

NOTES OF  
APPLIED STATISTICS

From Prof. Piercesare Secchi's lectures  
for the MSc in Mathematical Engineering

by Teo Bucci

Politecnico di Milano  
A.Y. 2021/2022



## LECTURE 2 22/02/2022

units: individuals, times, places

Feature space

$$\underline{x} = (x_1, x_2, \dots, x_p)^T \in \mathbb{R}^p$$

$p$ : number of features  
 $m$ : sample size

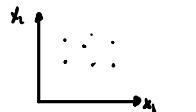
1	$x_{11}$	$x_{12}$	...	$x_{1p}$
	$x_{21}$	$x_{22}$	...	$x_{2p}$
	$\vdots$	$\vdots$		$\vdots$
$m$	$x_{m1}$	$x_{m2}$	...	$x_{mp}$

$\times$  DATA MATRIX DATA FRAME  $m \times p$

### ↳ Row Perspective

$$\mathbf{X} = \begin{bmatrix} \underline{x}_1 \\ \underline{x}_2 \\ \vdots \\ \underline{x}_m \end{bmatrix} \quad \underline{x}_i \in \mathbb{R}^p$$

profile of unit  $i$



### ↳ Column Perspective

$$\mathbf{X} = [x_1 \dots x_p] \quad x_i \in \mathbb{R}^m$$

Sample of feature  $i$



$Y$ : special variable

Goal: explain the variability of  $Y$  through the variability of  $\underline{x}$

↳ SUPERVISED PROBLEMS both the features  $\underline{x} = (\dots)$  and the target  $Y$  are obs. statistical unit

GENERAL PROBLEM OF STAT. LEARNING: supervised setting  $? f: \mathbb{R}^p \rightarrow \mathbb{R}$  s.t.  $f(\underline{x})$  is the "best" explanation of  $Y$  in terms of  $\underline{x}$

argmin <sub>$f$</sub>   $\mathbb{E}[(Y-f(\underline{x}))^2]$

PROP argmin <sub>$f$</sub>   $\mathbb{E}[(Y-f(\underline{x}))^2] = \mathbb{E}[Y|\underline{x}]$

PROOF  $\mathbb{E}[(Y-f(\underline{x}))^2] = \mathbb{E}[(Y-\mathbb{E}[Y|\underline{x}]) - (f(\underline{x}) - \mathbb{E}[Y|\underline{x}])]^2 =$   
 $= \mathbb{E}[(Y-\mathbb{E}[Y|\underline{x}])^2] + \mathbb{E}[(f(\underline{x}) - \mathbb{E}[Y|\underline{x}])^2] - 2\mathbb{E}[(Y-\mathbb{E}[Y|\underline{x}])(f(\underline{x}) - \mathbb{E}[Y|\underline{x}])]$

$$\mathbb{E}[\mathbb{E}[W|z]] = \mathbb{E}[W] \quad \begin{matrix} \text{does not} \\ \text{depend on } z \end{matrix} \quad \begin{matrix} \text{take } f = \mathbb{E}[Y|\underline{x}] \end{matrix}$$

$$(*) = \mathbb{E}\left[\underbrace{\mathbb{E}[(Y-\mathbb{E}[Y|\underline{x}])(f(\underline{x}) - \mathbb{E}[Y|\underline{x}])]}_{=0} | \underline{x}\right] = 0$$

$$= \mathbb{E}[Y|\underline{x}] - \mathbb{E}[Y|\underline{x}] = 0$$

■

SIGMA-FIELD generated by  $\underline{x}$ : when you know  $\underline{x}$ , you know everything that depends on  $\underline{x}$  (SIGMA-SPACE)

$$Y = \mathbb{E}[Y|\underline{x}] + \varepsilon$$

unknown based on the info carried by  $\underline{x}$

### ↳ BASIC MODEL FOR SL

$$Y = f(\underline{x}) + \varepsilon \quad \text{Var}(\varepsilon) = \mathbb{E}[\varepsilon^2] - \mathbb{E}[\varepsilon]^2 = \mathbb{E}[\varepsilon^2]$$

$$f(\underline{x}) = \mathbb{E}[Y|\underline{x}]$$

$\varepsilon$  r.v. s.t.  $\mathbb{E}[\varepsilon] = 0$

$\varepsilon$  indep of  $\underline{x}$ , something I can't explain in terms of  $\underline{x}$

$\underline{x} \xrightarrow{\hat{f}}$  estimate of  $f$

↳ known once  $\underline{x}$  are known

## LECTURE 3 24/02/2022

### SUPERVISED LEARNING

Basic model

$$\mathbb{R} \ni Y = f(\underline{x}) + \varepsilon \quad \underline{x} \in \mathbb{R}^p, \varepsilon \in \mathbb{R}, \mathbb{E}[\varepsilon] = 0$$

$$f: \mathbb{R}^p \rightarrow \mathbb{R}$$

$$\varepsilon \perp \underline{x} \quad \mathbb{E}[Y|\underline{x}] = f(\underline{x})$$

GOAL Use data to learn  $f$ : generate an estimate  $\hat{f}$  of  $f$

totally known when  $\underline{x}$  is known

$$\hat{f}: \underline{x} \rightarrow \hat{f}$$

Q: Properties of  $\hat{f}$ ?

Q: Variability of  $\hat{f}$ ?



"Not so long"  
 $\hat{f} \in \{f_\theta : \theta \in \mathbb{R}^q\}$   
 We approximate  $f$  with  
 a member from this  
 family close enough ( $f \approx f_\theta$ )

$\Rightarrow$  Use  $\mathbf{x}$  to estimate  $\hat{\theta}$  so that  $\hat{f}(\mathbf{x}) = f_\theta(\mathbf{x})$

Ex

$$f_\theta(\mathbf{x}) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$

$$\theta = (\beta_0, \dots, \beta_p) \in \mathbb{R}^{p+1}$$

$$\mathbf{x} \rightarrow \hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$$

LINEAR

Then variables are NOT necessarily the ones received with the dataset

Received:  $w_1, \dots, w_k$

Variables:  $x_1 = h_1(w_1, \dots, w_k)$   
 $\vdots$   
 $x_k = h_k(w_1, \dots, w_k)$

$$\Rightarrow f_\theta(w_1, \dots, w_k) = \beta_0 + \beta_1 h_1(w_1, \dots, w_k) + \dots + \beta_p h_p(w_1, \dots, w_k)$$

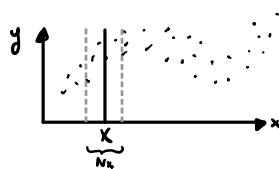
Ex

$$Y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 \log x$$

still LINEAR as long as  $\beta_1, \dots, \beta_k$  are known

↳ this is where PRIOR ENGINEERING / DOMAIN knowledge kicks in! 😊

Ex (KNN - K Nearest Neighbors)



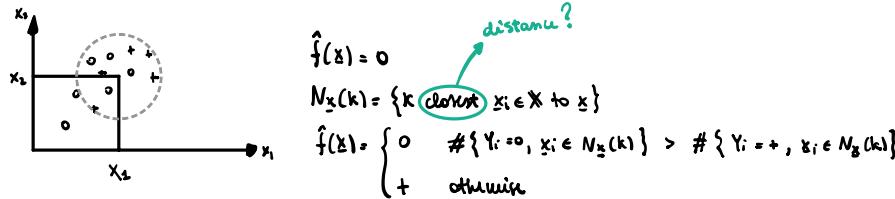
$$\hat{f}(\mathbf{x}) = ?$$

$$N_{\mathbf{x}}(k) = \{k \text{ closest } \mathbf{x}_i \in \mathbf{X} \text{ to } \mathbf{x}\}$$

$$\hat{f}(\mathbf{x}) = \frac{1}{k} \sum_{\mathbf{x}_i \in N_{\mathbf{x}}(k)} Y_i$$

$$\text{Var}(\hat{f}(\mathbf{x})) = \frac{1}{k} \text{Var}(E)$$

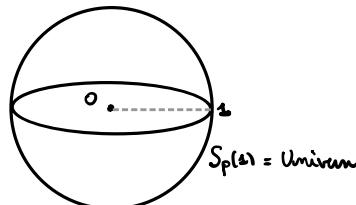
Ex (KNN for classification)



### Cure of dimensionality

↳ There is a "continuity" idea,  $f$  needs to be regular enough

↳ If  $p$  is large: nothing's around



$S_p(r) =$  sphere of radius  $r$  in  $\mathbb{R}^p$  centered in  $0$

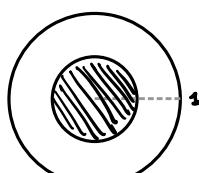
Friends are uniformly distributed in  $S_p(r)$

Q: How far do you have to go to reach 10% of your friends?

$$p=1 \quad \text{---} \quad -r \quad 0 \quad r \quad \text{---} \quad 10\%$$

$$\frac{\mu(S_1(r))}{\mu(S_1(z))} = 0,1 = \frac{2r}{2} = r \Rightarrow r = 0,1$$

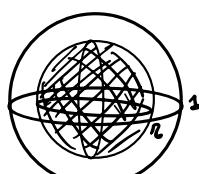
$p=2$



$$0,1 = \frac{\mu(S_2(r))}{\mu(S_2(z))} = \frac{\pi r^2}{\pi z^2} = r^2 \Rightarrow r = \sqrt{0,1} = 0,31 \quad \text{You have to travel 30% MORE!}$$

1

$p=3$



$$0,1 = \frac{\mu(S_3(r))}{\mu(S_3(z))} = \frac{\frac{4}{3}\pi r^3}{\frac{4}{3}\pi z^3} = r^3 \Rightarrow r = \sqrt[3]{0,1} = 0,46$$

$$p = 100 \Rightarrow \pi = (0.1)^{1/100} = 0.97$$

↳  $p=1 m=100$  Same information  
 ↳  $p=50 m=(100)^{50}$  but in the real universe there are  $\approx 10^{82}$  atoms 😊  
 ↳ Too much!

How to tackle the overfit?

- ↳ Use structured models
- ↳ reduce dimensionality, PCA, ICA

How to evaluate  $\hat{f}$ ?

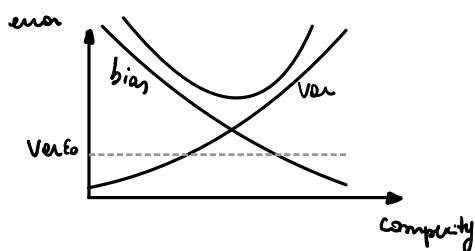
Error when predicting  $Y$  for a new unit  $(Y_0, \underline{x}_0)$

test error

$\begin{aligned} & \text{Random} \\ & \mathbb{E}_{\underline{x}}[(Y_0 - \hat{f}(\underline{x}_0))^2] \\ & \quad \text{expectation w.r.t. } \underline{x} \text{ is known} \\ & \quad \underline{x} \text{ is known} \end{aligned}$	$\begin{aligned} & Y_0 = f(\underline{x}_0) + \epsilon_0 \\ & \quad \epsilon_0 \text{ is known} \end{aligned}$
--	--

$$\begin{aligned} & = \mathbb{E}_{\underline{x}}[(f(\underline{x}_0) - \hat{f}(\underline{x}_0) + \epsilon_0)^2] \quad \mathbb{E}[\epsilon_0] = 0 \\ & = \mathbb{E}_{\underline{x}}[(f(\underline{x}_0) - \hat{f}(\underline{x}_0))^2 + \epsilon_0^2 + 2\epsilon_0(f(\underline{x}_0) - \hat{f}(\underline{x}_0))] \\ & = (f(\underline{x}_0) - \hat{f}(\underline{x}_0))^2 + \text{Var}(\epsilon_0) \end{aligned}$$

this is the part we can reduce      No control on this irreducible



Ex KNN

$$\hat{f}(\underline{x}) = \frac{1}{k} \sum_{i \in N_k(\underline{x})} y_i \quad \text{Var}(\hat{f}(\underline{x})) = \frac{1}{k} \text{Var}(\epsilon) \downarrow k \uparrow$$

but bias  $\uparrow k \uparrow$

Read Chap 2 LSS (Intro to Stat Learning)

## LECTURE 4 1/3/2022

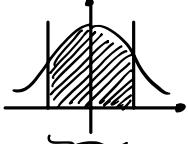
$$S_{11} = \frac{1}{m} \sum_{i=0}^m (\underline{x}_{i1} - \bar{\underline{x}}_1)(\underline{x}_{i1} - \bar{\underline{x}}_1) \quad \text{where } \bar{\underline{x}} = \begin{bmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_p \end{bmatrix} \text{ vector of sample means}$$

$S = [s_{11} \dots s_{pp}]$   $p \times p$  symmetric matrix of real numbers (sample covariance matrix)

↳ Unit measure is a problem  $\Rightarrow r_{jk} = \frac{s_{jk}}{\sqrt{s_{jj} \cdot s_{kk}}} \in [-1, 1]$  Standardize!

$$R = \begin{bmatrix} 1 & r_{12} & \dots & r_{1p} \\ \ddots & \ddots & \ddots & \ddots \end{bmatrix} \quad \text{correlation matrix}$$

Ex Height  $\underline{x}_1 = 170$  cm,  $\sqrt{s_{11}} = 3$  cm



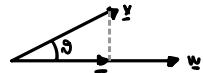
$$\text{Chubeychow: } \mathbb{P}[x \in [\bar{x}_i - k\sqrt{s_{ii}}, \bar{x}_i + k\sqrt{s_{ii}}]] \geq 1 - \frac{k^2}{k^2} \quad k=2$$

$\pm 2$ std dev	95%	Chubeychow
$\pm 2$ std dev	99%	Gaussian

## "Geometry" in $\mathbb{R}^m$

$$x, w \in \mathbb{R}^m, \langle x, w \rangle = x'w, \|x\| = \sqrt{\langle x, x \rangle} = \sqrt{\sum_{i=1}^m x_i^2}$$

$$\cos \theta = \frac{\langle x, w \rangle}{\|x\| \|w\|} = \frac{\sum x_i w_i}{\sqrt{\sum x_i^2} \sqrt{\sum w_i^2}}$$



$$\pi_{x \perp w} = \|x\| \cos \theta \cdot \frac{w}{\|w\|} = \|x\| \frac{x'w}{\|x\| \|w\|} \cdot \frac{w}{\|w\|} = \frac{x'w}{w'w} w = \left( \frac{w w'}{w'w} \right) x$$

operator that projects

### Data

$$\mathcal{L}(\underline{z}) = \{(\underline{c}, \dots, \underline{c}): c \in \mathbb{R}\}$$

$$\underline{z} = \begin{pmatrix} 1 \\ z_1 \\ \vdots \\ z_m \end{pmatrix}$$

No variability in this space

They all have the same MEAN but the approx is BAD for many of them

$$\begin{aligned} \underline{e}_1 &= \begin{pmatrix} x_{11} \\ x_{1m} \end{pmatrix} \\ \underline{d}_1 &= \underline{e}_1 - \bar{x}_1 \cdot \underline{1} \quad \text{MINIMIZE THIS!} \\ \pi_{\underline{e}_1 \perp \underline{z}} &= \frac{1 - \bar{x}_1}{m} \underline{e}_1 = \frac{1}{m} \sum_{i=1}^m (x_{i1} - \bar{x}_1) \underline{1} = \begin{pmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_1 \end{pmatrix} \end{aligned}$$

$$\|\underline{d}_1\| = \sqrt{\sum d_i^2} = \sqrt{\sum (x_{i1} - \bar{x}_1)^2}$$

$$\Rightarrow \|\underline{d}_1\|^2 = m S_m$$

$$\text{Commonly } \underline{d}_1 = \underline{e}_1 - \bar{x}_1 \cdot \underline{1} \in \mathcal{L}^\perp(\underline{z})$$

$$\begin{aligned} \theta = 0 &\Rightarrow \underline{d}_1 \in \mathcal{L}(\underline{d}_2) \quad \underline{d}_1 = k \underline{d}_2 \quad \text{same information carried} \\ \theta = \frac{\pi}{2} &\Rightarrow \text{nothing in } \underline{d}_2 \text{ carries info. about } \underline{d}_1 \text{ and vice versa} \end{aligned}$$

$$\underline{d}_2 = \underline{e}_2 - \bar{x}_2 \cdot \underline{1} \in \mathcal{L}^\perp(\underline{z})$$

$$\begin{array}{ccc} \text{information of } \underline{d}_2 \text{ carried by } \underline{d}_2 & & \text{information of } \underline{d}_2 \text{ carried by } \underline{d}_1 \\ \text{information of } \underline{d}_1 \text{ carried by } \underline{d}_2 & & \end{array}$$

$$\cos \theta = \frac{\underline{d}_1 \cdot \underline{d}_2}{\|\underline{d}_1\| \|\underline{d}_2\|} = \frac{\frac{1}{m} \sum (x_{i1} - \bar{x}_1)(x_{i2} - \bar{x}_2)}{\sqrt{\sum (x_{i1} - \bar{x}_1)^2} \sqrt{\sum (x_{i2} - \bar{x}_2)^2}} = r_{12} \in [-1, 1] \quad \text{OK}$$

### Data

$$\underline{X} = \begin{bmatrix} \underline{x}_1 \\ \vdots \\ \underline{x}_m \end{bmatrix} \quad \underline{x}_i \in \mathbb{R}^p$$

$x_1, \dots, x_m$  iid realizations of  $\underline{x}$  random vector of  $\mathbb{R}^p$

$$\mathbb{P}[\underline{x} \in B] = \nu_{\underline{x}}(B) \quad \forall B \in \mathcal{B}(\mathbb{R}^p)$$

↳ Borel

$$\text{It might happen } \nu_{\underline{x}}(B) = \int_B f(\underline{x}) d\underline{x} \quad \text{↳ density}$$

Ex

$$f(\underline{x}) = \frac{1}{\sqrt{(2\pi)^p |\Lambda \Sigma|}} \exp \left\{ -\frac{1}{2} (\underline{x} - \mu)' \Sigma^{-1} (\underline{x} - \mu) \right\} \quad \sum \in \mathbb{R}^{p \times p}, \mu \in \mathbb{R}^p \quad \mathbb{E}[\underline{x}] = \mu \quad \mathbb{E}[(\underline{x} - \mu)(\underline{x} - \mu)'] = \Sigma = [v_{ij}] \quad i, j = 1, \dots, p \quad v_{ij} = \text{Cov}(x_i, x_j)$$

$$\hookrightarrow c \in \mathbb{R}^p \quad \mathbb{E}[c' \underline{x}] = \mathbb{E}[\sum c_i x_i] = \sum c_i \mathbb{E}[x_i] = c' \mu$$

$$(R \exists) \text{Var}[c' \underline{x}] = c' \Sigma c \quad (\mathbb{R} \mathbb{R})$$

$$\text{Var}(c_1 x_1 + c_2 x_2) = c_1^2 \text{Var} x_1 + c_2^2 \text{Var} x_2 + 2c_1 c_2 \text{Cov}(x_1, x_2)$$

$$\hookrightarrow C \in \mathbb{R}^{p,p} \quad \mathbb{E}[C \underline{x}] = C \mu$$

$$\text{Var}(C \underline{x}) = C \Sigma C'$$

### Methods of moments

$$\hat{\mu} = \frac{1}{m} \sum_{i=1}^m x_i = \bar{x} \quad \hat{\Sigma} = \frac{1}{m} \sum_{i=1}^m (x_i - \bar{x})(x_i - \bar{x})' = S_m$$

- What are the properties of these two machines?

$$\textcircled{1} \quad \mathbb{E}[\bar{x}] = \mu \quad (\text{unbiased})$$

$$\text{Cov}(\bar{x}) = \frac{1}{m} \Sigma$$

$$\textcircled{2} \quad \mathbb{E}[S_m] = \frac{m-1}{m} \sum \quad \text{not unbiased} \rightsquigarrow S = \frac{1}{m-2} \sum_{i=1}^m (x_i - \bar{x})(x_i - \bar{x})' \quad \text{is unbiased}$$

## LECTURE 5 3/3/2022

$$\underline{X} = \begin{bmatrix} \underline{x}_1 \\ \vdots \\ \underline{x}_m \end{bmatrix} = [\underline{e}_1 \dots \underline{e}_p] \quad \underline{x}_i \in \mathbb{R}^p, \underline{e}_i \in \mathbb{R}^m$$

↳ realizations of  $x_1, \dots, x_m$  iid  $\sim \mu, \Sigma$

→ Deviation vector  $\underline{d}_i$

$$\begin{array}{ccc} \underline{e}_i & & \\ \underline{d}_i \in \mathcal{L}^\perp(\underline{z}) & & \\ \underline{x}_i - \bar{x} & & \underline{d}(\underline{z}) \end{array}$$

estimated by

$$\bar{x} = \frac{1}{m} \sum_{i=1}^m \underline{x}_i \text{ for } \mu$$

$$S_m = \frac{1}{m-2} \sum_{i=1}^m (\underline{x}_i - \bar{x})(\underline{x}_i - \bar{x})' \text{ for } \Sigma$$

Biased  $\rightsquigarrow \frac{1}{m-1}$  unbiased

$$\begin{aligned} \underline{d}_i &= \underline{e}_i - \bar{x}_i \cdot \underline{1} \\ &= \underline{e}_i - \frac{\underline{1} \cdot \underline{e}_i}{\underline{1} \cdot \underline{1}} \underline{e}_i = \left( I - \frac{\underline{1} \cdot \underline{1}'}{\underline{1}' \underline{1}} \right) \underline{e}_i \end{aligned}$$

Orthogonal projection on  $\mathcal{L}^{\perp}(1)$

$$\begin{aligned} \mathbf{d} &= [d_1 \dots d_p] \text{ matrix of obs. vectors} \\ &= \left( I - \frac{\underline{1} \cdot \underline{1}'}{\underline{1}' \underline{1}} \right) \mathbf{X} \end{aligned}$$

$$S = \frac{1}{m-1} \begin{bmatrix} \sum (x_{i1} - \bar{x}_1)^2 & \sum (x_{i1} - \bar{x}_1)(x_{i2} - \bar{x}_2) \dots \\ \vdots & \ddots \end{bmatrix} = \frac{1}{m-1} \begin{bmatrix} d_1' d_1 & d_1' d_2 \dots & d_1' d_p \\ \vdots & \ddots & \vdots \\ d_p' d_1 & \dots & d_p' d_p \end{bmatrix} = \frac{1}{m-1} \mathbf{d}' \mathbf{d} = \underbrace{\frac{1}{m-1} \mathbf{X}' \left( I - \frac{\underline{1} \cdot \underline{1}'}{\underline{1}' \underline{1}} \right)' \left( I - \frac{\underline{1} \cdot \underline{1}'}{\underline{1}' \underline{1}} \right) \mathbf{X}}_{\mathbf{d}' \mathbf{d}}$$

$$\Rightarrow S = \frac{1}{m-1} \mathbf{X}' \left( I - \frac{\underline{1} \cdot \underline{1}'}{\underline{1}' \underline{1}} \right) \mathbf{X}$$

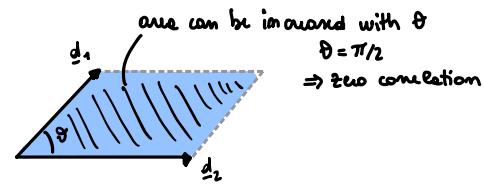
$\hookrightarrow$  Generalized variance :  $\det(S)$  (sample)  
 $\det(\Sigma)$  (population)

$\hookrightarrow$  Total variance :  $\text{tr}(S)$  (sample)  
 $\text{tr}(\Sigma)$  (population)

$\left[ \begin{smallmatrix} 1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \dots & 1 \end{smallmatrix} \right] \frac{\det \Sigma}{\text{tr } \Sigma} \right\}_{\text{REC}}$

Ex

$$p=2 \quad S = \frac{1}{m-1} \cdot \begin{bmatrix} d_1' d_1 & d_1' d_2 \\ d_2' d_1 & d_2' d_2 \end{bmatrix} = \frac{1}{m-1} \begin{bmatrix} \|d_1\|^2 & \|d_1\| \|d_2\| \cos \theta \\ \|d_1\| \|d_2\| \cos \theta & \|d_2\|^2 \end{bmatrix}$$



$$\text{Gen. Var.}(S) = \det(S) = \frac{1}{(m-1)^2} [\|d_1\|^2 \|d_2\|^2 - \|d_1\|^2 \|d_2\|^2 \cos^2 \theta] = \frac{1}{(m-1)^2} [\|d_1\|^2 \|d_2\|^2 \sin^2 \theta] \propto \text{Area}^2(\text{parallelogram}(d_1, d_2))$$

Prop  $\det(S) = 0 \iff d_1 \dots d_p \text{ are linearly dependent}$

Proof ( $\Leftarrow$ )  $\exists \underline{c} \neq \underline{0}, \underline{c} \in \mathbb{R}^p \quad \underline{d}' \underline{c} = \underline{0}$   
 $\sum \underline{d}_i c_i = \underline{0}$

$$S = \frac{1}{m-1} \mathbf{d}' \mathbf{d} \quad S \underline{c} = \frac{1}{m-1} \underbrace{\mathbf{d}' \mathbf{d} \underline{c}}_{=0} = \underline{0} \quad \Rightarrow \text{cols of } S \text{ are lin. dep.} \Rightarrow \det(S) = 0$$

$$(\Rightarrow) \quad \exists \underline{c} \in \mathbb{R}^p, \underline{c} \neq \underline{0} \quad \text{s.t.} \quad S \underline{c} = \underline{0} \quad \Rightarrow \quad \frac{1}{m-1} \mathbf{d}' \mathbf{d} \underline{c} = \underline{0} \quad \Rightarrow \quad \underbrace{\frac{1}{m-1} \underline{c}' \mathbf{d}' \mathbf{d} \underline{c}}_{\|\mathbf{d}\|^2 \underline{c}^2} = 0 \quad \Rightarrow \quad \|\mathbf{d} \underline{c}\|^2 = 0 \quad \Rightarrow \quad d_1 \dots d_p \text{ lin. dep.}$$

$$\text{Hence if } \det(S) = 0 \quad d \underline{c} = \underline{0} \quad \underline{0} \neq \underline{c} = \begin{pmatrix} c_1 \\ c_2 \\ \vdots \\ c_p \end{pmatrix} \quad \text{Wlog } c_1 \neq 0 \quad \sum_{i=1}^p c_i d_i = \underline{0} \quad \Rightarrow \quad d_1 = - \sum_{i=2}^p \frac{c_i}{c_1} d_i$$

$$\Rightarrow x_1 - \bar{x}_1 = - \sum_{i=2}^p \frac{c_i}{c_1} (x_i - \bar{x}_i) \quad \blacksquare$$

$\det(S) = 0$  never happens  $\Rightarrow \det(S) \approx 0$  is enough  
 $\hookrightarrow$  collinearity between variables (trouble!)

$\mathbf{X}$ :  $m \times p$  matrix

Corollary If  $p \geq m \Rightarrow \det(S) = 0$   
more features than observations

Proof  $d_1 \dots d_p \in \mathcal{L}^{\perp}(1) \iff$  at most  $m-1$  systems of linear indep. vectors  
 $\dim(\mathcal{L}^{\perp}(1)) \leq m-1$

If  $p \geq m > m-1 \Rightarrow d_1 \dots d_p$  are lin. dep.



Prop  $S'$  is positive semi-def

If  $\text{Det}(S) \neq 0 \Rightarrow S$  is positive def

proof We need to prove:  $\forall \xi \in \mathbb{R}^p \quad \xi' S \xi \geq 0$

$$\text{But } \frac{1}{m-1} \xi' d' d \xi = \frac{1}{m-1} \|d\xi\|^2 \geq 0$$

equal

If  $S$  is not pos. def.  $\Rightarrow \exists \xi \neq 0$  s.t.  $\xi' S \xi = 0 \Rightarrow \frac{1}{m-1} \xi' d' d \xi = 0 \Rightarrow \|d\xi\|^2 = d\xi = 0 \Rightarrow d \xi = 0 \Rightarrow d_1 \dots d_p \text{ lin. dep} \Rightarrow \text{Det}(S) = 0$

REMARK

$S'$  symmetric, real-valued  $\Rightarrow \exists (\lambda_i, e_i) \in \mathbb{R}^p$  s.t.  $S = \sum_{i=1}^p \lambda_i e_i e_i'$   
moreover  $e_i' e_j = \begin{cases} 0 & i \neq j \\ 1 & i=j \end{cases}$

thus  $\{e_1, \dots, e_p\}$  orthonormal system for  $\mathbb{R}^p$

They satisfy  $S \xi_i = \lambda_i \xi_i$   $i = 1 \dots p$

Convention

Also, since it's positive semidef we can order them:  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$   
( $>$  if positive def.)

Notation

$$P = [e_1 \dots e_p] \Rightarrow S = P \Lambda P'$$

$$\Lambda = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_p \end{bmatrix} \quad \text{If } \text{Det}(S) \neq 0 \Rightarrow S^{-1} = \sum_{i=1}^p \frac{1}{\lambda_i} e_i e_i' = P \Lambda^{-1} P$$

$$\text{If } \text{Det}(S) = 0 \Rightarrow \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k > \underbrace{\lambda_{k+1} = 0 = \dots = \lambda_p}_{\text{ignore them}} \Rightarrow S^{-1} = \sum_{i=1}^k \frac{1}{\lambda_i} e_i e_i'$$

If  $\text{Det}(S) \neq 0$ :  $d_{S^{-1}}(x, y) = \sqrt{(x-y)' S^{-1} (x-y)}$  Mahalanobis distance

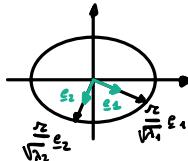
If  $S = I$  this is the Euclidean distance

$$S^{-1}(x) = \{z \in \mathbb{R}^p : (x-z)' (x-z) \leq n^2\}$$

$$S^{-1}(x) = \{z \in \mathbb{R}^p : (x-z)' S^{-1} (x-z) \leq n^2\} \rightsquigarrow \text{Ellips}$$

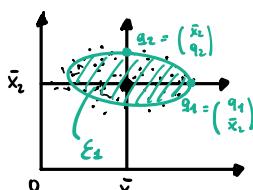
B p x p pos. def  $x \in \mathbb{R}^p$ :  $x' B x = 1$

$$B = \sum \lambda_i e_i e_i'$$



Ex

$$p=2 \quad S = \begin{bmatrix} S_{11} & 0 \\ 0 & S_{22} \end{bmatrix} \quad S_{11} > S_{22}$$



$$d_{\text{eucl.}}(q_1, \bar{x}) = d_{\text{eucl.}}(q_2, \bar{x})$$

$\hookrightarrow$  very little significative

$$\hookrightarrow \text{standardize: } d_{\text{eucl.}}(\text{Std}(q_1), \text{Std}(\bar{x})) = d_{\text{eucl.}}\left(\begin{pmatrix} \frac{q_1 - \bar{x}_1}{\sqrt{S_{11}}} & 0 \end{pmatrix}\right) = \frac{|q_1 - \bar{x}_1|}{\sqrt{S_{11}}}$$

$$d_{\text{eucl.}}(\text{Std}(q_2), \text{Std}(\bar{x})) = d_{\text{eucl.}}\left(\begin{pmatrix} 0 & \\ \frac{q_2 - \bar{x}_2}{\sqrt{S_{22}}} & 0 \end{pmatrix}\right) = \frac{|q_2 - \bar{x}_2|}{\sqrt{S_{22}}}$$

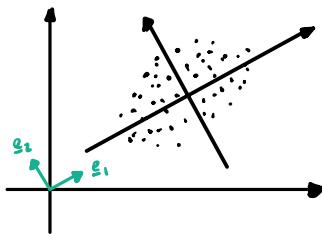
IDEA

The euclidean distance is right after standardization

$$d_{\text{eucl.}}(\text{Std}(q), \text{Std}(\bar{x})) = d_{\text{eucl.}}\left(\begin{pmatrix} \frac{q_1 - \bar{x}_1}{\sqrt{S_{11}}} & 0 \\ \frac{q_2 - \bar{x}_2}{\sqrt{S_{22}}} & 0 \end{pmatrix}\right) = \sqrt{\left(\frac{q_1 - \bar{x}_1}{\sqrt{S_{11}}}\right)^2 + \left(\frac{q_2 - \bar{x}_2}{\sqrt{S_{22}}}\right)^2}$$

$$\text{Vol}(\Sigma(S)) = k_p \sqrt{\prod_{i=1}^p \lambda_i} = \alpha \sqrt{\text{Det}(S)}$$

What if ...  $\underline{S}$  not diagonal



## LECTURE 6 8/3/2022

$S$  cov. matrix

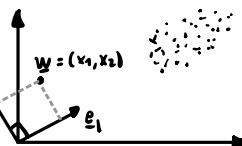
$$\det(S) > 0$$

$$S = \sum_{i=1}^p \lambda_i \underline{\varepsilon}_i \underline{\varepsilon}_i'$$

$$= P \Lambda P'$$

$$P = [\underline{\varepsilon}_1 \dots \underline{\varepsilon}_p]$$

$$\Lambda = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_p \end{bmatrix}$$



New system

$$\tilde{w} = \begin{pmatrix} \underline{\varepsilon}_1' \underline{w} \\ \underline{\varepsilon}_2' \underline{w} \\ \vdots \\ \underline{\varepsilon}_p' \underline{w} \end{pmatrix} = P' \underline{w}$$

$$\|\underline{w}\|^2 = \underline{w}' \underline{w} = \tilde{w}' \tilde{w} = \underline{w}' P' P \underline{w}$$

$$\text{The data } \tilde{\mathbf{x}} \text{ becomes: } \tilde{\mathbf{x}} = \begin{bmatrix} (P' \mathbf{x})' \\ (P' \mathbf{x}_m)' \end{bmatrix} = \mathbf{X} P$$

$$\tilde{\mathbf{x}} = \mathbf{X} P$$

$$S = P \Lambda P'$$

The covariance becomes

$$S = \frac{1}{m-1} \tilde{\mathbf{x}}' \left( I - \frac{1 \ 1'}{1 \ 1'} \right) \tilde{\mathbf{x}}$$

$$\tilde{S} = \frac{1}{m-1} \tilde{\mathbf{x}}' \left( I - \frac{1 \ 1'}{1 \ 1'} \right) \tilde{\mathbf{x}} = \frac{1}{m-1} P' \mathbf{X}' \left( I - \frac{1 \ 1'}{1 \ 1'} \right) \mathbf{X} P = P' S P = P' P \Lambda P' P = \Lambda = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_p \end{bmatrix}$$

## PCA: PRINCIPAL COMPONENT ANALYSIS

Opt. problem: does  $\exists \underline{a}^* \in \mathbb{R}^p$  s.t.  $\text{Var}(\underline{a}^* \underline{x}) = \sup_{\substack{\underline{a} \in \mathbb{R}^p \\ \|\underline{a}\|=1}} \text{Var}(\underline{a} \underline{x})$ ?

Not well posed: if  $\underline{a}^*$  is solution  $\Rightarrow \text{Var}(100 \underline{a}^* \underline{x}) = 100 \text{Var}(\underline{a}^* \underline{x}) > \text{Var}(\underline{a}^* \underline{x}) \Rightarrow \underline{a}^* \text{ not solution}$

We ask  $\|\underline{a}\| = 1$

$\hookrightarrow \underline{a}$  is supposed to be a direction

REMARK: this is not fitting the regression line when we want to predict one variable in terms of the other, here we are looking for the subspace which is closest (best approximates) the data.

prop B pos. semidef.  $p \times p$  ( $B = \sum_{i=1}^p \lambda_i \underline{\varepsilon}_i \underline{\varepsilon}_i'$ ), then

$$1) \sup_{\substack{\underline{a} \in \mathbb{R}^p \\ \underline{a}' \underline{a} = 1}} \frac{\underline{a}' B \underline{a}}{\underline{a}' \underline{a}} = \lambda_1 \quad \text{argmax } \underline{x} = \underline{\varepsilon}_1$$

$$2) \sup_{\substack{\underline{a} \in \mathbb{R}^p \\ \underline{a} \perp \underline{\varepsilon}_1}} \frac{\underline{a}' B \underline{a}}{\underline{a}' \underline{a}} = \lambda_2$$

$\vdots$

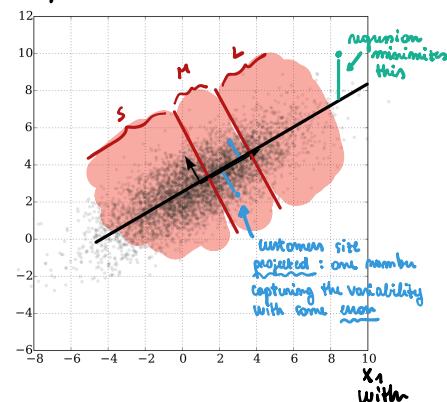
$$p) \sup_{\substack{\underline{a} \in \mathbb{R}^p \\ \underline{a} \perp \underline{\varepsilon}_1 \dots \underline{\varepsilon}_{p-1}}} \frac{\underline{a}' B \underline{a}}{\underline{a}' \underline{a}} = \lambda_p = \inf_{\substack{\underline{a} \in \mathbb{R}^p \\ \underline{a}' \underline{a} = 1}} \frac{\underline{a}' B \underline{a}}{\underline{a}' \underline{a}} \quad \text{argmax } \underline{x} = \underline{\varepsilon}_p$$

$$\text{proof} \quad \frac{\underline{a}' B \underline{a}}{\underline{a}' \underline{a}} = \frac{\underline{a}' P \Lambda P' \underline{a}}{\underline{a}' \underline{a}} = \frac{\underline{a}' \Lambda \underline{a}}{\underline{a}' \underline{a}} = \frac{\sum_{i=1}^p \lambda_i y_i^2}{\sum y_i^2} \leq \lambda_1 \frac{\sum y_i^2}{\sum y_i^2} = \lambda_1$$

but if  $\underline{a} = \underline{\varepsilon}_1$ :  $\frac{\underline{a}' B \underline{a}}{\underline{a}' \underline{a}} = \frac{\underline{a}' \Lambda \underline{a}}{\underline{a}' \underline{a}} = \lambda_1$  the sup is reached

$$\text{Step 2) } \underline{a} \perp \underline{\varepsilon}_1 \quad \frac{\underline{a}' B \underline{a}}{\underline{a}' \underline{a}} = \frac{\underline{a}' P \Lambda P' \underline{a}}{\underline{a}' \underline{a}} = \frac{\underline{a}' \Lambda \underline{a}}{\underline{a}' \underline{a}} = \frac{\sum_{i=2}^p \lambda_i y_i^2}{\sum y_i^2} \leq \lambda_2 \frac{\sum y_i^2}{\sum y_i^2} = \lambda_2$$

$x_2$  length



$$\underline{y} = \underline{P}' \underline{a} = \begin{bmatrix} \underline{e}_1' \\ \underline{e}_2' \\ \vdots \\ \underline{e}_p' \end{bmatrix} \underline{a} = \begin{bmatrix} 0 \\ y_2 \\ \vdots \\ y_p \end{bmatrix}$$

$$\text{but if } \underline{a} = \underline{e}_2 (\perp \underline{e}_2) : \frac{\underline{e}_2' \underline{B} \underline{e}_2}{\underline{e}_2' \underline{e}_2} = \frac{\lambda_2 \underline{e}_2' \underline{e}_2}{\underline{e}_2' \underline{e}_2} = \lambda_2$$

$$\text{Step 1} \quad \frac{\underline{a}' \underline{B} \underline{a}}{\underline{a}' \underline{a}} = \frac{\underline{a}' \underline{P} \Lambda \underline{P} \underline{a}}{\underline{a}' \underline{a}} = \frac{\underline{y}' \Lambda \underline{y}}{\underline{y}' \underline{y}} = \frac{\sum \lambda_i y_i^2}{\sum y_i^2} \geq \lambda_p \frac{\sum y_i^2}{\sum y_i^2} = \lambda_p \quad \blacksquare$$

## Back to PCA

$$\sup_{\substack{\underline{a} \in \mathbb{R}^p \\ \|\underline{a}\|=1}} \text{Var}(\underline{a}' \underline{x}) = \sup_{\substack{\underline{a} \in \mathbb{R}^p \\ \|\underline{a}\|=1}} \frac{\underline{a}' \Sigma \underline{a}}{\underline{a}' \underline{a}} = \lambda_1 \quad \Sigma = \underline{P} \Lambda \underline{P}' = \sum_{i=1}^p \lambda_i \underline{e}_i \underline{e}_i'$$

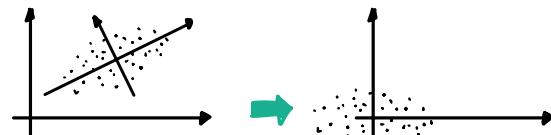
$$\text{We want } \text{Cov}(\underline{a}' \underline{x}, \underline{e}_1' \underline{x}) \geq 0 \quad \text{Cov}(\underline{a}' \underline{x}, \underline{e}_1' \underline{x}) = \underline{a}' \Sigma \underline{e}_1 = \lambda_1 \underline{a}' \underline{e}_1 = 0 \Leftrightarrow \underline{a}' \underline{e}_1 = 0 \Leftrightarrow \underline{a} \perp \underline{e}_1$$

$$\text{In general } \text{Cov}(\underline{a}' \underline{x}, \underline{b}' \underline{x}) = \underline{a}' \Sigma \underline{b}$$

$$C = \begin{bmatrix} \underline{a}' \\ \underline{b}' \end{bmatrix} \quad \text{Cov}(C \underline{x}) = C \Sigma C' = \begin{bmatrix} \underline{a}' \\ \underline{b}' \end{bmatrix} \Sigma \begin{bmatrix} \underline{a} & \underline{b} \end{bmatrix} = \begin{bmatrix} \underline{a}' \Sigma \underline{a} & \underline{a}' \Sigma \underline{b} \\ \underline{b}' \Sigma \underline{a} & \underline{b}' \Sigma \underline{b} \end{bmatrix}$$

$$\text{So } \sup_{\substack{\underline{a} \in \mathbb{R}^p \\ \|\underline{a}\|=1}} \text{Var}(\underline{a}' \underline{x}) = \lambda_1, \quad \text{In general: } \sup_{\substack{\underline{a} \in \mathbb{R}^p \\ \|\underline{a}\|=1}} \text{Var}(\underline{a}' \underline{x}) = \lambda_{k+1}$$

$$\text{Cov}(\underline{a}' \underline{x}, \underline{e}_j' \underline{x}) = 0 \quad \forall j=1 \dots k$$



DEF For  $k=1 \dots p$ ,  $\underline{e}_k$  identifies the direction of the  $k$ -th principal component

$$Y_k = \underline{e}_k' \underline{x} \quad (Y_k = \underline{e}_k (\underline{x} - \underline{\mu}))$$

$$Y_k = e_{1k} x_1 + e_{2k} x_2 + \dots + e_{pk} x_p \quad \text{or it will be } Y = \underline{P}' (\underline{x} - \underline{\mu}) \quad \text{so that}$$

$$\mathbb{E}[Y] = 0$$

$$\text{Cov}(Y) = \underline{P}' \Sigma \underline{P} = \underline{P}' \underline{\Lambda} \underline{P}' \underline{P} = \underline{\Lambda}$$

$$\text{PROP} \quad \text{Cov}(Y_k, X_i) = e_{ik} \cdot \frac{\sqrt{\lambda_k}}{\sqrt{\sigma_{ii}}} \quad \Sigma = \begin{bmatrix} \sigma_{11} & \sigma_{12} & \dots & \sigma_{1p} \\ \sigma_{21} & \sigma_{22} & \dots & \sigma_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{p1} & \sigma_{p2} & \dots & \sigma_{pp} \end{bmatrix}$$

$$\text{Proof} \quad \text{Cov}(Y_k, X_i) = \frac{\text{Cov}(Y_k, X_i)}{\sqrt{\lambda_k} \sqrt{\sigma_{ii}}} = \frac{\lambda_k e_{ik}}{\sqrt{\lambda_k} \sqrt{\sigma_{ii}}} = e_{ik} \frac{\sqrt{\lambda_k}}{\sqrt{\sigma_{ii}}}$$

$$\text{Cov}(Y_k, X_i) = \text{Cov}(\underline{e}_k' \underline{x}, \underline{e}_i' \underline{x}) = \underline{e}_k' \sum_{\text{over } i} \underline{e}_i = \underline{e}_k' \sum \underline{e}_i = \lambda_k \underline{e}_i' \underline{e}_k = \lambda_k e_{ik}$$

Is variability between variables VERY different?

$\hookrightarrow$  NO  $\Rightarrow$  PCA

$\hookrightarrow$  YES  $\Rightarrow$  Standardize  $\Rightarrow$  PCA

} may get different principal components

## LECTURE 7 10/3/2022

### REMARK

$$\text{Gen.Var}(\underline{x}) = \det \Sigma = \prod_{i=1}^p \lambda_i = \det \underline{\Lambda} = \text{Gen.Var}(Y)$$

$$\text{Total.Var}(\underline{x}) = \text{tr} \Sigma = \text{tr} \underline{\Lambda} = \text{Tot.Var}(Y)$$

### INTERPRETATION

$$\text{Cor}(Y_i, X_k) = \frac{e_{ki}}{\sqrt{\lambda_k}}$$

Interpret the loadings only if the variances are of comparable sizes

tells how much  $X_k$  is important in generating  $Y_i$   
(don't forget that also  $X_k$  plays an important role)

$\hookrightarrow$  what does it mean? "ART", depends on the problem and on how good you are ...



## PCA of std. variables

$$\underline{x} \stackrel{\text{std}}{\sim} \mathcal{N}^n(\underline{\mu}, \Sigma) = \underline{z}$$

$\Sigma = \begin{bmatrix} \sigma_{11} & 0 \\ 0 & \sigma_{pp} \end{bmatrix}$

$$\underline{z} = \left( \frac{x_1 - \mu_1}{\sqrt{\sigma_{11}}}, \dots, \frac{x_p - \mu_p}{\sqrt{\sigma_{pp}}} \right)^T$$

$$\text{Cov}(z) = \Sigma^{-1/2} \Sigma \Sigma^{-1/2} = \rho = \text{cov}(x)$$

$\Leftrightarrow \rho = \sum_{i=1}^p \lambda_i e_i e_i^T$  (NOT the same  $\lambda_i, e_i$  of  $\Sigma$ !)

