right)$$

$$\Rightarrow P(x \in \{\}) = \frac{\text{Vol}(\{\})}{V(d)} \leq \frac{\text{Vol}(\{\})}{V(d)} = \frac{2\alpha \cdot V(d-1)}{V(d)}.$$

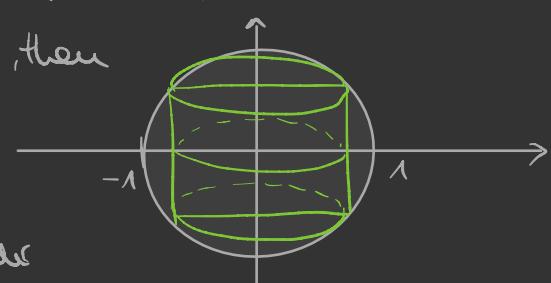
One more geometrical idea:



It is fully contained, since  $(\frac{1}{\sqrt{d-1}})^2 + (\sqrt{1 - \frac{1}{d-1}})^2 = 1$ .

Create a cylinder around the square, then

$$V(d) \geq 2 \underbrace{\frac{1}{\sqrt{d-1}}}_{\text{radius of cylinder}} \cdot \underbrace{\left(1 - \frac{1}{d-1}\right)^{\frac{d-1}{2}}}_{\text{base of cylinder}} V(d-1)$$



using  $V(R \cdot B^{d-1}) = V(B^{d-1})R^d$  with  $R$  radius.

With  $(1-\xi)^a \geq 1-a \cdot \xi$ , for  $\xi$  small and  $a$  large, we get  
that (Bernoulli'sche Ungleichung?)

$$(1 - \frac{1}{d-1})^{\frac{d-1}{2}} \geq 1 - (\frac{d-1}{2}) \frac{1}{d-1} = \frac{1}{2}.$$

Simplifying further, we get that

$$V(d) \geq \frac{1}{\sqrt{d-1}} \cdot V(d-1)$$

denominator  $\geq l b$   
 $\downarrow$

$$\frac{\text{numerator}}{\text{denominator}} \leq \frac{\text{numerator}}{l b}$$

With this we get

$$P(X \in \{\}) = \frac{2\alpha \cdot V(d-1)}{V(d)} \leq \frac{2\alpha \cdot \cancel{V(d-1)}}{\cancel{\frac{1}{\sqrt{d-1}} \cdot V(d-1)}} = 2\sqrt{d-1}\alpha.$$

□

Remark 2.9: Lemma 2.8 shows that independently of  $d$ , the  $X_i = \langle X, u_i \rangle$  component of a random unit vector  $X$  w.r.t. any orthogonal basis  $\{u_1, \dots, u_d\}$  is bounded away from zero with overwhelming probability.

Remark 2.10: Identifying  $\mathbb{C}$  with  $\mathbb{R}^2$ , we can extend the result of lemma 2.8 to random unit vectors in  $\mathbb{C}^d$  as follows:

The probability that, for a randomly chosen unit vector  $z \in \mathbb{C}^d$ ,  $|z_1| \geq \kappa > 0$  holds is at least  $1 - 2\alpha \sqrt{2d-1}$ .

# Worlesung 7

12.05.

Theorem 2.1 (speed of convergence of the power method)

Let  $A \in \mathbb{K}^{I \times J}$  and  $x \in \mathbb{K}^I$  be a random, unit length vector.

Let  $V$  be the space spanned by the left singular vectors of  $A$  corresponding to singular values  $\sigma_k \geq (1-\varepsilon)\sigma_1, k \in \{1, \dots, I\}$ .

Let  $m \in \mathbb{N}$ ,  $m \in \Omega\left(\frac{\ln(d/\varepsilon)}{\varepsilon}\right)$  for  $\beta \geq 1/2$ . Let  $w^*$  be the unit vector after  $m$  iterations of the power method, i.e.

$$w^* := \frac{(AA^H)^m \cdot x}{\|(AA^H)^m x\|_2}$$

Then, the probability that  $w^*$  has a component of at most  $\Theta\left(\frac{\varepsilon}{\alpha d^\beta}\right)$  orthogonal to  $V$  is  $\geq 1 - 2\sqrt{2d-1}$ .

error of the approximation

Note: For  $\alpha$  small, we get a better probability, but a worse approximation.

