

Session 7: Spectral Theorem, PCA and SVD

Optimization and Computational Linear Algebra for Data Science

Contents

1. The Spectral Theorem
2. Principal Component Analysis
3. Singular Value Decomposition

Midterm

- ❖ The Midterm exam is in 2 weeks.
- ❖ **Scope:** everything that we have seen so far (this week's video included).
- ❖ **Knowing is not enough!** You need to practice: review problems available on the course's webpage.
- ❖ Past years midterms also available, with solutions.
- ❖ **Important:** when working on a problem, take **at least** 10min on it before looking at the solution (in case you are stuck).
- ❖ The midterm is open books/notes, but **if you think that you need them for the exam, this probably means that you are not prepared enough.**

The Spectral Theorem

The spectral theorem

$$\text{Id}x = 1 \cdot x$$

Theorem

Let $A \in \mathbb{R}^{n \times n}$ be a **symmetric** matrix. Then there is a **orthonormal** basis of \mathbb{R}^n composed of eigenvectors of A .

$$P = \begin{pmatrix} | & & | \\ v_1 & \dots & v_n \\ | & & | \end{pmatrix}$$

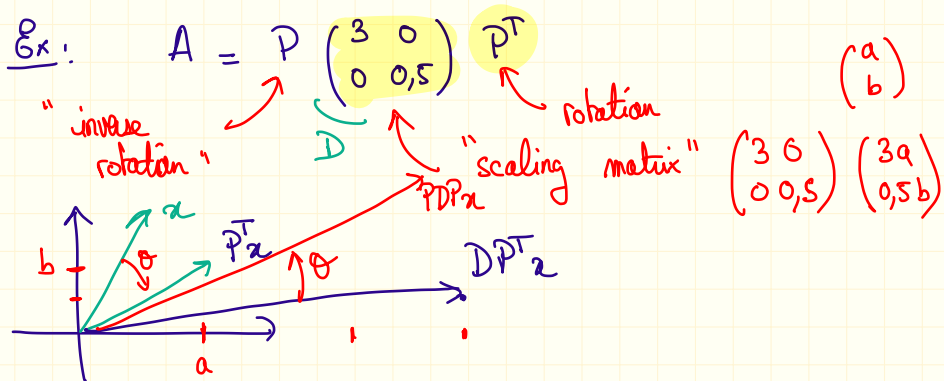
$$D = \begin{pmatrix} \lambda_1 & & (0) \\ & \ddots & \\ (0) & & \lambda_n \end{pmatrix}$$

Theorem (Matrix formulation)

Let $A \in \mathbb{R}^{n \times n}$ be a **symmetric** matrix. Then there exists an **orthogonal** matrix P and a **diagonal** matrix D of sizes $n \times n$ such that

$$A = PDP^T.$$

Geometric interpretation



The Theorem behind PCA

Theorem

by
Spec.
Thm

Let A be a $n \times n$ symmetric matrix and let $\lambda_1 \geq \dots \geq \lambda_n$ be its n eigenvalues and v_1, \dots, v_n be an associated orthonormal family of eigenvectors. Then

$$\lambda_1 = \max_{\substack{v \in \mathbb{R}^n \\ \|v\|=1}} \underbrace{v^T A v}_{\text{circled}}$$

and

$$v_1 = \arg \max_{\|v\|=1} v^T A v.$$

Moreover, for $k = 2, \dots, n$:

$$\lambda_k = \max_{\|v\|=1, v \perp v_1, \dots, v_{k-1}} v^T A v, \quad \text{and} \quad v_k = \arg \max_{\|v\|=1, v \perp v_1, \dots, v_{k-1}} v^T A v.$$

$k=2$

$$\lambda_2 = \max_{\substack{\|v\|=1 \\ v \perp v_1}} v^T A v$$

$$v_2 = \arg \max_{\substack{\|v\|=1 \\ v_2 \perp v_1}} v^T A v$$

Proof

- let $\sigma \in \mathbb{R}^n$ such that $\|\sigma\|=1$. Let $(\alpha_1 \dots \alpha_n)$ be the coordinates of σ in $B = (\sigma_1, \dots, \sigma_n)$.
- $$\begin{aligned} \underline{A\sigma} &= A(\alpha_1 \sigma_1 + \dots + \alpha_n \sigma_n) \\ &= \alpha_1 A\sigma_1 + \dots + \alpha_n A\sigma_n \\ &= \alpha_1 \lambda_1 \sigma_1 + \dots + \alpha_n \lambda_n \sigma_n \rightarrow \end{aligned}$$

