

EXPLORATING A MULTIVARIATE DATASET

Prediction problem

$\underline{x} \in \mathbb{R}^p$, $y \in \mathbb{R}$: we want to use \underline{x} to predict y .

What is the best function $f(\underline{x})$: $f: \mathbb{R}^p \rightarrow \mathbb{R}$ to predict y ?

"best": $f(\underline{x}) = \arg \min \mathbb{E}[(y - f(\underline{x}))^2] := \arg \min \text{MSE}$
(mean square error)

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

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

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

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

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

$$\hat{f}: \text{estimation of } f \text{ through } \underline{X} : \quad \underline{X} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1p} & y_1 \\ x_{21} & x_{22} & \dots & x_{2p} & y_2 \\ \vdots & & & & \\ x_{n1} & x_{n2} & \dots & x_{np} & y_n \end{bmatrix}$$

We want to see how good \hat{f} is in estimating new predictions:

$\underline{x}_0 \in \mathbb{R}^p$: new observation ($\notin \underline{X}$) for which we want to predict y_0

We want to know about the error:

$$\mathbb{E}_{\underline{X}}[(y_0 - \hat{f}(\underline{x}_0))^2 | \underline{X} = \underline{x}] := \underbrace{\mathbb{E}_{\underline{X}}[(y_0 - \hat{f}(\underline{x}_0))^2]}_{\text{we have one realization}}$$

$$\text{since } y_0 = f(\underline{x}_0) + \varepsilon_0$$

$$\implies \underbrace{\mathbb{E}_{\underline{X}}[(y_0 - f(\underline{x}_0))^2]}_{\text{this estimate the error of the prediction given } \underline{X} = \underline{x}, \text{ what if we want to consider all the possible realizations of } \underline{X}} = (f(\underline{x}_0) - \hat{f}(\underline{x}_0))^2 + \underbrace{\text{Var}(\varepsilon_0)}_{\text{this is irreducible: it's the part of } y \text{ not explainable through } \underline{x}}$$

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

$$\implies \mathbb{E}[\mathbb{E}_{\underline{X}}[(y_0 - \hat{f}(\underline{x}_0))^2]] = \text{Var}(\varepsilon_0) + \text{Var}(\hat{f}(\underline{x}_0)) + \underbrace{(f(\underline{x}_0) - \mathbb{E}[\hat{f}(\underline{x}_0)])^2}_{:= \text{BIAS}^2}$$

(how far is the model from what we want to estimate)

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

\implies LOCAL AVERAGE (fail with p large: CURSE OF DIMENSIONALITY)

\implies 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:

$$\mathbb{X} = [x_1 \ x_2 \ \dots \ x_p]$$

Every column x_j has n realizations : $y_j = [x_{1j}, x_{2j}, \dots, x_{nj}]^\top$

- mean of x_j : $\bar{x}_j := \frac{1}{n} \sum_{i=1}^n x_{ij}$

- variance of x_j : $s_{jj} := \frac{1}{n} \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2$

- covariance of x_j, x_k : $\text{Cov}(x_j, x_k) := \frac{1}{n} \sum_{i=1}^n (x_{ij} - \bar{x}_j)(x_{ik} - \bar{x}_k) := s_{kj}$

- covariance matrix (for x_1, \dots, x_p) :

$$S := \begin{bmatrix} s_{11} & s_{12} & \dots & s_{1p} \\ s_{21} & s_{22} & \dots & s_{2p} \\ \vdots & & & \\ s_{p1} & s_{p2} & \dots & s_{pp} \end{bmatrix} \in \mathbb{R}^{p \times p}$$

- 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) :

$$r := \begin{bmatrix} r_{11} & \dots & r_{1p} \\ \vdots & & \\ r_{p1} & \dots & r_{pp} \end{bmatrix}$$

Whatever the distribution of x_j is, we can say (Chebyshev) :

$$\Pr(\bar{x}_j - k\sqrt{s_{jj}} \leq x_j \leq \bar{x}_j + k\sqrt{s_{jj}}) \geq 1 - \frac{1}{k^2} \quad \forall k \neq 0$$

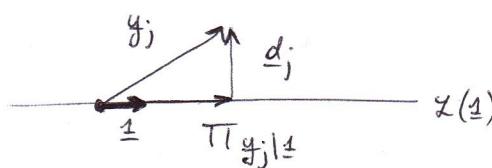
alternatively:

$$\Pr(|x_j - \mu| \leq k\sigma) \geq 1 - \frac{1}{k^2}$$

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

Geometrical interpretation:

- one variable



d_j = deviation vector

$\mathcal{L}(1)$ = Space with "no statistic" ($v \in \mathcal{L}(1) \Rightarrow v = k \cdot 1$)

$\Pi_{y_j | 1}$ = projection of y_j on $\mathcal{L}(1)$

Generic projection of v on $\mathcal{L}(w)$:

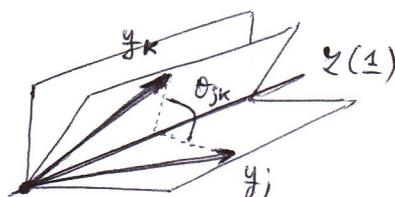
$$\Pi_{v | \mathcal{L}(w)} = \frac{w \cdot w^\top}{w^\top \cdot w} \cdot v$$

$$\Rightarrow \Pi_{y_j | \mathcal{L}(1)} = \frac{1 \cdot 1^\top}{1^\top \cdot 1} y_j = \bar{x}_j \cdot 1$$

$$\Rightarrow d_j = y_j - \bar{x}_j \cdot 1 \quad \Rightarrow \|d_j\| = \sqrt{n} \sqrt{s_{jj}}$$

this is why s_{jj} talks about the variability
(Note: d_j captures the informations that the mean doesn't)

- two variables



$$\begin{aligned} y_j &= \bar{x}_j \cdot 1 + d_j \\ y_k &= \bar{x}_k \cdot 1 + d_k \end{aligned}$$

$$\cos(\theta_{jk}) = \frac{\underline{d}_j^T \underline{d}_k}{\|\underline{d}_j\| \cdot \|\underline{d}_k\|} = \frac{s_{jk}}{\sqrt{s_{jj} s_{kk}}} = r_{jk}$$

- $\theta_{jk} = 0 \Rightarrow r_{jk} = 1 \Rightarrow \underline{d}_j \in \mathcal{L}(\underline{d}_k)$
- $\theta_{jk} = \frac{\pi}{2} \Rightarrow r_{jk} = 0 \Rightarrow \underline{d}_k \perp \underline{d}_j$

- by rows:

$$\underline{X} = \begin{bmatrix} \underline{x}_1^T \\ \vdots \\ \underline{x}_n^T \end{bmatrix} \quad \text{where } \underline{x}_j = \text{realization of } \underline{X_j}$$

$\underline{x}_1, \dots, \underline{x}_n \stackrel{iid}{\sim} \underline{X} \in \mathbb{R}^p$: we merge with random vectors

$$\underline{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_p \end{bmatrix}$$

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

- 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})

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

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

- correlation matrix (for \underline{X}): (ρ)

$$V = \text{diag}(\sigma_1, \dots, \sigma_p) \Rightarrow V^{1/2} = \text{diag}(\sqrt{\sigma_1}, \dots, \sqrt{\sigma_p})$$

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

$$\rho := V^{-1/2} \Sigma V^{1/2}$$

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

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

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

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

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

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

Estimators

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

n rows $\Rightarrow n$ random vectors $\underline{x}_j \in \mathbb{R}^p$ ($\underline{x}_1, \dots, \underline{x}_n \stackrel{iid}{\sim} \underline{X} \in \mathbb{R}^p$)

- Estimator for $\underline{\mu}$:

$$\bar{\underline{X}} = \frac{1}{n} \sum_{j=1}^n \underline{x}_j \xrightarrow{\text{realization}} \bar{\underline{x}} = \begin{bmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_p \end{bmatrix}$$

Assuming $\underline{x}_1, \dots, \underline{x}_n \stackrel{iid}{\sim} \underline{X}$, $\mathbb{E}[\underline{X}] = \underline{\mu}$, $\text{Cov}(\underline{X}) = \Sigma$:

- $\mathbb{E}[\bar{\underline{X}}] = \underline{\mu}$ (unbiased)

- $\text{Cov}(\bar{\underline{X}}) = \frac{1}{n} \Sigma$

- Estimator for Σ :

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

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

from now on:

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

Considering that:

$$\underline{d}_j = \underline{y}_j - \frac{\underline{1} \cdot \underline{1}^T}{\underline{1}^T \underline{1}} \underline{y}_j = (\mathbf{I} - \frac{\underline{1} \cdot \underline{1}^T}{\underline{1}^T \underline{1}}) \underline{y}_j \quad (\underline{d}_j = \underline{y}_j - \pi_{\underline{y}_j} \underline{1})$$

$$\mathbf{d} = [\underline{d}_1, \dots, \underline{d}_p] = (\mathbf{I} - \frac{\underline{1} \cdot \underline{1}^T}{\underline{1}^T \underline{1}}) \mathbf{X} \quad := \text{deviation matrix}$$

$$\begin{aligned} S &= \frac{1}{n} \mathbf{d}^T \mathbf{d} \\ &= \frac{1}{n-1} \mathbf{X}^T (\mathbf{I} - \frac{\underline{1} \cdot \underline{1}^T}{\underline{1}^T \underline{1}}) \mathbf{X} \end{aligned}$$

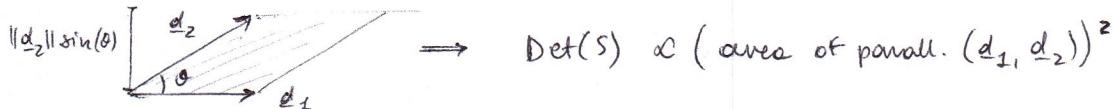
Variability in a multivariate sense

In \mathbb{R} the variability is explained by variance.

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

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

$$p=2: S = \frac{1}{n-1} \mathbf{d}^T \mathbf{d} \Rightarrow \text{Det}(S) = \|\underline{d}_1\|^2 \|\underline{d}_2\|^2 \sin^2 \theta \left(\frac{1}{n-1}\right)^2$$



- $\text{det}(S) \uparrow \Rightarrow \text{area} \uparrow$: can increase because of θ (max area with $\theta = \pi/2$) or because of $\|\underline{d}_1\| / \|\underline{d}_2\|$
- $\text{det}(S) = 0 \Rightarrow$ either:
 - $(\|\underline{d}_1\| \vee \|\underline{d}_2\|) = 0$
 - $\theta = 0$

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

$$p \neq 2: \text{Det}(S) \propto (\text{vol}(\text{parallelotope } (\underline{d}_1, \dots, \underline{d}_p))^2$$

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

$$p=2: \text{Tr}(S) = \frac{1}{n-1} (\|\underline{d}_1\|^2 + \|\underline{d}_2\|^2)$$

$$p \neq 2: \text{Tr}(S) = \frac{1}{n-1} (\|\underline{d}_1\|^2 + \dots + \|\underline{d}_p\|^2)$$

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

Note: $\text{Det}(S) = 0 \iff \underline{d}_1, \dots, \underline{d}_p$ are linearly dependent
 $(\exists c \neq 0 \text{ s.t. } c_1 \underline{d}_1 + \dots + c_p \underline{d}_p = 0)$

Consequence: $\text{Det}(S) = 0 \Rightarrow \exists k \text{ s.t. } \underline{d}_k = - \sum_{\substack{i=1 \\ i \neq k}}^p \frac{c_i}{c_k} \underline{d}_i$

$$\Rightarrow \underline{y}_k = \bar{\underline{x}}_k \underline{1} - \sum_{\substack{i=1 \\ i \neq k}}^p \frac{c_i}{c_k} (\underline{y}_i - \bar{\underline{x}}_i \underline{1})$$

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

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 $\det(S) > 0$ is pos. def. ($\lambda_i > 0 \forall i$)

Assuming S positive definite ($\det(S) > 0$):

(and an order $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p > 0$)

$$\underline{x}, \underline{y} \in \mathbb{R}^P : d_{S^{-1}}^2(\underline{x}, \underline{y}) = (\underline{x} - \underline{y})^\top S^{-1} (\underline{x} - \underline{y}) := \text{Mahalanobis distance (squared)}$$

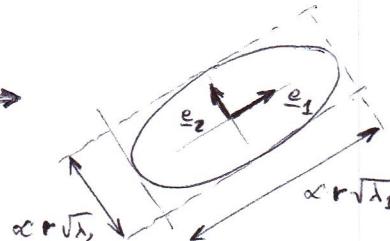
Note: $S^{-1} = \sum_{i=1}^P \frac{1}{\lambda_i} e_i e_i^\top$

