

# Principal Component Analysis (PCA)

Aim: to create new variables to summarize the data. We would like to have less variables than previously.  
→ be careful: no clustering with PCA!

Remember:

the data  $X$ :

$$\begin{matrix} & & x_{11} & x_{12} & \dots & x_{1p} \\ n \text{ rows} & & x_{n1} & x_{n2} & & x_{np} \end{matrix}$$

$P$  columns

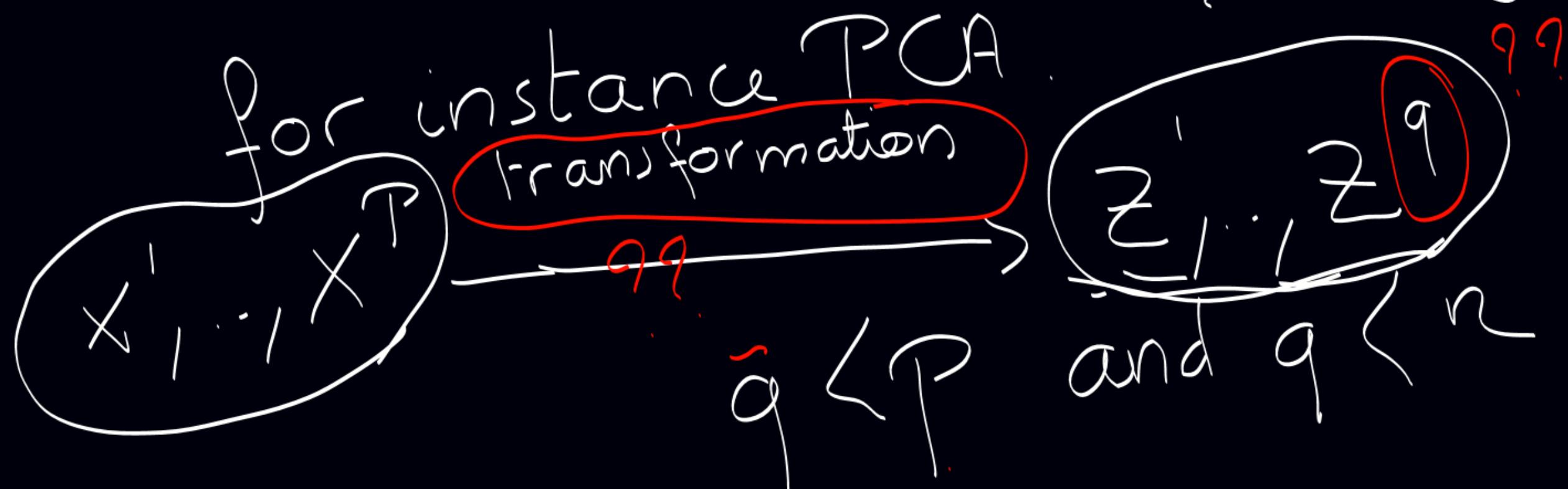
Now, we are often in the case where  $P > n$ .

If we want to perform linear regression

$\rightarrow \hat{P}$  because  $P > n \Rightarrow$  unicity

A solution:

We are going to reduce the number of variables by using



Example:

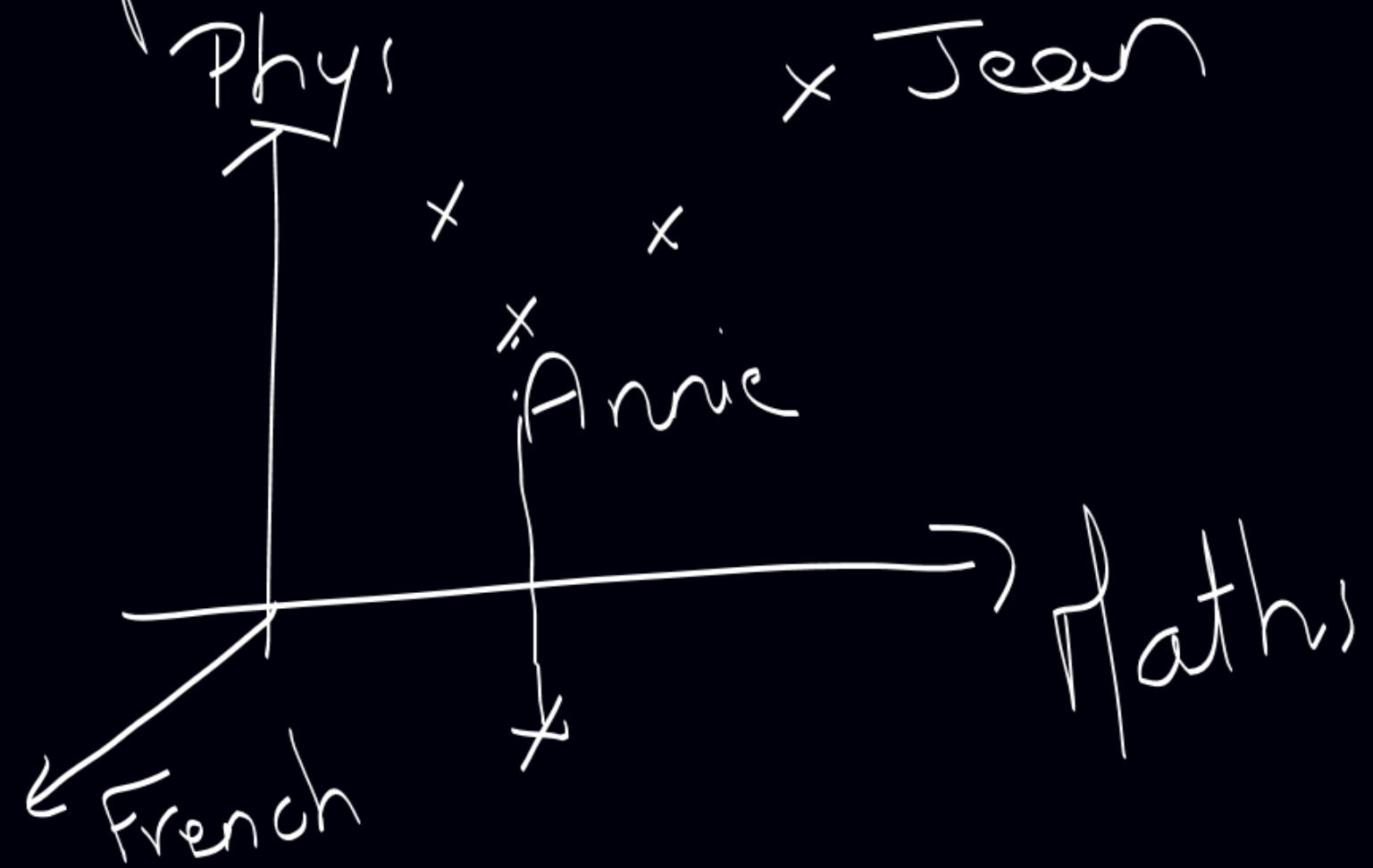
variables

Students	Maths	Phys.	French	English
Jean	6.00	6.00	5.00	5.50
	8.00	8.00	8.00	8.00
Alan		4.00	11.00	9.50
	6.00		15.50	15.00
Annie	14.50	14.50	12.00	12.50
Dionique	14.00	14.00	5.50	7.00
		10.00	14.00	11.50
Didier	11.00	7.00	8.50	9.50
Andre	5.50	12.50		
Pierre	13.00		12.50	12.00
Brigitte	9.00	9.50		
Emilie				

An individual is a point of  $\mathbb{R}^4$

→ I am not able to make a

drawing to represent the  
individuals



# 10] Computation of usual statistics

	Mean	Standard deviation	Min.	Max.
Maths	9.67	3.37	5.50	14.50
Phys.	9.83	2.99	6.00	14.50
French	10.22	3.47	5.00	15.50
English	10.06	2.81	5.50	15.00

↳ homogeneity among the variable !