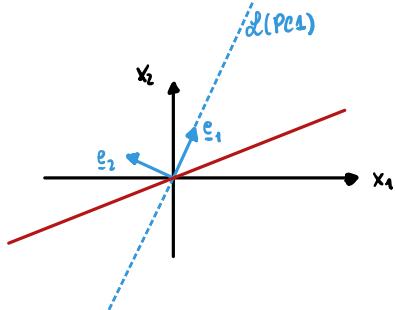
$\Rightarrow$  PCA of  $\rho$

Ex

$$x \in \mathbb{R}^2 \quad \Sigma = \begin{bmatrix} 1 & 4 \\ 4 & 100 \end{bmatrix} \Rightarrow \rho = \begin{bmatrix} 1 & 0.4 \\ 0.4 & 1 \end{bmatrix} \quad \mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

• PCA on  $\Sigma$ :

$$\begin{aligned} \lambda_1 &= 100.46 & e_1 &= (0.04, 0.999)^T \quad (e_1 \perp e_2) \\ \lambda_2 &= 0.84 & e_2 &= (0.999, -0.04)^T \end{aligned}$$



• linear space identified by PC1 ( $e_1$ )

$$0.999 x_1 - 0.04 x_2 = 0$$

$\Leftrightarrow x_2 = \frac{0.999}{0.04} x_1 \approx 25 x_1$

• Regression line

$$\left( \frac{y - \mu_y}{\sqrt{\sigma_{yy}}} = 10 \frac{x - \mu_x}{\sqrt{\sigma_{xx}}} \right)$$

$$\Leftrightarrow \frac{x_2}{10} = 0.4 \frac{x_1}{1} \Rightarrow x_2 = 4x_1$$

• PCA on  $\rho$ :

$$\begin{aligned} \lambda_1 &= 1.4 & e_1 &= (\sqrt{1/2}, \sqrt{1/2})^T \\ \lambda_2 &= 0.6 & e_2 &= (\sqrt{1/2}, -\sqrt{1/2})^T \end{aligned}$$

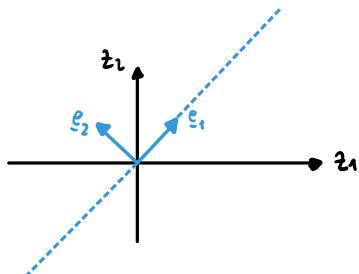
•  $z_2 = z_1$

$$\frac{x_2}{10} = x_1 \Rightarrow x_2 = 10 x_1$$

•  $z_2 = 0.4 z_1$

$$\frac{x_2}{10} = 0.4 x_1 \Rightarrow x_2 = 4 x_1$$

the regression line  
is SCALE-INVARIANT



In general

$$\Sigma = \begin{bmatrix} 1 & r \\ r & 1 \end{bmatrix} \quad \text{reg line } x_2 = r x_1$$

PC1       $x_2 = x_1$

HOWEVER  $\mu, \Sigma$  are unknown, but we can estimate them with data!

$$\hat{\mu} = \bar{x} \quad \hat{\Sigma} = \frac{1}{m-1} \sum_{i=1}^m (x_i - \bar{x})(x_i - \bar{x})^T$$

$$x_i \mapsto y_i = \begin{bmatrix} e_1^T x_i \\ \vdots \\ e_k^T x_i \end{bmatrix}$$

$$\mathbb{X} = \begin{bmatrix} x_1 & \cdots & x_p \\ x_{11} & \cdots & x_{1p} \\ x_{21} & \cdots & x_{2p} \\ \vdots & & \vdots \\ x_{m1} & \cdots & x_{mp} \end{bmatrix} \mapsto$$

$$\begin{bmatrix} y_1 & y_2 & \cdots & y_p \\ e_1^T x_1 & \vdots & e_1^T x_m \\ e_2^T x_1 & \vdots & e_2^T x_m \\ \vdots & & \vdots \\ e_k^T x_1 & \vdots & e_k^T x_m \end{bmatrix}$$

only keep first k components

What is the right k? ART...



## OBSERVATION

↳ for  $k=1, 2, 3$  you can PLOT

↳  $\text{Var}(Y_i) = \lambda_i$

↳  $\lambda_1$  explains  $\frac{\lambda_1}{\sum \lambda_i}$  of the tot. var.

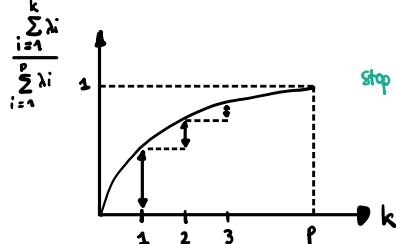
↳  $\lambda_1 + \lambda_2$  explain  $\frac{\lambda_1 + \lambda_2}{\sum \lambda_i}$  of the tot. var.  
⋮

Stop at 0.80 (80%) rule of thumb

↳ For PCA on  $\mathbf{y}$

$$\hat{\lambda} = \frac{\sum \lambda_i}{p} = \frac{tr \mathbf{S}}{p} = \frac{p}{p} = 1 \rightarrow \text{take the components for which the eigenvalues are } \lambda_i > 1$$

↳



Stop when adding one dimension doesn't explain a valuable amount of variability

## DIFFERENT PERSPECTIVE ON PCA

↳ optimal empirical orthonormal basis

Problem: Find linear space  $\mathcal{L}$  of dim.  $k < p$  "closest" to the points.

$\mathcal{L} = \text{span}(\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_k)$  o.m. basis

Find  $\mathbf{q}_1, \dots, \mathbf{q}_k$  s.t.  $\underbrace{\sum_{i=1}^m \|(\mathbf{x}_i - \bar{\mathbf{x}}) - \sum_{j=1}^k q_j (q_j^\top (\mathbf{x}_i - \bar{\mathbf{x}}))\|^2}_{(*)}$  is minimum

$$\begin{aligned} (*) &= \left\| \mathbf{y}_i - \sum_{j=1}^k q_j (q_j^\top \mathbf{y}_i) \right\|^2 = \left( \mathbf{y}_i - \sum_{j=1}^k q_j (q_j^\top \mathbf{y}_i) \right)^\top \left( \mathbf{y}_i - \sum_{j=1}^k q_j (q_j^\top \mathbf{y}_i) \right) \\ &= \mathbf{y}_i^\top \mathbf{y}_i - 2 \sum_{j=1}^k \mathbf{y}_i^\top q_j q_j^\top \mathbf{y}_i + \sum_{j=1}^k (q_j^\top \mathbf{y}_i)^2 \\ &= \mathbf{y}_i^\top \mathbf{y}_i - \sum_{j=1}^k (q_j^\top \mathbf{y}_i)^2 \end{aligned}$$

Hence  $(*) = \sum_{i=1}^m (\mathbf{y}_i^\top \mathbf{y}_i - \sum_{j=1}^k (q_j^\top \mathbf{y}_i)^2)$  to be min

$\Leftrightarrow \sum_{i=1}^m \sum_{j=1}^k (q_j^\top \mathbf{y}_i)^2$  to be max

$$\sum_i \sum_j q_j^\top \mathbf{y}_i \mathbf{y}_i^\top q_j = \sum_j q_j^\top \underbrace{\sum_i \mathbf{y}_i \mathbf{y}_i^\top}_{\text{denote } S} q_j = (m-1) \sum_j q_j^\top S q_j \text{ to be max}$$

$$\Leftrightarrow K=1 \Rightarrow \max_{\substack{\|\mathbf{q}_1\|=1 \\ \mathbf{q}_1 \in \mathbb{R}^p}} \mathbf{q}_1^\top S \mathbf{q}_1 = \lambda_1 \quad \mathbf{q}_1 = \mathbf{e}_1 \quad S = \sum_{i=1}^p \lambda_i \mathbf{e}_i \mathbf{e}_i^\top$$

↳ by induction:  $\mathbf{q}_1 = \mathbf{e}_1, \dots, \mathbf{q}_k = \mathbf{e}_k$

## APPROXIMATION ERROR

$$\sum_{i=1}^m \|(\mathbf{x}_i - \bar{\mathbf{x}}) - \sum_{j=1}^k q_j q_j^\top (\mathbf{x}_i - \bar{\mathbf{x}})\|^2 = \underbrace{\sum_{i=1}^m (\mathbf{x}_i - \bar{\mathbf{x}})' (\mathbf{x}_i - \bar{\mathbf{x}})}_{CS} - \underbrace{(m-1) \sum_{j=1}^k q_j^\top S q_j}_{(m-1) \sum_{j=1}^k \lambda_j} = (m-1) \sum_{i=1}^p \lambda_i - (m-1) \sum_{i=1}^p \lambda_i = (m-1) \sum_{i=k+1}^p \lambda_i$$

$$(2) = \text{tr} \left( \sum_{i=1}^m (\mathbf{x}_i - \bar{\mathbf{x}})' (\mathbf{x}_i - \bar{\mathbf{x}}) \right) = \sum_{i=1}^m \text{tr} [(\mathbf{x}_i - \bar{\mathbf{x}})' (\mathbf{x}_i - \bar{\mathbf{x}})] = \sum_{i=1}^m \text{tr} [(\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})'] = \text{tr} \left[ \sum_{i=1}^m (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})' \right] = (m-1) \text{tr}(S) = (m-1) \sum_{i=1}^p \lambda_i$$

left out from PCA



## SVD ( Singular Value Decomposition)

$$\mathbf{X} = \mathbf{U} \mathbf{D} \mathbf{V} \quad \mathbf{D} = \begin{bmatrix} d_1 & & 0 \\ & \ddots & \\ 0 & & d_p \end{bmatrix}$$

$m \times n$   $m \times p$   $p \times p$   
unitary      unitary

$$\mathbf{S} = \frac{1}{m-1} \mathbf{X}' \mathbf{X} = \frac{1}{m-1} \mathbf{V}' \mathbf{D} \underbrace{\mathbf{U}' \mathbf{U}}_{\mathbf{I}} \mathbf{D} \mathbf{V} = \frac{1}{m-1} \mathbf{V}' \mathbf{D}' \mathbf{D} \mathbf{V} = \frac{1}{m-1} \mathbf{P} \Lambda \mathbf{P} \quad \text{you can reach the same result!}$$

## LECTURE 8 11/3/2021

### Gaussian distribution

$\mathbb{R}^p \ni \underline{x} \sim N_p(\mu, \Sigma)$  if  $P[\underline{x} \in B] = \int_B \varphi_{\underline{x}}(t) dt \quad \forall B \in \mathbb{R}^p$  Borel

$\mathbb{R}^p \ni \underline{x} \sim N_p(\mu, \Sigma)$  if  $\mathbb{R}^p \ni \underline{x} \sim N_p(\mu, \Sigma)$  is still gaussian

#### PROP

- ↳  $E[\underline{x}] = \mu$
- ↳  $Cov(\underline{x}) = \Sigma$

$$\varphi_{\underline{x}}(t) = \frac{1}{\sqrt{(2\pi)^p \det \Sigma}} \exp \left\{ -\frac{1}{2} (\underline{x} - \mu)' \Sigma^{-1} (\underline{x} - \mu) \right\}$$

### THEOREM

$$\underline{x} \sim N_p(\mu, \Sigma) \iff \forall \underline{a} \in \mathbb{R}^q, \underline{a}' \underline{x} \sim N_q(\underline{a}' \mu, \underline{a}' \Sigma \underline{a})$$

Proof un characteristic function

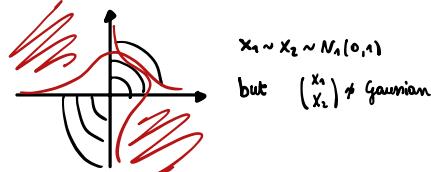
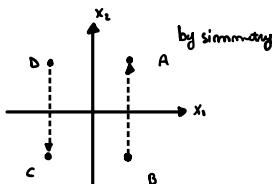
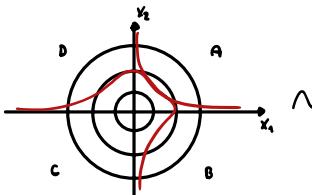
### COROLLARY

$$\underline{x} \sim N_p(\mu, \Sigma) \Rightarrow x_i \sim N_1(\mu_i, \sigma_{ii}) \quad \forall i=1, \dots, p$$

Proof  $x_i = \underline{a}_i' \underline{x} \quad \underline{a}_i = (0, \dots, 0, 1, 0, \dots, 0)$   
 Then  $\underline{x} \sim N_p(\mu, \Sigma)$   
 $\therefore x_i \sim N_1(\mu_i, \sigma_{ii})$   
 $= N_1(\mu_i, \sigma_{ii})$

### OBS

$$\begin{aligned} \underline{x} &\sim N_p(0, I) \\ x_1 &\sim N_1(0, 1) \\ x_2 &\sim N_1(0, 1) \end{aligned}$$



### PROP

$$\underline{x} \sim N_p(\mu, \Sigma) \Rightarrow A\underline{x} \sim N_q(A\mu, A\Sigma A')$$

Proof We need to prove  $\forall \underline{a} \in \mathbb{R}^q \quad \underline{a}' (A\underline{x}) \sim N_q(\underline{a}' (A\mu), \underline{a}' (A\Sigma A') \underline{a})$

$$(\underline{a}' A) \underline{x} = (A' \underline{a})' \underline{x} \stackrel{\text{Then}}{\sim} N_1 \left( \underbrace{(A' \underline{a})' \mu}_{\underline{a}' (A\mu)}, \underbrace{(A' \underline{a})' \Sigma (A' \underline{a})}_{\underline{a}' (A\Sigma A') \underline{a}} \right) \quad \blacksquare$$

Prop  $\underline{x} \sim N_p(\mu, \Sigma) \quad d \in \mathbb{R}^p \Rightarrow \underline{x} + d \sim N_p(\mu + d, \Sigma)$

Proof Exercise.

OBS  $\underline{z}_1, \dots, \underline{z}_p$  iid  $\sim N_1(0, 1)$

$$\underline{z} = (z_1, \dots, z_p)'$$

$$\underline{t} = (t_1, \dots, t_p)' \quad \varphi_{\underline{z}}(\underline{t}) = \prod_{i=1}^p \varphi_{z_i}(t_i) = \prod_{i=1}^p \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} t_i^2} = \frac{1}{\sqrt{(2\pi)^p}} \exp \left[ -\frac{1}{2} \sum_{i=1}^p t_i^2 \right] \quad \Rightarrow \underline{z} \sim N_p(0, I)$$

$\mu \in \mathbb{R}^p, \Sigma \in \mathbb{R}^{p,p}$  pos. def

$$\underline{x} = \sum^{1/2} \underline{z} + \mu \sim N_p(\mu, \Sigma) \quad \text{If you can generate random } \overset{\text{gauss.}}{\text{distribution}} \text{ in 1D, you can in any D.}$$

$$\text{OBS} \quad \mathbb{R}^p \ni \underline{x} = \begin{pmatrix} x_1 \in \mathbb{R}^q \\ x_2 \in \mathbb{R}^{p-q} \end{pmatrix}, \quad \mu = \begin{pmatrix} \mu_1 \in \mathbb{R}^q \\ \mu_2 \in \mathbb{R}^{p-q} \end{pmatrix} \quad \Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}$$

$$\underline{\text{PROP}} \quad (\underline{x}_1, \underline{x}_2) = \underline{x} \sim N_p \left( \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \right) \Rightarrow \underline{x} \sim N_p(\underline{\mu}, \Sigma)$$

Proof ...

$$\underline{\text{PROP}} \quad (\underline{x}_1, \underline{x}_2) = \underline{x} \sim N_p \left( \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \right) \quad \text{If } \Sigma_{12} = 0 \Rightarrow \underline{x}_1 \perp\!\!\!\perp \underline{x}_2$$

Proof Show that  $\psi_{\underline{x}}(\underline{t}) = \psi_{x_1}(t_1) \cdot \psi_{x_2}(t_2)$

THEOREM  $\underline{x}_1 | \underline{x}_2 = \underline{x}_2 \sim N_q \left( \mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (\underline{x}_2 - \mu_2), \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \right)$

Ex  $p=2$

$$\mathbb{R}^2 \ni \underline{x} = \begin{pmatrix} y \\ x \end{pmatrix} \sim N_2 \left( \begin{pmatrix} \mu_y \\ \mu_x \end{pmatrix}, \begin{pmatrix} \sigma_{yy} & \sigma_{yx} \\ \sigma_{xy} & \sigma_{xx} \end{pmatrix} \right) \quad Y \sim N_1(\mu_y, \sigma_{yy})$$

When conditioning on  $\underline{x}$ :

- If  $\rho \approx 1 \Rightarrow (*) \approx 0$  so I will be pretty sure on  $y$
- If  $\rho \approx 0 \Rightarrow (*) \approx \sigma_{yy}$  my information doesn't change when knowing what's carried by  $x$

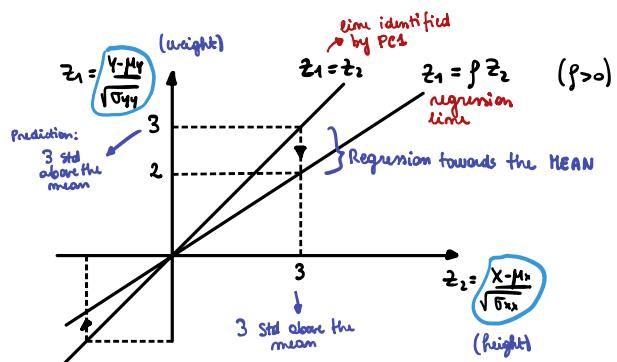
I know this even BEFORE observing  $x$

$$\mathbb{E}[Y | X=x] = \mu_y + \rho \frac{\sqrt{\sigma_{yy}}}{\sqrt{\sigma_{xx}}} (x - \mu_x)$$

$$\hookrightarrow \text{regression line: } y = \mu_y + \rho \frac{\sqrt{\sigma_{yy}}}{\sqrt{\sigma_{xx}}} (x - \mu_x)$$

$$\frac{y - \mu_y}{\sqrt{\sigma_{yy}}} = \rho \frac{x - \mu_x}{\sqrt{\sigma_{xx}}} \quad \begin{matrix} \text{slope} \\ \text{standardized variables} \end{matrix}$$

$$\begin{aligned} Y | X=x &\sim N_1 \left( \underbrace{\mu_y + \sigma_{yx} \cdot \sigma_{xx}^{-1} (x - \mu_x)}_{(*)}, \underbrace{\sigma_{yy} - \sigma_{xy} \sigma_{xx}^{-1} \sigma_{yx}}_{\sigma_{yy}(1 - \frac{\sigma_{xy}^2}{\sigma_{xx} \sigma_{yy}})} \right) \\ &= N \left( \mu_y + \rho \frac{\sqrt{\sigma_{yy}}}{\sqrt{\sigma_{xx}}} (x - \mu_x), \sigma_{yy}(1 - \rho^2) \right) \end{aligned}$$



## LECTURE 9 15/3/2022

trust? difficult

$\underline{x}_i$  realization of  $\underline{X}_i$ ,  $\underline{x}_1, \dots, \underline{x}_m$  iid  $\sim N_p(\underline{\mu}, \Sigma)$

Obvious candidates as estimators:

- for  $\underline{\mu}$ :  $\frac{1}{m} \sum_{i=1}^m \underline{x}_i$
- for  $\Sigma$ :  $\frac{1}{m-1} \sum_{i=1}^m (\underline{x}_i - \bar{\underline{x}})(\underline{x}_i - \bar{\underline{x}})^T$  (Methods of moments)

For gaussian we also have ML

### Maximum Likelihood

Suppose you toss a coin 5 times: HTHT

Model:  $X_1, X_2, X_3, X_4, X_5$  iid  $\sim \text{Be}(p)$   $p?$

$$P(\text{What has been observed}) = P_p(\text{HTHT}) = p^2(1-p)^3 \Rightarrow L(p) = p^2(1-p)^3$$

$$\text{In general } L(p | x_1, \dots, x_m) = p^{\sum x_i} (1-p)^{m - \sum x_i}$$

$$\hat{p} = \underset{p}{\operatorname{argmax}} L(p) = \underset{p}{\operatorname{argmax}} \log L(p)$$

func. of  $p$

$$\begin{aligned} L'(p) &= 2p - 5p^4 = 0 \\ &\Rightarrow p = 0 \vee \hat{p} = \sqrt[3]{\frac{2}{5}} \end{aligned}$$

A believes  $p=0.4$   $L(0.4) =$  [ ]

B believes  $p=0.5$   $L(0.5) =$  [ ]

$$L(\underline{\mu}, \Sigma | \underline{x}_1, \dots, \underline{x}_m) = \prod_{i=1}^m \frac{1}{\sqrt{(2\pi)^p \det \Sigma}} \exp \left\{ -\frac{1}{2} (\underline{x}_i - \underline{\mu}) \Sigma^{-1} (\underline{x}_i - \underline{\mu}) \right\}$$

$$(\hat{\underline{\mu}}, \hat{\Sigma}) = \underset{\substack{\underline{\mu} \in \mathbb{R}^p \\ \Sigma \text{ pos def}}}{\operatorname{argmax}} L(\underline{\mu}, \Sigma) = \dots = \left( \bar{\underline{x}}, \frac{m-1}{m} S \right)$$

## Invariance principle of MLE

Suppose  $\theta \in \mathbb{R}^k$  is a parameter, suppose  $\hat{\theta}$  is MLE for  $\theta$ . Consider  $h: \mathbb{R}^k \rightarrow \mathbb{R}^j$  a mapping. Then  $h(\hat{\theta}) = \widehat{h(\theta)}$

$$\text{Ex } \sum_i \lambda_i \tilde{x}_i \tilde{\varepsilon}_i, \quad \sum_i \tilde{\lambda}_i \tilde{x}_i \tilde{\varepsilon}_i \stackrel{\text{MLE for } \Sigma}{\Rightarrow} \tilde{\lambda}_i = \lambda_i \quad \tilde{\varepsilon}_i = \hat{\varepsilon}_i$$

$$\text{PROP } \bar{X} = \frac{1}{m} \sum_{i=1}^m x_i \sim N_p(\mu, \frac{1}{m} \Sigma)$$

$$\text{Proof } \bar{X} = \begin{bmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_m \end{bmatrix} \in \mathbb{R}^{mp} \sim N_{mp}\left(\begin{pmatrix} \mu \\ \vdots \\ \mu \end{pmatrix}, \underbrace{\begin{pmatrix} \Sigma & & \\ & \ddots & \\ & & \Sigma \end{pmatrix}}_{m \text{ times}}\right) \quad \text{Let } A = [\underbrace{I_{pp} \cdots I_{pp}}_{m \text{ times}}] \in \mathbb{R}^{P \times mp}$$

$$\Rightarrow A \bar{X} = \left( \sum_{i=1}^m x_{i1}, \dots, \sum_{i=1}^m x_{ip} \right)' \Rightarrow \frac{1}{m} A \bar{X} = \bar{X} \sim N_p\left(\frac{1}{m} A \begin{pmatrix} \mu \\ \vdots \\ \mu \end{pmatrix}, \frac{1}{m^2} A \Sigma A'\right) \dots$$

$$\text{DEF } \underline{z}_1, \dots, \underline{z}_m \stackrel{\text{iid}}{\sim} N_p(0, \Sigma) \quad \Sigma \text{ pos-def } (\det \Sigma > 0) \quad \text{Then } \sum_{i=1}^m \underline{z}_i \underline{z}_i' \sim \text{Wishart}(\Sigma, m)$$

(we don't write it, it's clear from the context)

## Properties

$$\begin{aligned} 1) \quad & A_1 \sim \text{Wishart}(\Sigma, m_1) \\ & A_2 \sim \text{Wishart}(\Sigma, m_2) \\ & A_1 \perp\!\!\!\perp A_2 \end{aligned} \Rightarrow A_1 + A_2 \sim \text{Wishart}(\Sigma, m_1 + m_2)$$

$$\text{Proof } A_1 = \sum_{i=1}^{m_1} \underline{z}_i \underline{z}_i', \quad A_2 = \sum_{i=1}^{m_2} \underline{z}_i \underline{z}_i' \Rightarrow \text{we order them } \underline{z}_1, \dots, \underline{z}_{m_1}, \underline{z}_1, \dots, \underline{z}_{m_2} \\ \underline{w}_1, \dots, \underline{w}_{m_1}, \underline{w}_{m_1+1}, \dots, \underline{w}_{m_1+m_2} \Rightarrow \underline{w}_1, \dots, \underline{w}_{m_1+m_2} \sim N_p(0, \Sigma)$$

$$\Rightarrow A_1 + A_2 = \sum_{i=1}^{m_1} \underline{z}_i \underline{z}_i' + \sum_{i=m_1+1}^{m_1+m_2} \underline{z}_i \underline{z}_i' = \sum_{i=1}^{m_1+m_2} \underline{w}_i \underline{w}_i' \sim \text{Wishart}(\Sigma, m_1 + m_2)$$

$$2) \quad C \in \mathbb{R}^{k \times p} \quad \left. \begin{array}{l} \\ A \sim \text{Wishart}(\Sigma, m) \end{array} \right\} \Rightarrow CAC' \sim \text{Wishart}(C\Sigma C', m)$$

$$\text{Proof } A = \sum_{i=1}^m \underline{z}_i \underline{z}_i' \Rightarrow CAC' = \sum_{i=1}^m C \underline{z}_i \underline{z}_i' C' = \sum_{i=1}^m \underline{w}_i \underline{w}_i' \quad \text{where } \underline{w}_i = C \underline{z}_i \sim N_k(0, C\Sigma C')$$

$$3) \quad \sigma^2 \in \mathbb{R}, \sigma^2 > 0 \quad \left. \begin{array}{l} \\ A \sim \text{Wishart}(\Sigma, m) \end{array} \right\} \Rightarrow \sigma^2 A \sim \text{Wishart}(\sigma^2 \Sigma, m)$$

$$\text{Proof } \sigma^2 A = \sum_{i=1}^m \sigma \underline{z}_i \underline{z}_i' \sigma \quad \underline{w}_i = \sigma \underline{z}_i \sim N_p(0, \sigma^2 \Sigma) \quad \text{use (2)}$$

$$4) \quad A \sim \text{Wishart}(\Sigma, m) \quad \Sigma = [\sigma^2] \quad (p=1) \\ A = \sum_{i=1}^m \underline{z}_i \underline{z}_i' = \sum_{i=1}^m \underline{z}_i^2 \Rightarrow \frac{1}{\sigma^2} A = \sum_{i=1}^m \frac{\underline{z}_i}{\sigma} \frac{\underline{z}_i}{\sigma} = \sum_{i=1}^m \underline{w}_i \underline{w}_i' \quad \underline{w}_i \sim N(0, 1) \Rightarrow \frac{1}{\sigma^2} A \sim \chi^2(m) \Rightarrow "A \sim \sigma^2 \chi^2(m)"$$

The Wishart is a multivariate extension of  $\chi^2$ .

$$\text{Now, if } A \sim \text{Wishart}(\Sigma, m), \quad 0 \neq \Sigma \in \mathbb{R}^p, \quad \Sigma \in \mathbb{R}^{p \times p} \stackrel{(2)}{\Rightarrow} \Sigma' A \Sigma \sim \text{Wishart}\left(\underbrace{\Sigma' \Sigma}_{>0}, m\right) \stackrel{(4)}{\sim} \Sigma' \Sigma \sim \chi^2(m)$$

$$\Rightarrow \Sigma' A \Sigma \sim (\Sigma' \Sigma) \chi^2(m)$$

$$\Rightarrow \Sigma' A \Sigma \sim \frac{1}{\Sigma' \Sigma} \sim \chi^2(m)$$

## THEOREM

$$x_1, \dots, x_m \stackrel{\text{iid}}{\sim} N_p(\mu, \Sigma) \Rightarrow \sum_{j=1}^m (x_j - \bar{x})(x_j - \bar{x})' \sim \text{Wishart}(\Sigma, m-1)$$

COROLLARY From + (3)

$$S = \frac{1}{m-1} \sum_{i=1}^m (x_i - \bar{x})(x_i - \bar{x})' \sim \text{Wishart}\left(\frac{1}{m-1} \Sigma, m-1\right)$$

$$\hat{\Sigma} = \frac{m-1}{m} S \sim \text{Wishart}\left(\frac{1}{m} \Sigma, m-1\right)$$

## Summary (THEOREM)

$X_1, \dots, X_m \stackrel{iid}{\sim} N_p(\mu, \Sigma)$  :

$$1) \bar{X} \sim N_p(\mu, \frac{1}{m}\Sigma)$$

$$2) (m-1)S \sim Wishart(\Sigma, m-1)$$

$$3) \bar{X} \perp\!\!\!\perp S$$

## LECTURE 10 17/3/2022

### THEOREM (LLN) (weak)

$\mathbb{R}^p \ni X_1, \dots, X_m \stackrel{iid}{\sim} \mu, \Sigma$

$$\bar{X}_m = \frac{1}{m} \sum_{i=1}^m X_i \xrightarrow{P} \mu$$

$$S = \frac{1}{m-1} \sum_{i=1}^m (X_i - \bar{X})(X_i - \bar{X})' \xrightarrow{P} \Sigma$$

### THEOREM CLT

$$\sqrt{m}(\bar{X} - \mu) \sim AN_p(0, \Sigma)$$

$$\sqrt{m}(\bar{X}_m - \mu) \xrightarrow{d} N_p(0, \Sigma)$$

### INFERENCE ON $\mu$

#### ► CASE $m \gg p$

Then CLT yields  $\sqrt{m}(\bar{X} - \mu) \xrightarrow{d} N_p(0, \Sigma) \Rightarrow \frac{\sqrt{m}(\bar{X} - \mu)}{\sqrt{m}} \Sigma^{-1}(\bar{X} - \mu) \sqrt{m} \sim \chi^2(p)$   
 ↪ pivotal ( $\chi^2$  doesn't depend on  $\Sigma, \mu$ )

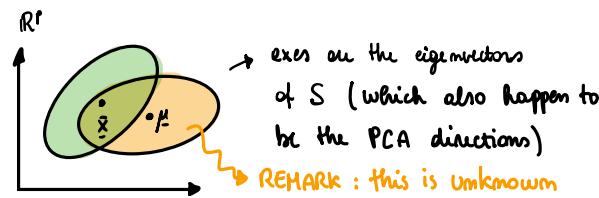
$\Sigma$  is not known, LLN  $S \sim \Sigma$  as  $m \rightarrow \infty$  thus  $m(\bar{X} - \mu) S^{-1}(\bar{X} - \mu) \sim \chi^2(p)$  pivotal

$$\begin{aligned} P(m(\bar{X} - \mu) S^{-1}(\bar{X} - \mu) < \chi_{1-\alpha}^2(p)) &= 1-\alpha \\ \Rightarrow P(\frac{1}{m} S^{-1}(\bar{X} - \mu) \leq \chi_{1-\alpha}^2(p)) &= 1-\alpha \end{aligned}$$

$$\begin{aligned} P(\bar{X} \in \mathcal{E}_{\frac{1}{m} S^{-1}}^{1-\alpha}(\mu)) &= 1-\alpha \\ P(\mu \in \mathcal{E}_{\frac{1}{m} S^{-1}}^{1-\alpha}(\bar{X})) &= 1-\alpha \end{aligned}$$

different centre  $\left\{ \begin{array}{l} \mathcal{E}_{\frac{1}{m} S^{-1}}^{1-\alpha}(\mu) = \{y \in \mathbb{R}^p : m(y - \mu)' S^{-1} (y - \mu) \leq \chi_{1-\alpha}^2(p)\} \\ \mathcal{E}_{\frac{1}{m} S^{-1}}^{1-\alpha}(\bar{X}) = \{y \in \mathbb{R}^p : m(y - \bar{X})' S^{-1} (y - \bar{X}) \leq \chi_{1-\alpha}^2(p)\} \end{array} \right.$

$$\bar{X} \in \mathcal{E}_{\frac{1}{m} S^{-1}}^{1-\alpha}(\mu) \Leftrightarrow \mu \in \mathcal{E}_{\frac{1}{m} S^{-1}}^{1-\alpha}(\bar{X})$$



DEF Confidence region of level  $1-\alpha$ ,  $\alpha \in (0, 1)$ , for  $\mu$ :

$$CR_{1-\alpha}(\mu) = \{y \in \mathbb{R}^p : m(y - \bar{X})' S^{-1} (y - \bar{X}) \leq \chi_{1-\alpha}^2(p)\}$$

►  $(1-\alpha)\%$  of the times this subset will cover the mean  $\mu$

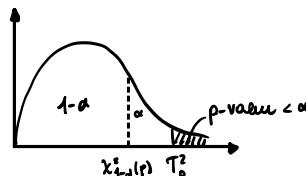
► Obv. I don't know whether I guessed correctly or not.

We can now build tests and the complement of the confidence region: the rejection region.

#### TESTING

$$H_0: \mu = \mu_0 \text{ vs } H_1: \mu \neq \mu_0$$

Evidence against  $H_0$ :  $d_{\frac{1}{m} S^{-1}}(\bar{X}, \mu_0)$  large



If  $H_0$  is TRUE  $\Rightarrow m(\bar{X} - \mu_0)' S^{-1} (\bar{X} - \mu_0) \sim \chi^2(p)$

Reject  $H_0$  if  $T_0^2 = m(\bar{X} - \mu_0)' S^{-1} (\bar{X} - \mu_0) > \chi_{1-\alpha}^2(p)$

↳ Rejection Region  $\{ \text{data} : T_0^2 > \chi_{1-\alpha}^2(p) \} = \{ \text{data} : p\text{-value} < \alpha \}$

► CASE  $m$  SMALL  $\Rightarrow$  assume Gaussianity

$$\underline{x}_1, \dots, \underline{x}_m \stackrel{iid}{\sim} N_p(\mu, \Sigma)$$

$$m(\bar{\underline{x}} - \mu)^\top S^{-1}(\bar{\underline{x}} - \mu) \sim ?$$

DEF  $W \sim \chi^2(m)$ ,  $Y \sim \chi^2(n)$ ,  $W \perp\!\!\!\perp Y \Rightarrow \frac{Y/m}{W/m} \sim F(m, n)$  Fisher's distribution

REMARK  $t = \bar{z}/\sqrt{\frac{W}{m}}$ ,  $\bar{z} \sim N(0, 1)$ ,  $W \sim \chi^2(m)$ ,  $\bar{z} \perp\!\!\!\perp W \Rightarrow t \sim t(m)$

$$\text{Then } t^2 = \frac{\bar{z}^2}{W/m} \quad \frac{\bar{z}^2 \sim \chi^2(1)}{W \sim \chi^2(m)} \Rightarrow t^2 \sim F(1, m)$$

REMARK  $F(m, n) \xrightarrow[m \rightarrow \infty]{\mathcal{L}} \frac{1}{m} \chi^2(m)$

$$\begin{aligned} \underline{Y} &\sim \chi^2(m) & \frac{W}{m} &\xrightarrow[m \rightarrow \infty]{} ? & W = \sum_{i=1}^m z_i^2 &\xrightarrow[m \rightarrow \infty]{\mathcal{L}N(0, 1)} \mathbb{E}[z_i^2] = 1 & \text{Var}(z_i) = \mathbb{E}[z_i^2] - (\mathbb{E}[z_i])^2 = 1 - 0 = 1 \\ \underline{W} &\sim \chi^2(m) \end{aligned}$$

### HOTELLING'S THEOREM (1931)

$$\left. \begin{array}{l} \underline{X} \sim N_p(\mu, \Sigma) \text{ det } \Sigma > 0 \\ W \sim \text{Wish} \left( \frac{1}{m} \Sigma, m \right) \\ \underline{X} \perp\!\!\!\perp W \end{array} \right\} \Rightarrow \frac{m-p+1}{mp} (\bar{\underline{x}} - \mu)^\top W^{-1} (\bar{\underline{x}} - \mu) \sim F(p, m-p+1)$$

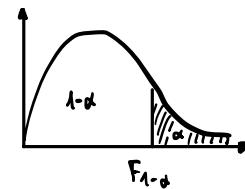
COROLLARY  $\underline{x}_1, \dots, \underline{x}_m \stackrel{iid}{\sim} N_p(\mu, \Sigma)$ ,  $\text{det } \Sigma > 0 \Rightarrow m(\bar{\underline{x}} - \mu)^\top S^{-1}(\bar{\underline{x}} - \mu) \sim \frac{(m-1)p}{m-p} F(p, m-p)$  pivotal

$$\begin{array}{l} \text{Proof} \quad \sqrt{m}(\bar{\underline{x}} - \mu) \sim N_p(0, \Sigma) \text{ (exact)} \\ S \sim \text{Wishart} \left( \frac{1}{m-1} \Sigma, m-1 \right) \\ \bar{\underline{x}} \perp\!\!\!\perp S \end{array} \left. \begin{array}{l} \xrightarrow[m=m-1]{\text{Hotelling}} \\ \xrightarrow[1]{\text{Hotelling}} \end{array} \right\} \frac{m-1-p+1}{(m-1)p} m(\bar{\underline{x}} - \mu)^\top S^{-1}(\bar{\underline{x}} - \mu) \sim F(p, m-1-p+1)$$

### LECTURE 11 18/3/2022

Def Hotelling's  $T^2$  statistic :  $T^2 = m(\bar{\underline{x}} - \mu)^\top S^{-1}(\bar{\underline{x}} - \mu)$

CONFIDENCE REGIONS Fix  $\alpha \in (0, 1)$   $P(T^2 \leq \frac{(m-1)p}{m-p} F_\alpha(p, m-p)) = 1-\alpha$



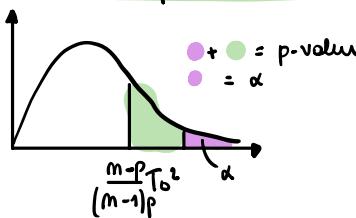
$$\text{CR}_{1-\alpha}(\mu) = \{ \eta \in \mathbb{R}^p : (\eta - \bar{\underline{x}})^\top S^{-1}(\eta - \bar{\underline{x}}) \leq \frac{(m-1)p}{m-p} F_\alpha(p, m-p) \} \quad \text{radius from the "m large" case}$$

#### TESTING

$$H_0: \mu = \mu_0 \text{ vs } H_1: \mu \neq \mu_0$$

$$\text{If } H_0 \text{ true } \Rightarrow T_0^2 \sim m(\bar{\underline{x}} - \mu_0)^\top S^{-1}(\bar{\underline{x}} - \mu_0) \sim \frac{(m-1)p}{m-p} F(p, m-p)$$

Reject  $H_0$  if  $T_0^2 > \frac{(m-1)p}{m-p} F_\alpha(p, m-p)$



REM Don't forget about  $\frac{(m-p)}{(m-1)p}$  !

► The CR identifies all the  $\mu_0$  for we can't reject  $H_0$

What happens as  $m \rightarrow \infty$ ?

$$\frac{(m-1)p}{m-p} F_{1-\alpha}(p, m-p) \xrightarrow[m \rightarrow \infty]{} \chi^2_{1-\alpha}(p) \quad \text{as we expected from the previous case}$$

## INFERENCE FOR LINEAR COMBINATIONS OF $\mu$

$\underline{a} \in \mathbb{R}^p$   $\underline{a}'\mu = a_1\mu_1 + \dots + a_p\mu_p$  (my portfolio and not all the market stocks)

$$\text{Ex } \underline{a} = \underline{\alpha}_i = (0 \dots \underset{i}{1} \dots 0) \Rightarrow \underline{a}'\mu = \mu_i$$

$$\underline{a} = (0 \dots 0 \underset{i}{1} 0 \dots 0 \underset{j}{-1} 0 \dots 0) \Rightarrow \underline{a}'\mu = \mu_i - \mu_j$$

### ► CASE "ONE-AT-THE-TIME"

►  $\underline{a}'\bar{X}$  (unbiased) estimator for  $\underline{a}'\mu$

► As before : •  $m$  small  $\Rightarrow$  Gaussianity  
•  $m$  large  $\Rightarrow$  CLT to have approximately Gaussianity

► If  $X_1, \dots, X_m \stackrel{iid}{\sim} N_p(\mu, \Sigma) \Rightarrow \bar{X} \sim N_p(\mu, \frac{1}{m}\Sigma) \Rightarrow \underline{a}'\bar{X} \sim N_1(\underline{a}'\mu, \frac{1}{m}\underline{a}'\Sigma\underline{a}) \Rightarrow$

► If  $m$  is small (assume Gaussianity)

We know  $(m-1)S \sim \text{Wish}(\Sigma, m-1)$

then  $(m-1)\underline{a}'S\underline{a} \sim \underbrace{(\underline{a}'\Sigma\underline{a})}_{\in \mathbb{R}} \chi^2(m-1) \sim \frac{(m-1)\underline{a}'S\underline{a}}{\underline{a}'\Sigma\underline{a}} \sim \chi^2(m-1)$  Monotone  $\bar{X} \perp S$  and so are their transformations, thus

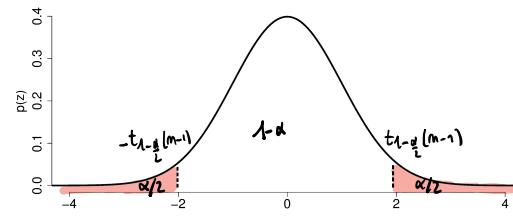
$$\Rightarrow \frac{\frac{\underline{a}'\bar{X} - \underline{a}'\mu}{\sqrt{\underline{a}'\Sigma\underline{a}}} \sqrt{m}}{\sqrt{\frac{(m-1)\underline{a}'S\underline{a}}{\underline{a}'\Sigma\underline{a}} \cdot \frac{1}{m-1}}} \sim t(m-1) \Rightarrow \frac{\underline{a}'\bar{X} - \underline{a}'\mu}{\sqrt{\underline{a}'S\underline{a}}} \sqrt{m} \sim t(m-1) \text{ pivotal}$$

**REMARK**  
 $N(0,1) \sim t(m)$   
 $\sqrt{\frac{\chi^2(m)}{m}}$  (II)

From this

$$\bullet \quad P\left[\frac{|\underline{a}'\bar{X} - \underline{a}'\mu|}{\sqrt{\underline{a}'S\underline{a}}} \sqrt{m} < t_{1-\frac{\alpha}{2}}(m-1)\right] = 1-\alpha = P\left[\underline{a}'\mu \in \left[\underline{a}'\bar{X} \pm t_{1-\frac{\alpha}{2}}(m-1) \sqrt{\frac{\underline{a}'S\underline{a}}{m}}\right]\right] \quad \forall \underline{a} \in \mathbb{R}^p$$

$$\bullet \quad CI_{1-\alpha}(\underline{a}'\mu) = \left[\underline{a}'\bar{X} \pm t_{1-\frac{\alpha}{2}}(m-1) \sqrt{\frac{\underline{a}'S\underline{a}}{m}}\right]$$



Ex  $\underline{a} = \underline{\alpha}_i$

$$CI_{1-\alpha}(\mu_i) = \left[\bar{X}_i \pm t_{1-\frac{\alpha}{2}}(m-1) \sqrt{\frac{S_{ii}}{m}}\right]$$

$$\underline{a} = (0 \dots 0 \underset{i}{1} 0 \dots 0) \quad CI_{1-\alpha}(\mu_i - \mu_j) = \left[\bar{X}_i - \bar{X}_j \pm t_{1-\frac{\alpha}{2}}(m-1) \sqrt{\frac{S_{ii} - 2S_{ij} + S_{jj}}{m}}\right]$$

We can do TESTING

$H_0: \underline{a}'\mu \leq \delta_0$ (Default pollution)	$H_1: \underline{a}'\mu > \delta_0$ (We want to prove that it's lower)	Reject if $t_{1-\alpha/2}(m-1) > t_{1-\alpha}(m-1)$
$t_0 = \frac{\underline{a}'\bar{X} - \delta_0}{\sqrt{\underline{a}'S\underline{a}}} \sqrt{m} \sim t(m-1)$		$H_0: \underline{a}'\mu \geq \delta_0$

### ► CASE "SIMULTANEOUS"

In the above we did NOT prove that  $P\left(\underline{a}'\mu \in \left[\underline{a}'\bar{X} \pm t_{1-\frac{\alpha}{2}}(m-1) \sqrt{\frac{\underline{a}'S\underline{a}}{m}}\right], \forall \underline{a} \in \mathbb{R}^p\right) = 1-\alpha$

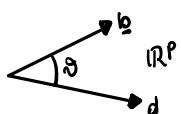
↳ the above would be "no matter what coin I toss I get H 1/2 of the times"

↳ vs "if I toss all the coins together I get all H with 1/2 probability"

What should go here to make the statement true?

New problem:  $P\left(\underline{a}'\mu \in \left[\underline{a}'\bar{X} \pm ? \sqrt{\frac{\underline{a}'S\underline{a}}{m}}\right], \forall \underline{a} \in \mathbb{R}^p\right) = 1-\alpha$

### LINEAR ALGEBRA RECALL



$$\frac{b'd}{\|b\| \|d\|} = \cos \theta$$

$$0 \leq \frac{(b'd)^2}{\|b\|^2 \|d\|^2} = \cos^2 \theta \leq 1$$

$$\Rightarrow (b'd) \leq \|b\|^2 \|d\|^2 \quad \forall b, d \in \mathbb{R}^p$$

"=" if  $b \in L(d)$

/ GS Inequality



Let  $B \in \mathbb{R}^{p \times p}$  pos. def.

Prop  $\forall b, d \in \mathbb{R}^p, (b'd)^2 \leq (b'Bb)(d'B^{-1}d)$

Proof  $(b'd)^2 = (\underbrace{b' B \cancel{B^{-1}} d}_w)^2 \stackrel{\text{CS}}{\leq} \|w\|^2 \|d\|^2 = (b' B \cancel{B^{-1}} B \cancel{B^{-1}} d) = (b' B \cancel{B^{-1}} b)(d' B^{-1} d)$  eq. holds if  $b' B^{-1} b \in L(B^{-1} d)$   
 $w = Bd$   $b \in L(B^{-1} d)$

Prop  $B \in \mathbb{R}^{p \times p}$  pos. def.,  $d \in \mathbb{R}^p \Rightarrow \max_{x \in \mathbb{R}^p, x \neq 0} \frac{(x'd)^2}{x'Bx} = d'B^{-1}d$

Proof CS:  $(x'd)^2 \leq (x'Bx)(d'B^{-1}d)$

$$\text{if } x \neq 0 \Rightarrow x'Bx > 0 \quad (\text{B pos. def}) \Rightarrow \frac{(x'd)^2}{x'Bx} \leq d'B^{-1}d \quad \forall x \neq 0$$

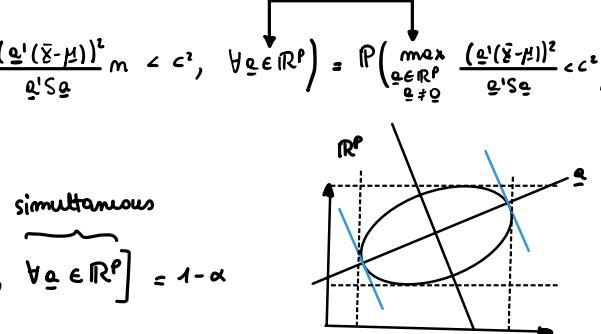
equality if  $x \in L(B^{-1}d)$

For  $a \in \mathbb{R}^p$   $\max_{a \in \mathbb{R}^p, a \neq 0} \frac{(a'(\bar{x}-\mu))^2}{a'Sa} = m(\bar{x}-\mu)'S^{-1}(\bar{x}-\mu) \sim \frac{(m-1)p}{m-p} F(p, m-p)$

we know it's distribution  
(clos to Hotelling's theorem)

Problem Find  $c > 0$  s.t.  $1-\alpha = P\left(\frac{|a'(\bar{x}-\mu)|}{\sqrt{a'Sa}} \sqrt{m} < c, \forall a \in \mathbb{R}^p\right) = P\left(\frac{(a'(\bar{x}-\mu))^2}{a'Sa} m < c^2, \forall a \in \mathbb{R}^p\right) = P\left(\max_{\substack{a \in \mathbb{R}^p \\ a \neq 0}} \frac{(a'(\bar{x}-\mu))^2}{a'Sa} < c^2\right)$

$$\Rightarrow c^2 = \frac{(m-1)p}{m-p} F_{1-\alpha}(p, m-p) \quad \text{☺}$$



We can finally conclude:  $P\left[\frac{|a'(\bar{x}-\mu)|}{\sqrt{a'Sa}} \sqrt{m} < \sqrt{\frac{(m-1)p}{m-p} F_{1-\alpha}(p, m-p)}, \forall a \in \mathbb{R}^p\right] = 1-\alpha$

$$\text{which is the same as } P\left[a'\mu \in \left[\bar{x} \pm \sqrt{\frac{(m-1)p}{m-p} F_{1-\alpha}(p, m-p)} \sqrt{\frac{a'Sa}{m}}\right], \forall a \in \mathbb{R}^p\right] = 1-\alpha$$

from which we easily need the simultaneous confidence interval

$$\text{SimCI}_{1-\alpha}(\mu) = \left[\bar{x} \pm \sqrt{\frac{(m-1)p}{m-p} F_{1-\alpha}(p, m-p)} \sqrt{\frac{a'Sa}{m}}\right]$$

## LECTURE 12 22/3/2022

Thus now can take all the possible linear combination we want and we get a confidence interval that is globally correct  $(1-\alpha)\%$  of times! Indeed all of the interval are correct with probability  $1-\alpha$ : all of the interval cover the linear combination  $(1-\alpha)\%$  of the time they are used!

