



# Fundamentos de Aprendizaje Automático y Reconocimiento de Patrones

Actividades en clase: Práctico 7

Graciana Castro, Martín Schmidt, Federico Lecumberry, Guillermo Carbajal

Instituto de Ingeniería Eléctrica

2024

# Objetivos

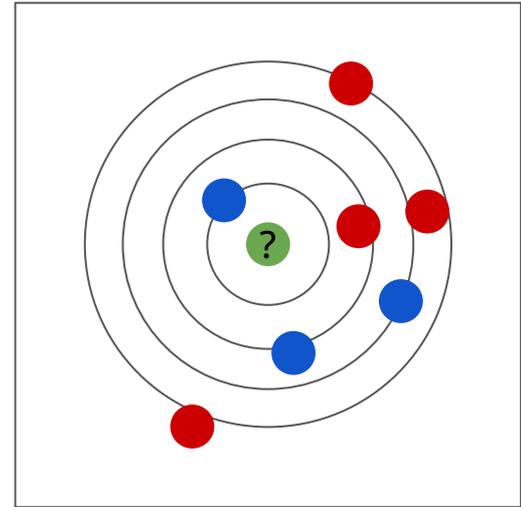
1. Aplicar *k-vecinos más cercanos* para clasificar entre las 10 categorías de dígitos. Optimizar el parámetro  $k$  utilizando la biblioteca *scikit-learn*.
2. Estimar densidades de probabilidad mediante el método de *Ventanas de Parzen*.
3. Implementar el algoritmo de clustering *k-means* y aplicarlo en datos sintéticos. Analizar su funcionamiento.
4. Realizar *agrupamiento de datos* mediante *Mezcla de Gaussianas*. Implementarlo utilizando el esquema *Expectation-Maximization* para encontrar los parámetros. Comparar este agrupamiento con el de *k-means*.

# Tabla de contenido

- Ejercicio 1: k-NN
- Ejercicio 2: ventanas de Parzen
- Ejercicio 3: k-means
- Ejercicio 4: mezcla de Gaussianas

# Ejercicio 1: k-NN

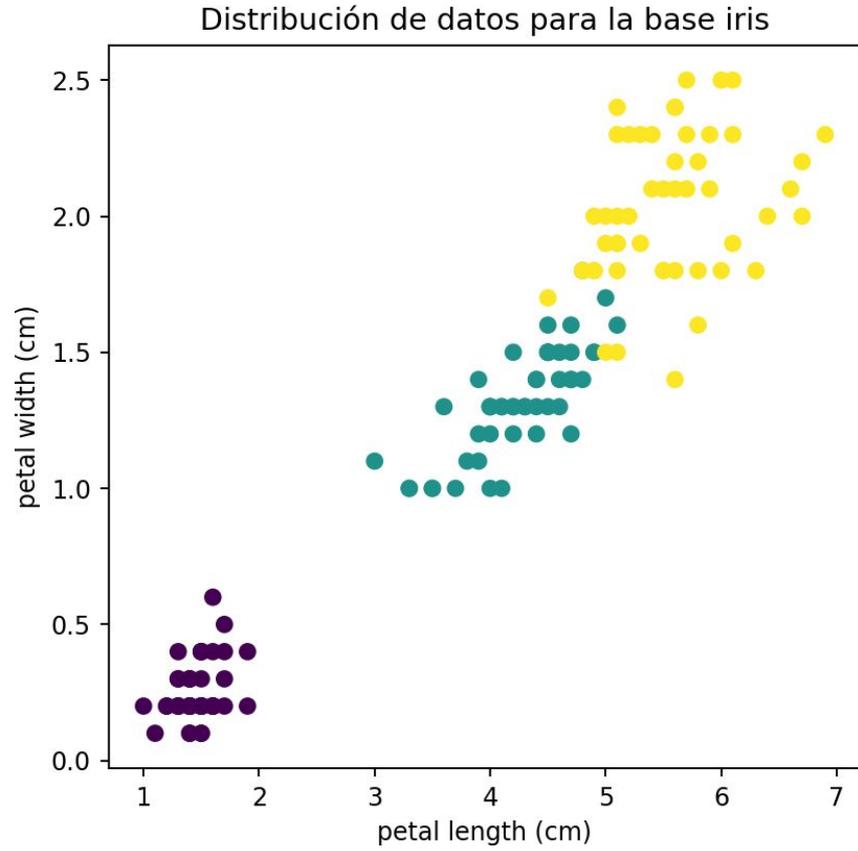
- Clasificación de dígitos utilizando  $k$ -NN
- Optimización del hyperparámetro  $k$ 
  - mediante conjunto de validación
  - mediante validación cruzada



