

EXPLORING A MULTIVARIATE DATASET

Prediction problem

$X \in \mathbb{R}^p, Y \in \mathbb{R}$: we want to use X to predict Y .

What is the best function $f(X)$: $f: \mathbb{R}^p \rightarrow \mathbb{R}$ to predict Y ?
 # best! : $f(\underline{x}) = \arg \min \mathbb{E}[(Y - f(\underline{x}))^2] := \text{arg min MSE}$ (mean square error)

$$\text{If } f(\underline{x}) = k \implies k = \arg \min \mathbb{E}[(Y - k)^2] = \mathbb{E}[Y] \quad (*)$$

$$\mathbb{E}[(Y - f(\underline{x}))^2] = \mathbb{E}[(Y - \mathbb{E}[Y|\underline{x}])^2] + \underbrace{\mathbb{E}[(\mathbb{E}[Y|\underline{x}] - f(\underline{x}))^2]}_{\text{constant wrt. } f(\underline{x})}$$

$$\implies f(\underline{x}) = \arg \min \text{MSE} = \arg \min \mathbb{E}[(\mathbb{E}[Y|\underline{x}] - f(\underline{x}))^2] \stackrel{(*)}{=} \mathbb{E}[Y|\underline{x}] \text{ best guess of } Y \text{ once we know } X$$

$$\text{The model is: } \begin{cases} Y = f(\underline{x}) + \varepsilon \\ f(\underline{x}) = \mathbb{E}[Y|\underline{x}] \\ \mathbb{E}[\varepsilon] = 0 \end{cases}$$

Suppose now that we have a model \hat{f} . How good is the model?

$$\hat{f} : \text{estimation of } f \text{ through } \hat{X} : \quad \hat{X} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix}$$

We want to see how good \hat{f} is in estimating new predictions: $\underline{x}_0 \in \mathbb{R}^p$: new observation ($\notin \hat{X}$) for which we want to predict y_0

We want to know about the error:

$$\mathbb{E}[(y_0 - \hat{f}(\underline{x}_0))^2 | \hat{X}] = \mathbb{E}[\varepsilon^2] = \mathbb{E}[\mathbb{E}[\varepsilon^2 | \hat{X}]]$$

we have one realization

$$\text{since } y_0 = f(x_0) + \varepsilon_0 \implies \mathbb{E}[\varepsilon^2 | \hat{X}] = (f(\underline{x}_0) - \hat{f}(\underline{x}_0))^2 + \underbrace{\text{Var}(\varepsilon_0)}_{\text{this is irreducible: it's the part of } Y \text{ not explainable through } \hat{X}}$$

this estimate the error of the prediction given $\hat{X} = X$, what if we want to consider all the possible realizations of X ?

$$\implies \mathbb{E}[\mathbb{E}_{\hat{X}}[(y_0 - \hat{f}(\underline{x}_0))^2]] = \text{Var}(\varepsilon_0) + \text{Var}(\hat{f}(\underline{x}_0)) + \frac{(f(\underline{x}_0) - \mathbb{E}[\hat{f}(\underline{x}_0)])^2}{\text{Bias}^2} \quad (\text{how far is the model from what we want to estimate})$$

How to estimate f ? (How to get \hat{f} ?)

LOCAL AVERAGE

- we react with:
 - reduce of dimensionality (PCA) (data driven)
 - parametric models (knowledge is necessary)

Geometry of the data

We can explore the dataset by rows or by columns:

- by columns:
 - $\hat{X} = [x_1 \ x_2 \ \dots \ x_p]$
 - Every column X_j has n realizations : $\hat{y}_j = [x_{1j}, x_{2j}, \dots, x_{nj}]^T$
 - mean of X_j : $\bar{x}_{ij} := \frac{1}{n} \sum_{i=1}^n x_{ij}$
 - variance of X_j : $S_{jj} := \frac{1}{n} \sum_{i=1}^n (x_{ij} - \bar{x}_{ij})^2$
 - covariance of X_j, X_k : $\text{Cov}(X_j, X_k) := \frac{1}{n} \sum_{i=1}^n (x_{ij} - \bar{x}_{ij})(x_{ik} - \bar{x}_{ik}) := S_{kj}$
 - covariance matrix (for X_1, \dots, X_p):
 - $S := \begin{bmatrix} S_{11} & S_{12} & \dots & S_{1p} \\ S_{21} & S_{22} & \dots & S_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ S_{p1} & S_{p2} & \dots & S_{pp} \end{bmatrix} \in \mathbb{R}^{p \times p}$
 - correlation of X_j, X_k : $r_{kj} = \frac{S_{kj}}{\sqrt{S_{kk} S_{jj}}} = \text{Corr}(x_j, x_k) \in [-1, 1]$
 - correlation matrix (for X_1, \dots, X_p):
 - $R := \begin{bmatrix} r_{11} & r_{12} & \dots & r_{1p} \\ r_{21} & r_{22} & \dots & r_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ r_{p1} & r_{p2} & \dots & r_{pp} \end{bmatrix}$
- whatever the distribution of X_j is, we can say (Chebyshev) :
- $\Pr(X_j - k\sqrt{S_{jj}} \leq X_j \leq \bar{x}_{ij} + k\sqrt{S_{jj}}) \geq 1 - \frac{1}{k^2} \quad \forall k \neq 0$
- alternatively:
- $\Pr(|X_j - \mu| \leq k\sigma) \geq 1 - \frac{1}{k^2}$

Thanks to Chebyshev, whenever we have the mean and the standard deviation we can construct an interval.

Geometrical interpretation:

- one variable



Geometric projection of \underline{y} on $\underline{z}(1)$:

$$\text{Proj}_{\underline{z}(1)} \underline{y} = \frac{\underline{w} \cdot \underline{w}^T}{\underline{w}^T \underline{w}} \cdot \underline{y} \quad \underline{z}(1) = \text{Space with "no statistic"} \\ (\underline{z} \in \mathcal{L}(\underline{z}) \Rightarrow \underline{z} = k \cdot \underline{z})$$

Projection of y_j on $\mathcal{L}(\underline{z})$

$$\text{Proj}_{\underline{z}(1)} y_j = \frac{\underline{z} \cdot \underline{z}^T}{\underline{z}^T \underline{z}} \cdot y_j = \bar{x}_{ij} \quad (\text{how far is the model from what we want to estimate})$$

- two variables



$$y_j = \bar{x}_{ij} \underline{z} + d_j$$

$$y_j = \bar{x}_{ik} \underline{z} + d_k$$

$$d_j = \alpha_{jk} \underline{z}$$

$$\cos(\theta_{jk}) = \frac{d_j^\top d_k}{\|d_j\| \|d_k\|} = \frac{s_{jk}}{\sqrt{s_{jj} s_{kk}}} = r_{jk}$$

- $\theta_{jk} = 0 \Rightarrow r_{jk} = 1 \Rightarrow d_j \in \mathcal{L}(d_k)$
- $\theta_{jk} = \frac{\pi}{2} \Rightarrow r_{jk} = 0 \Rightarrow d_k \perp d_j$

by rows:

$$\underline{X} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \quad \text{where } x_j \text{ realization of } X_j \\ X_1, \dots, X_n \text{ iid } X \in \mathbb{R}^p : \text{ we measure with random vectors } \underline{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_p \end{bmatrix}$$

- mean: $\mathbb{E}[\underline{x}] = \begin{bmatrix} \mathbb{E}[x_1] \\ \vdots \\ \mathbb{E}[x_p] \end{bmatrix} = \mu = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_p \end{bmatrix}$

covariance of X_j and X_k :

$$\sigma_{jk} = \mathbb{E}[(X_j - \mu_j)(X_k - \mu_k)] = [\Sigma]_{jk}$$

covariance matrix (for \underline{X})

$$\Sigma := [s_{ij}] \in \mathbb{R}^{p \times p}$$

$$\Sigma := \mathbb{E}[(X - \mu)(X - \mu)^\top]$$

correlation matrix (for \underline{X}): (§)

$$V = \text{diag}(\sigma_1, \dots, \sigma_p) \Rightarrow V^{-\frac{1}{2}} = \text{diag}\left(\frac{1}{\sqrt{\sigma_1}}, \dots, \frac{1}{\sqrt{\sigma_p}}\right) \\ V^{-\frac{1}{2}} = \text{diag}\left(\frac{1}{\sigma_1}, \dots, \frac{1}{\sigma_p}\right)$$

$$\beta := V^{-\frac{1}{2}} \Sigma V^{\frac{1}{2}}$$

linear combinations of the components of \underline{X} :

$$c \in \mathbb{R}^p : \mathbb{E}[c^\top \underline{X}] = c^\top \mu$$

$$\text{Var}(c^\top \underline{X}) = c^\top \Sigma c$$

k linear combinations of the component of \underline{X} :

$$C \in \mathbb{R}^{k \times p} : \mathbb{E}[C \underline{X}] = C\mu$$

$$\text{Cov}(C \underline{X}) = C \Sigma C^\top$$

Estimators

We consider \underline{X} and we look at it by rows: Can we estimate μ and Σ ? We assume every row to be a realization of a random vector:

n rows \rightarrow in random vectors $\underline{x}_1, \dots, \underline{x}_n \text{ iid } X \in \mathbb{R}^p$

Estimator for μ :

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n \underline{x}_i \quad \text{realization} \rightarrow \bar{X} = \begin{bmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_p \end{bmatrix}$$

Assuming X_1, \dots, X_n iid X , $\mathbb{E}[\underline{x}] = \mu$, $\text{cov}(\underline{x}) = \Sigma$:

- $\mathbb{E}[\bar{X}] = \mu$ (unbiased)
- $\text{Cov}(\bar{X}) = \frac{1}{n} \Sigma$

Estimator for Σ :

$$S = \frac{1}{n} \sum_{i=1}^n (\underline{x}_i - \bar{X})(\underline{x}_i - \bar{X})^\top \quad \text{realization} \rightarrow S = \begin{bmatrix} s_{11} & \cdots & s_{1p} \\ \vdots & \ddots & \vdots \\ s_{p1} & \cdots & s_{pp} \end{bmatrix}$$

$$\mathbb{E}[S] = \frac{n-1}{n} \Sigma \rightarrow \frac{n-1}{n} S \text{ is unbiased for } \Sigma$$

From now on:

$$S = \frac{1}{n-1} \sum_{i=1}^n (\underline{x}_i - \bar{X})(\underline{x}_i - \bar{X})^\top$$

Considering that:

$$d_j = \underline{x}_j - \frac{\bar{x} \cdot \bar{x}^\top}{\bar{x}^\top \bar{x}} \bar{x}_j = \left(I - \frac{\bar{x} \cdot \bar{x}^\top}{\bar{x}^\top \bar{x}} \right) \bar{x}_j \quad (d_j = \underline{x}_j - \frac{\bar{x} \cdot \bar{x}^\top}{\bar{x}^\top \bar{x}} \bar{x}_j)$$

$$d = [d_1, \dots, d_p] = \left(I - \frac{\bar{x} \cdot \bar{x}^\top}{\bar{x}^\top \bar{x}} \right) \underline{X} \quad := \text{deviation matrix}$$

$$S = \frac{1}{n} d^\top d$$

$$= \frac{1}{n-1} \underline{X}^\top (I - \frac{\bar{x} \cdot \bar{x}^\top}{\bar{x}^\top \bar{x}}) \underline{X}$$

Variability in a multivariate sense

In \mathbb{R}^p the variability is explained by variance.
In \mathbb{R}^p the variance is not enough.

Generalized variance: $\text{Det}(S)$

$$p=2 : S = \frac{1}{n-1} d^\top d \Rightarrow \text{Det}(S) = \|d_1\|^2 \|d_2\|^2 \sin^2 \theta \left(\frac{1}{n-1} \right)^2$$

$$\|d_2\| \sin(\theta) \quad \Rightarrow \quad \text{Det}(S) \propto \text{area of parallelogram } (d_1, d_2)$$

$$\bullet \text{det}(S) \uparrow \Rightarrow \text{area } \uparrow : \text{can increase because of } \theta \text{ (max area with } \theta = \pi/2)$$

$$\text{or because of } \|d_1\| / \|d_2\| = 0$$

$$\bullet \text{det}(S) = 0 \Rightarrow \text{either: } \bullet (\|d_1\| = 0 \text{ or } \|d_2\| = 0) \quad \bullet \theta = 0$$

In any case, it doesn't mean that $\#$ variability

$$p=4 : \text{Det}(S) \propto (\text{vol parallelepiped } (d_1, \dots, d_4))^2$$

Total variance: $\text{Tr}(S)$

$$p=2 : \text{Tr}(S) = \frac{1}{n-1} (\|d_1\|^2 + \|d_2\|^2)$$

$$p=4 : \text{Tr}(S) = \frac{1}{n-1} (\|d_1\|^2 + \|d_2\|^2 + \|d_3\|^2 + \|d_4\|^2)$$

(capturing the sum of the marginal variabilities of the variables)

Note: $\text{Det}(S) = 0 \iff \begin{cases} d_1, \dots, d_p \text{ are linearly dependent} \\ (S \neq 0 \text{ st. } c_1 d_1 + \dots + c_p d_p = 0) \end{cases}$

Consequence: $\text{Det}(S) = 0 \iff \begin{cases} d_1, \dots, d_p \text{ are linearly independent} \\ (S \neq 0 \text{ st. } c_1 d_1 + \dots + c_p d_p \neq 0) \end{cases}$

$$\text{Consequence for } k : \quad \exists k \text{ st. } d_k = - \sum_{i \neq k} \frac{c_i}{c_k} d_i$$

$$\Rightarrow y_k = \bar{x}_k - \sum_{i \neq k} \frac{c_i}{c_k} (y_i - \bar{x}_i)$$

