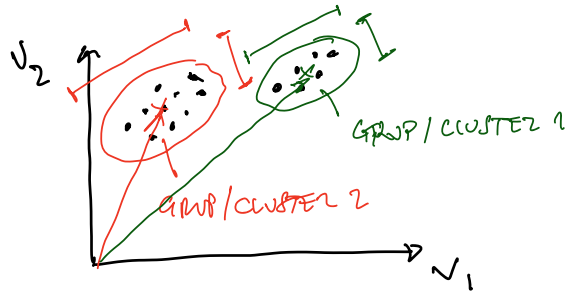


Tema 2: TÈCNiques D'AGUPAMENT DE DADES (CLUSTERING ALGORITHMS)

↳ No-SUPERVISAT

$X_{NOBS \times NVAR}$ (dades numèriques)

Ex: $NOBS = 16$
 $NVAR = 2$



$NOBS = 10^4$
 $NVAR = 18$ } $\rightarrow ?$

Algorithms
 ↙ ↘
 K-means Agrupament
 (K-centroids) Agglomeration

Models de
mescla
Gaussian
(GMM)

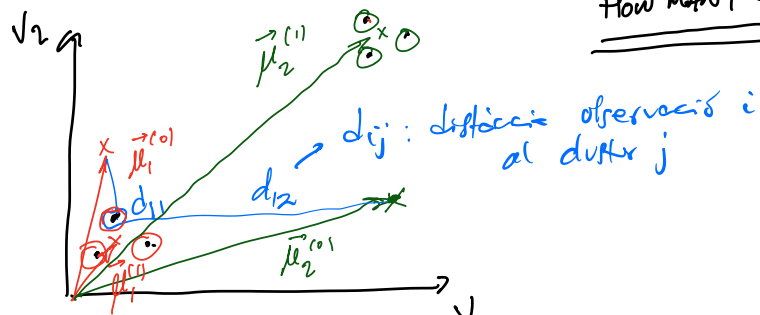
distàncies
entre observacions

distribució
probabilitat
joint

How MANY clusters?

① K-MEANS : Algorisme iteratiu (iteratiu)

Ex:



(1) Establir # clusters que volem identificar (k)

(2) Definir centroides aleatoris (k)

per cadascun dels k grups

$\{\bar{\mu}_1^{(0)}, \bar{\mu}_2^{(0)}, \dots, \bar{\mu}_k^{(0)}\}$

centroides clusters a iteració 0

goto (3)
 until
 $|\Delta \hat{\mu}_i| < \epsilon$

[3] Assignar cada observació al cluster
 que tingui el centroide més proper.
 [4] Recalcular centroides
 $\{\vec{\mu}_1^{(1)}, \vec{\mu}_2^{(1)}, \dots, \vec{\mu}_k^{(1)}\}$

clustering final

$$Y = \{y_1, y_2, \dots, y_{n_{obs}}\}$$

Assign
 cluster
 per cada
 observació

$$y_i \in \{1, \dots, k\}$$

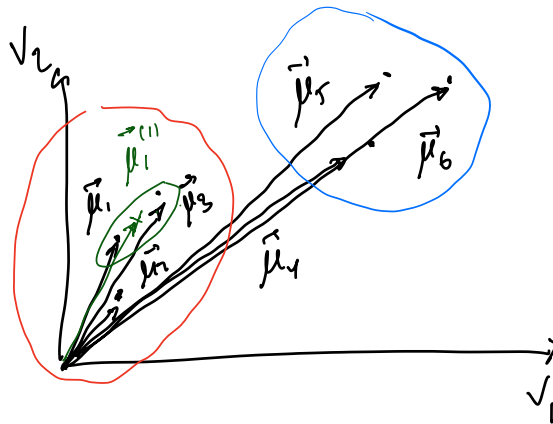
$K?$ → pairplot: v_1 v_2 $v_3 \dots v_{n_{obs}}$
 v_1 v_2 v_3 v_4

