

# Numerical data description

## Part of Exploratory Data Analysis

Location: "where"



- Sample mean
- Median
- Mode

Scale: "spread"



- Sample variance / standard deviation
- Median absolute deviation (MAD)
- Interquartile range (IQR)

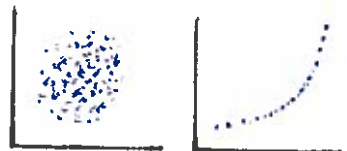
Shape: "form"

- Sample skewness
- Galton's measure of skewness
- Sample kurtosis
- Robust kurtosis



Dependence: "relation"

- Sample covariance / correlation
- Kendall's  $\tau$



## Multivariate measures

Now:  $d$ -dim. random vector  $\underline{X} = \begin{pmatrix} X_1 \\ \vdots \\ X_d \end{pmatrix}$

$n$  observations

→ Observation matrix  $X = \begin{pmatrix} X_{11} & \dots & X_{1n} \\ X_{21} & \dots & \vdots \\ \vdots & \ddots & \vdots \\ X_{d1} & \dots & X_{dn} \end{pmatrix}$

$X_{ij} = j$ -th observation of  $i$ -th attribute

Attention: here we have  $d \times n$  observation matrices  $D$

In MLPR it was the other way

Both are common in literature

### Notation:

$X_{i:}$ :  $i$ -th row vector of  $X$

$$X_{i:} = (X_{i1}, \dots, X_{in})$$

$X_{:j}$ :  $j$ -th column vector of  $X$

$$X_{:j} = \underline{X}_j = \begin{pmatrix} X_{1j} \\ X_{2j} \\ \vdots \\ X_{dj} \end{pmatrix}$$

- Sample covariance matrix:

$$\text{cov}(X) = \begin{pmatrix} \text{cov}(x_1, x_1) & \text{cov}(x_1, x_2) & \dots & \text{cov}(x_1, x_d) \\ \vdots & \ddots & & \vdots \\ \text{cov}(x_d, x_1) & \dots & & \text{cov}(x_d, x_d) \end{pmatrix}$$

Estimate of covariance matrix

$$\text{Cov}[\underline{x}] = \mathbb{E}[(\underline{x} - \mathbb{E}[\underline{x}])(\underline{x} - \mathbb{E}[\underline{x}])^T]$$

Symmetric with variances on the diagonal

Positive semi-definite:

$$\begin{aligned} \underline{w}^T \text{Cov}[\underline{x}] \underline{w} &= \mathbb{E}[\underbrace{\underline{w}^T (\underline{x} - \mathbb{E}[\underline{x}])}_{\text{scalar}} \underbrace{(\underline{x} - \mathbb{E}[\underline{x}])^T \underline{w}}_{\text{scalar}}] \\ &= \mathbb{E}[(\underline{w}^T (\underline{x} - \mathbb{E}[\underline{x}]))^2] \geq 0 \end{aligned}$$

Eigenvalue decomposition:

$$\text{Cov}[\underline{x}] = U \Lambda U^T$$

where

$\Lambda$  diagonal matrix with eigenvalues  $\lambda_1, \dots, \lambda_d$

$U$  orthonormal matrix with eigenvectors

Total variance of  $d$  attributes  
= sum of eigenvalues

$$\begin{aligned}\sum_{i=1}^d \text{Var}[x_i] &= \text{trace}(\text{Cov}[\underline{x}]) \\ &= \text{trace}(\underline{U} \underline{\Lambda} \underline{U}^T) \\ &= \text{trace}(\underline{\Lambda} \underline{U}^T \underline{U}) \\ &= \text{trace}(\underline{\Lambda}) = \sum_{i=1}^d \lambda_i\end{aligned}$$

Covariance of linear transformation:

$$\text{Cov}[\underline{A}\underline{x} + \underline{b}] = \underline{A} \text{Cov}[\underline{x}] \underline{A}^T$$