"Neighborhood" in Mahalanobis' distance:

$$E_{r^2, S^{-1}}(\underline{x}) = \{ \underline{x} \in \mathbb{R}^P : (\underline{x} - \underline{x})^\top S^{-1} (\underline{x} - \underline{x}) \leq r^2 \}$$

Graphically ($P=2$):

$$S^{-1} = \sum_{i=1}^2 \frac{1}{\lambda_i} e_i e_i^\top \Rightarrow$$



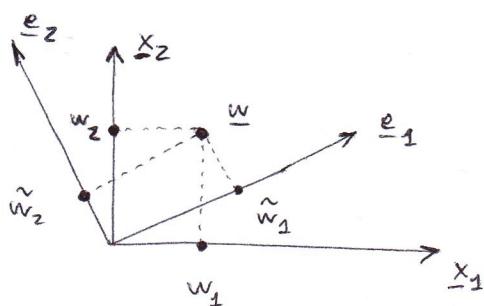
The Mahalanobis' distance is the right one for capturing the distance between data.
(The euclidean distance is not good, in fact the Mahalanobis' distance is the Euclidean distance after we standardized the data)

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 // to the axes of the system

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



| old system | new system |
|--|--|
| $\underline{w} = [x_1 \dots x_p]^\top$ | $\tilde{\underline{w}} = [e_1^\top \underline{w} \dots e_p^\top \underline{w}] = P^\top \underline{w}$ |
| $\underline{X} = \begin{bmatrix} \underline{x}_1^\top \\ \vdots \\ \underline{x}_n^\top \end{bmatrix}$ | $\tilde{\underline{X}} = \underline{X} P$ |
| $S = \frac{1}{n-1} \underline{X} \left(\underline{I} - \frac{1}{n-1} \underline{1} \underline{1}^\top \right) \underline{X}^\top$ | $\tilde{S} = P^\top S P = \Lambda$ |

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

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

There is always a reference system for which the coordinates are uncorrelated: correlation is NOT a property of the data

Do we lose information? No:

PRINCIPAL COMPONENT ANALYSIS (PCA)

Let's refer to the model (not the data) :

- \underline{X} random vector in \mathbb{R}^P :
- $\mathbb{E}[\underline{X}] = \underline{\mu}$
- $\text{Cov}(\underline{X}) = \Sigma$

Problem : Find $\underline{a} \in \mathbb{R}^P$ s.t. $\|\underline{a}\|=1$ and $\text{Var}(\underline{a}^T \underline{X})$ is maximized

(i.e. find the direction \underline{a} s.t. the variability of the projection of \underline{X} on \underline{a} is maximum)

$$\max_{\underline{a} \in \mathbb{R}^P: \|\underline{a}\|=1} \text{Var}(\underline{a}^T \underline{X}) = \max_{\underline{a} \in \mathbb{R}^P: \|\underline{a}\|=1} \underline{a}^T \Sigma \underline{a} = \max_{\underline{a} \in \mathbb{R}^P} \frac{\underline{a}^T \Sigma \underline{a}}{\underline{a}^T \underline{a}}$$

Geometrical lemma:

$$B \in \mathbb{R}^{P \times P} \text{ positive semidefinite}, \quad B = \sum_{i=1}^P \lambda_i \underline{e}_i \underline{e}_i^T : \quad (\lambda_1 \geq \lambda_2 \geq \dots)$$

$$\Rightarrow \left\{ \begin{array}{l} 1. \max_{\underline{X} \in \mathbb{R}^P} \frac{\underline{X}^T B \underline{X}}{\underline{X}^T \underline{X}} = \lambda_1, \quad \text{arg max}(\dots) = \underline{e}_1 \\ 2. \max_{\substack{\underline{X} \in \mathbb{R}^P \\ \underline{X} \perp \underline{e}_1}} \frac{\underline{X}^T B \underline{X}}{\underline{X}^T \underline{X}} = \lambda_2, \quad \text{arg max}(\dots) = \underline{e}_2 \\ \vdots \\ p. \max_{\substack{\underline{X} \in \mathbb{R}^P \\ \underline{X} \perp \underline{e}_1, \dots, \underline{X} \perp \underline{e}_{p-1}}} \frac{\underline{X}^T B \underline{X}}{\underline{X}^T \underline{X}} = \lambda_p = \min_{\substack{\underline{X} \in \mathbb{R}^P \\ \underline{X} \perp \underline{e}_1, \dots, \underline{X} \perp \underline{e}_{p-1}}} \frac{\underline{X}^T B \underline{X}}{\underline{X}^T \underline{X}} \end{array} \right.$$

Back to PCA :

$$\max_{\underline{a} \in \mathbb{R}^P} \text{Var}(\underline{a}^T \underline{X}) = \max_{\underline{a} \in \mathbb{R}^P} \frac{\underline{a}^T \Sigma \underline{a}}{\underline{a}^T \underline{a}} = \lambda_1, \quad \text{arg max}(\dots) = \underline{e}_1$$

\Rightarrow First Principal Component (PC1) : $Y_1 = \underline{e}_1^T \underline{X}$ ($/ Y_1 = \underline{e}_1^T (\underline{X} - \underline{\mu})$)

Problem : Find $\underline{a} \in \mathbb{R}^P$ s.t. $\|\underline{a}\|=1$, $\text{Var}(\underline{a}^T \underline{X})$ is max and $\underline{a} \perp \underline{e}_1$

(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)

$$\max_{\substack{\underline{a} \in \mathbb{R}^P: \|\underline{a}\|=1 \\ \text{cov}(\underline{a}^T \underline{X}, \underline{e}_1^T \underline{X})=0}} \text{Var}(\underline{a}^T \underline{X}) = \max_{\substack{\underline{a} \in \mathbb{R}^P \\ \underline{a} \perp \underline{e}_1}} \frac{\underline{a}^T \Sigma \underline{a}}{\underline{a}^T \underline{a}}$$

$$\Rightarrow \max_{\substack{\underline{a} \in \mathbb{R}^P: \\ \underline{a} \perp \underline{e}_1}} \frac{\underline{a}^T \Sigma \underline{a}}{\underline{a}^T \underline{a}} = \lambda_2, \quad \text{arg max}(\dots) = \underline{e}_2$$

\Rightarrow Second Principal Component (PC2) : $Y_2 = \underline{e}_2^T \underline{X}$ ($/ Y_2 = \underline{e}_2^T (\underline{X} - \underline{\mu})$)

Generalized problem :

$$\max_{\substack{\underline{a} \in \mathbb{R}^P: \|\underline{a}\|=1 \\ \text{cov}(\underline{a}^T \underline{X}, \underline{e}_j^T \underline{X})=0 \quad j=1, \dots, k-1}} \text{Var}(\underline{a}^T \underline{X}) = \lambda_k, \quad \text{arg max}(\dots) = \underline{e}_k$$

\Rightarrow k^{th} Principal Component (PCk) : $Y_k = \underline{e}_k^T \underline{X}$ ($/ Y_k = \underline{e}_k^T (\underline{X} - \underline{\mu})$)

$$\underline{Y} := [Y_1 \dots Y_p]^\top = P^\top \underline{X} \quad := \text{vector of Principal Components}$$

- $E[\underline{Y}] = E[P^\top \underline{X}] = P^\top \underline{\mu}$
- $\text{Cov}(\underline{Y}) = \text{Cov}(P^\top \underline{X}) = P^\top \Sigma P = \Lambda$
 $\text{Cov}(Y_i, Y_j) = 0 \quad \forall i \neq j \quad \longrightarrow \quad \text{no correlation between the coordinates of } \underline{Y}$
 $\text{Cov}(Y_i, Y_i) = \text{Var}(Y_i) = \lambda_i \quad \forall i$
- Ordering: we ordered the components \Rightarrow the first component is the one with larger variability, the second is the one with the second larger var. and so on \Rightarrow we can capture the most of the variability just with the first components (the last express small variability)
- We're not losing variability:
 - generalized variance with \underline{Y} : $\text{Det}(\Lambda) = \prod_{i=1}^p \lambda_i = \text{Det}(\Sigma)$
 - total variance with \underline{Y} : $\text{Tr}(\Lambda) = \sum_{i=1}^p \lambda_i = \text{Tr}(\Sigma)$
- $Y_i = e_{1i} X_1 + \dots + e_{pi} X_p$
 $\{e_{1i}, \dots, e_{pi}\}$ are the loading (weights):
 $\text{corr}(Y_i, X_k) = \frac{e_{ki} \sqrt{\lambda_i}}{\sqrt{\sigma_{kk}}}$

PCA on standardized variables

(we're still working with the model, not data)

$$V := \text{diag}(\sigma_{11}, \dots, \sigma_{pp})$$

$$\underline{X} \longrightarrow \underline{Z} = V^{-1/2}(\underline{X} - \underline{\mu}) = \left[\frac{X_1 - \mu_1}{\sqrt{\sigma_{11}}}, \dots, \frac{X_p - \mu_p}{\sqrt{\sigma_{pp}}} \right]^\top$$

$$\underline{Z} : \bullet E[\underline{Z}] = 0$$

$$\bullet \text{Cov}(\underline{Z}) = V^{-1/2} \Sigma V^{-1/2} := \rho \quad := \text{Correlation matrix}$$

$$\rho = \sum_{i=1}^p \lambda_i e_i e_i^\top \quad (\text{eigenvalues/vectors of } \rho)$$

(covariance matrix of the standardized variables)

$$\implies Y_i = e_i^\top \underline{Z} = e_i^\top V^{-1/2}(\underline{X} - \underline{\mu})$$

$$(\text{Note: } \text{Tr}(\rho) = \sum_{i=1}^p \text{Var}(Y_i) = \sum_{i=1}^p \text{Var}(Z_i) = p)$$

$$\bullet \text{PCA}(\Sigma) \neq \text{PCA}(\rho)$$

- 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 PCA(ρ):

if X_1 is in km and X_2 in mm \Rightarrow $\text{Var}(X_1)$ will mask $\text{Var}(X_2) \Rightarrow \text{PCA}(\rho)$

PCA on data

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

Usually $\underline{\mu}$ and Σ are unknown, we use data:

$$\left. \begin{array}{l} S \text{ estimates } \Sigma \\ \underline{X} \text{ estimates } \underline{\mu} \end{array} \right\} \implies \text{PCA on } S$$

$$S = \sum_{i=1}^p \lambda_i e_i e_i^T \implies \text{PCA: } y_i \text{ projection on } e_i$$

$$\underset{i\text{-th stat. unit}}{\underset{x_i}{\xrightarrow{\text{PCA}}}} \quad y_i = \begin{bmatrix} e_1^T x_i \\ \vdots \\ e_p^T x_i \end{bmatrix} \quad \begin{array}{l} \leftarrow \text{"score" of } x_i \text{ on the 1st PC} \\ \leftarrow \text{"score" of } x_i \text{ on the } p\text{th PC} \end{array}$$

$$\bar{X} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \xrightarrow{\text{PCA}} \begin{bmatrix} y_1^T \\ \vdots \\ y_n^T \end{bmatrix} = \begin{bmatrix} y_{11} & y_{12} & \cdots & y_{1p} \\ y_{21} & y_{22} & \vdots & y_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ y_{n1} & y_{n2} & \cdots & y_{np} \end{bmatrix} \quad \begin{array}{l} \rightarrow \text{scores on the} \\ p\text{-th principal component} \end{array}$$

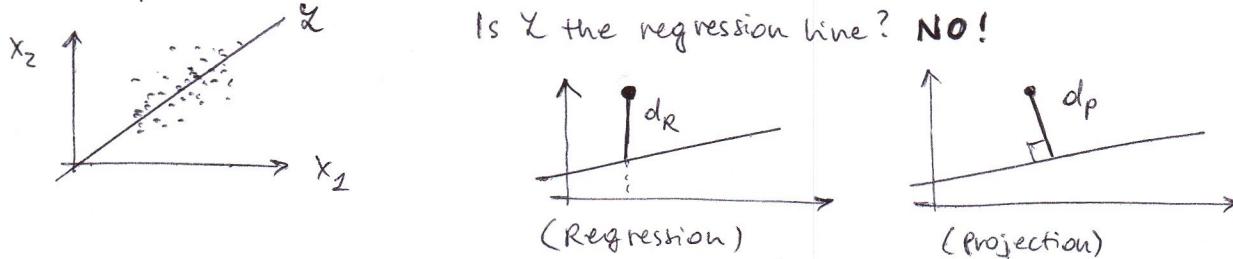
$y_1 \quad y_2 \quad \cdots \quad y_p$

- 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 x_1, \dots, x_k

Geometrical meaning of PCA

Problem: given $x_1, \dots, x_n \in \mathbb{R}^p$ find a linear space \mathcal{L} s.t. $\dim(\mathcal{L}) = k \leq p$ which best approximate the data (i.e. such that \mathcal{L} is "closest" to data).

