

TÉCNICAS DE AGUPAMIENTO (CLUSTERING)

DATOS NUMÉRICOS : $A_{\text{NOBS} \times \text{FEATURES}}$

$$= \begin{pmatrix} x_1 & x_2 & \dots & x_M \\ \vdots & \vdots & \ddots & \vdots \\ a_{ij} \in \mathbb{R} & & & \end{pmatrix} \begin{matrix} \text{obs 1} \\ \text{obs 2} \\ \vdots \\ \text{obs } n \end{matrix}$$

3 TÉCNICAS

↳ K-means

↳ aglomerativo / jerárquico

↳ Mezclas Gaussianas (Gaussian mixture models, GMM)

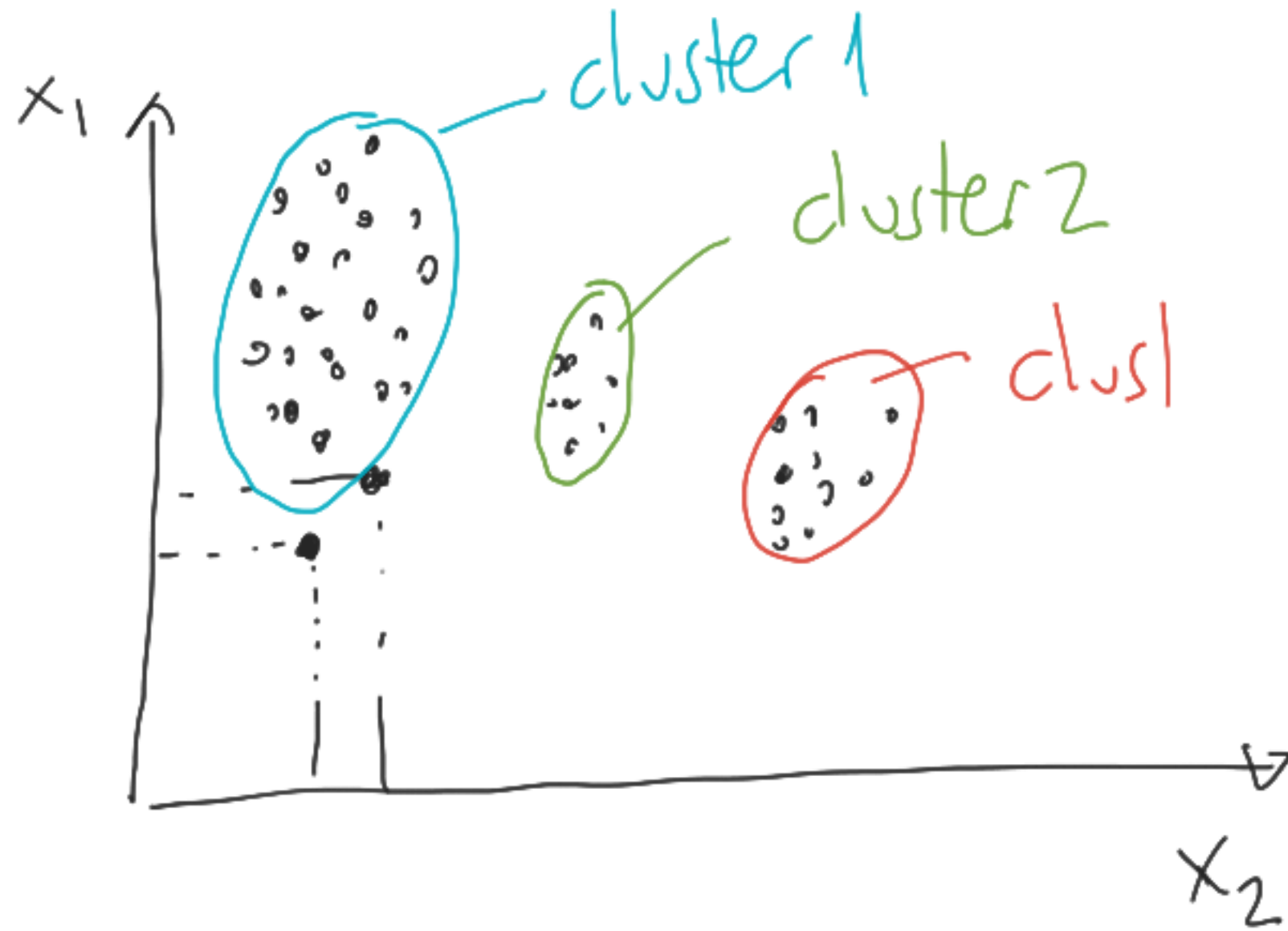
↓
modelo estadístico

similitud - cercanía

CLUSTERING — no supervision

ex: data 2D $\{x_1, x_2\}$

scatterplot



(no need to ~~ethique~~ ~~de classe~~
 $A, w_i, i=1, \dots, N$)