$\begin{pmatrix} \alpha_1 \lambda_1 \\ \vdots \\ \alpha_n \lambda_n \end{pmatrix}$ are the coords. of $A\sigma$ in B .
- $$\begin{aligned} \underline{\sigma^T A \sigma} &= \underline{\sigma} \cdot \underline{(A\sigma)} = \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{pmatrix} \cdot \begin{pmatrix} \alpha_1 \lambda_1 \\ \vdots \\ \alpha_n \lambda_n \end{pmatrix} \\ &= \underline{\alpha_1^2 \lambda_1 + \dots + \alpha_n^2 \lambda_n} \end{aligned}$$

Proof

Maximize

$$\alpha_1^2 \lambda_1 + \dots + \alpha_n^2 \lambda_n$$

subject to $\|v\|^2 = 1$

$v^T A v$

$$1 = \alpha_1^2 + \dots + \alpha_n^2$$

• Since $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$, the maximum is achieved for $\alpha_1 = 1$, $\alpha_2 = \dots = \alpha_n = 0$.

• This corresponds to $v = v_1$

• The corresponding value of $v^T A v$ is then $1^2 \cdot \lambda_1 = \lambda_1$

Proof

if now we maximize $\sigma^T A \sigma$ subject to $\begin{cases} \|\sigma\|=1 \\ \sigma \perp \sigma_1 \end{cases}$
we maximize $\underbrace{\langle \sigma, \sigma_1 \rangle}_{\alpha_1} = 0$

$$\cancel{\alpha_1^2} \lambda_1 + \alpha_2^2 \lambda_2 + \dots + \alpha_n^2 \lambda_n \quad \text{subj. to} \quad \begin{cases} \alpha_1^2 + \dots + \alpha_n^2 = 1 \\ \alpha_1 = 0 \end{cases}$$

The maximum is now achieved for

$$\begin{cases} \alpha_1 = 0 \\ \alpha_2 = 1 \\ \alpha_3 = \dots = \alpha_n = 0 \end{cases}$$

σ :

$$\begin{cases} \alpha_1 = \langle \sigma_1, \sigma \rangle \\ \vdots \\ \alpha_n = \langle \sigma_n, \sigma \rangle \end{cases}$$

Principal Component Analysis

Empirical mean and covariance

We are given a dataset of n points $a_1, \dots, a_n \in \mathbb{R}^d$

$$\underline{d = 1}$$

❖ Mean

$$\mu = \frac{1}{n} \sum_{i=1}^n a_i \in \mathbb{R}$$

❖ Variance

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (a_i - \mu)^2 \in \mathbb{R}$$

Empirical mean and covariance

We are given a dataset of n points $a_1, \dots, a_n \in \mathbb{R}^d$

$$\underline{d = 1}$$

Mean

$$\mu = \frac{1}{n} \sum_{i=1}^n a_i \in \mathbb{R}$$

Variance

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n \underline{(a_i - \mu)^2} \in \mathbb{R}$$

$$\underline{d \geq 2}$$

Mean

$$\mu = \frac{1}{n} \sum_{i=1}^n a_i \in \mathbb{R}^d$$

Covariance matrix

$$S = \frac{1}{n} \sum_{i=1}^n \underline{(a_i - \mu)(a_i - \mu)^T} \in \mathbb{R}^{d \times d}$$
$$= \frac{1}{n} \sum_{i=1}^n \underline{a_i a_i^T} \quad \text{if } \underline{\mu = 0.}$$