Consider $p=2$:



Here we're trying to minimize d_P , not d_R ($\neq d_P$)!

(Note that the regression depends on the system, the projection no)

Problem:

Find y_1, \dots, y_k (we identify \mathcal{L} with an orthonormal basis:
 $\mathcal{L} = \text{span}(y_1, \dots, y_k)$)

such that:

$$\underbrace{\sum_{i=1}^n \| (x_i - \bar{x}) - \sum_{j=1}^k y_j y_j^T (x_i - \bar{x}) \|^2}_{(*)} \text{ is minimum.}$$

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

$$\min (*) \iff \max \sum_{i=1}^n \sum_{j=1}^k (y_j^T (x_i - \bar{x}))^2$$

$$\iff \max (n-1) \sum_{j=1}^k y_j^T S y_j : \quad k=1 \implies \max (\cdot) = \lambda_1 \quad \text{arg max} (\cdot) = \frac{e_1}{\|e_1\|}$$

by induction on k :

$$y_1 = e_1, \dots, y_k = e_k \implies \text{The directions of the first } k \text{ principal components identify the closest linear space (of dim } k) \text{ to the data}$$

$$\text{Error of approximation: } (n-1) \sum_{j=k+1}^p \lambda_j$$

MULTIVARIATE GAUSSIAN DISTRIBUTION

$$\underline{X} \sim N_p(\mu, \Sigma)$$

$$f(\underline{x}) = \frac{1}{\sqrt{(2\pi)^p \det(\Sigma)}} e^{-\frac{1}{2} (\underline{x}-\mu)^\top \Sigma^{-1} (\underline{x}-\mu)}$$

- $\underline{X} \sim N_p(\mu, \Sigma) \implies E[\underline{X}] = \mu, \text{Cov}(\underline{X}) = \Sigma$
- $\underline{X} \sim N_p(\mu, \Sigma) \iff \underline{a}^\top \underline{X} \sim N_1(\underline{a}^\top \mu, \underline{a}^\top \Sigma \underline{a}) \quad \forall \underline{a} \in \mathbb{R}^p$
(A linear combination of the components is gaussian)
- $\underline{X} = [X_1, \dots, X_p]^\top \sim N_p(\mu, \Sigma), \Sigma = [\sigma_{ij}] \implies X_j \sim N_1(\mu_j, \sigma_{jj})$
- $\underline{X} \sim N_p(\mu, \Sigma), A \in \mathbb{R}^{q \times p} \implies A\underline{X} \sim N_q(A\mu, A\Sigma A^\top)$
- $\underline{X} \sim N_p(\mu, \Sigma), \underline{d} \in \mathbb{R}^p \implies \underline{X} + \underline{d} \sim N_p(\mu + \underline{d}, \Sigma)$
- $\underline{X} = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} \sim N_p\left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}\right) \implies X_1 \sim N_1(\mu_1, \Sigma_{11})$

So for example:

$$\begin{bmatrix} X_1 \\ X_2 \\ X_3 \end{bmatrix} \sim N \implies \begin{bmatrix} X_1 \\ X_3 \end{bmatrix} \sim N$$

$$X_1 \perp\!\!\!\perp X_2 \iff \Sigma_{12} = \Sigma_{21} = 0$$

- $\underline{X} = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} \sim N_p\left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}\right), \det(\Sigma) \neq 0 :$
 $\implies X_1 | X_2 = x_2 \sim N_1(\mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (x_2 - \mu_2), \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21})$
 $\text{Cov}(X_1 | X_2 = x_2) = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \perp\!\!\!\perp X_2$

\therefore Partial covariances

(even without knowing X_2 we already know how information will be modified)

(This generates the Regression Effect)

- $\underline{X} \sim N_p(\mu, \Sigma), \det(\Sigma) > 0 \implies (\underline{X}-\mu)^\top \Sigma^{-1} (\underline{X}-\mu) \sim \chi^2(p)$
- $\underline{X} \sim N_p(\mu, \Sigma), \det(\Sigma) = 0 \implies (\underline{X}-\mu)^\top \Sigma^{-1} (\underline{X}-\mu) \sim \chi^2(k)$
 $k = \text{rank}(\Sigma)$
 $\Sigma^{-1} = \sum_{i=1}^k \frac{1}{\lambda_i} e_i e_i^\top$
 $(\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k > 0 = \lambda_{k+1} = \dots = \lambda_p)$
- $\alpha \in (0, 1), \det(\Sigma) > 0 \implies P((\underline{X}-\mu)^\top \Sigma^{-1} (\underline{X}-\mu) \leq \chi^2_\alpha(p)) = 1-\alpha$



Estimators for μ and Σ

$$\bar{\underline{X}} = \begin{bmatrix} \bar{X}_1 \\ \bar{X}_2 \\ \vdots \\ \bar{X}_n \end{bmatrix} : \underline{x}_i \text{ realization of } X_i \implies \underline{X}_1, \dots, \underline{X}_n \stackrel{iid}{\sim} N_p(\mu, \Sigma)$$

$$\bullet \bar{\underline{X}} = \frac{1}{n} \sum_{i=1}^n \underline{X}_i$$

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

$$\bullet \hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (\underline{X}_i - \bar{\underline{X}})(\underline{X}_i - \bar{\underline{X}})^\top \quad (\text{biased, but MLE})$$

Distributions: assuming $\underline{X}_1, \dots, \underline{X}_n \stackrel{\text{iid}}{\sim} N_p(\mu, \Sigma)$

- $\bar{X} \sim N_p(\mu, \frac{1}{n} \Sigma)$
- $S \sim \text{Wish}(\frac{1}{n-1} \Sigma, n-1) \quad (n-1)S \sim \text{Wish}(\Sigma, n-1)$
- $\hat{\Sigma} \sim \text{Wish}(\frac{1}{n} \Sigma, n-1) \quad n\hat{\Sigma} \sim \text{Wish}(\Sigma, n-1)$

About Wishart distribution (matrix distribution!):

(Def.) $\underline{Z}_1, \dots, \underline{Z}_m$ iid $\sim N_p(\underline{0}, \Sigma)$, $\det(\Sigma) > 0$

$$\Rightarrow \sum_{i=1}^m \underline{Z}_i \underline{Z}_i^\top \sim \text{Wishart}(\Sigma, m) \quad (\text{the param p is in } \Sigma)$$

Properties:

- $A_1 \sim \text{Wish}(\Sigma, m_1)$, $A_2 \sim \text{Wish}(\Sigma, m_2)$, $A_1 \perp\!\!\!\perp A_2$
 $\Rightarrow A_1 + A_2 \sim \text{Wish}(\Sigma, m_1 + m_2)$
- $A \sim \text{Wish}(\Sigma, m)$, $C \in \mathbb{R}^{k \times p}$ const
 $\Rightarrow C A C^\top \sim \text{Wish}(C \Sigma C^\top, m)$
- $A \sim \text{Wish}(\Sigma, m)$, $\sigma^2 > 0$
 $\Rightarrow \sigma^2 A \sim \text{Wish}(\sigma^2 \Sigma, m)$
- $A \sim \text{Wish}(\Sigma, m)$, $\Sigma 1 \times 1 \Rightarrow A \sim \Sigma \chi^2(m)$

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

LLN

$\underline{X}_1, \dots, \underline{X}_n$ random vectors iid such that $E[\underline{X}_i] = \mu$, $\text{Cov}(\underline{X}) = \Sigma$ exists:

- $\bar{X} = \frac{1}{n} \sum_{i=1}^n \underline{X}_i \xrightarrow{\text{P}} \mu \text{ as } n \rightarrow \infty$
- $S = \frac{1}{n-1} \sum_{i=1}^n (\underline{X}_i - \bar{X})(\underline{X}_i - \bar{X})^\top \xrightarrow{\text{P}} \Sigma \text{ as } n \rightarrow \infty$

CLT

$\underline{X}_1, \dots, \underline{X}_n$ random vectors iid such that $E[\underline{X}_i] = \mu$, $\text{Cov}(\underline{X}) = \Sigma$ exists:

$$\Rightarrow \sqrt{n}(\bar{X} - \mu) \sim \underbrace{A N_p(\underline{0}, \Sigma)}$$

asymptotically normal
meaning: for large n one can approximate the distribution of

$$\sqrt{n}(\bar{X} - \mu) \text{ with a } N_p(\underline{0}, \Sigma)$$

In practice, for large n :

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

It doesn't mean that if the sample is large then it is gaussian. It means that if a sample is large enough then the sample mean is gaussian

INFERENCE FOR THE MEAN μ

- $X_1, \dots, X_n \in \mathbb{R}^p$ random vectors s.t.
- $(X_1, \dots, X_n) \stackrel{\text{iid}}{\sim} N_p(\mu, \Sigma)$
 - $E[\bar{X}] = \mu$
 - $\text{Cov}(\bar{X}) = \frac{1}{n} \Sigma$
 - $\det(\Sigma) > 0$

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \quad \text{pointwise estimator for } \mu$$

Two cases:

- n large ($n > p$):

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

Pivotal statistic distr.:

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

If Σ is not known:

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

Pivotal statistic distr.:

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

- Confidence regions:

$$\text{IP}(\text{d}_{S^{-1}}(\bar{X}, \mu) \leq \chi^2_\alpha(p)) = \text{IP}(\mu \in \mathcal{E}_{S^{-1}}^\alpha(\bar{X})) = 1 - \alpha$$

(i.e. the random ellipse centered in \bar{X} and shaped with S will be covering the true (fixed) value of μ , $(1-\alpha)$ -100% of the times that we generate it)

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

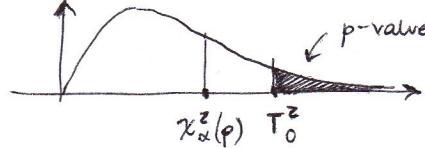
- Testing:

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

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

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

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



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

- n small

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

Pivotal statistic distr.:

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

Hotelling's theorem: $\bar{X} \sim N_p(\mu, \Sigma)$ ($\det(\Sigma) > 0$) $\Leftrightarrow W \sim \text{Wish}(\frac{1}{m}\Sigma, m)$

$$\Leftrightarrow \frac{m-p+1}{m} (\bar{X} - \mu)^T W^{-1} (\bar{X} - \mu) \sim F(p, m-p+1)$$

- Confidence regions:

$$\text{IP} \left(n(\bar{\boldsymbol{x}} - \boldsymbol{\mu})^T S^{-1} (\bar{\boldsymbol{x}} - \boldsymbol{\mu}) \leq \frac{(n-1)p}{n-p} F_{\alpha}(p, n-p) \right) = 1 - \alpha$$

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

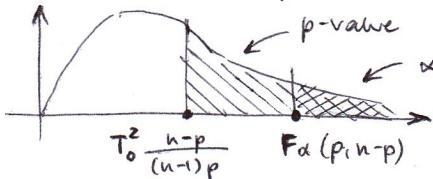
- Testing:

$$\begin{cases} H_0: \boldsymbol{\mu} = \boldsymbol{\mu}_0 \\ H_1: \boldsymbol{\mu} \neq \boldsymbol{\mu}_0 \end{cases}$$

$$T_0^2 = n(\bar{\boldsymbol{x}} - \boldsymbol{\mu}_0)^T S^{-1} (\bar{\boldsymbol{x}} - \boldsymbol{\mu}_0)$$

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

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



CI for linear combinations of the mean

$\underline{\boldsymbol{x}}_1, \dots, \underline{\boldsymbol{x}}_n \stackrel{iid}{\sim} N_p(\boldsymbol{\mu}, \Sigma)$, small n case

$\bar{\boldsymbol{x}} = \frac{1}{n} \sum_{i=1}^n \underline{\boldsymbol{x}}_i$ estimator for $\boldsymbol{\mu}$

- One-at-the-time CI($\boldsymbol{\mu}$) ($1! \alpha$)

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

$$\frac{\underline{a}^T \bar{\boldsymbol{x}} - \underline{a}^T \boldsymbol{\mu}}{\sqrt{\underline{a}^T S \underline{a}}} \sqrt{n} \sim t(n-1)$$

$$\Rightarrow \text{IP} \left(\underline{a}^T \boldsymbol{\mu} \in [\underline{a}^T \bar{\boldsymbol{x}} \pm t_{\frac{\alpha}{2}}(n-1) \sqrt{\frac{\underline{a}^T S \underline{a}}{n}}] \right) = 1 - \alpha \quad \forall \underline{a} \in \mathbb{R}^p$$