So if we have an ellipse of the  $CR_{1-\alpha}(\mu)$  then its projection along any direction  $a$  gives us the simultaneous confidence interval  $\text{SimCI}_{1-\alpha}(a'\mu)$

### BONFERRONI STRATEGY FOR SIMULTANEOUS CONFIDENCE INTERVAL

Using the previous method we can ALL possible linear combinations: we'd like to fix a FINITE number!

Let us fix  $a_1, \dots, a_k \in \mathbb{R}^p$ : we want the k CI for  $a_i'\mu, i=1, \dots, k$  with a simultaneous confidence of  $1-\alpha$

Consider  $CI_{1-\alpha}(a_i'\mu) = \left[\bar{x} \pm t_{1-\frac{\alpha}{k}(m-1)} \sqrt{\frac{a_i'Sa_i}{m}}\right]$  (one-at-the-time)

Then

$$P(A \cup B) = P(A) + P(B) - P(A \cap B) \leq P(A) + P(B) \quad \text{Bonferroni ineq.}$$

$$P\left[\bigcap_{i=1}^k [a_i'\mu \in CI_{1-\alpha}(a_i'\mu)]\right] = 1 - P\left[\bigcup_{i=1}^k [a_i'\mu \notin CI_{1-\alpha}(a_i'\mu)]\right] \geq 1 - \sum_{i=1}^k P[a_i'\mu \notin CI_{1-\alpha}(a_i'\mu)] = 1 - \sum_{i=1}^k \alpha = 1 - k\alpha$$

$$\Rightarrow \text{Bonf CI}_{1-\alpha}(a_i'\mu) = \left[\bar{x} \pm t_{1-\frac{\alpha}{k}(m-1)} \sqrt{\frac{a_i'Sa_i}{m}}\right] \quad \text{Notice that Bonf CI}_{1-\alpha}(a_i'\mu) \xrightarrow{k \rightarrow \infty} [-\infty, +\infty], \text{ good for small } k$$

Simultaneous Testing can be done but it's very conservative (with big data you'd never reject H<sub>0</sub>)

$$H_0: \begin{cases} H_{01} & a_1'\mu = \delta_1 \\ H_{02} & a_2'\mu = \delta_2 \\ \vdots & \vdots \\ H_{0k} & a_k'\mu = \delta_k \end{cases} \quad \text{vs} \quad H_1: \text{at least one } H_{0i} \text{ is false} \quad \Rightarrow \text{Reject H}_0 \text{ if for at least one } i: \frac{|a_i'\bar{x} - \delta_i|}{\sqrt{a_i'Sa_i}} \sqrt{m} > t_{1-\frac{\alpha}{k}(m-1)}$$



$$\text{Indeed } P[\text{Reject } H_0 \mid H_0 \text{ true}] = P\left[\bigcup_{i=1}^k \{t_i > t_{i-\frac{\alpha}{2k}}(m-i)\} \mid \bigcap_{i=1}^k H_0 \text{ true}\right] \leq \sum_{i=1}^k P[\dots] = \sum_{i=1}^k \frac{\alpha}{k} = \alpha$$

so the overall probability of rejecting one  $H_0$  when all are true is equal to  $\alpha$

**FALSE DISCOVERY RATE (FDR)**  
(Benjamini & Hochberg 1995)

citations  
2014 23'000  
2019 53'000  
2022 83'821 } very important strategy  
for sim. testing!

Let  $\Delta$  be any strategy for testing (e.g. Bonferroni)

Decisions taken following  $\Delta$

		Not reject $H_0$	Reject $H_0$	
		$U$	$V$ false discoveries	$k_0 = U + V$
True	$H_0$	$T$ missed discoveries	$S$ true discoveries	$k_0 - k_0 = T + S$
	$H_1$	$R$		$k - R$

We can observe that

With Bonferroni strategy you reject  $H_0$  at level  $\frac{\alpha}{k}$

Let  $I_0$  be the set of the true  $H_0$ 's  $|I_0| = k_0$

$$\underbrace{P[V \geq 1]}_{\text{Family Wise Error Rate (FWER)}} = P\left[\bigcup_{j \in I_0} \{\text{Reject } H_0\}\right] \leq \sum_{j \in I_0} P[\{\text{Reject } H_0\}] \leq \sum_{j \in I_0} \frac{\alpha}{k} = k_0 \frac{\alpha}{k} \leq k \frac{\alpha}{k} = \alpha$$

Family Wise Error Rate (FWER)

$$\text{We define } Q = \begin{cases} \frac{V}{R}, & R > 0, \Leftrightarrow V > 0 \\ 0, & R = 0, \Leftrightarrow V = 0 \end{cases} \Rightarrow \text{FALSE DISCOVERY RATE } \mathbb{E}[Q]$$

**REMARK**  $k = k_0$   $S = T = 0$   $V = R$  then all no discoveries to be made

$$Q = \{0, 1\}$$

$$\text{FDR} = \mathbb{E}[Q] = P[V > 0] = P[V \geq 1] = \text{FWER}$$

**REMARK**  $k > k_0$   $\oplus$

$$\begin{aligned} \text{if } V = 0 \Rightarrow Q = 0 \\ \text{if } V > 0 \Rightarrow Q = \frac{V}{R} \leq 1 \\ \text{moreover } \mathbb{E}[V > 0] = \begin{cases} 0 & \text{if } V = 0 \\ 1 & \text{if } V > 0 \end{cases} \end{aligned} \Rightarrow Q \leq \mathbb{E}[V > 0]$$

$$\text{thus if } \oplus \quad \text{FDR} = \mathbb{E}[Q] \leq \mathbb{E}[\mathbb{E}[V > 0]] = P[V > 0] = P[V \geq 1] = \text{FWER}$$

We have a way to control FDR!

**Conclusion:** Bonferroni is good but not usable, FDR is less good but usable!

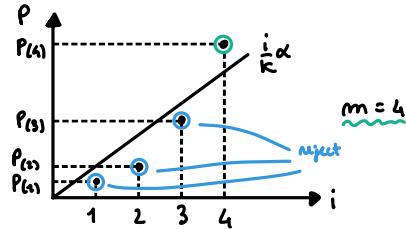
### STRATEGY FOR CONTROLLING FDR

Let  $p_i$  be the p-value of the single test  $H_{0i}$  vs  $H_{1i}$

We order the p-values  $p_{(1)} \leq p_{(2)} \leq \dots \leq p_{(k)}$   
 $\downarrow \quad \downarrow \quad \downarrow$   
 $H_{0(1)} \quad H_{0(2)} \quad H_{0(k)}$



$$m = \max \{i \in \{1, \dots, k\} : p_{(i)} \leq \frac{i}{k} \alpha\} \quad \alpha \in (0, 1) \quad \alpha, k \text{ fixed}$$



**THEOREM (B&H)** If  $p_1, \dots, p_k$  are independent  $\Rightarrow$  The strategy which rejects  $H_{0(i)}$  if  $i \leq m$  controls FDR at level  $\alpha$  ( $\text{FDR} \leq \alpha$ )

### THEOREM (BENJAMINI - YEKUTELI 2001)

- If  $p_1, \dots, p_k$  are positively correlated  $\Rightarrow$  B&H 95 strategy controls FDR at level  $\alpha$

- If  $p_1, \dots, p_k$  are negatively correlated  $\Rightarrow$  diff. procedure: Reject  $H_{0(i)}$  if  $i \leq m^* = \max \{j \in \{1, \dots, k\} : p_{(j)} \leq \frac{j}{C_{k,k}} \alpha\}$

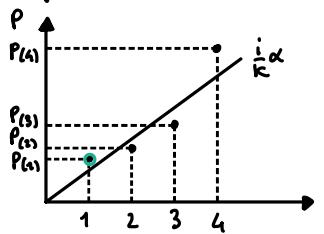


$$\text{where } C(k) = \sum_{j=1}^k \frac{1}{j}$$

Note: mixed cases are not covered: we need both ALL positively correlated, or ALL negatively correlated p-values!

### Efron's modification of B&H strategy

It might happen that  $P(i)$  is above the line, B&H would accept all



LECTURE 13 25/3/2022

## COMPARISON OF MEAN

### ► REPEATED MEASURES

e.g. person, family

For each  $i=1, \dots, m$  statistical unit, with  $p \geq 1$  features, we make  $q \geq 2$  measurements

$$X_{1i} = \begin{pmatrix} \text{height (time 1)} \\ \text{weight (time 1)} \\ \text{salary (time 1)} \end{pmatrix} \xrightarrow{\text{treatment}} X_{2i} = \begin{pmatrix} H & (\text{time 2}) \\ W & (\text{time 2}) \\ S & (\text{time 2}) \end{pmatrix} \quad \begin{matrix} \text{dependence between} \\ \text{No dependence between} \end{matrix}$$

### ► CASE $q=2$ (PAIRS)

$$\begin{aligned} R^P \ni \mu_1 &\leftarrow \left( \begin{matrix} X_{11} \\ X_{12} \\ \vdots \\ X_{1m} \end{matrix} \right), \left( \begin{matrix} X_{21} \\ X_{22} \\ \vdots \\ X_{2m} \end{matrix} \right), \dots, \left( \begin{matrix} X_{m1} \\ X_{m2} \\ \vdots \\ X_{mm} \end{matrix} \right) \} q=2 \\ R^P \ni \mu_2 &\leftarrow \left( \begin{matrix} X_{21} \\ X_{22} \\ \vdots \\ X_{2m} \end{matrix} \right) \end{aligned}$$

REMARK "pairs" have to be meaningful, it doesn't make sense to pair a group of French and Italians based on height

!

The main question is ("did the treatment have an effect?")

$$\left\{ H_0: \mu_1 = \mu_2 \text{ vs } H_1: \mu_1 \neq \mu_2 \right\} \Leftrightarrow \left\{ H_0: \mu_1 - \mu_2 = 0 \text{ vs } H_1: \mu_1 - \mu_2 \neq 0 \right\}$$

$$D_i = X_{1i} - X_{2i} \quad i=1, \dots, m$$

► Assume  $D_1, \dots, D_m$  iid  $\sim N_p(\underline{\delta}, \Sigma_D)$

$$\text{► Test statistic } \bar{D} = \frac{1}{m} \sum_{i=1}^m D_i \quad S_D = \frac{1}{m-1} \sum_{i=1}^m (D_i - \bar{D})(D_i - \bar{D})'$$

$$\Rightarrow m(\bar{D} - \underline{\delta})' S_D^{-1} (\bar{D} - \underline{\delta}) \sim \frac{(m-1)p}{m-p} F(p, m-p)$$

pivotal quantity

$$\bullet \text{Test } H_0: \underline{\delta} = \underline{\delta}_0 \text{ vs } H_1: \underline{\delta} \neq \underline{\delta}_0 \quad (\text{e.g. } \underline{\delta}_0 = 0)$$

$$T_0^2 = m(\bar{D} - \underline{\delta}_0)' S_D^{-1} (\bar{D} - \underline{\delta}_0) \Rightarrow \alpha \in (0, 1) \text{ reject if } T_0^2 > \frac{(m-1)p}{m-p} F_{1-\alpha}(p, m-p)$$

$$\bullet \text{Confidence Region } CR_{1-\alpha}(\underline{\delta}) = CR_{1-\alpha}(\mu_1 - \mu_2) = \{ \underline{\eta} \in \mathbb{R}^p : m(\bar{D} - \underline{\eta})' S_D^{-1} (\bar{D} - \underline{\eta}) \leq \frac{(m-1)p}{m-p} F_{1-\alpha}(p, m-p) \}$$

$$\bullet \text{Simultaneous Confidence Interval } \text{Sim CI}(\underline{\alpha}' \underline{\delta}) = [ \underline{\alpha}' \bar{D} \pm \sqrt{\frac{(m-1)p}{m-p} F_{1-\alpha}(p, m-p)} \sqrt{\frac{\underline{\alpha}' S_D^{-1} \underline{\alpha}}{m}} ]$$

$$\bullet \text{Bonferroni Simultaneous Confidence Interval } \text{Bonf CI}_{1-\alpha}(\mu_{ij} - \mu_{jk}) = [ \bar{D}_{ij} \pm t_{1-\alpha/2p} (m-1) \sqrt{\frac{S_{jj}}{m}} ]$$

### ► CASE $p=1$ (1 feature, but many repeated measures)

$$\text{For } i=1, \dots, m \quad p=1 \quad \underbrace{\{(X_{1i}, X_{2i}, \dots, X_{qi})'\}}_{q \geq 2} \sim N_q(\underline{\mu}, \Sigma) \quad \text{► Assume Gaussianity}$$

$$H_0: \mu_1 = \mu_2 = \dots = \mu_q \quad \text{vs} \quad H_1: \exists i, j \text{ s.t. } \mu_i \neq \mu_j$$

Different approach than before

DEF  $C \in \mathbb{R}^{(q-1) \times q}$  is called CONTRAST MATRIX if  $C = \begin{bmatrix} c_1 \\ \vdots \\ c_{q-1} \end{bmatrix}, c_i \in \mathbb{R}^q$

- 1)  $c_{11}, \dots, c_{q-1,1}$  are lin. indep.
- 2)  $c_i' \underline{1} = 0 \quad \forall i=1, \dots, q-1$  (i.e.  $c_i \perp \underline{1}$ )

$$H_0: C\mu = 0 \quad \text{vs} \quad H_1: C\mu \neq 0$$

$$\bar{X} = \frac{1}{m} \sum_{i=1}^m \underline{x}_i \quad \underline{x}_i = \begin{pmatrix} x_{1i} \\ x_{2i} \\ \vdots \\ x_{qi} \end{pmatrix} \stackrel{iid}{\sim} N_{q-1}(\mu, \Sigma) \quad C\bar{X} \sim N_{q-1}(C\mu, \frac{1}{m} C\Sigma C^T) \Rightarrow \sqrt{m}(C\bar{X} - C\mu) \sim N_{q-1}(0, C\Sigma C^T)$$

$$S = \frac{1}{m-1} \sum_{i=1}^m (\underline{x}_i - \bar{X})(\underline{x}_i - \bar{X})^T \Rightarrow (m-1)S \sim \text{Wishart}(\Sigma, m-1) \Rightarrow (m-1)CSC^T \sim \text{Wishart}(C\Sigma C^T, m-1)$$

$\Rightarrow$  by Hotelling's theorem  $m(C\bar{X} - C\mu)^T (CSC^T)^{-1} (C\bar{X} - C\mu) \sim \frac{(m-1)(q-1)}{m-q+1} F(q-1, m-q+1) \rightarrow$  pivotal quantity

• Test  $H_0: C\mu = 0$  at level  $\alpha$  if  $T_0^2 = m(C\bar{X})^T (CSC^T)^{-1} (C\bar{X}) > \frac{(m-1)(q-1)}{m-q+1} F_{1-\alpha}(q-1, m-q+1)$

...

⚠ I could choose many different basis

Example

$$C_1 = \begin{bmatrix} 1 & -1 & 0 & \cdots & \cdots & 0 \\ 1 & 0 & -1 & 0 & \cdots & 0 \\ \vdots & & & & \ddots & \vdots \\ 1 & \cdots & \cdots & \cdots & \cdots & -1 \\ R \curvearrowleft & & & & & \end{bmatrix} \quad C_2 = \begin{bmatrix} 1 & -1 & 0 & \cdots & \cdots & 0 \\ 0 & 1 & -1 & \cdots & \cdots & 0 \\ \vdots & & & & \ddots & \vdots \\ 0 & \cdots & \cdots & \cdots & \cdots & 1 & -1 \end{bmatrix}$$

Comparing with baseline

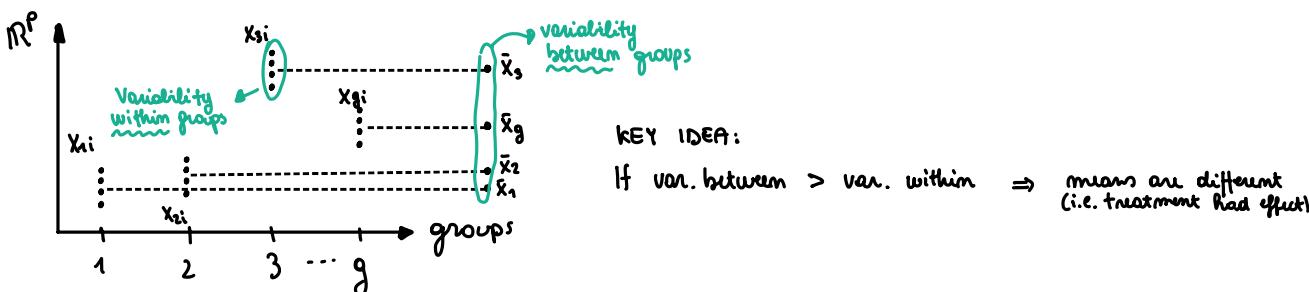
Compare with previous

Notice  $\exists B \in \mathbb{R}^{(q-1) \times (q-1)}$  s.t.  $\det(B) \neq 0$  and  $C_1 = BC_2 \Rightarrow T_0^2 = m(BC_2 \bar{X})^T (BC_2 S C_2^T B^T)^{-1} (BC_2 \bar{X})$   
 $\hookrightarrow$  maps a basis into another  $= m(C_2 \bar{X})^T B^T (B^{-1})^{-1} (C_2 S C_2^T) B^{-1} B (C_2 \bar{X})$   
 same!

## ► MANOVA (MULTI-VARIATE ANALYSIS OF VARIANCE)

$m_1 \dots m_g$  observations  
 $g$  indep. groups  
 $\underbrace{X_{11}, \dots, X_{1m_1}}_{\text{group } 1} \stackrel{iid}{\sim} N_p(\mu_1, \Sigma) \quad \underbrace{X_{21}, \dots, X_{2m_2}}_{\text{group } 2} \stackrel{iid}{\sim} N_p(\mu_2, \Sigma) \quad \vdots \quad \underbrace{X_{g1}, \dots, X_{gm_g}}_{\text{group } g} \stackrel{iid}{\sim} N_p(\mu_g, \Sigma) \quad p \text{ features}$

GOAL: make inference on the differences among groups:  $\mu_1, \dots, \mu_g$   
 $\Sigma$  have to be equal  $\left\{ \begin{array}{l} \cdot \text{ do testing for this} \\ \cdot \text{ if not, transform the data until yes} \end{array} \right.$



► CASE  $g=2, p \geq 1$  (2 groups, many features)

► Assume gaussianity

$$\begin{aligned} \bar{X}_1 &\sim N_p(\mu_1, \frac{1}{m_1} \Sigma) \Rightarrow \bar{X}_1 - \bar{X}_2 \sim N_p(\mu_1 - \mu_2, (\frac{1}{m_1} + \frac{1}{m_2}) \Sigma) \Rightarrow \left( \frac{1}{m_1} + \frac{1}{m_2} \right)^{-1/2} (\bar{X}_1 - \bar{X}_2 - (\mu_1 - \mu_2)) \sim N_p(0, \Sigma) \\ \bar{X}_2 &\sim N_p(\mu_2, \frac{1}{m_2} \Sigma) \end{aligned}$$

$$S_1 = \frac{1}{m_1-1} \sum_{i=1}^{m_1} (\underline{x}_{1i} - \bar{X}_1)(\underline{x}_{1i} - \bar{X}_1)^T \text{ est. of } \Sigma_1 \Rightarrow (m_1-1)S_1 \sim \text{Wishart}(\Sigma, m_1-1) \quad \text{II} \Rightarrow (m_1-1)S_1 + (m_2-1)S_2 \sim \text{Wishart}(\Sigma, m_1+m_2-2)$$

$$S_2 = \frac{1}{m_2-1} \sum_{i=1}^{m_2} (\underline{x}_{2i} - \bar{X}_2)(\underline{x}_{2i} - \bar{X}_2)^T \text{ est. of } \Sigma_2 \Rightarrow (m_2-1)S_2 \sim \text{Wishart}(\Sigma, m_2-1)$$

$$S_{\text{pooled}} = \frac{(m_1-1)S_1 + (m_2-1)S_2}{m_1+m_2-2} \quad \text{s.t. } (m_1+m_2-2) S_{\text{pooled}} \sim \text{Wishart}(\Sigma, m_1+m_2-2)$$

$$\bar{X}_1 \perp\!\!\!\perp S_{\text{pooled}} \perp\!\!\!\perp \bar{X}_2, \text{ since } S_1 \perp\!\!\!\perp \bar{X}_1, S_2 \perp\!\!\!\perp \bar{X}_2 \Rightarrow (m_1+m_2-2) S_{\text{pooled}} \perp\!\!\!\perp \left( \frac{1}{m_1} + \frac{1}{m_2} \right)^{-1/2} [(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)]$$

$$\text{Hotelling} \Rightarrow \left( \frac{1}{m_1} + \frac{1}{m_2} \right)^{-1} [(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)]^T S_{\text{pooled}}^{-1} [(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)] \sim \frac{(m_1+m_2-2)p}{m_1+m_2-1-p} F(p, m_1+m_2-1-p) \rightarrow \text{pivotal}$$

- Test  $H_0: \mu_1 - \mu_2 = \delta_0$  vs  $H_1: \mu_1 - \mu_2 \neq \delta_0$  We reject at level  $\alpha \in (0, 1)$  if  
 $T_0^2 = \left( \frac{1}{m_1} + \frac{1}{m_2} \right)^{-1} [(\bar{X}_1 - \bar{X}_2) - \delta_0]^T S_{\text{pooled}}^{-1} [(\bar{X}_1 - \bar{X}_2) - \delta_0] > \frac{(m_1 + m_2 - 2)p}{m_1 + m_2 - 1 - p} F_{1-\alpha}(p, m_1 + m_2 - 1 - p)$

### • Confidence Region

$$CR_{1-\alpha}(\beta_1 - \beta_2) = \left\{ \eta \in \mathbb{R}^p : \left( \frac{1}{m_1} + \frac{1}{m_2} \right)^{-1} [(\bar{X}_1 - \bar{X}_2) - \delta]^T S_{\text{pooled}}^{-1} [(\bar{X}_1 - \bar{X}_2) - \delta] < \frac{(m_1 + m_2 - 2)p}{m_1 + m_2 - 1 - p} F_{1-\alpha}(p, m_1 + m_2 - 1 - p) \right\}$$

- Simultaneous Confidence interval Sim CI<sub>1-α</sub>( $\beta_1 - \beta_2$ )  $\rightarrow \dots$
- Bonferroni Simultaneous Confidence interval Bonf CI<sub>1-α</sub>  $\rightarrow \dots$

### ► SPECIAL CASE BEHRENS-FISHER PROBLEM

$\Downarrow X_1, \dots, X_{1m_1} \stackrel{\text{iid}}{\sim} N_p(\mu_1, \Sigma_1)$   
 $X_2, \dots, X_{2m_2} \stackrel{\text{iid}}{\sim} N_p(\mu_2, \Sigma_2)$

$H_0: \Sigma_1 = \Sigma_2$  vs  $H_1: \Sigma_1 \neq \Sigma_2$  Here we don't want to reject!

- some tests based on Gaussianity (Extended Leven Test)
- Permutation tests (nonparametric)

If  $m_1, m_2$  are large enough we don't need Gaussianity

$$\Rightarrow \text{by CLT } \bar{X}_1 \sim N_p(\mu_1, \frac{1}{m_1} \Sigma_1) \Rightarrow \bar{X}_1 - \bar{X}_2 \sim N_p(\mu_1 - \mu_2, \frac{1}{m_1} \Sigma_1 + \frac{1}{m_2} \Sigma_2)$$

$$\bar{X}_2 \sim N_p(\mu_2, \frac{1}{m_2} \Sigma_2)$$

Therefore  $\left( \frac{1}{m_1} + \frac{1}{m_2} \right)^{-1} [(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)]^T \left[ \frac{1}{m_1} \Sigma_1 + \frac{1}{m_2} \Sigma_2 \right]^{-1} [(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)] \sim \chi^2(p)$

unknown but  
 $\Sigma_1 \rightarrow \Sigma_1$  in P and thus L  
 $\Sigma_2 \rightarrow \Sigma_2$

So our piv. quantity is

$$\left( \frac{1}{m_1} + \frac{1}{m_2} \right)^{-1} [(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)]^T \left[ \frac{1}{m_1} S_1 + \frac{1}{m_2} S_2 \right]^{-1} [(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)] \sim \chi^2(p)$$

## LECTURE 14 29/3/2022

► CASE  $g \geq 2, p=1$  (many groups, 1 treatment)  $\rightarrow$  ANOVA case, not multivariate (one-way)

$\hookrightarrow$  1 treatment

$\Downarrow \begin{cases} X_1, \dots, X_{1m_1} \sim N(\mu_1, \sigma^2) \\ \vdots \\ X_g, \dots, X_{gm_g} \sim N(\mu_g, \sigma^2) \end{cases}$

Goals:

- $H_0: \mu_1 = \dots = \mu_g$  (e.g. fertilizer has no effect) vs  $H_1: \exists i, j : \mu_i \neq \mu_j$  (e.g. has some effect)
- If reject  $H_0 \Rightarrow$  estimate the  $\mu_i$

New parametrization of the problem

$$\mu_i = \mu + \tau_i$$

overall mean      how diff. group  $i$  is from the mean

$g \mapsto g+1$  parameters ( $\mu, \tau_i$ )

$\Rightarrow$  we need one more constraint!  $\textcircled{2}$

Suppose  $m_1 + \dots + m_g = m$ ,  $\bar{X} = \frac{1}{m} \sum_{i=1}^g \sum_{j=1}^{m_i} X_{ij}$   $\mathbb{E}[\bar{X}] = \frac{1}{m} \sum_{i=1}^g \sum_{j=1}^{m_i} (\mu + \tau_i) = \frac{1}{m} \sum_{i=1}^g (m_i \mu + m_i \tau_i) = \mu + \frac{1}{m} \sum_{i=1}^g m_i \tau_i = \mu + \frac{1}{m} \sum_{i=1}^g m_i \tau_i$

$\bar{X}$  unbiased  $\Leftrightarrow \mathbb{E}[\bar{X}] = \mu \Leftrightarrow \sum_{i=1}^g m_i \tau_i = 0$   $\textcircled{1}$

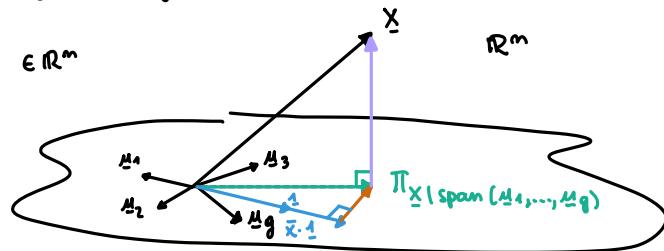
Est. for  $\tau_i$ ?  $\mathbb{E}[X_i - \bar{X}] = \mathbb{E}[\bar{X}_i] - \mu = \frac{1}{m_i} \sum_{j=1}^{m_i} \mathbb{E}[X_{ij}] - \mu = \mu + \tau_i - \mu = \tau_i$

mean of group  $i$       overall mean       $\mu + \tau_i$        $\mu + \tau_i$

## DECOMPOSITION OF VARIANCE

Let  $\mathbb{R}^m \ni \underline{x} = (x_{11}, x_{12}, \dots, x_{1m_1}, x_{21}, \dots, x_{2m_2}, \dots, x_{g1}, \dots, x_{gm_g})'$

$$\underline{\mu}_1 = \begin{pmatrix} 1 \\ 0 \\ 1 \\ 0 \end{pmatrix} \}_{M_1} \quad \underline{\mu}_2 = \begin{pmatrix} 0 \\ 1 \\ 0 \\ 1 \end{pmatrix} \}_{M_2} \quad \underline{\mu}_g = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{pmatrix} \}_{M_g}$$



### REMARK

- 1)  $\underline{\mu}_1, \dots, \underline{\mu}_g$  are lin. indep.
- 2)  $\underline{\mu}_i \cdot \underline{\mu}_j = m_i \delta_{ij} = m_j \delta_{ij} = \begin{cases} 0 & i \neq j \\ m_i & i = j \end{cases}$
- 3)  $\mathbb{R}^m \ni 1 \in \text{span}(\underline{\mu}_1, \dots, \underline{\mu}_g)$

$$\Pi_{\underline{x} \in \text{span}(\underline{\mu}_1, \dots, \underline{\mu}_g)} = \sum_{i=1}^g \frac{\underline{\mu}_i \cdot \underline{x}}{\underline{\mu}_i \cdot \underline{\mu}_i} \underline{\mu}_i = \sum_{i=1}^g \left( \frac{1}{m_i} \sum_{j=1}^{m_i} x_{ij} \right) \cdot \underline{\mu}_i = \sum_{i=1}^g \bar{x}_i \cdot \underline{\mu}_i$$

$$\Pi_{\underline{x} \in \mathbb{R}^m} = \underbrace{\bar{x} \cdot 1}_{\in \mathbb{R}}$$

$$\Pi_{\sum_{i=1}^g \bar{x}_i \cdot \underline{\mu}_i} = \frac{1}{1 \cdot 1} \left( \sum_{i=1}^g \bar{x}_i \cdot \underline{\mu}_i \right) = \sum_{i=1}^g \bar{x}_i \cdot \frac{1}{1 \cdot 1} \underline{\mu}_i = \frac{1}{m} \sum_{i=1}^g m_i \bar{x}_i \cdot 1 = \left( \frac{1}{m} \sum_{i=1}^g \sum_{j=1}^{m_i} x_{ij} \right) 1 = \bar{x} \cdot 1$$

$$\Rightarrow \underline{x} \cdot \bar{x} \cdot 1 + \sum_{i=1}^g (\bar{x}_i - \bar{x}) \underline{\mu}_i + \sum_{i=1}^g \sum_{j=1}^{m_i} (x_{ij} - \bar{x}_i) \underline{\mu}_i$$

all  
ORTHOGONAL!

$$\Rightarrow \|\underline{x}\|^2 = \|\bar{x} \cdot 1\|^2 + \left\| \sum_{i=1}^g (\bar{x}_i - \bar{x}) \underline{\mu}_i \right\|^2 + \left\| \underline{x} - \sum_{i=1}^g \bar{x}_i \underline{\mu}_i \right\|^2$$

$$\sum_{i=1}^g \sum_{j=1}^{m_i} x_{ij}^2 = \underbrace{\sum_{i=1}^g (\bar{x}_i - \bar{x})^2 m_i}_{SS_{\text{MEAN}}} + \underbrace{\sum_{i=1}^g (\bar{x}_i - \bar{x})^2 m_i}_{SS_{\text{TREATMENT}}} + \underbrace{\sum_{i=1}^g \sum_{j=1}^{m_i} (x_{ij} - \bar{x}_i)^2}_{SS_{\text{RESIDUALS}}}$$

Which is equivalent to

$$\sum_{i=1}^g \sum_{j=1}^{m_i} (x_{ij} - \bar{x})^2 = \underbrace{\sum_{i=1}^g (\bar{x}_i - \bar{x})^2 m_i}_{SS_{\text{CENTERED}}} + \underbrace{\sum_{i=1}^g \sum_{j=1}^{m_i} (x_{ij} - \bar{x}_i)^2}_{SS_{\text{TREATMENT}}} + \underbrace{\sum_{i=1}^g \sum_{j=1}^{m_i} (x_{ij} - \bar{x}_i)^2}_{SS_{\text{RESIDUALS}}}$$

$\hookrightarrow$  not captured by the lin.  
Comb. of  $\underline{\mu}_1, \dots, \underline{\mu}_g$  with  $\underline{x}$

IDEA: Reject  $H_0$  if  $\frac{SS_{\text{TREAT}}}{SS_{\text{RES}}} = \chi^2$  is large

If  $H_0$  is true,  $\chi^2 \sim ?$

- We know that  $\underline{x}$  is Gaussian  $\Rightarrow$  each term is Gaussian
- they are also orthogonal  $\Rightarrow$  they are independent
- In  $\chi^2$  we have the quotient of two squared moments  $\Rightarrow$  Fisher

$$\text{Focus on } \sum_{i=1}^g \sum_{j=1}^{m_i} (x_{ij} - \bar{x}_i)^2 = \sum_{i=1}^g (m_i - 1) S_i^2 \sim \sigma^2 \chi^2(m_i - 1) \sim \sigma^2 \chi^2(m-g)$$

$\hookrightarrow S_i^2 = \frac{1}{m_i - 1} \sum_{j=1}^{m_i} (x_{ij} - \bar{x}_i)$  est. of  $\sigma^2$   $\sim \sigma^2 \chi^2(m_i - 1)$  of group i

If  $H_0$  is true  $\Rightarrow \mu_1 = \mu_2 = \dots = \mu_g = \mu \Rightarrow \frac{1}{m-1} \sum_{i=1}^g \sum_{j=1}^{m_i} (x_{ij} - \bar{x})^2 = S^2$  is an estimator of  $\sigma^2$

$$\Rightarrow (m-1) S^2 = \sum_{i=1}^g \sum_{j=1}^{m_i} (x_{ij} - \bar{x}_i)^2 \sim \sigma^2 \chi^2(m-g)$$

Thus we can conclude:

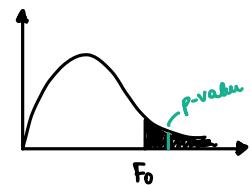
$$SS_{\text{centered}} = SS_{\text{TREAT}} + SS_{\text{RES}}$$

$$\sim \sigma^2 \chi^2(m-g-1) \sim \sigma^2 \chi^2(m-g)$$

$$F_0 = \frac{\sum_{i=1}^g m_i (\bar{X}_{ij} - \bar{X})^2 / (g-1)}{\sum_{i=1}^g \sum_{j=1}^{m_i} (\bar{X}_{ij} - \bar{X}_i)^2 / (m-g)} \sim F(g-1, m-g)$$

↓ pivotal

Reject  $H_0$  if  $F_0 > F_{1-\alpha}(g-1, m-g)$



### ► GENERAL CASE $p \geq 1, g \geq 2$

$$\mathbb{R}^p \ni X_{ij} = \mu + T_i + \varepsilon_{ij}$$

where

- $\varepsilon_{ij} \stackrel{iid}{\sim} N_p(0, \Sigma)$
- $\mu \in \mathbb{R}^p$
- $T_1, \dots, T_g \in \mathbb{R}^p$  s.t.  $\sum_{i=1}^g m_i T_i = 0$
- $\bar{\mu} = \frac{1}{m} \sum_{i=1}^g \sum_{j=1}^{m_i} X_{ij} = \bar{X}$  unbiased
- $\hat{T}_i = \frac{1}{m_i} \sum_{j=1}^{m_i} X_{ij} - \bar{X} = \bar{X}_i - \bar{X}$  unbiased

Goal:  $H_0: \mu_1 = \dots = \mu_g$  vs  $H_1: \exists i j \mu_i \neq \mu_j$

same as:  $H_0: T_1 = \dots = T_g = 0$  vs  $H_1: \exists T_i \neq 0$

Model

REMARK Fix  $k \in \{1, \dots, p\} \Rightarrow$  each  $X_{ijk} = \mu_k + T_{ik} + \varepsilon_{ijk} \stackrel{N_p(0, \Sigma_{kk})}{\sim}$  is an ANOVA model (salary)

Why don't we just "press the button"  $p$  times? CORRELATION!

Decomposition of covariance

$$\sum_{i=1}^g \sum_{j=1}^{m_i} (\bar{X}_{ij} - \bar{X})(\bar{X}_{ij} - \bar{X})' = \underbrace{\sum_{i=1}^g (\bar{X}_i - \bar{X})(\bar{X}_i - \bar{X})' m_i}_{\text{Cov}_B} + \underbrace{\sum_{i=1}^g \sum_{j=1}^{m_i} (\bar{X}_{ij} - \bar{X}_i)(\bar{X}_{ij} - \bar{X}_i)' m_i}_{\text{Cov}_W}$$

proof we  
 $\bar{X}_{ij} - \bar{X} = (\bar{X}_i - \bar{X}) + (\bar{X}_{ij} - \bar{X}_i)$

LECTURE 15 31/3/2022

By analogy with ANOVA we'd like consider  $B/W$  but we can't take the ratio of 2 matrices!

•  $B: \text{rank}(B) \leq S = \min(g-1, p)$  (If we have 3 groups,  $p=100$ ,  $B \in 100 \times 100$  with rank 2)

$$\begin{aligned} W &= \sum_{i=1}^g (m_i - 1) S_i \\ &\stackrel{\text{def}}{=} \frac{1}{m_i - 1} \sum_{j=1}^{m_i} (\bar{X}_{ij} - \bar{X}_i)(\bar{X}_{ij} - \bar{X}_i)' \text{ est. of } \Sigma \\ &\Downarrow \\ &\frac{1}{m-g} W = \frac{1}{m-g} \sum_{i=1}^g m_i S_i \text{ est. of } \Sigma \end{aligned}$$

• Wilks'  $\Lambda_W$  proposal: test statistic  $\Lambda_W = \frac{\det(W)}{\det(B+W)}$  large  $\Leftrightarrow B$  is not adding much to  $W$

⇒ Reject  $H_0$  if  $\Lambda_W$  is small

If  $H_0$  is true the distribution of  $\Lambda_W$  is known as long as

- $p \geq 1, g=2,3$
- $p=2, g \geq 1$
- as  $m \rightarrow \infty$ :  $\underbrace{-(m-1 - \frac{p+g}{2}) \log(\Lambda_W)}_{\text{reject H0 if } > \chi^2_\alpha(p(g-1))} \sim \chi^2(p(g-1))$  (Bartlett's)

Note that since we have a minus in front of the pivotal statistic, we are taking  $-\log$ : so we reject for small value of  $\Lambda_W$ , so we reject for values of  $-\log \Lambda_W$  which are big.



### Rejecting $H_0$

If  $H_0$  is rejected we need to compare confidence intervals for  $T_{ik} - T_{il}$  component-wisely using Bonferroni's Simultaneous  $i, k = 1, \dots, g$  to see when there was difference  $l = 1, \dots, p$

An estimator is  $\underbrace{(\bar{X}_{il} - \bar{X}_e)}_{\text{est of } \tau_{il}} - \underbrace{(\bar{X}_{ke} - \bar{X}_e)}_{\text{est of } \tau_{ke}} = \bar{X}_{il} - \bar{X}_{ke} \sim N(\tau_{il} - \tau_{ke}, \sigma_{\epsilon}^2 (\frac{1}{m_i} + \frac{1}{m_k}))$

↑ Not pivotal,  
estimate of  $\sigma_{\epsilon}^2$ :  $\frac{1}{m-g}$  W.E.

$p_{\epsilon}(g-1) \frac{1}{2}$  SBCI at level  $1-\alpha$

$$\text{SBCI}_{1-\alpha}(\tau_{il} - \tau_{ke}) = \left[ \bar{X}_{il} - \bar{X}_{ke} \pm t_{\frac{\alpha}{2}} \cdot \frac{s_e}{p_{\epsilon}(g-1)} \sqrt{\frac{N_{el}}{m-g} \left( \frac{1}{m_i} + \frac{1}{m_k} \right)} \right]$$

NB If  $p=1$  we get SBCI im ANOVA

### ► SHORT EXCURSUS ON TWO-WAY ANOVA

{ treatment 1:  $g$  levels  $\{1, 2, \dots, g\}$  factor 1  
treatment 2:  $b$  levels  $\{1, 2, \dots, b\}$  factor 2

We observe  $X_{ijk} \in \mathbb{R}$ :  
 $i = 1, \dots, g$   
 $j = 1, \dots, b$   
 $k = 1, \dots, m$  (balanced)  
↳ total sample size is  $m \cdot b \cdot g$

		treat 2							
		1	2	3	4	...	...	b	
treat 1	1								
	2			$X_{2ik}$					$X_{2ik}, k=1, \dots, m$
:									
g									

Model

$$X_{ijk} = \mu_{ij} + \varepsilon_{ijk} \quad \mu_{ij} = \mu + \tau_i + \beta_j + \gamma_{ij} \quad \begin{array}{l} \text{gb} \\ \text{1} \quad g \quad b \quad gb \\ \uparrow \quad \uparrow \quad \uparrow \quad \uparrow \\ \text{effect of I treat} \quad \text{effect of II treat} \quad \text{effect of interaction} \end{array}$$

$\stackrel{iid N(0, \sigma^2)}{\sim}$

Constraints

$$0 = \sum_{i=1}^g \tau_i \quad 0 = \sum_{j=1}^b \beta_j \quad \sum_{j=1}^b \gamma_{ij} = 0 \quad \underbrace{\sum_{i=1}^g \gamma_{ij} = 0}_{(g)} \quad \Rightarrow \text{in this way we don't overparametrize!}$$

$$(1) \qquad (1) \qquad \qquad (g) \qquad (b)$$

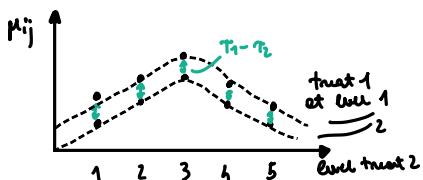
$$\sum_{j=1}^b \gamma_{ij} = 0 \quad \underbrace{(g+b-1)}_{(g+b-1)}$$

Estimators

$$\hat{\mu} = \bar{X} = \frac{1}{mgb} \sum_{ijk} X_{ijk} \quad \left. \begin{array}{l} \hat{\tau}_i = \bar{X}_{i..} - \bar{X} \\ \hat{\beta}_j = \bar{X}_{..j} - \bar{X} \\ \hat{\gamma}_{ij} = \bar{X}_{ij} - (\bar{X}_{i..} - \bar{X}) - (\bar{X}_{..j} - \bar{X}) = \bar{X}_{ij} - \bar{X}_{i..} - \bar{X}_{..j} + \bar{X} \\ \hat{\varepsilon}_{ijk} = \frac{1}{m} \sum_{ij} X_{ijk} : \boxed{\text{not}} \end{array} \right\}$$

### REMARK

Suppose  $\gamma_{ij} = 0 \quad \forall i, j \Rightarrow X_{ijk} = \underbrace{\mu + \tau_i + \beta_j}_{\mu_{ij}} + \varepsilon_{ijk}$

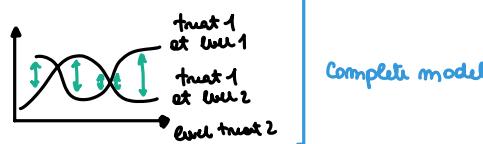


Additive model

$$\mu_{1j} = \mu + \tau_1 + \beta_j$$

$$\mu_{2j} = \mu + \tau_2 + \beta_j$$

If there are interactions  
TEST FOR INTERACTIONS!