A_i

→ no-supervised

• clustering

• reducción dimensional (PCA)

	x_1	x_2	...	x_M
obs 1	0	0		0
obs 2	0	0		0

...

→ supervised

• clasificación
• regresión

	atributos			
	x_1	x_2	...	x_M
obs 1	0	0		0
obs 2	0	0		0

etiqueta dese

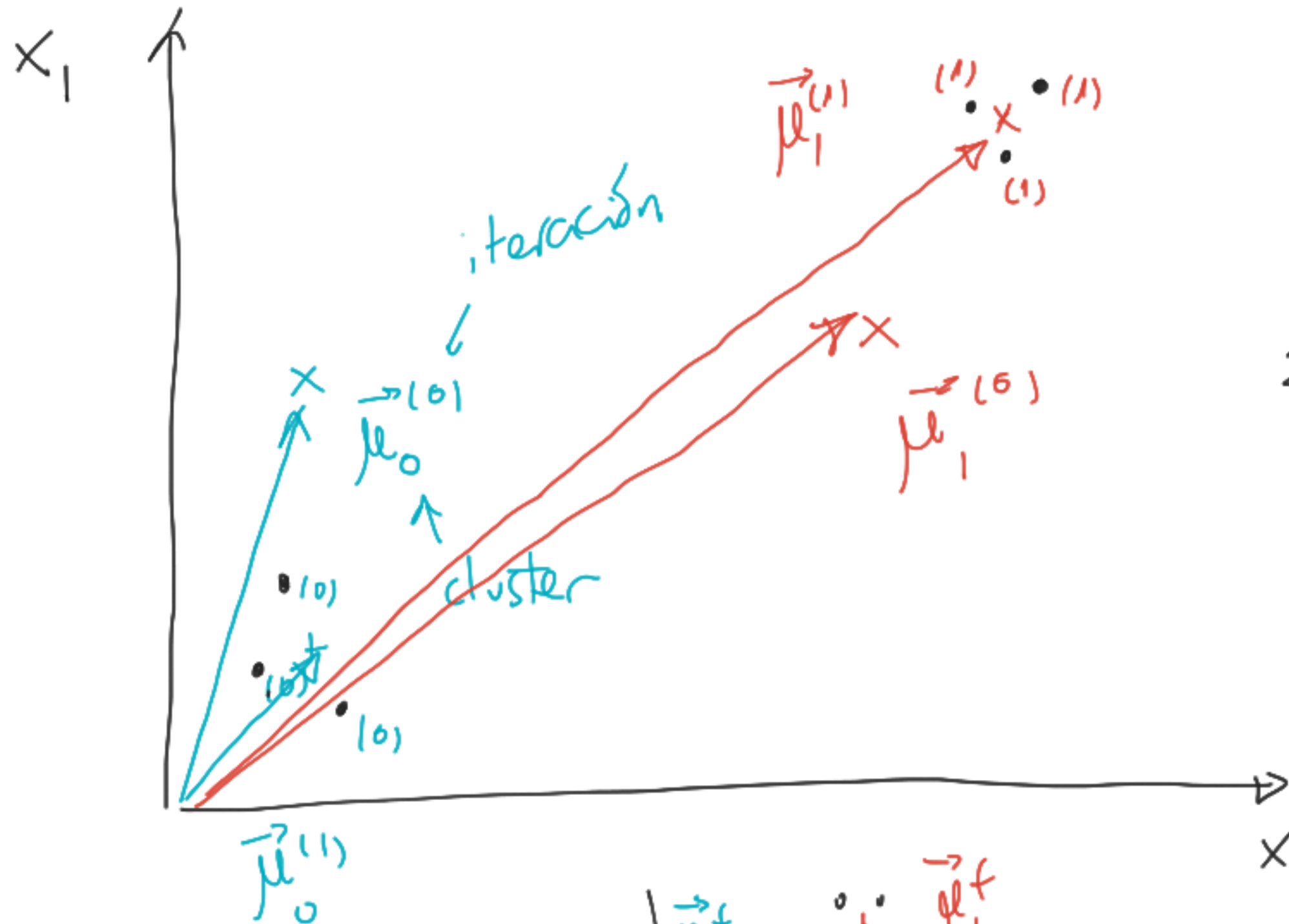
	w
+	'coche'
+	'moto'