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

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

Reject if:

$$\frac{|\underline{a}^T \bar{\boldsymbol{x}} - \delta_0|}{\sqrt{\underline{a}^T S \underline{a}}} \sqrt{n} > t_{\frac{\alpha}{2}}(n-1)$$

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

Reject if:

$$\frac{\underline{a}^T \bar{\boldsymbol{x}} - \delta_0}{\sqrt{\underline{a}^T S \underline{a}}} \sqrt{n} > t_{\alpha}(n-1)$$

Examples:

$$\underline{a} = [0 \dots 0 \ 1 \ 0 \ \dots \ 0]^T \Rightarrow \underline{a}^T \boldsymbol{\mu} = \mu_i$$

$$\underline{a} = [0 \dots 0 \ 1 \ 0 \ \dots \ 0 \ 1 \ 0 \ \dots \ 0]^T \Rightarrow \underline{a}^T \boldsymbol{\mu} = \mu_i - \mu_j$$

- Simultaneous CI($\underline{a}^T \boldsymbol{\mu}$) ($\infty \alpha$)

$$\text{IP} \left(\frac{|\underline{a}^T (\bar{\boldsymbol{x}} - \boldsymbol{\mu})|}{\sqrt{\underline{a}^T S \underline{a}}} \sqrt{n} \leq \sqrt{\frac{(n-1)p}{n-p} F_{\alpha}(p, n-p)}, \forall \underline{a} \in \mathbb{R}^p \right) = 1 - \alpha$$

$$\text{Sim CI}_{1-\alpha}(\underline{a}^T \boldsymbol{\mu}) = \left[\underline{a}^T \bar{\boldsymbol{x}} \pm \sqrt{\frac{(n-1)p}{n-p} F_{\alpha}(p, n-p)} \sqrt{\frac{\underline{a}^T S \underline{a}}{n}} \right]$$

- Bonferroni's correction ($\underline{a} \in \mathbb{R}^p$, small k)

Given $\underline{a}_1, \dots, \underline{a}_k \in \mathbb{R}^p$ find $CI(\underline{a}_1^\top \mu), \dots, CI(\underline{a}_k^\top \mu)$ with simultaneous confidence of $1-\alpha$, $\alpha \in (0,1)$:

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

Testing k assumptions simultaneously

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

Given $\underline{a}_1, \dots, \underline{a}_k \in \mathbb{R}^p$:

$$\begin{cases} H_0 : \begin{cases} \underline{a}_1^\top \mu = \delta_1 \\ \vdots \\ \underline{a}_k^\top \mu = \delta_k \end{cases} \\ H_1 : \exists i : \underline{a}_i^\top \mu \neq \delta_i \end{cases}$$

Bonferroni method:

Reject at level α if for at least one i : $\frac{|\underline{a}_i^\top \bar{x} - \delta_i|}{\sqrt{\underline{a}_i^\top S \underline{a}_i}} \sqrt{n} > t_{\frac{\alpha}{2k}}(n-1)$

Truth:

| | not-rejecting H_0 | rejecting H_0 |
|-------|-----------------------------|-----------------------------|
| H_0 | V | $V_{\text{false positive}}$ |
| H_1 | $T_{\text{false negative}}$ | S |

$K = R$
not rejected R
rejected

$K_0 = \# \text{ true hypothesis}$

$K - K_0 = \# \text{ false hypothesis}$

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

- False discovery rate : $FDR := \mathbb{E}\left[\frac{V}{R}\right]$

Property:

$$FDR \leq FWER$$

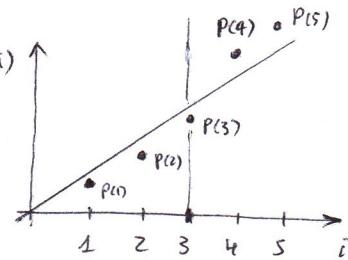
Methods:

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

1. For each of the k tests compute the p-value p_i
2. Order the p-values: $p_{(1)} \leq p_{(2)} \leq \dots \leq p_{(k)}$
3. $m := \max\{i \in \{1, \dots, k\} : p_{(i)} \leq \frac{i}{k} \alpha\}$
4. Reject $H_{0(1)}, \dots, H_{0(m)}$

Ex :

$$m = 3 :$$



Remember:

$$H_{0_1} \neq H_{0_{(1)}}$$

In this way, if the p-values are LL $\Rightarrow FDR \leq \alpha$

Comparing means of Gaussian distributions

- Paired data

Each unit is observed 2 times:

$$\underline{X}_{1i} = \begin{bmatrix} X_{1i1} \\ X_{1i2} \\ \vdots \\ X_{1ip} \end{bmatrix} \quad \underline{X}_{2i} = \begin{bmatrix} X_{2i1} \\ X_{2i2} \\ \vdots \\ X_{2ip} \end{bmatrix} \quad i = 1, \dots, n$$

\underline{X}_{1i} iid with mean μ_1

\underline{X}_{2i} iid with mean μ_2

Goal: inference on $\mu_1 - \mu_2$

$$\Rightarrow D_i := \underline{X}_{1i} - \underline{X}_{2i} \quad i = 1, \dots, n, \quad E[D_i] = \mu_1 - \mu_2$$

Goal: inference on $E[D_i]$

Assumptions: $D_1, \dots, D_n \stackrel{iid}{\sim} N_p(\underline{\delta}, \Sigma_D)$ (we're not assuming on distr. of $\underline{X}_{1i}/\underline{X}_{2i}$)

Pivotal statistic:

$$n(\bar{D} - \underline{\delta})^T S_D^{-1}(\bar{D} - \underline{\delta}) \sim \frac{(n-1)p}{n-p} F(p, n-p)$$

- CR $1-\alpha(\mu_1 - \mu_2) = \{\underline{\delta} \in \mathbb{R}^p; n(\bar{D} - \underline{\delta})^T S_D^{-1}(\bar{D} - \underline{\delta}) \leq \frac{(n-1)p}{n-p} F_\alpha(p, n-p)\}$
- $H_0: \mu_1 - \mu_2 = \underline{\delta}_0$ vs. $H_1: \mu_1 - \mu_2 \neq \underline{\delta}_0$
Reject at level α if: $n(\bar{D} - \underline{\delta}_0)^T S_D^{-1}(\bar{D} - \underline{\delta}_0) > \frac{(n-1)p}{n-p} F_\alpha(p, n-p)$
- $\text{Sim CI}_{1-\alpha}(\mu_{1i} - \mu_{2i}) = [\bar{D}_i \pm \sqrt{\frac{(n-1)p}{n-p} F_\alpha(p, n-p)} \sqrt{\frac{1}{n}} \sqrt{S_{D,ii}}]$
- Bonf Sim CI $1-\alpha = [\bar{D}_i \pm t_{\frac{\alpha}{2p}}(n-1) \sqrt{\frac{1}{n}} \sqrt{S_{D,ii}}]$

- Repeated univariate measures

Each unit is observed q times

$$\underline{X}_i = [X_{i1}, \dots, X_{iq}]^T \quad i = 1, \dots, n$$

Ex. n patients,

X_{ij} = blood pressure of patient i at time j

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

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

We introduce the contrast matrix C :

$$\text{examples: } C = \begin{bmatrix} 1 & -1 & 0 & \dots & 0 \\ 0 & 1 & -1 & \dots & 0 \end{bmatrix}, \quad C = \begin{bmatrix} 1 & 0 & \dots & -1 \\ 0 & 1 & \dots & 0 \end{bmatrix} \quad (*)$$

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

Pivotal statistic:

$$n(C\bar{X} - C\underline{\mu})^T (CSC)^{-1}(C\bar{X} - C\underline{\mu}) \sim \frac{(n-1)(q-1)}{n-q+1} F(q-1, n-q+1)$$

Reject at level α if:

$$n(C\bar{X} - C\underline{\mu})^T (CSC)^{-1}(C\bar{X} - C\underline{\mu}) > F_\alpha(q-1, n-q+1)$$

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

the second is testing $\mu_1 - \mu_q, \mu_2 - \mu_q, \dots$

M(ANOVA) ONE-WAY

Cases :

| | P # features | g # groups | Description |
|---|-------------------|-----------------|---|
| 1 | ≥ 1 | 2 | We have n patients and one treatment. We apply the treatment to n_2 patients and not to n_1 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 pat.s we apply level 1 : for n_g pat.s we apply level $g-1$. For each patients we have 1 feature |
| 3 | ≥ 1 | ≥ 2 | Like in the second case but for each patient we have p features. |

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

$$\begin{aligned} \underline{x}_{11}, \dots, \underline{x}_{1n_1} &\stackrel{\text{iid}}{\sim} N_p(\mu_1, \Sigma) \\ \underline{x}_{21}, \dots, \underline{x}_{2n_2} &\stackrel{\text{iid}}{\sim} N_p(\mu_2, \Sigma) \end{aligned} \quad \left. \begin{array}{l} \\ \parallel \end{array} \right.$$

Goal: inference on $\mu_1 - \mu_2$

$$\begin{aligned} \mu_1 &\leftarrow \bar{\underline{x}}_1 = \frac{1}{n_1} \sum_{i=1}^{n_1} \underline{x}_{1i} \sim N_p(\mu_1, \frac{1}{n_1} \Sigma) \\ \mu_2 &\leftarrow \bar{\underline{x}}_2 = \frac{1}{n_2} \sum_{i=1}^{n_2} \underline{x}_{2i} \sim N_p(\mu_2, \frac{1}{n_2} \Sigma) \end{aligned} \quad \left. \begin{array}{l} \\ \parallel \end{array} \right.$$

$$S_{\text{pooled}} := \frac{(n_1-1)S_1 + (n_2-1)S_2}{n_1 + n_2 - 2}$$

with $\left\{ \begin{array}{l} S_1 = \frac{1}{n_1-1} \sum_{i=1}^{n_1} (\underline{x}_{1i} - \bar{\underline{x}}_1)(\underline{x}_{1i} - \bar{\underline{x}}_1)^T \\ S_2 = \frac{1}{n_2-1} \sum_{i=1}^{n_2} (\underline{x}_{2i} - \bar{\underline{x}}_2)(\underline{x}_{2i} - \bar{\underline{x}}_2)^T \end{array} \right. \quad \begin{array}{l} (\Sigma \text{ in group 1}) \\ (\Sigma \text{ in group 2}) \end{array}$

Pivotal statistic:

$$\left(\frac{1}{n_1} + \frac{1}{n_2} \right)^{-1} \left[(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2) \right]^T S_{\text{pooled}}^{-1} \left[(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2) \right] \sim \frac{(n_1 + n_2 - 2)p}{n_1 + n_2 - 1 - p} F(p, n_1 + n_2 - 1 - p)$$

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:

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

Confidence region:

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

Note: for n_1 and n_2 very large:

$$\left[(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2) \right]^T \left(\frac{S_1}{n_1} + \frac{S_2}{n_2} \right)^{-1} \left[(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2) \right] \sim \chi^2(p)$$

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

$$\left. \begin{array}{l} X_{11}, \dots, X_{1n_1} \stackrel{\text{iid}}{\sim} N_1(\mu_1, \sigma^2) \\ \vdots \\ X_{g1}, \dots, X_{gn_g} \stackrel{\text{iid}}{\sim} N_1(\mu_g, \sigma^2) \end{array} \right\} \quad n = n_1 + \dots + n_g$$

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

Parametrization: $\mu_i = \underline{\mu} + \tau_i$ $i = 1, \dots, g$
overall mean

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

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

$$\underline{\mu} \leftarrow \bar{X} = \frac{1}{n} \sum_{i=1}^g \sum_{j=1}^{n_i} X_{ij}$$

$$\tau_i \leftarrow \bar{X}_i - \bar{X} = \left[\frac{1}{n_i} \sum_{j=1}^{n_i} X_{ij} - \bar{X} \right]$$

Variance decomposition:

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

Pivotal statistic:

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

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

Case 3: $p \geq 1, g \geq 2$

$$\begin{aligned} X_{11}, \dots, X_{1n_1} &\stackrel{iid}{\sim} N_p(\mu_1, \Sigma) \\ \vdots \\ X_{g1}, \dots, X_{gn_g} &\stackrel{iid}{\sim} N_p(\mu_g, \Sigma) \end{aligned} \quad \left. \right\} \text{II}$$