(there is a perfect linear relationship between the variable k and all the other variables, so the variable k is useless (is obtainable as function of the others))

Consequence: $\underline{X} \in \mathbb{R}^{n \times p} : \text{If } p \geq n \Rightarrow \text{Det}(S) = 0$

Note that: $S \in \mathbb{R}^{p \times p}$ real and symmetric $\Rightarrow S = \sum_{i=1}^p \lambda_i e_i e_i^\top$

$$\Rightarrow P := [e_1, \dots, e_p] : \Lambda := \text{diag}(\lambda_1, \dots, \lambda_p) : S = P \Lambda P^\top$$

$$\Rightarrow \begin{cases} \det(S) = \prod_{i=1}^p \lambda_i \\ \text{Tr}(S) = \sum_{i=1}^p \lambda_i \end{cases}$$

S is positive semi-definite ($\lambda_i \geq 0 \forall i$), if $\text{Det}(S) > 0$ is pos. def. ($\lambda_i > 0 \forall i$)
Assuming S positive definite ($\text{Det}(S) > 0$):
(and an order $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p > 0$)

$$x, y \in \mathbb{R}^p : d_{S^{-1}}(x, y) = (x - y)^\top S^{-1}(x - y) := \text{Mahalanobis distance (standardized)}$$

$$\text{Note: } S^{-1} = \sum_{i=1}^p \frac{1}{\lambda_i} e_i e_i^\top$$

"Neighborhood" in Mahalanobis' distance:

$$e_{r^2, S^{-1}}(\bar{x}) = \{x \in \mathbb{R}^p : (x - \bar{x})^\top S^{-1}(x - \bar{x}) \leq r^2\}$$

Graphically ($p=2$):



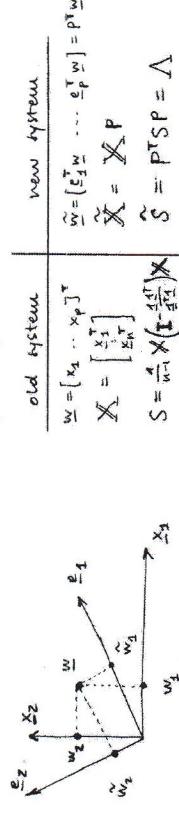
The Mahalanobis' distance is the right one for capturing the distance between data.
(The Euclidean distance is not good, in fact the Mahalanobis' distance is the Euclidean distance after we standardized the data)

Note that: it would be convenient to find a system s.t. S is diagonal w.r.t. that system. Graphically ($p=2$):



in this way the axes of the ellipse
are \perp to the axes of the system

\Rightarrow We introduce a system identified by the eigenvectors of S :



Observation: $\|w\| = \|\tilde{w}\|$

Do we lose information? No:

- $\text{Det}(S) = \prod_{i=1}^p \lambda_i = \text{Det}(\tilde{S})$
- $\text{Tr}(S) = \sum_{i=1}^p \lambda_i = \text{Tr}(\tilde{S})$

There is always a reference system
for which the coordinate vectors are uncorrelated
the data.

PRINCIPAL COMPONENT ANALYSIS (PCA)

Let's refer to the word "word" (not the data):
 \underline{x} random vector in \mathbb{R}^p :

$$\bullet \quad \mathbb{E}[\underline{x}] = \underline{\mu} \\ \bullet \quad \text{Cov}(\underline{x}) = \Sigma$$

Problem: Find $\underline{a} \in \mathbb{R}^p$ s.t. $\|\underline{a}\| = 1$ and $\text{Var}(\underline{a}^\top \underline{x})$ is maximized
(i.e. find the direction \underline{a} s.t. the variability of the projection of \underline{x} on \underline{a} is maximum.)

$$\max_{\underline{a} \in \mathbb{R}^p, \|\underline{a}\|=1} \text{Var}(\underline{a}^\top \underline{x}) = \max_{\underline{a} \in \mathbb{R}^p, \|\underline{a}\|=1} \underline{a}^\top \Sigma \underline{a} = \max_{\underline{a} \in \mathbb{R}^p} \frac{\underline{a}^\top \Sigma \underline{a}}{\underline{a}^\top \underline{a}}$$

Geometrical lemma:
 $B \in \mathbb{R}^{p \times p}$ positive semi-definite, $B = \sum_{i=1}^p \lambda_i e_i e_i^\top$:

$$\left\{ \begin{array}{l} 1. \quad \max_{\underline{x} \in \mathbb{R}^p} \frac{\underline{x}^\top B \underline{x}}{\underline{x}^\top \underline{x}} = \lambda_1 , \quad \text{arg max}(\cdot) = e_1 \\ 2. \quad \max_{\substack{\underline{x} \in \mathbb{R}^p \\ \underline{x} \perp e_1}} \frac{\underline{x}^\top B \underline{x}}{\underline{x}^\top \underline{x}} = \lambda_2 , \quad \text{arg max}(\cdot) = e_2 \\ \vdots \\ p. \quad \max_{\substack{\underline{x} \in \mathbb{R}^p \\ \underline{x} \perp e_1, \dots, \underline{x} \perp e_{p-1}}} \frac{\underline{x}^\top B \underline{x}}{\underline{x}^\top \underline{x}} = \lambda_p = \min_{\substack{\underline{x} \in \mathbb{R}^p \\ \underline{x} \perp \underline{e}_1, \dots, \underline{x} \perp \underline{e}_{p-1}}} \frac{\underline{x}^\top B \underline{x}}{\underline{x}^\top \underline{x}} \end{array} \right.$$

Back to PCA:

$$\max_{\underline{a} \in \mathbb{R}^p} \text{Var}(\underline{a}^\top \underline{x}) = \max_{\substack{\underline{a} \in \mathbb{R}^p \\ \|\underline{a}\|=1}} \frac{\underline{a}^\top \Sigma \underline{a}}{\underline{a}^\top \underline{a}} = \lambda_1 , \quad \text{arg max}(\cdot) = e_1$$

\Rightarrow First Principal Component (PC1): $\underline{Y}_1 = \underline{e}_1^\top \underline{x} \quad (\text{if } Y_1 = \underline{e}_1^\top \underline{x})$
problem: Find $\underline{a} \in \mathbb{R}^p$ s.t. $\|\underline{a}\|=1$, $\text{Var}(\underline{a}^\top \underline{x})$ is max and $\underline{a} \perp e_1$
(i.e. we want to find an other direction of max variability but we want the projection on this direction to be uncorrelated with the previous projection)

$$\max_{\substack{\underline{a} \in \mathbb{R}^p : \|\underline{a}\|=1 \\ \text{Cov}(\underline{a}^\top \underline{x}, \underline{e}_1^\top \underline{x})=0}} \text{Var}(\underline{a}^\top \underline{x}) = \max_{\substack{\underline{a} \in \mathbb{R}^p \\ \underline{a} \perp \underline{e}_1}} \frac{\underline{a}^\top \Sigma \underline{a}}{\underline{a}^\top \underline{a}}$$

$$\Rightarrow \text{Second Principal Component (PC2)}: \underline{Y}_2 = \underline{e}_2^\top \underline{x} \quad (\text{if } Y_2 = \underline{e}_2^\top \underline{x})$$

Generalized problem:

$$\max_{\substack{\underline{a} \in \mathbb{R}^p : \|\underline{a}\|=1 \\ \text{Cov}(\underline{a}^\top \underline{x}, \underline{e}_1^\top \underline{x})=0}} \text{Var}(\underline{a}^\top \underline{x}) = \lambda_k , \quad \text{arg max}(\cdot) = \underline{e}_k$$

\Rightarrow k-th Principal Component (PCk): $\underline{Y}_k = \underline{e}_k^\top \underline{x} \quad (\text{if } Y_k = \underline{e}_k^\top \underline{x})$

- $\underline{Y} := [Y_1 \dots Y_p]^T = P^T \underline{X}$:= vector of Principal Components
- $E[\underline{Y}] = E[P^T \underline{X}] = P^T \underline{\mu}$
- $Cov(\underline{Y}) = Cov(P^T \underline{X}) = P^T \Sigma P = \Lambda$
- $Cov(Y_i, Y_j) = 0 \quad \forall i \neq j$ no correlation between the coordinates of \underline{Y}
- $Cov(Y_i, Y_i) = Var(Y_i) = \lambda_i \quad \forall i$
- Ordering: we ordered the components \Rightarrow the first component is the one with larger variability, the second is the one with the second largest var. and so on \Rightarrow we can capture the most of the variability (not with the first components (the last express small variability))
- We're not losing variability:

 - generalized variance with \underline{Y} : $\det(\Lambda) = \prod_{i=1}^p \lambda_i = \det(\Sigma)$
 - total variance with \underline{Y} : $\text{Tr}(\Lambda) = \sum_i \lambda_i = \text{Tr}(\Sigma)$

- $Y_i = e_i^\top X + \epsilon_i$ ϵ_i 's are the loading (weights):

PCA on standardized variables
(we're still working with the model, not data)

$$V := \text{diag}(\sigma_1, \dots, \sigma_p)$$

$$\underline{X} \longrightarrow \underline{Z} = V^{-\frac{1}{2}} (\underline{X} - \underline{\mu}) = \left[\frac{x_1 - \mu_1}{\sigma_1}, \dots, \frac{x_p - \mu_p}{\sigma_p} \right]^T$$

- $E[Z] = 0$
- $Cov(Z) = V^{-\frac{1}{2}} \Sigma V^{-\frac{1}{2}} := S$:= covariance matrix (covariance matrix of the standardized variables)
- $\rho = \frac{1}{n} \sum_i e_i^\top \underline{e}_i^\top$ (eigenvalues/vectors of S)

$$(Note: \text{Tr}(S) = \sum_i \text{Var}(Y_i) = \sum_i \text{Var}(Z_i) = p)$$

- $\text{PCA}(\Sigma) \neq \text{PCA}(S)$

- If we think that there are variables with very different variabilities (maybe because of the units of measure or maybe because of the phenomena) then it would be better $\text{PCA}(S)$: If X_1 is in km and X_2 in mm \Rightarrow $\text{Var}(X_1)$ will mask $\text{Var}(X_2) \Rightarrow \text{PCA}(S)$

PCA on data

$$\underline{X} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} : \underline{x}_i \text{ realization of } X_i, \quad \underline{x}_1, \dots, \underline{x}_n \stackrel{iid}{\sim} \underline{X} \in \mathbb{R}^p$$

Usually Σ and Σ are unknown, we use data:
 Σ estimates \sum Σ estimates λ $\left\{ \begin{array}{l} \text{PCA on } S \\ \text{PCA on } \Sigma \end{array} \right.$

$$S = \sum_{i=1}^p \lambda_i e_i e_i^\top \xrightarrow{\text{PCA}} \underline{y}_i = \begin{bmatrix} e_1^\top \underline{x}_i \\ e_2^\top \underline{x}_i \\ \vdots \\ e_p^\top \underline{x}_i \end{bmatrix} \xrightarrow{\text{PCA}} \underline{y}_i = \begin{bmatrix} \underline{x}_1^\top \\ \vdots \\ \underline{x}_n^\top \end{bmatrix} = \begin{bmatrix} y_{11} & \dots & y_{1p} \\ \vdots & \ddots & \vdots \\ y_{n1} & \dots & y_{np} \end{bmatrix} = \begin{bmatrix} y_1 \\ \vdots \\ y_p \end{bmatrix}$$

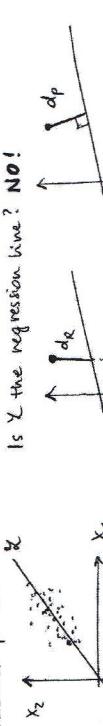
- PCA for categorical data is called "Correspondence analysis" and it's performed on the table of joint frequencies among variables

- we have also to observe the smaller λ_i : if $\lambda_p \approx 0$ it means that \exists linear relationship between $\underline{x}_1, \dots, \underline{x}_k$

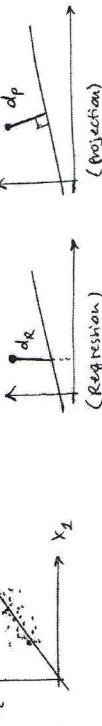
Geometrical meaning of PCA

problem: given $\underline{x}_1, \dots, \underline{x}_n \in \mathbb{R}^p$ find a linear space \underline{L} s.t. $\text{dim}(\underline{L}) = k \leq p$ which best approximate the data (i.e. such that \underline{L} is closest to data).

Consider $p=2$:



Is \underline{L} the regression line? No!



