

Exploratory Data Analysis

- Numerical data description (**robust**)
 - Visualisation [Labs]
- Understand structure of data

EDA linked to Preprocessing

- Standardisation
 - Centring $[C_n = I_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T]$
 - Scaling to unit variance
- Outlier detection / removal
 - Tukey's fences

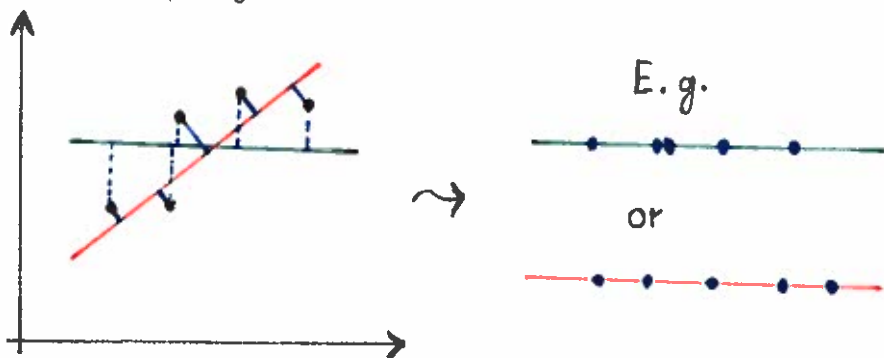


How to represent our data?

- to find and exploit structure
- to facilitate further processing

Quiz: Useful data direction

- Given: 2D data
- Represent using just 1D data
- Linear projection



Should we select the direction such that we:

- A) maximise the variance
- B) minimise the approximation error (squared)
- C) approximate the sample covariance matrix as well as possible (in the Frobenius norm)
- D) either A, B or C depending on the data structure
- E) do something else
- F) ... you don't know

Principal Component Analysis

First view: Find directions in data space along which data are maximally variable

→ "PC directions"

First PC direction: \underline{w}_1

First PC: Projected variable $z_1 = \underline{w}_1^T \underline{x}$

where \underline{x} : centred random variable
 $\text{Cov}[\underline{Ax} + \underline{b}] = \underline{A} \text{Cov}[\underline{x}] \underline{A}^T$

$$\text{Var}[z_1] = \text{Var}[\underline{w}_1^T \underline{x}] = \underline{w}_1^T \underline{\Sigma} \underline{w}_1$$

where $\underline{\Sigma} = \text{Cov}[\underline{x}]$

→ Optimisation problem:

$$\underset{\underline{w}_1}{\text{maximise}} \quad \underline{w}_1^T \underline{\Sigma} \underline{w}_1$$

$$\text{s.t. } \|\underline{w}_1\| = 1$$

EV decomposition: $\Sigma = U \Lambda U^T$

W.l.o.g. $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d \geq 0$

U orth. basis $\Rightarrow \underline{w}_1 = \sum_{i=1}^d \alpha_i \underline{u}_i = U \underline{\alpha}$

$$\Rightarrow \underline{w}_1^T \Sigma \underline{w}_1 = \underline{\alpha}^T U^T U \Lambda U^T U \underline{\alpha} = \underline{\alpha}^T \Lambda \underline{\alpha} = \sum_{i=1}^d \alpha_i^2 \lambda_i$$

$$\|\underline{w}_1\|^2 = \underline{w}_1^T \underline{w}_1 = \underline{\alpha}^T U^T U \underline{\alpha} = \sum_{i=1}^d \alpha_i^2 \stackrel{!}{=} 1$$

maximise $\alpha_1, \dots, \alpha_d$

Solution: $\underline{\alpha} = (1 \ 0 \ 0 \dots 0)^T$

$\Rightarrow \underline{w}_1 = U \begin{pmatrix} 1 \\ 0 \\ \vdots \end{pmatrix} = \underline{u}_1$ i.e. EV with largest eigenvalue

PC: $z_1 = \underline{w}_1^T \underline{x}$

Properties: $E[z_1] = E[\underline{w}_1^T \underline{x}] = \underline{w}_1^T E[\underline{x}] = 0$

$$\text{Var}[z_1] = \underline{w}_1^T \Sigma \underline{w}_1 = \sum_{i=1}^d \alpha_i^2 \lambda_i = \lambda_1$$

For observations $X = (\underline{x}_1 \ \underline{x}_2 \ \dots \ \underline{x}_n)$:

$(z_1)_i = \underline{w}_1^T \underline{x}_i$ also called PC (scores)

All n first PC scores: $\underline{z}_1^T = \underline{w}_1^T X$

Subsequent PC directions

Maximise variance of $z_2 = \underline{w}_2^T \underline{x}$

but also reveal something new:

Require \underline{w}_2 orthogonal to \underline{w}_1

→ Optimisation problem:

$$\text{maximise } \underline{w}_2^T \Sigma \underline{w}_2$$

$$\text{s.t. } \|\underline{w}_2\| = 1 \quad \text{and} \quad \underline{w}_2^T \underline{w}_1 = 0$$

Like before, using the EV decomposition
express $\underline{w}_2 = U \underline{b}$

$$\Rightarrow \underline{w}_2^T \underline{w}_1 = \underline{w}_2^T \underline{u}_1 = \underline{b}^T U^T \underline{u}_1 \stackrel{\text{orth.}}{=} \underline{b}^T \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} = b_1$$

\rightarrow Insert constraint $b_1 = 0$:

$$\text{maximise } \sum_{i=2}^d b_i^2 \lambda_i \quad \text{s.t.} \quad \sum_{i=2}^d b_i^2 = 1$$

b_2, \dots, b_d

Like before, solution $\underline{b} = (0 \mid 0 \dots 0)^T$

$$\underline{w}_2 = U \underline{b} = \underline{u}_2$$

$$\text{Var}[z_2] = \lambda_2$$

Also: z_1 and z_2 are uncorrelated:

$$E[z_1, z_2] = 0$$

Further PC directions:

$$\text{maximise } \underline{w}_m^T \Sigma \underline{w}_m$$

\underline{w}_m

$$\text{s.t. } \|\underline{w}_m\| = 1 \quad \text{and} \quad \underline{w}_m^T \underline{w}_i = 0, \quad i = 1, \dots, m-1$$

$$\rightarrow \underline{w}_m = \underline{u}_m, \quad \text{Var}[z_m] = \lambda_m$$

Total variance of z_1, \dots, z_k :

(i.e. variance "explained" by k PCs):

$$\sum_{m=1}^k \text{Var}[z_i] = \sum_{m=1}^k \lambda_m$$

Corresponding normalised measure:

$$\text{Fraction of variance explained} = \frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^d \lambda_i}$$