Complete model

- then:
- (A)  $H_0: \gamma_{ij} = 0 \quad \forall i, j$  vs  $H_1: \exists \gamma_{ij} \neq 0$
  - (B)  $H_0: \tau_i = 0 \quad \forall i$  vs  $H_1: \exists \tau_i \neq 0$
  - (C)  $H_0: \beta_j = 0 \quad \forall j$  vs  $H_1: \exists \beta_j \neq 0$

Decomposition of variance:

$\sum_{i=1}^g \sum_{j=1}^b \sum_{k=1}^m (\bar{x}_{ijk} - \bar{x})^2$	SS <sub>centered</sub>	Dof
$\sum_{i=1}^g (\bar{x}_{i\cdot} - \bar{x})^2 b m$	SS <sub>treat1</sub>	$g-1$
$\sum_{j=1}^b (\bar{x}_{\cdot j} - \bar{x})^2 g m$	SS <sub>treat2</sub>	$b-1$
$\sum_{i=1}^g \sum_{j=1}^b (\bar{x}_{ij\cdot} - \bar{x}_{i\cdot} - \bar{x}_{\cdot j} + \bar{x})^2 m$	SS <sub>interaction</sub>	$(g-1)(b-1)$
$\sum_{i=1}^g \sum_{j=1}^b \sum_{k=1}^m (\bar{x}_{ijk} - \bar{x}_{ij\cdot})^2$	SS <sub>residuals</sub>	$gb(m-1)$

if  $m=1$  (possibly) overfitting! To take into account variability we want  $m \geq 2$ . As  $m$  is greater, the "meter" is more precise

(A) Reject at level  $\alpha$  if  $\frac{\frac{1}{(g-1)(b-1)} SS_{interaction}}{\frac{1}{gb(m-1)} SS_{residuals}} > F_{1-\alpha}(1, (g-1)(b-1), gb(m-1))$

If  $H_0$  not rejected ( $H_0$  true)

the model is  $X_{ijk} = \mu + \tau_i + \beta_j + \epsilon_{ijk}$  additive! (simpler)  $\Rightarrow$  in this case  $SS_{interaction}$  is neglected and becomes part of  $SS_{res}$  (also the Dof becomes  $gb(m-1) + (g-1)(b-1)$ )

(B) Suppose the model is additive

(C similar) Reject at level  $\alpha$  if  $\frac{SS_{treat1}/(g-1)}{(SS_{res} + SS_{inter})/(gb(m-1) + (g-1)(b-1))} > F_{1-\alpha}(g-1, gb(m-1) - b - g + 1)$

Suppose the model is complete

Reject at level  $\alpha$  if  $\frac{SS_{treat1}/(g-1)}{(SS_{res})/(gb(m-1))} > F_{1-\alpha}(g-1, gb(m-1))$

## LECTURE 16 1/4/2022

### CLASSIFIERS

Each unit is represented by  $(x_1, \dots, x_p)^T \in \mathcal{X}$  (e.g.  $\mathbb{R}^p$ ) with features  $L \in \{1, \dots, g\}$

Goal: Learn  $\delta: \mathcal{X} \rightarrow \{1, \dots, g\}$

### Supervised classification

Training set

$$\mathbb{X} = \begin{bmatrix} x_{11} & \dots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_{m1} & \dots & x_{mp} \end{bmatrix} \begin{bmatrix} l_1 \\ \vdots \\ l_m \end{bmatrix} \quad \begin{array}{l} x_{ij} \in \mathbb{R} \\ l_i \in \{1, \dots, g\} \quad \forall i \end{array}$$

Learn  $\delta$  by using:

- training set
- model (discriminant analysis)

### Unsupervised classification

Training set

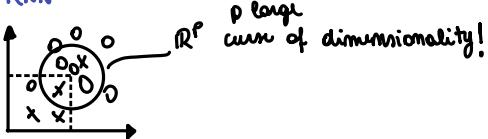
$$\mathbb{X} = \begin{bmatrix} x_{11} & \dots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_{m1} & \dots & x_{mp} \end{bmatrix} \begin{bmatrix} l_1 \\ \vdots \\ l_m \end{bmatrix} \quad \begin{array}{l} x_{ij} \in \mathbb{R} \\ l_i \in \{1, \dots, g\} \quad \forall i \end{array}$$

Estimate  $\hat{g}$  and  $\hat{l}_1, \dots, \hat{l}_m$  (cluster analysis)

then  $\Rightarrow$  supervised classification (i.e. find  $\delta: \mathcal{X} \rightarrow \{1, \dots, \hat{g}\}$ )

You assume they exist and they are unknown (with the labels nor how many groups)

## Ex KNN



## Ex Logistic regression

### Supervised classification

Ingredients:

- $\underline{x} \mid L=i \sim f_i(\underline{x})$  density on  $\mathbb{R}^p$

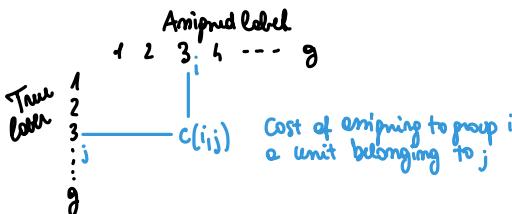
they need to be different

- prior distribution

$$p_i = P[L=i] \quad i = \{1, \dots, g\}$$

$$p_i \geq 0 \quad \sum_{i=1}^g p_i = 1$$

- Costs of misclassification



$$\begin{aligned} c(i,j) &\geq 0 \\ c(i,i) &= 0 \\ c(i,j) &\neq c(j,i) \text{ in general} \end{aligned}$$

REMARK  $\delta: \mathbb{R}^p \rightarrow \{1, 2, \dots, g\}$  is a class

$$R_i = \{\underline{x} \in \mathbb{R}^p : \delta(\underline{x}) = i\} = \delta^{-1}(i)$$

$$\{R_1, \dots, R_g\} \leftrightarrow \delta$$

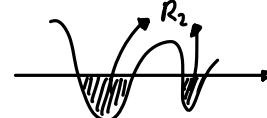
↳ partition of  $\mathbb{R}^p$

Optimality criterion  $\delta$  is optimal if it minimizes  $\underbrace{\mathbb{E}[\text{cost of misclassification}]}_{\text{ECM}} = \text{ECM}$

## Ex $g=2$

$$\underline{x} \in \mathbb{R}^p \quad \delta(\underline{x}) = \begin{cases} 1 & \text{if } \underline{x} \in R_1 \\ 2 & \text{if } \underline{x} \in R_2 \end{cases}$$

$$\begin{aligned} \text{ECM}(\delta) &= \int_{R_1} C(1,2) f_2(\underline{x}) p_2 d\underline{x} + \int_{R_2} C(2,1) f_1(\underline{x}) p_1 d\underline{x} \\ &= \int_{\mathbb{R}^p} C(1,2) f_2(\underline{x}) p_2 d\underline{x} - \int_{R_2} C(1,2) f_2(\underline{x}) p_2 d\underline{x} + \int_{R_2} C(2,1) f_1(\underline{x}) p_1 d\underline{x} \\ &= C(1,2) p_2 + \int_{R_2} [C(2,1) f_1(\underline{x}) p_1 - C(1,2) f_2(\underline{x}) p_2] d\underline{x} \end{aligned}$$



$$R_2 = \{\underline{x} \in \mathbb{R}^p \mid C(2,1) f_1(\underline{x}) p_1 \leq C(1,2) f_2(\underline{x}) p_2\}$$

$$R_1 = \{\underline{x} \in \mathbb{R}^p \mid C(2,1) f_1(\underline{x}) p_1 > C(1,2) f_2(\underline{x}) p_2\}$$

↳ the equal sign doesn't matter

I will classify to group 2 all units such that if I make a mistake I pay less than what I would pay if I were to classify the same units in the other group and make a mistake

General case:  $g \geq 2$

$$\text{ECM}(\delta) = \int_{R_1} \sum_{k \neq 1} C(1|k) f_k(\underline{x}) p_k d\underline{x} + \cdots + \overbrace{\int_{R_i} \sum_{k \neq i} C(i|k) f_k(\underline{x}) p_k d\underline{x} + \cdots + \int_{R_g} \sum_{k \neq g} C(g|k) f_k(\underline{x}) p_k d\underline{x}}$$

$$\hat{\delta} = \underset{\delta}{\text{argmin}} \text{ ECM}(\delta) \quad \hat{\delta} \in \{\hat{R}_1, \dots, \hat{R}_g\}$$

$$\hat{R}_1 = \left\{ \underline{x} \in \mathbb{R}^p : \sum_{k \neq 1} c(i|k) f_k(\underline{x}) p_k \leq \sum_{k \neq j} c(j|k) f_k(\underline{x}) p_k, \forall j \in \{2, \dots, g\} \right\}$$

$$\hat{R}_i = \left\{ \underline{x} \in \mathbb{R}^p : \sum_{k \neq i} c(i|k) f_k(\underline{x}) p_k \leq \sum_{k \neq j} c(j|k) f_k(\underline{x}) p_k, \forall j \in \{1, \dots, g\} \setminus \{i\} \right\}$$

(w constant)  
doesn't change!

REMARK  $\underline{x} \in R_i \Leftrightarrow \frac{\sum_{k \neq i} c(i|k) f_k(\underline{x}) p_k}{\sum_{t=1}^g f_t(\underline{x}) p_t} \leq \frac{\sum_{k \neq j} c(j|k) f_k(\underline{x}) p_k}{\sum_{t=1}^g f_t(\underline{x}) p_t} \quad \forall j \in \{1, \dots, g\} \setminus \{i\} \Leftrightarrow \textcircled{*}$

$$\frac{\frac{P(X=\underline{x}|L=k)P(L=k)}{\sum_{t=1}^g P(X=\underline{x}|L=t)P(L=t)}}{P(X=\underline{x})} = \frac{P(X=\underline{x}|L=k)P(L=k)}{P(X=\underline{x})} \stackrel{\text{Bayes}}{=} P(L=k|X=\underline{x})$$

$$\textcircled{*} \Leftrightarrow \sum_{k \neq i} c(i|k) P(L=k|X=\underline{x}) \leq \sum_{k \neq j} c(j|k) P(L=k|X=\underline{x}) \quad \forall j \in \{1, \dots, g\} \setminus \{i\}$$

### Special cases

1) All costs are equal:

$$c(i|j) = w > 0 \quad \forall i \neq j \quad (\text{default for a computer})$$

$$c(i|i) = 0$$

$$\Rightarrow \underline{x} \in R_i \Leftrightarrow \sum_{k \neq i} P(L=k|\underline{x}=\underline{x}) \leq \sum_{k \neq j} P(L=k|\underline{x}=\underline{x}) \quad \forall j \neq i$$

$$\Leftrightarrow 1 - P[L=i|\underline{x}=\underline{x}] \leq 1 - P[L=j|\underline{x}=\underline{x}] \quad \forall j \neq i$$

$$\Leftrightarrow P[L=i|\underline{x}=\underline{x}] \geq P[L=j|\underline{x}=\underline{x}] \quad \forall j \neq i \quad (\text{BAYES CLASSIFIER})$$

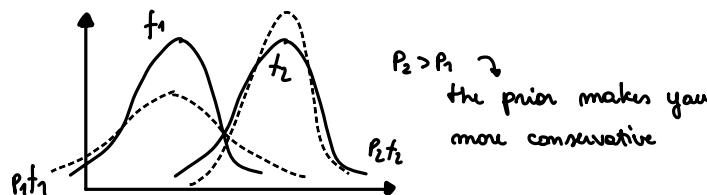
2)  $c(i|j) = w > 0 \quad \textcircled{*} \quad p_1 = p_2 = \dots = p_g = \frac{1}{g}$  (equally indifferent before looking at the date  
 $c(i|i) = 0$  of which class the unit belongs to)

$$\Rightarrow \underline{x} \in R_i \Leftrightarrow P(L=i|\underline{x}=\underline{x}) \geq P(L=j|\underline{x}=\underline{x})$$

↳ Abused in practice!

$$\Leftrightarrow \frac{f_i(\underline{x}) p_i}{\sum_{t \neq i} f_t(\underline{x}) p_t} \geq \frac{f_j(\underline{x}) p_j}{\sum_{t \neq j} f_t(\underline{x}) p_t}$$

$$\Leftrightarrow f_i(\underline{x}) \geq f_j(\underline{x}) \quad \forall j \neq i$$



### LECTURE 17 5/4/2022

### 3) QDA (Quadratic Discriminant analysis)

Hypothesis:

$$\cdot c(i|j) = w > 0 \quad \forall j \neq i$$

$$\rightarrow \cdot f_i(\underline{x}) = \frac{1}{\sqrt{(2\pi)^p |\Sigma_i|}} \exp \left\{ -\frac{1}{2} (\underline{x} - \mu_i)' \Sigma_i^{-1} (\underline{x} - \mu_i) \right\} \quad \underline{x} \in \mathbb{R}^p, \Sigma_i \text{ pos. def.}, \mu_i \in \mathbb{R}^p$$

$$\underline{x} \in R_i \Leftrightarrow f_i(\underline{x}) p_i \geq f_j(\underline{x}) p_j \quad \forall j \neq i$$

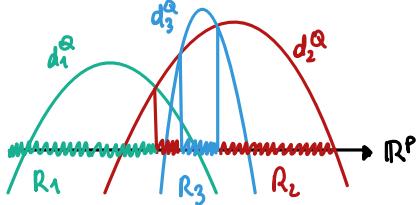
$$\Leftrightarrow p_i \frac{1}{\sqrt{(2\pi)^p |\Sigma_i|}} \exp \left\{ -\frac{1}{2} (\underline{x} - \mu_i)' \Sigma_i^{-1} (\underline{x} - \mu_i) \right\} \geq p_j \frac{1}{\sqrt{(2\pi)^p |\Sigma_j|}} \exp \left\{ -\frac{1}{2} (\underline{x} - \mu_j)' \Sigma_j^{-1} (\underline{x} - \mu_j) \right\} \quad \forall j \neq i$$



$$\Leftrightarrow \log p_i - \frac{1}{2} \log |\Sigma_i| - \frac{1}{2} \underbrace{(\bar{x} - \mu_i)' \Sigma_i^{-1} (\bar{x} - \mu_i)}_{\text{MALANOBIS DISTANCE}} \geq \log p_j - \frac{1}{2} \log |\Sigma_j| - \frac{1}{2} (\bar{x} - \mu_j)' \Sigma_j^{-1} (\bar{x} - \mu_j) \quad b_j + i$$

$$= d_j^Q(\bar{x}) \quad \text{QUADRATIC DISCRIMINANT FUNCTIONS}$$

$$\Leftrightarrow d_i^Q(x) \geq d_j^Q(x) \quad \forall j \neq i$$



#### 4) LDA (Linear Discriminant Analysis)

- $c(i|j) = w > 0 \quad \forall j \neq i$
  - $f_i(x)$  is  $N_p(\mu_i, \Sigma)$   SAME FOR ALL  $i = 1, \dots, g$  GROUPS

$$x \in R_i \Leftrightarrow f_i(x)p_i \geq f_j(x)p_j \quad \forall j \neq i$$

$$\Leftrightarrow \log p_i - \frac{1}{2} \log |\Sigma_i| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)' \Sigma_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) \geq \log p_j - \frac{1}{2} \log |\Sigma_j| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_j)' \Sigma_j^{-1} (\mathbf{x} - \boldsymbol{\mu}_j) \quad \forall j \neq i$$

$$\Leftrightarrow \log p_i - \frac{1}{2} x^T \Sigma^{-1} x + x^T \Sigma^{-1} \mu_i - \frac{1}{2} \mu_i^T \Sigma^{-1} \mu_i \geq \log p_j - \frac{1}{2} x^T \Sigma^{-1} x + x^T \Sigma^{-1} \mu_j - \frac{1}{2} \mu_j^T \Sigma^{-1} \mu_j \quad \forall j \neq i$$

no more quadratic!

$$\Leftrightarrow \underline{x}^T \Sigma^{-1} \mu_i + \log p_i - \frac{1}{2} \mu_i^T \Sigma^{-1} \mu_i \geq \underline{x}^T \Sigma^{-1} \mu_j + \log p_j - \frac{1}{2} \mu_j^T \Sigma^{-1} \mu_j \quad \forall j \neq i$$

$$x \in R_i \Leftrightarrow d_i(x) \geq d_j(x) \quad \forall j \neq i$$

VERY ROBUST!

**IMPORTANT** Gaussianity is not much important, instead focus on modifying data to have the same covariance matrix for all groups.

- Where is DATA and ? To estimate the unknown  $\mu_i, \Sigma_i, f_i(x)$   $i=1, \dots, g$  (training data)

$$\mathbb{X} = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix} \begin{bmatrix} l_1 \\ \vdots \\ l_m \end{bmatrix}$$

**WARNING** Don't use date to estimate  $\pi_i$ ,  $i=1, \dots, g$  unless you know what you're doing.

⇒ ok if sample is "random and large" enough to estimate the frequencies, however things with small frequency will not be many (dov.) and the classifier won't have a lot of information to learn how to classify them!

In LDA and QDA estimate  $\mu_i$  with  $\bar{x}_i = \frac{1}{m_i} \sum_{j:j_i=i} x_j$  (both)

$$\Sigma_i \text{ with } S_i = \frac{1}{m_i - 1} \sum_{j: l_j = i} (x_j - \bar{x}_i)(x_j - \bar{x}_i)^T \quad (\text{QDA})$$

$$\sum \text{ with } Spooled = \frac{1}{m-g} \sum_{i=1}^g (M_{i-1}) S_i \quad (LDA)$$

If data in each group are not many with respect to p  $\Rightarrow$  parametrize  $\Sigma_i, i=1, \dots, g$

## Ex (Naive Bayes classifier)

Consider QDA and assume:  $\Sigma_i = \begin{bmatrix} \sigma_{11}^{(i)} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_{pp}^{(i)} \end{bmatrix}$  indep components

$$d_i^{\text{LDA}}(x) = \log p_i - \frac{1}{2} \underbrace{\log |\Sigma_i|}_{\log \sum_{i=1}^g \pi_{ii}} - \frac{1}{2} (\underline{x} - \mu_i)' \Sigma_i^{-1} (\underline{x} - \mu_i)$$

$$\sum_{j=1}^g \frac{(x_j - \bar{x}_{ij})^2}{\pi_{jj}}$$

$\sigma_{ij}$  is estimated with  $\frac{1}{m_i-1} \sum_{j: l_j=1} (x_{ij} - \bar{x}_{ij})^2$

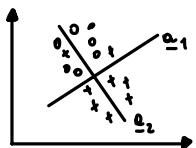
REMARK PCA and then classification is a bad recipe

### Fisher's argument for LDA

► Remove Gaussianity assumption!

Covariance between groups  $B = \frac{1}{g-1} \sum_{i=1}^g (\mu_i - \bar{\mu})(\mu_i - \bar{\mu})'$   $\bar{\mu} = \frac{1}{g} \sum_{i=1}^g \mu_i$

Covariance within groups



GOAL:

Find  $\underline{a} \in \mathbb{R}^p$  that maximizes Variability between groups

$$\text{Reparametrize the problem } \underline{u} = \Sigma^{1/2} \underline{a} \quad (\underline{a} = \Sigma^{-1/2} \underline{u}) \Rightarrow \frac{\underline{a}' B \underline{a}}{\underline{a}' \Sigma \underline{a}} = \frac{\underline{u}' \Sigma^{-1/2} B \Sigma^{-1/2} \underline{u}}{\underline{u}' \Sigma^{-1/2} \Sigma^{-1/2} \underline{u}} = \frac{\underline{u}' \Sigma^{-1/2} B \Sigma^{-1/2} \underline{u}}{\underline{u}' \underline{u}}$$

$$\underset{\underline{u} \in \mathbb{R}^p}{\operatorname{argmax}} \frac{\underline{u}' \Sigma^{-1/2} B \Sigma^{-1/2} \underline{u}}{\underline{u}' \underline{u}} = \underline{e}_1 \Rightarrow \underline{a}_1 = \Sigma^{-1/2} \underline{e}_1$$

$$\text{If } \Sigma^{-1/2} B \Sigma^{-1/2} = \sum_{i=1}^g \lambda_i \underline{e}_i \underline{e}_i' \quad S = \operatorname{rank} B = \min \{g-1, p\}$$

$$\text{on im PCA} \quad \underline{a}_1 = \Sigma^{-1/2} \underline{e}_1 \rightarrow \underline{a}_1' \underline{x} \quad \text{1st discriminant score (Fisher's score)} \\ \vdots \quad \vdots \\ \underline{a}_S = \Sigma^{-1/2} \underline{e}_S \quad \underline{a}_S' \underline{x} \quad \text{sth} \quad "$$

$$A = \begin{bmatrix} \underline{a}_1' \\ \vdots \\ \underline{a}_S' \end{bmatrix} \Rightarrow \operatorname{Cov}(\underline{a}_i' \underline{x}, \underline{a}_j' \underline{x}) = \underline{a}_i' \Sigma \underline{a}_j = \underline{e}_i' \Sigma^{-1/2} \Sigma \Sigma^{-1/2} \underline{e}_j = \underline{e}_i' \underline{e}_j = \begin{cases} 0 & i \neq j \\ 1 & i = j \end{cases} \Rightarrow \operatorname{Cov}(A \underline{x}) = I$$

$$\underline{x} \xrightarrow{A} \begin{bmatrix} \underline{a}_1' \underline{x} \\ \vdots \\ \underline{a}_k' \underline{x} \end{bmatrix} \left. \right\} \text{take } k \quad k=2,3 \quad (\text{for making pictures})$$

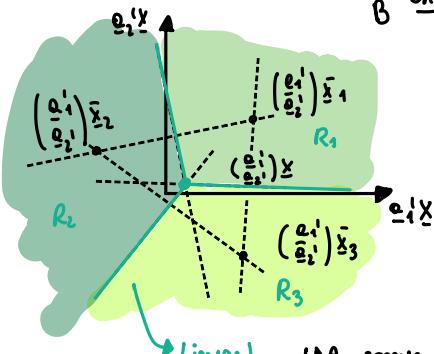
► How to build a classifier out of this?

Training data

$$\mathbb{X} = \begin{bmatrix} \underline{x}_1' \\ \vdots \\ \underline{x}_m' \end{bmatrix} \begin{bmatrix} \underline{e}_1 \\ \vdots \\ \underline{e}_m \end{bmatrix} \quad \mu_i \xrightarrow{\text{est.}} \bar{\underline{x}}_i$$

$$\Sigma \xrightarrow{\text{est.}} \text{Spooler} = \frac{1}{m-g} \sum_{i=1}^g (m_i-1) S_i$$

$$B \xrightarrow{\text{est.}} \hat{B} = \frac{1}{g-1} \sum_{i=1}^g m_i (\bar{\underline{x}}_i - \bar{\underline{x}})(\bar{\underline{x}}_i - \bar{\underline{x}})' \quad \bar{\underline{x}} = \frac{1}{m} \sum \bar{\underline{x}}_i$$



$$\underline{x} \in R_i \Leftrightarrow \sum_{j=1}^k (\underline{a}_j' \underline{x} - \underline{a}_j' \bar{\underline{x}}_i)^2 \leq \sum_{j=1}^k (\underline{a}_j' \underline{x} - \underline{a}_j' \bar{\underline{x}}_t)^2 \quad \forall j \neq i$$

[Logistic regression useful when distrib. is unknown]

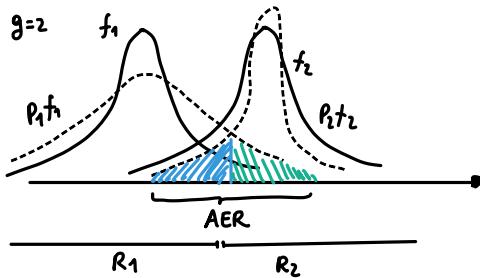
Linear! LDA comes out even without Gaussianity

## EVALUATING A CLASSIFIER

Let  $\delta: \mathbb{R}^p \rightarrow \{1, \dots, g\}$  be a classifier

$$AER(\delta) = \sum_{k=1}^g \int_{R_k} f_k(x) p_k dx + \dots + \sum_{k=g}^g \int_{R_k} f_k(x) p_k dx$$

ACTUAL ERROR RATE  $\rightarrow$  unknown ( $f, p$ )



## NON-PARAMETRIC APPROACH TO ESTIMATE AER

- Naive approach

$$\mathbb{X} = \begin{bmatrix} \mathbb{x}_1 \\ \vdots \\ \mathbb{x}_m \end{bmatrix} \begin{bmatrix} l_1 \\ \vdots \\ l_m \end{bmatrix} \quad \delta(\mathbb{x}_i) = \hat{l}_i \quad \varepsilon_i = \begin{cases} 0 & l_i = \hat{l}_i \\ 1 & l_i \neq \hat{l}_i \end{cases} \Rightarrow APER = \frac{1}{m} \sum_{i=1}^m \varepsilon_i \quad \text{APPARENT ERROR RATE (TRAINING ERROR)}$$

## CONFUSION MATRIX ( $g=2$ )

		True	
		1	2
Predicted	1	$m_{11}$	$m_{12}$
	2	$m_{21}$	$m_{22}$

APER =  $\frac{m_{12} + m_{21}}{m}$

→ useful to see the weak points of the classifier

- Use a test set

$$\tilde{\mathbb{X}} = \begin{bmatrix} \tilde{\mathbb{x}}_1 \\ \vdots \\ \tilde{\mathbb{x}}_m \end{bmatrix} \begin{bmatrix} \tilde{l}_1 \\ \vdots \\ \tilde{l}_m \end{bmatrix} \quad \text{Apply } \delta \text{ (learnt from } \mathbb{X} \text{) to } \tilde{\mathbb{X}} \quad \delta(\tilde{\mathbb{x}}_i) = \hat{l}_i \quad \varepsilon_i = \begin{cases} 0 & \hat{l}_i = \tilde{l}_i \\ 1 & \hat{l}_i \neq \tilde{l}_i \end{cases} \quad \widehat{AER}(\delta) = \frac{1}{N} \sum_{i=1}^N \varepsilon_i$$

## In between? CROSS-VALIDATION

- L10 - Leave one out

$$\mathbb{X} = \begin{bmatrix} \mathbb{x}_1 \\ \mathbb{x}_2 \\ \vdots \\ \mathbb{x}_m \end{bmatrix} \begin{bmatrix} l_1 \\ l_2 \\ \vdots \\ l_m \end{bmatrix}$$

FOR  $i = 1, \dots, m$   
 $\mathbb{X}_{-i}$  = without datum  $i$   
 learn a classifier  $\delta_{-i}$  from  $\mathbb{X}_{-i}$   
 test  $\delta_{-i}$  on  $i$     $\varepsilon_i = \begin{cases} 0 & \delta_{-i}(\mathbb{x}_i) = l_i \\ 1 & \delta_{-i}(\mathbb{x}_i) \neq l_i \end{cases}$   
 END

$$\widehat{AER}(\delta) = \frac{1}{m} \sum_{i=1}^m \varepsilon_i$$

Small bias  $\Downarrow$

High variance  $\Updownarrow$

it's a mean of NON-INDEPENDENT measurements

- k-fold cross validation

( $k \leq m$ , im "practically"  $k=5, 10$  (?)

you take out more data  $\Rightarrow k$  parts



$$ER_i = \frac{1}{m_i} \sum_j \varepsilon_j$$

$\hookrightarrow$  # of  $j$  s.t. they belong to part- $i$

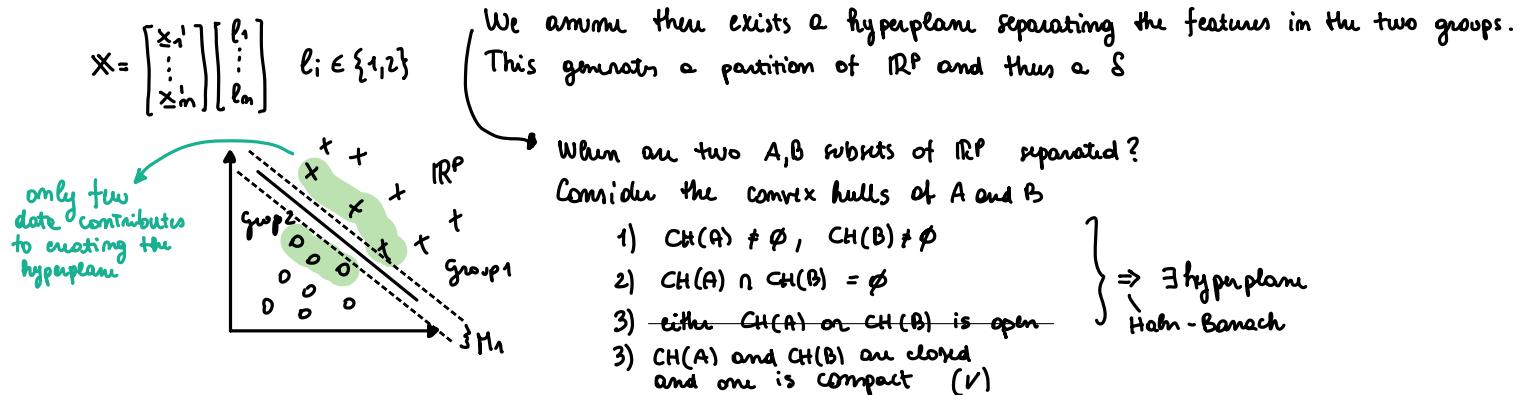
$$\widehat{AER}(\delta) = \frac{1}{M} \sum_{i=1}^k m_i ER_i$$

OBS  $k=m \Rightarrow L10$

Before splitting  $\Rightarrow$  Permute data! To avoid building classifiers of linked variables (depending if the original dataset was ordered)  
 ↳  $m!$  permutations  
 ↳  $(k=m)$ , run  $d$  permutations and get different estimates  $\widehat{AER}_1(S), \dots, \widehat{AER}_d(S) \rightarrow$  Mean and confidence interval 2-test

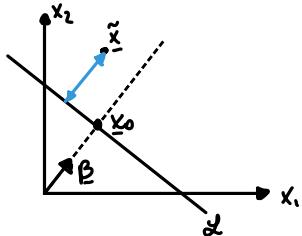
## Short exercises on separating Hyperplanes and SVM (Supported Vector Machine) ["Elements of Statistical Learning"]

Idea ( $g=2$ )



Hyperplanes

$\mathcal{L}$  hyperplane in  $\mathbb{R}^p$ : affine subspace of dim  $p-1$



Let  $\beta \perp \mathcal{L}, \|\beta\| = 1, \beta \in \mathbb{R}^p \quad \underline{x}_0 \in \text{span}(\beta) \cap \mathcal{L} \quad \beta_0 = \|\underline{x}_0\|$

$$\underline{x} \in \mathcal{L} \Leftrightarrow \Pi_{\underline{x} \mid \beta} = \underline{x}_0 \Leftrightarrow \underbrace{\underline{x}' \beta}_{x_1 \beta_1 + \dots + x_p \beta_p} = \beta_0$$

$$y_i := \begin{cases} 1 & \underline{x}_i \in \text{Group 1} \\ -1 & \underline{x}_i \in \text{Group 2} \end{cases}$$

$$y_i (\underline{x}_i \beta - \beta_0) \geq 0$$

$\underline{x} \notin \mathcal{L} : \quad \underline{x}' \beta > \beta_0 \quad \underline{x}' \beta - \beta_0$  is the distance (with sign) between  $\underline{x}$  and  $\mathcal{L}$

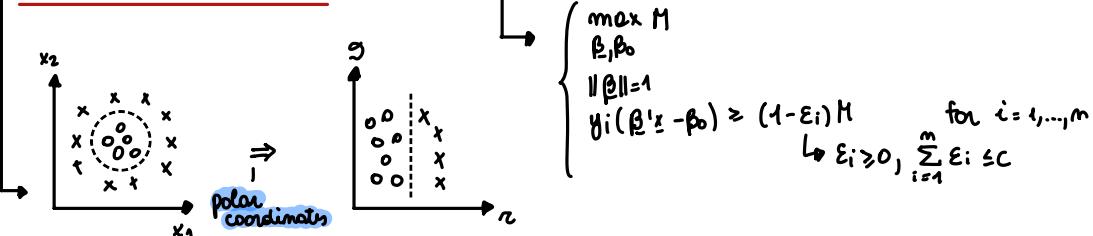
Sol: see ref.

$$M_1 = \min \{ y_i (\underline{x}_i \beta - \beta_0) : i = 1, \dots, m \} \Rightarrow \text{Find } \beta \in \mathbb{R}^p \text{ s.t. } \boxed{\beta = \arg \max_{\|\beta\|=1, \beta \in \mathbb{R}^p} M_1}$$

What if group 1 and 2 are not separated by a hyperplane?

- Increase the dimension  $p$  of the embedding space  $\rightarrow$  the curse of dimensionality plays in our favour
- Solve the "soft" problem allowing some overlapping

### LECTURE 19 8/4/2022



### UNSUPERVISED CLASSIFICATION (CLUSTER ANALYSIS)

► You don't know the labels nor how many they are

► Goal: - Estimate  $g$   
 - Estimate  $\{l_1, \dots, l_m\}$   $\Rightarrow$  at the end you have a training set

$$\mathbb{X} = \begin{bmatrix} \underline{x}_1 \\ \vdots \\ \underline{x}_m \end{bmatrix} \begin{bmatrix} \hat{l}_1 \\ \vdots \\ \hat{l}_m \end{bmatrix}$$

► Two lines of attack

- parametric modeling

→ based on ML (Maximum Likelihood)

→ EM

- nonparametric approach based on distances

► Idea: two units belonging to the same cluster are more **similar** than two units belonging to different clusters.

$$d: \mathbb{R}^p \times \mathbb{R}^p \rightarrow [0, +\infty)$$

$$d(\underline{x}, \underline{x}) = 0 \quad \forall \underline{x} \in \mathbb{R}^p$$

$$d(\underline{x}, \underline{y}) = 0 \Leftrightarrow \underline{x} = \underline{y} \quad (\text{stronger, also } \Rightarrow)$$

pseudo-metric

$$d(\underline{x}, \underline{y}) = d(\underline{y}, \underline{x}) \quad \forall \underline{x}, \underline{y} \in \mathbb{R}^p$$

$$d(\underline{x}, \underline{y}) \leq d(\underline{x}, \underline{z}) + d(\underline{z}, \underline{y}) \quad \forall \underline{x}, \underline{y}, \underline{z} \in \mathbb{R}^p$$

$$(d(\underline{x}, \underline{y}) \leq \max \{d(\underline{x}, \underline{z}) + d(\underline{z}, \underline{y})\})$$

distance / metric

ultra-metric

Example of metrics

→ Better to standardize!

Euclidean

$$- d(\underline{x}, \underline{y}) = \sqrt{(\underline{x}-\underline{y})'(\underline{x}-\underline{y})}$$

$$- d(\underline{x}, \underline{y}) = \sqrt{(\underline{x}-\underline{y})' Z^{-1} (\underline{x}-\underline{y})}$$

$$- d(\underline{x}, \underline{y}) = \left( \sum_{i=1}^p |x_i - y_i|^{1/m} \right)^{1/m}$$

$$- d(\underline{x}, \underline{y}) = \sum_{i=1}^p \frac{|x_i - y_i|}{x_i + y_i} \quad \underline{x}, \underline{y} \in (\mathbb{R}^+)^p \quad \text{Canberra}$$

Minkowski ( $\Sigma$  usually unknown)

$\ell^m$  distances  $\rightarrow m=2$  Euclidean

$m=1$  Manhattan

$$\underline{x}, \underline{y} \in \{0, 1\}^p$$

$$\underline{x} = (0 \ 1 \ 1 \ 0 \ 0 \ 0 \ 1)$$

$$d_{\text{Euclid.}}(\underline{x}, \underline{y}) = \sqrt{\sum_{i=1}^p (x_i - y_i)^2} = \sqrt{\# \text{ discordances}}$$

$= \sqrt{c+b}$

		$\underline{y}$
	0	1
x	a	b

$\Rightarrow$  we can build a table

$(1, 0) \rightarrow \text{blue}$

$(0, 1) \rightarrow \text{green}$

$(0, 0) \rightarrow \text{not green nor blue} \rightarrow \text{brown}$

You can transform categorical variables

Ex  $x \in \{\text{blue, green, brown}\} \Rightarrow x_1 <_0^1 \text{blue} \quad x_2 <_0^1 \text{green} \quad x_3 <_0^1 \text{brown}$

Ex Mix:  $d = \lambda d^{\text{quant}} + (1-\lambda)d^{\text{cat.}}$

$$\text{Receive } \underline{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix} \xrightarrow{d} D = \begin{bmatrix} 0 & & & \\ d_{21} & 0 & & \\ d_{31} & d_{23} & 0 & \dots \\ \vdots & & & \end{bmatrix}$$

Distance matrix

How to measure the distance between SUBSETS?

Let  $U, V$  be finite subsets of  $\mathbb{R}^p$

$$- d(U, V) = \min \{ d(\underline{x}, \underline{y}) : \underline{x} \in U, \underline{y} \in V \}$$

SINGLE LINKAGE ( $d(\text{Italy, Switzerland}) = 0$ )

$$- d(U, V) = \max \{ d(\underline{x}, \underline{y}) : \underline{x} \in U, \underline{y} \in V \}$$

COMPLETE LINKAGE

$$- d(U, V) = \frac{1}{|U|+|V|} \sum_{\substack{\underline{x} \in U \\ \underline{y} \in V}} d(\underline{x}, \underline{y}) \quad \text{AVERAGE LINKAGE}$$

$$- d(U, V) = d(\text{centre of } U, \text{ centre of } V) \quad \text{also to be defined!} \quad \text{CENTRED LINKAGE}$$

Hierarchical Agglomerative Clustering

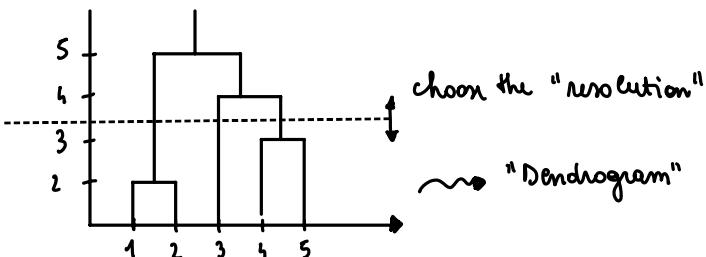
- Set a distance and a linkage

- Every unit is a cluster

- Iterate until convergence

- aggregate the two clusters which are closest

<u>EX</u>	$m=5$	$D:$	$\{1,2\}$	$\{1,2\}$	$\{1,2\}$	$\{1,2\}$
		1 2 3 4 5	{1,2} 3 4 5	{1,2} 3 {4,5}	{1,2} {4,5} 0	{1,2} {3,4,5} 0
		1 0	0	0	0	0
		2 0	2 0	3 0	3 5 0	5 0
		3 6 5 0	3 5 0	4 9 4 0	4 8 0	5 8 3 0
		4 10 9 4 0	4 9 4 0	5 8 3 0	5 8 3 0	5 8 3 0
		5 9 8 5 3 0	5 8 3 0	5 8 3 0	5 8 3 0	5 8 3 0



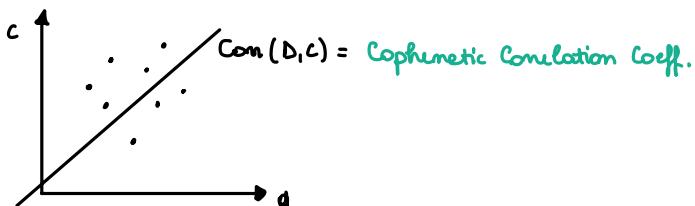
OBS Jitter the data  $\Rightarrow$  add noise  
 Careful when using single linkage  $\Rightarrow$  chain effect  
 CL, AL  $\Rightarrow$  ellipsoidal clusters

## LECTURE 19 21/04/2022

New distance matrix

	1	2	3	4	5	cophenetic distances (ultrametric)
1	0					
2	2	0				
3	5	5	0			
4	5	5	4	0		
5	5	5	4	3	0	

For every couple of units  $x_i, x_j$  in  $X$  (training set) you have 2 distances  $\begin{cases} d_{ij} \text{ distance in } D \\ c_{ij} \text{ distance in } C \end{cases}$



## Ward agglomerative algorithm

$X = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}$  We want to divide it in  $\underbrace{C_1, \dots, C_k}_{k}$  clusters  $C_i \cap C_j = \emptyset \quad \forall i \neq j, \bigcup_{i=1}^k C_i = X$

for  $i = 1, \dots, k$

$$\bar{x}_j = \frac{1}{\#C_j} \sum_{x_i \in C_j} x_i$$

$$ESS_j = \sum_{x_i \in C_j} \|x_i - \bar{x}_j\|^2$$

variability around the centroid

• When every unit is a cluster  $\Rightarrow ESS = 0$

iterate:  
 aggregate two clusters which cause the minimum increment of ESS

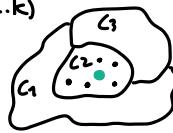
## k-Means

$\mathbb{R}^p, d, \mathbb{X}$

Let  $C_1, \dots, C_k$  be clusters partitioning  $\mathbb{X}$

DEF  $\bar{x}_j$  CENTROID of  $C_j$  ( $j=1 \dots k$ )

$$\bar{x}_j = \underset{\bar{x} \in \mathbb{R}^p}{\operatorname{argmin}} \sum_{x_i \in C_j} d^2(x_i, \bar{x})$$



$\bar{x}_j$  may not belong to  $\mathbb{X}$

→ If  $d$  is euclidean this is the barycenter

⇒ "Find  $C_1, \dots, C_k$  st.  $\sum_{j=1}^k \sum_{x_i \in C_j} d^2(x_i, \bar{x}_j)$  is minimum"  $\rightsquigarrow$  Grundy solution with k-means!

### • Initialization Step

either: specify  $C_1, \dots, C_k$  randomly (distributing  $k$  labels randomly among the  $m$  units)

or: specify  $\bar{x}_1, \dots, \bar{x}_k$  centroids in  $\mathbb{R}^p$  randomly

### • Iterate until convergence

→ for  $j=1, \dots, k$  compute  $C_j$  so that at the end we have  $\bar{x}_1, \dots, \bar{x}_k$  ] Difficult step! Minimization  $\Rightarrow$  Way out

→ for all  $x_i$  attribute it to cluster  $C_j$  if  $d^2(x_i, \bar{x}_j) = \min \{ d^2(x_i, \bar{x}_j), j=1, \dots, k \}$

Convergence  $\Leftrightarrow$  Step 1 doesn't change

$$\begin{aligned} \bar{x}_j &= \underset{\substack{\text{modified} \\ \mathbb{X} \in \mathbb{X}}}{\operatorname{argmin}} \sum_{x_i \in C_j} d^2(x_i, \bar{x}) \\ &\text{only points of the training set} \\ &\text{Medoid} \\ &\text{(k-medoid)} \end{aligned}$$

## Label switching problem

Suppose  $k=3$   $n=2$

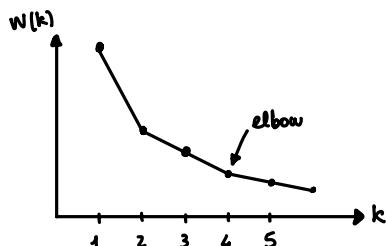
	1	2	3
1	*		
2		*	
3	*		

They give you the same  
cluster, but with different names  
(in  $p=2$  easy, pictures!  
in  $p=100$  permute the matrix...)

## How to choose $k$ ?

One way to do it

$$W(k) = \sum_{j=1}^k \sum_{x_i \in C_j} d^2(x_i, \bar{x}_j)$$



## Graphical representation

- Fisher's score

- Multidimensional scaling (MDS)

↳ Raw data  $\Rightarrow$  New representation

$\mathbb{R}^p, d \quad \mathbb{R}^q$  ( $q \leq p$ ) desired

$$\Delta = [d_{ij}] \quad \mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}$$

$$\mathbb{Y} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} \text{ s.t. } \text{euclid.}(y_i, y_j) \approx d_{ij}$$

- Classical MDS

Find  $y_1, \dots, y_m$  in  $\mathbb{R}^q$  s.t.  $\sum_{ij} (d_{ij} - \delta_{ij})^2$  is min (no unique sol!)

OBS if  $d$  is Euclidean distance in  $\mathbb{R}^p \Rightarrow$  PCA span  $(e_1, \dots, e_q)$

- Kruskal's MDS

$$\text{STRESS} = \frac{\sum_{i,j} (\theta(d_{ij}) - \delta_{ij})^2}{\sum_{i,j} \delta_{ij}^2} \quad \text{minimize w.r.t. } \theta$$

## LECTURE 20 26/4/2022

$Y \in \mathbb{R}$  target

$\underline{x} = (x_1, \dots, x_p)^\top \in \mathbb{R}^p$  features

$\underline{y} = f(\underline{x}) + \varepsilon \quad \varepsilon \perp \!\!\! \perp \underline{x}$

↳ How to capture it? • Use prior knowledge

• Data

$\sigma$ -field (everything you know when you know  $\underline{x}$ )

Optimal  $f$  w.r.t. MSE  $\Rightarrow f(\underline{x}) = \mathbb{E}[Y | \underline{x}]$

We know an analytical expression when both  $\underline{x}, Y$  are Gaussian

- parametric model

$$\hat{f}(\underline{x}) = \beta_0 + \beta_1 \underline{x}_1^2 + \beta_2 \underline{x}_2 + \beta_3 \underline{x}_3 \quad \text{PRIOR KNOWLEDGE}$$

- totally data driven LINEAR MODELS

- CART - Random forests  
- NN - DL  
Good for predictions but it's black-box

BETTER WITH LOTS OF DATA  
(the features already explain everything)  
without need of transformations

## CART (Classification and Regression Trees)

Idea:  $f$  is piecewise constant

$$\hat{f}(\underline{x}) = \sum_{j=1}^J c_j \mathbb{1}_{\{\underline{x} \in R_j\}}$$

finite partition



$$c_j = \bar{y}_j = \text{mean of } Y \text{ in } R_j$$

$$\text{Goal: find } J \text{ and } \{R_1, \dots, R_J\} \text{ s.t. } \sum_{j=1}^J \sum_{i: \underline{x}_i \in R_j} \|y_i - \bar{y}_j\|^2 \text{ is min}$$

⚠ WARNING overfitting ⚡

CART is a greedy algorithm for solving ↑

Iterate the Step

Step: find  $S_1$  cutoff such that  $\underbrace{\sum_{i=1}^n (y_i - \bar{y}_i)^2}_{\text{Variability of } Y} - \left[ \underbrace{\sum_{x_{ij} \leq S_1} (y_i - \bar{y}_1)^2}_{R_1} + \underbrace{\sum_{x_{ij} > S_1} (y_i - \bar{y}_2)^2}_{R_2} \right]$  is max.

Variability of  $Y$  after the cut

Repeat for  $i = 2, \dots, p$

Having  $S_1, \dots, S_p$ , choose  $S_j^*$  that maximizes ↑

$$R_1 = \{ \underline{x} : x_{ij} \leq S_j^* \} \quad R_2 = \underline{x}^c$$

Stop splitting if the number of units in a region is less than  $m$  ( $m = 5, 10$ , choose wisely)

## How to control overfitting?

► prune the tree

↳ penalize the number of leaves

$$\min \left\{ \sum_{j=1}^J \sum_{i: \underline{x}_i \in R_j} (y_i - \bar{y}_j)^2 + \alpha J \right\} \quad \alpha > 0$$

PENALIZATION PARAMETER  
(usually found through cross-validation)

## Linear Models for regression

$$\mathbb{X} = \begin{bmatrix} x_{11} & \dots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_{m1} & \dots & x_{mp} \end{bmatrix} \begin{bmatrix} Y \\ Y_1 \\ \vdots \\ Y_m \end{bmatrix} \Rightarrow \text{Design matrix } \mathbb{Z} = \begin{bmatrix} 1 & z_{11} & \dots & z_{1n} \\ 1 & z_{21} & \dots & z_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & z_{m1} & \dots & z_{mn} \end{bmatrix} \begin{bmatrix} Y \\ Y_1 \\ \vdots \\ Y_m \end{bmatrix}$$

$$z_i = g_i(x_1, \dots, x_p)$$

$\hookrightarrow$  known  $i=1, \dots, n$

$$\hat{f}(z) = \beta_0 + \beta_1 z_1 + \dots + \beta_n z_n \quad z = (1 z_1 \dots z_n)^T \in \mathbb{R}^{n+1}$$

LM:  $Y = \mathbb{Z}' \beta + \varepsilon$

$$\varepsilon \perp \vDash \mathbb{E}[\varepsilon] = 0$$

$$\text{Var}(\varepsilon) = \sigma^2 I$$

Model for the experiment generating the data:

$$Y = (Y_1 \dots Y_m)^T \in \mathbb{R}^m$$

$$\mathbb{Z} \in \mathbb{R}^{m \times (n+1)}$$

$$Y = \mathbb{Z}\beta + \varepsilon$$

$$\text{OLS } \varepsilon \perp \mathbb{Z} \quad \mathbb{E}[\varepsilon] = 0 \quad \text{Cov}(\varepsilon) = \sigma^2 I \quad (\text{no Gaussianity assumption here!})$$

Ordinary Least Squares method

Find  $\beta$  s.t.  $\|Y - \mathbb{Z}\beta\|^2$  is min

PROP If  $\mathbb{Z}$  is full-rank ( $\text{rank } \mathbb{Z} = n+1$ )

$$\Rightarrow \hat{\beta} = \underset{\beta \in \mathbb{R}^{n+1}}{\text{arg min}} \|Y - \mathbb{Z}\beta\|^2 = (\mathbb{Z}' \mathbb{Z})^{-1} \mathbb{Z}' Y$$

proof  $\mathbb{Z}\beta = \beta_0 1 + \beta_1 e_1 + \dots + \beta_n e_n$

$$\mathbb{Z} = \begin{bmatrix} 1 & e_1 & \dots & e_n \end{bmatrix}$$

$(n+1) \times (n+1)$  full rank

$$\mathbb{Z}' \mathbb{Z} = \sum_{i=1}^{n+1} \lambda_i e_i e_i' \quad \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_{n+1} > 0$$

$$\downarrow$$

$$(\mathbb{Z}' \mathbb{Z})^{-1} = \sum_{i=1}^{n+1} \frac{1}{\lambda_i} e_i e_i'$$

$$\mathcal{L}(\mathbb{Z}) \ni q_i := \frac{1}{\sqrt{\lambda_i}} \mathbb{Z} e_i \quad i=1, \dots, n+1$$

$$q_i' q_j = \frac{1}{\sqrt{\lambda_i \lambda_j}} e_i' \mathbb{Z}' \mathbb{Z} e_j = \frac{\lambda_i}{\sqrt{\lambda_i \lambda_j}} e_i' e_j = \begin{cases} 1 & i=j \\ 0 & i \neq j \end{cases} \Rightarrow \text{o.m.b. for } \mathcal{L}(\mathbb{Z})$$

$$\Pi_{Y \mid \mathcal{L}(\mathbb{Z})} = \sum_{i=1}^{n+1} \Pi_{Y \mid q_i} q_i = \sum_{i=1}^{n+1} q_i q_i' \frac{1}{q_i' q_i} Y = \sum_{i=1}^{n+1} \frac{1}{\lambda_i} \mathbb{Z} e_i e_i' \mathbb{Z}' Y = \mathbb{Z} \left( \sum_{i=1}^{n+1} \frac{1}{\lambda_i} e_i e_i' \right) \mathbb{Z}' Y = \mathbb{Z} \underbrace{(\mathbb{Z}' \mathbb{Z})^{-1} \mathbb{Z}'}_{\mathbf{H}} Y$$

REMARK If  $\underbrace{\text{rank } \mathbb{Z}}_k < n+1$

We can't invert it, but we can use the generalized Moore-Penrose inverse

## LECTURE 21 2/5/2022

### Ex (ANOVA)

$$\Rightarrow \begin{cases} X_{11}, \dots, X_{1m_1} \stackrel{iid}{\sim} N_1(\mu_1, \sigma^2) \\ \vdots \\ X_{g1}, \dots, X_{gm_g} \stackrel{iid}{\sim} N_g(\mu_g, \sigma^2) \end{cases} \quad Y = (X_{11} X_{12} \dots X_{1m_1} X_{21} \dots X_{gm_g}) \in \mathbb{R}^{m_1 + \dots + m_g \times 1}$$

Take as design matrix  $\mathbb{Z} = \begin{bmatrix} 1 & 1 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & \dots & 1 \\ 1 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & \dots & 0 \end{bmatrix} \begin{matrix} m_1 \\ m_2 \\ \vdots \\ m_g \end{matrix}$

$\updownarrow \quad \mathbb{Z} \quad \mathbb{Z}' \quad m$

$\beta = (\mu, \tau_1, \dots, \tau_g) \in \mathbb{R}^{g+1}$

$$Y = Z\beta + \varepsilon \quad \text{means} \quad X_{ij} = \mu + \tau_i + \varepsilon_{ij} \quad \begin{array}{l} i=1, \dots, g \\ j=1, \dots, m_i \end{array} \quad \text{if } \varepsilon = (\varepsilon_1 \dots \varepsilon_{gm}) \in \mathbb{R}^m \sim N_m(0, \sigma^2 I)$$

► Notice  $Z$  is NOT full-rank

$$\dim(\mathcal{L}(Z)) = g = n < n+1$$

► We impose  $\sum_{i=1}^g m_i \tau_i = 0$  to avoid overparametrization

$$\rightarrow \tau_g = -\frac{1}{m_g} \sum_{i=1}^{g-1} m_i \tau_i \quad \text{consider treatment at (arbitrary!) level } g$$

► Choose  $Z = [\underline{z}_0 \ \underline{z}_1 \ \dots \ \underline{z}_g]$  it's full-rank!

$$\left[ \begin{array}{c} 1 \\ \vdots \\ 1 \end{array} \right] \left[ \begin{array}{c} 1 \\ \vdots \\ 1 \\ 0 \\ \vdots \\ 0 \\ -\frac{m_1}{m_g} \\ \vdots \\ -\frac{m_{g-1}}{m_g} \end{array} \right] \left\{ \begin{array}{c} m_1 \\ \vdots \\ m_{g-1} \\ m_g \end{array} \right\} \left[ \begin{array}{c} 0 \\ \vdots \\ 1 \\ -\frac{m_{g-1}}{m_g} \\ \vdots \\ -\frac{m_{g-1}}{m_g} \end{array} \right] \left\{ \begin{array}{c} m_{g-1} \\ m_g \end{array} \right\}$$

How to estimate  $\beta_0, \dots, \beta_n$  and  $\sigma^2$ ?