PCA

- ❖ We are given a dataset of n points $a_1, \dots, a_n \in \mathbb{R}^d$, where d is «large».
- ❖ **Goal:** represent this dataset in lower dimension, i.e. find $\tilde{a}_1, \dots, \tilde{a}_n \in \mathbb{R}^k$ where $k \ll d$.
- ❖ Assume that the dataset is centered: $\sum_{i=1}^n a_i = 0$.
- ❖ Then, S can be simply written as:

"cov. matrix
without the $1/n$ "

$$S = \sum_{i=1}^n a_i a_i^T = A^T A.$$

$$= BB^T$$

where A is the $n \times d$ "data matrix":

$$B = \begin{pmatrix} 1 & & \\ a_1 & \dots & a_n \\ 1 & & \end{pmatrix} = A^T$$

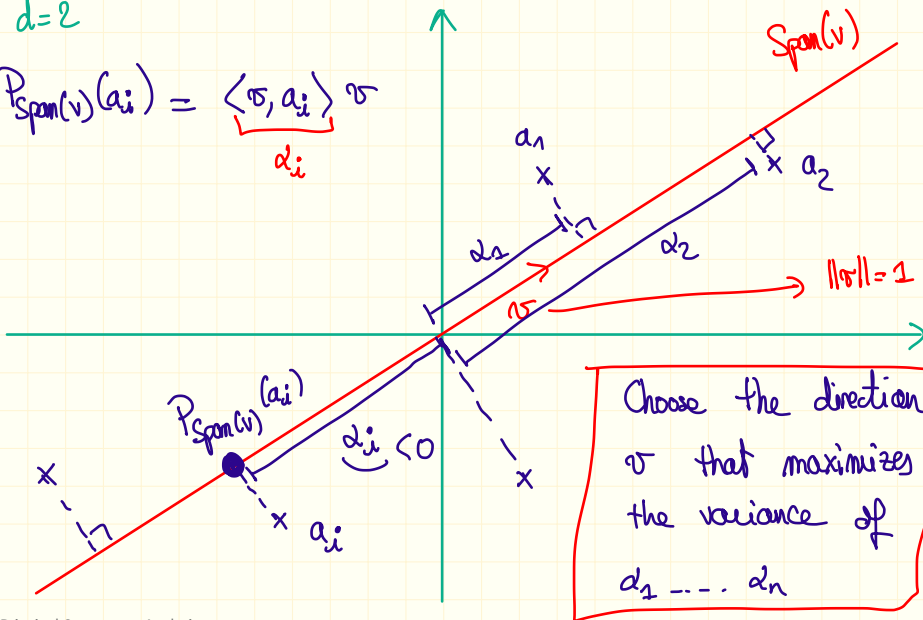
$$A = \begin{pmatrix} -a_1^T - \\ \vdots \\ -a_n^T - \end{pmatrix}$$



Direction of maximal variance

$d=2$

$$P_{\text{Span}(v)}(a_i) = \underbrace{\langle v, a_i \rangle}_{d_i} v$$



Direction of maximal variance

Mean : $\frac{a_1 + \dots + a_n}{n} = \frac{\langle \sigma, a_1 \rangle + \dots + \langle \sigma, a_n \rangle}{n}$

$= \langle \sigma, \frac{a_1 + \dots + a_n}{n} \rangle = 0$

Variance : $\frac{1}{n} \sum_{i=1}^n \langle \sigma, a_i \rangle^2 = \frac{1}{n} \sum_{i=1}^n \sigma^T a_i a_i^T \sigma$

$\langle \sigma, a_i \rangle = \sigma^T a_i$
 $= a_i^T \sigma$

$= \frac{1}{n} \sigma^T \left(\sum_{i=1}^n a_i a_i^T \right) \sigma$

$= \frac{1}{n} \boxed{\sigma^T S \sigma}$

Direction of maximal variance

Good news: $S = A^T A$ is symmetric. $\leftarrow S = A^T A \leftarrow \begin{matrix} \cdot \text{sym.} \\ \cdot \text{PSD} \end{matrix}$

Spectral Theorem: let $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$ be the eigenvalues of S and (v_1, \dots, v_d) an associated orthonormal basis of eigenvectors. $\rightarrow 0$

• By the theorem we saw before a ^{unit} vector v that maximises $v^T S v$ is $v = v_1$

• and the corresponding "variance" $v^T S v$ is equal to λ_1 .