Coefficients of correlation

4 rows	1	1	1
	1	1	1
		1	1
			1

4 columns

$$\text{Var}(x^i) = c \times \sum_{R=1}^n (x_R^i - \bar{x}^i)^2$$

$$\text{Cov}(x^i, x^j) = c \times \sum_{R=1}^n (x_R^i - \bar{x}^i)(x_R^j - \bar{x}^j)$$

$$\text{cor}(x^i, x^j) = \frac{c \times \sum_{R=1}^n (x_R^i - \bar{x}^i)(x_R^j - \bar{x}^j)}{\sqrt{(\sum (x_R^i - \bar{x}^i)^2)(\sum (x_R^j - \bar{x}^j)^2)}}$$

Let  $X$  be the table with the data.

$$X = \begin{pmatrix} x_{11} \\ \vdots \\ x_{1P} \\ \vdots \\ x_{n1} \\ \vdots \\ x_{nP} \end{pmatrix} \leftarrow \begin{array}{l} \text{observations} \\ \text{for the individual } 1 \\ \uparrow \\ \text{observations for the variable 1} \end{array}$$

two pre processings

\* To center the data : always done!

$$X_c = \begin{pmatrix} x_{11} - \bar{x}_{\cdot 1} & x_{12} - \bar{x}_{\cdot 2} & \dots & x_{1p} - \bar{x}_{\cdot p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} - \bar{x}_{\cdot 1} & x_{n2} - \bar{x}_{\cdot 2} & \dots & x_{np} - \bar{x}_{\cdot p} \end{pmatrix}$$

$$\begin{pmatrix} a & b & c & d \\ a & b & c & d \\ \vdots & & & \\ a & b & c & d \end{pmatrix} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} (a \ b \ c \ d)$$

where  $\bar{x}_{\cdot j} = \frac{1}{n} \sum_{k=1}^n x_{kj}$

empirical mean of  
the variable number  $j$

Rk: The empirical mean  
of each variable of  $X_c$  is

equal to 0.

$$\frac{1}{n} \sum_{j=1}^n (x_{0j} - \bar{x}_{0j}) = \frac{1}{n} \sum_{j=1}^n x_{0j} - \frac{1}{n} \sum_{j=1}^n \bar{x}_{0j}$$
$$= \bar{x}_{0j} - \cancel{\bar{x}} \cancel{\times} \cancel{\bar{x}_{0j}} = 0$$

2nd preprocessing : normalization

(not always done)

$$X_n = \begin{pmatrix} \frac{x_{11} - x_{\cdot 1}}{s_1} & \frac{x_{12} - x_{\cdot 2}}{s_2} & \frac{x_{1p} - x_{\cdot p}}{s_p} \\ \frac{x_{21} - x_{\cdot 1}}{s_1} & \frac{x_{22} - x_{\cdot 2}}{s_2} & \frac{x_{2p} - x_{\cdot p}}{s_p} \\ \vdots & \vdots & \vdots \\ \frac{x_{n1} - x_{\cdot 1}}{s_1} & \frac{x_{n2} - x_{\cdot 2}}{s_2} & \frac{x_{np} - x_{\cdot p}}{s_p} \end{pmatrix}$$