Modello:

$$X_{ij} = \mu + \tau_i + \varepsilon_{ij} \quad \begin{array}{l} i=1, \dots, g \quad j=1, \dots, n_i \\ \sum_{i=1}^g n_i \tau_i = 0 \end{array} \\ \varepsilon_{ij} \sim N(0, \Sigma)$$

$$\mu \leftarrow \bar{x}$$

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

$$\text{Goal: } \begin{cases} H_0: \tau_1 = \tau_2 = \dots = \tau_g = 0 \\ H_1: \exists \tau_i \neq 0 \end{cases}$$

Covariance decomposition:

$$\underbrace{\sum_{i=1}^g \sum_{j=1}^{n_i} (X_{ij} - \bar{x})(X_{ij} - \bar{x})^\top}_{\text{covariability that we would compute under } H_0} = \underbrace{\sum_{i=1}^g (\bar{x}_i - \bar{x})(\bar{x}_i - \bar{x})^\top n_i}_{\therefore = B \text{ between}} + \underbrace{\sum_{i=1}^g \sum_{j=1}^{n_i} (X_{ij} - \bar{x}_i)(X_{ij} - \bar{x}_i)^\top}_{? = W \text{ within}}$$

covariability that we would compute under H_0

: = B
between

? = W
within

We want again see " $\frac{B}{W}$ " to decide whether to reject or not

Proposals:

| Λ | reject H_0 if Λ is .. |
|--|---------------------------------|
| WILKS: $\Lambda_W = \frac{\text{Det}(W)}{\text{Det}(W+B)}$ | small |
| LAWLEY-HOTELING $\Lambda_{LH} = \text{Tr}(BW^{-1})$ | large |
| PILLAI: $\Lambda_p = \text{Tr}(B \cdot (B+W)^{-1})$ | large |

(interpretations on notes)

(*) when all the n_i are big enough:

Barlett's approximation:

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

Reject at level α if:

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

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

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

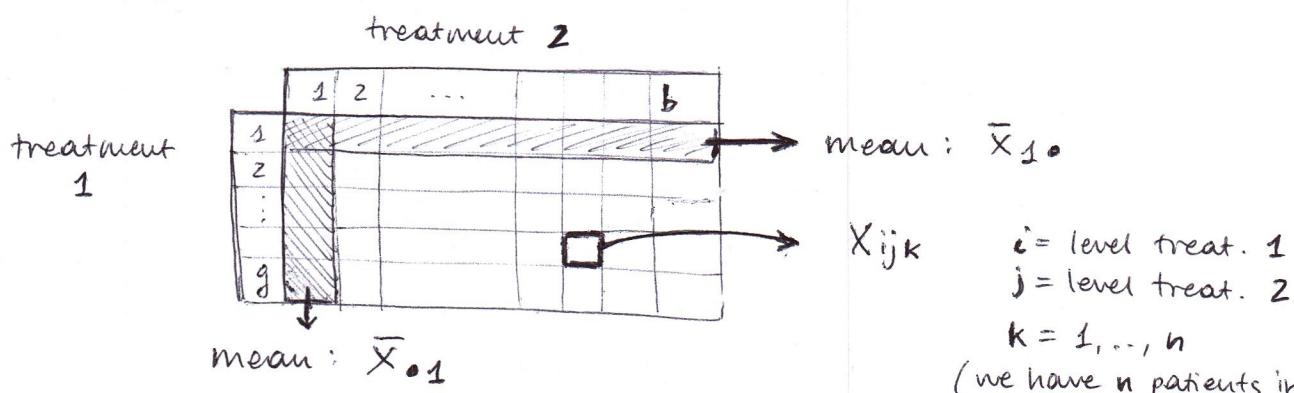
→ Bonferroni CI ($\tau_{ie} - \tau_{ke}$) $k, i = 1, \dots, g$
 $e = 1, \dots, p$

- $\bar{x}_{ie} - \bar{x}_{ke} \sim N(\tau_{ie} - \tau_{ke}, \frac{1}{n_i} \sigma_{ee} + \frac{1}{n_k} \sigma_{ee})$
- $\frac{1}{n-g} W$ estimate $\Sigma \Rightarrow \frac{Wee}{n-g}$ estimate σ_{ee}

$$\text{Bonferroni CI}_{1-\alpha}(\tau_{ie} - \tau_{ke}) = \left[\bar{x}_{ie} - \bar{x}_{ke} \pm t \frac{\alpha}{pg(g-1)} \sqrt{\frac{Wee}{n-g} \left(\frac{1}{n_k} + \frac{1}{n_i} \right)} \right]$$

(M)ANOVA TWO-WAYS

We have two treatments



$$\text{Model: } X_{ijk} = \mu + \underbrace{\tau_i}_{\substack{\text{treat.} \\ 1}} + \underbrace{\beta_j}_{\substack{\text{treat.} \\ 2}} + \underbrace{\gamma_{ij}}_{\text{interactions}} + \varepsilon_{ijk} \quad \begin{array}{l} i=1, \dots, g \\ j=1, \dots, b \\ k=1, \dots, n \end{array}$$

Constraints:

$$\sum_{i=1}^g \tau_i = 0$$

$$\sum_{j=1}^b \beta_j = 0$$

$$\sum_{i=1}^g \gamma_{ij} = \sum_{j=1}^b \gamma_{ij} = 0$$

Decomposition of variance:

$$\sum_{i=1}^g \sum_{j=1}^b \sum_{k=1}^n (x_{ijk} - \bar{x})^2 = \sum_{i=1}^g bn (\bar{x}_{i\cdot} - \bar{x})^2 + \text{SS treat 1}$$

$$+ \sum_{j=1}^b gn (\bar{x}_{\cdot j} - \bar{x})^2 + \text{SS treat 2}$$

$$+ \sum_{i=1}^g \sum_{j=1}^b (\bar{x}_{ij} - \bar{x}_{i\cdot} - \bar{x}_{\cdot j} + \bar{x})^2 n + \text{SS interactions}$$

$$+ \sum_{i=1}^g \sum_{j=1}^b \sum_{k=1}^n (x_{ijk} - \bar{x}_{ij})^2 \text{ SS residuals}$$

in this sum up
in the additive
model
(with no interact.)

$$\begin{cases} H_0 : \gamma_{ij} = 0 \\ H_1 : \exists \gamma_{ij} \neq 0 \end{cases} \rightarrow \text{Reject at level } \alpha \text{ if:}$$

$$\frac{\frac{\text{SS interactions}}{(g-1)(b-1)}}{\frac{\text{SSres}}{gb(n-1)}} > F_\alpha((g-1)(b-1), gb(n-1))$$

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

$$\begin{cases} H_0 : \tau_1 = \tau_2 = \dots = \tau_g = 0 \\ H_1 : \exists \tau_i \neq 0 \end{cases} \rightarrow \text{Reject at level } \alpha \text{ if:}$$

$$\frac{\frac{\text{SS treat 1}}{g-1}}{\frac{\text{SS residuals}}{gb(n-g)}} > F_\alpha(g-1, gb(n-g))$$

CLASSIFICATION

Supervised

$$X = \left[\begin{array}{cccc|c} x_1 & x_2 & \dots & x_p & L \\ x_{11} & x_{12} & \dots & x_{1p} & l_1 \\ \vdots & & & & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} & l_n \end{array} \right]$$

Goal: Find $\delta: X \rightarrow \{1, \dots, g\}$
 $(X = [x_1, \dots, x_p] \in X)$
 (DISCRIMINANT ANALYSIS)

Unsupervised

$$X = \left[\begin{array}{cccc} x_1 & x_2 & \dots & x_p \\ x_{11} & x_{12} & \dots & x_{1p} \\ \vdots & & & \\ x_{n1} & x_{n2} & \dots & x_{np} \end{array} \right]$$

we believe there are hidden labels

- Goal:
- estimate tables (CLUSTER ANALYSIS)
 - Find $\delta: X \rightarrow \{1, \dots, g\}$ (DISCRIMINANT ANALYSIS)

Ingredients: (SUPERVISED CLASSIFICATION)

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

Optimality criterium:

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

Goal: $\min ECM(\delta)$ (Expected Cost of Misclassification of δ)

- $g = 2$

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

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

Explanation:

we put x in R_1 if :

$$\text{cost}\left(\begin{array}{l} x \text{ belonging to 2} \\ x \text{ attributed to 1} \end{array}\right) \leq \text{cost}\left(\begin{array}{l} x \text{ belonging to 1} \\ x \text{ attributed to 2} \end{array}\right)$$

We're choosing where to put x basing on the minimum cost that we have to pay if we're wrong.

- $g > 2$

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

$$R_2 = \{x \in \mathbb{R}^p : \sum_{k \neq 2} c(2|k)f_k(x)p_k \leq \sum_{k \neq j} c(j|k)f_k(x)p_k, j=1, 3, \dots, g\}$$

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

Explanation:

we put x in R_1 if :

$$\text{cost}\left(\begin{array}{l} x \text{ belonging to } k (k \in \{2, \dots, g\}) \\ x \text{ attributed to 1} \end{array}\right) \leq \text{cost}\left(\begin{array}{l} x \text{ belonging to } j (j \neq k) \\ x \text{ attributed to 1} \end{array}\right)$$

We chose to put x in R_1 if the cost of misclassification is minimum

$$\delta(x) = i \iff x \in R_i$$

Aj

Optimal classifier:

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

($i \in \{1, \dots, g\}$)

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

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

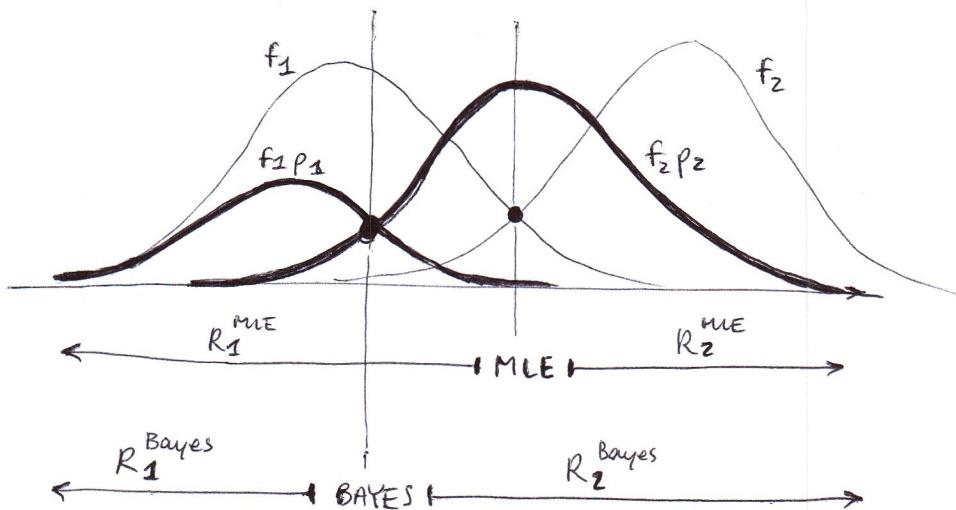
- BAYES CLASSIFIER: all costs are constant (and equal)

$$\delta(\underline{x}) = i \iff \Pr(L=j | \underline{X}=\underline{x}) \leq \Pr(L=i | \underline{X}=\underline{x}) \quad \forall j \neq i$$

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

- MLE CLASSIFIER: $c(ilj) = \text{const } \forall i \neq j$, $p_1 = p_2 = \dots = p_g = \frac{1}{g}$

$$\delta(\underline{x}) = i \iff f_j(\underline{x}) \leq f_i(\underline{x}) \quad \forall j \neq i$$



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)

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

$\pi_k := \frac{c_k p_k}{\sum c_j p_j}$ are acting like priors

BAYES
CLASSIFIER

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

Special Bayes classifiers

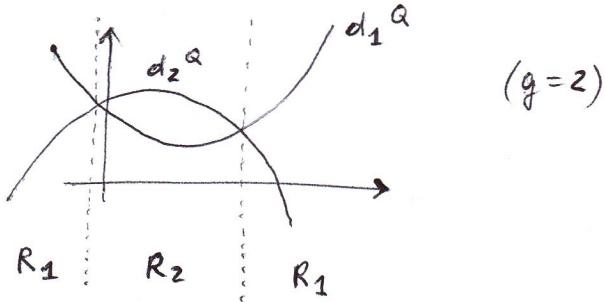
- **QDA** (Quadratic Discriminant Analysis)

$$\underline{x} | L=i \sim N_p(\mu_i, \Sigma_i)$$

→ we can substitute the generic $P(L=j | \underline{x} = \underline{x}) \leq P(L=i | \underline{x} = \underline{x})$:

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

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



- **LDA** (Linear Discriminant Analysis)

$$\underline{x} | L=i \sim N_p(\mu_i, \Sigma) \quad (\Sigma_1 = \Sigma_2 = \dots = \Sigma_g = \Sigma)$$