Here we're trying to minimize dp , not d_L ($\neq dp$)!
(Note that the regression depends on the system, the projection no)

Problem:

$$\text{find } \underline{q}_1, \dots, \underline{q}_k \quad (\text{we identify } \underline{L} \text{ with an orthonormal basis})$$

$$\underline{L} = \text{span}(\underline{q}_1, \dots, \underline{q}_k)$$

such that:

$$\sum_{i=1}^n \|(\underline{x}_i - \underline{\Sigma}) - \sum_{j=1}^k q_j q_j^\top (\underline{x}_i - \underline{\Sigma})\|^2 \text{ is minimum.}$$

$(*) = \text{(distance between every single datum (after centering) and the projection on } \underline{L} \text{ of the datum (after centering))}^2$

$$\min(*) \iff \max \sum_{i=1}^k \sum_{j=1}^k q_j q_j^\top (\underline{x}_i - \underline{\Sigma})^2$$

$$\iff \max (n-1) \sum_{i=1}^k q_i^\top S q_i$$

by induction on k :

$$q_1 = \underline{e}_1, \dots, q_k = \underline{e}_k \quad \text{The directions of the first } k \text{ principal components identify the closest linear space (of dim } k \text{) to the data}$$

$$\text{Error of approximation: } (n-1) \sum_{i=1}^k \lambda_i$$

MULTIVARIATE GAUSSIAN DISTRIBUTION

$$\begin{aligned} X &\sim N_p(\mu, \Sigma) \\ f(x) &= \frac{1}{\sqrt{(2\pi)^p \det(\Sigma)}} e^{-\frac{1}{2}(x-\mu)^\top \Sigma^{-1}(x-\mu)} \end{aligned}$$

- $\underline{X} \sim N_p(\mu, \Sigma)$
 - $\underline{X} \sim N_p(\mu, \Sigma) \iff E[\underline{X}] = \mu, \quad \text{Cov}(\underline{X}) = \Sigma$
 - $\underline{X} \sim N_p(\mu, \Sigma) \iff \underline{q}^\top \underline{X} \sim N_q(\underline{\mu}^\top \underline{\mu}, \underline{\Sigma}^\top \underline{\Sigma}) \quad \forall \underline{q} \in \mathbb{R}^q$
(\underline{q} : linear combination of the components is gaussian)
 - $\underline{X} = [X_1, \dots, X_p]^\top \sim N_p(\mu, \Sigma), \quad \Sigma = [\sigma_{ij}] \Rightarrow X_i \sim N_2(\mu_i, \sigma_{ii})$
 - $\underline{X} \sim N_p(\mu, \Sigma), \quad A \in \mathbb{R}^{q \times p} \Rightarrow A\underline{X} \sim N_q(A\mu, A\Sigma A^\top)$
 - $\underline{X} \sim N_p(\mu, \Sigma), \quad d \in \mathbb{R}^p \Rightarrow \underline{X} + \underline{d} \sim N_p(\mu + \underline{d}, \Sigma)$
 - $\underline{X} = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} \sim N_p\left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_1 & \Sigma_{12} \\ \Sigma_{21} & \Sigma_2 \end{bmatrix}\right) \Rightarrow X_1 \sim N_2(\mu_1, \Sigma_1)$
- for example:
- $$\begin{bmatrix} X_1 \\ X_2 \end{bmatrix} \sim N \Rightarrow \begin{bmatrix} X_1 \\ X_3 \end{bmatrix} \sim N$$
- $$X_1 \perp\!\!\!\perp X_2 \iff \Sigma_{12} = \Sigma_{21} = 0$$

$$\underline{X} = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} \sim N_p\left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_1 & \Sigma_{12} \\ \Sigma_{21} & \Sigma_2 \end{bmatrix}\right), \quad \det(\Sigma) \neq 0 :$$

$$\begin{aligned} \rightarrow X_1 | X_2 = x_2 &\sim N_q(\mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (x_2 - \mu_2), \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}) \\ \text{Cov}(X_1 | X_2 = x_2) &= \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \perp\!\!\!\perp X_2 \end{aligned}$$

:= Partial covariances

even without knowing X_2 we already know how informations will be modified

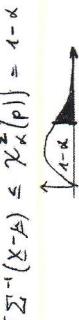
(This generates the Regression Effect)

- $\underline{X} \sim N_p(\mu, \Sigma), \quad \det(\Sigma) > 0 \Rightarrow (\underline{X} - \mu)^\top \Sigma^{-1} (\underline{X} - \mu) \sim \chi^2(p)$
- $\underline{X} \sim N_p(\mu, \Sigma), \quad \det(\Sigma) = 0 \Rightarrow (\underline{X} - \mu)^\top \Sigma^{-1} (\underline{X} - \mu) \sim \chi^2(k)$

$$K = \text{rank}(\Sigma)$$

$$\Sigma = \sum_{i=1}^k \frac{1}{\lambda_i} \underline{e}_i \underline{e}_i^\top$$

$$(\lambda_1 > \lambda_2 > \dots > \lambda_K > 0 = \lambda_{K+1} = \dots = \lambda_p)$$



It doesn't mean

that if the sample is large than it is gaussian. It means that if a sample is large enough then the sample mean is gaussian

statistically normal

measuring: for large n one can approximate the distribution of $\sqrt{n}(\bar{X} - \mu)$ with a $N_p(0, \Sigma)$

In practice, for large n :

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

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

- $\bar{X} \sim N_p(\mu, \frac{1}{n} \Sigma)$
 - $S \sim \text{Wish}(\frac{1}{n-1} \Sigma, n-1)$
 - $\hat{\Sigma} \sim \text{Wish}(\frac{1}{n} \Sigma, n-1)$
- About Wishart distribution (matrix distribution):
- (Def.) $\bar{X}_1, \dots, \bar{X}_m \stackrel{iid}{\sim} N_p(\underline{0}, \Sigma), \quad \det(\Sigma) > 0$
- $$\Rightarrow \sum_{i=1}^m \bar{X}_i \bar{X}_i^\top \sim \text{Wishart}(\Sigma, m)$$
- (the parameter p is in Σ)

Properties:

- $A_1 \sim \text{Wish}(\Sigma, m_1), \quad A_2 \sim \text{Wish}(\Sigma, m_2), \quad A_2 \perp\!\!\!\perp A_1$
- $\Rightarrow A_1 + A_2 \sim \text{Wish}(\Sigma, m_1 + m_2)$
- $A \sim \text{Wish}(\Sigma, m), \quad C \in \mathbb{R}^{k \times p} \text{ const}$
- $\Rightarrow C A C^\top \sim \text{Wish}(C \Sigma C^\top, m)$
- $A \sim \text{Wish}(\Sigma, m), \quad \sigma^2 > 0$
- $\Rightarrow \sigma^2 A \sim \text{Wish}(\sigma^2 \Sigma, m)$
- $\bullet A \sim \text{Wish}(\Sigma, m), \quad \Sigma \perp\!\!\!\perp A \sim \Sigma \text{Exp}(\mu)$

Note that: $\bar{X} \perp\!\!\!\perp S$ and are sufficient statistics (i.e. if the data is generated by a gaussian distribution (no matter m) then all we need to know is \bar{X}, S)

LLN

- X_1, \dots, X_n random vectors iid such that $E[X_i] = \mu, \quad \text{Cov}(X_i) = \Sigma$ exists:
- $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{P} \mu$ as $n \rightarrow \infty$
 - $S = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})^\top \xrightarrow{P} \Sigma$ as $n \rightarrow \infty$

CLT

- X_1, \dots, X_n random vectors iid such that $E[X_i] = \mu, \quad \text{Cov}(X_i) = \Sigma$ exists:
- $$\Rightarrow \sqrt{n}(\bar{X} - \mu) \xrightarrow{D} N_p(0, \Sigma)$$

It doesn't mean

that if the sample is large than it is gaussian. It means that if a sample is large enough then the sample mean is gaussian

statistically normal

measuring: for large n one can approximate the distribution of $\sqrt{n}(\bar{X} - \mu)$ with a $N_p(0, \Sigma)$

In practice, for large n :

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

Estimators for μ and Σ

$$\hat{X} = \begin{bmatrix} X_1^\top \\ \vdots \\ X_n^\top \end{bmatrix} : \quad X_i \text{ realization of } X_i \Rightarrow X_1, \dots, X_n \stackrel{iid}{\sim} N_p(\mu, \Sigma)$$

- $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ (unbiased)
- $S = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})^\top$ (biased, but MVE)
- $\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})^\top$

INFERENCE FOR THE MEAN μ

$X_1, \dots, X_n \in \mathbb{R}^p$ random vectors s.t. $E[X] = \mu$
 $(X_1, \dots, X_n) \stackrel{iid}{\sim} \bar{X}$ $Cov(\frac{X}{n}) = \Sigma$
 $\text{Det } (\Sigma) > 0$

$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$: pointwise estimator for μ

Two cases:

• n large ($n \gg p$):

CLT: $\sqrt{n}(\bar{X} - \mu) \sim N_p(0, \Sigma)$

Pivotal statistic dist.:

$$n(\bar{X} - \mu)^T S^{-1}(\bar{X} - \mu) \sim \chi^2(p)$$

If Σ is not known:

$$\text{LLN: } S \xrightarrow{n \rightarrow \infty} \Sigma$$

Pivotal statistic dist.:

$$n(\bar{X} - \mu)^T S^{-1}(\bar{X} - \mu) \sim \chi^2(p)$$

• Confidence regions:

$P(d_{S^{-1}}(\bar{X}, \mu) \leq \chi^2_\alpha(p)) = P(\mu \in \Sigma^{-1}(\bar{X})) = 1 - \alpha$
 (i.e. the random ellipse centered in \bar{X} and shaped with S
 will be covering the true (fixed) value of μ with $(1 - \alpha) \cdot 100\%$
 of the times that we generate it)

$$CR_{1-\alpha}(\mu) = \{\eta \in \mathbb{R}^p : n(\eta - \bar{X})^T S^{-1}(n\eta - \bar{X}) \leq \chi^2_\alpha(p)\}$$

Testing :

$$\begin{cases} H_0: \mu = \mu_0 \\ H_1: \mu \neq \mu_0 \end{cases}$$

Test statistic: $T_0^2 = n(\bar{X} - \mu_0)^T S^{-1}(\bar{X} - \mu_0)$

If H_0 is true $\Rightarrow T_0^2 \sim \chi^2(p)$

Rejection Region $_\alpha$ = { $T_0^2 > \chi^2_\alpha(p)$ }



$$T_0^2 > \chi^2_\alpha(p) \iff \text{p-value} \leq \alpha$$

• n small

Here we have to add: $X_1, \dots, X_n \stackrel{iid}{\sim} N_p(\mu, \Sigma)$, $\text{det}(\Sigma) > 0$
 Pivotal statistic dist.:

$$n(\bar{X} - \mu)^T S^{-1}(\bar{X} - \mu) \sim \frac{(n-1)p}{n-p} F(p, n-p)$$

Hotteling's theorem: $\underline{X} \sim N_p(\mu, \Sigma)$ ($\text{det}(\Sigma) > 0$) $\perp\!\!\!\perp W \sim \text{Wish}(\frac{1}{m}\Sigma, m)$
 $\Rightarrow \frac{1}{m-p+1} (\underline{X} - \mu)^T W^{-1} (\underline{X} - \mu) \sim F(p, m-p+1)$

• Confidence regions:

$$P(n(\bar{X} - \mu)^T S^{-1}(\bar{X} - \mu) \leq \frac{(n-1)p}{n-p} F_\alpha(p, n-p)) = 1 - \alpha$$

$$CR_{1-\alpha}(\mu) = \{\eta \in \mathbb{R}^p : n(\bar{X} - \mu)^T S^{-1}(\bar{X} - \mu) \leq \frac{(n-1)p}{n-p} F_\alpha(p, n-p)\}$$

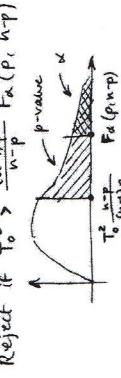
• Testing:

$$\begin{cases} H_0: \mu = \mu_0 \\ H_1: \mu \neq \mu_0 \end{cases}$$

$$T_0^2 = n(\bar{X} - \mu_0)^T S^{-1}(\bar{X} - \mu_0)$$

If H_0 is true $\Rightarrow T_0^2 \sim \frac{(n-1)p}{n-p} F_\alpha(p, n-p)$

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



CI for linear combinations of the mean

$X_1, \dots, X_n \stackrel{iid}{\sim} N_p(\mu, \Sigma)$, small n case
 $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ estimator for μ

• One-at-the-time CI($\underline{a}^T \mu$) ($1 \leq q \leq p$)

Let $\underline{a} \in \mathbb{R}^p$, $CI_{1-\alpha}(\underline{a}^T \mu)$?

$$\frac{\underline{a}^T \bar{X} - \underline{a}^T \mu}{\sqrt{\underline{a}^T \Sigma \underline{a}}} \sqrt{n} \sim t(n-1)$$

$\Rightarrow P(\underline{a}^T \mu \in [\underline{a}^T \bar{X} \pm t_{\alpha/2}(n-1) \sqrt{\frac{\underline{a}^T \Sigma \underline{a}}{n}}]) = 1 - \alpha$

$$\Rightarrow CI_{1-\alpha}(\underline{a}^T \mu) = [\underline{a}^T \bar{X} \pm t_{\alpha/2}(n-1) \sqrt{\frac{\underline{a}^T \Sigma \underline{a}}{n}}]$$

(Testing: $\begin{cases} H_0: \underline{a}^T \mu = \delta_0 \\ H_1: \underline{a}^T \mu \neq \delta_0 \end{cases}$)

Reject if:

$$\frac{|\underline{a}^T \bar{X} - \delta_0|}{\sqrt{\underline{a}^T \Sigma \underline{a}}} \sqrt{n} > t_{\alpha/2}(n-1)$$

Examples:

$$\begin{cases} \underline{a} = [0 \dots 0 \dots 0]^T \\ \underline{a} = [0 \dots 0 \dots 1 \dots 0 \dots 0]^T \end{cases} \Rightarrow \begin{cases} \underline{a}^T \mu = \mu_1 \\ \underline{a}^T \mu = \mu_i - \mu_j \end{cases}$$