### Ordinary Least Squares (OLS)

Assume  $Z$  full rank:  $\text{rank } Z = n+1 \leq m$

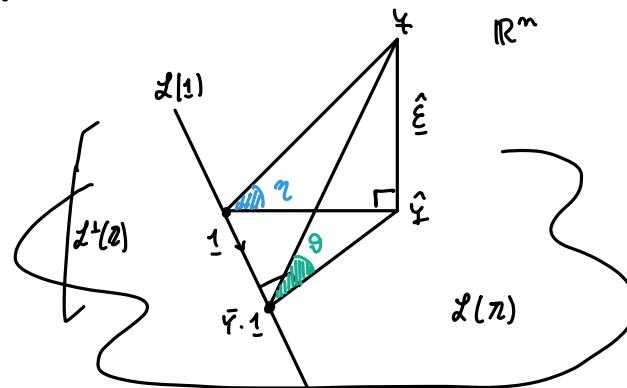
$\hat{Y} = \Pi_{\mathcal{L}(Z)} Y$  fitted values

$\hat{\varepsilon} = Y - \hat{Y}$  residuals

$$Y = \hat{Y} + \hat{\varepsilon} \Rightarrow \hat{Y} \perp \hat{\varepsilon}$$

$$\Pi_{\mathcal{L}(Z)} Y = \bar{Y} \cdot 1 \quad \text{left adjoint}$$

$$\Pi_{\mathcal{L}(Z)} Y = \frac{1 \cdot 1^T}{1^T 1} \hat{Y} = \frac{1 \cdot 1^T}{1^T 1} H Y = \frac{1 \cdot 1^T}{1^T 1} Y = \bar{Y} \cdot 1$$



$$\Rightarrow \|Y - \bar{Y} \cdot 1\|^2 = \|\hat{Y} - \bar{Y} \cdot 1\|^2 + \|\hat{\varepsilon}\|^2$$

$$\sum (y_i - \bar{y})^2 = \sum (y_i - \bar{y})^2 + \sum \hat{\varepsilon}_i^2$$

centered  
on the mean

$$R^2 = 1 - \frac{\sum \hat{\varepsilon}_i^2}{\sum (y_i - \bar{y})^2} = 1 - \sin^2(\theta) = \cos^2 \theta \in [0, 1]$$

•  $R^2 = 1 \Rightarrow \theta = 0 \Rightarrow Y = \hat{Y}$  perfect fit

•  $R^2 = 0 \Rightarrow \theta = \pi/2 \Rightarrow \hat{Y} = \bar{Y} \cdot 1$  predict with mean

► "Models through the origin"

$$Z = \begin{bmatrix} 1 & z_1 & \dots \\ \vdots & \vdots & \ddots \end{bmatrix} \quad \beta = (\beta_0 \ \beta_1 \dots \ \beta_g) \quad \Rightarrow \bar{Y} = \bar{Y}$$

you cannot project on  $\mathcal{L}(Z)$  since it's outside  $\mathcal{L}(Z)$   
 $R^2$  cannot be computed, or it can, but  $\star$  doesn't hold anymore, it might be also greater than 1!

You can only use  $\sum \hat{\varepsilon}_i^2 = \sum \hat{y}_i^2 + \sum \hat{\varepsilon}_i^2$

$$\Rightarrow \text{invent } \tilde{R}^2 = \frac{\sum \hat{y}_i^2}{\sum y_i^2} = \cos^2 \eta \in [0, 1]$$

► To take into account the model complexity

$$R^2_{adj} = 1 - \frac{\frac{1}{m-(n+1)} \sum_{i=1}^n \hat{\epsilon}_i^2}{\frac{1}{m-1} \sum_{i=1}^m (y_i - \bar{y})^2}$$

## Properties of $\hat{\beta}$ and $\hat{\epsilon}$

Assume full-rank, rank  $Z = n+1 \leq m$

- 1)  $E[\hat{\beta}] = \beta$  (unbiased)
- 2)  $Cov(\hat{\beta}) = \sigma^2 (Z'Z)^{-1}$
- 3)  $E[\hat{\epsilon}] = 0$
- 4)  $Cov(\hat{\epsilon}) = \sigma^2 (I - H)$
- 5)  $E[\hat{\epsilon}' \hat{\epsilon}] = E\left[\sum_{i=1}^m \hat{\epsilon}_i^2\right] = \sigma^2 (m - (n+1)) \stackrel{\text{cov}}{\Rightarrow} \frac{\hat{\epsilon}' \hat{\epsilon}}{m - (n+1)} = \frac{\sum \hat{\epsilon}_i^2}{m - (n+1)}$  unbiased for  $\sigma^2$

### Proof

$$\begin{aligned}
 1) E[\hat{\beta}] &= E[(Z'Z)^{-1} Z' \underline{y}] = (Z'Z)^{-1} Z' E[\underline{y}] = \beta \\
 2) Cov(\hat{\beta}) &= Cov((Z'Z)^{-1} Z' \underline{y}) = (Z'Z)^{-1} Z' (I\sigma^2) Z (Z'Z)^{-1} = \sigma^2 (Z'Z)^{-1} \\
 3) E[\hat{\epsilon}] &= E[(I - H)\underline{y}] = E[\underline{y}] - E[H\underline{y}] = Z\beta - Z\beta = 0 \\
 4) Cov(\hat{\epsilon}) &= Cov((I - H)\underline{y}) = (I - H)\sigma^2 \underline{y} (I - H)' = \sigma^2 (I - H)(I - H)' = \sigma^2 (I - H)^2 = \sigma^2 (I - H) \\
 5) E[\hat{\epsilon}' \hat{\epsilon}] &= \underbrace{\text{tr}[E[\hat{\epsilon}' \hat{\epsilon}]]}_{1 \times 1 \text{ matrix}} = E[\text{tr}(\hat{\epsilon}' \hat{\epsilon})] = E[\text{tr}(\hat{\epsilon} \hat{\epsilon}')] = E[\text{tr}(\underbrace{\hat{\epsilon} \hat{\epsilon}'}_{\text{matrix}} (I - H)' \underline{y} (I - H))] = E[\text{tr}((I - H) \underline{y} \hat{\epsilon}' (I - H)')] \\
 &= \text{tr}((I - H) E[\hat{\epsilon} \hat{\epsilon}'] (I - H)') = \sigma^2 \text{tr}((I - H)(I - H)') \quad (I - H) \hat{\epsilon} \\
 &= \sigma^2 \underbrace{\text{tr}(I - H)}_{\text{tr}(I) - \text{tr}(H)} = m - \text{tr}(Z(Z'Z)^{-1} Z') = m - \text{tr}(Z'Z (Z'Z)^{-1}) = m - \text{tr}(I_{n+1}) = m - (n+1) \\
 &= \sigma^2 (m - (n+1))
 \end{aligned}$$

The det. must be zero  
(This will have an impact)

Corollary  $S^2 = \frac{\sum \hat{\epsilon}_i^2}{m - (n+1)}$  unbiased for  $\sigma^2$

Obs  $Cov(\hat{\beta}) = \sigma^2 (Z'Z)^{-1}$

↳ are correlated, we can control the design matrix (e.g. PCA) Danger zone when regressors are

Obs Design of Experiments

From now on  $\epsilon \sim N_m(0, \sigma^2 I)$

Prop Assume  $Z$  full rank and  $\sigma^2$

- 1)  $\hat{\beta}$  and  $\hat{\sigma}^2 = \frac{\sum \hat{\epsilon}_i^2}{m}$  are MLE of  $\beta$  and  $\sigma^2$
- 2)  $\hat{\beta} \sim N_{n+1}(\beta, \sigma^2 (Z'Z)^{-1})$
- 3)  $\hat{\epsilon} \sim N_m(0, \sigma^2 (I - H))$
- 4)  $\hat{\epsilon} \perp \hat{\beta}$
- 5)  $\hat{\epsilon}' \hat{\epsilon} = \sum_{i=1}^m \hat{\epsilon}_i^2 \sim \sigma^2 \chi^2(m - (n+1))$

Proof:

- For 1): Just compute derivative of log likelihood

- For 2), 3) and 4) Note that:

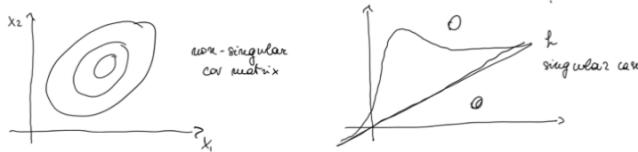
$$\begin{bmatrix} \hat{\beta} \\ \hat{\epsilon} \end{bmatrix} = \begin{bmatrix} (Z'Z)^{-1} Z' \\ I - Z(Z'Z)^{-1} Z' \end{bmatrix} \underline{y} \text{ and } \underline{y} \sim N_n(0, \sigma^2 I)$$

Then we just need to verify:  $Cov \begin{bmatrix} \hat{\beta} \\ \hat{\epsilon} \end{bmatrix} = \begin{bmatrix} (Z'Z)^{-1} & 0 \\ 0 & I - H \end{bmatrix}$  and so since we have two zero blocks off diagonal than the two vectors are independent!



- For 5):  $\hat{\epsilon} \sim \mathcal{N}_n(\underline{0}, \sigma^2(I - H))$  but:  
 $(I - H)$  is a singular matrix, indeed  $\det(I - H) = 0$  as its rank is  $n - (r + 1)$  and not  $n$ !

So we have a singular Gaussian distribution: it's a Gaussian distribution defined on a sub-space of Lebesgue measure equal to zero:



So the error vector  $\hat{\epsilon}$  is a non-singular Gaussian, whereas the residual vector  $\hat{\epsilon}$  is a singular Gaussian constrained on a linear sub-space!

To prove this we need to compute:  $\hat{\epsilon}^T(I - H)^{-1}\hat{\epsilon}$  but  $I - H$  is singular!

So we can compute the Moore-Penrose Generalised inverse so:  $\hat{\epsilon}^T(I - H)^{\dagger}\hat{\epsilon} \sim \sigma^2\chi^2(n - (r + 1))$  and this is the Mahalanobis Distance of a Gaussian from its mean. Then:  $\hat{\epsilon}^T\hat{\epsilon} \sim \sigma^2\chi^2(n - (r + 1))$

## LECTURE 22 5/5/2022

$$2) \Rightarrow \frac{1}{S^2} (\hat{\beta} - \beta)^T Z^T Z (\hat{\beta} - \beta) \sim \chi^2(r+1) \quad \text{if } S^2 \text{ is }$$

$$3) \Rightarrow \frac{1}{S^2} \hat{\epsilon}^T \hat{\epsilon} \sim \chi^2(m - (r+1)) \quad \text{if } S^2 \text{ is }$$

General Penrose inverse

$$\frac{\frac{1}{S^2} (\hat{\beta} - \beta)^T Z^T Z (\hat{\beta} - \beta) / (n+1)}{\frac{1}{S^2} \hat{\epsilon}^T \hat{\epsilon} / (m - (r+1))} \sim F(r+1, m - (r+1))$$

$$CR_{1-\alpha}(\beta) = \left\{ \eta \in \mathbb{R}^{r+1} : \frac{1}{S^2} (\hat{\beta} - \eta)^T Z^T Z (\hat{\beta} - \eta) \leq (r+1) F_{1-\alpha}(r+1, m - (r+1)) \right\}$$

$$\frac{1}{S^2} \hat{\epsilon}^T \hat{\epsilon} \sim \chi^2(m - (r+1)) \rightarrow \frac{(m - (r+1)) S^2}{S^2} \sim \chi^2(m - (r+1)) \rightarrow CI_{1-\alpha}(\beta) = \left[ \frac{(m - (r+1)) S^2}{\chi^2_{\alpha/2}(m - (r+1))}, \frac{(m - (r+1)) S^2}{\chi^2_{1-\alpha/2}(m - (r+1))} \right]$$

Linear combinations,  $\underline{\alpha} \in \mathbb{R}^{r+1}$

$$\underline{\alpha}^T \beta = \alpha_0 \beta_0 + \dots + \alpha_r \beta_r \sim N_1(\underline{\alpha}^T \beta, \underline{\alpha}^T (\underline{\alpha}^T \underline{\alpha})^{-1} \underline{\alpha}) \rightarrow \frac{\underline{\alpha}^T (\hat{\beta} - \beta)}{\underline{\alpha}^T (\underline{\alpha}^T \underline{\alpha})^{-1} \underline{\alpha}} \sim N(0, 1) \quad \text{if} \quad \underline{\alpha}^T \hat{\epsilon} \sim \underline{\alpha}^T \underline{\alpha} \chi^2(m - (r+1)) \rightarrow \frac{\underline{\alpha}^T (\hat{\beta} - \beta)}{\underline{\alpha}^T (\underline{\alpha}^T \underline{\alpha})^{-1} \underline{\alpha}} \sim t(m - (r+1))$$

$$CI_{1-\alpha}(\underline{\alpha}^T \beta) = \left[ \underline{\alpha}^T \hat{\beta} \pm S \sqrt{\underline{\alpha}^T (\underline{\alpha}^T \underline{\alpha})^{-1} \underline{\alpha}} t_{1-\alpha/2}(m - (r+1)) \right]$$

$H_0: \underline{\alpha}^T \beta = \gamma_0 \quad \text{vs} \quad \underline{\alpha}^T \beta \neq \gamma_0$

$$\text{Reject at } \alpha \in [0, 1] \text{ if } \frac{|\underline{\alpha}^T \hat{\beta} - \gamma_0|}{S \sqrt{\underline{\alpha}^T (\underline{\alpha}^T \underline{\alpha})^{-1} \underline{\alpha}}} > t_{1-\alpha/2}(m - (r+1))$$

Ex (Done in all software for linear regression)

$H_0: \beta_i = 0 \quad \text{vs} \quad \beta_i \neq 0$

$$\text{Reject if } \frac{|\hat{\beta}_i - 0|}{S \sqrt{\text{diag}((Z^T Z)^{-1})}} > t_{1-\alpha/2}(m - (r+1))$$

Sim. CI for  $\underline{\alpha}^T \beta$

$$\text{From the Maximum Likelihood: } \max_{\underline{\alpha} \in \mathbb{R}^{r+1}} \frac{1}{S^2} \frac{(\underline{\alpha}^T (\hat{\beta} - \beta))^2}{\underline{\alpha}^T (\underline{\alpha}^T \underline{\alpha})^{-1} \underline{\alpha}} = \frac{1}{S^2} (\hat{\beta} - \beta)^T Z^T Z (\hat{\beta} - \beta) \sim (r+1) F(r+1, m - (r+1))$$

$$\Rightarrow \text{Sim. } CI_{1-\alpha}(\underline{\alpha}^T \beta) = \left\{ \underline{\alpha}^T \beta \pm \sqrt{\underline{\alpha}^T (\underline{\alpha}^T \underline{\alpha})^{-1} \underline{\alpha}} \sqrt{S^2(r+1) F_{1-\alpha}(r+1, m - (r+1))} \right\}$$

NOT GIVEN BY R!  
BUT VERY MEANINGFUL

C  $p \times (n+1)$  not random

$H_0: C\beta = 0$  vs  $C\beta \neq 0$

$$\left[ \begin{array}{c} C_{11}\beta_0 + C_{12}\beta_1 + \dots + C_{1(n+1)}\beta_n \\ C_{21}\beta_0 + \dots + C_{2(n+1)}\beta_n \end{array} \right] = \left[ \begin{array}{c} 0 \\ \vdots \\ 0 \end{array} \right]$$

$$C\hat{\beta} \sim N_p(C\beta, \sigma^2 C(\mathbf{Z}'\mathbf{Z})^{-1}C)$$

$$\text{Under } H_0: \frac{\frac{1}{p} (\hat{C\beta})' (\sigma^2 C(\mathbf{Z}'\mathbf{Z})^{-1}C)^{-1} (\hat{C\beta})}{\hat{\xi}' \hat{\xi}} \sim \chi^2(p) \sim F(p, m-(n+1))$$

$$\text{So } \underbrace{\frac{1}{S^2} (\hat{C\beta})' [C(\mathbf{Z}'\mathbf{Z})^{-1}C]^{-1} (\hat{C\beta})}_{\text{reject at level } \alpha \text{ if this is } > p F_{1-\alpha}(p, m-(n+1))} \sim p F(p, m-(n+1))$$

Special case

$$Y = \beta_0 + \beta_1 Z_1 + \dots + \beta_n Z_n + \varepsilon \quad H_0: \beta_1 = \beta_2 = \dots = \beta_{n-1} = 0$$

can I take out  
simultaneously the  
last  $p$ ?

$$\Rightarrow \text{Choose } C = \begin{bmatrix} 0 & \dots & 0 & 1 \\ \vdots & \ddots & \vdots & \vdots \\ 0 & \dots & 0 & 1 \end{bmatrix} = \begin{bmatrix} 0 & I_p \end{bmatrix} \quad p \times (n+1)$$

↳ We can SHUFFLE to test the subset we want!

$$\begin{array}{c} Y \\ \downarrow \\ L(Z_1) \quad L(Z) \\ \curvearrowright \quad \curvearrowright \\ \mathbf{Z} = [Z_1 \quad Z_2] \\ \uparrow \quad \tilde{p} \\ \text{simplified} \end{array}$$

Reject if (D) large:

$$\frac{\hat{\xi}_1' \hat{\xi}_1 - \hat{\xi}' \hat{\xi}}{p S^2} \sim F(p, m-(n+1))$$

Very special case

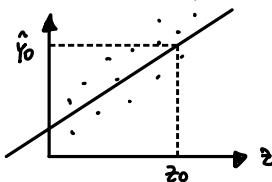
$$H_0: \beta_1 = \dots = \beta_n = 0 \quad (\text{all but } \beta_0)$$

Should I throw everything in the trash bin?

$$Z_1 = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \quad \hat{Y}_1 = \bar{Y} \cdot 1$$

$$\frac{\hat{\xi}_1' \hat{\xi}_1}{S^2 n} = \frac{\sum (Y_i - \bar{Y})^2}{n S^2} = \frac{\sum (\hat{Y}_i - \bar{Y})^2 / n}{\sum \hat{Y}_i^2 / n} \sim F(1, m-(n+1))$$

## PREDICTION



Gauss-Markov theorem Given  $\alpha \in \mathbb{R}^{n+1}$  given  $\alpha \in \mathbb{R}^{n+1}$

The optimal estimator of  $\alpha' \beta$  is  $\hat{\alpha}' \hat{\beta}$  (min. variance)

among those which are

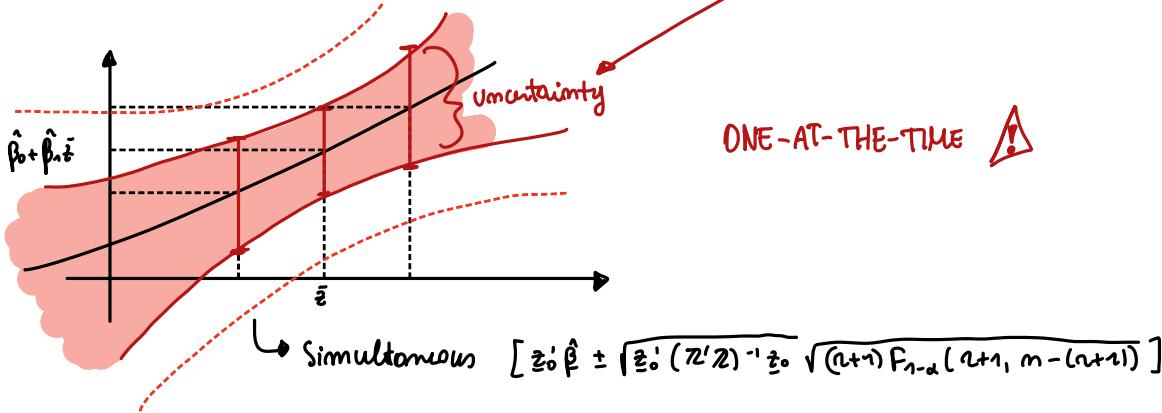
- unbiased
- linear functions of  $\gamma$

$$\frac{\hat{\alpha}' \hat{\beta} - \hat{\alpha}' \beta}{S \sqrt{\hat{\alpha}' (\mathbf{Z}' \mathbf{Z})^{-1} \hat{\alpha}}} \sim t(m-(n+1))$$

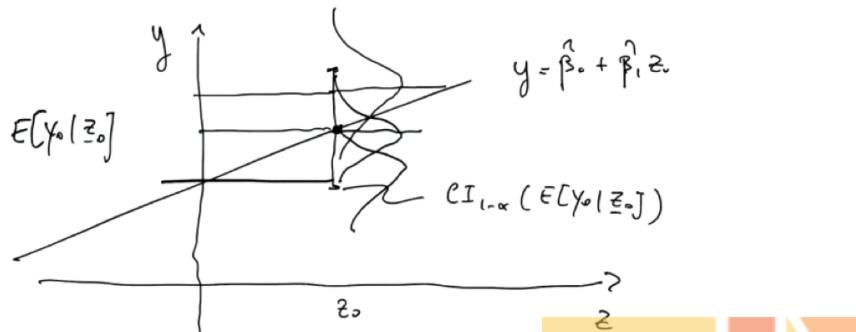
PIVOTAL

$$\Rightarrow CI_{1-\alpha}(\hat{\alpha}' \hat{\beta}) = \left[ \hat{\alpha}' \hat{\beta} \pm S \sqrt{\hat{\alpha}' (\mathbf{Z}' \mathbf{Z})^{-1} \hat{\alpha}} t_{\alpha/2}(m-(n+1)) \right] \quad \alpha \in (0, 1)$$

Uncertainty



In the above we have seen the prediction for  $\mathbb{E}[y_0 | \underline{z}_0]$ . Well, what about  $y_0$ ?



Can we find  $I$  such that  $P[y_0 \in I | \underline{z}_0] = 1 - \alpha$  ?  
PREDICTION INTERVAL

Obviously this PI will be larger than the CI for the mean: with the PI we are uncertain about the centre of the distribution and there is an extra variability due to the fact that  $y$  is different from its mean!

We know that:  $\underline{z}_0' \hat{\beta} \sim \mathcal{N}_1(\underline{z}_0' \beta, \sigma^2 \underline{z}_0' (\mathbb{Z}' \mathbb{Z})^{-1} \underline{z}_0)$  this is in the training set!

Then  $y_0$  is a new statistical unit: not in the training set! Thus:  $y_0 \perp \underline{z}_0' \hat{\beta}$  since  $\epsilon_0 \perp \epsilon$

Thus:  $y_0 - \underline{z}_0' \hat{\beta} \sim \mathcal{N}(0, \sigma^2(1 + \underline{z}_0' (\mathbb{Z}' \mathbb{Z})^{-1} \underline{z}_0))$  and  $\hat{\epsilon}' \hat{\epsilon} \sim \sigma^2 \chi^2(n - (r + 1))$  Moreover:  $\hat{\epsilon}' \hat{\epsilon} \perp y_0 - \underline{z}_0' \hat{\beta}$  since  $\hat{\epsilon}' \hat{\epsilon} \perp \epsilon$  and  $\hat{\epsilon}' \hat{\epsilon} \perp \underline{z}_0' \hat{\beta}$

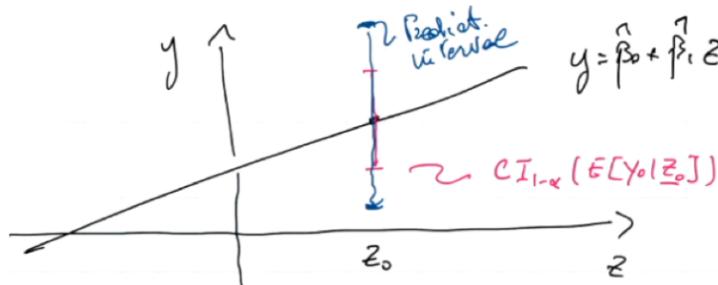
Thus  $\frac{y_0 - \underline{z}_0' \hat{\beta}}{\sqrt{\sigma^2(1 + \underline{z}_0' (\mathbb{Z}' \mathbb{Z})^{-1} \underline{z}_0)}} \sim \mathcal{N}(0, 1)$  and  $\sqrt{\frac{\hat{\epsilon}' \hat{\epsilon}}{\sigma^2(n - (r + 1))}}$  is a chi-squared. Thus their quotient is a  $t$ -student that is:

$$\frac{\frac{y_0 - \underline{z}_0' \hat{\beta}}{\sqrt{\frac{\hat{\epsilon}' \hat{\epsilon}}{\sigma^2(n - (r + 1))}}}}{\sqrt{\frac{\hat{\epsilon}' \hat{\epsilon}}{\sigma^2(n - (r + 1))}}} = \frac{y_0 - \underline{z}_0' \hat{\beta}}{S \sqrt{1 + \underline{z}_0' (\mathbb{Z}' \mathbb{Z})^{-1} \underline{z}_0}} \sim t(n - (r + 1))$$

So the prediction interval of probability  $1 - \alpha$  is:

$$PI_{1-\alpha}(y_0) = \left[ \underline{z}_0' \hat{\beta} \pm S \sqrt{1 + \underline{z}_0' (\mathbb{Z}' \mathbb{Z})^{-1} \underline{z}_0} t_{\alpha/2}(n - (r + 1)) \right]$$

We see that the prediction interval is larger than the confidence interval:



## GLS : Generalized Least Squares

$$\hat{\beta} = \underset{\beta \in \mathbb{R}^{n+1}}{\operatorname{argmin}} \|y - Z\beta\|^2$$

OLS

$$\hat{\beta} = \underset{\beta \in \mathbb{R}^{n+1}}{\operatorname{argmin}} (y - Z\beta)' W^{-1} (y - Z\beta) \quad W \text{ maxm pos. def}$$

GLS (if  $W=I$  we have OLS)

$$(y - Z\beta)' W^{-1} (y - Z\beta) = (W^{-1/2} y - W^{-1/2} Z\beta)' (W^{-1/2} y - W^{-1/2} Z\beta)$$

$$= \underbrace{\|W^{-1/2} y - W^{-1/2} Z\beta\|^2}_{\tilde{y} \quad \tilde{Z}} \Rightarrow \hat{\beta} = (\tilde{Z}' \tilde{Z})^{-1} \tilde{Z}' \tilde{y} = (Z' W^{-1} Z)^{-1} Z' W^{-1} y$$

We relax the assumption of LM

$$\operatorname{Cov}(\varepsilon) = \sigma^2 \Sigma \neq \sigma I \quad \Rightarrow \quad \hat{\beta} = \underset{\beta \in \mathbb{R}^{n+1}}{\operatorname{argmin}} (y - Z\beta)' \Sigma^{-1} (y - Z\beta)$$

$$\hat{\beta} = (Z' \Sigma^{-1} Z)^{-1} Z' \Sigma^{-1} y$$

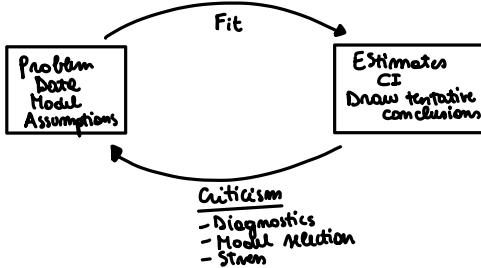
EX

$$\Sigma = \begin{bmatrix} m_1 & & \\ & m_2 & \\ & & \ddots \\ & & & m_n \end{bmatrix} \quad \operatorname{Var}(y_i) \propto m_i \quad y_i = \sum_{j=1}^{m_i} T_{ij}$$

"Weighted Least Squares" → You trust less some observations

↳ EX  $y_i = \frac{1}{m_i} \sum_j T_{ij} \quad \operatorname{Var}(y_i) \propto \frac{1}{m_i} \quad \Sigma = \begin{bmatrix} 1/m_1 & \dots & 1/m_n \end{bmatrix}$

## DIAGNOSTICS



- 1) Residual Analysis
- 2) Influential Cases
- 3) Collinearity

### 1) Residual Analysis

Model Fitted Model

$$y = Z\beta + \varepsilon \quad \hat{y} = Z\hat{\beta} + \hat{\varepsilon} \quad \hat{\beta} = (Z' Z)^{-1} Z' y$$

$$\mathbb{E}[\varepsilon] = 0 \quad \mathbb{E}[\hat{\varepsilon}] = 0 \quad \hat{\varepsilon} = y - \hat{y} = \pi y_1 \pi_{Z^{\perp}}(z)$$

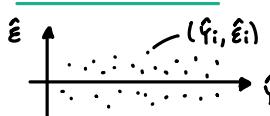
$$\operatorname{Cov}(\varepsilon) = I\sigma^2 \quad \operatorname{Cov}(\hat{\varepsilon}) = \sigma^2 (I - \pi)$$

$$\varepsilon \sim N_m(0, \sigma^2 I)$$

$$\hat{\varepsilon} \in \mathbb{R}^m \quad \hat{\varepsilon} \in \underline{Z^{\perp}}(Z) \subseteq \mathbb{R}^m$$

$$\dim = m - (n+1)$$

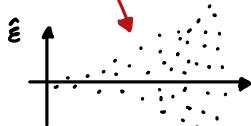
Plot the residuals



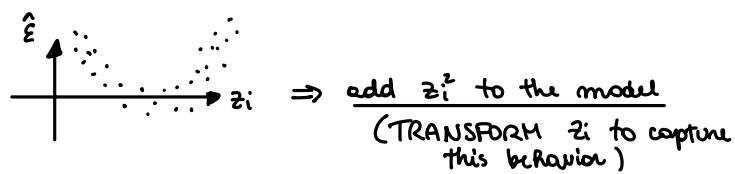
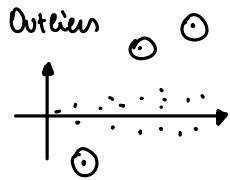
We don't want PATTERNS!

⚠️ HETEROGENEITY

You can't trust all your observations in the same way



⇒ Weighted Least Squares, Variance Stabilizing Transformations



$$\text{Var}(\hat{e}_i) = \sigma^2$$

$$\text{Var}(\hat{z}_i) = \sigma^2(1-R_{ii})$$

$$\Rightarrow \frac{\hat{e}_i}{S\sqrt{1-R_{ii}}}$$

STUDENTISED RESIDUALS

[for carefully designed experiments  $R_{ii} \approx 0$  so they are the same.]

The variances are different

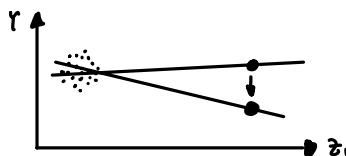
You can control  
 $\rightarrow H = \sigma^2(Z'Z)^{-1}Z'$

## 2) Influential Cases

$R_{ii}$  leverage

$$\in [0,1]$$

If  $R_{ii} \uparrow 1 \Rightarrow \text{Var}(\hat{e}_i) \rightarrow 0 \Rightarrow \hat{e}_i \rightarrow 0$  Fishy, you know this before collecting data!



[meaningful explanation in recording]

Hold out unit  $i=1, \dots, n$  from training set. Does the model change significantly?

$$\text{Cook's distance } d(\hat{\beta}^{(i)}, \hat{\beta}) = D_i = \frac{(\hat{\beta}^{(i)} - \hat{\beta})'(Z'Z)(\hat{\beta}^{(i)} - \hat{\beta})}{S^2(n+1)} = \left( \frac{\hat{e}_i}{S\sqrt{1-R_{ii}}} \right)^2 \frac{R_{ii}}{1-R_{ii}} \cdot \frac{1}{n+1}$$

## 3) Collinearity

$$\text{Cov}(\hat{\beta}) = \sigma^2(Z'Z)^{-1}$$

$$\text{Var}(\hat{\beta}_j) = \frac{\sigma^2}{\sum_{i=1}^n (z_{ij} - \bar{z}_j)^2} \cdot \frac{1}{1-R_j^2}$$

coefficient of determination when  $z_j$  is regressed on  $z_1, \dots, \cancel{z_j}, \dots, z_r$   
 $R_j \uparrow 1 \Rightarrow \text{Var}(\hat{\beta}_j) \uparrow \infty$

Variance Inflation Factor  
"VIF"

LECTURE 24 10/5/2022

## Model Selection

$$Y = Z\beta + \varepsilon$$

$[1 \ z_1 \ \dots \ z_n]^T \rightarrow$  we have 2 possibilities for each model ( $2^n$  models)

$$" (1+1)^n = \sum_{k=0}^n \binom{n}{k} 1 \cdot 1^{n-k} = \sum_{k=0}^n \binom{n}{k}$$

- Brute force / Griddy  $\Rightarrow$  we check them ALL (pretty expensive)

For  $k=0, \dots, n$

- fit the  $\binom{n}{k}$  models with  $k$  regressors
- Compute  $R^2$  and choose the one with maximum  $R^2$

$$R_0^2 \leq R_1^2 \leq \dots \leq R_n^2$$



Otherwise we could use  $R^2_{adj}$  (not monotonic), AIC, BIC to choose the best  $k$

► No intelligence, we're just checking them all and choosing according to some criteria (doable only if the models are not too big)

## • Forward / backward selection

↳ 1. select one variable (best model)

iterate (step k)

add variable corresponding to maximum increase of  $\frac{\text{SS}(M_k) - \text{SS}(M_{k-1})}{\text{SS}(M_k)}$

stop if the F-test fails (you cannot reject H<sub>0</sub>)

$$F = \frac{\frac{SS_{reg}(M_k) - SS_{reg}(M_{k+1})}{m - (k+1)}}{\frac{SS_{reg}(M_{k+1})}{m - (k+1)}}$$

### Model selection & collinearity

Note:  $\underline{y} = \underline{Z}\underline{\beta} + \underline{\epsilon}$  is the model for the observed data from OLS we get:  $\hat{\underline{\beta}} = (\underline{Z}^T \underline{Z})^{-1} \underline{Z}^T \underline{y}$  assuming  $\underline{Z}$  has full rank!

Then the fitted model for the mean of  $\underline{y}$  is:  $y_0 = \underline{z}_0^T \hat{\underline{\beta}}$

The fitted model always go through the baricentre of the observed data: for instance suppose  $\underline{z}_0 = \frac{\underline{Z}^T \underline{1}}{\underline{1}^T \underline{1}} = [1 \ \bar{z}_1 \ \dots \ \bar{z}_r]^T$  so it's the vectors of the mean of the regressors!

Then:  $y_0 = \underline{z}_0^T \hat{\underline{\beta}} = \frac{1^T \underline{Z}}{n} (\underline{Z}^T \underline{Z})^{-1} \underline{Z}^T \underline{y} = \frac{1}{n} \underline{1}^T H \underline{y} = \frac{1}{n} (H \underline{1})^T \underline{y} = \frac{1}{n} \underline{1}^T \underline{y} = \bar{y}$  Hence:  $\bar{y} = \hat{\beta}_0 + \hat{\beta}_1 \bar{z}_1 + \dots + \hat{\beta}_r \bar{z}_r \implies$

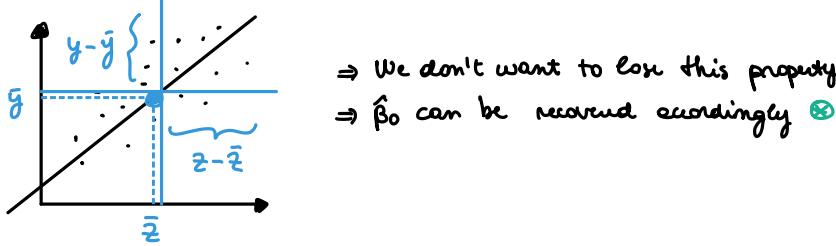
$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{z}_1 - \dots - \hat{\beta}_r \bar{z}_r$  So a different representation for the fitted model is:

$$y_0 - \bar{y} = \hat{\beta}_1(z_1 - \bar{z}_1) + \dots + \hat{\beta}_r(z_r - \bar{z}_r)$$

So instead of working with the original variables, we can first centre the data on the baricentre:  $\underline{z} - \bar{\underline{z}}$  and  $\underline{y} - \bar{\underline{y}}$  and then fit the model!

So **centering**: from  $\underline{y}$  we consider  $\mathbb{R}^n \ni \underline{y}^* = [\underline{y}_1 - \bar{\underline{y}}, \dots, \underline{y}_n - \bar{\underline{y}}]^T$  and from  $\underline{Z}$  we consider:  $\begin{bmatrix} z_{11} - \bar{z}_1 & \dots & z_{1r} - \bar{z}_r \\ \vdots & & \vdots \\ z_{n1} - \bar{z}_1 & \dots & z_{nr} - \bar{z}_r \end{bmatrix} = \underline{Z}^*$

which is an  $n \times r$  matrix!



OLS becomes

$$\left\{ \begin{array}{l} \hat{\underline{\beta}}^* = \underset{\underline{\beta} \in \mathbb{R}^r}{\text{argmin}} \| \underline{y}^* - \underline{Z}^* \underline{\beta} \|^2 \\ \hat{\beta}_0 = \bar{y} - \hat{\beta}_1^* \bar{z}_1 - \dots - \hat{\beta}_r^* \bar{z}_r \\ \hat{\beta}_i = \hat{\beta}_i^* \quad i=1,\dots,r \end{array} \right.$$

From now on, assume  $\underline{Z}$   $m \times n$  centred,  $\underline{y} \in \mathbb{R}^m$  centred.

### 1) PCA regression

The problem of collinearity is due to the fact that the regressors aren't orthogonal, thus we can transform them so that we get new orthogonal regressors: we just do PCA of the dependent variables!

$\underline{Z} = [\underline{z}_1, \dots, \underline{z}_r]$  with  $\underline{z}_i \in \mathbb{R}^n$  (Remember: these are already centred variables!) Then we perform PCA of  $\underline{Z}$  so that we get:  $PC_1, \dots, PC_r$

We can then check for an elbow and reduce the dimensionality: we consider  $PC_1, \dots, PC_k$  with  $k \leq r$

$$\underline{Z}^* = [PC_1 \dots PC_k] \quad \text{fit with this } \underline{Z}^* \quad \underline{y} = \underline{Z}^* \underline{\beta} + \underline{\epsilon}, \quad \underline{\beta} \in \mathbb{R}^k$$

$$\begin{aligned} \cdot \quad PC_1 &= e_{11} z_1 + \dots + e_{1k} z_n \\ \vdots & \\ \cdot \quad PC_k &= e_{k1} z_1 + \dots + e_{nk} z_n \end{aligned}$$

$$y_0 = \hat{\beta}_1 PC_1 + \dots + \hat{\beta}_k PC_k$$

Substituting this into the fitted model:

$$y_0 = \underbrace{z_1(\hat{\beta}_1 + \dots + \hat{\beta}_k)}_{\hat{\epsilon}_1} + \underbrace{z_2(\hat{\beta}_1 + \dots + \hat{\beta}_k)}_{\hat{\epsilon}_2} + \dots + \underbrace{z_r(\hat{\beta}_1 + \dots + \hat{\beta}_k)}_{\hat{\epsilon}_r} = z_1 \hat{y}_1 + \dots + z_r \hat{y}_r$$

We solved the problem of collinearity since the principal components are orthogonal, but there is no guarantee that the principal component are good directions for prediction: we might have thrown away the information of  $\mathbb{Z}$  correlated with  $y$

Note: there are algorithms that do dimension reduction without loosing too much correlation (e.g: slice-inverse regression).

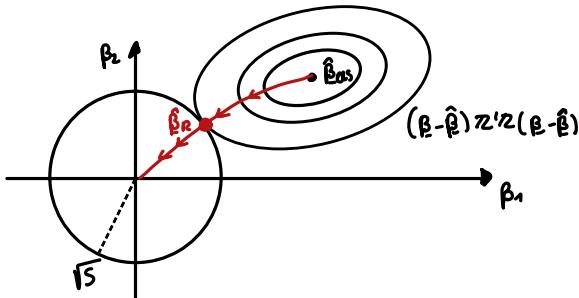
## 2) Ridge regression (Hoerl & Kennard 70's)

Collinearity  $\Rightarrow \text{Var}(\hat{\beta})$  explodes

Idea: constrain  $\|\hat{\beta}\|^2$  to be small

$$\text{OLS: } \hat{\beta} = \underset{\beta}{\text{argmin}} \|\mathbf{y} - \mathbf{Z}\beta\|^2 = \underset{\beta}{\text{argmin}} \|\underbrace{\mathbf{y} - \hat{\mathbf{y}}}_{\hat{\epsilon}} - \underbrace{\mathbf{Z}(\beta - \hat{\beta})}_{\epsilon \in \mathcal{Z}(\mathbf{Z})}\|^2 = \underset{\beta}{\text{argmin}} [\|\hat{\epsilon}\|^2 + \|\mathbf{Z}(\beta - \hat{\beta})\|^2] = \underset{\beta}{\text{argmin}} \|\mathbf{Z}(\beta - \hat{\beta})\|^2 = \underset{\beta}{\text{argmin}} (\beta - \hat{\beta})' \mathbf{Z}' \mathbf{Z} (\beta - \hat{\beta})$$

RIDGE:  $\left\{ \begin{array}{l} \underset{\beta}{\text{argmin}} (\beta - \hat{\beta})' \mathbf{Z}' \mathbf{Z} (\beta - \hat{\beta}) \\ \|\beta\|^2 \leq s \end{array} \right.$



OBS

- $\hat{\beta}_R$  is biased
  - For suitable choice of  $s$   $\mathbb{E}[\|\hat{\beta}_R - \beta\|^2] \xrightarrow{\text{Variance}} < \mathbb{E}[\|\hat{\beta}_{\text{OLS}} - \beta\|^2]$
- ↳ How to choose  $s$ ?

Lagrangian is  $\underset{\beta}{\text{argmin}} [\|\mathbf{Z}(\beta - \hat{\beta})\|^2 + \lambda \|\beta\|^2]$   $\lambda = \lambda(s)$   $\lambda \uparrow s \downarrow$

$$\hat{\beta}_R = (\mathbf{Z}' \mathbf{Z} + \lambda \mathbf{I})^{-1} \mathbf{Z}' \mathbf{y}$$

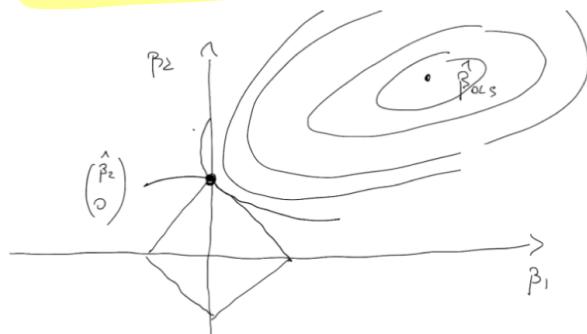
$$\mathbb{E}[\hat{\beta}_R] = (\mathbf{Z}' \mathbf{Z} + \lambda \mathbf{I})^{-1} \mathbf{Z}' \mathbf{Z} \beta \quad (\lambda = 0 \text{ OLS})$$

Collinearity ✓

Model Selection ✗

## 3) Lasso regression (Tibshirani)

It's similar to ridge but here we change the constraint so that we have the situation represented below:



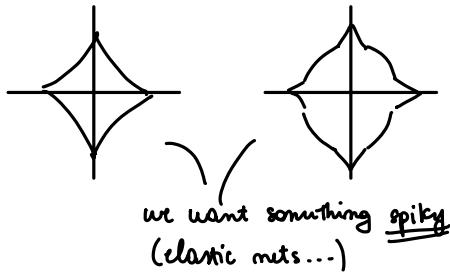
So instead of spherical constraints we take a diamond constraint: so that we can hope that the tangent point with the contour point will be on one of the vertices so that one of the  $\hat{\beta}_i$  will be exactly zero!

Not only we shrink the  $\hat{\beta}$  but we also select the variables: feature selection!



$$\text{LASSO: } \begin{cases} \underset{\beta}{\text{argmin}} (\beta - \hat{\beta})' Z' Z (\beta - \hat{\beta}) \\ \|\beta\|_1 \leq s \quad (\sum |\beta_i| \leq s) \end{cases}$$

Collinearity ✓  
Model Selection ✓  
Analytical solution X  $\Rightarrow$  numerical solutions !



## LECTURE 25 12/5/2022

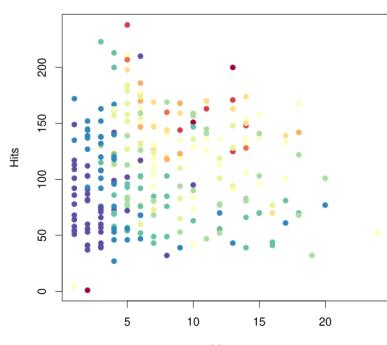
### Tree-based Methods

- Here we describe *tree-based* methods for regression and classification.
- These involve *stratifying* or *segmenting* the predictor space into a number of simple regions.
- Since the set of splitting rules used to segment the predictor space can be summarized in a tree, these types of approaches are known as *decision-tree* methods.