iris

pl	pw	sl	sw	class
1.3	2.4	5.1		{0, 1, 2}

K-means

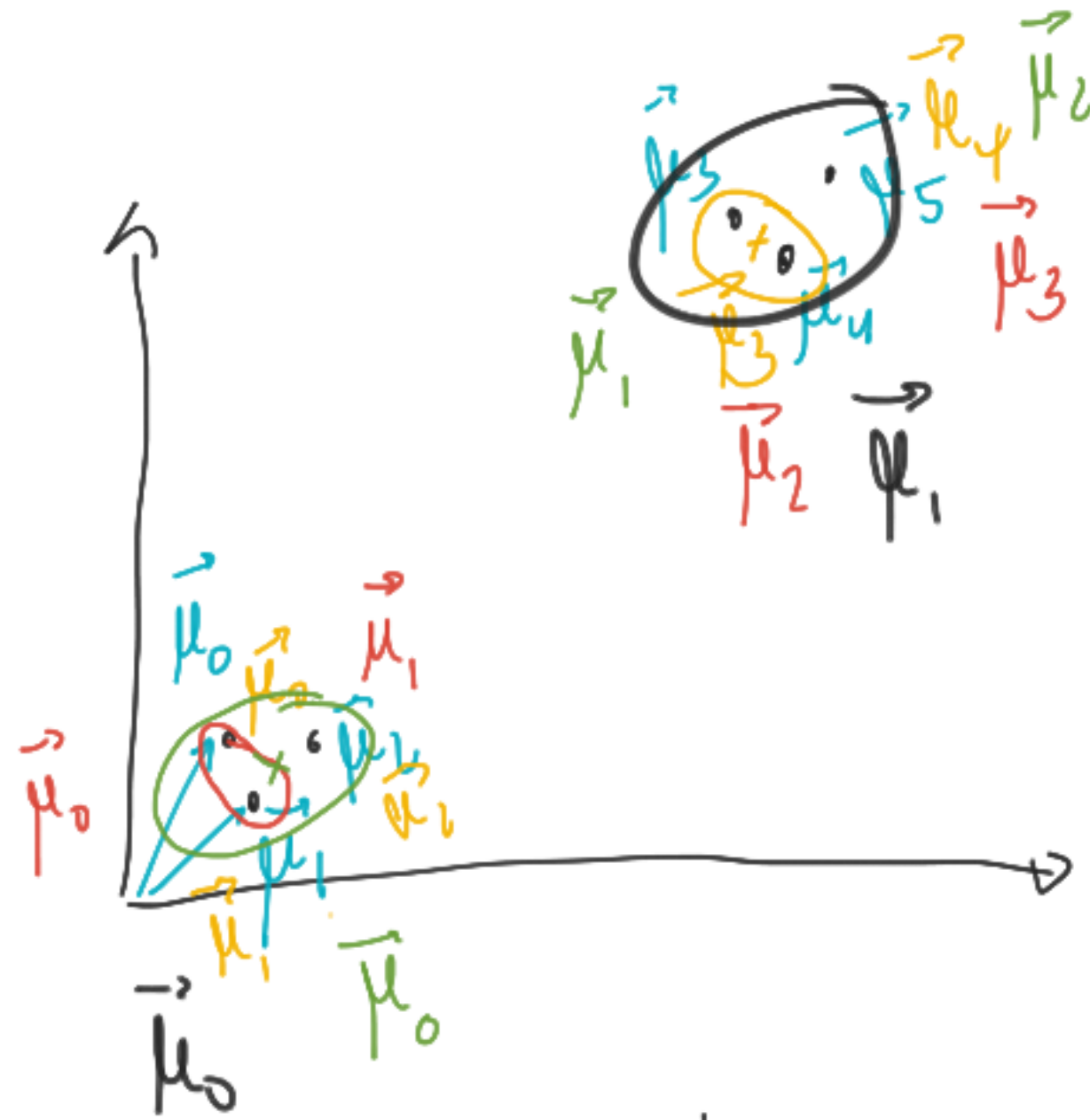
→ algoritmo iterativo basado en distancias entre observaciones



resultado :