Proof: Let  $A = \sum_{k=1}^r \sigma_k u_k v_k^H$  be the SVD of  $A$ . If  $r \leq n = |I|$ , complete  $\{u_1, \dots, u_r\}$  to an orthonormal basis  $\{u_1, \dots, u_r, \dots, u_n\}$  of  $\mathbb{K}^n$ . Then for  $x \in \mathbb{K}^n$ :

$$x = \sum_{i=1}^n \langle x, u_i \rangle u_i$$

Since

$$AA^H = \sum_{k=1}^r \sigma_k u_k v_k^H \cdot \sum_{i=1}^n \sigma_i v_i u_i^H$$

$$\begin{aligned} \text{Cauchy-} \\ \text{product} \quad &= \sum_{k=1}^r \sum_{j=0}^{r-1} \underbrace{\sigma_k u_k v_k^H}_{\delta_{k,k-j}} \cdot \underbrace{\sigma_{k-j} v_{k-j} u_{k-j}^H}_{\delta_{k,k-j} \text{ due to orthogonality}} \\ &= \sum_{k=1}^r \sigma_k^2 u_k u_k^H \end{aligned}$$

we get that

$$(AA^H)^2 = \left( \sum_{k=1}^r \sigma_k^2 u_k u_k^H \right)^2$$

$$\begin{aligned}
&= \left( \sum_{k=1}^r \sum_{j=0}^{k-1} \sigma_k^2 u_k u_k^H \cdot \sigma_{k-j}^2 u_{k-j} u_{k-j}^H \right) \\
&= \left( \sum_{k=1}^r \sigma_k^{2.2} u_k u_k^H \right).
\end{aligned}$$

By induction, this yields

$$(AA^H)^m = \sum_{k=1}^r \sigma_k^{2m} u_k u_k^H.$$

Setting  $\sigma_k = 0$  for  $r+1 \leq k \leq n$ , we write the extended sum

$$(AA^H)^m = \sum_{k=1}^n \sigma_k^{2m} u_k u_k^H.$$

and hence for  $x \in \mathbb{K}^d$  picked at random

$$(AA^H)_x^m = \sum_{k=1}^n \sigma_k^{2m} u_k \langle x, u_k \rangle.$$

We note that since  $u_k$ 's are fixed, but  $x$  was random, it is as we would pick  $\langle x, u_k \rangle$  at random.

By Lemma 2.8 and Remark 2.9 and 2.10, the probability that  $|\langle x, u_k \rangle| \geq \alpha > 0$  is at least  $1 - 2\alpha \sqrt{2d-1}$ .

Suppose now that  $\sigma_1, \sigma_2, \dots, \sigma_r$  are the singular values of  $A$  that are  $\geq (1-\varepsilon)\sigma_1$ , i.e.

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r \geq (1-\varepsilon)\sigma_1 \geq \sigma_{r+1} \geq \dots \geq \sigma_n.$$

For  $\sigma_1 > \sigma_2 \exists \varepsilon > 0$  small enough s.t.  $\sigma_\varepsilon = 0$ .  $\square$

By Pythagoras-Fourier theorem, we have

$$\begin{aligned}
\|(AA^H)_x^m\|_2^2 &= \left\| \sum_{k=1}^n \sigma_k^{2m} \langle x, u_k \rangle u_k \right\|_2^2 = \\
&\stackrel{\text{w.l.o.g.}}{=} \sum_{k=1}^n \sigma_k^{4m} |\langle x, u_k \rangle|^2 \\
&\geq \sigma_1^{4m} |\langle x, u_1 \rangle|^2 \\
&\stackrel{\text{w.l.o.g.}}{\geq} \alpha^2 \sigma_1^{4m}
\end{aligned}$$

whp  
with high prob.

with probability at least  $1 - 2\alpha \sqrt{2d-1}$ ,

which gives us a lower bound for  $\|(AA^H)^m x\|_2^2$ .

