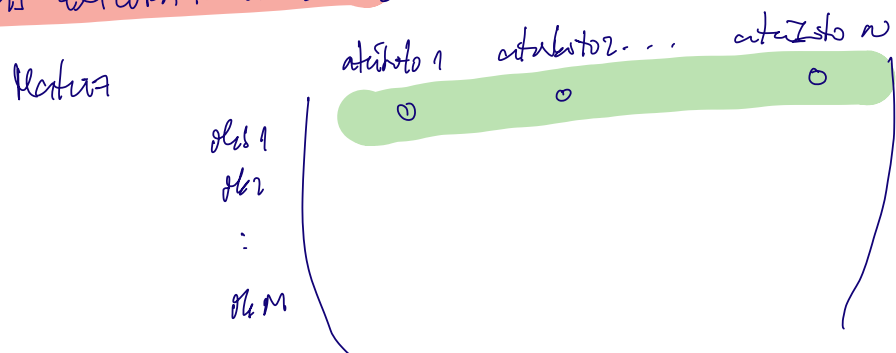


1. Intro AI

2. ANÁLISIS EXPLORATORIO DATOS



pandas → .csv
pandas → .excel → pandas dataframe

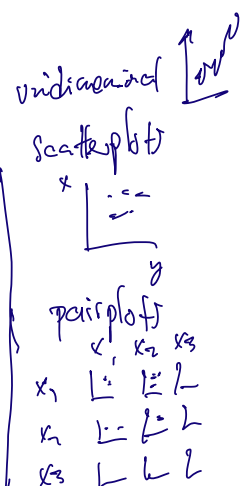
- filter
- order
- statistics ...

sklearn

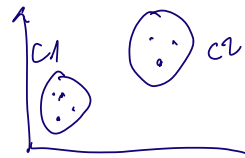
imputation data
(missing, NaNs...)

seaborn

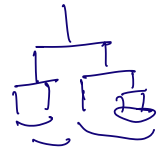
visualizer



3. CLUSTERING



- k-MEANS → simple
- AGGLOMERATIVE → dendrogram
- GAUSSIAN MIXTURE MODELS → How many clusters?



...
 selección > clustering

mpg, irr...

↳

Indices
 Pearson
 selección modelos
 BIC, AIC

4. REDUCCIÓN DE DIMENSIONALIDAD

↳ ANÁLISIS DE COMPONENTES PRINCIPALES (PCA)

Matriz Datos

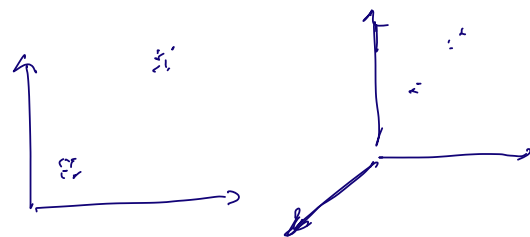
$X_{NROWS \times NFEAT}$

IMPUTACIÓN

CLUSTERING

+ w class blocks
 CLASIFICACIÓN

$X_{NROWS \times NFEAT}$

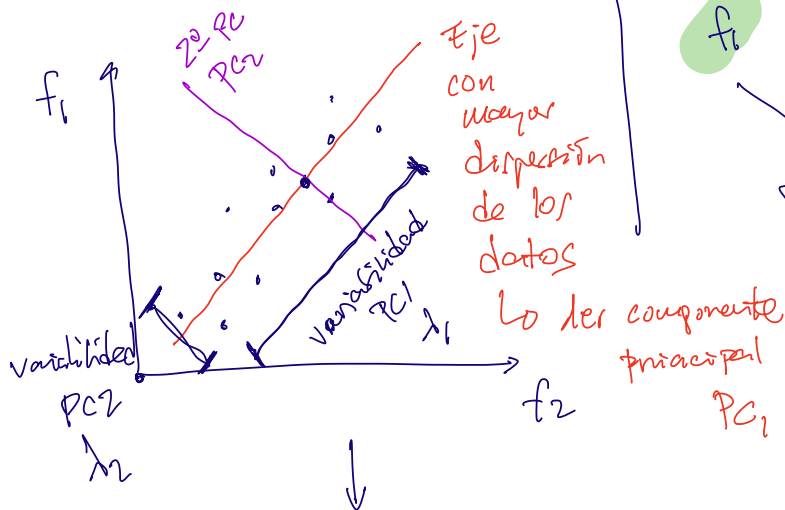


$X_{NROWS \times d}$

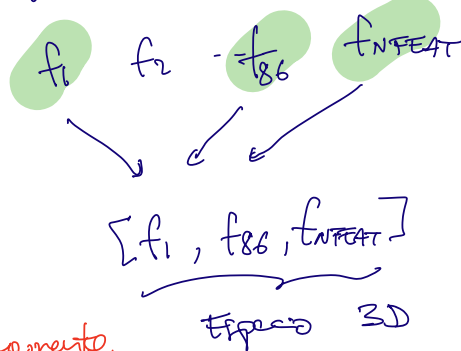
↓
 dimensión reducida
 $d < NFEAT$

CUANTO MAYOR LA DIMENSIONALIDAD
 (# FEATURES) → MÁS DATOS
 SON NECESARIOS PARA
 IDENTIFICAR PATRONES

Ex: PCA 2D NFEAT=2



SELECCIÓN ATRIBUTOS



NFEAT = 1

2d PCA 2d
(f1, f2) → (PC1, PC2)

Desplazamiento } transf. lineal
Rotación } sistema de coord.