1. Establecer el # clusters que se desea identificar
 $N_{CLUSTERS} = 2$
2. Valores iniciales de los centroides (aleatorio)
 $\bar{\mu}_0^{(0)}, \bar{\mu}_1^{(0)}$
3. Calcular distancia entre cada observación y los centroides y asociarla al más cercano
4. Recalcular los centroides
5. Repetir pasos 3 y 4 hasta ^{que} los cent.

aglomerativo jerárquico



→ 1. Indicar #clusters ($N=2$)

2. Inicialización: tanto centroides como observaciones

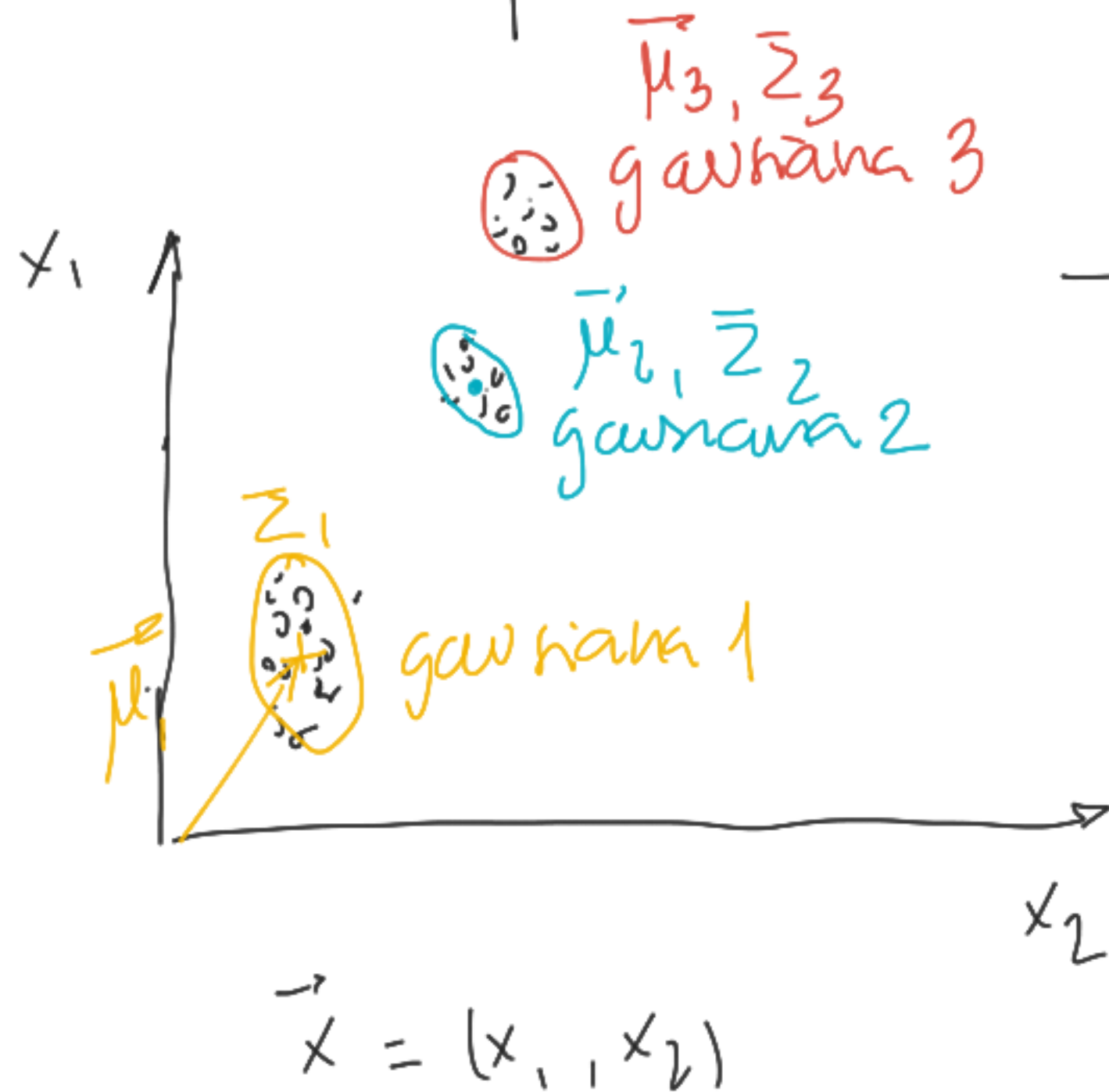
3. Agrupar centroides por distancia

4. Repetir hasta que solo queden $N=2$

DENDROGRAMA!

Mezclas gaussianas (GMM)

modelo probabilístico distribución observaciones