② AGUPAMENT AGLOMERATIU

↳ Hierarchical clustering

↳ agglomeratiu

↳ divisiu



initial clusters = # observations

↓
 agrupar clusters
 fins que
 no hi ha
 més guany

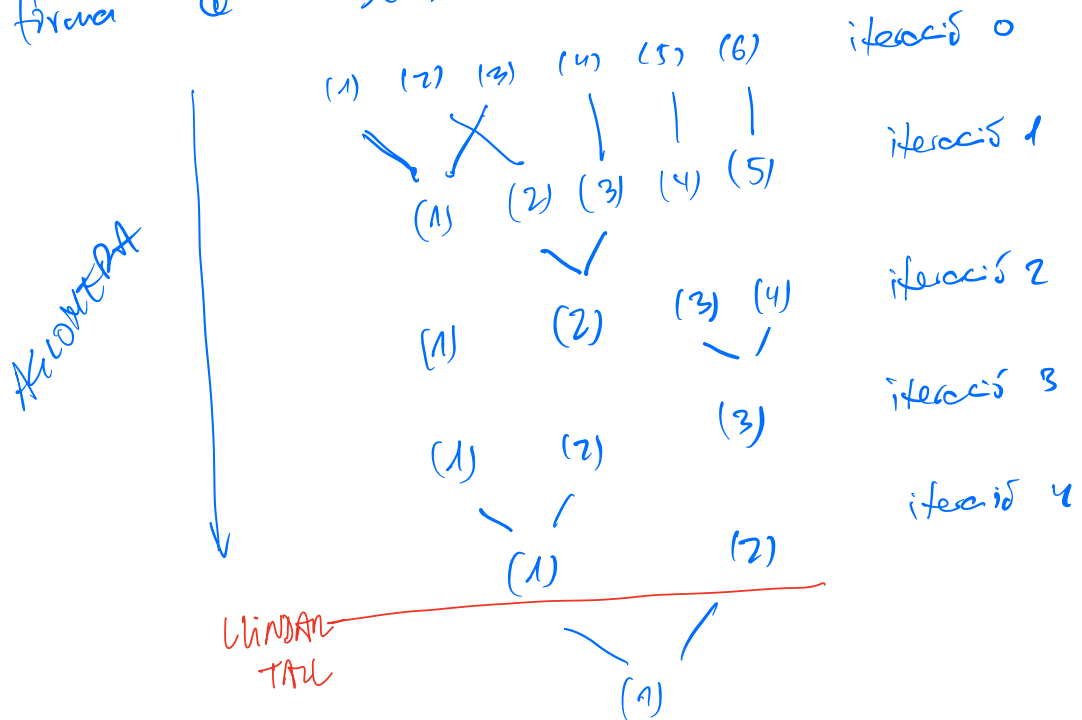
Ⓚ

- (1) Indices # clusters K
- (2) Initializer #clusters \rightarrow # observations
- (3) Agrupos els 2 clusters mit proper

$$\begin{array}{lcl} (\vec{\mu}_1, \vec{\mu}_3) & \rightarrow & \vec{\mu}_1^{(1)} \\ \vec{\mu}_2 & \rightarrow & \vec{\mu}_2^{(1)} \\ \vec{\mu}_3 & \rightarrow & \times \\ \vec{\mu}_4 & \rightarrow & \vec{\mu}_3^{(1)} \\ \vec{\mu}_5 & \rightarrow & \vec{\mu}_4^{(1)} \\ \vec{\mu}_6 & \rightarrow & \vec{\mu}_5 \end{array} \quad \left. \vphantom{\begin{array}{lcl} (\vec{\mu}_1, \vec{\mu}_3) & \rightarrow & \vec{\mu}_1^{(1)} \\ \vec{\mu}_2 & \rightarrow & \vec{\mu}_2^{(1)} \\ \vec{\mu}_3 & \rightarrow & \times \\ \vec{\mu}_4 & \rightarrow & \vec{\mu}_3^{(1)} \\ \vec{\mu}_5 & \rightarrow & \vec{\mu}_4^{(1)} \\ \vec{\mu}_6 & \rightarrow & \vec{\mu}_5 \end{array}} \right\} \begin{array}{l} 5 \text{ clusters} \\ (\#obs - 1) \end{array}$$

(2) Repeat until #clusters = k

la différence d'agrupements es pot representar gràficament
en forma de DENDROGRAMA



③ MODELS DE MESCLA GAUSSIANA (GAUSSIAN MIXTURE MODEL)

$$X_{N \times N \times N \times N} = \left(\text{---} \right) \rightarrow \vec{x} \in \mathbb{R}^{N \times N}$$

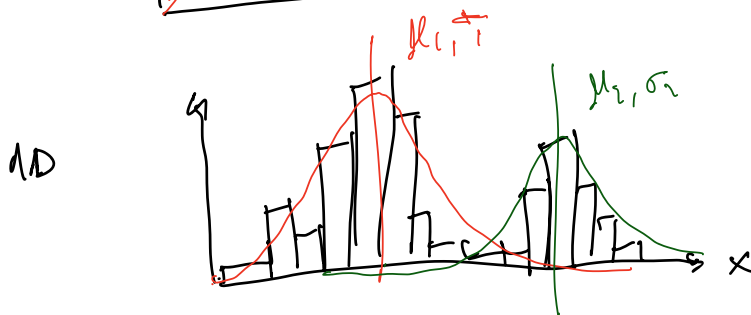
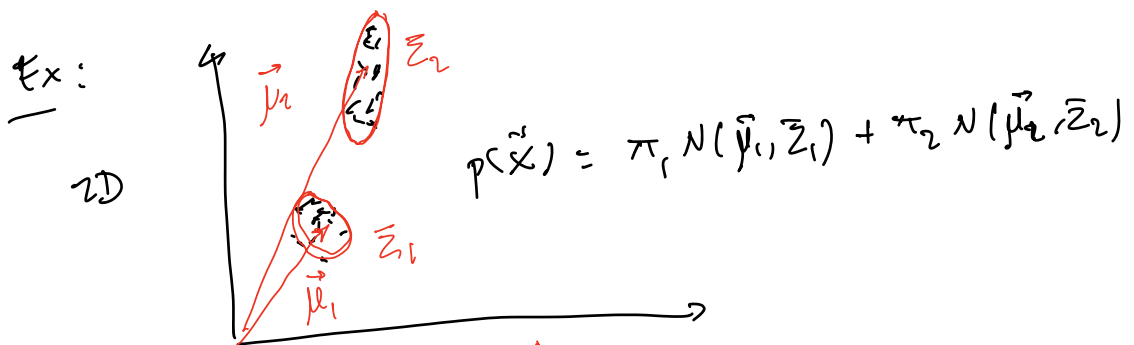
(numèriques)