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

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

Fisher's argument for LDA

Suppose $\underline{x} | L=i \sim (\mu_i, \Sigma)$ (no gaussianity)

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

$$\Rightarrow \begin{cases} a_1 = \Sigma^{-1/2} e_1 \\ \vdots \\ a_s = \Sigma^{-1/2} e_s \end{cases} \quad \bar{\mu} = \frac{1}{g} \sum_{i=1}^g \mu_i$$

where e_j are the eigenvectors of $\Sigma^{-1/2} B \Sigma^{-1/2}$ ($B = \frac{1}{g-1} \sum_{i=1}^g (\mu_i - \bar{\mu})(\mu_i - \bar{\mu})^T$)

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

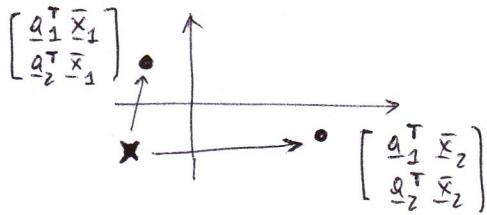
Classifier:

1. $\bar{x}_i \rightarrow [a_1^T \bar{x}_i, \dots, a_k^T \bar{x}_i, \dots, a_s^T \bar{x}_i]^T$ we consider only the first k projections

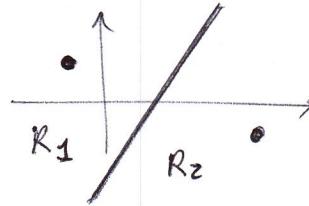
2. observation $\underline{x} \rightarrow [a_1^T \underline{x}, \dots, a_k^T \underline{x}]^T$

3. $R_i = \{ \underline{x} \in \mathbb{R}^p : \sum_{j=1}^k (a_j^T \underline{x} - a_j^T \bar{x}_i)^2 \leq \sum_{j=1}^k (a_j^T \underline{x} - a_j^T \bar{x}_h)^2, h=1, \dots, g \}$

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



\underline{x} closest to the mean of group 1 $\Rightarrow \delta(\underline{x}) = 1$



Estimate of the parameters

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

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

- LDA : $\hat{\mu}_k = \bar{\underline{x}}_k$

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

Evaluating a classifier

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

- non-parametric estimate

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

| | | attributed L | |
|------------|---|----------------|----------|
| | | 1 | 2 |
| actual L | 1 | n_{11} | n_{12} |
| | 2 | n_{21} | n_{22} |

$$\hat{AER}(\delta) = \frac{n_{12} + n_{21}}{n}$$

- leave-one-out cross validation

For $i=1, \dots, n$

1. we take out " i " from the training set
2. we train δ on $\mathcal{X}-i \Rightarrow \delta_{-i}$
3. $\delta_{-i}(\underline{x}_i) = \hat{l}_i$
4. $\varepsilon_i = \begin{cases} 1 & \text{if } \hat{l}_i \neq l_i \\ 0 & \text{if } \hat{l}_i = l_i \end{cases}$

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

(the final considered classifier is δ , not δ_{-i})

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

K-fold cross-validation algorithm:

0. Set $k < n$ and randomly split the units of the training set in k parts (randomly = permute the rows (randomly) and then split in k parts)

For $j = 1, \dots, k$

1. Take out part j from training set
2. Train δ on \mathcal{X} -part j $\Rightarrow \delta_{\text{-part } j}$
3. Apply $\delta_{\text{-part } j}$ to the part j :

$$\text{Err}_j = \frac{1}{n_j} \sum_{\{i \in \text{part } j\}} \varepsilon_i$$

$$\varepsilon_i = \begin{cases} 1 & \delta_{\text{-part } j}(x_i) \neq l_i \\ 0 & \delta_{\text{-part } j}(x_i) = l_i \end{cases}$$

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

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

$$\Rightarrow \begin{cases} \mathbb{E}[\hat{\text{AER}}(\delta)] = \frac{1}{B} \sum_{j=1}^B \hat{\text{AER}}_j(\delta) \\ \text{Var}(\hat{\text{AER}}(\delta)) = \frac{1}{B-1} \sum_{j=1}^B (\hat{\text{AER}}_j(\delta) - \mathbb{E}[\hat{\text{AER}}(\delta)])^2 \end{cases}$$

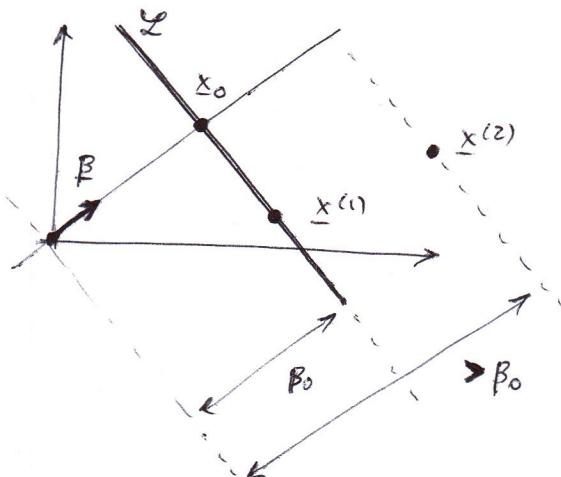
$$CI_{1-\alpha}(\mathbb{E}[\hat{\text{AER}}(\delta)]) = \left[\mathbb{E}[\hat{\text{AER}}(\delta)] \pm \sqrt{\frac{\text{Var}(\hat{\text{AER}}(\delta))}{B}} z_{\frac{\alpha}{2}} \right] \quad (\text{CLT})$$

Support Vector Machines (SVM)

If we have a set (suppose $g=2$):



• Ǝ SEPARATING HYPERPLANE



To identify the hyperplane we introduce a vector β s.t. $\beta \in \mathbb{R}^p$, $\beta \perp \mathcal{L}$, $\|\beta\| = 1$

$$x_0 \in \text{span}(\beta) \cap \mathcal{L}$$

$$\beta_0 = \|x_0\|$$

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

$$\underline{x}^{(1)} \in \mathcal{L} \iff \Pi_{\underline{x}^{(1)}}|\mathcal{L}(\beta) = \underline{x}_0$$

$$\iff \beta^T \underline{x}^{(1)} = \beta_0$$

For example, $\underline{x}^{(2)} \notin \mathcal{L} \Rightarrow \beta^T \underline{x}^{(2)} > \beta_0$.

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

let say:

$$\begin{array}{ccc} \text{plus} & \Rightarrow & \left\{ \begin{array}{l} \beta^T \underline{x} - \beta_0 > 0 \Rightarrow \underline{x} \in \text{Plus} \\ \beta^T \underline{x} - \beta_0 < 0 \Rightarrow \underline{x} \in \text{minus} \end{array} \right. \end{array}$$

Suppose we have 2 groups:

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

\rightarrow Distance between \underline{x}_i and \mathcal{L} : $y_i (\beta^T \underline{x}_i - \beta_0)$:

$$\text{Let } M := \min \{ y_i (\beta^T \underline{x}_i - \beta_0) : i=1, \dots \}$$



Optimal separating plane:

$$\mathcal{L} \text{ (to } \beta \text{ and } \beta_0 \text{) s.t. } \max_{\beta_0, \beta} M$$

under the constraints:

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

• \mathcal{L} SEPARATING HYPERPLANE

We can either:

- Allow some overlapping:

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

- Use non-linear boundaries

UNSUPERVISED CLASSIFICATION

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

Dissimilarity functions:

- $d(\underline{x}, \underline{y}) = \sqrt{(\underline{x}-\underline{y})^T (\underline{x}-\underline{y})} = \sqrt{\sum_{i=1}^p (x_i - y_i)^2}$ Euclidean
- $d_{\Sigma^{-1}}(\underline{x}, \underline{y}) = \sqrt{(\underline{x}-\underline{y})^T \Sigma^{-1} (\underline{x}-\underline{y})}$ Mahalanobis
- $d(\underline{x}, \underline{y}) = (\sum_{i=1}^p |x_i - y_i|^m)^{1/m}$ ℓ^m distance (Minkowski)

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

cambene

Dissimilarity matrix:

$$X = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \longrightarrow d_{ij} = d(x_i, x_j) \longrightarrow D = \begin{bmatrix} 0 & d_{12} & d_{13} & \dots \\ d_{21} & 0 & d_{23} & \dots \\ \vdots & \vdots & \ddots & \ddots \end{bmatrix}$$

Dissimilarity between clusters:

U, V two sets of points: $d(U, V) = ?$

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

Hierarchical Agglomerative Clustering Algorithm

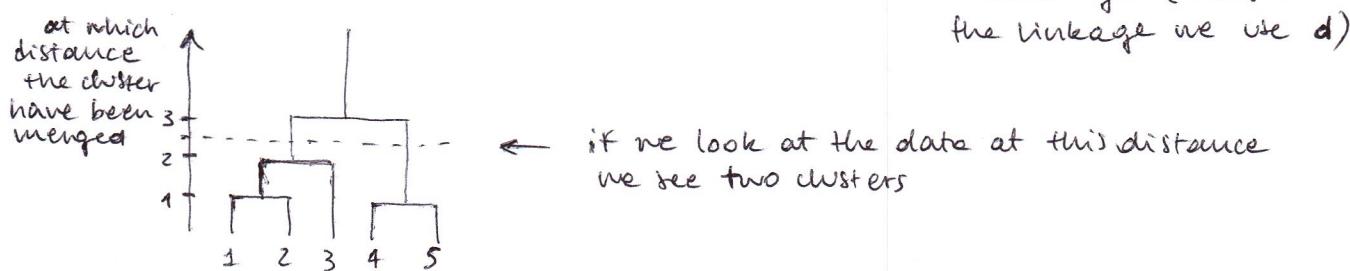
- Set a d and a linkage
- Initialization: every unit is a cluster
- Until convergence repeat:

1. merge the two cluster which are less similar
2. compute the new distance matrix



Graphical representation:

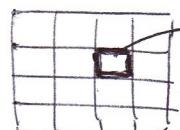
Dendrogram



Another representation of the Dendrogram:

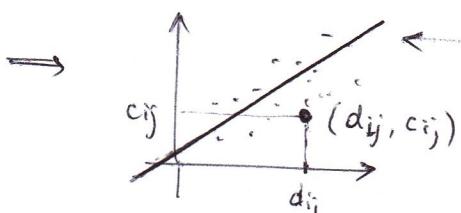
Cophenetic distance

Dendrogram



level at which i and j have been merged

For every couple of elements: $(x_i, x_j) \xrightarrow{i} (d_{ij}, c_{ij}) \quad i, j = 1, \dots, n$



Cophenetic correlation coefficient

$$|CPCC| = |\text{Corr}(C, D)|$$

$$= |\text{Corr}\{(d_{ij}, c_{ij}) : i, j = 1, \dots, n\}|$$

(the closer to 1 the better the coph. matrix is repr. the dendrogr.)

Ward's method for Hierarchical clustering

Suppose we split data into k groups.

Let C_1, \dots, C_k be clusters.

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

↑
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 .

(Method 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 :

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

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

K-means algorithm:

- Initialization:
 - randomly create C_1, \dots, C_k in the training set (\Rightarrow step 1)
 - randomly assign k centroids: $\bar{x}_1, \dots, \bar{x}_k$ (\Rightarrow step 2)
- Iterate until convergence:

1. For $j = 1, \dots, k$ compute the centroids of $C_j \Rightarrow \bar{x}_1, \dots, \bar{x}_k$

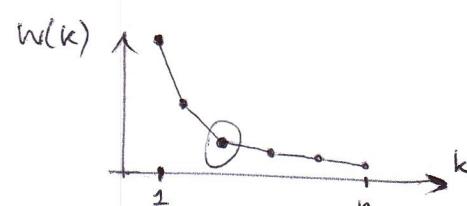
2. For each stat. unit x_i :

assign x_i to cluster C_j if:

$$d^2(x_i, \bar{x}_j) = \min \{ d^2(x_i, \bar{x}_1), \dots, d^2(x_i, \bar{x}_k) \}$$

How to choose k ?