where  $s_j^2 = \frac{1}{n} \sum_{k=1}^n (x_{kj} - \bar{x}_{\cdot j})^2$

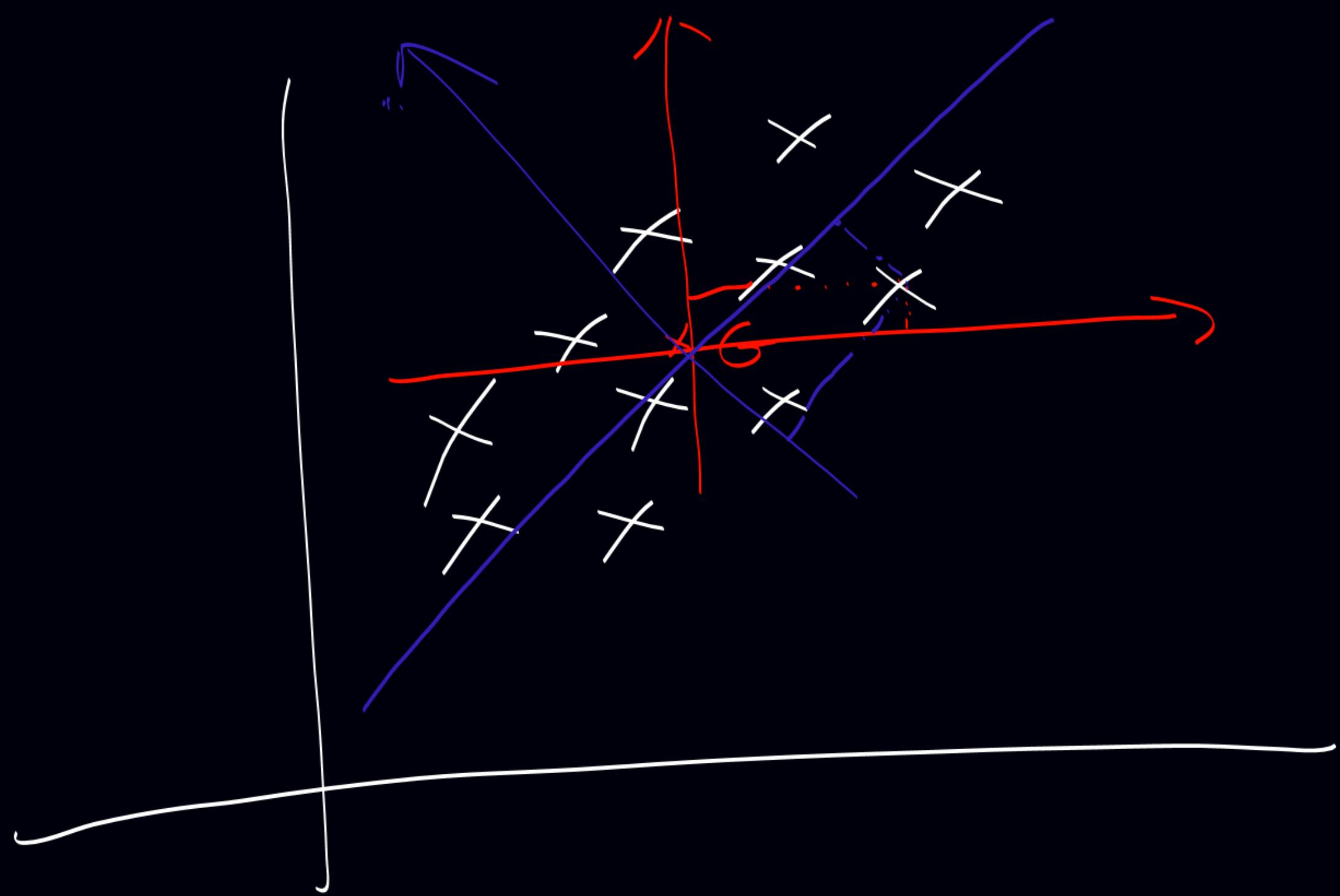
Matrix of data  $X$

$X_c$ : the centered version

$$V = \frac{1}{n} X_c^T X_c \rightarrow \text{the covariance matrix}$$

$X_n$ : the normalized version

$$R = \frac{1}{n} X_n^T X_n \rightarrow \text{correlation matrix}$$



Compute the eigenvalues  
and the eigenvectors of

$V$  and  $R$

$$V = \frac{1}{n} \mathbf{x}_c \mathbf{x}_c^T \quad \text{and} \quad R = \frac{1}{n} \mathbf{x}_n \mathbf{x}_n^T$$

↳ eigen

$$\text{cov}(X^i, X^j) = \frac{1}{n} \sum_{k=1}^n (x_{ki} - \bar{x}_{\cdot i})(x_{kj} - \bar{x}_{\cdot j})$$

$$\left( \frac{1}{n} X^t X \right)_{ij} = \frac{1}{n} \times \text{product between the}$$

the column  $i$  and  
the column  $j$  of  $X^t$

$\times \mathbb{R}^3$

$$A : (x, y, z) = x \times (1, 0, 0) + y \times (0, 1, 0) + z \times (0, 0, 1)$$
$$\{(1, 0, 0), (0, 1, 0), (0, 0, 1)\}$$

$e_1 = (1, 0, 0)$   
 $e_2 = (0, 1, 0)$   
 $e_3 = (0, 0, 1)$

all the points  
of  $\mathbb{R}^3$  can be  
written as a linear  
combination of  $\{e_1, e_2, e_3\}$

The family  $\{e_1, e_2, e_3\}$  generates  $\mathbb{R}^3$ .