$$p(\vec{x}) = \sum_{i=1}^{Ng} \pi_i \cdot N(\vec{\mu}_i, \Sigma_i)$$

pdf observacions pesos centroides matriu covariància

$N \times 1$ $N \times N$

descriu la distribució de les observacions a l'espai $N \times N$ dimensional com una suma ponderada de gaussianes multivariades



$$p(x) = \pi_1 N(\mu_1, \sigma_1) + \pi_2 N(\mu_2, \sigma_2)$$

$$p(\vec{x}) = \sum_{i=1}^{Ng} \pi_i N(\vec{\mu}_i, \Sigma_i) \rightarrow \text{paràmetres a estimar.}$$

estimació ~~parametres~~:
 Algorisme d'expectació -
Maximització
 (EM-algorithm)

VARIABLE LATENT
 ↓
 assignació de observacions → Gaussianes
 ↓
 etiqueta de clustering

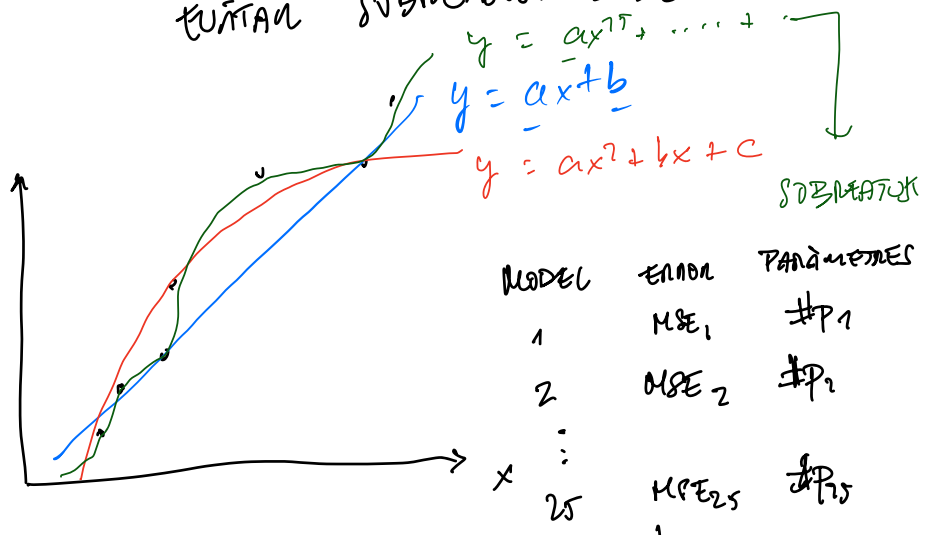
Dades → $\begin{bmatrix} \text{GMM} \\ \text{NA} \end{bmatrix}$ → clustering variable

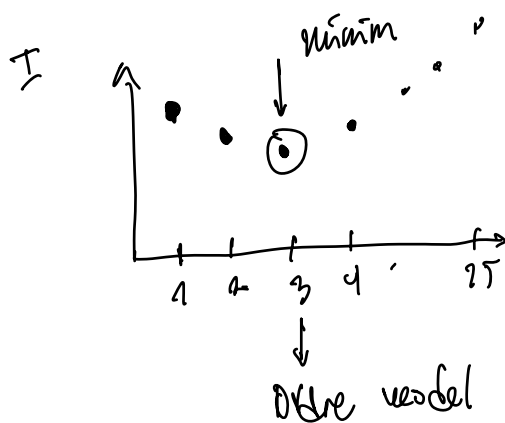
How many clusters? → Model Selection
 Índex de penalització (Bayesian Information Criterion)
 Akaike i.c.)

$I \sim$ error model - complexitat model
 (MSE) (# parameters)

EVITAR SOBREATUST DADES

ex: regressió polinòmica
 y



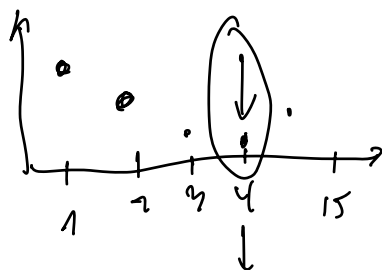


↓

Index	parameters
model	Index
1	I_1
2	I_2
⋮	⋮
25	I_{25}

enum + Relaxed models

model	nbr	#parameters	I
1 gaussian			I_1
2 gaussian			I_2
⋮			⋮
15 gaussian			I_{15}



clusters = # gaussianes.