# Data preprocessing

Prepare data for further analysis

- Standardisation: make variables comparable

Often: normalise to have mean 0 and variance 1

linear transformation

→ express using matrix operations

Given:  $X = (x_{:1} \dots x_{:n}) = (\underline{x}_1 \dots \underline{x}_n)$

Sample mean: d-dim. vector  $\bar{X} = (\bar{x}_1, \dots, \bar{x}_d)$   
with  $\bar{x}_i = \frac{1}{n} \sum_{j=1}^n x_{ij}$

Centred observations:  $\tilde{x}_i = x_i - \bar{X}$

→ matrix  $\tilde{X} = (\tilde{x}_1 \dots \tilde{x}_n)$

$$= (\underline{x}_1 - \bar{X} \dots \underline{x}_n - \bar{X})$$

$$= X - \underbrace{(\bar{X} \dots \bar{X})}_{n\text{-times}}$$

Centering matrix :  $C_n = I_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T$   
where  $I_n$  = identity matrix ( $n \times n$ )

and  $\mathbf{1}_n = (1 \ 1 \dots 1)^T$  ( $n \times 1$ )

It holds:  $\tilde{X} = X C_n$

Proof:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{n} X \mathbf{1}_n$$

$$\bar{x} \mathbf{1}_n^T = \begin{pmatrix} \bar{x}_1 \\ \bar{x}_2 \\ \vdots \\ \bar{x}_d \end{pmatrix} (1 \ 1 \dots 1) = \begin{pmatrix} \bar{x}_1 & \dots & \bar{x}_1 \\ \bar{x}_2 & & \vdots \\ \vdots & & \vdots \\ \bar{x}_d & \dots & \bar{x}_d \end{pmatrix} = \underbrace{\begin{pmatrix} \bar{x} & \dots & \bar{x} \end{pmatrix}}_{n\text{-times}}$$

$$\Rightarrow (\bar{x} \dots \bar{x}) = X \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T$$

$$\begin{aligned} \tilde{X} &= X - (\bar{x} \dots \bar{x}) = X - X \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T \\ &= X (I_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T) = X C_n \quad \square \end{aligned}$$

Properties:

- Idempotence:  $C_n C_n = C_n$  "removing mean twice"

- Multiplication from right removes row mean

- Multiplication from left removes column mean

$$C_n \underline{v} = \begin{pmatrix} v_1 - \bar{v} \\ v_2 - \bar{v} \\ \vdots \\ v_n - \bar{v} \end{pmatrix}$$

Expression of sample covariance matrix

$$\begin{aligned}\text{cov}(X) &= \frac{1}{n} \sum_{i=1}^n (\underline{x}_i - \underline{\bar{x}})(\underline{x}_i - \underline{\bar{x}})^T = \frac{1}{n} \sum_{i=1}^n \tilde{\underline{x}}_i \tilde{\underline{x}}_i^T \\ &= \frac{1}{n} \tilde{X} \tilde{X}^T = \frac{1}{n} X C_n X^T\end{aligned}$$

Scaling to unit variance:

$$\text{Recall: } \text{Cov}[A\underline{x} + \underline{b}] = A \text{Cov}[\underline{x}] A^T$$

$$\text{Set } A = \text{diag}\left(\frac{1}{\text{std}(X)}\right)$$

where  $\text{diag}(\underline{v}) =$  diagonal matrix  
with  $\underline{v}$  on diagonal

$$\text{and } \text{std}(X) = (\text{std}(x_1), \dots, \text{std}(x_d))^T$$

$\rightarrow \text{cov}(AX)$  has ones on diagonal  
(i.e., unit variance)

Note: correlations are still present

# Outlier detection and removal

Outlier: observation that is unusual compared to others

Could occur due to:

- measurement error  $\rightarrow$  remove
- wrong assumptions  $\rightarrow$  keep

Indication: e.g. mean  $\gg$  median

Often used: Tukey's fences

$$[Q_1 - k \cdot IQR, Q_3 + k \cdot IQR]$$

where  $k \geq 0$  often  $k = 1.5$

Observations outside often labelled as outliers