• The dimensionally-reduced dataset is then $\tilde{x}_1 = \langle v_1, x_1 \rangle \dots x_n = \langle v_1, x_n \rangle$

2nd direction of maximal variance

- We would like to find another vector v , such that the variance of $\langle v, a_1 \rangle \dots \langle v, a_n \rangle$ is large.

- Goal: find v that maximizes

$$\boxed{v^T S v \quad \text{subject to} \quad \begin{cases} \|v\| = 1 \\ v \perp v_1 \end{cases}}$$

- An optimal vector is $\boxed{v = v_2}$ $\lambda_1 > \lambda_2 > \dots$ $\lambda_2 = \lambda_3$

$$\tilde{a}_1 = \begin{pmatrix} \langle v_1, a_1 \rangle \\ \langle v_2, a_1 \rangle \end{pmatrix} \quad \text{---} \quad \tilde{a}_n = \begin{pmatrix} \langle v_1, a_n \rangle \\ \langle v_2, a_n \rangle \end{pmatrix}$$

1st princ. comp. of a_1

j^{th} direction of maximal variance

- The « j^{th} direction of maximal variance » is v_j since v_j is solution of

$$\text{maximize } \underline{v^T S v}, \quad \text{subject to } \|v\| = 1, \quad v \perp v_1, v \perp v_2, \dots, v \perp v_{j-1}.$$

- The dimensionally reduced dataset is then

$$\underbrace{\begin{pmatrix} \langle v_1, a_1 \rangle \\ \langle v_2, a_1 \rangle \\ \vdots \\ \langle v_k, a_1 \rangle \end{pmatrix}}_{\text{column 1}}, \underbrace{\begin{pmatrix} \langle v_1, a_2 \rangle \\ \langle v_2, a_2 \rangle \\ \vdots \\ \langle v_k, a_2 \rangle \end{pmatrix}}_{\text{column 2}}, \underbrace{\begin{pmatrix} \langle v_1, a_3 \rangle \\ \langle v_2, a_3 \rangle \\ \vdots \\ \langle v_k, a_3 \rangle \end{pmatrix}}_{\text{column 3}} \cdots \underbrace{\begin{pmatrix} \langle v_1, a_n \rangle \\ \langle v_2, a_n \rangle \\ \vdots \\ \langle v_k, a_n \rangle \end{pmatrix}}_{\text{column n}} \cdot \mathbb{R}^k$$

Recap

- ① Center your dataset \rightarrow get a dataset such that $\sum_{i=1}^n a_i = 0$.
- ② Compute the covariance matrix $S = \sum_{i=1}^n a_i a_i^T$
($= A^T A$)
- ③ Compute the eigenvalues $\lambda_1 \rightarrow \lambda_d$ of S
and associated eigenvectors $as_1 \rightarrow as_d$ of S .
- ④ Sort eigenvalues / eigenvectors.
- ⑤ Select some k .
- ⑥ Compute $\tilde{a}_1 \dots \tilde{a}_n$.

Which value of k should we take?

1st way

When using the k first principal components we capture a fraction

$f_k =$

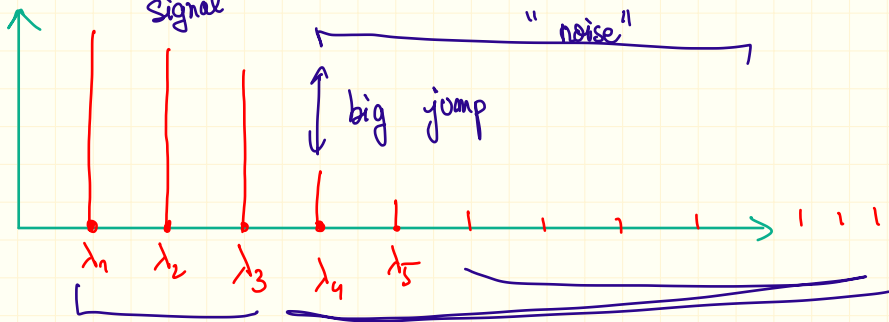
$$\frac{\lambda_1 + \lambda_2 + \dots + \lambda_k}{\lambda_1 + \lambda_2 + \dots + \lambda_k + \dots + \lambda_d}$$

of the total variance.