inconveniente: requiere elegir # observaciones



Suma ponderada de Gaussianas

$$P(\vec{x}) = \sum_{i=1}^3 \pi_i \cdot N(\vec{\mu}_i, \Sigma_i) =$$

pesos

$$= \pi_1 \cdot N(\vec{\mu}_1, \Sigma_1) + \pi_2 \cdot N(\vec{\mu}_2, \Sigma_2) + \pi_3 \cdot N(\vec{\mu}_3, \Sigma_3)$$

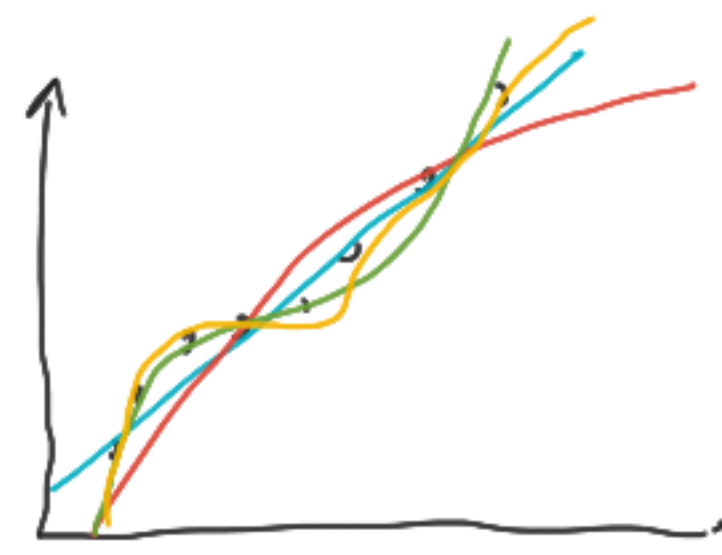
modelo GMM N Gaussianos en d dimensiones

Gaussianos	parámetros gaussianos	pesos
1	$\vec{\mu}_1 (1 \times d)$, $\bar{\Sigma}_1 (d \times d)$	π_1
2	$\vec{\mu}_2 (1 \times d)$, $\bar{\Sigma}_2$	π_2
.	.	.
.	.	.
N	$\vec{\mu}_N$, $\bar{\Sigma}_N$	π_N

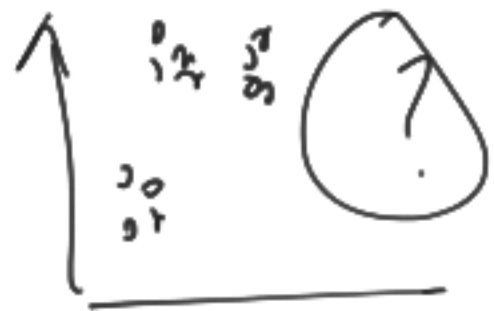
ESTIMAR TODOS LOS PARÁMETROS!