To be able to say that  $\{e_1, e_2, e_3\}$  is a basis of  $\mathbb{R}^3$ , the linear combination should be unique.

A way to show this is to prove that if  $(0,0,0) = \lambda_1 e_1 + \lambda_2 e_2 + \lambda_3 e_3 \Rightarrow \lambda_1 = \lambda_2 = \lambda_3 = 0$

$A \Rightarrow$  eigenvalues  $\lambda_1, \dots, \lambda_p$   
eigenvectors  $u_1, \dots, u_p$

$\exists \lambda_i \neq \lambda_j \quad u_i \perp u_j$

$$\sum_k u_{ik} \times u_{jk} = 0$$
$$k=1$$

$\sum \lambda_i = \lambda_j$   
 $u_i$  and  $u_j$  is not necessary  
orthogonal

There exists a technic named  
Gram-Schmidt that allow to  
create  $v_i$  and  $v_j$  eigen vectors of A  
associated to the eigenvalue  $\lambda_i$   
such that  $v_i \perp v_j$

$$U_1 = \begin{pmatrix} 1, 0 \\ 0, 1 \end{pmatrix} \quad \geq \quad V_1 = \begin{pmatrix} 1, 0 \\ 0, 1 \end{pmatrix}$$

$$U_2 = \begin{pmatrix} 1, 1 \\ 1, 1 \end{pmatrix} \quad \geq \quad V_2 = \begin{pmatrix} 0, 1 \\ 1, 0 \end{pmatrix}$$

$$V = P D P^{-1}$$

with  $P$  composed  
 & the eigenvectors  
 $P^{-1} = t(P)$

Now what is T\$ scores?

$$X_C = \begin{pmatrix} & \text{ind1} \\ & \text{ind2} \\ & \vdots \\ & \text{indn} \end{pmatrix}$$

$X_C \% * \% E1\$ \text{ vectors}$

$E1\$ \text{ vectors}$

$$\left( v_1 \mid v_2 \mid v_3 \mid v_4 \right)$$

Rba:

$A$  is a matrix  
 $\lambda$  is an eigenvalue for  $A$  if  
there exists a vector  $X$  which is  
not null if  $AX = \lambda X$   $X$ : eigenvector  
associated to  $\lambda$