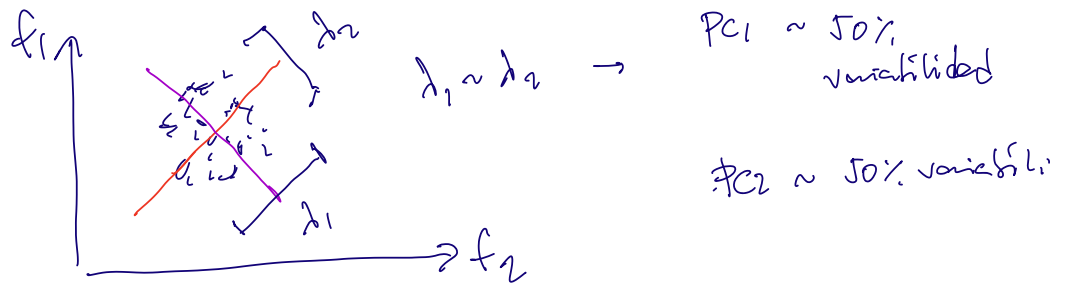
$$\begin{pmatrix} PC_1 \\ PC_2 \end{pmatrix} = \begin{pmatrix} \alpha_{11} & \alpha_{12} \\ \alpha_{21} & \alpha_{22} \end{pmatrix} \begin{pmatrix} f_1 \\ f_2 \end{pmatrix}$$

$$PC_1 = \alpha_{11} f_1 + \alpha_{12} f_2$$

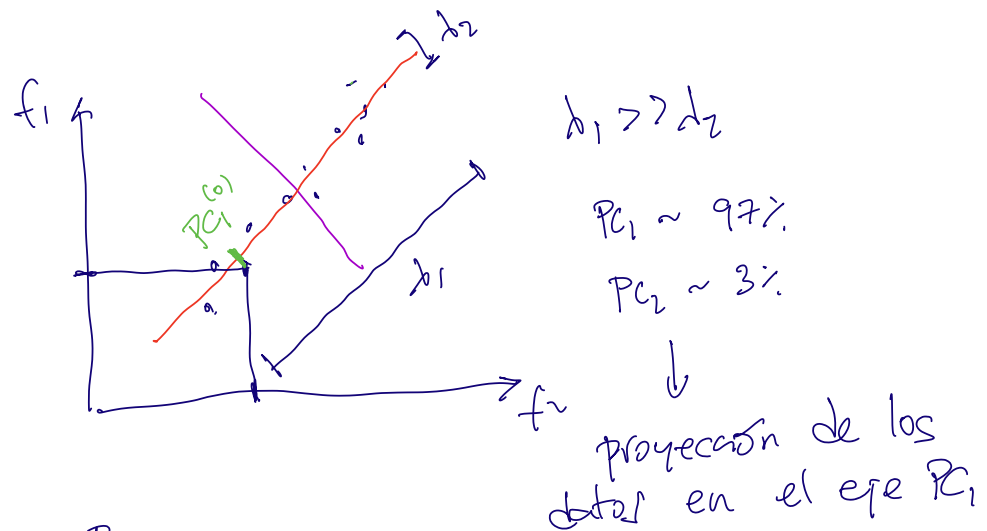
↓
La nueva variable es combinación lineal de f1 y f2

cuánta variabilidad explica el $PC_1 = \frac{\lambda_1}{(\lambda_1 + \lambda_2)}$

$$PC_2 = \frac{\lambda_2}{(\lambda_1 + \lambda_2)}$$



NO PUEDO REDUCIR LA
DIMENSIONALIDAD
PORQUE AUNAS HAY
CORRELACIÓN ENTRE f_1 y f_2



$(f_1, f_2) \rightarrow PC_1$
(1d)

Matemáticamente:

① ESCALAR
DATOS

$X_{\text{NOBS} \times \text{NFEAT}}$
estandarizar

$$X_S = (X - \bar{X}) / \sigma$$

datos
escalados

↓
media de
cada
FEATURE
centrado

↑
escalado

2) matrix covarianza

$$C = \frac{1}{N_{OBS}-1} \cdot \underbrace{\left(\frac{X-\bar{X}}{\sigma} \right)^+}_{(NFEAT \times N_{OBS})} \cdot \underbrace{\left(\frac{X-\bar{X}}{\sigma} \right)}_{(N_{OBS} \times NFEAT)}$$