• Simultaneous CI($\underline{a}^T \mu$) ($\alpha \geq 1 - \alpha$)

$$P\left(\frac{|\underline{a}^T(\bar{X} - \mu)|}{\sqrt{\underline{a}^T \Sigma \underline{a}}} \sqrt{n} \leq \sqrt{\frac{(n-1)p}{n-p}} F_\alpha(p, n-p), \forall \underline{a} \in \mathbb{R}^p\right) = 1 - \alpha$$

$$\text{Sim } CI_{1-\alpha}(\underline{a}^T \mu) = \left[\underline{a}^T \bar{X} \pm \sqrt{\frac{(n-1)p}{n-p}} F_\alpha(p, n-p) \sqrt{\frac{\underline{a}^T \Sigma \underline{a}}{n}} \right]$$

- Bonferroni's correction ($k \leq$, break K)

Given $\alpha_1, \dots, \alpha_k \in \mathbb{R}^p$ fixed $CI(\alpha_i^\top \mu)$, ..., $CI(\alpha_k^\top \mu)$ with simultaneous confidence of $1-\alpha$, $\alpha \in (0,1)$:

$$\text{sim } CI_{1-\alpha}(\alpha_i^\top \mu) = \left[\bar{\alpha}_i^\top \bar{\mu} \pm t_{\frac{\alpha}{2k}}(n-1) \sqrt{\frac{\alpha_i^\top \Sigma \alpha_i}{n}} \right]$$

Testing K assumptions simultaneously

$X_1, \dots, X_n \stackrel{iid}{\sim} N_p(\mu, \Sigma)$ $\det(\Sigma) > 0$

Given $\alpha_1, \dots, \alpha_K \in \mathbb{R}^p$:

$$\begin{cases} H_0 : \left\{ \begin{array}{l} \alpha_1^\top \mu = \beta_1 \\ \vdots \\ \alpha_K^\top \mu = \beta_K \end{array} \right. \\ H_A : \exists i : \alpha_i^\top \mu \neq \beta_i \end{cases}$$

Bonferroni method:

$$\text{reject at level } \alpha \text{ if for at least one } i: \frac{|\alpha_i^\top \bar{\Sigma} - \beta_i|}{\sqrt{\alpha_i^\top \Sigma \alpha_i}} \sqrt{n} > t_{\frac{\alpha}{2K}(n-1)}$$

		not-rejecting H_0		rejecting H_0
		H_0	H_A	
Truth:	H_0	U	V	$\mu_0 = \# \text{ true hypotheses}$
	H_A	T false	S	$K - \mu_0 = \# \text{ false hypotheses}$

$K - R$
* not rejected
* rejected

- Family-wise error rate : $FWER := P(V \geq 1) =$ probability that we'll reject at least one of the true null hypotheses

- False discovery rate : $FDR := E[\frac{V}{K}]$

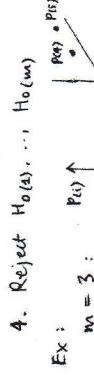
Property:

$$FDR \leq FWER$$

Methods:

- Bonferroni guarantees : $FWER \leq \alpha$ but with large K it becomes a problem
- Benjamini and Hochberg strategy for controlling FDR:

- For each of the K tests compute the p-value p_i
- Order the p-values: $p_{(1)} \leq p_{(2)} \leq \dots \leq p_{(K)}$
- $m := \max\{i \in \{1, \dots, K\} : p_{(i)} \leq \frac{\alpha}{k}\}$
- Reject $H_0(1), \dots, H_0(m)$



Ex: $m = 3$: $p(1) \uparrow$ $p(2) \uparrow$ $p(3) \uparrow$
p-value order:

$$H_0(3) \neq H_0(2)$$

reject at level α :

$$H_0(3) \neq H_0(2)$$

In this way, if the p-values are \ll α \rightarrow $FDR \leq \alpha$

Comparing means of Gaussian distributions

- Paired data

Each unit is observed 2 times:

$$\begin{aligned} \underline{X}_{1i} &= \begin{bmatrix} X_{11i} \\ X_{12i} \\ \vdots \\ X_{1Ki} \end{bmatrix} & \underline{X}_{2i} &= \begin{bmatrix} X_{21i} \\ X_{22i} \\ \vdots \\ X_{2Ki} \end{bmatrix} \\ \underline{X}_{1i} &\text{ iid with mean } \mu_1 \\ \underline{X}_{2i} &\text{ iid with mean } \mu_2 \end{aligned}$$

Goal: inference on $\mu_1 - \mu_2$
Assumptions: $\underline{D}_i := \underline{X}_{1i} - \underline{X}_{2i} \stackrel{iid}{\sim} N_p(\underline{\delta}, \Sigma_D)$

(we're not assuming Σ_D is known)

Pivotal statistic:

$$n(\bar{\Delta} - \underline{\delta})^\top S_D^{-1}(\bar{\Delta} - \underline{\delta}) \sim \frac{(n-1)p}{np} F(p, n-p)$$

- $CR_{1-\alpha}(\mu_1 - \mu_2) = \{\underline{\delta} \in \mathbb{R}^p : n(\bar{\Delta} - \underline{\delta})^\top S_D^{-1}(\bar{\Delta} - \underline{\delta}) \leq \frac{(n-1)p}{np} F_\alpha(p, n-p)\}$
- $H_0 : \mu_1 - \mu_2 = \underline{\delta}_0 \text{ vs. } H_1 : \mu_1 - \mu_2 \neq \underline{\delta}_0$
Reject at level α if: $n(\bar{\Delta} - \underline{\delta}_0)^\top S_D^{-1}(\bar{\Delta} - \underline{\delta}_0) > \frac{(n-1)p}{np} F_\alpha(p, n-p)$
- Sim CI $\alpha(\mu_1 - \mu_2) = \left[\bar{\Delta}_i \pm \sqrt{\frac{(n-1)p}{np} F_\alpha(p, n-p)} \sqrt{\frac{1}{n} \sum_{j=1}^n \Sigma_{jj}} \right]$
- Confidence interval: $\bar{\Delta}_i \pm t_{\frac{\alpha}{2p}}(n-1) \sqrt{\frac{1}{n} \sum_{j=1}^n \Sigma_{jj}}$

- Repeated univariate measure

Each unit is observed q times
 $\underline{X}_i = [X_{i1}, \dots, X_{iq}]^\top$
 $i = 1, \dots, n$
Ex. n patients,
 $X_{ij} = \text{blood pressure of patient } i \text{ at time } j$

$$\mathbb{E}[\underline{X}_i] = \underline{\mu} = [\mu_1, \dots, \mu_q]^\top$$

Goal: $H_0 : \mu_1 = \mu_2 = \dots = \mu_q \text{ vs. } H_1 : \exists i, j \text{ s.t. } \mu_i \neq \mu_j$
We introduce the contrast variable C ,

examples: $C = \begin{bmatrix} 1 & -1 & 0 & \dots & -1 \\ 0 & 1 & -1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \end{bmatrix}$, $c = \begin{bmatrix} \alpha_1 & \dots & \alpha_q \end{bmatrix}$

Test: $H_0 : C \underline{\mu} = 0 \text{ vs. } H_1 : C \underline{\mu} \neq 0$

Pivotal statistic:

$$n(C \bar{\Sigma} - C \underline{\mu})^\top (C S C)^{-1} (C \bar{\Sigma} - C \underline{\mu}) \sim \frac{(n-1)(q-1)}{n-q+1} F(q-1, n-q+1)$$

Reject at level α if:

$$n(C \bar{\Sigma} - C \underline{\mu})^\top (C S C)^{-1} (C \bar{\Sigma} - C \underline{\mu}) > F_\alpha(q-1, n-q+1)$$

(*) the first is testing $\mu_1 - \mu_2, \mu_2 - \mu_3, \dots$
the second is testing $\mu_2 - \mu_1, \mu_3 - \mu_1, \dots$

M(ANOVA) ONE-WAY

Cases:

P	# Features	# groups	Description
1	≥ 1	2	We have n patients and one treatment. We apply the treatment to n_1 patients and not to n_2 patients. For each patient we have p features
2	1	≥ 2	We have n patients and one treatment. We apply the treatment at different levels: for n_1 patients we don't apply for n_2 pairs we apply level q . For each patients we have q feature
3	≥ 1	≥ 2	Like in the second case but for each patient we have p features.

Case 1 :

$$\begin{aligned} \bar{X}_{11}, \dots, \bar{X}_{1n_1} &\stackrel{\text{iid}}{\sim} N_p(\mu_1, \Sigma) \\ \bar{X}_{21}, \dots, \bar{X}_{2n_2} &\stackrel{\text{iid}}{\sim} N_p(\mu_2, \Sigma) \end{aligned} \quad \text{Goal: Inference on } \mu_1 - \mu_2$$

$$\begin{aligned} \mu_1 &\leftarrow \bar{X}_1 = \frac{1}{n_1} \sum_{i=1}^{n_1} \bar{X}_{1i} \sim N_p(\mu_1, \frac{1}{n_1} \Sigma) \\ \mu_2 &\leftarrow \bar{X}_2 = \frac{1}{n_2} \sum_{i=1}^{n_2} \bar{X}_{2i} \sim N_p(\mu_2, \frac{1}{n_2} \Sigma) \\ S_{\text{pooled}} &:= \frac{(n_1-1)S_1 + (n_2-1)S_2}{n_1+n_2-2} \end{aligned}$$

$$\text{with } \begin{cases} S_1 = \frac{1}{n_1-1} \sum_{i=1}^{n_1} (\bar{X}_{1i} - \bar{\bar{X}}_1)(\bar{X}_{1i} - \bar{\bar{X}}_1)^T & (\Sigma \text{ in group 1}) \\ S_2 = \frac{1}{n_2-1} \sum_{i=1}^{n_2} (\bar{X}_{2i} - \bar{\bar{X}}_2)(\bar{X}_{2i} - \bar{\bar{X}}_2)^T & (\Sigma \text{ in group 2}) \end{cases}$$

Partial statistic:

$$(\frac{1}{n_1} + \frac{1}{n_2})^{-1} [(\bar{X}_1 - \bar{\bar{X}}_2) - (\mu_1 - \mu_2)]^T S_{\text{pooled}}^{-1} [\bar{X}_1 - \bar{\bar{X}}_2] - (\mu_1 - \mu_2) \sim$$

$$\sim \frac{(n_1+n_2-2)p}{n_1+n_2-2-p} F(p, n_1+n_2-2-p)$$

$$\text{Testing: } \begin{cases} H_0: \mu_1 - \mu_2 = \delta_0 \\ H_1: \mu_1 - \mu_2 \neq \delta_0 \end{cases}$$

Reject at level α if:

$$(\frac{n_1}{n_1+n_2})^{-1} [(\bar{X}_1 - \bar{\bar{X}}_2) - \delta_0]^T S_{\text{pool}}^{-1} [(\bar{X}_1 - \bar{\bar{X}}_2) - \delta_0] > \frac{(n_1+n_2-2)p}{n_1+n_2-2-p} F_\alpha(p, n_1+n_2-2-p)$$

Confidence region:

$$CR_{1-\alpha}(\mu_1 - \mu_2) = \left\{ \underline{\delta} \in \mathbb{R}^p : \left(\frac{1}{n_1} + \frac{1}{n_2} \right)^{-1} [(\bar{X}_1 - \bar{\bar{X}}_2) - \delta] \leq \underline{\delta} \right\}$$

Note: for n_1 and n_2 very large:

$$[[(\bar{X}_1 - \bar{\bar{X}}_2) - (\mu_1 - \mu_2)]^* \left(\frac{S_1}{n_1} + \frac{S_2}{n_2} \right)^{-1} [(\bar{X}_1 - \bar{\bar{X}}_2) - (\mu_1 - \mu_2)]] \sim \mathcal{N}^2(p)$$

Case 2 : $p = 1, q \geq 2$

$$\begin{aligned} X_{11}, \dots, X_{1n_1} &\stackrel{\text{iid}}{\sim} N_1(\mu_1, \sigma^2) \\ X_{21}, \dots, X_{qn_2} &\stackrel{\text{iid}}{\sim} N_1(\mu_2, \sigma^2) \end{aligned} \quad \text{HOMO}$$

$$\begin{cases} H_0: \mu_1 = \mu_2 = \dots = \mu_q \\ H_1: \exists i \neq j : \mu_i \neq \mu_j \end{cases}$$

$$\text{Parametrization: } \mu_i = \mu + \tau_i \quad i = 1, \dots, q$$

overall mean

$$\Rightarrow X_{ij} = \mu + \tau_i + \varepsilon_{ij} \quad \varepsilon_{ij} \sim N_1(0, \sigma^2)$$

$$\text{Constraint: } \sum_{i=1}^q n_i \tau_i = 0$$

$$\begin{aligned} \mu &\leftarrow \bar{X} = \frac{1}{n} \sum_{i=1}^q \sum_{j=1}^{n_i} X_{ij} \\ \tau_i &\leftarrow \bar{X}_i - \bar{X} = \left[\frac{1}{n_i} \sum_{j=1}^{n_i} X_{ij} - \bar{X} \right] \end{aligned}$$

Variance decomposition:

$$\sum_{i=1}^g \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)^2 = \underbrace{\sum_{i=1}^g (\bar{x}_i - \bar{x})^2 n_i}_{SS_{\text{treatment}}} + \underbrace{\sum_{i=1}^g \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)^2}_{SS_{\text{residuals}}}$$

Pivotal statistic:

$$F_0 = \frac{SS_{\text{treatment}} / (g-1)}{SS_{\text{residual}} / (n-g)} \sim F(g-1, n-g)$$

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

Proposals:

Λ	reject H_0 if Λ is ..
WILKS: $\Lambda_W = \frac{\det(\mathbf{W})}{\det(\mathbf{W} + \mathbf{B})}$	small
MANLEY-HOTELING: $\Lambda_M = \text{Tr}(\mathbf{B}\mathbf{W}^{-1})$	large
PIKAI: $\Lambda_P = \text{Tr}(\mathbf{B}(\mathbf{B} + \mathbf{W})^{-1})$	large

(interpretations on words)

(*) when all the n_i are big enough:
Barlett's approximation:

$$-\left(n-1 - \frac{p+q}{2}\right) \log \Lambda_W \sim \chi^2(p(q-1))$$

Reject at level α if:

$$-\left(n-1 - \frac{p+q}{2}\right) \log \Lambda_W > \chi^2(p(q-1))$$

If we reject H_0 : at what level the effect made effect?