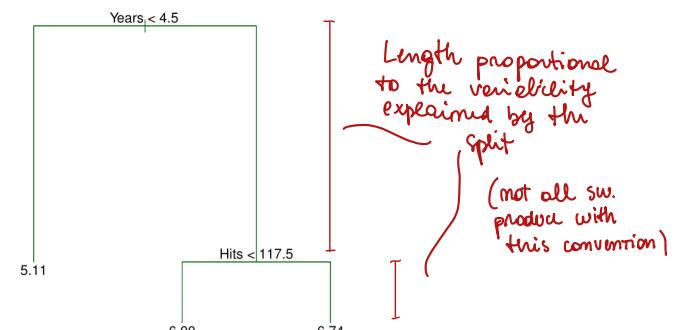
### Pros and Cons

- Tree-based methods are simple and useful for interpretation.
- However they typically are not competitive with the best supervised learning approaches in terms of prediction accuracy.
- Hence we also discuss *bagging*, *random forests*, and *boosting*. These methods grow multiple trees which are then combined to yield a single consensus prediction.
- Combining a large number of trees can often result in dramatic improvements in prediction accuracy, at the expense of some loss interpretation.
- Decision trees can be applied to both regression and classification problems.
- We first consider regression problems, and then move on to classification.

Baseball salary data: how would you stratify it?  
Salary is color-coded from low (blue, green) to high (yellow, red)



Decision tree for these data

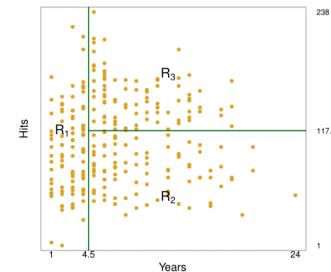


## Details of previous figure

- For the Hitters data, a regression tree for predicting the log salary of a baseball player, based on the number of years that he has played in the major leagues and the number of hits that he made in the previous year.
- At a given internal node, the label (of the form  $X_j < t_k$ ) indicates the left-hand branch emanating from that split, and the right-hand branch corresponds to  $X_j \geq t_k$ . For instance, the split at the top of the tree results in two large branches. The left-hand branch corresponds to  $\text{Years} < 4.5$ , and the right-hand branch corresponds to  $\text{Years} \geq 4.5$ .
- The tree has two internal nodes and three terminal nodes, or leaves. The number in each leaf is the mean of the response for the observations that fall there.

## Results

- Overall, the tree stratifies or segments the players into three regions of predictor space:  $R_1 = \{X \mid \text{Years} < 4.5\}$ ,  $R_2 = \{X \mid \text{Years} \geq 4.5, \text{Hits} < 117.5\}$ , and  $R_3 = \{X \mid \text{Years} \geq 4.5, \text{Hits} \geq 117.5\}$ .



6 / 51

7 / 51

## Terminology for Trees

- In keeping with the *tree* analogy, the regions  $R_1$ ,  $R_2$ , and  $R_3$  are known as *terminal nodes*.
- Decision trees are typically drawn *upside down*, in the sense that the leaves are at the bottom of the tree.
- The points along the tree where the predictor space is split are referred to as *internal nodes*.
- In the hitters tree, the two internal nodes are indicated by the text `Years<4.5` and `Hits<117.5`.

## Interpretation of Results

- Years** is the most important factor in determining **Salary**, and players with less experience earn lower salaries than more experienced players.
- Given that a player is less experienced, the number of **Hits** that he made in the previous year seems to play little role in his **Salary**.
- But among players who have been in the major leagues for five or more years, the number of **Hits** made in the previous year does affect **Salary**, and players who made more **Hits** last year tend to have higher salaries.
- Surely an over-simplification, but compared to a regression model, it is easy to display, interpret and explain

## Details of the tree-building process

- We divide the predictor space — that is, the set of possible values for  $X_1, X_2, \dots, X_p$  — into  $J$  distinct and non-overlapping regions,  $R_1, R_2, \dots, R_J$ .
- For every observation that falls into the region  $R_j$ , we make the same prediction, which is simply the mean of the response values for the training observations in  $R_j$ .

## More details of the tree-building process

- In theory, the regions could have any shape. However, we choose to divide the predictor space into high-dimensional rectangles, or *boxes*, for simplicity and for ease of interpretation of the resulting predictive model.
- The goal is to find boxes  $R_1, \dots, R_J$  that minimize the RSS, given by

$$\sum_{j=1}^J \sum_{i \in R_j} (y_i - \hat{y}_{R_j})^2,$$

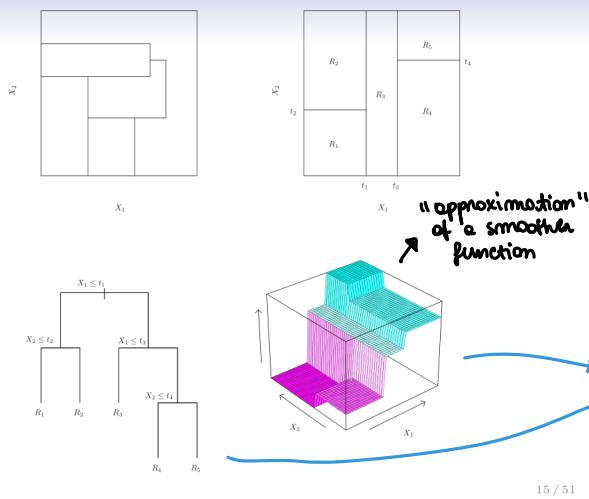
where  $\hat{y}_{R_j}$  is the mean response for the training observations within the  $j$ th box.

## Details— Continued

- We first select the predictor  $X_j$  and the cutpoint  $s$  such that splitting the predictor space into the regions  $\{X \mid X_j < s\}$  and  $\{X \mid X_j \geq s\}$  leads to the greatest possible reduction in RSS.
- Next, we repeat the process, looking for the best predictor and best cutpoint in order to split the data further so as to minimize the RSS within each of the resulting regions.
- However, this time, instead of splitting the entire predictor space, we split one of the two previously identified regions. We now have three regions.
- Again, we look to split one of these three regions further, so as to minimize the RSS. The process continues until a stopping criterion is reached; for instance, we may continue until no region contains more than five observations.

## Predictions

- We predict the response for a given test observation using the mean of the training observations in the region to which that test observation belongs.
- A five-region example of this approach is shown in the next slide.



Details of previous figure

**Top Left:** A partition of two-dimensional feature space that could not result from recursive binary splitting.

**Top Right:** The output of recursive binary splitting on a two-dimensional example.

**Bottom Left:** A tree corresponding to the partition in the top right panel.

**Bottom Right:** A perspective plot of the prediction surface corresponding to that tree.

15 / 51

## Pruning a tree

- The process described above may produce good predictions on the training set, but is likely to **overfit** the data, leading to poor test set performance. *Why?*
- A smaller tree with fewer splits (that is, fewer regions  $R_1, \dots, R_J$ ) might lead to lower variance and better interpretation at the cost of a little bias.
- One possible alternative to the process described above is to grow the tree only so long as the decrease in the RSS due to each split exceeds some (high) threshold.
- This strategy will result in smaller trees, but is too **short-sighted**: a seemingly worthless split early on in the tree might be followed by a very good split — that is, a split that leads to a large reduction in RSS later on.

## Choosing the best subtree

- The tuning parameter  $\alpha$  controls a trade-off between the subtree's complexity and its fit to the training data.
- We select an optimal value  $\hat{\alpha}$  using cross-validation.
- We then return to the full data set and obtain the subtree corresponding to  $\hat{\alpha}$ .

## Pruning a tree— continued

- A better strategy is to grow a very large tree  $T_0$ , and then **prune** it back in order to obtain a **subtree**.
- Cost complexity pruning** — also known as **weakest link pruning** — is used to do this.
- we consider a sequence of trees indexed by a nonnegative tuning parameter  $\alpha$ . For each value of  $\alpha$  there corresponds a subtree  $T \subset T_0$  such that

$$\sum_{m=1}^{|T|} \sum_{x_i \in R_m} (y_i - \hat{y}_{R_m})^2 + \alpha |T|$$

is as small as possible. Here  $|T|$  indicates the number of terminal nodes of the tree  $T$ ,  $R_m$  is the rectangle (i.e. the subset of predictor space) corresponding to the  $m$ th terminal node, and  $\hat{y}_{R_m}$  is the mean of the training observations in  $R_m$ .

## Summary: tree algorithm

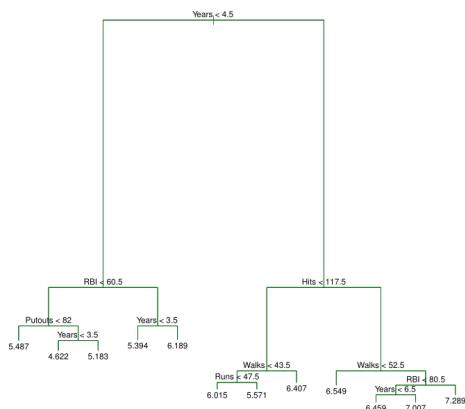
- Use recursive binary splitting to grow a large tree on the training data, stopping only when each terminal node has fewer than some minimum number of observations.
- Apply cost complexity pruning to the large tree in order to obtain a sequence of best subtrees, as a function of  $\alpha$ .
- Use K-fold cross-validation to choose  $\alpha$ . For each  $k = 1, \dots, K$ :
  - Repeat Steps 1 and 2 on the  $\frac{K-1}{K}$ th fraction of the training data, excluding the  $k$ th fold.
  - Evaluate the mean squared prediction error on the data in the left-out  $k$ th fold, as a function of  $\alpha$ .
- Average the results, and pick  $\alpha$  to minimize the average error.
- Return the subtree from Step 2 that corresponds to the chosen value of  $\alpha$ .

## Baseball example continued

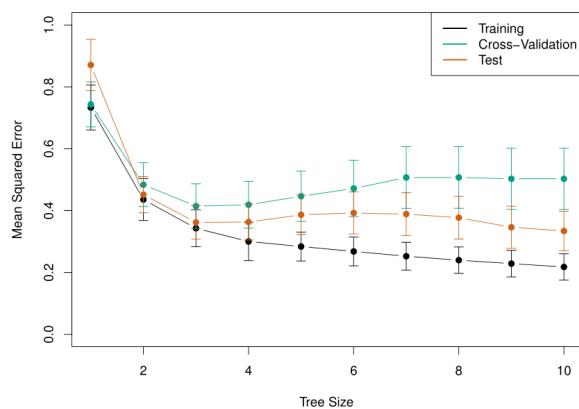
- First, we randomly divided the data set in half, yielding 132 observations in the training set and 131 observations in the test set.
- We then built a large regression tree on the training data and varied  $\alpha$  in order to create subtrees with different numbers of terminal nodes.
- Finally, we performed six-fold cross-validation in order to estimate the cross-validated MSE of the trees as a function of  $\alpha$ .

21 / 51

## Baseball example continued



## Baseball example continued



22 / 51

23 / 51

## Classification Trees

- Very similar to a regression tree, except that it is used to predict a qualitative response rather than a quantitative one.
- For a classification tree, we predict that each observation belongs to the *most commonly occurring class* of training observations in the region to which it belongs.

## Details of classification trees

- Just as in the regression setting, we use recursive binary splitting to grow a classification tree.
- In the classification setting, RSS cannot be used as a criterion for making the binary splits.
- A natural alternative to RSS is the *classification error rate*. This is simply the fraction of the training observations in that region that do not belong to the most common class:

$$E = 1 - \max_k (\hat{p}_{mk}).$$

Here  $\hat{p}_{mk}$  represents the proportion of training observations in the  $m$ th region that are from the  $k$ th class.

- However classification error is not sufficiently sensitive for tree-growing, and in practice two other measures are preferable.

## Gini index and Deviance

- The *Gini index* is defined by

$$G = \sum_{k=1}^K \hat{p}_{mk}(1 - \hat{p}_{mk}),$$

a measure of total variance across the  $K$  classes. The Gini index takes on a small value if all of the  $\hat{p}_{mk}$ 's are close to zero or one.

- For this reason the Gini index is referred to as a measure of node *purity* — a small value indicates that a node contains predominantly observations from a single class.
- An alternative to the Gini index is *cross-entropy*, given by

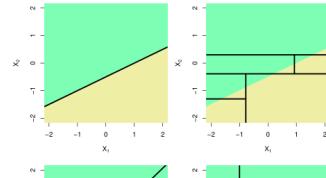
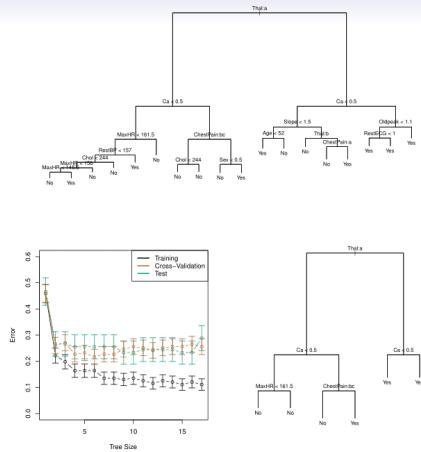
$$D = - \sum_{k=1}^K \hat{p}_{mk} \log \hat{p}_{mk}.$$

- It turns out that the Gini index and the cross-entropy are very similar numerically.

## Example: heart data

- These data contain a binary outcome *HD* for 303 patients who presented with chest pain.
- An outcome value of *Yes* indicates the presence of heart disease based on an angiographic test, while *No* means no heart disease.
- There are 13 predictors including *Age*, *Sex*, *Chol* (a cholesterol measurement), and other heart and lung function measurements.
- Cross-validation yields a tree with six terminal nodes. See next figure.

## Trees Versus Linear Models



Top Row: True linear boundary; Bottom row: true non-linear boundary.

Left column: linear model; Right column: tree-based model

28 / 51

## Advantages and Disadvantages of Trees

- ▲ Trees are very easy to explain to people. In fact, they are even easier to explain than linear regression!
- ▲ Some people believe that decision trees more closely mirror human decision-making than do the regression and classification approaches seen in previous chapters.
- ▲ Trees can be displayed graphically, and are easily interpreted even by a non-expert (especially if they are small).
- ▲ Trees can easily handle qualitative predictors without the need to create dummy variables.
- ▼ Unfortunately, trees generally do not have the same level of predictive accuracy as some of the other regression and classification approaches seen in this book.

However, by aggregating many decision trees, the predictive performance of trees can be substantially improved. We introduce these concepts next.

30 / 51

$$\begin{aligned}
 X &= \mu + \varepsilon \\
 \mathbb{E}[\varepsilon] &= 0 & \mathbb{E}[\bar{X}_m] &= \mu \\
 X_1: & \mathbb{E}[X_1] = \mu & \text{Var}(\bar{X}_m) &= \sigma^2/m \\
 & \text{Var}(X_1) = \sigma^2 & X_1 \rightarrow f_1 & \text{reduces variance?} \\
 X_2 \dots X_m \text{ iid} & & X_2 \rightarrow f_2 & \dots \\
 & & X_B \rightarrow f_B & \dots \\
 & & \downarrow \text{They are less} & \text{Bootstrap!} \\
 & & \text{various} & \\
 \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} \rightarrow f & & \text{Sample with replacement} & \\
 & & \left[ \begin{bmatrix} x_1^{*} \\ x_2^{*} \\ \vdots \\ x_m^{*} \end{bmatrix} \begin{bmatrix} y_1^{*} \\ y_2^{*} \\ \vdots \\ y_m^{*} \end{bmatrix} \rightarrow f^{*} \right] & \text{Repeat } B \text{ times} \\
 & & f_1^{*}, f_2^{*}, \dots, f_B^{*} & \xrightarrow{\text{I have something}} \text{that mimics the dist. of the regressors} \\
 & & \frac{1}{B} \sum_{i=1}^B f_i^{*} & \\
 \text{Is it good?} & & & \\
 \text{► these are NOT independent since they are generated from the same dataset} & & & \\
 P[u \notin X^*] &= (1 - \frac{1}{m})^m \rightarrow e^{-1} & & \\
 \text{With prob. } 2/3 \text{ the same units in the training set} & & & \\
 \hookrightarrow \text{the } 1/3 \text{ data you didn't use} \Rightarrow \text{TEST} & & &
 \end{aligned}$$

## Bagging

- *Bootstrap aggregation*, or *bagging*, is a general-purpose procedure for reducing the variance of a statistical learning method; we introduce it here because it is particularly useful and frequently used in the context of decision trees.
- Recall that given a set of  $n$  independent observations  $Z_1, \dots, Z_n$ , each with variance  $\sigma^2$ , the variance of the mean  $\bar{Z}$  of the observations is given by  $\sigma^2/n$ .
- In other words, *averaging a set of observations reduces variance*. Of course, this is not practical because we generally do not have access to multiple training sets.

## Bagging—continued

- Instead, we can bootstrap, by taking repeated samples from the (single) training data set.
- In this approach we generate  $B$  different bootstrapped training data sets. We then train our method on the  $b$ th bootstrapped training set in order to get  $f^{*b}(x)$ , the prediction at a point  $x$ . We then average all the predictions to obtain

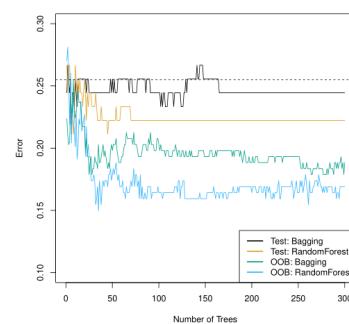
$$\hat{f}_{\text{bag}}(x) = \frac{1}{B} \sum_{b=1}^B \hat{f}^{*b}(x).$$

This is called *bagging*.

## Bagging classification trees

- The above prescription applied to regression trees
- For classification trees: for each test observation, we record the class predicted by each of the  $B$  trees, and take a *majority vote*: the overall prediction is the most commonly occurring class among the  $B$  predictions.

## Bagging the heart data



33 / 51

34 / 51



## Details of previous figure

Bagging and random forest results for the Heart data.

- The test error (black and orange) is shown as a function of  $B$ , the number of bootstrapped training sets used.
- Random forests were applied with  $m = \sqrt{p}$ .
- The dashed line indicates the test error resulting from a single classification tree.
- The green and blue traces show the OOB error, which in this case is considerably lower

## Out-of-Bag Error Estimation

- It turns out that there is a very straightforward way to estimate the test error of a bagged model.
- Recall that the key to bagging is that trees are repeatedly fit to bootstrapped subsets of the observations. One can show that on average, each bagged tree makes use of around two-thirds of the observations.
- The remaining one-third of the observations not used to fit a given bagged tree are referred to as the *out-of-bag* (OOB) observations.
- We can predict the response for the  $i$ th observation using each of the trees in which that observation was OOB. This will yield around  $B/3$  predictions for the  $i$ th observation, which we average.
- This estimate is essentially the LOO cross-validation error for bagging, if  $B$  is large.

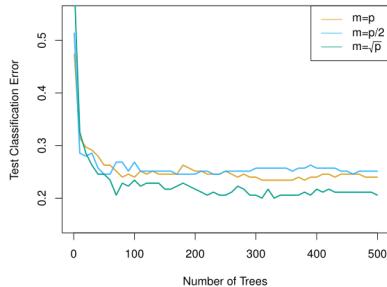
## Random Forests

- Random forests* provide an improvement over bagged trees by way of a small tweak that *decorrelates* the trees. This reduces the variance when we average the trees.
- As in bagging, we build a number of decision trees on bootstrapped training samples.
- But when building these decision trees, each time a split in a tree is considered, *a random selection of  $m$  predictors* is chosen as split candidates from the full set of  $p$  predictors. The split is allowed to use only one of those  $m$  predictors.
- A fresh selection of  $m$  predictors is taken at each split, and typically we choose  $m \approx \sqrt{p}$  — that is, the number of predictors considered at each split is approximately equal to the square root of the total number of predictors (4 out of the 13 for the Heart data).

## Example: gene expression data

- We applied random forests to a high-dimensional biological data set consisting of expression measurements of 4,718 genes measured on tissue samples from 349 patients.
- There are around 20,000 genes in humans, and individual genes have different levels of activity, or expression, in particular cells, tissues, and biological conditions.
- Each of the patient samples has a qualitative label with 15 different levels: either normal or one of 14 different types of cancer.
- We use random forests to predict cancer type based on the 500 genes that have the largest variance in the training set.
- We randomly divided the observations into a training and a test set, and applied random forests to the training set for three different values of the number of splitting variables  $m$ .

## Results: gene expression data



## Details of previous figure

- Results from random forests for the fifteen-class gene expression data set with  $p = 500$  predictors.
- The test error is displayed as a function of the number of trees. Each colored line corresponds to a different value of  $m$ , the number of predictors available for splitting at each interior tree node.
- Random forests ( $m < p$ ) lead to a slight improvement over bagging ( $m = p$ ). A single classification tree has an error rate of 45.7%.

## Boosting

- Like bagging, boosting is a general approach that can be applied to many statistical learning methods for regression or classification. We only discuss boosting for decision trees.
- Recall that bagging involves creating multiple copies of the original training data set using the bootstrap, fitting a separate decision tree to each copy, and then combining all of the trees in order to create a single predictive model.
- Notably, each tree is built on a bootstrap data set, independent of the other trees.
- Boosting works in a similar way, except that the trees are grown *sequentially*: each tree is grown using information from previously grown trees.

## Boosting algorithm for regression trees

- Set  $\hat{f}(x) = 0$  and  $r_i = y_i$  for all  $i$  in the training set.
- For  $b = 1, 2, \dots, B$ , repeat:
  - Fit a tree  $\hat{f}^b$  with  $d$  splits ( $d + 1$  terminal nodes) to the training data  $(X, r)$ .
  - Update  $\hat{f}$  by adding in a shrunken version of the new tree:

$$\hat{f}(x) \leftarrow \hat{f}(x) + \lambda \hat{f}^b(x).$$

- Update the residuals,

$$r_i \leftarrow r_i - \lambda \hat{f}^b(x_i).$$

- Output the boosted model,

$$\hat{f}(x) = \sum_{b=1}^B \lambda \hat{f}^b(x).$$



## What is the idea behind this procedure?

- Unlike fitting a single large decision tree to the data, which amounts to *fitting the data hard* and potentially overfitting, the boosting approach instead *learns slowly*.
- Given the current model, we fit a decision tree to the residuals from the model. We then add this new decision tree into the fitted function in order to update the residuals.
- Each of these trees can be rather small, with just a few terminal nodes, determined by the parameter  $d$  in the algorithm.
- By fitting small trees to the residuals, we slowly improve  $f$  in areas where it does not perform well. The shrinkage parameter  $\lambda$  slows the process down even further, allowing more and different shaped trees to attack the residuals.

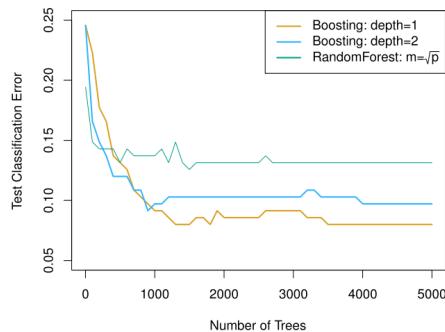
## Boosting for classification

- Boosting for classification is similar in spirit to boosting for regression, but is a bit more complex. We will not go into detail here, nor do we in the text book.
- Students can learn about the details in *Elements of Statistical Learning, chapter 10*.
- The R package `gbm` (gradient boosted models) handles a variety of regression and classification problems.

43 / 51

44 / 51

## Gene expression data continued



## Details of previous figure

- Results from performing boosting and random forests on the fifteen-class gene expression data set in order to predict *cancer* versus *normal*.
- The test error is displayed as a function of the number of trees. For the two boosted models,  $\lambda = 0.01$ . Depth-1 trees slightly outperform depth-2 trees, and both outperform the random forest, although the standard errors are around 0.02, making none of these differences significant.
- The test error rate for a single tree is 24%.

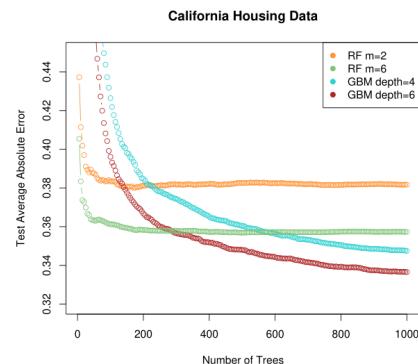
45 / 51

46 / 51

## Tuning parameters for boosting

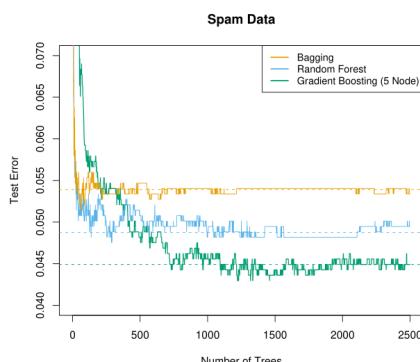
1. The *number of trees*  $B$ . Unlike bagging and random forests, boosting can overfit if  $B$  is too large, although this overfitting tends to occur slowly if at all. We use cross-validation to select  $B$ .
2. The *shrinkage parameter*  $\lambda$ , a small positive number. This controls the rate at which boosting learns. Typical values are 0.01 or 0.001, and the right choice can depend on the problem. Very small  $\lambda$  can require using a very large value of  $B$  in order to achieve good performance.
3. The *number of splits*  $d$  in each tree, which controls the complexity of the boosted ensemble. Often  $d = 1$  works well, in which case each tree is a *stump*, consisting of a single split and resulting in an additive model. More generally  $d$  is the *interaction depth*, and controls the interaction order of the boosted model, since  $d$  splits can involve at most  $d$  variables.

## Another regression example



from *Elements of Statistical Learning, chapter 15*.

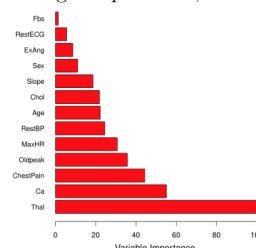
## Another classification example



from *Elements of Statistical Learning, chapter 15*.

## Variable importance measure

- For bagged/RF regression trees, we record the total amount that the RSS is decreased due to splits over a given predictor, averaged over all  $B$  trees. A large value indicates an important predictor.
- Similarly, for bagged/RF classification trees, we add up the total amount that the Gini index is decreased by splits over a given predictor, averaged over all  $B$  trees.



Variable importance plot for the *Heart* data

49 / 51

50 / 51



- Decision trees are simple and interpretable models for regression and classification
- However they are often not competitive with other methods in terms of prediction accuracy
- Bagging, random forests and boosting are good methods for improving the prediction accuracy of trees. They work by growing many trees on the training data and then combining the predictions of the resulting ensemble of trees.
- The latter two methods— random forests and boosting—are among the state-of-the-art methods for supervised learning. However their results can be difficult to interpret.

51 / 51

## LECTURE 16/5/2022

### MIXED EFFECT MODELS

#### ► MOTIVATIONS

##### General overview

- Regression models are a very common statistical tool.  
A very important assumption of regression models is the **independence and homoscedasticity** of the observations.
- *How often such assumptions are met in the real world?* Seldom, to be optimistic.  
=> a more general framework to properly account for complex structure present into the data is needed.
- Modeling general dependence and heteroscedasticity among observations concerns the **design of the variance covariance matrix** of the errors.

Ex: It is often the case that statistical units are **correlated and/or grouped within a hierarchy**:

- Measurements of the same individuals over time
- Patients sharing the same GP, hospital of admission, district, etc.
- Pupils grouped into classes, schools, districts, etc.
- Units belonging to close districts.



##### General overview

- Observations **within groups are more similar (possibly correlated)** than observations between groups
  - Independence hypothesis does not hold
  - Ignoring the dependence structure induces biased estimates of parameters
- **Linear Mixed Models – LMMs** (a.k.a. **Random Effects Models** or **Multilevel/Hierarchical models**) are a generalization of traditional regression models (**Fixed Effect Models**), adding to the linear predictor (a.k.a. fixed component) a random component.
- LMMs allow to **disentangle the contribution of different kind of dependencies among observations**, focusing on inference related to the presence of groups.



## Motivating examples



### 1. ARMD - Age-Related Macular Degeneration Trial

- The ARMD data arise from a randomized multi-center clinical trial comparing an experimental treatment (interferon- $\alpha$ ) versus placebo for patients diagnosed with ARMD.
- We focus on the comparison between placebo and the highest dose (6 million units daily) of interferon- $\alpha$ .
- Patients with macular degeneration progressively lose vision.
- Visual acuity of each of 240 patients was assessed at baseline and at four post-randomization timepoints (4, 12, 24, and 52 weeks).
- Visual acuity was evaluated based on patient's ability to read lines of letters on standardized vision charts. The charts display lines of five letters of decreasing size, which the patient must read from top (largest letters) to bottom (smallest letters).
- In our analyses, we will focus on the visual acuity defined as the total number of letters correctly read.

=> longitudinal data in the form of up to five visual acuity measurements collected at different, but common to all patients, timepoints.



Long & short format



## Motivating examples



### 2. Multi center trials

- Many trials are multi-center, due to the inadequate number of patients in a single center (rare disease) or to shared studies and protocols.
  - Often the analysis ignores the center from which the data were obtained, making the implicit assumption that centers are identical. But they are not: differences usually arise in a) case mix, b) overall success in recruitment and outcome and c) relative benefit (hospital effect).
  - Observational multicenter study (**N=48 hospitals**) on primary lung cancer.
  - Investigation of clinical factors associated to 3y death for any cause in **802 pts** undergoing surgery.
- **Goals:**
- Risk stratification of the population under study  
=> association between mortality and features measured at baseline/entrance
  - Assessment of the grouper effect and scenario analysis



## Motivating examples



### 3. Schools

- Nowadays, Italian (but also international) students are tested by means of standardized tests given at different grades.
  - The **school dataset** collects information about the tests of 1000 students enrolled in 50 different primary schools.
  - For each student, beside the test result, we observe the gender, the socioeconomic index and the anonymous id of the school in which he/she is enrolled.
  - Being enrolled in a school might have a relevant effect on student test scores (e.g. very bad/good teachers) → students are not independent
  - **n = 1000 students** within **N = 50 schools**
  - Investigation of the association between student-level characteristics and student test scores
  - Investigation of the school effect on student test scores
- **Goals:**
- Prediction of student test scores  
=> association between test score and gender and socioeconomic index of the student
  - Assessment of the school effect and scenario analysis



- ▶ LM
- ▶ LM 2.0 relaxing homogeneity
- ▶ LM 3.0 relaxing independence
- ▶ LMM

- LMs are used to quantify the relationship between a dependent variable and a set of covariates with the use of a linear function depending on a (possibly) small number of regression parameters.
- LMs are suitable for analyzing data involving **independent observations with a homogenous variance** (e.g., standard linear regression, ANOVA/ANCOVA models).
- The classical LM for **independent, normally distributed observations**  $y_j, j = 1, \dots, n$  with a constant variance can be specified in a variety of ways. A commonly used specification is:

Model equation  
at the level of  
observations

$$y_j = \beta_0 + \beta_1 z_{1j} + \dots + \beta_p z_{pj} + \varepsilon_j = \mathbf{z}_j^t \boldsymbol{\beta} + \varepsilon_j$$

$$\varepsilon_j \sim N(0, \sigma^2) \quad j \perp j' \text{ independent errors}$$



Note that  $E[y_j] = \mu_j = \mathbf{z}_j^t \boldsymbol{\beta}$

$$V[y_j] = V[\varepsilon_j] = \sigma^2$$

13



- LMs are used to quantify the relationship between a dependent variable and a set of covariates with the use of a linear function depending on a (possibly) small number of regression parameters.
- LMs are suitable for analyzing data involving **independent observations with a homogenous variance** (e.g., standard linear regression, ANOVA/ANCOVA models).
- The classical LM for **independent, normally distributed observations**  $y_j, j = 1, \dots, n$  with a constant variance can be specified in a variety of ways. A commonly used specification is:

Model equation  
for all data

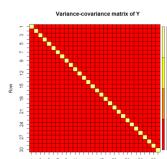
$$\mathbf{y} = \mathbf{Z}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

$$\boldsymbol{\varepsilon} \sim N(\mathbf{0}, R) \quad R = \sigma^2 I_n \in \mathbb{R}^{n \times n}$$

$$\mathbf{y} = (y_1, \dots, y_n)' \quad \boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)'$$

$$\mathbf{Z} = \begin{pmatrix} \mathbf{z}_1 \\ \mathbf{z}_2 \\ \vdots \\ \mathbf{z}_n \end{pmatrix} = \begin{pmatrix} 1 & z_{11} & \dots & z_{p1} \\ \dots & \dots & \dots & \dots \\ 1 & z_{1n} & \dots & z_{pn} \end{pmatrix} = (\mathbf{1}' \quad \mathbf{z}_1' \quad \dots \quad \mathbf{z}_p') \in \mathbb{R}^{n \times (p+1)}$$

14



## LM - estimation

- **Goal:** finding estimates of a set of parameters  $\boldsymbol{\beta}$  and  $\sigma^2$ .
- Estimation via least squares (OLS)
- However, OLS is less suitable for more complex LMs, including LMMs.

$$\hat{\boldsymbol{\beta}}_{OLS} = (\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{y}$$

Note: OLS estimate does not require the normality assumption

Valid under the assumption of uncorrelated residual errors

$$\hat{\sigma}_{OLS}^2 = \frac{1}{n-(p+1)} (\mathbf{y} - \mathbf{Z}\hat{\boldsymbol{\beta}}_{OLS})'(\mathbf{y} - \mathbf{Z}\hat{\boldsymbol{\beta}}_{OLS})$$

Note: OLS estimates are **unbiased**



15

## LM - estimation

### Maximum Likelihood (ML) estimation

- The likelihood function for a Normal LM is

$$L_{full}(\boldsymbol{\beta}, \sigma^2; \mathbf{y}) = (2\pi\sigma^2)^{-n/2} \prod_{j=1}^n \exp \left[ -\frac{(y_j - \mathbf{z}'_j \boldsymbol{\beta})^2}{2\sigma^2} \right]$$
$$l_{full}(\boldsymbol{\beta}, \sigma^2; \mathbf{y}) = -\frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_{j=1}^n (y_j - \mathbf{z}'_j \boldsymbol{\beta})^2$$

- Maximization of  $l_{full}$  provides the ML estimates of unknown parameters

$$\hat{\boldsymbol{\beta}}_{ML} = (Z'Z)^{-1} Z' \mathbf{y} \quad \text{Same as OLS}$$

$$\hat{\sigma}_{ML}^2 = \frac{1}{n} (\mathbf{y} - Z\hat{\boldsymbol{\beta}}_{ML})' (\mathbf{y} - Z\hat{\boldsymbol{\beta}}_{ML}) = \frac{1}{n} \sum_{j=1}^n (y_j - \mathbf{z}'_j \hat{\boldsymbol{\beta}}_{ML})^2 \quad \text{Biased!}$$

16



## LM - estimation

### Restricted Maximum Likelihood (REML) estimation

- To obtain an unbiased estimate for  $\sigma^2$ , an approach that is orthogonal to the estimation of  $\boldsymbol{\beta}$  is needed. This is possible considering the likelihood function based on a set of  $n-(p+1)$  independent contrast of  $\mathbf{y}$ .

$$l_{REML}(\sigma^2; \mathbf{y}) = -\frac{n-(p+1)}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_{j=1}^n r_j^2 \Rightarrow \hat{\sigma}_{REML}^2 = \frac{1}{n-(p+1)} \sum_{j=1}^n r_j^2$$

where  $r_j$  are the residuals.

- OLS estimators are equivalent to the REML estimates.

This equivalence holds for classical LM with independent, homoscedastic errors, NOT for more complex formulations.

- The REML objective function does not allow to directly estimate the coefficients  $\boldsymbol{\beta}$ . The ML formula should be used in this case. Also in this case, the equivalence between ML and REML estimates is true for the classical ML only.

17



## ► LM 2.0 relaxing homogeneity

## LM 2.0

- In the previous case, we formulated the classical LM for independent observations.

⇒ key assumptions:

1. observations are independent and normally distributed with a constant, i.e., homogeneous variance
2. the expected value of the observations can be expressed as a linear function of covariates.

- We now relax the homoscedasticity assumption and allow for the observations to be heteroscedastic, i.e., to have different variances, while retaining the assumption that the observations are independent and normally distributed.

⇒ LMs with heterogeneous variance

- Important concept: variance function

- The new setting will ask for suitable estimation methods (e.g WLS, GLS and IRLS estimation).

19



➤ In the classical LM with homogeneous variance, the variance of the dependent variable is  $V[y_j] = \sigma^2$

➤ We now relax the constant variance assumption and assume that  $V[y_j] = \sigma_j^2$

Therefore:

Model equation  
at the level of  
observations

$$y_j = \beta_0 + \beta_1 z_{1j} + \dots + \beta_p z_{pj} + \varepsilon_j = \mathbf{z}_j^t \boldsymbol{\beta} + \varepsilon_j \quad j = 1, \dots, n$$

$$\varepsilon_j \sim N(0, \sigma_j^2) \quad j \perp j' \text{ independent errors}$$



Note that

$$E[y_j] = \mu_j = \mathbf{z}_j^t \boldsymbol{\beta}$$

$$V[y_j] = V[\varepsilon_j] = \sigma_j^2$$



- ❖ The model contains  $n+p+1 > n$  parameters => the model is not identifiable!
- ❖ It may become identifiable if we impose additional constraints on the residual variances
  1. Assume known variance weights
  2. More general way: to represent variances more parsimoniously as a function of a small set of parameters => variance function.

20



## LM 2.0 – Variance Function

➤ A more general and flexible way to introduce variance heterogeneity is by means of a **variance function**

$$\lambda(\delta, \mu; v)$$

which assumes positive values and is continuous and differentiable with respect to  $\delta$ .

➤ Therefore:

$$V[\varepsilon_j] = \sigma^2 \lambda^2(\delta, \mu_j; v_j)$$

where ->  $v_j$  is a **vector of (known) covariates** defining the variance function for observation  $j$ ,  
->  $\delta$  contains a **small set of variance parameters**, common to all observations.

Note that:

- ❖ because the function  $\lambda(\cdot)$  involves  $\mu_j$ , it in fact depends on  $\boldsymbol{\beta}$ , too
- ❖ the parameter  $\sigma$  should be interpreted as a scale parameter only and no more as residual error standard deviation.

21



## LM 2.0 – Variance Function

Therefore:

Model equation  
at the level of  
observations

$$y_j = \beta_0 + \beta_1 z_{1j} + \dots + \beta_p z_{pj} + \varepsilon_j = \mathbf{z}_j^t \boldsymbol{\beta} + \varepsilon_j$$

$$\varepsilon_j \sim N(0, \sigma^2 \lambda_j^2)$$

where  $\lambda_j^2 = \lambda^2(\delta, \mu_j; v_j)$  and  $j \perp j'$  independent errors

Note that

$$E[y_j] = \mu_j = \mathbf{z}_j^t \boldsymbol{\beta}$$

$$V[y_j] = V[\varepsilon_j] = \sigma^2 \lambda_j^2 \longrightarrow$$

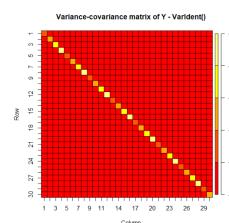
Model equation  
at the level of  
observations

$$\mathbf{y} = \mathbf{Z}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

$$\boldsymbol{\varepsilon} \sim N(\mathbf{0}, R) \quad R = \sigma^2 \Lambda \Lambda \quad \text{where } \Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$$

$$R \in \mathbb{R}^{n \times n}$$

22



## LM 2.0 – Variance Function

Examples of the variance function  $\lambda(\cdot)$ :

1. Known weights,  $\lambda(\cdot) = \lambda(v)$
2. Variance functions depending on  $\delta$  but not on  $\mu$  =>  $\lambda(\cdot) = \lambda(\delta; v)$
3. Variance functions depending on  $\delta$  and  $\mu$  =>  $\lambda(\cdot) = \lambda(\delta, \mu; v)$
4. Variance functions depending on  $\mu$  but not on  $\delta$  =>  $\lambda(\cdot) = \lambda(\mu; v)$

➤ We will symbolically refer to groups 2–4 as  $\langle\delta\rangle$ ,  $\langle\delta, \mu\rangle$ , and  $\langle\mu\rangle$ -group, respectively.

➤ Specification of an LM with heterogeneous variance is very general and encompasses all four groups of variance functions => the use of a variance function from any of the aforementioned groups does not pose difficulties in terms of the model specification.

However, in models involving variance functions from groups  $\langle\delta, \mu\rangle$  or  $\langle\mu\rangle$ , the parameters  $\beta$  are shared by the mean and variance structures (*mean-variance models*) => require different estimation approaches and inference techniques, as compared to the models involving known weights or variance functions from the  $\langle\delta\rangle$ -group.

23



## LM 2.0 – Variance Function

**Table 7.1** A summary of the parts of Chap. 7 that contain the information about particular groups of variance functions and the corresponding estimation methods

Group	Arguments	Examples	Estimation algorithm	Section
Known weights	— —	varFixed(·)	WLS	7.4.1
$\langle\delta\rangle$	— —	Table 7.2	ML/REML	7.4.2
$\langle\delta, \mu\rangle$	— +	Table 7.3	ML/REML-based GLS	7.8.1.1
$\langle\mu\rangle$	+ —	Table 7.4	IRLS	7.8.1.2

**Table 7.2** Examples of variance functions from the  $\langle\delta\rangle$ -group<sup>a</sup>

Function $\lambda(\cdot)$	$\lambda_i$	Description
varPower( $\delta; v_i, s_i$ )	$ v_i ^{\delta_i}$	Power of a variance covariate $v_i$
varExp( $\delta; v_i, s_i$ )	$\exp(v_i \delta_i)$	Exponent of a variance covariate
varConstPower( $\delta; v_i, s_i$ )	$\delta_{1,s_i} +  v_i ^{\delta_{2,s_i}}$	Constant plus power variance function $\delta_{1,s_i} > 0$
varIdent( $\delta; s_i$ )	$\delta_i$	Different variances per stratum $\delta_i \equiv 1, \delta_i > 0$ for $s \neq 1$

$s_i$  = stratum the  $j$ -th obs belongs to

**Table 7.3** Examples of variance functions from the  $\langle\delta, \mu\rangle$ -group<sup>a</sup>

Function $\lambda(\cdot)$	$\lambda_i$	Description
varPower( $\delta, \mu_i; s_i$ )	$ \mu_i ^{\delta_i}$	Power of $ \mu_i $
varExp( $\delta, \mu_i; s_i$ )	$\exp(\mu_i \delta_i)$	Exponent of $\mu_i$
varConstPower( $\delta, \mu_i; s_i$ )	$\delta_{1,s_i} +  \mu_i ^{\delta_{2,s_i}}$	Constant plus power variance function $\delta_{1,s_i} > 0$

**Table 7.4** Examples of variance functions from the  $\langle\mu\rangle$ -group<sup>a</sup>

Function $\lambda(\cdot)$	$\lambda_i$	Description
varPower( $\mu_i; s_i, \delta$ )	$ \mu_i ^{\delta_i}$	Power of $ \mu_i , \delta_i$ known
varExp( $\mu_i; s_i, \delta$ )	$\exp(\mu_i \delta_i)$	Exponent of $\mu_i, \delta_i$ known
varConstPower( $\mu_i; s_i, \delta$ )	$\delta_{1,s_i} +  \mu_i ^{\delta_{2,s_i}}$	Constant plus power variance function, $\delta_{1,s_i} > 0, \delta_{1,s_i}$ and $\delta_{2,s_i}$ known

## LM 2.0 – Estimation

### Profiled Likelihood

➤ The log-likelihood function for the LM 2.0 model is:

$$l_{full}(\beta, \sigma^2, \delta) = -\frac{n}{2} \log(\sigma^2) - \frac{1}{2} \sum_{j=1}^n \log(\lambda_j^2) - \frac{1}{2\sigma^2} \sum_{j=1}^n \lambda_j^{-2} (\mathbf{y}_j - \mathbf{x}'_j \beta)^2$$

Note that it depends on  $\delta$  through  $\lambda$ , then in the case of  $\lambda=1$  we step back to the ML of the usual LM case.

➤ Profiling of a likelihood function can be done in a variety of ways. Here we follow the profiling approach implemented in the `gls()` function of the `nlme` package.

=> We first profile out the  $\beta$  parameters, and then, we profile out  $\sigma^2$ .

➤ The advantage of using the function is that it does not depend on all the initial parameters.  
Thus, optimization of the function is performed in a parameter space of a lower dimension.

Assume that  $\delta$  is known. Then,

$$\begin{aligned} & \text{maximize } l_{full} \text{ wrt } \beta \text{ for any value of } \delta \Rightarrow \hat{\beta}(\delta) \Rightarrow l_{ML}^*(\sigma^2, \delta) \\ & \Rightarrow \text{maximize } l_{ML}^* \text{ wrt } \sigma^2 \text{ for any value of } \delta \Rightarrow \hat{\sigma}_{ML}^2(\delta) \Rightarrow l_{ML}^*(\delta) \end{aligned}$$

25



## LM 2.0 – Estimation

### Profiled Likelihood

- ML estimates of unknown parameters of LM 2.0:

$$\hat{\beta}_{ML} = \hat{\beta}(\hat{\delta}_{ML}) = (\sum_{j=1}^n \hat{\lambda}_j^{-2} \mathbf{x}_j \mathbf{x}'_j)^{-1} \sum_{j=1}^n \hat{\lambda}_j^{-2} \mathbf{x}_j y_j \quad \hat{\sigma}_{ML}^2 = \hat{\sigma}^2(\hat{\delta}_{ML}) = \sum_{j=1}^n \hat{\lambda}_j^{-2} \hat{r}_j^2 y_j$$

$$\hat{\lambda}_j = \lambda(\hat{\delta}_{ML}; v_j) \quad \hat{r}_j = r_j(\hat{\delta}_{ML})$$

- Biased => REML approach

$$l_{REML}(\sigma^2, \delta) = -\frac{n-(p+1)}{2} \log(\sigma^2) - \frac{1}{2} \sum_{j=1}^n \log(\lambda_j^2) - \frac{1}{2\sigma^2} \sum_{j=1}^n \lambda_j^{-2} r_j^2 - \frac{1}{2} \log \left[ \det \left( \sum_{j=1}^n \lambda_j^{-2} \mathbf{x}_j \mathbf{x}'_j \right) \right]$$

$$\Rightarrow \hat{\sigma}_{REML}^2(\hat{\delta}) = \sum_{j=1}^n \hat{\lambda}_j^{-2} \hat{r}_j^2 / (n-p-1)$$

- By maximization of  $l_{REML}$  with respect to  $\delta$ , we obtain an estimator  $\hat{\delta}_{REML}$  of  $\delta$ .

This is used to yield an estimator  $\hat{\beta}_{REML}$  of  $\beta$ .