Choose k such that $f_k \geq 80\%$

Which value of k should we take?

2nd way: Plot the eigenvalues of S in decreasing order



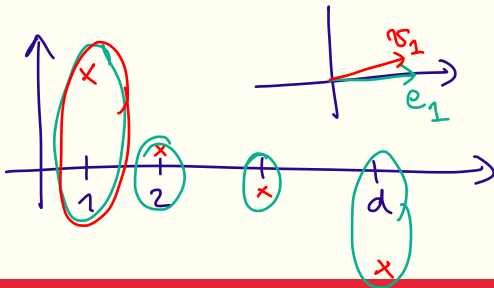
→ select $k = 3$.

$$v_1 = \begin{pmatrix} 0,8 \\ 0,1 \\ 0,2 \end{pmatrix}$$

$$a_1 = \begin{pmatrix} a_{11} \\ \vdots \\ a_{nd} \end{pmatrix}$$

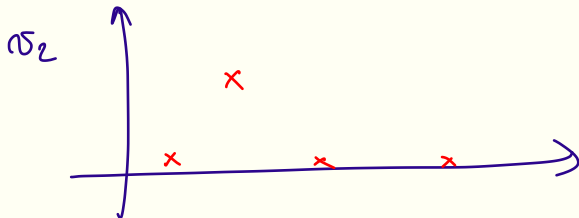
10 000

entries
of v_1



$$a_1 = \begin{pmatrix} f_1 \\ f_2 \\ \vdots \\ f_d \end{pmatrix}$$

Singular Value Decomposition



$$f_1 = \langle a_1, e_1 \rangle$$

Singular values/vectors

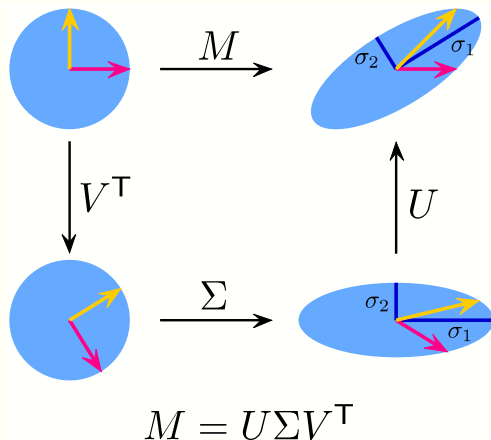
Singular Value decomposition

Theorem

Let $A \in \mathbb{R}^{n \times m}$. Then there exists two orthogonal matrices $U \in \mathbb{R}^{n \times n}$ and $V \in \mathbb{R}^{m \times m}$ and a matrix $\Sigma \in \mathbb{R}^{n \times m}$ such that $\Sigma_{1,1} \geq \Sigma_{2,2} \geq \dots \geq 0$ and $\Sigma_{i,j} = 0$ for $i \neq j$, that verify

$$A = U\Sigma V^T.$$

Geometric interpretation of $U\Sigma V^T$



Questions?

$$d \times d \left(A^T A \stackrel{?}{=} \sum_{i=1}^n a_i a_i^T \right) d \times d$$

$$A = \begin{pmatrix} \text{---} a_1^T \text{---} \\ \text{---} a_n^T \text{---} \end{pmatrix}$$

$$\underline{\underline{(A^T A)_{k,l}}} = \sum_{i=1}^n \overbrace{(A^T)_{k,i}}^{A_{i,k}} A_{i,l}$$

$$= \sum_{i=1}^n \underbrace{(a_i)_k (a_i)_l}_{\substack{\text{---} a_i^T \text{---} \\ \left(\begin{smallmatrix} a_i \end{smallmatrix} \right) \left(\begin{smallmatrix} d \times d \end{smallmatrix} \right)}} = (a_i a_i^T)_{k,l}$$

$$= \left(\sum_{i=1}^n \overbrace{a_i a_i^T}^{d \times d} \right)_{k,l}$$

$$y = ax + b$$