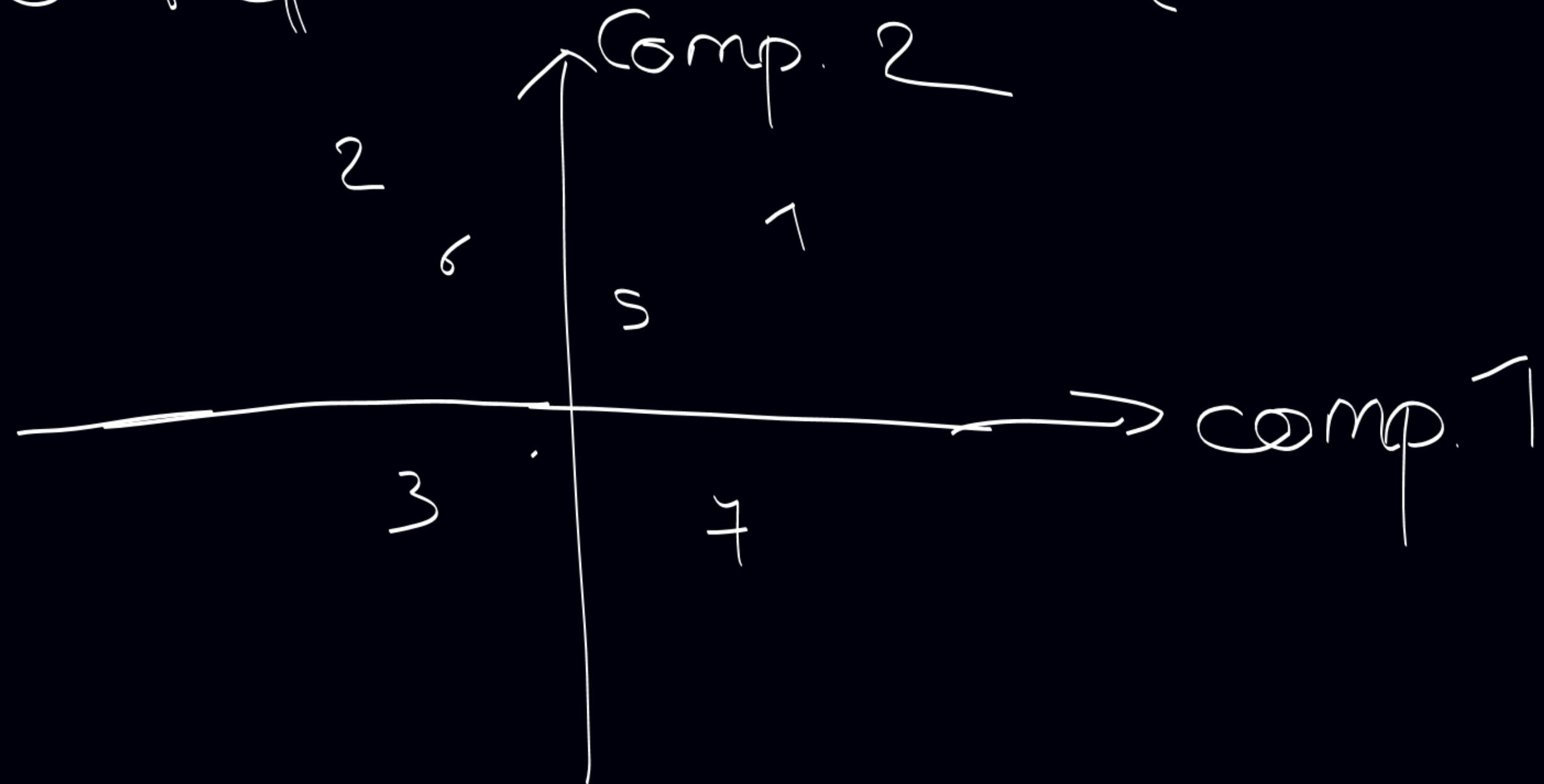
Let  $\mu$  be a real

if  $X$  is an eigenvector for  $A$   
associated to  $\lambda$

$\Rightarrow \mu X$  is also an eigenvector  
for  $A$  associated to  $\lambda$

In particular  $-X$  is an eigenvector

To represent the data



# notion of Inertia

↳ associated to the notion of

distance

let  $u_i = \begin{pmatrix} u_{i1} \\ \vdots \\ u_{ip} \end{pmatrix}$   $u_j = \begin{pmatrix} u_{j1} \\ \vdots \\ u_{jp} \end{pmatrix}$

$d(u_i, u_j) = \sum_{k=1}^p (x_{ik} - x_{jk})^2$

↓ scalar product

$\langle \overrightarrow{ou_i}, \overrightarrow{ou_j} \rangle = u_i \cdot u_j$

A diagram illustrating the geometric interpretation of the cosine formula. It shows a point  $O$  at the origin. Two vectors originate from  $O$ : one pointing towards a point  $u_i$  and another pointing towards a point  $u_j$ . The angle between these two vectors is labeled  $\alpha_{ij}$ .

$$\cos(\alpha_{ij}) = \frac{\langle \vec{Ou}_i, \vec{Ou}_j \rangle}{\|\vec{Ou}_i\| \times \|\vec{Ou}_j\|}$$

with  $\|\vec{Ou}_i\|^2 = \langle \vec{Ou}_i, \vec{Ou}_i \rangle$

$$X_C = \begin{pmatrix} x_{11} - x_{\cdot 1} \\ x_{12} - x_{\cdot 2} \\ \vdots \\ \vdots \\ x_{n1} - x_{\cdot 1} \\ x_{n2} - x_{\cdot 2} \\ \vdots \\ x_{np} - x_{\cdot p} \end{pmatrix}$$

$$U_1 = \begin{pmatrix} x_{11} - x_{\cdot 1} \\ x_{12} - x_{\cdot 2} \\ \vdots \\ x_{1P} - x_{\cdot P} \end{pmatrix}$$