26

## ► LM 3.0 relaxing independence

## LM 3.0

- The essential assumption for the LMs considered before was that the **observations** collected during the study were **independent of each other**.
- This **assumption is restrictive** in studies which use sampling designs that lead to correlated data.
  - Studies collecting measures over time, i.e., in a longitudinal fashion;
  - Designs involving hierarchies or grouping (e.g., cluster-randomization clinical trials; in studies collecting spatially correlated data, etc.)
- Note that for such designs, the distinction between sampling units (e.g., subjects in a longitudinal study) and analysis units (e.g., time-specific measurements) is important.
- We now consider a class of **more general LMs** that allow **relaxing the assumptions of independence**, namely **LMs with fixed effects and correlated residual errors for grouped data**, or simply as **LMs for correlated data**.
- These models can be viewed as an example of population-averaged models, i.e., models in which the parameters are interpreted as quantifying effects of covariates on the marginal mean value of the dependent variable for the entire population.
- Important concept: ***correlation structure***.



28

## LM 3.0

- We now introduce LM with fixed effects and correlated residual errors for grouped data with hierarchical structure.
  - Single-level of grouping, with N groups (levels of a grouping factor) indexed by  $i$  ( $i = 1, \dots, N$ )
  - $n_i$  observations per group indexed by  $j$  ( $j = 1, \dots, n_i$ ).

- Then, for the group  $i$

$$\mathbf{y}_i = \mathbf{Z}_i \boldsymbol{\beta} + \boldsymbol{\varepsilon}_i \quad \boldsymbol{\varepsilon}_i \sim N_{n_i}(\mathbf{0}, R_i) \quad R_i = \sigma^2 \mathbf{R}_i$$

**Model equation  
at the level of group  
of observations**

$$\mathbf{y}_i = (y_{i1}, \dots, y_{in_i})' \quad \boldsymbol{\varepsilon}_i = (\varepsilon_{i1}, \dots, \varepsilon_{in_i})'$$

Multiple level of  
grouping is possible,  
but requires suitable  
design of Z matrix  
and software  
specifications.

$$\mathbf{Z}_i = \begin{pmatrix} \mathbf{z}_{i1} \\ \mathbf{z}_{ij} \\ \vdots \\ \mathbf{z}_{in_i} \end{pmatrix} = \begin{pmatrix} 1 & z_{i11} & \dots & z_{ip1} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & z_{i1n_i} & \dots & z_{ipn_i} \end{pmatrix} = (\mathbf{1}' \quad \mathbf{z}'_{i1} \quad \dots \quad \mathbf{z}'_{ip})$$

Note that

$$E[y_{ij}] = \mu_{ij} = \mathbf{z}_{ij}' \boldsymbol{\beta} \quad V[\mathbf{y}_i] = V[\boldsymbol{\varepsilon}_i] = \sigma^2 \mathbf{R}_i$$



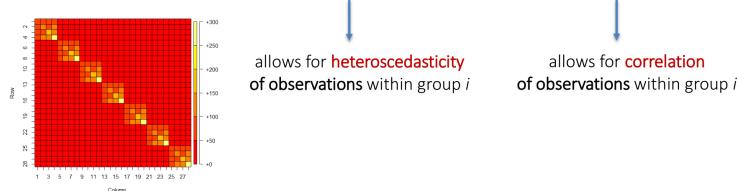
29



## LM 3.0 – Correlation function

- This model is not identifiable in its most general form due to i) the non uniqueness of the  $R_i$  representation and ii) the model potentially involves too many unknown parameters.  
=> additional constraints

$$R_i = \sigma^2 R_i = \sigma^2 \Lambda_i C_i \Lambda_i \quad \text{where} \quad \Lambda_i = (\lambda_{i1}, \dots, \lambda_{in_i}) \quad \text{and} \quad C_i = \text{corr matrix}$$



- Correlation function (general version):  $\text{Corr}(\varepsilon_{ij}, \varepsilon'_{ij}) = h[d(t_{ij}, t'_{ij}), \boldsymbol{\varrho}]$  where
- $\boldsymbol{\varrho}$  is a vector of correlation parameters
  - $d(\cdot)$  is a distance function of vectors of position variables  $t_{ij}$  and  $t'_{ij}$  corresponding to,  $\varepsilon_{ij}$  and  $\varepsilon'_{ij}$
  - $h[\cdot]$  is a continuous function with respect to  $\boldsymbol{\varrho}$ , ranging between -1 and 1 and s.t.  $h(0, \boldsymbol{\varrho}) = 1$ .

30



## LM 3.0 – Correlation function

- By assuming various distances and correlation functions, a variety of correlation structures can be obtained.
- The correlation structures can be classified into two main groups:

### 1. "Serial" structures

=> correlation structures which are defined in the context of time-series or longitudinal data.

Ex: Autocorrelation =>  $\text{Corr}(\varepsilon_{ij}, \varepsilon'_{ij}) = \rho$

### 2. "Spatial" structures

=> correlation structures which are defined in the context of spatial data.

**Table 10.1** Examples of serial and spatial correlation structures

Correlation structure	Function $h(\cdot, \cdot)$	Comment
Serial	(Auto)correlation function	
<i>corCompSymm</i> <sup>a</sup>	$h(k, \boldsymbol{\varrho}) \equiv \varrho^k$	$k = 1, 2, \dots;  \varrho  < 1$
<i>corARI</i>	$h(k, \boldsymbol{\varrho}) \equiv \varrho^k$	$k = 0, 1, \dots;  \varrho  < 1$
<i>corCAR1</i>	$h(s, \boldsymbol{\varrho}) \equiv \varrho^s$	$s \geq 0; \varrho \geq 0$
<i>corSymm</i>	$h(d(j, j'), \boldsymbol{\varrho}) \equiv \varrho_{j, j'}$	$j < j';  \varrho_{j, j'}  < 1$
Spatial	Correlation function	
<i>corExp</i>	$h(s, \boldsymbol{\varrho}) \equiv e^{-s/\varrho}$	$s \geq 0; \varrho > 0$
<i>corGaus</i>	$h(s, \boldsymbol{\varrho}) \equiv e^{-(s/\varrho)^2}$	$s \geq 0; \varrho > 0$
<i>corLin</i>	$h(s, \boldsymbol{\varrho}) \equiv (1 - s/\varrho)I(s < \varrho)$	$s \geq 0; \varrho > 0$
<i>corRatio</i>	$h(s, \boldsymbol{\varrho}) \equiv 1 - (s/\varrho)^2 / \{1 + (s/\varrho)^2\}$	$s \geq 0; \varrho > 0$
<i>corSpher</i>	$h(s, \boldsymbol{\varrho}) \equiv [1 - 1.5(s/\varrho) + 0.5(s/\varrho)^3]I(s < \varrho)$	$s \geq 0; \varrho > 0$

<sup>a</sup>The names of the structures follow the convention used in the **nlme** package

31



## LM 3.0 - Estimation

- Estimation carried out via WLS or likelihood based methods mentioned before.

32



## Why LMMs

- Linear mixed effects models (LMMs) are models in which a random component appears in the model function.
- Up to LM 3.0, we modeled the possible dependence among observations acting on the structure of the variance-covariance matrix of the errors.  
This enabled us to represent a wide range of possible dependence between observations in a flexible way, in particular the one related to the presence of groups/strata.  
=> PB: we have to assume exactly the kind of dependence existing among observations.
- LMMs enable us to disentangle the effect of the presence of groups from the general kind of dependence occurring among observations.  
=> they allow for partitioning of the global variance and for inference on the population of groupers the actual groups come from.
- More in general, LMMs represent the most popular and effective approach for the analysis of grouped data, enabling
  - o the estimation of group level effect
  - o the modelling of the dependence between observations not deriving by assumptions directly made on the var-cov matrix of the errors, but driven by assumptions on between groups variability.



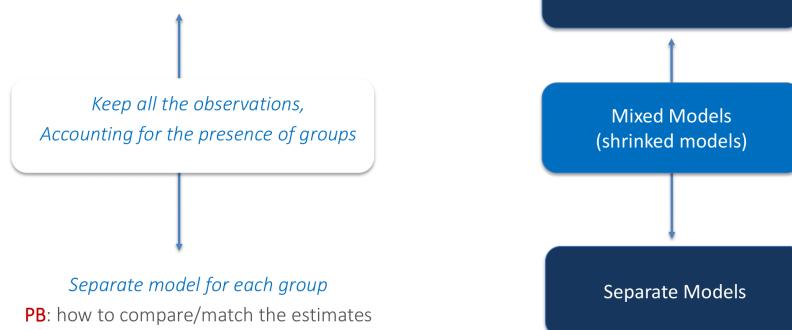
34



## LMM – Pooled, separate or to mixed effects estimates?

*Back to independent observations averaging over the groups*

PB: low number of observations and loosing information



## LMM – group-level specification

- For hierarchical data with a single level of grouping (N groups indexed by  $i = 1, \dots, N$ , with  $n_i$  observations per group indexed by  $j = 1, \dots, n_i$ ), we can formulate the classical LMM at a given level of a grouping factor as follows:

Model equation  
at the level of group  
of observations

$$\mathbf{y}_i = \mathbf{Z}_i \boldsymbol{\beta} + \mathbf{W}_i \mathbf{b}_i + \boldsymbol{\varepsilon}_i$$

$$\mathbf{b}_i \sim N_q(\mathbf{0}, D) \quad \boldsymbol{\varepsilon}_i \sim N_{n_i}(\mathbf{0}, R_i) \quad \text{with } \mathbf{b}_i \perp \boldsymbol{\varepsilon}_i \quad \forall i \neq i'$$

$$D = \sigma^2 D \quad \text{and} \quad R_i = \sigma^2 R_i \quad R_i \text{ and } D \text{ positive-definite}$$

$$\mathbf{Z}_i = \begin{pmatrix} \mathbf{z}_{i1} \\ \vdots \\ \mathbf{z}_{ij} \\ \vdots \\ \mathbf{z}_{in_i} \end{pmatrix} = \begin{pmatrix} 1 & z_{i11} & \dots & z_{ip1} \\ \dots & \dots & \dots & \dots \\ 1 & z_{i1n_i} & \dots & z_{ipn_i} \end{pmatrix} = (\mathbf{1}' \quad \mathbf{z}'_{i1} \quad \dots \quad \mathbf{z}'_{ip}) \quad Z_i \in \mathbb{R}^{n_i \times (p+1)}$$

$$\mathbf{b}_i = (b_{i0}, b_{i1}, \dots, b_{iq})' \quad \mathbf{W}_i = (\mathbf{1}' \quad \mathbf{w}'_{i1} \quad \dots \quad \mathbf{w}'_{iq}) \quad W_i \in \mathbb{R}^{n_i \times (q+1)}$$

- LMMs in their general form are not unique. To make it identifiable we will specify the structure of the matrix  $R_i$  in terms of a set of parameters for a variance function and a correlation matrix as in LM 3.0.

36



## LMM – group-level specification

- Generalization to multilevel LMMs is possible, though notationally more complex.
- Let  $\mathbf{y} = (\mathbf{y}'_1, \dots, \mathbf{y}'_N)'$  be the vector of all  $n = \sum_{i=1}^N n_i$  observed values of the dependent variable.  
Let  $\mathbf{b} = (\mathbf{b}'_1, \dots, \mathbf{b}'_N)'$  and  $\boldsymbol{\varepsilon} = (\boldsymbol{\varepsilon}'_1, \dots, \boldsymbol{\varepsilon}'_N)'$  be the vectors containing all the  $N \times q$  random effects and  $n$  errors, respectively.

Model equation  
for all data

$$\mathbf{y} = \mathbf{Z}\boldsymbol{\beta} + \mathbf{W}\mathbf{b} + \boldsymbol{\varepsilon}$$

$$\mathbf{b} \sim N_{N(q+1)}(\mathbf{0}, \sigma_b^2 \mathbf{D}) \quad \boldsymbol{\varepsilon} \sim N_n(\mathbf{0}, \sigma_\varepsilon^2 \mathbf{R}) \quad \text{with } \mathbf{b}_i \perp \boldsymbol{\varepsilon}_i \quad \mathbf{b}_i \perp \mathbf{b}_{i'}$$

where

$$\mathbf{D} = \mathbf{I}_N \otimes \mathbf{D} = \begin{pmatrix} \mathbf{D} & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \mathbf{D} \end{pmatrix} \quad \text{and} \quad \mathbf{R} = \begin{pmatrix} \mathbf{R}_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \mathbf{R}_N \end{pmatrix}$$

LMM with *nested structure* for random effects. It is possible to formulate LMMs with non-block-diagonal matrices  $\mathbf{W}$ ,  $\mathbf{D}$ , and  $\mathbf{R}$  (crossed random effects).

$$\mathbf{Z} = \begin{pmatrix} \mathbf{Z}_1 \\ \vdots \\ \mathbf{Z}_i \\ \vdots \\ \mathbf{Z}_N \end{pmatrix} \in \mathbb{R}^{n \times (p+1)} \quad \mathbf{W}_i = \begin{pmatrix} \mathbf{W}_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \mathbf{W}_N \end{pmatrix} \in \mathbb{R}^{n \times N(q+1)}$$

37



## LMM – conditional and unconditional distributions

### UNCONDITIONAL DISTRIBUTION of the RANDOM EFFECTS

- The unconditional distribution of the random effects is a **multivariate normal distribution** with zero mean and variance-covariance matrix  $\mathbf{D}$ .

$$D(\sigma^2, \boldsymbol{\theta}_D) = \sigma^2 D(\boldsymbol{\theta}_D)$$

where  $\boldsymbol{\theta}_D$  is a vector of parameters, which represent the (scaled by  $\sigma^2$ ) variances and covariances of the elements of  $\mathbf{b}$ .

- Note that the matrix  $\mathbf{D}$  is parameterized using a vector of parameters  $\boldsymbol{\theta}_D$ .

- **General case:** any two elements of the vector  $\mathbf{b}$  can be correlated and there are no restrictions imposed on the matrix  $\mathbf{D}$ , except that it is positive-definite and symmetric  
=> general structure for  $\mathbf{D}$  with  $\boldsymbol{\theta}_D$  containing  $q(q+1)/2$  distinct elements corresponding to  $q$  variances and  $q(q-1)/2$  covariances of the random effects included in  $\mathbf{b}$ .  
-> Although  $q$  is typically small, estimating all the parameters may be difficult if the sample size  $n$  is limited.
- **Simplified structure of the matrix  $\mathbf{D}$ :** all elements of the vector  $\mathbf{b}_i$  are independent  
=> diagonal form for  $\mathbf{D}$  with  $\boldsymbol{\theta}_D$  containing  $q$  distinct elements corresponding to  $q$  variances  
Plausibility of the assumption will depend on the data at hand.

38



## LMM – conditional and unconditional distributions

### CONDITIONAL DISTRIBUTION of the OBSERVATIONS given the RANDOM EFFECTS

- For the classical LMMs, the conditional distribution  $f_{\mathbf{y}|\mathbf{b}}(\mathbf{y}_i|\mathbf{b}_i)$  is multivariate normal, with

$$E[\mathbf{y}_i|\mathbf{b}_i] = \boldsymbol{\mu}_i = \mathbf{Z}_i \boldsymbol{\beta} + \mathbf{W}_i \mathbf{b}_i \quad \text{with } \boldsymbol{\mu}_i = (\mu_{i1}, \dots, \mu_{in_i}) \\ V[\mathbf{y}_i|\mathbf{b}_i] = V[\boldsymbol{\varepsilon}_i|\mathbf{b}_i] = \sigma^2 \mathbf{R}_i$$

Thus, conditionally on the (unknown) values of the random effects, the mean value of the dependent-variable is defined by a linear combination of the fixed effects and random effects, respectively.

- In their most general form, LMMs are not identifiable => constraining like in LM 3.0.

$$V[\varepsilon_{ij}|\mathbf{b}_i] = \sigma^2 \mathbf{R}_{ij} = \sigma^2 \lambda^2(\mu_{ij}, \boldsymbol{\delta}; \boldsymbol{\nu}_{ij})$$

39



## LMM – conditional and unconditional distributions

### MARGINAL DISTRIBUTION of the OBSERVATIONS

- The marginal distribution  $f_y(\mathbf{y}_i)$  of the observations is obtained by “integrating out” the random effects.

$$f_y(\mathbf{y}_i) = \int f_{y,b}(\mathbf{y}_i, \mathbf{b}_i) d\mathbf{b}_i = \int f_{y|b}(\mathbf{y}_i | \mathbf{b}_i) f_b(\mathbf{b}_i) d\mathbf{b}_i$$

- Given that  $f_{y,b}(\mathbf{y}_i, \mathbf{b}_i)$  and  $f_b(\mathbf{b}_i)$  are densities of multivariate normal distributions, the marginal distribution of  $\mathbf{y}$  is also multivariate normal and it can be derived analytically.

$$E[\mathbf{y}_i] = Z_i \boldsymbol{\beta}$$

$$V[\mathbf{y}_i] = \sigma^2 [W_i D W_i' + R_i] = \sigma^2 [W_i D(\boldsymbol{\theta}_D) W_i' + R_i(\boldsymbol{\theta}_R, v_i)]$$

- The **marginal mean value** is defined by a linear combination of the vectors of covariates as for the LMs.

- The **marginal variance-covariance matrix** consists of two components.

- The first one is contributed by the random effects.

=> *We are partitioning the global variance!*

- The second one is related to the residuals.

40



## LMM – Estimation

- Numerical methods are always needed in the case of LMMs

- Again, different approaches are possible, leading to possibly different estimates:

- ML --> produce variance estimates biased downwards to some degree
- REML – Restricted (Residual) ML --> produce unbiased estimates marginalizing wrt nuisance parameters

- Model fitting has 3 distinctive components

- Estimates of fixed effects
- Estimates of random effects
- Estimates of variance parameters

- **WARNING:** variance components are nonnegative by definition.

Nevertheless, methods for estimating them may underestimate variance values, and when real values are close to 0 this may result in unrealistic negative estimates. This may happen when:

- the # of random effects ( $q$ ) is small
- the # of obs per random effect is small



## LMM – Inference on parameters

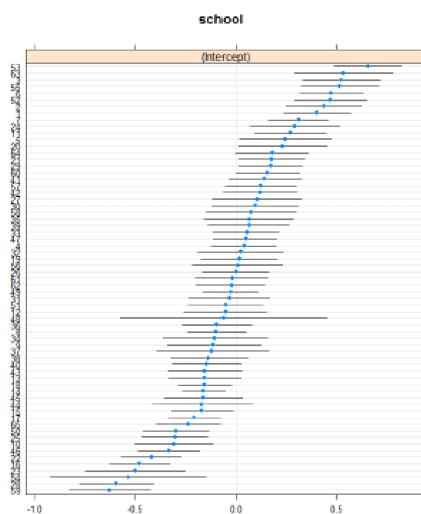
- **Profiled confidence intervals** for the fixed effects and for the standard deviations of random effects and residuals are provided.

- Points and intervals estimates for the random effects can be visualized using a **dotplot**

--> shows the point and interval estimates for the random effects, ordering them and highlighting which are significantly different from 0.

- To get the final estimate for the outcome we add to the contribution of each covariate in the linear predictor the estimate of the random effect the statistical unit belongs to  
--> It may increase or decrease the estimated value for the outcome.

- For each group, we may assess if the corresponding effect is **positive or negative**, or defining a threshold for labelling them as **outlier**



## LMM – a simple guide to their use

- Compute the amount of variability the presence of the grouper accounts for  
=> **PVRE** (Proportion of Variance due to Random Effect)

$$PVRE = \frac{\hat{\sigma}_b^2}{\hat{\sigma}_b^2 + \hat{\sigma}_e^2}$$

Random intercept case

- Use the estimate of  $\sigma_b^2$  for performing **scenario analysis**  
=> what if belonging to a  $-2\hat{\sigma}_b^2$  or a  $+2\hat{\sigma}_b^2$  group, given the same unit conditions?

- We may **cluster the random effects estimates** and use the medoids of the groups as representative elements for **quantifying the effect of the «communities of groupers»**.

- We may use the point estimates of  $N$  random effects as a new response to be modeled through group-specific covariates OR we may introduce group-level covariates within the same model  
=> ( $Z \neq W$ )

43



## LMMs – Example Multicenter Randomized Trial



- The trial was **randomised double blind comparison of 3 treatments for hypertension** and has been reported by Hall et al. (1991). **Diastolic Blood Pressure (DBP) after treatment** was the primary endpoint.
- One treatment was a **new drug (A)** and the **other two (B and C)** were standard drugs for controlling hypertension (A=Carvedilol, B=Nifedipine, C=Atenolol).
- **29 centres** participated in the trial and patients were randomised in order of entry.
- 2 pre treatments and 4 post treatment visits were made as follows:
  - Visit 1 (week 0): Measurements were made to determine whether pts met eligibility criteria for the trial. Pts who did so received a placebo treatment for 1 week, after which they returned for a second visit
  - Visit 2 (week 1): Measurements were repeated and pts who still satisfied the eligibility criteria were entered into the study and randomized to receive one of the 3 treatments
  - Visit 3-6 (weeks 3,5,7,9): Measurements were repeated at 4 post-treatment visits.
- **311 pts** were assessed to entry into the study. Of these, **288** were suitable and randomised to receiving one of the 3 treatments. **30 pts dropped out** of the study prior to Visit 6.
- Measurements of cardiac function, laboratory values, and adverse events were recorded at each visit.



## LMMs – Example Multicenter Randomized Trial



- Focus on a toy-example subsample      **N = 3 hospitals**  
**n = 9 observations**  
**3 covariates => p+1 = 4**  
**Random intercept only => q+1 = 1**
- Fit a Linear Mixed Model for dependent (grouped) observations with homoscedastic residuals

$$\mathbf{y} = (176, 194, 156, 150, 150, 160, 150, 160, 160)'$$

Centre	Treatment	Pre-treatment systolic BP	Post-treatment systolic BP
1	A	178	176
1	A	168	194
1	B	196	156
1	B	170	150
2	A	165	150
2	B	190	160
3	A	175	150
3	A	180	160
3	B	175	160

$$\Rightarrow \mathbf{Z} = \begin{pmatrix} 1 & 178 & 1 & 0 \\ 1 & 168 & 1 & 0 \\ 1 & 196 & 0 & 1 \\ 1 & 170 & 0 & 1 \\ 1 & 165 & 1 & 0 \\ 1 & 190 & 0 & 1 \\ 1 & 175 & 1 & 0 \\ 1 & 180 & 1 & 0 \\ 1 & 175 & 0 & 1 \end{pmatrix},$$



EX1: Random intercept only – homoscedastic residuals

EX: Multicentre Randomized Trial

(1)

N = 3 hospitals

m = 9 observations  
(group 1 = 4, group 2 = 2, group 3 = 3)

p + 1 = 4 covariates

q + 1 = 1 RANDOM INTERCEPT ONLY

> Group Level Formulation

$$\underline{y}_i = \underline{Z}_i \underline{\beta}^T + W_i \underline{b}_{io} + \underline{\varepsilon}_i \rightarrow \underline{\varepsilon}_i \sim N_{m_i}(0, \sigma_e^2 R_i)$$

$$\underline{b}_{io} \sim N_3(0, \sigma_b^2)$$

$$E[\underline{y}_i] = \underline{Z}_i \underline{\beta}^T$$

$$\underline{\beta} = (\beta_0, \beta_1, \beta_2, \beta_3)$$

$$V[\underline{y}_i] = W_i \sigma_b^2 W_i^T + \sigma_e^2 R_i$$

$$i=1 = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix} + \sigma_e^2 \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

$$W_i \in \mathbb{R}^{m_i \times (p+1)} = m_i \times 4$$

$$i=2 = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix} + \sigma_e^2 \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

$$i=3 = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix} + \sigma_e^2 \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

$$W_i \in \mathbb{R}^{m_i \times (q+1)} = m_i \times 1$$

Homoschedasticity  
(but could generalize)

$$= \begin{bmatrix} \sigma_b^2 & \sigma_b^2 & \sigma_b^2 & \sigma_b^2 \\ \sigma_b^2 & \sigma_b^2 & \sigma_b^2 & \sigma_b^2 \\ \sigma_b^2 & \sigma_b^2 & \sigma_b^2 & \sigma_b^2 \\ \sigma_b^2 & \sigma_b^2 & \sigma_b^2 & \sigma_b^2 \end{bmatrix} + \begin{bmatrix} \sigma_e^2 & \sigma_e^2 & 0 \\ 0 & \sigma_e^2 & 0 \\ 0 & 0 & \sigma_e^2 \end{bmatrix}$$

$$= \begin{bmatrix} \sigma_b^2 + \sigma_e^2 & \sigma_b^2 & \sigma_b^2 & \sigma_b^2 \\ \sigma_b^2 & \sigma_b^2 + \sigma_e^2 & \sigma_b^2 & \sigma_b^2 \\ \sigma_b^2 & \sigma_b^2 & \sigma_b^2 + \sigma_e^2 & \sigma_b^2 \\ \sigma_b^2 & \sigma_b^2 & \sigma_b^2 & \sigma_b^2 + \sigma_e^2 \end{bmatrix} \in \mathbb{R}^{m_i \times m_i}$$

> All Data Formulation

$$\underline{y} = \underline{Z} \underline{\beta}^T + W \underline{b}^T + \underline{\varepsilon} \rightarrow \underline{\varepsilon} \sim N_9(0, \sigma_e^2 R) \quad R \in \mathbb{R}^{m \times m} = 9 \times 9$$

$$\underline{b} \sim N_N(0, D) \quad D \in \mathbb{R}^{N \times N} = 3 \times 3$$

$$D = \begin{bmatrix} \sigma_b^2 & & \\ & \sigma_b^2 & 0 \\ 0 & 0 & \sigma_b^2 \end{bmatrix}$$

$$\underline{\beta} = (\beta_0, \beta_1, \beta_2, \beta_3)$$

$$\underline{Z} \in \mathbb{R}^{n \times (p+1)} = 9 \times 4$$

$$W \in \mathbb{R}^{m \times N(q+1)} = 9 \times 3(4)$$

$$\underline{Z} = \begin{bmatrix} 1 & Z_{111} & Z_{211} & Z_{311} \\ 1 & \vdots & \vdots & \vdots \\ 1 & Z_{112} & Z_{212} & Z_{312} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & Z_{119} & Z_{219} & Z_{319} \end{bmatrix}$$

$$W = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{bmatrix}$$

$$\underline{b}_o = [b_{o1}, b_{o2}, b_{o3}]$$

$$\Rightarrow E[\underline{y}] = \underline{Z} \underline{\beta}^T$$

$$V[\underline{y}] = W D W^T + \sigma_e^2 R$$

$$= \begin{bmatrix} \text{---} & \text{---} & \text{---} \\ \boxed{(\star)} & \text{---} & \text{---} \\ \text{---} & \text{---} & \text{---} \end{bmatrix} \quad \leftarrow$$

This structure arise from the assumptions on fixed and random effects carried out in the LMM formulation

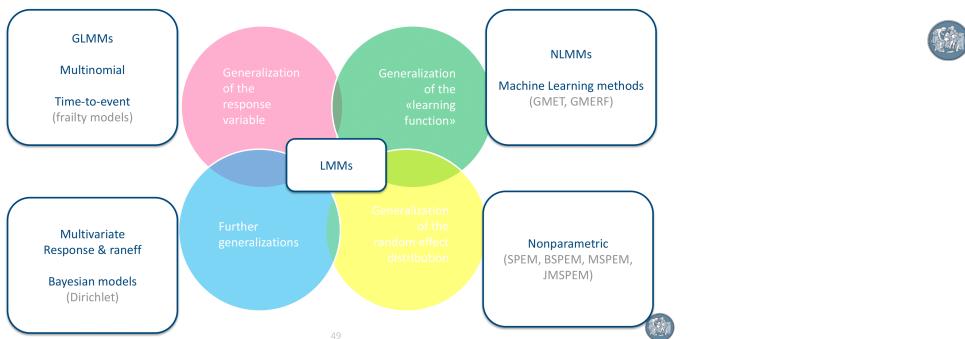
We might design something similar also working on variance and correlation function of LM3.0, but then we loose the opportunity of making inference on groups and disentangle the sources.

- observations between groups are uncorrelated ( $\Rightarrow$  independent in M setting)
- observations within groups correlated with dependence induced by random effects

## Take home messages

- Mind the variance! Modeling variability is what makes LMs still a powerful tool in the statistical learning landscape.
- The model employing random effects implies a marginal normal distribution which is similar to distributions considered in the context of LMs 3.0, but with the variance-covariance matrix of  $\gamma$  of a very specific parametric form.
- LMs 3.0 (fixed effects and correlated residual errors) are less restrictive than LMMs.  
=> LMs 3.0 are more flexible than LMMs, but they **do not allow making inference about the variability that may be related to different levels of the data hierarchy**.
- Neglecting the correlation structure among the observations leads to a big loss in the information carried by the data.
- LMMs allow to extract information at all the levels of the hierarchy and to disentangle the source of variation in the response each level of hierarchy is responsible for (identifying the important ones!).

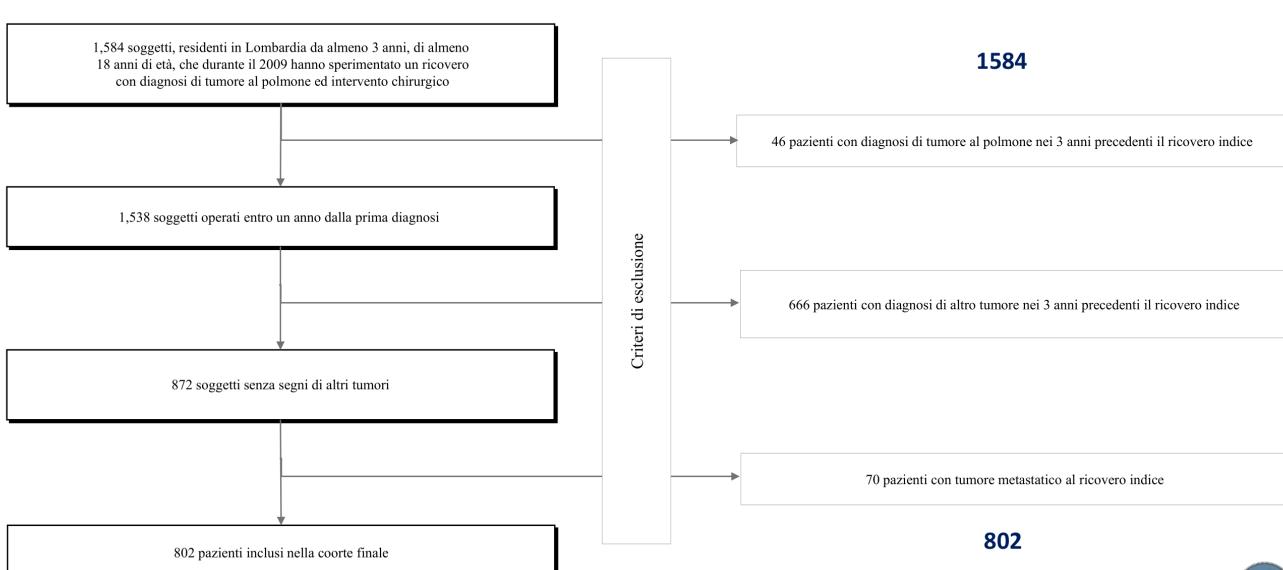
## Beyond LMMs



## Case study I – Scenario analysis on hospital effect in a multi center trial

- Observational multicenter study (N=48 hospitals) on **lung cancer**.
- Investigation of **clinical factors associated to 3y death for any cause** in pts undergoing surgery for lung cancer.
- **Goals:**
  - Risk stratification of the population under study  
(association between mortality and features measured at baseline/entrance)
  - Assessment of the grouper effect
  - Scenario analysis

## Case study I – Cohort selection



## Case study I – Hospital volume



48 hospitals, different volume and outcomes

min = 1 -- max = 109



Volume	Numero di ospedali	Numero di pz
<10	27	87
10 – 30	11	206
31 – 50	7	289
>50	3	220



## Case study I – GLMM with single level of grouping



$$\text{logit}(p_{ij}) = \log\left(\frac{p_{ij}}{1 - p_{ij}}\right) = \beta_0 + \beta_1 z_{1ij} + \dots + \beta_p z_{pij} + b_i$$

➤  $b_i$  random effect  $\sim N(0, \sigma_b^2)$

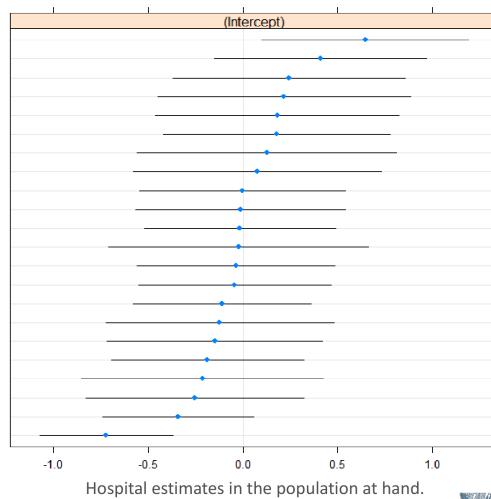
➤  $b_i$  takes the same value for each observation within the same group and different values in different groups.

➤ It can be interpreted as *the effect of being hospitalized in a given structure*.

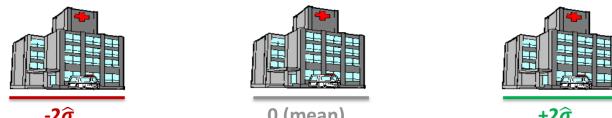
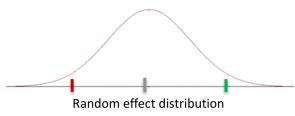
➤ Given the patient-specific features, hospital have a significant effect on the log odds of survival probability

% of total variability accounted for by the grouper

VPC = 4.6%



## Case study I – GLMM with single level of grouping



Man  
65 years  
MCS = 1  
NO adjuvant Chemio  
**Low complexity surgery**



<b>0.73</b>	<b>0.86</b>	<b>0.93</b>
<b>0.45</b>	<b>0.64</b>	<b>0.80</b>

Man  
65 years  
MCS = 1  
NO adjuvant Chemio  
**High complexity surgery**





➤ The 3y overall mortality after surgery for lung cancer depends on individual characteristics AND on the hospital the patients are admitted to (protocols? Physicians ability?)

➤ Each time you have data with grouped structure, using LMMs allows for

- ❖ proper modeling of the **hierarchical structure**
- ❖ better management of **missing data** and unbalancing between groups
- ❖ **variance partitioning**
- ❖ **inference** on grouper(s) population => **scenario analysis**



## Case study II



Develop an **EM algorithm** for **semi-parametric mixed-effects models** for hierarchical data and apply it to **INVALSI data** as an **unsupervised clustering** tool for Italian schools.

Research questions:

- Are there **differences across the effects of Italian schools** on their student achievements?
- Is it possible to identify **groups of schools** that behave in similar/different ways?
- Do these groups of schools differ in terms of their **characteristics**?



## Case study II



Mixed-effects linear model (two-level):

$$\mathbf{y}_i = \boldsymbol{\beta} \mathbf{X}_i + \mathbf{b}_i \mathbf{Z}_i + \boldsymbol{\epsilon}_i \quad \text{for } i = 1, \dots, N$$

where

$\mathbf{y}_i$  is the  $n_i$ -dimensional vector of **response variable** in the  $i$ -th group;  
 $\mathbf{X}_i$  and  $\mathbf{Z}_i$  are the  $(n_i \times (p+1))$  and  $(n_i \times (q+1))$  matrices of **fixed and random covariates** respectively;  
 $\boldsymbol{\beta}$  is the  $(p+1)$ -dimensional vector of **fixed coefficients**;  
 $\mathbf{b}_i$  is the  $(q+1)$ -dimensional vector of **random coefficients** in the  $i$ -th group;  
 $\boldsymbol{\epsilon}_i$  is the  $(n_i)$ -dimensional vector of **errors**,  $\boldsymbol{\epsilon}_i \sim \mathcal{N}(0, \sigma^2_{\epsilon})$ .

- Parametric framework  $\rightarrow \mathbf{b} \sim \mathcal{N}(0, \boldsymbol{\Sigma}_b)$
- Non-parametric framework  $\rightarrow \mathbf{b} \sim \text{discrete distribution } P^*$





$P^*$  can be interpreted as the **mixing distribution** that generates the density of the stochastic model in (★).

The ML estimator  $\hat{P}^*$  of the random effects distribution  $P^*$  can be expressed as a

- **set of points**  $(\mathbf{c}_1, \dots, \mathbf{c}_M)$ , where  $M \leq N$  and  $\mathbf{c}_m \in \mathbb{R}^2$  for  $m = 1, \dots, M$ ,
- and a **set of weights**  $(w_1, \dots, w_M)$ , where  $\sum_{m=1}^M w_m = 1$  and  $w_m \geq 0$  for each  $m = 1, \dots, M$ , that represents the proportion of groups belonging to each cluster  $m$ .

61



Besides the **estimate of  $\mathbf{M}$** , the joint estimation of  $\sigma^2, \beta, (\mathbf{c}_1, \dots, \mathbf{c}_M)$  and  $(w_1, \dots, w_M)$  in model (★) is performed through the **maximization of the likelihood**, mixture by the discrete distribution of the random effects:

$$L(\beta, \sigma^2 | \mathbf{y}) = p(\mathbf{y} | \beta, \sigma^2) = \sum_{m=1}^M \frac{w_m}{(2\pi\sigma^2)^{\frac{\sum_{i=1}^N n_i}{2}}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^N \sum_{j=1}^{n_i} (y_{ij} - c_{0m} - \beta x_{ij} - c_{1m} z_{ij})^2 \right\},$$

with respect to the fixed coefficient  $\beta$ , the error variance  $\sigma^2$  and the random effects distribution  $(\mathbf{c}_m, w_m)$ , for  $m = 1, \dots, M$ .

→ **EM algorithm** ([SPEM algorithm](#)).

62



The algorithm is inspired by the one proposed by Azzimonti *et. al* (2013), but with the following improvements:

- the **initialization of the support points** has been reviewed;
- the **maximization** step is made in **closed form**;
- the **covariates are group specific** (they can be different in terms of number of observations and assumed values).

The **support reduction** of the discrete distribution stands on 2 criteria:

- two points  $\mathbf{c}_m$  and  $\mathbf{c}_k$  closer than a fixed threshold  $D$  **collapse to a unique point**  $\mathbf{c}_{m,k} = \frac{\mathbf{c}_m + \mathbf{c}_k}{2}$  with weight  $w_{mk} = w_m + w_k$ ;
- we remove points with a weight **w lower than a fixed threshold  $\tilde{w}$** , parameterising the remaining weights.

63





The update steps of the SPEM algorithm are the following:

- $w_m^{(up)} = p(\mathbf{b} = \mathbf{c}_m) = \frac{1}{N} \sum_{i=1}^N W_{im}$  for  $m = 1, \dots, M$
- $(\beta^{(up)}, \sigma^2)^{(up)} = \arg \max_{\beta, \sigma^2} \sum_{m=1}^M \sum_{i=1}^N W_{im} \ln p(\mathbf{y}_i | \beta, \sigma^2, \mathbf{c}_m)$
- $\mathbf{c}_m^{(up)} = \arg \max_{\mathbf{c}} \sum_{i=1}^N W_{im} \ln p(\mathbf{y}_i | \beta, \sigma^2, \mathbf{c})$   $m = 1, \dots, M$

where

- $W_{im} = \frac{w_m p(\mathbf{y}_i | \beta, \sigma^2, \mathbf{c}_m)}{\sum_{k=1}^M w_k p(\mathbf{y}_i | \beta, \sigma^2, \mathbf{c}_k)} = p(\mathbf{b}_i = \mathbf{c}_m | \mathbf{y}_i, \beta, \sigma^2)$
- $p(\mathbf{y}_i | \beta, \sigma^2, \mathbf{c}_m) = \frac{1}{(2\pi\sigma^2)^{\frac{n_i}{2}}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{j=1}^{n_i} (y_{ij} - \beta x_{ij} - c_{0m} - c_{1m} z_{ij})^2 \right\}$

SPEM algorithm → R package

64



## Case study II – INVALSI data



### ■ Students

- ▶ Maths **INVALSI test scores at grade 8 (MATH8)**,
- ▶ Maths **INVALSI test scores at grade 6 (MATH6)**,
- ▶ socio-economical index (**ESCS**).



### ■ Schools

- ▶ Information about **school-body composition** (e.g. percentages of females/immigrants, average ESCS, private/public),
- ▶ information about **school principal** (gender, age, education),
- ▶ **managerial practices** of the school.

■ **Size:** 6,572 students within 462 schools.

65



## Case study II – INVALSI data



Semi-parametric two-level model for **students** (level 1) nested within **schools** (level 2):

$$\begin{aligned} \mathbf{y}_i &= c_{0m} + c_{1m}\mathbf{z}_i + \beta\mathbf{x}_i + \epsilon_i \quad i = 1, \dots, N \quad m = 1, \dots, M \\ \epsilon_i &\sim \mathcal{N}(\mathbf{0}, \sigma_\epsilon^2) \quad \text{ind.} \end{aligned} \quad (\heartsuit)$$

where

$N$  is the total number of schools,

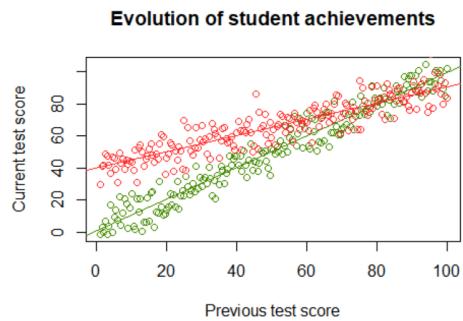
$\mathbf{y}_i$  is the **maths score at grade 8** of students attending the  $i$ -th school,

$\mathbf{z}_i$  is the **maths score at grade 6** of students attending the  $i$ -th school,

$\mathbf{x}_i$  is the **socio-economic index “ESCS”** of students attending in the  $i$ -th school.

66





Example of distributions of student achievements within two schools. **Previous score** represents the student score before entering the school, **current score** represents the student score at the end of the period within the school.

**How does the relation change across the two schools?**

67

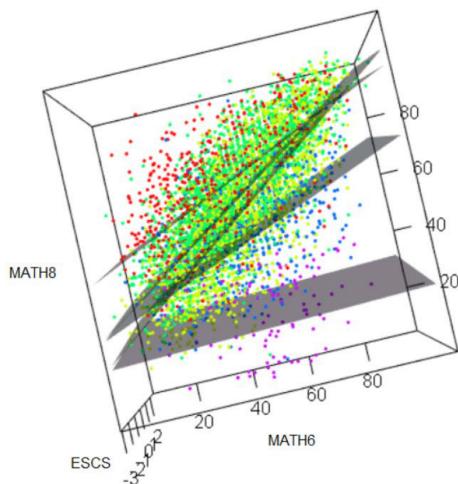


The SPEM algorithm, considering  $D = 0.5$  and  $\tilde{w} = 0.01$ , identifies  $M=5$  distinct clusters of schools:

Cluster	$\beta$	$c_0$	$c_1$	$w$
Cluster 1	1.417	46.028	0.454	12.2%
Cluster 2	1.417	22.579	0.707	39.6%
Cluster 3	1.417	30.293	0.648	37.5%
Cluster 4	1.417	31.207	0.393	8.8%
Cluster 5	1.417	25.359	0.027	1.9%

Table: ML estimates of fixed and random coefficients of model (♥) obtained by the SPEM algorithm.

68



Plot of data with the 5 identified regression planes. Colors represent the 5 clusters.

69





To explore a posteriori the clusters we apply the following **multinomial logit model**, for each school  $i$ ,  $i = 1, \dots, N$  and each cluster  $m = \{1, 4, 5\}$ :

$$\ln\left(\frac{P(Y_i = m)}{P(Y_i = C_{ref})}\right) = \beta_{m0} + \sum_{q=1}^Q \beta_{mq} V_{iq}.$$

where

- $Y_i$  represents the **cluster of belonging** of school  $i$ ;
- $C_{ref}$  is the **reference cluster** (union of clusters 2 and 3);
- $V$  is the  $N \times q$  matrix of **school level variables**;
- $\beta_m$  is the vector of coefficients for cluster  $m$ .

70



School level variables	cluster 1	cluster 4	cluster 5
Private School	0.884	-9.187***	-6.147***
Scientific education (yes=1) of the school principal	-0.135	0.171	-6.019***
Central Italy	0.744	0.648	15.691***
Southern Italy	1.201.	1.200.	14.687***

Table: Results of the multinomial logit model. Asterisks denote the significance of the coefficients.

71



- The SPEM algorithm for hierarchical data can be used as a tool to perform **in-built clustering** in many classification problems
- Contrarily to existing methods, **it does not need to fix a priori the number of mass points** of the random effects discrete distribution
- It represents a **novelty and a value-added** both in the **non-parametric framework** of mixed-effects models and in the **research about school effectiveness**
- When applied to INVALSI data, it identifies **differences across the impacts of Italian schools** on their students achievements

72



# R corner

Useful references to R packages and R functions to implement mixed-effects models

## R packages and available codes

### R packages employed in the R laboratory session:

1. **nlme**: fixed and mixed-effects regression models with homoscedastic/heteroscedastic and independent/correlated residuals
2. **lme4**: mixed-effects regression models only with homoscedastic and independent residuals (but with a bunch of accessories)
3. **insight**: extract information from a mixed-effects model (e.g., variance decomposition)

R codes relative to "PoliMi literature" about mixed-effects models:

1. **GMET()** function: supplementary material of 'Fontana, L., Masci, C., Ieva, F. and Paganoni, A.M. (2021). Performing learning analytics via generalized mixed-effects trees'. *Data*, 6, 74.'
2. **GMERF()** function: supplementary material of 'Pellagatti M., Masci, C., Ieva, F. and Paganoni, A.M. (2021). Generalized mixed-effects random forest: A flexible approach to predict university student dropout'. *Statistical Analysis and Data Mining: The ASA Data Science Journal*, 14(3), 241-257.'
3. **SPEM()** function: supplementary material of 'Masci, C., Ieva, F. and Paganoni, A.M. (2018). Semi-parametric mixed-effects models for unsupervised classification of Italian schools'. *Journal of the Royal Statistical Society: Series A (Statistics in Society)* 182.4, pp.'
4. **MSPEM()** function: supplementary material of 'Masci, C., Ieva, F. and Paganoni, A.M. (2022). Semiparametric multinomial mixed-effects models: a university students profiling tool.' *The Annals of Applied Statistics* - in press.'