(i.e. we want to compare Σ_i and Σ_k component-wise
so we discover if level i produced a different effect
than level k and on what component)

→ Bonferroni CI ($\tau_{ii} - \tau_{ik}$) $\quad k, i = 1, \dots, q$
 $i = 1, \dots, p$

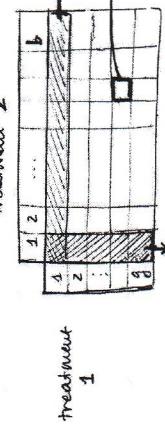
$$\bullet \bar{x}_{it} - \bar{x}_{ik} \sim N(\tau_{ii} - \tau_{ik}, \frac{1}{n_i} \sigma_{ii} + \frac{1}{n_k} \sigma_{kk})$$

$$\bullet \frac{1}{n_i} \text{W estimate } \Sigma \Rightarrow \frac{\text{W}i}{n_i} \text{ estimate } \sigma_{ii}$$

$$\text{Bonferroni CI}_{1-\alpha}(\tau_{ii} - \tau_{ik}) = \left[\bar{x}_{ie} - \bar{x}_{ke} \pm t_{\alpha/2} \sqrt{\frac{n_i \sigma_{ii}}{n_i - g} + \frac{1}{n_k} \left(\frac{n_i \sigma_{ii}}{n_i - g} + \frac{1}{n_k} \right)} \right]$$

(M) ANOVA TWO-WAYS

We have two treatments



$i = \text{level treat. 1}$
 $k = 1, \dots, n$
(we have n patients in this cell)

Covariance decomposition:

$$\sum_{i=1}^g \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)^2 = \underbrace{\sum_{i=1}^g (\bar{x}_i - \bar{x})^2 n_i}_{\text{between}} + \underbrace{\sum_{i=1}^g \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)^2}_{\text{within}}$$

covariability that we would carry out under H_0

$H_0: \Sigma_1 = \Sigma_2 = \dots = \Sigma_g = \Sigma$

$H_1: \exists \Sigma_i \neq \Sigma$

treatment:

$$\bar{x}_i \leftarrow \bar{x}_i - \bar{x}$$

Goal: $H_0: \Sigma_1 = \Sigma_2 = \dots = \Sigma_g = \Sigma$

$H_1: \exists \Sigma_i \neq \Sigma$

ANOVA

$$\sum_{i=1}^g \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)^2 = \underbrace{\sum_{i=1}^g (\bar{x}_i - \bar{x})^2 n_i}_{\text{between}} + \underbrace{\sum_{i=1}^g \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)^2}_{\text{within}}$$

covariance decomposition:

$$\sum_{i=1}^g \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)(\bar{x}_j - \bar{x})^T n_i = \underbrace{\sum_{i=1}^g (\bar{x}_i - \bar{x})^2 n_i}_{\text{between}} + \underbrace{\sum_{i=1}^g \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)(\bar{x}_j - \bar{x})^T}_{\text{within}}$$

we want again see $\frac{B}{n}$ to decide whether to reject or not

$B = \text{B}_{\text{between}}$

$n = \text{N}$

$\bar{x}_{ij} = \text{mean in cell } (i, j)$

$\bar{x}_i = \text{mean in row } i$

$\bar{x}_j = \text{mean in column } j$

$\bar{x} = \text{mean overall}$

Model: $X_{ijk} = \mu + \underbrace{\tau_i}_{\text{treat. interactions}} + \underbrace{\beta_j}_{\text{treat. interactions}} + \underbrace{\gamma_{ij}}_{\text{treat. interactions}} + \varepsilon_{ijk}$

Constraints:

$$\sum_{i=1}^q \tau_i = 0 \quad \sum_{j=1}^b \beta_j = 0 \quad \sum_{i=1}^q \gamma_{ij} = 0$$

Decomposition of variance:

$$\begin{aligned} \sum_{i=1}^q \sum_{j=1}^b \sum_{k=1}^n (\bar{x}_{ijk} - \bar{x})^2 &= \sum_{i=1}^q b n (\bar{x}_{i\cdot} - \bar{x})^2 + \\ &\quad + \sum_{j=1}^b q n (\bar{x}_{\cdot j} - \bar{x})^2 + \\ &\quad + \sum_{i=1}^q \sum_{j=1}^b \left(\bar{x}_{ij\cdot} - \bar{x}_{i\cdot} - \bar{x}_{\cdot j} + \bar{x} \right)^2 n + \text{SS interactions} \\ &\quad \left. \begin{array}{l} \text{this term up} \\ \text{in the additive} \\ \text{model} \\ \text{(with no interaction)} \end{array} \right) \end{aligned}$$

$$\begin{cases} H_0: \gamma_{ij} = 0 \\ H_1: \exists \gamma_{ij} \neq 0 \end{cases} \rightarrow \frac{\text{Reject at level } \alpha \text{ if:}}{\frac{\frac{(g-1)(b-1)}{\text{SS interactions}}}{\frac{g(b-1)}{\text{SSres}}}} \rightarrow F_\alpha \left(\frac{g-1}{g(b-1)}, g(b-1) \right)$$

If we accept $H_0 \Rightarrow$ additive model, and so:

$$\begin{cases} H_0: \tau_1 = \tau_2 = \dots = \tau_q = 0 \\ H_1: \exists \tau_i \neq 0 \end{cases} \rightarrow \frac{\text{Reject at level } \alpha \text{ if:}}{\frac{\frac{g-1}{\text{SS treat. 1}}}{\frac{g-1}{\text{SS interactions}}}} \rightarrow F_\alpha \left(\frac{g-1}{g(b-1)}, g(b-1) \right)$$

(SUPERVISED CLASSIFICATION)

Input features:

- $X|L=i \sim f_i(\underline{x})$ = distribution of the features in group 1
- $P(L=i) = p_i$ = prior probabilities
- $c(i,j) = \text{cost of misclassification}$ ($c(i,j) = \text{what we pay if we attribute to } i \text{ a unit belonging to } j$)

Optimality criterion:

- (Note that instead of defining d we define a partition R_i)

Goal: min $\text{ECM}(\delta)$ (Expected cost of misclassification of δ)

- $g = 2$

$$R_1 = \{ \underline{x} \in \mathbb{R}^p : c(1|2)f_2(\underline{x})p_2 \leq c(2|1)f_1(\underline{x})p_1 \}$$

$$R_2 = \{ \underline{x} \in \mathbb{R}^p : c(2|1)f_2(\underline{x})p_2 \leq c(1|2)f_1(\underline{x})p_1 \}$$

Explanation:

we put \underline{x} in R_2 if:

$$\text{cost} \left(\frac{\underline{x} \text{ belonging to } 2}{\underline{x} \text{ attributed to } 1} \right) \leq \text{cost} \left(\frac{\underline{x} \text{ belonging to } 1}{\underline{x} \text{ attributed to } 2} \right)$$

We're choosing where to put \underline{x} basing on the minimum cost that we have to pay if we're wrong

- $g \geq 2$

$$R_1 = \{ \underline{x} \in \mathbb{R}^p : \sum_{k=2}^g c(1|k)f_k(\underline{x})p_k \leq \sum_{k \neq j} c(j|k)f_k(\underline{x})p_k, j = 2, \dots, g \}$$

$$R_2 = \{ \underline{x} \in \mathbb{R}^p : \sum_{k=2}^g c(2|k)f_k(\underline{x})p_k \leq \sum_{k \neq j} c(j|k)f_k(\underline{x})p_k, j = 2, \dots, g \}$$

$$R_i = \{ \underline{x} \in \mathbb{R}^p : \sum_{k \neq i} c(i|k)f_k(\underline{x})p_k \leq \sum_{k \neq j} c(j|k)f_k(\underline{x})p_k, j \neq i \}$$

Unsupervised

$$\mathbf{X} = \begin{bmatrix} x_1 & x_2 & \dots & x_p & L \\ x_{11} & x_{12} & \dots & x_{1p} & l_1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} & l_n \end{bmatrix}$$

- Goal: • estimate tables (CUSTER ANALYSIS)
 $\{X = [x_1, \dots, x_p] \in X\}$
 (DISCRIMINANT ANALYSIS)

CLASSIFICATION

Supervised

$$\text{cost} \left(\frac{\underline{x} \text{ belonging to } k (k \neq l, g) }{\underline{x} \text{ attributed to } j} \right) \leq \text{cost} \left(\frac{\underline{x} \text{ belonging to } j (j \neq k) }{\underline{x} \text{ attributed to } j} \right)$$

- we chose to put \underline{x} in R_2 if the cost of misclassification is minimum

$$\delta(\underline{x}) = i \quad \longleftrightarrow \quad \underline{x} \in R_i$$

Optimal classifier:

$$\delta(x) = i \iff \left[\sum_{k \neq i} c(i|k) f_k(x) p_k \leq \sum_{k \neq j} c(j|k) f_k(x) p_k \quad \forall j \neq i \right]$$

$$\iff \left[\sum_{k \neq i} c(i|k) P(L=k | X=x) \leq \sum_{k \neq j} c(j|k) P(L=k | X=x) \quad \forall j \neq i \right]$$

(i.e. the expected posterior cost \leq all the other expected costs)

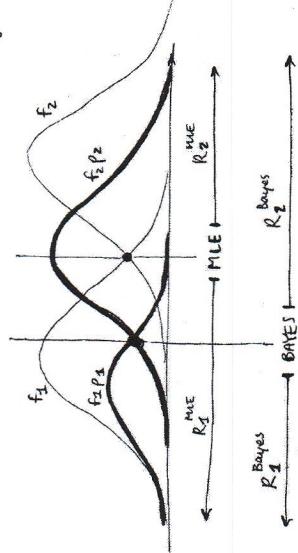
• BAYES CLASSIFIER : all costs are constant (and equal)

$$f_i(x) = i \iff P(L=i | X=x) \leq P(L=j | X=x) \quad \forall j \neq i$$

(i.e. we attribute x to i if the posterior probability of belonging to group i is maximum)

$$f_i(x) = \text{const} \quad \forall i, \quad p_1 = p_2 = \dots = p_g = \frac{1}{g}$$

$$f_j(x) \leq f_i(x) \quad \forall j \neq i$$



Bayes classifier

Bayes classifier is more flexible than it seems.

If we have costs (so Bayes assumptions seem to be violated) we can modify the priors in order to take into account the costs and then we go back to Bayes classifier. (Costs and priors play a similar role)

e.g. $c(i|k) = c_k \geq 0$ (we pay c_k for every unit belonging to k not attributed to k (whenever it goes))

$$\pi_k := \frac{c_k p_k}{\sum c_l p_l}$$

$$\Rightarrow \pi_i = \frac{c_i p_i}{\sum c_j p_j} \leq \frac{c_i p_i}{\sum c_k p_k} \leq \frac{c_i p_i}{\sum c_k p_k} \leq \frac{c_i p_i}{c_k p_k} = \frac{c_i}{c_k} \pi_k$$

Special Bayes classifiers

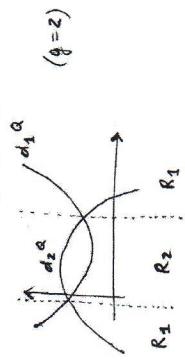
• QDA (Quadratic Discriminant Analysis)

$$X | L=i \sim N_p(\mu_i, \Sigma_i)$$

\rightarrow we can substitute the generic $P(L=i | X=x) \leq P(L=j | X=x)$:

$$R_i = \{x \in \mathbb{R}^p : d_i^Q(x) \geq d_j^Q(x), \quad j=1, \dots, g\}$$

$$d_i^Q(x) = \log(p_i) - \frac{1}{2} \log(\det(\Sigma_i)) - \frac{1}{2} (x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i)$$



• LDA (Linear Discriminant Analysis)

$$X | L=i \sim N_p(\mu_i, \Sigma)$$

$$R_i = \{x \in \mathbb{R}^p : d_i^L(x) \geq d_j^L(x), \quad j=1, \dots, g\}$$

$$d_i^L(x) = \log(p_i) + \mu_i^T \Sigma^{-1} x - \frac{1}{2} \mu_i^T \Sigma \mu_i$$

Fisher's argument for LDA

Suppose $X | L=i \sim (\mu_i, \Sigma)$ (no gaussianity)

Goal: we want to find the direction \mathbf{a} which maximizes the variability BETWEEN groups w.r.t. the variability WITHIN groups

$$\mathbf{a} = \begin{cases} \mathbf{a}_1 = \sum_{i=1}^n \mathbf{z}_i \\ \mathbf{a}_2 = \sum_{i=1}^n \mathbf{z}_i^2 \end{cases} \quad \rightarrow \quad \bar{\mu}_1 = \frac{1}{g} \sum_{i=1}^g \mu_i$$

$$\text{Estimates: } \hat{\mu}_i = \bar{X}_i, \quad \hat{\Sigma} = \frac{1}{n-g} \sum_{i=1}^g (n_i-1) S_i$$

Classifier:
1. $\bar{X}_i \rightarrow [\mathbf{a}_1^T \bar{X}_i, \dots, \mathbf{a}_k^T \bar{X}_i]^T$
2. observation $x \rightarrow [\mathbf{a}_1^T x, \dots, \mathbf{a}_k^T x]^T$

we consider only the first k projections
BAYES CLASSIFIER

$$R_i = \{x \in \mathbb{R}^p : \sum_{j \neq i} f_k(x) \pi_k \leq \sum_{k \neq j} f_k(x) \pi_k \quad j=1, \dots, g\}$$

↓

We are assigning \underline{x} to the closest mean \bar{x}_i :



Evaluation of the parameters

$$\bullet \text{ QDA} : \hat{\mu}_k = \frac{1}{n_k} \sum_{i: l_i=k} \underline{x}_i = \bar{x}_k$$

$$\hat{\Sigma}_k = \frac{1}{n_k - 1} \sum_{i: l_i=k} (\underline{x}_i - \bar{x}_k)(\underline{x}_i - \bar{x}_k)^T = S_k$$

$$\bullet \text{ LDA} : \hat{\mu}_k = \bar{x}_k$$

$$\hat{\Sigma} = \frac{1}{n-g} \sum_{k=1}^g S_k (n_k - 1)$$

Evaluating a classifier

We have to estimate the Actual Error Rate of δ ($AER(\delta)$).