Inertia with respect to the mean

point of  $\bar{x}_c$

$$\begin{aligned}\bar{I}_o &= \frac{1}{n} \sum_{i=1}^n d^2(o_i, u_i) \\ &= \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_{\cdot j})^2\end{aligned}$$

$$= \sum_{j=1}^P \left( \frac{1}{n} \sum_{i=1}^n (x_{ij} - \bar{x}_{\cdot j})^2 \right)$$

Variance of  $\bar{V}_j$

$$\bar{I}_0 = \sum_{j=1}^P V[\bar{V}_j] = \text{tr}(\Sigma)$$

where  $\Sigma$  : matrix of  
variance-covariance

Let  $\Delta$  be an axis such that

$$o \in \Delta$$

$$\underline{\pi}_{\Delta} = \frac{1}{n} \sum_{i=1}^n d^2(h_{\Delta,i}, u_i)$$

with  $h_{\Delta,i}$ : projection of  $u_i$   
on  $\Delta$



Let  $V$  denotes a subspace of  $\mathbb{R}^P$   
and  $V^*$  \_\_\_\_\_ the orthogonal

Rk:  $o \in V$  supplementary of  $V$  in  $\mathbb{R}^P$

$$\underline{\underline{I}}_V = \frac{1}{n} \sum_{i=1}^n d(h_{V,i}, v_i)^2$$

with  $h_{V,i}$ : projection of  $v_i$   
on  $V$

We can prove that:

$$d^2(h_{V,i}, o_i) + d^2(h_{V^*,i}, o_i) = d^2(o, o_i)$$
$$= d^2(o, h_{V,i}) + d^2(o, h_{V^*,i})$$

Because of this:

generalization

$$\overline{I}_o = \overline{I}_V + \overline{I}_{V^*} = \overline{I}_{A_1} + \overline{I}_{A_2} + \dots + \overline{I}_{A_p}$$

all  $A_i$  are with dimension 1  
and orthogonal

1st Step : determination of  $\Delta_1$

→ idea:  $\Delta_1$  should be such that

$\Delta_1$  is as close as possible to  
the cloud of points.

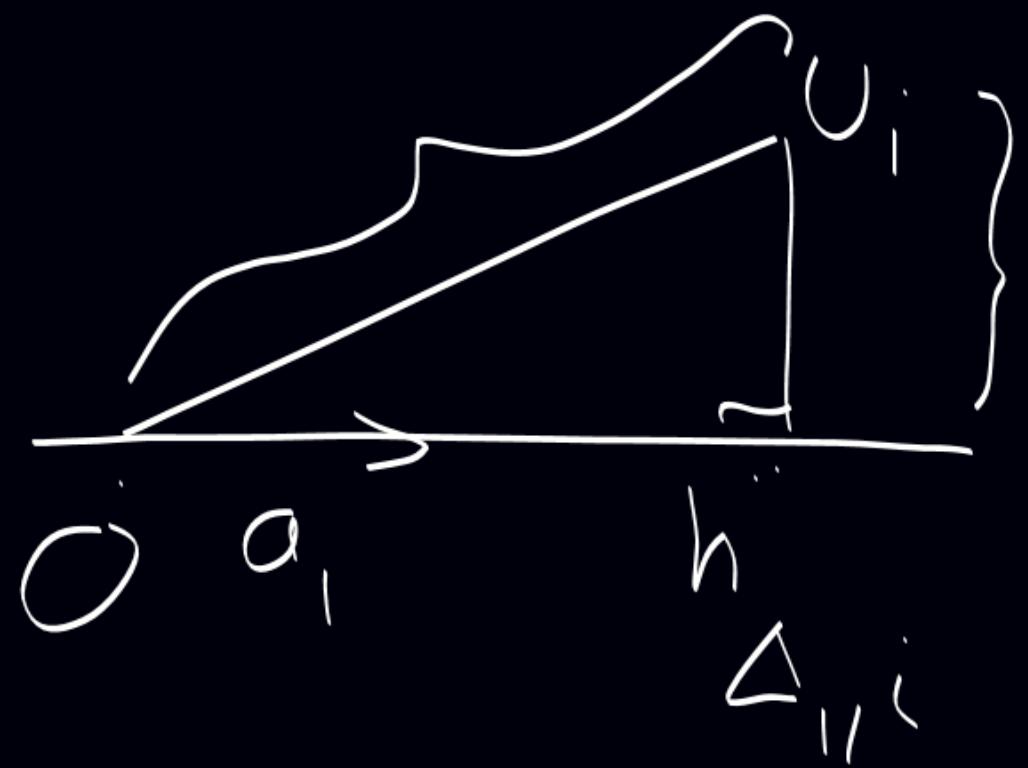
we search  $\Delta_1$  such that :