74



## References

### References

#### Journal papers

- Masci, C., Ieva, F. and Paganoni, A.M. (2022). Semiparametric multinomial mixed-effects models: a university students profiling tool.' *The Annals of Applied Statistics* - in press
- Cannistra' M., Masci, C., Ieva, F., Agasisti T. and Paganoni, A.M. (2021). 'Early-predicting dropout of university students: an application of innovative machine learning and multilevel statistical techniques'. *Studies in Higher Education*, 1-22, DOI:10.1080/03075079.2021.2018415.
- Fontana, L., Masci, C., Ieva, F. and Paganoni, A.M. (2021). Performing learning analytics via generalized mixed-effects trees'. *Data*, 6, 74.
- Masci, C., Ieva, F., Agasisti, T. and Paganoni, A.M. (2021). Evaluating class and school effects on the joint achievements in different subjects: a bivariate semiparametric mixed-effects model'. *Computational Statistics*, 36, pages 2337–2377.
- Pellagatti M., Masci, C., Ieva, F. and Paganoni, A.M. (2021). Generalized mixed-effects random forest: A flexible approach to predict university student dropout'. *Statistical Analysis and Data Mining: The ASA Data Science Journal*, 14(3), 241-257.
- Gasperoni, F.; Ieva, F.; Paganoni, A.M.; Jackson, C.; Sharples, L. (2020) Evaluating the effect of healthcare providers on the clinical path of Heart Failure patients through a novel semi-Markov multi-state model. *BMC Health Services Research*, 20 (1): 1-18.
- Gasperoni, F., Ieva, F., Paganoni, A.M., Jackson C., Sharples L.D. (2019) Nonparametric frailty Cox models for hierarchical time-to-event data *Biostatistics*, 21 (3): 531-544.
- Masci, C., Ieva, F. and Paganoni, A.M. (2018). Semi-parametric mixed-effects models for unsupervised classification of Italian schools'. *Journal of the Royal Statistical Society: Series A (Statistics in Society)* 182.4, pp.
- Masci, C., Agasisti, T. and Johnes, G. (2018). Student and school performance across countries: A machine learning approach'. *European Journal of Operational Research*, 269(3), pp. 1072-1085. 313-3142.
- Masci, C., Ieva, F., Agasisti, T. and Paganoni, A.M. (2017). Bivariate multilevel models for the analysis of mathematics and reading pupils' achievements'. *Journal of Applied Statistics* 44:7, pp. 1296-1317.
- Masci, C., Ieva, F., Agasisti, T. and Paganoni, A. M. (2016). Does class matter more than school? Evidence from a multilevel statistical analysis on Italian junior secondary school students'. *Socio-Economic Planning Sciences* (54), pp. 47-57.
- Guglielmi, A., Ieva, F., Paganoni, A.M., Ruggeri, F., Soriano, J. (2014). Semiparametric Bayesian modeling for the classification of patients with high observed survival probabilities. *Journal of the Royal Statistical Society - Series C*, 63 (1): 25-46
- Grieco, N., Ieva, F., Paganoni, A.M. (2012). Performance assessment using mixed effects models: a case study on coronary patient care. *IMA Journal of Management Mathematics*, 23(2), 117-131
- Guglielmi, A., Ieva, F., Paganoni, A.M., Ruggeri, F. (2012). A Bayesian random-effects model for survival probabilities after Acute Myocardial Infarction. *Chilean Journal of Statistics*, 3(1): 1-15.
- Ieva, F., Paganoni, A.M. (2010). Multilevel models for clinical registers concerning STEMI patients in a complex urban reality: a statistical analysis of MOMI2 survey. *Communications in Applied and Industrial Mathematics*, 1(1), 128-147

76



## References

#### Books

- Galecki, A., & Burzykowski, T. (2013). Linear mixed-effects model using R. Springer, New York, NY.
- Pinheiro, J., & Bates, D. (2006). Mixed-effects models in S and S-PLUS. Springer science & business media.
- Gelman, A., Hill, J. (2006) Data Analysis Using Regression and Multilevel/Hierarchical Models. Cambridge University Press.
- Goldstein, H. (2003) Multilevel Statistical Models. Third edition. London

#### Codes

- Pinheiro J, Bates D, R Core Team (2022). *nlme: Linear and Nonlinear Mixed Effects Models*. R package version 3.1-157, <https://CRAN.R-project.org/package=nlme>.
- Bates D, Mächler M, Bolker B, Walker S (2015). "Fitting Linear Mixed-Effects Models Using lme4." *Journal of Statistical Software*, 67(1), 1-48. doi: [10.18637/jss.v067.i01](https://doi.org/10.18637/jss.v067.i01)
- Lüdecke D, Waggoner P, Makowski D (2019). "insight: A Unified Interface to Access Information from Model Objects in R." *Journal of Open Source Software*, 4(38), 1412. doi:[10.21105/joss.01412](https://doi.org/10.21105/joss.01412).

77



## LECTURE 26 20/5/2022

### Spatial Statistics

#### GEOSTATISTICS

- $s_1, \dots, s_m$  spatial locations  $\in D$  (fixed)
- $z_{s_1}, \dots, z_{s_m}$  observations collected at those locations  
↳ correlated
- Aims
  - Model and estimate the spatial dependence
  - Predictions: what happens in  $s_0$ ?
  - Estimate the mean of  $z_s$



► LATTICE DATA Goal: - clustering - modelling the mean

► SPATIAL POINT PROCESSES  $s_1, \dots, s_m$  random locations Goal: study the dist. of  $s_1, \dots, s_m$

We consider  $z_{s_1}, \dots, z_{s_m}$  scalar

Model  $\{z_s, SED\}$

$$z_s = m_s + \delta_s$$

drift                      individual

$$\mathbb{E}[z_s] = m_s \quad \mathbb{E}[\delta_s] = 0$$

### STATIONARITY

► STRONG STATIONARITY  $(z_{s_1}, \dots, z_{s_m}) \sim (z_{s_1+h}, \dots, z_{s_m+h}) \quad \forall s_1, \dots, s_m \in D \quad \forall h \in \mathbb{R}^d \quad \forall m \geq 1$

► 2nd ORDER STATIONARITY

- $\mathbb{E}[z_s] = m \quad \forall s \in D$
- $\text{Cov}(z_{s_1}, z_{s_2}) = C(s_1 - s_2)$



Properties:

- Positive semidefinite  $\sum_i \lambda_i \lambda_j C(s_i - s_j) \geq 0 \quad \forall \lambda_i, \lambda_j \in \mathbb{R} \quad \forall s_1, \dots, s_m \in D$
- Symmetric  $C(-h) = C(h)$
- Bounded  $|C(h)| \leq \underbrace{C(0)}_{\text{variance}}$

VARIOGRAM  $2\gamma(s_1 - s_2) = \text{Var}(z_{s_1} - z_{s_2}) = \mathbb{E}[(z_{s_1} - z_{s_2})^2] - (\mathbb{E}[z_{s_1}] - \mathbb{E}[z_{s_2}])^2$

Relation between  $2\gamma(\cdot)$  and  $C(\cdot)$

$$\begin{aligned} 2\gamma(h) &= \text{Var}(z_{s_1} + z_{s_2} - 2\text{cov}(z_{s_1}, z_{s_2})) \\ &= 2C(0) - 2C(h) \end{aligned}$$

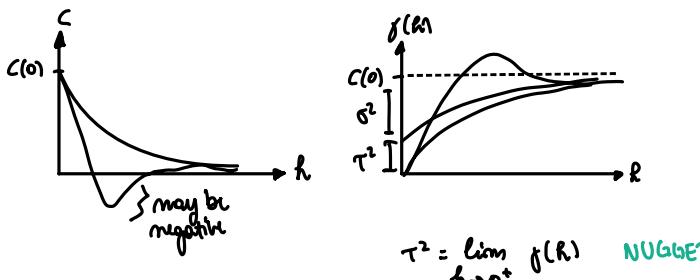
$\frac{\parallel s_1 - s_2 \parallel}{h}$

### Properties

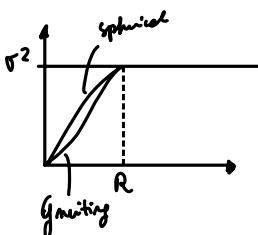
- conditional negative def.  $\sum_i \lambda_i \lambda_j \gamma(s_i - s_j) \leq 0 \quad \forall s_i, s_j \in D \quad \forall \lambda_i \in \mathbb{R} \quad \sum_i \lambda_i = 0$
- symm  $\gamma(-h) = \gamma(h) \quad \forall h \in \mathbb{R}^d$
- Non negative  $\gamma(h) \geq 0 \quad \forall h \in \mathbb{R}^d$
- Zero at origin  $\gamma(0) = 0$
- Sub quadratic growth  $\lim_{\parallel h \parallel \rightarrow \infty} \frac{2\gamma(h)}{\parallel h \parallel^2} = 0$

Isotropy: second order stationarity + directional homogeneity

- $\mathbb{E}[z_s] = m$
- $\text{Cov}(z_{s_1}, z_{s_2}) = C(\parallel h \parallel) \Rightarrow 2\gamma(\parallel h \parallel)$



$$R: f(h) = \sigma^2 + \tau^2 \quad h \geq R$$



Valid model:  $f$  s.t. all variogram prop. are satisfied

## LECTURE 27 23/5/2022

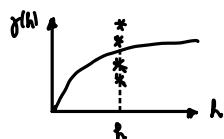
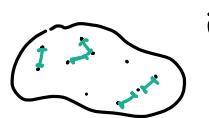
$$C(R) = \underbrace{C(0)}_{\text{firm } f(h)} - f(R) \quad f \xrightarrow{\sim} C$$

Estimates for  $f$  are more robust.

### Estimation of the variogram $\hat{f}$

$$\begin{aligned} ? \quad 2\hat{f}(R) &= \text{Var}(z_{s_1} - z_{s_2}) \quad \|s_1 - s_2\| = R \\ &= \mathbb{E}[(z_{s_1} - z_{s_2})^2] - (\mathbb{E}[z_{s_1}] - \mathbb{E}[z_{s_2}])^2 \end{aligned} \quad \xrightarrow{\text{under stationarity}}$$

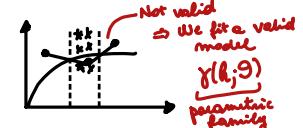
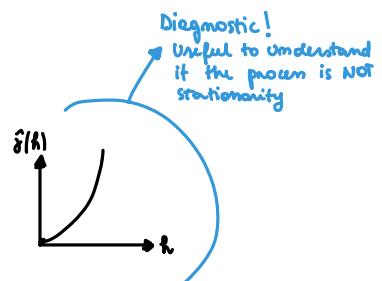
If we're not under stationarity,  $2\hat{f} = \mathbb{E}[(z_{s_1} - z_{s_2})^2]$  will estimate  $2\hat{f}(R) + (m_{s_1} - m_{s_2})^2$



$$\begin{aligned} 2\hat{f}(R) &= \frac{1}{|N(R)|} \sum_{(i,j) \in N(R)} (z_{s_i} - z_{s_j})^2 \\ N(R) &= \{(i,j) : \|s_i - s_j\| \leq R\} \end{aligned}$$

Binned Empirical estimator

Exact  $R$  is too strict



### Spatial predictions (Kriging)

$s_1, \dots, s_m \in D$  locations

$z_{s_1}, \dots, z_{s_m}$  observations

$z_{s_0}$ ?

In general find  $f(z_{s_1}, \dots, z_{s_m})$  measurable that minimizes  $\mathbb{E}[(z_{s_0} - f(z))^2]$

⇒ Solution  $\mathbb{E}[z_{s_0} | z]$

properties: 1) Unbiased  $\mathbb{E}[z_{s_0} | z] = \mathbb{E}[z_{s_0}]$  residuals

2) Orthogonal  $\mathbb{E}[z_{s_0} | z] \perp z_{s_0} - \mathbb{E}[z_{s_0} | z]$

3) Data interpolation If  $s_0 \in \{s_1, \dots, s_m\}$ ,  $s_0 = s_i \Rightarrow \mathbb{E}[z_{s_0} | z] = z_{s_i}$

USA school:  
assume it and go on [ 4) If  $\{z_s, s \in D\}$  is Gaussian  $\mathbb{E}[z_{s_0} | z] = \lambda_0^* + \sum_{i=1}^m \lambda_i^* z_{s_i}$

→ We won't look at the BEST function if we don't have Gaussianity, we'll restrict ourselves to a smaller family

Best Linear Unbiased Predictor

→ Kriging:  $z_{s_0} = \lambda_0 + \sum_i \lambda_i z_{s_i} \Rightarrow$  we look for  $\lambda_0, \dots, \lambda_m$  the min  $\mathbb{E}[(z_{s_0} - z_{s_0}^*)^2]$  s.t.  $\mathbb{E}[z_{s_0}^*] = \mathbb{E}[z_{s_0}]$

BLUP



3 types:

- 1) Simple Kr: Rp  $E[z_s]$  known everywhere (unless in practice)
- 2) Ordinary Kr: Rp 2nd order stationarity but unknown mean
- 3) Universal Kr: relax 2nd order stationarity

$c(R)$  known

## ► Ordinary Kr

Assumptions:

2nd order stationarity & isotropy

$$m = E[z_s]$$

$$C(R) = \text{Cov}(z_{s1}, z_{s2}) \quad \|s_1 - s_2\| = R$$

Problem:

find  $\lambda_0, \dots, \lambda_m \in \mathbb{R}$  st.  $\min E[(z_{s0} - (\lambda_0 + \sum_i \lambda_i z_{si}))^2]$

$$\text{st. } E[(\lambda_0 + \sum_i \lambda_i z_{si})] = m \quad \sum_i \lambda_i = 1$$

$$\text{umb. } \Rightarrow \lambda_0 + \sum_i \lambda_i E[z_{si}] = m$$

$$\lambda_0 + \sum_i \lambda_i m = m$$

$$\mid \text{impose uniform umbenndnung: } \begin{cases} \lambda_0 = 0 \\ \sum_i \lambda_i = 1 \end{cases}$$

Not a convex combination because we don't have  $0 \leq \lambda_i \leq 1$

Objective functional  $\Phi(\lambda, \zeta) = E[(z_{s0} - \sum_i \lambda_i z_{si})^2] + 2 \zeta (\sum_i \lambda_i - 1)$

$$= \text{Var}(z_{s0} - \sum_i \lambda_i z_{si}) + 2 \zeta (\sum_i \lambda_i - 1)$$

$$= \underbrace{\text{Var}(z_{s0})}_{c(0)} + \sum_i \sum_j \lambda_i \lambda_j \text{Cov}(z_{si}, z_{sj}) - 2 \sum_i \lambda_i \text{Cov}(z_{s0}, z_{si}) + 2 \zeta (\sum_i \lambda_i - 1)$$

$$= c(0) + \sum_i \sum_j \lambda_i \lambda_j C(\|s_i - s_j\|) - 2 \sum_i \lambda_i C(\|s_0 - s_i\|) + 2 \zeta (\sum_i \lambda_i - 1)$$

$$\begin{cases} \frac{\partial \Phi}{\partial \lambda_i} = 0 = \cancel{2} \sum_j \lambda_j C(\|s_i - s_j\|) - \cancel{2} C(\|s_0 - s_i\|) + \cancel{2} \zeta \quad i = 1, \dots, m \\ \frac{\partial \Phi}{\partial \zeta} = 0 = 2 (\sum_i \lambda_i - 1) \end{cases}$$

$$(\Sigma)_{ij} = C(\|s_i - s_j\|) \text{ cov. matrix}$$

$$\rightarrow \begin{cases} \sum_j \lambda_j C(\|s_i - s_j\|) + \zeta = C(\|s_0 - s_i\|) \\ \sum_i \lambda_i = 1 \end{cases} \quad \text{OK system}$$

$$\Leftrightarrow \left( \begin{array}{cc} \sum_i & 1 \\ 1 & 0 \end{array} \right) \begin{pmatrix} \lambda \\ \zeta \end{pmatrix} = \begin{pmatrix} C(s_0) \\ 1 \end{pmatrix} \quad \begin{matrix} (s_0)_i = C(\|s_0 - s_i\|) \\ \text{vector of covariances} \end{matrix}$$

## ► Universal Kr

Assumptions

$$z_s = m_s + \delta_s \quad \text{SDE}$$

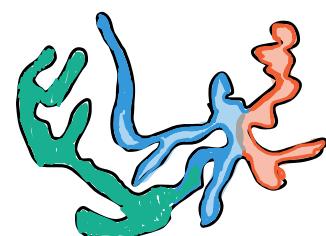
$$E[\delta_s] = 0$$

$$\text{Cov}(\delta_{s1}, \delta_{s2}) = C(R) \quad R = \|s_1 - s_2\|$$

$$m_s = E[z_s] \quad \text{drift}$$

$$= \sum_{l=0}^L a_l f_l(s) \quad \begin{matrix} \text{unknown} \\ \text{regressors known} \\ \text{coeffs.} \end{matrix}$$

$$(m_s = \underbrace{a_0}_{\text{ex}} + a_1 x + a_2 y) \quad (x, y) = s$$



Problem

Find  $\lambda_0, \dots, \lambda_m \in \mathbb{R}$  that solve  $\min E[(z_{s0} - \lambda_0 - \sum_i \lambda_i z_{si})^2]$  st.  $E[(\lambda_0 + \sum_i \lambda_i z_{si})] = m_{s0} = \sum_{l=0}^L a_l f_l(s_0)$

$$\lambda_0 - \sum_i \lambda_i E[z_{si}] = \sum_{l=0}^L a_l f_l(s_0)$$

$$\lambda_0 - \sum_i \lambda_i \sum_l a_l f_l(s_i) = \sum_l a_l f_l(s_0)$$

$$\lambda_0 - \sum_l a_l \sum_i \lambda_i f_l(s_i) = \sum_l a_l f_l(s_0) \rightarrow \begin{cases} \lambda_0 = 0 \\ \sum_i \lambda_i f_l(s_i) = f_l(s_0) \quad l = 0, \dots, L \end{cases}$$

$$\Rightarrow \dots \Rightarrow \begin{pmatrix} \sum_l & F \\ F' & 0 \end{pmatrix} \begin{pmatrix} \lambda \\ \zeta \end{pmatrix} = \begin{pmatrix} 0 \\ f_0 \end{pmatrix}$$

## LECTURE 28 24/5/2022

To estimate  $\Sigma$  we can estimate  $2\hat{f}(\cdot)$  from the residuals  $\delta_{s_1}, \dots, \delta_{s_m}$

1) Empirical estimate

$$\hat{f}(R) = \frac{1}{2|N(R)|} \sum_{(i,j) \in N(R)} (\delta_{s_i} - \delta_{s_j})^2$$

2) Fit a valid model  $f(h; \underline{\theta})$ , get  $\hat{\delta}$

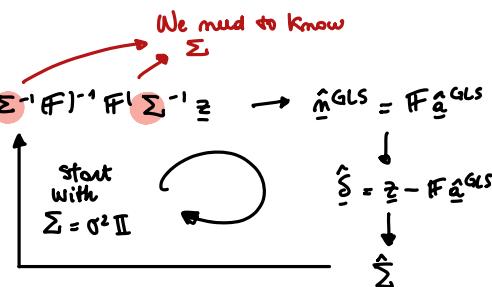
However the  $\delta$ 's are unknown! We first have to estimate them

$$z_{s_i} = \sum_{j \neq i} a_j f_j(s_i) + \delta_{s_i}, \quad s_i \in D$$

$$\underline{z} = F\underline{a} + \underline{\delta} \quad \text{Cov}(\underline{\delta}) = \Sigma$$

$$\Rightarrow \underset{\underline{a} \in R^{L+M}}{\text{argmin}} (\underline{z} - F\underline{a})^T \Sigma^{-1} (\underline{z} - F\underline{a}) = \hat{\underline{a}}^{\text{GLS}} = \dots = (F^T \Sigma^{-1} F)^{-1} F^T \Sigma^{-1} \underline{z} \rightarrow \hat{\underline{a}}^{\text{GLS}} = F^T \Sigma^{-1} \underline{z}$$

BLUE for  $\underline{a}$



## LECTURE 29 26/5/2022

- Informally, **functional data** are entities that can be described through a function, e.g., a curve, a surface, an image
- A **functional dataset** consists of a sample of functional observations
- Even though observations are actually discrete and affected by noise, the observed values reflect a **smooth variation of the phenomenon**. One might be interested not only in **point-wise** values, but also in **differential properties** of the data

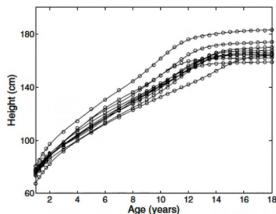


Figure 1.1. The heights of 10 girls measured at 31 ages. The circles indicate the unequally spaced ages of measurement.

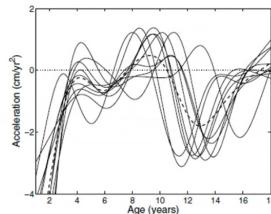


Figure 1.2. The estimated accelerations of height for 10 girls, measured in centimeters per year. The heavy dashed line is the cross-sectional mean, and is a rather poor summary of the curves.

Ramsay Silverman 2005 Springer

### Books:

- Ramsay, J.O. and Silverman, B.W. (2005). *Functional Data Analysis*, Springer, 2nd ed.
- Ramsay, J.O. and Silverman, B.W. (2002). *Applied Functional Data Analysis*, Springer.
- Ramsay, J.O., Hooker, G. and Graves, S. (2009). *Functional Data Analysis with R and Matlab*, Springer.
- Ferraty, F. and Vieu, P. (2006). *Nonparametric Functional Data Analysis: Theory and Practice*, Springer.
- Horvath, L. and Kokoszka P. (2012). *Inference for Functional Data with Applications*, Springer.
- Kokoszka P. and Reimherr, M. (2017). *Introduction to Functional Data Analysis*. Chapman & Hall

### Introductory paper:

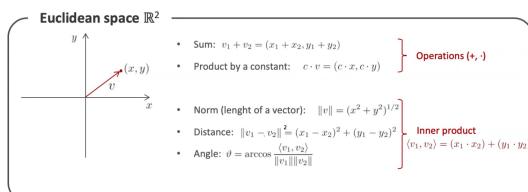
- Sørensen, H., Goldsmith, J., Sangalli, L.M. (2013), "An introduction with medical applications to functional data analysis". *Statistics in Medicine*, 32, pp. 5222–5240.

### Software:

- R package fda (corresponding Matlab code available from <http://www.psych.mcgill.ca/misc/fda/>)
- R package Refund
- Matlab code PACE
- R package mgcv
- R package fdakma (alignment and clustering)
- R package fdaPDE (functional data over complex multidimensional domains)

The notion of **Hilbert space** generalizes the concept of Euclidean space to spaces of any (even infinite) dimension

- Vectorial structure (linear combinations)
- Distance, angles, projections (measure of dependence, best approximations)



### $L^2$ : space of real-valued square-integrable functions

- Sum:  $(f_1 + f_2)(t) = f_1(t) + f_2(t)$
- Product by a constant:  $(c \cdot f)(t) = c \cdot f(t)$

### Operations ( $\star, \cdot$ )

$$\int_0^1 f(t)^2 dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

$$\int_0^1 (f_1(t) \cdot f_2(t)) dt$$

$$\int_0^1 (f_1(t) - f_2(t))^2 dt$$

Let  $H$  be a linear space. An inner product on  $H$  is a bilinear, symmetric, positive definite form

$$\langle \cdot, \cdot \rangle : H \times H \rightarrow \mathbb{R}$$

that satisfies

- (i)  $\langle \lambda x + y, z \rangle = \lambda \langle x, z \rangle + \langle y, z \rangle \quad \forall \lambda \in \mathbb{R}, \forall x, y, z \in H$
- (ii)  $\langle x, y \rangle = \langle y, x \rangle \quad \forall x, y \in H$
- (iii)  $\langle x, x \rangle \geq 0 \quad \forall x \in H$
- (iv)  $\langle x, x \rangle = 0 \iff x = 0$

In particular:

- The inner product allows to measure lengths and angles
- It allows to define orthogonality: two vectors in  $H$  are orthogonal if  $\langle x, y \rangle = 0$
- The inner product induces a norm and a metric
- The inner product allows generalizing the Pythagoras' Theorem:

$$\|x + y\|^2 = \|x\|^2 + \|y\|^2 \text{ if and only if } \langle x, y \rangle = 0$$

A (real) Hilbert space  $H$  is an inner product space that is complete, in the norm induced by the inner product.

- A Hilbert space is complete in the sense that it contains all the limit points of its Cauchy sequences;
- A Hilbert space is separable if it contains a dense countable subset
- Useful properties:

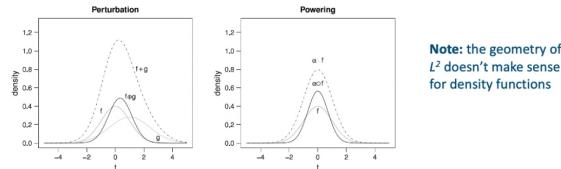
- In a Hilbert space one has the notion of orthogonal projection and of best approximations
- A Hilbert space  $H$  is separable iff it has an orthonormal basis  $\{u_n\}_{n \in \mathbb{N}}$
- If  $H$  is separable Hilbert space,  $\{u_n\}_{n \in \mathbb{N}}$  is an orthonormal basis and  $x \in H$  then

$$x = \sum_{n=1}^{\infty} \langle x, u_n \rangle u_n. \quad \text{Basis expansion}$$

## B<sup>2</sup>: space of density functions on a closed interval $I$ , with $\log$ in $L^2$

- Equivalence relation:  $f, g$  are equivalent if they are proportional (*scale invariance*)

- Sum (perturbation):  $(f \oplus g)(t) = \frac{f(t)g(t)}{\int_I f(s)g(s) ds}$
- Product by a constant (powering):  $(\alpha \odot f)(t) = \frac{f(t)^{\alpha}}{\int_I f(s)^{\alpha} ds}, \quad t \in I$ .
- Inner product:  $\langle f, g \rangle_B = \frac{1}{2\eta} \int_I \int_I \ln \frac{f(t)}{f(s)} \ln \frac{g(t)}{g(s)} dt ds$
- Norm:  $\|f\|_B = \left[ \frac{1}{2\eta} \int_I \int_I \ln^2 \frac{f(t)}{f(s)} dt ds \right]^{1/2}$



- Let  $H$  be a Hilbert space, whose points are functions defined on a closed interval  $T = [t_{min}, t_{max}]$  (e.g., range of time during which the data are collected)
- Hereafter, we will always consider functional data in Hilbert spaces

### Definition 1:

A **functional random variable** is a random element defined on a probability space  $(\Omega, \mathcal{F}, \mathbb{P})$  with values in  $H$

$$X : \Omega \rightarrow H$$

### Definition 2:

A **functional datum**  $x$  is a realization of a functional random variable, i.e., for  $\omega \in \Omega$ ,

$$x = X(\omega) : T = [t_{min}, t_{max}] \rightarrow \mathbb{R}$$

### Definition 3:

A **functional dataset** is a collection of functional data.

Let  $X : \Omega \rightarrow H$  be a functional random variable in  $H$ .

We assume  $\mathbb{E}[\|X\|_H^4] < \infty$

### Definition 4:

We call Fréchet mean of  $X$  the (unique) element  $\mu$  of  $H$  that solves

$$\underset{x \in H}{\operatorname{arginf}} \mathbb{E}[\|X - x\|_H^2].$$

- If  $H=L^2$  the Fréchet mean coincides a.e. with the point-wise mean

$$\mathbb{E}[X(t)] = \mu(t), \quad t \in T$$

- In any  $H$ , we can estimate the mean via the sample estimator

$$\bar{X} = \frac{1}{N} \sum_{i=1}^N X_i$$

In  $H=L^2$ , this is the point-wise sample mean



Let  $X : \Omega \rightarrow H$  be a zero-mean functional random variable in  $H$ , such that  $\mathbb{E}[\|X\|_H^4] < \infty$

**Definition 5:**  
We call covariance operator of  $X$  the operator from  $H$  to  $H$  defined as

$$Cx = \mathbb{E}[\langle X, x \rangle X], \quad x \in H$$

- If  $H=L^2$  the covariance operator can be equivalently defined through a kernel operator

$$[Cx](t) = \int_T c(s, t)x(s)d(s), \quad x \in L^2$$

where the covariance kernel is precisely the point-wise covariance

$$c(s, t) = \mathbb{E}[X(s)X(t)]$$

- In  $H=\mathbb{R}^p$ , the covariance operator coincides with the linear operator defined by the covariance matrix

Let  $X : \Omega \rightarrow H$  be a zero-mean functional random variable in  $H$ , such that  $\mathbb{E}[\|X\|_H^4] < \infty$

**Definition 5:**  
We call covariance operator of  $X$  the operator from  $H$  to  $H$  defined as

$$Cx = \mathbb{E}[\langle X, x \rangle X], \quad x \in H$$

- In any  $H$ , the covariance operator can be estimated through the sample covariance operator

$$Sx = \frac{1}{N} \sum_{i=1}^N \langle X_i, x \rangle X_i, \quad x \in H$$

- If  $H=L^2$ , one can use the alternative definition

$$[Sx](t) = \int_T \hat{c}(s, t)x(s)d(s), \quad x \in L^2 \quad \hat{c}(s, t) = \frac{1}{N} \sum_{i=1}^N X_i(s)X_i(t)$$

$$z_i = f(s_i) + \epsilon_i \quad i = 1, \dots, n \quad \epsilon_1, \dots, \epsilon_n \text{ i.i.d. } E[\epsilon_i] = 0, \text{Var}[\epsilon_i] = \sigma^2$$

$\psi_1, \dots, \psi_K$ :  $K$  basis functions

Basis expansion

$$\hat{f}(s) = \sum_{k=1}^K \hat{c}_k \psi_k(s) \quad \rightarrow \text{Find } \hat{c}_k, k = 1, \dots, K \text{ (i.e., find } \hat{f}) \text{ by minimizing}$$

$$\begin{aligned} \text{SSE} &= \sum_{j=1}^n (z_j - f(s_j))^2 = \sum_{j=1}^n \left( z_j - \sum_{k=1}^K c_k \psi_k(s_j) \right)^2 \\ \Psi &= \begin{bmatrix} \psi_1(s_1) & \psi_2(s_1) & \dots & \psi_K(s_1) \\ \psi_1(s_2) & \psi_2(s_2) & \dots & \psi_K(s_2) \\ \vdots & \vdots & & \vdots \\ \psi_1(s_n) & \psi_2(s_n) & \dots & \psi_K(s_n) \end{bmatrix} \end{aligned}$$

$$\begin{aligned} z &= (z_1, \dots, z_n)^t \\ f &= (f(s_1), \dots, f(s_n))^t \\ c &= (c_1, \dots, c_K)^t \\ \epsilon &= (\epsilon_1, \dots, \epsilon_n)^t \end{aligned}$$

curve fitting

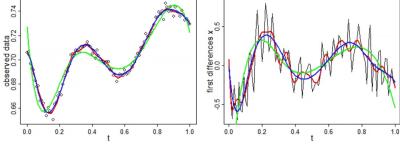
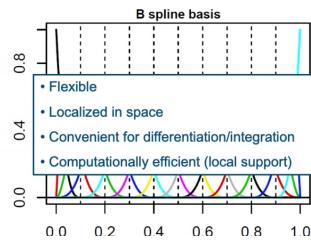
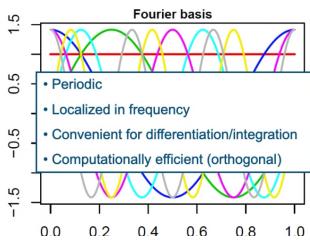
$$\underline{z} = \Psi \underline{c} + \underline{\epsilon}$$

$$\text{SSE} = (z - \Psi c)^t (z - \Psi c)$$

$$\hat{c} = (\Psi^t \Psi)^{-1} \Psi^t z$$

$$\dot{z} = \dot{f} = \Psi \hat{c} = \Psi (\Psi^t \Psi)^{-1} \Psi^t z = S z$$

$$df = K = \text{tr}(S) = \text{tr}(S^t S) = \text{tr}(2S - S^t S)$$



Choose the number K of basis

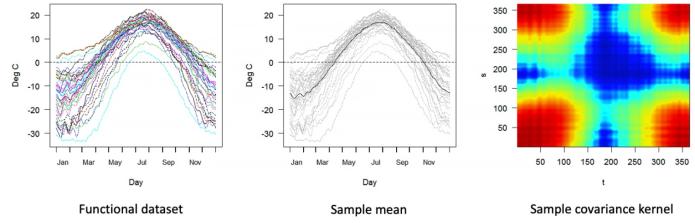
$K << N$

$K = 7$

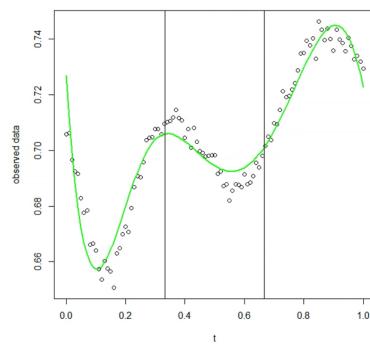
$K = 12$

$K = 30$

Dataset of Temperatures in Canada (35 observations)



1) Use a functional space with only few dimensions (few basis)  $K << n$



$K = 7$   
(spline of order 5)



Distributional properties of the estimator:

$$\psi = (\psi_1, \dots, \psi_K)^t$$

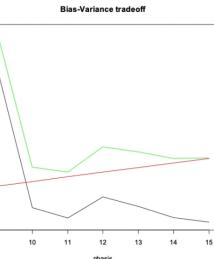
$$f(s) = \psi(s)^t \hat{c} = \psi(s)^t (\Psi^t \Psi)^{-1} \Psi^t z$$

$$E[\hat{f}(s)] = f(s) - E[\hat{f}(s)]$$

$$\text{Bias}[\hat{f}(s)] = f(s) - E[\hat{f}(s)]$$

$$\text{Var}[\hat{f}(s)] = E[\{\hat{f}(s) - E[\hat{f}(s)]\}^2] = \sigma^2 \psi(s)^t (\Psi^t \Psi)^{-1} \psi(s)$$

$$\text{MSE}[\hat{f}(s)] = E[\{\hat{f}(s) - f(s)\}^2] = \text{Bias}^2[\hat{f}(s)] + \text{Var}[\hat{f}(s)]$$



## ► Cross-validation

Distributional properties of the estimator:

$$\psi = (\psi_1, \dots, \psi_K)^t$$

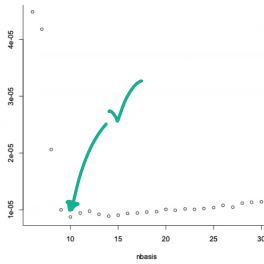
$$\hat{f}(s) = \psi(s)^t \hat{c} = \psi(s)^t (\Psi^t \Psi)^{-1} \Psi^t z$$

$$E[\hat{f}(s)] = f(s) - E[\hat{f}(s)]$$

$$\text{Bias}[\hat{f}(s)] = f(s) - E[\hat{f}(s)]$$

$$\text{Var}[\hat{f}(s)] = E[\{\hat{f}(s) - E[\hat{f}(s)]\}^2] = \sigma^2 \psi(s)^t (\Psi^t \Psi)^{-1} \psi(s)$$

$$\text{MSE}[\hat{f}(s)] = E[\{\hat{f}(s) - f(s)\}^2] = \text{Bias}^2[\hat{f}(s)] + \text{Var}[\hat{f}(s)]$$



Selection of the number of basis  $K$  minimizing

$$\begin{aligned} GCV(K) &= \frac{1}{n(1 - \text{tr}(S)/n)^2} (z - \hat{z})^t (z - \hat{z}) \\ &= \left(\frac{n}{n - K}\right) \left(\frac{1}{n - K}\right) (z - \hat{z})^t (z - \hat{z}) \end{aligned}$$

$$\begin{aligned} \hat{f}'(s) &= \sum_{k=1}^K \hat{c}_k \psi'_k(s) & \hat{f}''(s) &= \sum_{k=1}^K \hat{c}_k \psi''_k(s) \end{aligned}$$

Smoothing requires special care when the curve estimate is asked, not only to provide a good smoothing of the data, but also to reflect the features of the curve that are represented by its derivatives

Curve derivatives (or their functions) are very often

- objects of analysis
- helpful for further processing and analysis of the data (curve alignment/clustering)

2) Use a rich functional space but with regularization

$K \sim n$

$$\text{SSE}_\lambda = \text{SSE} + \lambda \int (f''(s))^2 ds \quad \text{penalize curvature} \quad (\text{potentially a different class})$$

$$\{R_\psi\}_{(k,l)} \quad (k, l)-\text{entry} : \int \psi''_k(s) \psi''_l(s) ds$$

$$\text{SSE}_\lambda = \text{SSE} + \lambda \mathbf{c}^t R_\psi \mathbf{c}$$

$$\hat{c}_\lambda = (\Psi^t \Psi + \lambda R_\psi)^{-1} \Psi^t z$$

$$\hat{z} = \hat{f} = \Psi (\Psi^t \Psi + \lambda R_\psi)^{-1} \Psi^t z = S z \quad \text{Sub-projection operator}$$

$$df = \text{tr}(S) < K \quad (\text{or } df = \text{tr}(S^t S) \text{ or } df = \text{tr}(2S - S^t S))$$

Distributional properties of the estimator:

$$\hat{f}(s) = \psi(s)^t \hat{c} = \psi(s)^t (\Psi^t \Psi + \lambda R_\psi)^{-1} \Psi^t z$$

$$E[\hat{f}(s)] = \psi(s)^t (\Psi^t \Psi + \lambda R_\psi)^{-1} \Psi^t f$$

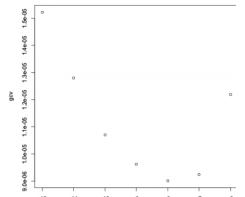
$$\text{Bias}[\hat{f}(s)] = f(s) - E[\hat{f}(s)]$$

$$\begin{aligned} \text{Var}[\hat{f}(s)] &= E[\{\hat{f}(s) - E[\hat{f}(s)]\}^2] = \\ &= \sigma^2 \psi(s)^t (\Psi^t \Psi + \lambda R_\psi)^{-1} (\Psi^t \Psi) (\Psi^t \Psi + \lambda R_\psi)^{-1} \psi(s) \end{aligned}$$

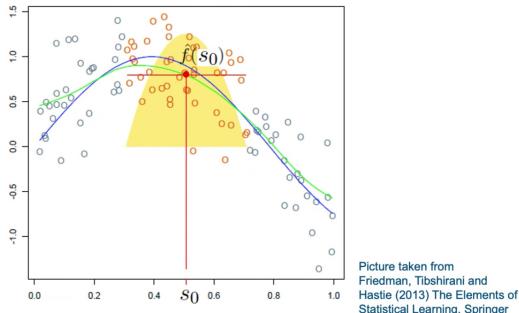
$$\text{MSE}[\hat{f}(s)] = E[\{\hat{f}(s) - f(s)\}^2] = \text{Bias}^2[\hat{f}(s)] + \text{Var}[\hat{f}(s)]$$

Selection of the smoothing parameter  $\lambda$  minimizing

$$\begin{aligned} GCV(\lambda) &= \frac{1}{n(1 - \text{tr}(S)/n)^2} (z - \hat{z})^t (z - \hat{z}) \\ &= \left(\frac{n}{n - \text{tr}(S)}\right) \left(\frac{1}{n - \text{tr}(S)}\right) (z - \hat{z})^t (z - \hat{z}) \end{aligned}$$



Local polynomial regression (kernel smoother)



At each abscissa  $s_0$ , find  $(c_0, \dots, c_L)$  that minimize

$$\sum_{i=1}^n \text{Kern}_h(s_0, s_i) [(z_i - \sum_{l=0}^L c_l (s_0 - s_i)^l)^2] \quad \text{where } \text{Kern}_h(s_0, s_i) = D\left(\frac{|s_0 - s_i|}{h}\right)$$

## Positive functions:

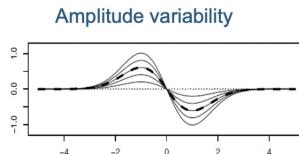
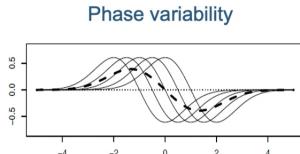
$$f(s) = e^{W(s)} \quad \text{where } W(s) = \sum_k c_k \psi_k(s)$$

Estimate  $f$  by minimizing

$$\sum_{i=1}^n (z_i - e^{W(s_i)})^2 + \lambda \int W''(s) ds$$

## Increasing functions:

$$f(s) = C + \int_{s_0}^s \exp\{W(t)\} dt \quad \text{where } W(s) = \sum_k c_k \psi_k(s)$$



**Phase variability:** different curves exhibit more or less the same features but that these features occur at different times or space locations for different statistical units.

If not taken properly into account, the misalignment acts as a confounding factor and may blur subsequent analyses.

Landmarks: significant (univocally identifiable) shape-events in a curve, e.g. crossings of zero, peaks, valleys, points of inflection.

$c_1, \dots, c_N$ , where  $c_i : [0, T] \rightarrow \mathbb{R}^d$

Suppose

- $L$  landmarks; for the  $i$ -th curve, located at  $t_{i1}, \dots, t_{iL}$
- a template curve  $c_0$  is available with landmark locations  $t_{01}, \dots, t_{0L}$   
If not, we can define  $t_{0j}$  as the average of the  $t_{ij}$ 's

Warping function for the  $i$ -th curve: any strictly increasing function  $h_i$  s.t.

- $h_i(0) = 0$
- $h_i(t_{0j}) = t_{ij}$ , for  $j = 1, \dots, L$
- $h_i(T) = T$

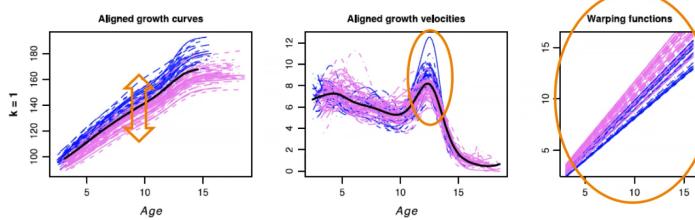
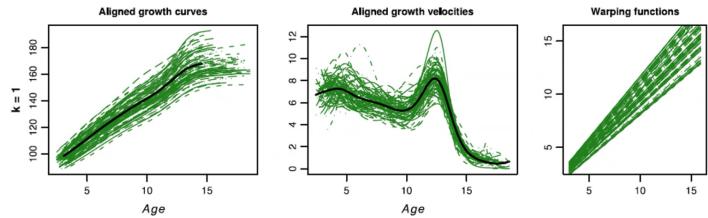
►  $(0, 0), (t_{01}, t_{i1}), \dots, (t_{0L}, t_{iL}), (T, T)$  : interpolated by a piece-wise line, a polygon or higher order monotone splines (strictly increasing)

### Registration of a set of functions

Find suitable warping functions  $h_1(t), \dots, h_N(t)$  such that  $c_1(h_1(t)), \dots, c_N(h_N(t))$  are the most similar.

The functions  $h_i$  should be increasing; they capture the phase variability. Amplitude variability is the remaining variability in vertical direction among the aligned curves.

In some cases, time or location is merely shifted from curve to curve, for example, because the measurements are started at random time points. For these situations, it is natural to use  $h_i(t) = t + \delta t$ . In other situations, phase variation is a matter of dilation, in which case  $h_i(t) = \alpha t$  is a natural choice of warping function. In yet other situations, the time or space deformation is more complex.



Once the biological clocks are aligned

the height of boys  
stochastically dominates the  
one of girls for any registered  
biological age

Aligned growth velocities  
boys have a more  
pronounced growth  
during puberty (more  
prominent growth  
velocity peak)

Warping functions  
Neat separation  
of boys and girls  
in the phase.  
The biological  
clocks of boys  
and girls run at  
different speeds

**Goal of Alignment:**  
Decoupling Phase and Amplitude Variability



**Goal of K-mean Clustering:**  
Decoupling Within and Between-cluster (Amplitude) Variability



**Goal of K-mean Alignment:**  
Identifying Phase Variability, Within-cluster Amplitude Variability  
and Between-cluster Amplitude Variability

## LECTURE 30 30/5/2022

### Off: 21 hidden step: Principal Component Analysis



104

Courtesy of P. Secchi

**Problem:** Given a dataset of  $N$  zero-mean multivariate observations in  $\mathbb{R}^p, X_1, \dots, X_N$  find the orthonormal directions  $a_1, \dots, a_p$  of maximum variability (for the dataset).

Equivalently, for  $k=1, \dots, p$ , find:  $a_k = \operatorname{argmax}_{a \in \mathbb{R}^p} \operatorname{Var}(a'X)$   
subject to:  $a'a = 1, a'j'a = 0 \text{ for } j < k$

- We can re-write the problem as

$$a_k = \operatorname{argmax}_{a \in \mathbb{R}^p} \frac{1}{N} \sum_{i=1}^N (a'X_i)^2$$

subject to:  $a'a = 1, a'j'a = 0 \text{ for } j < k$

or, equivalently

$$a_k = \operatorname{argmax}_{a \in \mathbb{R}^p} \frac{1}{N} \sum_{i=1}^N \langle a, X_i \rangle^2$$

subject to:  $\|a\| = 1, \langle a_j, a \rangle = 0 \text{ for } j < k$

**Note 1.** We assume  $N > p$  and absence of collinearity, i.e. the data matrix is full rank.

**Note 2.** If  $X_1, \dots, X_N$  are not zero-mean, they can be centered by subtracting the (sample) mean. For unbiasedness, divide by  $N-1$  instead of  $N$ .

**Problem:** Given a dataset of  $N$  zero-mean functional observations in  $H, X_1, \dots, X_N$ , find the directions of maximum variability (in  $H$ ) of the dataset, i.e., for  $k=1, \dots, N$ , find  $\xi_k$  maximizing

$$\frac{1}{N} \sum_{i=1}^N \langle \xi, X_i \rangle_H^2$$

subject to:  $\|\xi\| = 1, \langle \xi_j, \xi \rangle_H = 0 \text{ for } j < k$

- We look for an orthonormal system in  $H$  maximizing the variability of the corresponding projections

Indeed,  $\langle \xi, X_i \rangle_H$  is the projection of  $X_i$  along the direction  $\xi$  (i.e., a «direction» in  $H$ ).  
Note that  $\langle \xi, X_i \rangle_H$  is a scalar, hence  $\frac{1}{N} \sum_{i=1}^N \langle \xi, X_i \rangle_H^2$  is a sample variance in the usual sense.

**Note 1.** If the data are not zero-mean, they can be centered by subtracting the (sample) mean.  $N$  should then be replaced by  $N-1$ .

**Note 2.** If data are linearly independent and centered on the sample mean, we can find at most  $N-1$  principal components.

**Problem:** Given a dataset of  $N$  zero-mean functional observations in  $H, X_1, \dots, X_N$ , find the directions of maximum variability (in  $H$ ) of the dataset, i.e., for  $k=1, \dots, N$ , find  $\xi_k$  maximizing

$$\frac{1}{N} \sum_{i=1}^N \langle \xi, X_i \rangle_H^2$$

subject to:  $\|\xi\| = 1, \langle \xi_j, \xi \rangle_H = 0 \text{ for } j < k$

- We look for an orthonormal system in  $H$  maximizing the variability of the corresponding projections

Indeed,  $\langle \xi, X_i \rangle_H$  is the projection of  $X_i$  along the direction  $\xi$  (i.e., a «direction» in  $H$ ).  
Note that  $\langle \xi, X_i \rangle_H$  is a scalar, hence  $\frac{1}{N} \sum_{i=1}^N \langle \xi, X_i \rangle_H^2$  is a sample variance in the usual sense.

**Note 1.** If the data are not zero-mean, they can be centered by subtracting the (sample) mean.  $N$  should then be replaced by  $N-1$ .

**Note 2.** If data are linearly independent and centered on the sample mean, we can find at most  $N-1$  principal components.

$$a_k = \operatorname{argmax}_{a \in \mathbb{R}^p} \frac{1}{N} \sum_{i=1}^N \langle a, X_i \rangle^2$$

subject to:  $\|a\| = 1, \langle a_j, a \rangle = 0 \text{ for } j < k$

**Solution:** Call  $S$  the sample covariance matrix of  $X_1, \dots, X_N$ . Then, the **principal components** are found as the eigenvectors of the matrix  $S$ ; for  $k=1, \dots, p$ , they solve the eigen-equation

$$Se_k = \lambda_k e_k$$

The eigenvalue  $\lambda_k$  associated with the eigenvector  $e_k$  represents the variability along the direction  $e_k$ .

**Note.** We call  $u_{ik}$  the projection of the observation  $X_i$  along the direction  $e_k$ , i.e.,

$$u_{ik} = \langle X_i, e_k \rangle = X_i' e_k$$