Now consider the component of  $(AA^H)^m x$ , which is orthogonal w.r.t.  $V = \text{span}\{u_1, \dots, u_{r_\varepsilon}\}$ . We write

$$(AA^H)^m x = \sum_{k=1}^{r_\varepsilon} \Theta_k^{2m} \langle x, u_k \rangle u_k \in V$$

$$+ \sum_{k=r_\varepsilon+1}^n \Theta_k^{2m} \langle x, u_k \rangle u_k \in V^\perp.$$

Then

$$\left\| \sum_{k=r_\varepsilon+1}^n \Theta_k^{2m} |\langle x, u_k \rangle| \right\|_2^2 \stackrel{\text{P.F.T.}}{=} \sum_{k=r_\varepsilon+1}^n \Theta_k^{4m} |\langle x, u_k \rangle|^2$$

$$\stackrel{\text{by construction}}{\leq} (1-\varepsilon)^{4m} \Theta_1^{4m} \sum_{k=r_\varepsilon+1}^n |\langle x, u_k \rangle|^2$$

$$\leq (1-\varepsilon)^{4m} \Theta_1^{4m} \leq 1$$

Since  $\sum_{k=r_\varepsilon+1}^n |\langle x, u_k \rangle|^2 \leq \sum_{k=1}^n |\langle x, u_k \rangle|^2 = \|x\|_2^2 = 1$  by the Pythagoras-Fourier theorem.

In summary, we get

$$\frac{\|\mathcal{P}_{V^\perp}((AA^H)^m x)\|_2^2}{\|(AA^H)^m x\|_2^2} \stackrel{\text{upper bound}}{\leq} \frac{\Theta_1^{4m} (1-\varepsilon)^{4m}}{\alpha^2 \cdot \Theta_1^{4m}} = \frac{(1-\varepsilon)^{4m}}{\alpha^2}.$$

Using the Taylor expansion at  $\varepsilon=0$ , we get  $(1-\varepsilon) \approx e^{-\varepsilon}$  and hence we can rewrite

$$\frac{(1-\varepsilon)^{4m}}{\alpha^2} = \Theta(\alpha^{-2} e^{-4\varepsilon m})$$

$$m = C \cdot \left( \frac{\ln(d^\beta/\varepsilon)}{\varepsilon} \right) = \Theta(\alpha^{-2} \exp \left( \underbrace{-4 \cdot \varepsilon \cdot C \cdot \ln(d^\beta/\varepsilon) \cdot \varepsilon^{-1}}_{= \ln \left( \left( \frac{d^\beta}{\varepsilon} \right)^{-4C} \right)} \right))$$

$$= \Theta\left(\frac{\varepsilon}{\alpha^2 \cdot d^{4\beta C}}\right).$$

## Relation between left & right singular vectors:

$v$ 's are right singular vectors

$$\frac{(A^T A)^m x}{\|(A^T A)^m x\|_2} \longrightarrow v_1 v_1^T x$$

and  $u$ 's are left singular vectors

$$\frac{(A A^T)^m x}{\|(A A^T)^m x\|_2} \longrightarrow u_1 u_1^T x$$

i.e.

$$A^T A = \sum_{k=1}^d \sigma_k^2 v_k v_k^T.$$

and

$$A A^T = \sum_{k=1}^d \sigma_k^2 u_k u_k^T.$$

Having  $v_1$ , then  $u_1 = \frac{A v_1}{\|A v_1\|}$  , or

having  $u_1$ , then  $v_1 = \frac{A^T u_1}{\|A^T u_1\|}$ .

Calculation of  $\sigma_2, \dots, \sigma_d$ , once we have  $\sigma_1$ :

Compute  $A_1 = A - \sigma_1 u_1 v_1^T$  "deflation of  $A$ "

and again execute the SVD-powers method again.

No Problem: not as stable as it can be.