- $o \in \Delta_1$
- $I_{\Delta_1}$  minimum  $\Leftrightarrow I^*_{\Delta_1}$  maximum

$\Delta_1$  is defined by a unitary vector  
of  $\mathbb{R}^P$  denoted  $\vec{a}_1$

To determine  $\vec{a}_1$  we should have  
 $\overline{\Delta}_1^*$  maximum and  $\|\vec{a}_1\|^2 = 1$

$$\left\langle \overrightarrow{Oa_1}, \overrightarrow{oh} \Delta_{1,i} \right\rangle^2 = d^2(O, h \Delta_{1,i})$$



$$= \left\langle \overrightarrow{Oa_1}, \overrightarrow{Ou_i} \right\rangle^2 \\ = a_1 u_i + u_i a_1$$

$$\hat{\boldsymbol{\alpha}}_1 = \frac{1}{n} \sum_{i=1}^n (\boldsymbol{t}_{\boldsymbol{a}_1} \boldsymbol{U}_i^T \boldsymbol{U}_i \boldsymbol{a}_1)$$

$$= \boldsymbol{t}_{\boldsymbol{a}_1} \left( \frac{1}{n} \sum_{i=1}^n \boldsymbol{U}_i^T \boldsymbol{U}_i \right) \boldsymbol{a}_1$$

$\Sigma$  : matrix of variance  
covariance

Optimization  $P^b$ :

maximize  $t_{\alpha_1} \sum \alpha_i$  with  $t_{\alpha_1} \alpha_i = 1$

$$g(\alpha_1) = t_{\alpha_1} \sum \alpha_i - \lambda \left( \alpha_1 \alpha_1 - 1 \right)$$

$$\frac{\partial g(\alpha_1)}{\partial \alpha_1} = 2 \sum \alpha_i - 2 \lambda \alpha_1$$

$$\underline{\text{Rk}} : \alpha_1 = (a_{11}, \dots, a_{1P})$$
$$\frac{\partial g(\alpha_1)}{\partial \alpha_1} = \begin{pmatrix} \frac{\partial g}{\partial a_{11}} \\ \vdots \\ \frac{\partial g}{\partial a_{1P}} \end{pmatrix}$$

$$(1) \left| \begin{array}{l} \sum a_i = \lambda a_i \Leftrightarrow a_i \text{ is an eigen-} \\ \text{-vector of } \Sigma \text{ associated} \\ t_{a_i a_i} = 1 \text{ to the eigenvalue } \lambda. \end{array} \right.$$

$$(2) \Rightarrow \underbrace{t_{a_i} \sum a_i}_{\text{quantity that}} = \lambda t_{a_i a_i} = \lambda \quad \left| \Rightarrow \lambda \text{ is} \right. \\ \text{we want to maximize} \quad \left. \begin{array}{l} \text{The biggest} \\ \text{eigenvalue} \\ \text{of } \Sigma \end{array} \right.$$

In conclusion :

$\alpha_i$  is a unitary eigenvector  
associated to the biggest  
eigenvalue of  $\Sigma$ .

olympic2 .Txt

TODSTI .Pdf

Project1 .pdf → groups

& 4 persons!