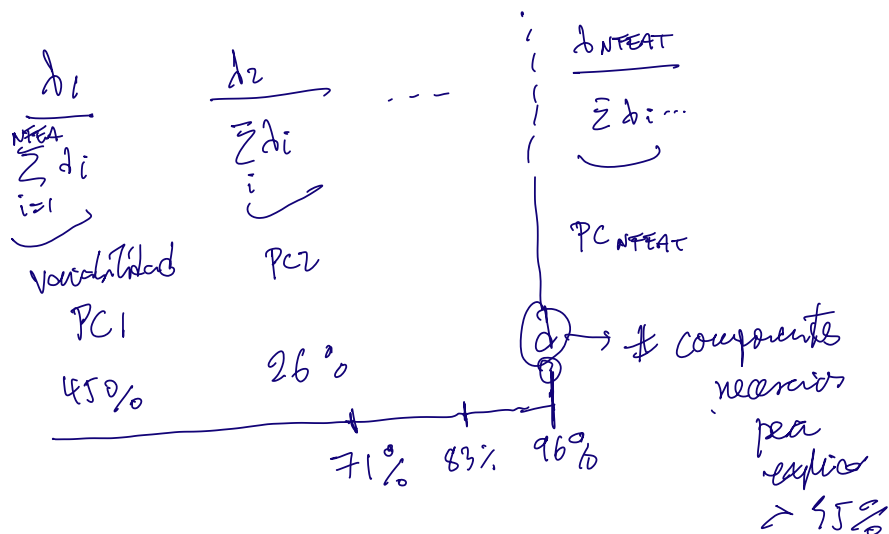
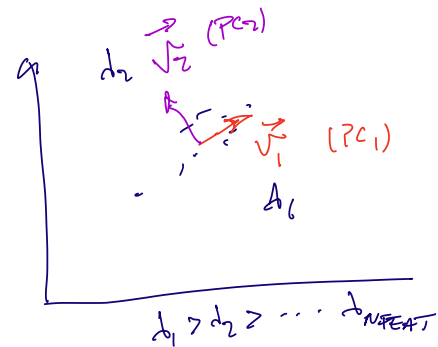
$NFEAT \times NFEAT$

3) Diagonalizar C

↓
valores propios $\lambda_i \rightarrow$ variabilidad explicada PC_i
vectores propios $\vec{v}_i \rightarrow$ dirección PC_i
 $i = 1 \dots NFEAT$

$$C \cdot \vec{v} = \lambda \cdot \vec{v}$$

4) Cuantos PC_i son necesarios para explicar el 95% variabilidad



$$(f_1 \dots f_{NFEAT}) \rightarrow (PC_1 \dots PC_d)$$

IMPORTANCIA ESCALA DATOS:

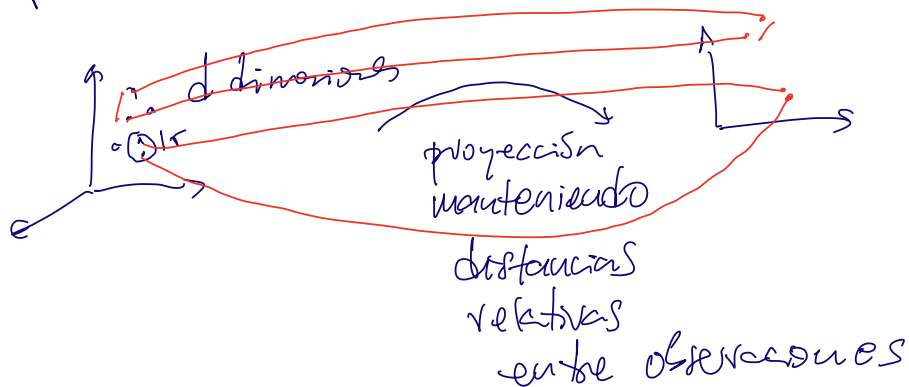
$$X = \begin{pmatrix} f_1(\mu m) & f_2(\mu m) \\ 1 & 33 \\ 0.5 & 89 \\ 0.3 & 124 \\ 1.2 & 96 \\ 5.4 & 348 \\ 1.8 & 430 \end{pmatrix} \rightarrow XS \begin{pmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} & \begin{bmatrix} 1 \\ 1 \end{bmatrix} \end{pmatrix}$$

$\begin{matrix} \uparrow \\ \sigma_1 \end{matrix}$
 $\begin{matrix} \uparrow \\ \text{rango } \sigma_2 \end{matrix}$

SELECCIÓN ATRIBUTOS

5. TÉCNICAS VISUALIZACIÓN DATOS MULTIDIMENSIONALES

Espacio atributos dimensión $d = N_{ATR}$



→ MULTIDIMENSIONAL SCALING (MDS)

→ t-SNE