How MANY clusters? → Selección Modelos

GMM + Índice de parsimonia

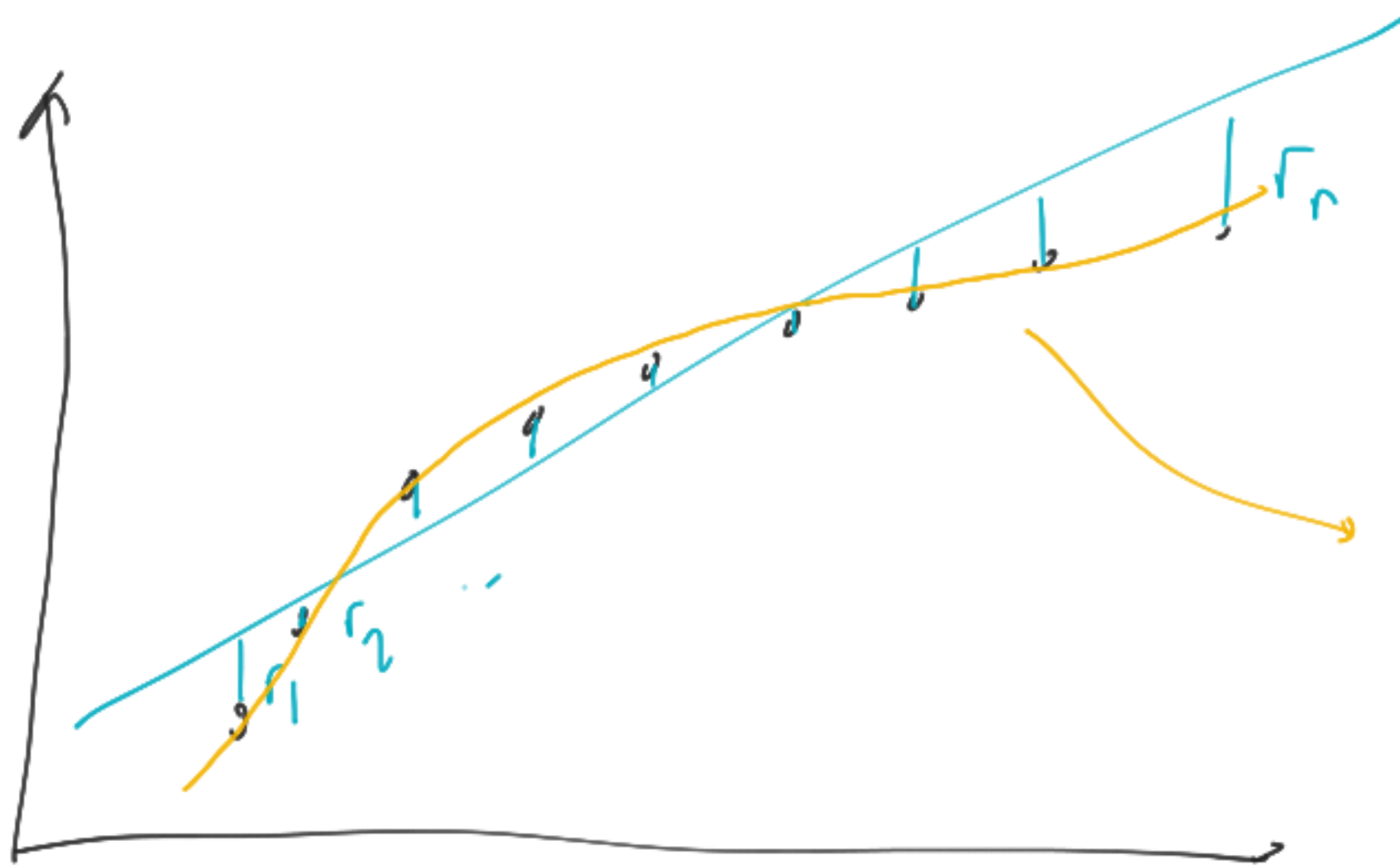


$$BIC = -\text{Error} + \text{Complejidad (Resídeo)} \quad \text{Error vs Complejidad (A parámetros)}$$



		<u>BIC</u>	$y_{\text{model}} - y_{\text{real}}$
GMM	1 gaussian	700	<div style="text-align: center;"> <u> </u> parsimonia </div>
GMM	2 gaussianes	600	
GMM	<u>3 gaussianes</u>	<u>500</u>	
GMM	4 gaussianes	550	

→ menor BIC modelo
óptimo → 3 clusters



$$ax + b \rightarrow (a, b)$$

$$\sum r_i^2 \rightarrow \text{error}$$

$$ax^4 + bx^3 + cx^2 + dx + e$$

$$(a, b, c, d, e)$$