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



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 $E[Y|X=x] = f(x)$
 (= Regression function)

Two basically approach:

1. Totally data-driven (non-parametric)
2. Parametric approach (model based)

→ CART
 classification and
 regression trees

Linear Models

$$\underline{X} \rightarrow \underline{Z} = \begin{bmatrix} z_1 & \dots & z_r \\ 1 & z_{11} & \dots & z_{1r} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & z_{n1} & \dots & z_{nr} \end{bmatrix} \in \mathbb{R}^{n \times r}$$

z_1, \dots, z_r
 known functions of X_1, \dots, X_p

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

Model for Y :

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

$$\beta = [\beta_0, \beta_1, \dots, \beta_r]^T \in \mathbb{R}^{r+1}$$

$$\varepsilon \in \mathbb{R}^n \text{ s.t. } \begin{cases} E[\varepsilon] = 0 \\ \text{Cov}(\varepsilon) = \sigma^2 I \end{cases}$$

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

- $\varepsilon_1, \dots, \varepsilon_n$:
- uncorrelated
 - same variance
 - mean 0
 - independent of the z 's

Estimating β and σ^2

(i.e. fitting the model)

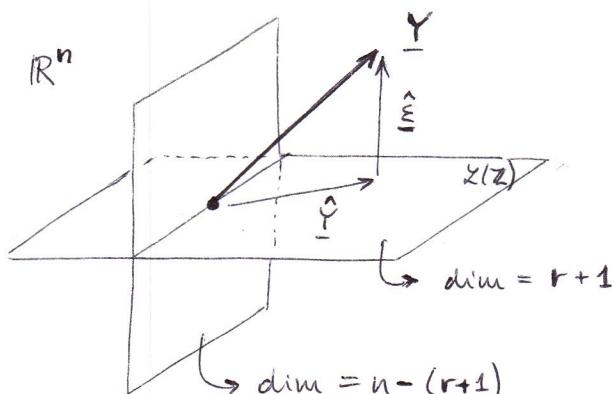
OLS : Ordinary Least Squares

$$\hat{\beta} = \arg \min_{\beta} \|Y - Z\beta\|^2$$

$$Z\hat{\beta} = \hat{Y} \quad \text{fitted values}$$

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

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)

$$\hat{Y} = Z(Z^T Z)^{-1} Z^T Y := HY \quad H = P_{\text{col}(Z)}$$

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

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

- rank(Z) = $k < r+1 < n$ ($\dim(Z(Z)) = k$)

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

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

$$(Z^T Z)^{-1} = \sum_{i=1}^k \frac{1}{\lambda_i} e_i e_i^T$$

where

$$\left\{ \begin{array}{l} Z^T Z = \sum_{i=1}^{r+1} \lambda_i e_i e_i^T \\ \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k = 0 = \lambda_{k+1} = \dots = \lambda_{r+1} \end{array} \right.$$

- rank(Z) = $r+1 = n$ ($Z(Z) = \mathbb{R}^n$)

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

Coefficients of determination

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

- Decomposition of variance:

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

CSS_{TOT}
(centered
sum of squares)

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

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

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

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

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

From now on: $\Sigma \sim N_n(0, \sigma^2 I)$:

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

Inference:

$$\left\{ \begin{array}{l} \frac{1}{S^2} (\hat{\beta} - \beta)^T (\mathbf{Z}^T \mathbf{Z}) (\hat{\beta} - \beta) \sim F_{r+1, n-(r+1)} \\ (S^2 = \frac{\hat{\Sigma}^T \hat{\Sigma}}{n - (r+1)}) \end{array} \right.$$

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

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

Simultaneous CI for $a^T \beta$:

$$\text{Sim CI}_{1-\alpha}(a^T \beta) = \left[a^T \hat{\beta} \pm \sqrt{a^T (\mathbf{Z}^T \mathbf{Z})^{-1} a} \sqrt{S^2(r+1) F_{1-\alpha}(r+1, n-(r+1))} \right]$$

Simultaneous CI_{1-α} for
any linear combination of β's

Special case

$$\text{Sim CI}_{1-\alpha}(\beta_i) = \left[\hat{\beta}_i \pm \sqrt{\text{diag}_i(\mathbf{Z}^T \mathbf{Z})^{-1}} \sqrt{S^2(r+1) F_{\alpha}(r+1, n-(r+1))} \right]$$

Testing the β 's

$$\begin{cases} H_0: C\hat{\beta} = \underline{0} \\ H_1: C\hat{\beta} \neq \underline{0} \end{cases} \quad : \quad \frac{1}{s^2} (\underline{C}\hat{\beta})^\top (\underline{C}(\underline{Z}^\top \underline{Z})^{-1} \underline{C}^\top)^{-1} (\underline{C}\hat{\beta}) \sim p F(p, n-(r+1))$$

Reject H_0 at level α if: $\frac{1}{s^2} (\underline{C}\hat{\beta})^\top (\underline{C}(\underline{Z}^\top \underline{Z})^{-1} \underline{C}^\top)^{-1} (\underline{C}\hat{\beta}) > p F_{1-\alpha}(p, n-(r+1))$

Special case I:

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

We're comparing $\underline{Y} = \underline{Z}\underline{\beta} + \underline{\varepsilon}$ vs. $\underline{Y} = \underline{Z}_1\underline{\beta}_1 + \underline{\varepsilon}_1$

$$\rightarrow \frac{SS_{\text{res}}(\underline{Z}_1) - SS_{\text{res}}(\underline{Z})}{s^2_p} \sim F(p, n-(r+1)) \quad \left(s^2 = \frac{\underline{\varepsilon}^\top \underline{\varepsilon}}{n-(r+1)} \right)$$

$$\text{where } SS_{\text{res}}(\underline{Z}_1) - SS_{\text{res}}(\underline{Z}) = \underline{\varepsilon}_1^\top \underline{\varepsilon}_1 - \underline{\varepsilon}^\top \underline{\varepsilon}$$

Special case II:

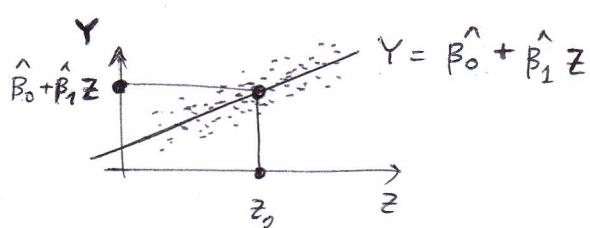
$$\begin{cases} H_0: \beta_1 = \beta_2 = \dots = \beta_r = 0 \\ H_1: \exists \beta_j \neq 0 \end{cases}$$

$$\rightarrow \frac{SS_{\text{res}}(\underline{Z}_1) - SS_{\text{res}}(\underline{Z})}{s^2_r} = \frac{\frac{\sum_{i=1}^n (Y_i - \bar{Y})^2 - \sum_{i=1}^n \hat{\varepsilon}_i^2}{r}}{\frac{\sum_{i=1}^n \hat{\varepsilon}_i^2}{n-(r+1)}} \sim F(r, n-(r+1))$$

Prediction

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

$$\underline{z}_0 = [1 \ z_{01} \ \dots \ z_{0r}]^\top \rightarrow Y_0 ?$$



$$\begin{cases} Y_0 = \underline{z}_0^\top \underline{\beta} + \varepsilon_0 & \varepsilon_0 \sim N(0, \sigma^2) \\ \mathbb{E}[Y_0] = \underline{z}_0^\top \underline{\beta} & (\varepsilon_0 \perp \varepsilon) \end{cases}$$

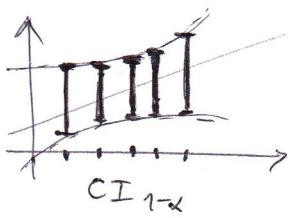
$$\frac{\underline{z}_0^\top \hat{\underline{\beta}} - \underline{z}_0 \underline{\beta}}{s \sqrt{\underline{z}_0 (\underline{Z}^\top \underline{Z})^{-1} \underline{z}_0}} \sim t(n-(r+1))$$

$$CI_{1-\alpha}(\underline{z}_0^\top \underline{\beta}) = \left[\underline{z}_0^\top \hat{\underline{\beta}} \pm s \sqrt{\underline{z}_0^\top (\underline{Z}^\top \underline{Z})^{-1} \underline{z}_0} t_{1-\frac{\alpha}{2}}(n-(r+1)) \right]$$

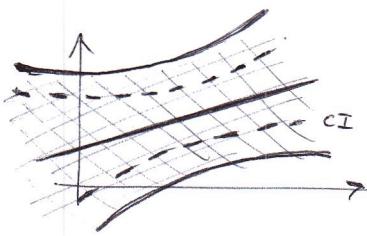
$$\text{Sim } CI_{1-\alpha}(\underline{z}_0^\top \underline{\beta}) = \left[\underline{z}_0^\top \hat{\underline{\beta}} \pm s \sqrt{\underline{z}_0^\top (\underline{Z}^\top \underline{Z})^{-1} \underline{z}_0} \sqrt{(r+1) F_{1-\alpha}(r+1, n-(r+1))} \right]$$

With $CI_{1-\alpha}(\underline{z}_0^\top \hat{\beta})$ we generate a CI for any \underline{z}_0 but fixed.

With $\text{sim } CI_{1-\alpha}(\underline{z}_0^\top \hat{\beta})$ we generate a bound of CI (so for any \underline{z}_0) :



only the vertical line has a level of $1-\alpha$



Note: this is the prediction of $E[Y_0 | \underline{z}_0]$, not Y_0 !

$$\frac{Y_0 - \underline{z}_0^\top \hat{\beta}}{S \sqrt{1 + \underline{z}_0^\top (\underline{Z}^\top \underline{Z})^{-1} \underline{z}_0}} \sim t(n-(r+1))$$

$$PI_{1-\alpha}(Y_0) = [\underline{z}_0^\top \hat{\beta} \pm S \sqrt{1 + \underline{z}_0^\top (\underline{Z}^\top \underline{Z})^{-1} \underline{z}_0} t_{1-\frac{\alpha}{2}(n-(r+1))}]$$

prediction interval
of probability $1-\alpha$

$$P(Y_0 \in PI_{1-\alpha}(Y_0)) = 1-\alpha$$

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

GLS : $\hat{\beta} = \underset{\beta}{\text{argmin}} (\underline{Y} - \underline{Z}\beta)^\top W^{-1} (\underline{Y} - \underline{Z}\beta)$
 Generalized least squares $\hat{\beta} = (\underline{Z}^\top W^{-1} \underline{Z})^{-1} \underline{Z}^\top W^{-1} \underline{Y}$

Collinearity

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

For every column :

$$\underline{z}_j := \gamma_0 + \gamma_1 \underline{z}_1 + \dots + \gamma_{j-1} \underline{z}_{j-1} + \gamma_{j+1} \underline{z}_{j+1} + \dots + \gamma_r \underline{z}_r$$

We substitute the j th column of \underline{Z} with \underline{z}_j (\uparrow) and we perform the linear model analysis again, getting R_j^2 .

$$\Rightarrow VIF_j = \frac{1}{1 - R_j^2}$$

Variance Inflation Factor

If $VIF(\beta_j) > 5/10 \Rightarrow$ probably β_j can be expressed as a linear combination of the other regressors

Collinearity and variable selection

Let's work with the problem "centered":

$$\underline{Z} \rightarrow \underline{Z}^* = \begin{bmatrix} z_{11} - \bar{z}_1 & \cdots & z_{1r} - \bar{z}_r \\ z_{21} - \bar{z}_1 & \cdots & z_{2r} - \bar{z}_r \\ \vdots & \ddots & \vdots \\ z_{n1} - \bar{z}_1 & \cdots & z_{nr} - \bar{z}_r \end{bmatrix} := \underline{Z}$$

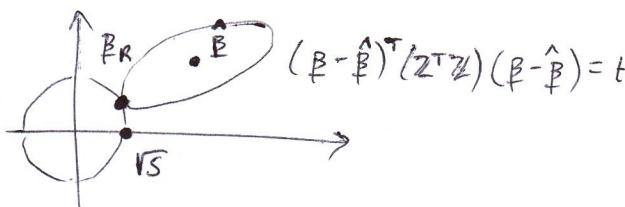
$$\underline{Y} \rightarrow \underline{Y}^* = [y_1 - \bar{Y}, \dots, y_n - \bar{Y}]^T := \underline{Y}$$

PCA Regression

- $\underline{Z} = [\underline{z}_1, \dots, \underline{z}_r]$
- PCA on $\underline{Z} \Rightarrow \text{PC}_1, \dots, \text{PC}_r$
- Reduce dimensionality: $\text{PC}_1, \dots, \text{PC}_k \quad k \leq r$
 $\underline{Z}^* := [\underline{\text{PC}}_1, \dots, \underline{\text{PC}}_k]$
- Fit $\underline{Y} = \underline{Z}^* \underline{\beta} + \underline{\varepsilon}$

$$\begin{aligned} \Rightarrow Y_0 &= \underline{z}_0^T \hat{\underline{\beta}} \\ &= \hat{\beta}_1 \text{PC}_1 + \dots + \hat{\beta}_k \text{PC}_k \\ &= \hat{\beta}_1 (e_{11} \underline{z}_1 + \dots + e_{1r} \underline{z}_r) + \dots + \hat{\beta}_k (e_{1k} \underline{z}_1 + \dots + e_{rk} \underline{z}_r) \\ &= \underline{z}_1 (\hat{\beta}_1 e_{11} + \dots + \hat{\beta}_k e_{1k}) + \dots + \underline{z}_r (\hat{\beta}_1 e_{r1} + \dots + \hat{\beta}_k e_{rk}) \\ &= \hat{\delta}_0 \underline{z}_1 + \dots + \hat{\delta}_r \underline{z}_r \end{aligned}$$

RIDGE Regression



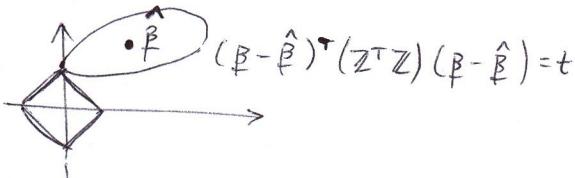
(not scale invariant!)