- non-parametric estimate
we apply δ to the training set and we compute the confusion matrix (suppose we have $g=2$):

observed \underline{x}			
\underline{x}	\underline{x}_2	\underline{x}_1	\underline{x}_2
actual \underline{x}	n_{11}	n_{12}	n_{21}
	n_{22}		n_{12}

- leave-one-out cross validation
For $i=1, \dots, n$

- we take out \underline{x}_i from the training set
- we train δ on $\mathcal{X}_{-i} \Rightarrow \delta_{-i}$
- $\delta_{-i}(\underline{x}_i) = \hat{l}_i$

- $\epsilon_i = \begin{cases} 1 & \text{if } \hat{l}_i \neq l_i \\ 0 & \text{if } \hat{l}_i = l_i \end{cases}$

$$\Rightarrow \hat{AER}(\delta) = \frac{\sum_{i=1}^n \epsilon_i}{n}$$

K-fold cross validation

Equal to leave-one-out but we take out a set of k elements and not only one element. This is to reduce the variance of $AER(\delta)$ estimated with leave-one-out

K-fold cross-validation algorithm:

- Set $k \in \mathbb{N}$ and randomly split the units of the training set in k parts (randomly = permute the rows randomly) and then split in k parts)

For $j=1, \dots, k$

- Take out part j from training set
- Train δ on $\mathcal{X}_{-part j} \Rightarrow \delta_{-part j}$
- Apply $\delta_{-part j}$ to the part j :

$$Err_j = \frac{1}{n_j} \sum_{i \in \text{part } j} \epsilon_i$$

$$\epsilon_i = \begin{cases} 1 & \text{if } \delta_{-part j}(\underline{x}_i) \neq l_i \\ 0 & \text{if } \delta_{-part j}(\underline{x}_i) = l_i \end{cases}$$

$$\Rightarrow \hat{AER}(\delta) = \frac{1}{n} \sum_{j=1}^k n_j Err_j$$

Note: by initializing B times we got: $\hat{AER}_1(\delta), \dots, \hat{AER}_B(\delta)$

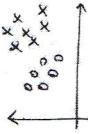
$$\rightarrow \hat{E}[\hat{AER}(\delta)] = \frac{1}{B} \sum_{j=1}^B \hat{AER}_j(\delta)$$

$$\rightarrow \text{Var}(\hat{AER}(\delta)) = \frac{4}{B-1} \sum_{j=1}^B (\hat{AER}_j(\delta) - \hat{E}[\hat{AER}(\delta)])^2$$

$$CIT_{1\alpha}(\hat{E}[\hat{AER}(\delta)]) = \left[\hat{E}[\hat{AER}(\delta)] \pm \sqrt{\frac{\text{Var}(\hat{AER}(\delta))}{B}} \cdot Z_{\alpha/2} \right] \text{ (CLT)}$$

Support Vector Machines (SVM)

(if we have a set (suppose $g=2$))



- SEPARATING HYPERPLANE

$$\begin{aligned} \text{To identify the hyperplane we} \\ \text{introduce a vector } \beta \text{ s.t. } \beta \perp \Sigma, \|\beta\|=1 \\ \Sigma_0 \in \text{span}(\beta) \cap \Sigma \\ \beta_0 = \|\Sigma_0\| \end{aligned}$$



$AER(\delta)$ estimated with leave-one-out

A generic point \underline{x} can (or not) be on Σ :

$$\underline{x} \in \Sigma \iff \nabla_{\underline{x}} \cdot \nabla_{\underline{x}}(\zeta(\underline{p})) = \underline{x}_0$$

$$\iff \underline{\beta}^T \underline{x} \in \Sigma \iff \underline{\beta}^T \underline{x} = \beta_0$$

For example, $\underline{x} \in \Sigma \iff \underline{\beta}^T \underline{x} \geq \beta_0 \iff \underline{\beta}^T \underline{x} > \beta_0$.

$\underline{\beta}^T \underline{x} - \beta_0$ measures (with sign) the distance between \underline{x} and Σ

let say: $\begin{cases} \text{plus} & \rightarrow \begin{cases} \underline{\beta}^T \underline{x} - \beta_0 > 0 & \Rightarrow \underline{x} \in \text{plus} \\ \underline{\beta}^T \underline{x} - \beta_0 < 0 & \Rightarrow \underline{x} \in \text{minus} \end{cases} \\ \text{minus} & \rightarrow \end{cases}$

Suppose we have 2 groups:

$$\begin{cases} y_i = 1 & \text{if } \ell_i = 1 \\ y_i = -1 & \text{if } \ell_i = 2 \end{cases} \quad (= \text{"plus"} = 1 \quad \text{"minus"} = 2)$$

Distance between \underline{x}_i and Σ : $y_i (\underline{\beta}^T \underline{x}_i - \beta_0) \geq 1$

Optimal separating plane:

$$\Sigma \text{ (so } \underline{\beta} \text{ and } \beta_0 \text{) s.t.} \max_{\underline{\beta}, \beta_0} M$$

under the constraints:

$$\begin{cases} \|\underline{\beta}\| = 1 \\ y_i (\underline{\beta}^T \underline{x}_i - \beta_0) \geq 1 \quad i = 1, \dots, n \end{cases}$$

- SEPARATING HYPERPLANE
we can either:

• Allow some overlapping:

$$\max_{\underline{\beta}, \beta_0} M \quad \text{s.t.} \quad \begin{cases} \|\underline{\beta}\| = 1 \\ y_i (\underline{\beta}^T \underline{x}_i - \beta_0) \geq M(1 - \varepsilon_i) \\ \varepsilon_i \geq 0 \\ \sum \varepsilon_i \leq c \end{cases} \quad \leftarrow \text{Budget constraint}$$

- Use non-linear boundaries

UNSUPERVISED CLASSIFICATION

Idea: units belonging to the same group are more similar than units belonging to other groups.

Dissimilarity functions:

$$1. d(\underline{x}, \underline{y}) = \sqrt{(\underline{x} - \underline{y})^T (\underline{x} - \underline{y})} = \sqrt{\sum_{i=1}^p (x_i - y_i)^2} \quad \text{Euclidean}$$

$$2. d_{\Sigma^{-1}}(\underline{x}, \underline{y}) = \sqrt{(\underline{x} - \underline{y})^T \Sigma^{-1} (\underline{x} - \underline{y})} \quad \text{Mahalanobis}$$

$$3. d(\underline{x}, \underline{y}) = \left(\sum_{i=1}^p |x_i - y_i|^m \right)^{1/m} \quad \text{L}^m \text{ distance (Minkowski)}$$

$$4. d(\underline{x}, \underline{y}) = \sum_{i=1}^p \frac{|x_i - y_i|}{x_i + y_i} \quad \text{cosine}$$

$$\text{Dissimilarity matrix:} \quad \underline{X} = \begin{bmatrix} \underline{x}_1^T \\ \vdots \\ \underline{x}_n^T \end{bmatrix} \quad \rightarrow \quad d_{ij} = d(\underline{x}_i, \underline{x}_j) \quad \rightarrow \quad D = \begin{bmatrix} 0 & d_{12} & \dots \\ d_{12} & 0 & \dots \\ \vdots & \vdots & \ddots \end{bmatrix}$$

Dissimilarity between clusters:

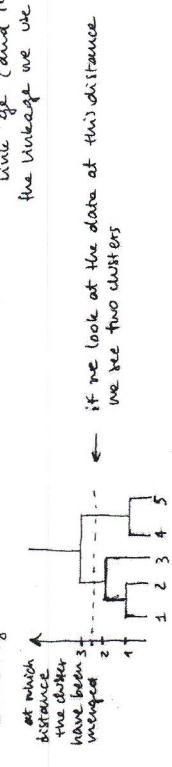
$\underline{U}, \underline{V}$ two sets of points: $d(\underline{U}, \underline{V}) = ?$

- Single linkage: $d(\underline{U}, \underline{V}) = \min \{ d(\underline{x}, \underline{y}) : \underline{x} \in \underline{U}, \underline{y} \in \underline{V} \}$
- Complete linkage: $d(\underline{U}, \underline{V}) = \max \{ d(\underline{x}, \underline{y}) : \underline{x} \in \underline{U}, \underline{y} \in \underline{V} \}$
- Average linkage: $d(\underline{U}, \underline{V}) = \frac{1}{\#\underline{U} \#\underline{V}} \sum_{\underline{x} \in \underline{U}} \sum_{\underline{y} \in \underline{V}} d(\underline{x}, \underline{y})$

Hierarchical Agglomerative Clustering Algorithm

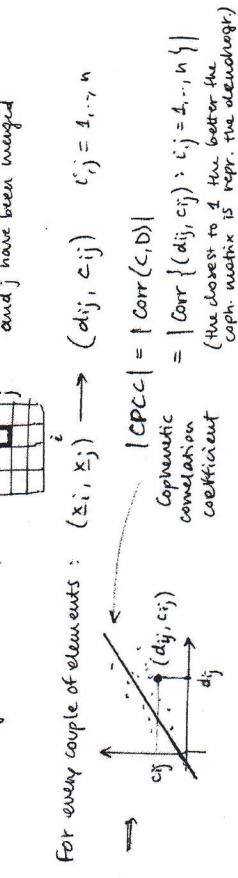
- Set a d and a linkage
- Initialization: every unit is a cluster
- Until convergence repeat:
 1. merge the two cluster which are less similar
 2. compute the new distance matrix

Geographical representation:
Dendrogram



- Here we use the link of (and in the linkage we use it)
- Another representation of the dendrogram:

Dendrogram →
Cophenetic distance



Ward's method for Hierarchical clustering

Suppose we split data into k groups.
Let C_1, \dots, C_k be clusters.

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

\downarrow
barycenter of
the cluster

$$ESS := ESS_1 + \dots + ESS_k$$

At the next iteration we merge the two cluster which generate the minimum increase of ESS.
(It had focus on the minimization of the loss of informations due to merging clusters)

K-means : non-hierarchical method of clustering

Def. Given a cluster $C_j \subseteq$ training set, we call centroid of C_j :

$$\bar{\bar{x}}_j = \arg \min_{\bar{x} \in \mathbb{R}^p} \sum_{x_i \in C_j} d^2(\bar{x}_i, \bar{x})$$

Optimal clustering: find C_1, \dots, C_k s.t. $\sum_{j=1}^k \left[\sum_{x_i \in C_j} d^2(\bar{x}_i, \bar{\bar{x}}_j) \right]$ is min.

k-means algorithm:

- Initialization:
 - randomly create C_1, \dots, C_k in the training set (\Rightarrow step 1)
 - randomly assign k centroids: $\bar{\bar{x}}_1, \dots, \bar{\bar{x}}_k$ (\Rightarrow step 2)
- Iterate until convergence:
 - for $j = 1, \dots, k$ compute the centroids of C_j $\Rightarrow \bar{\bar{x}}_1, \dots, \bar{\bar{x}}_k$
 - for each stat. unit x_i :
 - assign x_i to cluster C_j if: $d^2(x_i, \bar{\bar{x}}_j) = \min \{ d^2(x_i, \bar{\bar{x}}_1), \dots, d^2(x_i, \bar{\bar{x}}_k) \}$

How to choose k ?