For  $A \in \mathbb{R}^{n \times d}$ :

$$A^T A \underset{d \times d}{=} \begin{array}{c|c} \xrightarrow{\quad n \quad} & \xleftarrow{\quad d \quad} \\ \boxed{\phantom{0}} & \boxed{\phantom{0}} \end{array} \in \mathbb{R}^{d \times d}$$

$$\Rightarrow v_1 \in \mathbb{R}^d$$

$$A A^T \underset{n \times n}{=} \begin{array}{c|c} \xleftarrow{\quad d \quad} & \xrightarrow{\quad n \quad} \\ \boxed{\phantom{0}} & \boxed{\phantom{0}} \end{array} \in \mathbb{R}^{n \times n}$$

$$\Rightarrow u_1 \in \mathbb{R}^n$$

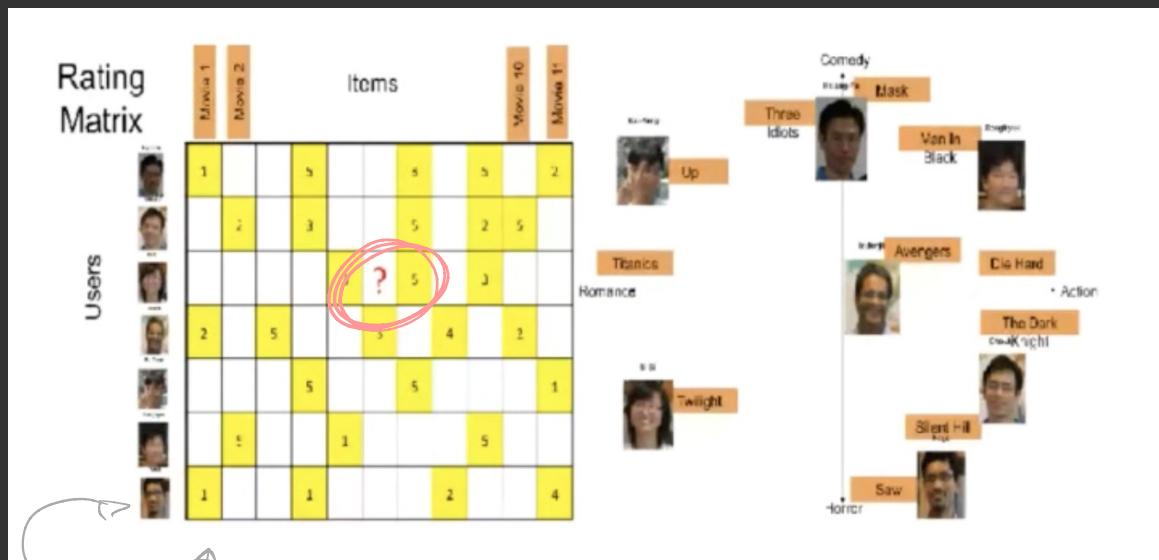
# Vorlesung 8

13.05.2020

## 2.5. Applications of SVD

### 2.5.1. Principal Component Analysis

Example: Customers Product Matrix



↑ customer n bought product 1 a single time.

Goal: complete this matrix

Problem: Let  $A$  with elements  $a_{ij}$  represent the probability of customer  $i$  to buy product  $j$ .

Assume:  $\Rightarrow k$  underlying basic factors like age, income, family...  
to determine customer's behavior (weighted combination)  
 $\Rightarrow$  look for these factors in our matrix.

$\Rightarrow$  characterize the  $i$ th customer's behavior by

$$u_i = (u_{l,i})_{l=1,\dots,k} \in \mathbb{R}^k \quad \left. \begin{array}{l} \text{weights for each basic factor } l. \\ \text{how relevant are the factors for a customer?} \end{array} \right\}$$

$\Rightarrow$  vector of probabilities  $v_l = (v_{l,j})_{j=1,\dots,d}$   $\quad \left. \begin{array}{l} \text{how does each factor effect the prob. to buy a product?} \end{array} \right\}$

representing the probability of purchasing product  $j$ , when only looking on basic factor  $l$ .

$\Rightarrow \exists U \in \mathbb{C}^{n \times k}$  and  $V \in \mathbb{C}^{k \times d}$  matrices s.t.

$$A = UV^T$$

↳ often no equality here

Theu

- » Noise is often contained in the reformation of A
  - » Don't know how large the K should be chosen

## Real-Life Example : Netflix - Recommendation

How can we find  $\Phi_{\mathcal{K}}, \Psi_{\mathcal{K}} \in \Omega^c$ ?