Problem :

$$\begin{cases} \underset{\underline{\beta}}{\text{arg min}} \| \underline{Z}(\underline{\beta} - \hat{\underline{\beta}}) \|^2 \\ \| \underline{\beta} \|^2 \leq s \\ \hat{\underline{\beta}} = (\underline{Z}^T \underline{Z})^{-1} \underline{Z}^T \underline{Y} \end{cases}$$

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

LASSO Regression



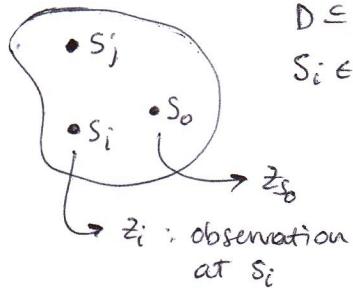
(λ with cross validation and prediction error for \underline{Y})

Problem :

$$\begin{cases} \underset{\underline{\beta}}{\text{arg min}} (\underline{\beta} - \hat{\underline{\beta}})^T (\underline{Z}^T \underline{Z}) (\underline{\beta} - \hat{\underline{\beta}}) \\ \| \underline{\beta} \|_1 \leq s \\ \hat{\underline{\beta}} \text{ OLS} \end{cases}$$

$$\underset{\underline{\beta}}{\text{arg min}} \| \underline{Y} - \underline{Z}\underline{\beta} \| + \lambda \| \underline{\beta} \|_1$$

SPATIAL DATA



$$D \subseteq \mathbb{R}^d$$

$$s_i \in D$$

$$s_1, \dots, s_n \in D \subseteq \mathbb{R}^d \text{ fixed}$$

$$z_{s_1}, \dots, z_{s_n} \text{ r.r. observations}$$

- Goals:
- Spatial dependence estimation
 - Prediction (spatial) : kriging

$$s_1, \dots, s_n \in D \text{ sites}$$

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

Assumptions:

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

Definitions:

- $m_s := \mathbb{E}[z_s]$ 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:
 1. $\mathbb{E}[z_s] = m \quad \forall s \in D$
 2. $\text{Cov}(z_{s_1}, z_{s_2}) = \underbrace{C(s_1 - s_2)}_{:= \text{COVARIogram}} \quad s_1, s_2 \in D$

\therefore COVARIogram

Properties:

- Symmetry: $C(-h) = C(h) \quad \forall h \in \mathbb{R}^d$
- Boundness: $|C(h)| \leq C(0) \quad \forall h \in \mathbb{R}^d$
- Pos. def.: $\sum_{i,j} \lambda_i \lambda_j C(s_i - s_j) \geq 0 \quad \forall \lambda_i, \lambda_j \in \mathbb{R}, \quad \forall s_i, s_j \in D$

Def. VARIOGRAM:
(under 2nd order of
stationarity) : $2\gamma(s_1 - s_2) := \text{Var}(z_{s_1} - z_{s_2}) = \mathbb{E}[(z_{s_1} - z_{s_2})^2]$

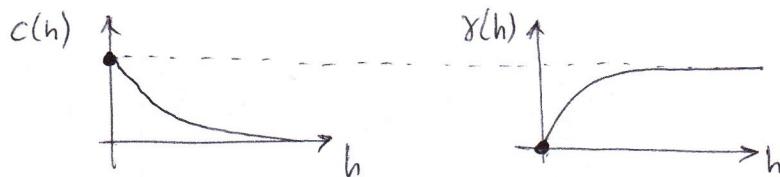
SEMI-VARIOGRAM: $\gamma(s_1 - s_2) := C(0) - C(s_1 - s_2)$

Properties:

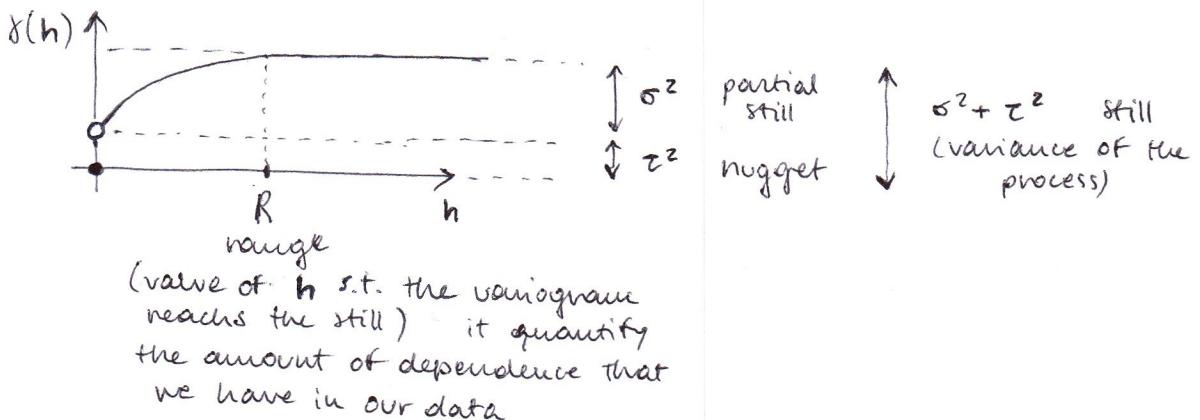
- Symmetry: $\gamma(-h) = \gamma(h)$
- Null at 0: $\gamma(0) = 0$
- Conditional negative definite: $\sum_{i,j} \lambda_i \lambda_j \gamma(s_i - s_j) \leq 0$

$$\begin{aligned} \lambda_i, \lambda_j &\in \mathbb{R} \\ \sum \lambda_i &= 0 \\ \forall s_i, s_j &\in D \end{aligned}$$

Examples:



Structural properties:



Def. (ISOTROPY): A field is isotropic if:
(under 2nd order of stationarity)

$$\text{Cov}(z_{s_i}, z_{s_j}) = C(\|s_i - s_j\|)$$

$$s_i, s_j \in D$$

Estimate Spatial Dependence (= estimate the variogram)

- Empirical estimate:

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

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

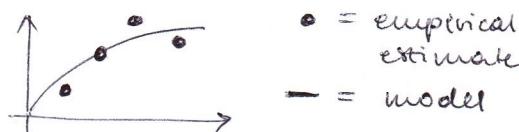
$$N_2(h) = \{(i,j) : \|s_i - s_j\| = \Delta h\}$$

← since we don't have too much data



- Parametric model:

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

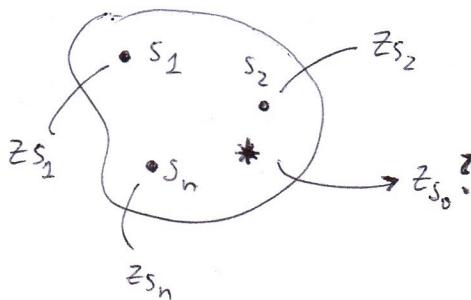


$$\text{For example: } \hat{\theta} = \arg \min_{\theta \in \Theta} \sum_{k=1}^K (\delta(h_k) - \gamma(h_k; \theta))^2 \text{ (OLS)}$$

Properties:

- δ_1, δ_2 valid variograms $\Rightarrow \delta_1 + \delta_2$ valid variogram
- $\alpha > 0$, δ_1 valid variogram $\Rightarrow \alpha \delta_1$ 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:

$$z_{s_0}^* = f(z_{s_1}, \dots, z_{s_n})$$

$$f = \operatorname{argmin} \mathbb{E}[(z_{s_0} - f(z_{s_1}, \dots, z_{s_n}))^2]$$

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 predictor
- unbiased

- Problem (kriging):

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

where $\lambda_0^*, \dots, \lambda_n^*$ solve?

$$\min_{\lambda_0, \lambda_1, \dots, \lambda_n \in \mathbb{R}} \mathbb{E}[(z_{s_0} - (\lambda_0 + \sum_i \lambda_i z_{s_i}))^2]$$

$$\text{such that: } \mathbb{E}[\lambda_0 + \sum_i \lambda_i z_{s_i}] = \mathbb{E}[z_{s_0}]$$

• Ordinary Kriging (• unknown mean
• stationarity 2nd order)

Assumptions: • unknown mean $\mathbb{E}[z_s] := m \quad \forall s \in D$

• stationarity 2nd order

• $\operatorname{Cov}(z_{s_1}, z_{s_2}) = C(s_1 - s_2) \quad \forall s_1, s_2 \in D$ known

(we actually estimate it:
 $\hat{s} \rightarrow \hat{\lambda} \rightarrow \hat{\Sigma}$)

Solution:

$$\text{Solve: } \begin{bmatrix} \Sigma & \mathbf{1} \\ \mathbf{1}^\top & 0 \end{bmatrix} \begin{bmatrix} \underline{\lambda} \\ \mu \end{bmatrix} = \begin{bmatrix} \sigma_0 \\ 1 \end{bmatrix} \implies \underline{\lambda}^* \implies z_{s_0}^* = \underline{\lambda}^{*\top} \underline{z}$$

where:

$$\Sigma = \operatorname{Cov}(\underline{z}) : \Sigma_{ij} = \operatorname{Cov}(z_{s_i}, z_{s_j}) = C(s_i - s_j)$$

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

μ = Lagrangian multiplier

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

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

• Universal Kriging ($\begin{array}{l} \text{: unknown mean} \\ \text{: non stationarity} \end{array}$)

Model :
$$\left\{ \begin{array}{l} z_s = m_s + \delta_s \quad s \in D \\ m_s = E[z_s] \quad \text{drift} \\ \delta_s = z_s - m_s \quad \text{residual} \end{array} \right.$$

- Assumptions :
- $m_s = \sum_{j=1}^L a_j f_j(s)$:
 - a_j unknown coeff.
 - $f_j(s)$ known regressors
 - δ_s random field ($\{\delta_s, s \in D\}$)
2nd order stationary random field for the residuals
 - $E[\delta_s] = 0$
 - $\text{Cov}(\delta_{s_1}, \delta_{s_2}) = \text{Cov}(z_{s_1}, z_{s_2}) = C(s_1 - s_2)$

Solution :

Solve :
$$\begin{bmatrix} \Sigma & F \\ F^T & 0 \end{bmatrix} \begin{bmatrix} \lambda \\ \mu \end{bmatrix} = \begin{bmatrix} \Sigma_0 \\ f_0 \end{bmatrix} \Rightarrow \underline{\lambda^*} \Rightarrow z_{s_0}^* = \underline{\lambda^{*\top}} \underline{z}$$

where :

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

F = design matrix

D = matrix of zeros

$$\Sigma_{0i} = \text{Cov}(z_{s_0}, z_{s_i})$$

f_0 = design vector : $f_{0j} = f_j(s_0)$: value of the j -th regressor at the new location

Estimations :

If Σ is unknown we should estimate $\hat{\gamma} \rightarrow \hat{C} \rightarrow \hat{\Sigma}$.
We know how to estimate $\hat{\gamma}$ from stationary data, here we have to use the residuals (since the data are not stationary)

- $\delta_{s_1}, \dots, \delta_{s_n}$ given $\Rightarrow \hat{\gamma} \rightarrow \hat{C} \rightarrow \hat{\Sigma}$
 - $\delta_{s_1}, \dots, \delta_{s_n}$ unknown $\Rightarrow \hat{\delta}_{s_i} \rightarrow \hat{\gamma} \rightarrow \hat{C} \rightarrow \hat{\Sigma}$
- $\hat{\delta}_{s_i} ? : \left\{ \begin{array}{l} \hat{\delta}_{s_i} = z_{s_i} - \hat{m}_{s_i} \\ \hat{m}_{s_i} = \sum_j \hat{a}_j f_j(s_i) \end{array} \right.$

where \hat{a} is estimated with an iterative procedure

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