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



Multidimensional scaling

Given the distances* among n statistical units, look for the k -dimensional representation of the n statistical units s.t. the distances* among the representations of the n units are as close as possible to the original distances* among the n units.
* distances = dissimilarities

REGRESSION

General goal: explain Y in terms of X .
More specifically the regression is concentrated in estimating $E[Y | X = x] = f(x)$
(f := Regression function)

Two basically approach:

1. Totally data-driven (non-parametric) \rightarrow classification and regression trees
2. Parametric approach (model based)

Linear Models

$$X \longrightarrow Z = \begin{bmatrix} 1 & z_1 & \dots & z_r \\ & \vdots & & \vdots \\ 1 & z_{n1} & \dots & z_{nr} \end{bmatrix} \in \mathbb{R}^{n \times r}$$

linear model: $E[Y | z_1, \dots, z_r] = \beta_0 + \beta_1 z_1 + \dots + \beta_r z_r$

Model for Y :

$$\begin{aligned} Y &= Z \beta + \varepsilon \\ \beta &= [\beta_0, \beta_1, \dots, \beta_r]^T \in \mathbb{R}^{r+1} \\ \varepsilon &\in \mathbb{R}^n \text{ s.t. } \begin{cases} E[\varepsilon] = 0 \\ \text{Cov}(\varepsilon) = \sigma^2 I \end{cases} \end{aligned}$$

i.e. $Y_i = \beta_0 + \beta_1 z_{i1} + \dots + \beta_r z_{ir} + \varepsilon_i$

z_{i1}, \dots, z_{in} :

- unrelated
- same variance
- mean 0
- independent of the z 's

Estimating β and σ^2

(i.e. fitting the model)

OLS : Ordinary least squares

$$\begin{aligned} \hat{\beta} &= \arg \min_{\beta} \| Y - Z \beta \|^2 \\ Z \hat{\beta} &= \hat{Y} \\ &\quad \text{fitted values} \\ Y - \hat{Y} &= \hat{\varepsilon} \\ &\quad \text{residuals} \end{aligned}$$

Goal: find $\hat{Y} \in \mathcal{L}(Z)$ closest to Y
(closest in the sense of Euclidean distance in \mathbb{R}^n)

Note: we can't take $r+1$ too large because we need degrees of freedom for the estimation of variance

Cases:

- rank(Z) = $r+1 < n$ (Z full rank)

$$\hat{Y} = Z(Z^T Z)^{-1} Z^T Y := HY \quad H = \text{IT}_{n,1} Z^{(r)}$$

$$\hat{\beta} = (Z^T Z)^{-1} Z^T Y$$

$$\hat{\Sigma} = (I-H)Y$$

$$\hat{Y} = Z(Z^T Z)^{-1} Z^T Y \quad (\text{dim}(Z/Z)) = k$$

$$\hat{Y} = Z(Z^T Z)^{-1} Z^T Y \quad (\exists! \hat{Y})$$

$$\hat{\beta} = (Z^T Z)^{-1} Z^T Y \quad (\exists! \hat{\beta})$$

$$(Z^T Z)^{-1} = \sum_{i=1}^k \frac{1}{\lambda_i} e_i e_i^T \quad \text{where } Z^T Z = \sum_{i=1}^{r+1} \lambda_i e_i e_i^T$$

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k = 0 = \lambda_{k+1} = \dots = \lambda_n$$

$$\hat{Y} = Y, \quad \hat{\Sigma} = 0$$

Coefficients of determination

$$\sum_{i=1}^n Y_i^2 = \underbrace{\sum_{i=1}^n \hat{Y}_i^2}_{SS_{\text{TOT}}} + \underbrace{\sum_{i=1}^n \hat{\epsilon}_i^2}_{SS_{\text{RES}}}$$

Decomposition of variance:

$$\sum_{i=1}^n (Y_i - \bar{Y})^2 = \underbrace{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}_{CSS_{\text{TOT}}} + \underbrace{\sum_{i=1}^n \hat{\epsilon}_i^2}_{SS_{\text{RES}}}$$

(centered sum of squares)

$$\Rightarrow R^2 = 1 - \frac{\sum_{i=1}^n \hat{\epsilon}_i^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} = 1 - \frac{SS_{\text{RES}}}{CSS_{\text{TOT}}}$$

Remember that it doesn't hold if $\hat{\Sigma} \notin \mathbb{R}$ (in that case: $R^2 = 1 - \frac{\|\hat{\Sigma}\|^2}{\|\Sigma\|^2}$)

$$\Rightarrow R^{\text{adj}} = 1 - \frac{\sum_{i=1}^n (Y_i - \bar{Y})^2}{n-1}$$

Properties of $\hat{\beta}$ and $\hat{\Sigma}$:

- $E[\hat{\beta}] = \beta$
- $\text{Cov}(\hat{\beta}) = \sigma^2 (Z^T Z)^{-1}$
- $E[\hat{\Sigma}] = 0$
- $\text{Cov}(\hat{\Sigma}) = \sigma^2 (I-H)$
- $E[\hat{\Sigma}^T \hat{\Sigma}] = E\left[\sum_{i=1}^n \hat{\epsilon}_i^2\right] = \sigma^2 (n-(r+1))$
- $\Rightarrow \Sigma^2 := \frac{\hat{\Sigma}^T \hat{\Sigma}}{n-(r+1)} \text{ is s.t. } E[\Sigma^2] = \sigma^2$

From now on: $\hat{\Sigma} \sim N_n(\beta, \sigma^2 I)$:

- $\hat{\beta}$ and $\hat{\Sigma}^2 = \frac{\hat{\Sigma}^T \hat{\Sigma}}{n}$ are i.i.d.
- $\hat{\beta} \sim N_{r+1}(\beta, \sigma^2 (Z^T Z)^{-1})$
- $\hat{\Sigma} \sim N_n(0, \sigma^2 (I-H))$
- $\hat{\Sigma}^T \hat{\Sigma} = \sum_{i=1}^n \hat{\epsilon}_i^2 \sim \sigma^2 \chi^2(n-(r+1))$

Inference:

$$\begin{cases} \frac{1}{\sigma^2} (\hat{\beta} - \beta)^T (Z^T Z) (\hat{\beta} - \beta) \sim (r+2) F(r+1, n-(r+1)) \\ \text{CRITICAL}(\alpha) = \left\{ \hat{\beta} \in \mathbb{R}^r : \frac{1}{\sigma^2} (\hat{\beta} - \beta)^T (Z^T Z) (\hat{\beta} - \beta) \leq F_{\alpha}(r+1, n-(r+1)) \cdot (r+1) \right\} \end{cases}$$

$$\begin{cases} \frac{(n-(r+1)) \Sigma^2}{\sigma^2} \sim \chi^2(n-(r+1)) \\ CI_{1-\alpha}(\sigma^2) = \left\{ \sigma^2 \in \mathbb{R} : \frac{(n-(r+1)) \Sigma^2}{\chi^2_{1-\alpha}(n-(r+1))} \leq \sigma^2 \leq \frac{(n-(r+1)) \Sigma^2}{\chi^2_{\alpha}(n-(r+1))} \right\} \end{cases}$$

Simultaneous CI for $\hat{\beta}^T \hat{\beta}$:

$$\text{Sim CI}_{1-\alpha}(\hat{\beta}^T \hat{\beta}) = \underbrace{\left[\hat{\beta}^T \hat{\beta} \pm \sqrt{\frac{1}{\sigma^2} (Z^T Z)^{-1} \alpha} \cdot \sqrt{S^2(r+1) F_{1-\alpha}(r+1, n-(r+1))} \right]}_{\text{Simultaneous CI for every linear combination of } \beta_i}$$

Special case
Sim CI_{1-\alpha}(\beta_i) = [\hat{\beta}_i \pm \sqrt{\text{diag}[Z^T Z]^{-1} \cdot S^2(r+1) F_{1-\alpha}(r+1, n-(r+1))}]

$$S^2(Z^T Z)^{-1} = \text{cov}(\beta)$$

$$= \text{VCOV}(\beta)$$

Testing the β 's

$$\begin{cases} H_0: \beta_1 = 0 \\ H_1: \beta_1 \neq 0 \end{cases} : \frac{1}{S^2} (\hat{\beta}_1)^T (C(Z^T Z)^{-1} C^T)^{-1} (\hat{\beta}) \sim F(p, n-(r+1))$$

Reject H_0 : at level α if: $\frac{1}{S^2} (\hat{\beta}_1)^T (C(Z^T Z)^{-1} C^T)^{-1} (\hat{\beta}) > F_{r+1, n-(r+1)}$

Special case I:

$$\begin{cases} H_0: \beta_1 = \beta_2 = \dots = \beta_{r-1} = 0 \\ H_1: \exists \beta_j \neq 0 \end{cases} \quad (\text{testing } p \text{ parameters})$$

we're comparing $Y = Z\beta + \varepsilon$ vs. $Y = Z_1\beta_1 + \varepsilon_1$
 $\Rightarrow \frac{SS_{\text{res}}(Z_1) - SS_{\text{res}}(Z)}{S^2 p}$ ~ $F(r, n-(r+1))$

where $SS_{\text{res}}(Z_1) - SS_{\text{res}}(Z) = \varepsilon_1^T \varepsilon_1 - \varepsilon^T \varepsilon$

Special case II:

$$\begin{cases} H_0: \beta_1 = \beta_2 = \dots = \beta_r = 0 \\ H_1: \exists \beta_j \neq 0 \end{cases}$$

$$\Rightarrow \frac{SS_{\text{res}}(Z_1) - SS_{\text{res}}(Z)}{S^2 r} = \frac{\sum_{i=1}^n (Y_i - \bar{Y})^2 - \sum_{i=1}^n \hat{Y}_i^2}{r} \sim F(r, n-(r+1))$$

Prediction

$$Y = Z\beta + \varepsilon$$

$$Z_0 = [1 \ Z_1 \ \dots \ Z_{r-1}]^T \rightarrow Y_0? \quad \begin{cases} Y_0 = Z_0^T \hat{\beta} + \varepsilon_0 \\ E[Y_0] = Z_0^T \beta \end{cases} \quad (\varepsilon_0 \perp \varepsilon)$$

$$Y = \hat{\beta}_0 + \hat{\beta}_1 Z_1 \quad \hat{\beta}_0 + \hat{\beta}_1 Z_1$$

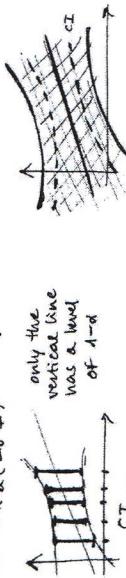


$$\frac{Z_0^T \hat{\beta} - Z_0^T \beta}{S \sqrt{Z_0^T (Z^T Z)^{-1} Z_0}} \sim t(n-(r+1))$$

$$CI_{1-\alpha}(Z_0^T \hat{\beta}) = \left[Z_0^T \hat{\beta} \pm S \sqrt{Z_0^T (Z^T Z)^{-1} Z_0} t_{1-\alpha/2}(n-(r+1)) \right]$$

$$\text{Sim } CI_{1-\alpha}(Z_0^T \hat{\beta}) = \left[Z_0^T \hat{\beta} \pm S \sqrt{Z_0^T (Z^T Z)^{-1} Z_0} \sqrt{(r+1) F_{r+1, n-(r+1)}} \right]$$

With $CI_{1-\alpha}(Z_0^T \hat{\beta})$ we generate a CI for any Z_0 but fixed.
 With $\text{Sim } CI_{1-\alpha}(Z_0^T \hat{\beta})$ we generate a band of CI (so for any Z_0)



Note: this is the prediction of $E[Y|Z_0]$, not Y_0 !

$$\frac{Y_0 - Z_0^T \hat{\beta}}{S \sqrt{1 + Z_0^T (Z^T Z)^{-1} Z_0}} \sim t(n-(r+1))$$

$$PI_{1-\alpha}(Y_0) = \left[Z_0^T \hat{\beta} \pm S \sqrt{1 + Z_0^T (Z^T Z)^{-1} Z_0} t_{1-\alpha/2}(n-(r+1)) \right]$$

prediction interval
of probability $1-\alpha$
 $P(Y_0 \in PI_{1-\alpha}(Y_0)) = 1-\alpha$

(Remember: they're one at the time!)

$$\begin{matrix} \text{GLS} & : & \hat{\beta} = \underset{\text{Generalized least squares}}{\arg \min}_{\beta} (Y - Z\beta)^T W^{-1} (Y - Z\beta) \end{matrix}$$

Collinearity

$$\hat{\beta} = (Z^T Z)^{-1} Z Y \quad : \quad \text{if } (Z^T Z) \text{ is close to be singular we have a problem due to collinearity (one regressor can be expressed as linear combination of the others)}$$

For every column:

$$\varepsilon_j := Y_0 + \delta_1 Z_1 + \dots + \delta_{j-1} Z_{j-1} + \delta_j Z_j + \dots + \delta_r Z_r$$

We substitute the j th column of Z with ε_j and we perform the linear model analysis again, getting R_j^2 .

$$\Rightarrow VIF_j = \frac{1}{1-R_j^2} \quad \begin{matrix} \text{If } VIF(\beta_j) > 5/10 \rightarrow \text{probably } \beta_j \text{ can be expressed as a linear combination of the other regressors} \\ \text{Variance Inflation Factor} \end{matrix}$$

Collinearity and variable selection

Let's work with the problem "centered":

$$\mathbb{Z} \rightarrow \mathbb{Z}^* = \begin{bmatrix} z_{11} - \bar{z}_1 & \dots & z_{1r} - \bar{z}_r \\ z_{21} - \bar{z}_2 & \dots & z_{2r} - \bar{z}_r \\ \vdots & & \vdots \\ z_{n1} - \bar{z}_1 & \dots & z_{nr} - \bar{z}_r \end{bmatrix}$$

$$Y \rightarrow Y^* = [Y_1 - \bar{Y}, \dots, Y_n - \bar{Y}]^T$$

PCA Regression

- $\mathbb{Z} = [\mathbb{Z}_1, \dots, \mathbb{Z}_r]$
- PCA on $\mathbb{Z} \Rightarrow PC_1, \dots, PC_r$
- Reduce dimensionality: $PC_1, \dots, PC_k \quad k \leq r$
- Fit $\hat{Y} = \mathbb{Z}^* \beta + \varepsilon$

$$\begin{aligned} \hat{Y}_0 &= \mathbb{Z}_0^* \hat{\beta} \\ &= \hat{\beta}_1 PC_1 + \dots + \hat{\beta}_k PC_k \\ &= \hat{\beta}_1 (e_{11} \hat{\beta}_1 + \dots + e_{1k} \hat{\beta}_k) + \dots + \hat{\beta}_k (e_{1k} \hat{\beta}_1 + \dots + e_{1k} \hat{\beta}_k) \\ &\vdots \\ &= \hat{\beta}_1 (e_{11} \hat{\beta}_1 + \dots + e_{1k} \hat{\beta}_k) + \dots + \hat{\beta}_r (e_{1r} \hat{\beta}_1 + \dots + e_{1r} \hat{\beta}_k) \\ &\vdots \\ &= \hat{\beta}_0 \mathbb{Z}_0 + \dots + \hat{\beta}_r \mathbb{Z}_r \end{aligned}$$

RIDGE Regression

Problem:

$$\begin{cases} \text{argmin}_{\beta} \| \mathbb{Z} (\beta - \hat{\beta}) \|^2 \\ \|\beta\|^2 \leq S \end{cases}$$

$$\hat{\beta} = (\mathbb{Z}^T \mathbb{Z})^{-1} \mathbb{Z}^T Y$$

$$\hat{\beta} = (\mathbb{Z}^T \mathbb{Z} + \lambda I)^{-1} \mathbb{Z}^T Y$$

(not state invariant!)

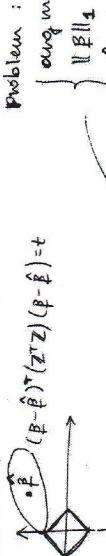
$$\begin{cases} \text{var}(z_j - \bar{z}_j) = \mathbb{E}[(z_j - \bar{z}_j)^2] \\ \text{cov}(z_i, z_j) = \mathbb{E}[(z_i - \bar{z}_i)(z_j - \bar{z}_j)] \end{cases}$$

$$\text{Def. VARIOGRAM: } \gamma(z_j - z_i) := \text{var}(z_j - \bar{z}_j) - \text{var}(z_i - \bar{z}_i)$$

$$\begin{cases} \text{symmetry: } c(-h) = c(h) \quad \forall h \in \mathbb{R}^d \\ \text{boundness: } |c(h)| \leq C(h) \quad \forall h \in \mathbb{R}^d \\ \text{pos. def.: } \sum_{i,j} \lambda_i \lambda_j c(z_i - z_j) \geq 0 \quad \forall \lambda_i, \lambda_j \in \mathbb{R} \end{cases}$$

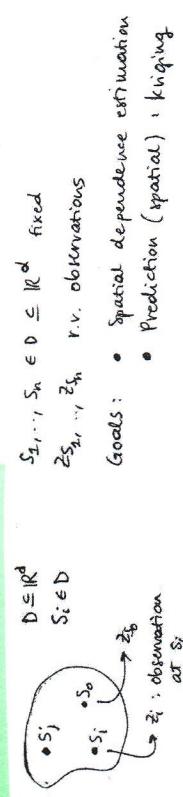
$$\text{Def. COVARIOTRAM: } \gamma(z_j - z_i) := \frac{1}{n} \sum_{s_i, s_j} (z_i - \bar{z}_i)(z_j - \bar{z}_j)$$

LASSO Regression



(with cross validation and prediction error for γ)

Spatial Data



$s_1, \dots, s_n \in D \subseteq \mathbb{R}^d$ fixed
 z_{s_1}, \dots, z_{s_n} r.v. observations

Goals:
• Spatial dependence estimation
• Prediction (spatial): knowing

$s_1, \dots, s_n \in D$ sites

z_{s_1}, \dots, z_{s_n} observations : $\{z_s, s \in D\}$

Assumptions:

- $E[z_s] < \infty \quad \forall s \in D$
- $\text{var}(z_s) < \infty \quad \forall s \in D$

Definitions:

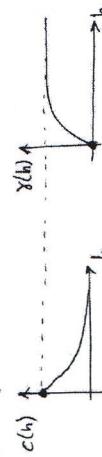
- $m_s := E[z_s]$ spatial mean of $\{z_s, s \in D\}$
- $C(s_1, s_2) = \text{cov}(z_{s_1}, z_{s_2})$ covariance function
- $\{z_s, s \in D\}$ second order stationary if:
 - $E[z_s] = m \quad \forall s \in D$
 - $\text{cov}(z_{s_1}, z_{s_2}) = \frac{C(s_1 - s_2)}{\text{COVARIOTRAM}} \quad s_1, s_2 \in D$

Properties:

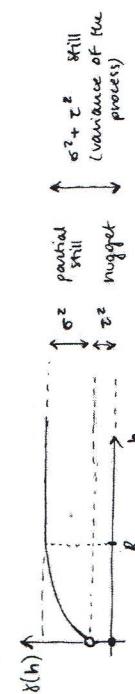
- symmetry: $c(-h) = c(h) \quad \forall h \in \mathbb{R}^d$
- boundness: $|c(h)| \leq C(h) \quad \forall h \in \mathbb{R}^d$
- pos. def.: $\sum_{i,j} \lambda_i \lambda_j c(z_i - z_j) \geq 0 \quad \forall \lambda_i, \lambda_j \in \mathbb{R}, \forall z_i, z_j \in D$

- VARIOGRAM: $\gamma(z_j - z_i) := \text{var}(z_j - \bar{z}_j) - \text{var}(z_i - \bar{z}_i) = \mathbb{E}[(z_j - \bar{z}_j)^2] - \mathbb{E}[(z_i - \bar{z}_i)^2]$
- COVARIOTRAM: $\gamma(z_j - z_i) := C(0) - C(z_j - z_i)$
- Properties:
 - symmetry: $\delta(-h) = \delta(h)$
 - null at 0: $\delta(0) = 0$
 - conditional negative definite: $\sum_{i,j} \lambda_i \lambda_j \gamma(z_i - z_j) \leq 0 \quad \forall \lambda_i, \lambda_j \in \mathbb{R}$

Examples:



Structural properties:

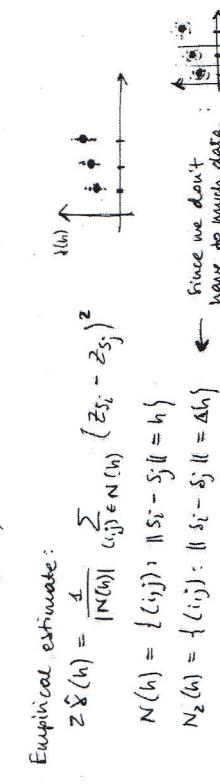


$\delta(h) = \begin{cases} \sigma^2 & \text{partial still} \\ \tau^2 & \text{nugget} \\ R & \text{range} \end{cases}$

(value of h s.t. the variogram reaches the still) if quantity the amount of dependence that we have in our data

Def. (isotropy): A field is isotropic if: $\text{cov}(z_{s_i}, z_{s_j}) = c(\|s_i - s_j\|)$
(under 2nd order of stationarity)

Estimate Spatial Dependence
(= estimate the variogram)



• Empirical estimate:

$$z\hat{\delta}(h) = \frac{1}{|\mathcal{N}(h)|} \sum_{(i,j) \in \mathcal{N}(h)} (z_{s_i} - z_{s_j})^2$$

$$\mathcal{N}(h) = \{(i,j) : \|s_i - s_j\| = h\}$$

$$\mathcal{N}_2(h) = \{(i,j) : \|s_i - s_j\| = \Delta h\}$$

Since we don't have so much data

• Parametric model:

$\delta(h; \theta)$: we have to estimate θ (best parameter vector to fit the empirical estimate)

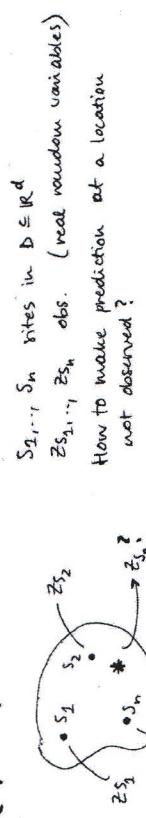


For example: $\hat{\theta} = \arg \min_{\theta \in \Theta} \sum_{k=1}^K (\delta(h_k) - \delta(h_k; \theta))^2$ (OLS)

Properties:

- δ_1, δ_2 valid variograms $\rightarrow \chi_1 + \chi_2$ valid variogram
- $\alpha > 0, \chi_1$ valid variogram $\rightarrow \alpha \chi_1$ valid variogram

(Spatial) Predictor



s_1, \dots, s_n sites in $D \subseteq \mathbb{R}^d$
 z_{s_1}, \dots, z_{s_n} obs. (real random variables)

How to make prediction at a location not observed?

Problem:

$$\begin{aligned} z_{s_0}^* &= f(z_{s_1}, \dots, z_{s_n}) \\ f &= \arg \min \mathbb{E}[(z_s - f(z_1, \dots, z_n))^2] \end{aligned}$$

Solution:

$$f(z_{s_1}, \dots, z_{s_n}) = \mathbb{E}[z_{s_0} | z_{s_1}, \dots, z_{s_n}]$$

problem: if data are not gaussian f can be highly non-linear.
we need an alternative way:

KRIGING : it looks for the best linear unbiased predictor

Problem (kriging):

$$z_{s_0}^* = \lambda_0^* + \sum_i \lambda_i^* z_{s_i}$$

where $\lambda_0^*, \dots, \lambda_n^*$ value,

$$\left\{ \begin{array}{l} \min_{\lambda_0, \lambda_1, \dots, \lambda_n \in \mathbb{R}} \mathbb{E}[(z_{s_0} - (\lambda_0 + \sum_i \lambda_i z_{s_i}))^2] \\ \text{such that: } \mathbb{E}[\lambda_0 + \sum_i \lambda_i z_{s_i}] = \mathbb{E}[z_{s_0}] \end{array} \right.$$

• Ordinary kriging (stationarity 2nd order)

Assumption:

- unknown mean $\mathbb{E}[z_s] := m \quad \forall s \in D$
- stationary 2nd order
- $\text{cov}(z_{s_1}, z_{s_2}) = C(s_1 - s_2) \quad \forall s_1, s_2 \in D$ known (we actually estimate it: $\hat{C} \rightarrow \hat{C} \rightarrow \hat{\Sigma}$)

Solution:

$$\text{solve: } \left[\begin{array}{c} \Sigma \frac{1}{n} \\ \vdots \\ 1 \end{array} \right] \left[\begin{array}{c} \lambda_0 \\ \vdots \\ \lambda_n \end{array} \right] = \left[\begin{array}{c} 0 \\ \vdots \\ 1 \end{array} \right] \Rightarrow \lambda^* \Rightarrow z_{s_0}^* = \lambda^* z_{s_0}$$

where:

$$\Sigma = \text{cov}(\bar{z}) : \quad \Sigma_{ij} = \text{cov}(z_{s_i}, z_{s_j}) = C(s_i - s_j)$$

$$\sigma_{0i} = \text{cov}(z_{s_0}, z_{s_i}) = C(s_0 - s_i)$$

μ = lagrangian multiplier

Note: we can also get the variance (of ordinary kriging)

$$\sigma_{\text{OK}}^2(s_0) = C(0) - \sum_i \lambda_i^* C(s_0 - s_i) - \mu^* = \mathbb{E}[(z_{s_0} - z_{s_0}^*)^2]$$

Universal Kriging (: unknown mean)

$$\text{Model : } \begin{cases} \hat{z}_s = m_s + \delta_s & \text{S.E.D} \\ m_s = \mathbb{E}[z_s] & \text{drift} \\ \delta_s = z_s - m_s & \text{residual} \end{cases}$$

Assumptions :

- $m_s = \sum_j a_j f_j(s) : \bullet a_j$ unknown coeff.
- $f_j(s)$ known regressors

* δ_s random field ($\{\delta_s, s \in \Omega\}$)

2nd order stationarity random field for the residuals

- $\mathbb{E}[\delta_s] = 0$
- $\text{Cov}(\delta_{s_1}, \delta_{s_2}) = \text{C}(s_1 - s_2)$

Solution:

$$\text{solve : } \begin{bmatrix} \Sigma & \mathbf{F}^T \\ \mathbf{F} & 0 \end{bmatrix} \begin{bmatrix} \lambda \\ \mu \end{bmatrix} = \begin{bmatrix} \mathbf{f}_0 \\ \mu \end{bmatrix} \implies \lambda^* \implies \mathbf{z}_{s_0}^* = \lambda^{*T} \mathbf{z}$$

where :

$$\Sigma = \text{Cov}(\mathbf{z})$$

\mathbf{F} = design matrix

$\mathbf{0}$ = matrix of zeros

$$\mathbf{f}_0 = \text{Cov}(\mathbf{z}_{s_0}, \mathbf{z}_{s_0})$$

$$\mathbf{f}_0 = \text{design vector} : f_0 = f_j(s_0) : \text{value of the } j\text{-th regressor at the new location}$$

Estimations:

If Σ is unknown we should estimate $\hat{\Sigma} \rightarrow \hat{\mathbf{C}} \rightarrow \hat{\Sigma}$

We know how to estimate $\hat{\mathbf{C}}$ from stationary data, here we have to use the residuals (since the data are not stationary)

$$\begin{aligned} \delta_{s_1}, \dots, \delta_{s_n} \text{ given} &\implies \hat{\delta} \rightarrow \hat{\mathbf{C}} \rightarrow \hat{\Sigma} \\ \delta_{s_1}, \dots, \delta_{s_n} \text{ unknown} &\implies \hat{\delta} \rightarrow \hat{\mathbf{C}} \rightarrow \hat{\Sigma} \\ \delta_{s_i} ? &: \begin{cases} \delta_{s_i} = z_{s_i} - \hat{m}_{s_i} \\ \hat{m}_{s_i} = \sum_j \hat{a}_j f_j(s_i) \end{cases} \quad \text{where } \hat{\mathbf{a}} \text{ is estimated with an iterative procedure} \end{aligned}$$

Note: Here we can't get a real true underestimate kriging variance (since it's based on Σ and we only have $\hat{\Sigma}$)