$$k = 1, 3, 5, 7$$

## Ejercicio 2: ventanas de Parzen

- Datos: Flores de la base Iris



## Ejercicio 2: ventanas de Parzen

- Datos: Flores de la base Iris
- Estimación de densidades

$$P(\mathbf{x}) = \frac{1}{M} \sum_{i=1}^M K(\mathbf{x}, \mathbf{z}_i)$$

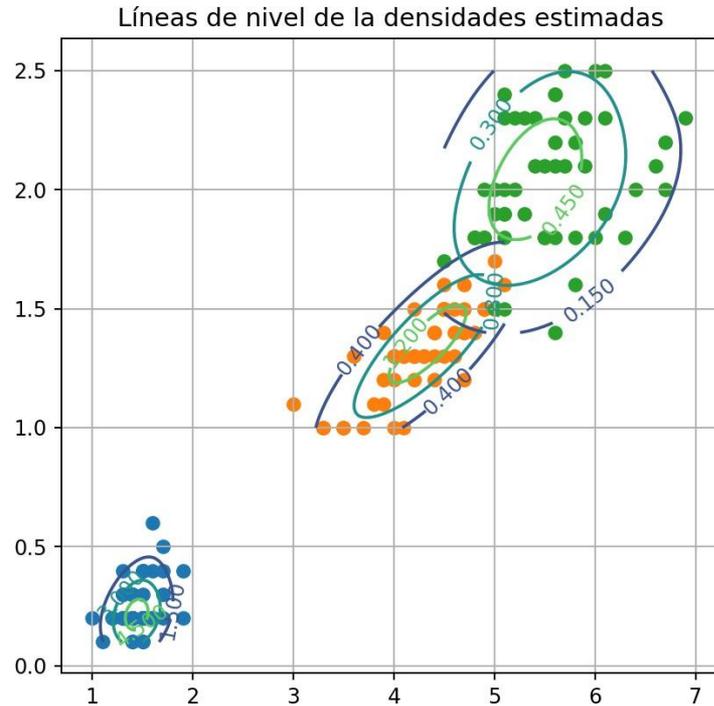
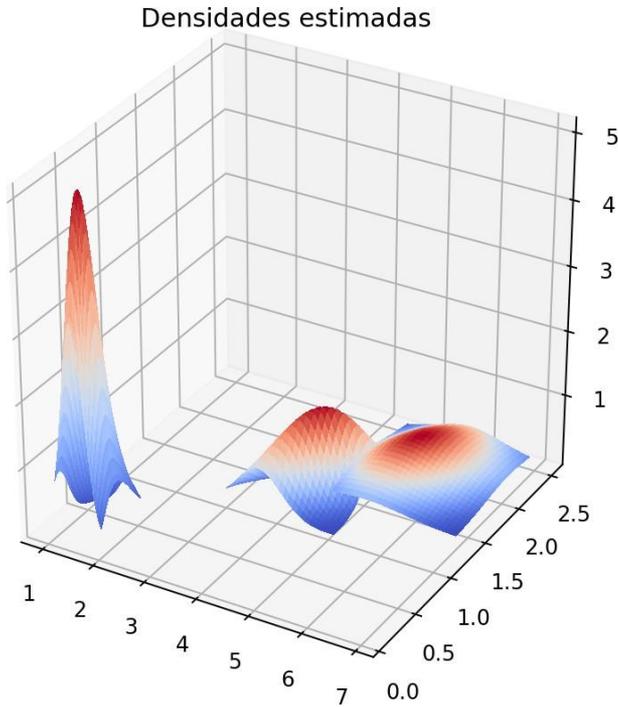
- Se utilizará un kernel gaussiano

$$K(\mathbf{x}, \mathbf{z}_i) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left( -\frac{1}{2} (\mathbf{x} - \mathbf{z}_i)^T \Sigma^{-1} (\mathbf{x} - \mathbf{z}_i) \right)$$

con  $\Sigma$  de la forma  $\Sigma = r^2 \Sigma_c$  siendo  $\Sigma_c$  la matriz de covarianza de los datos.

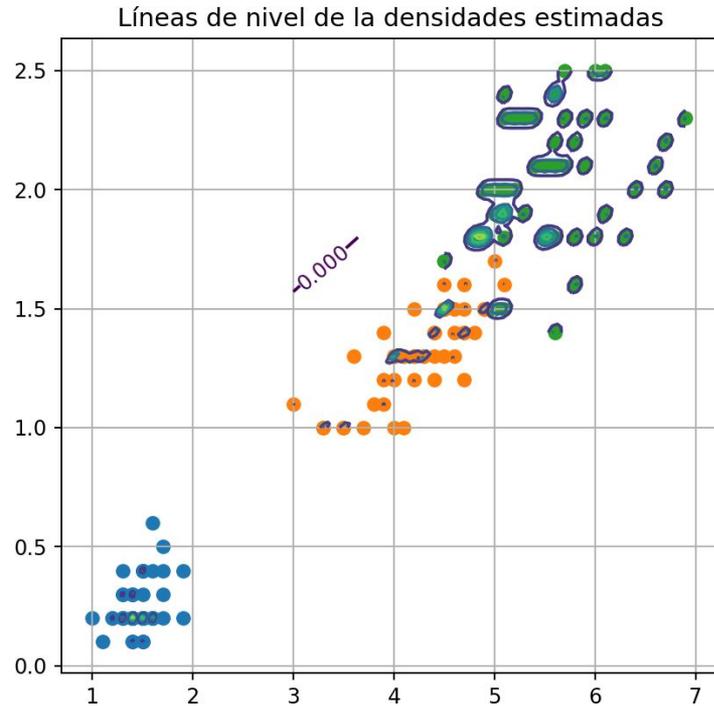
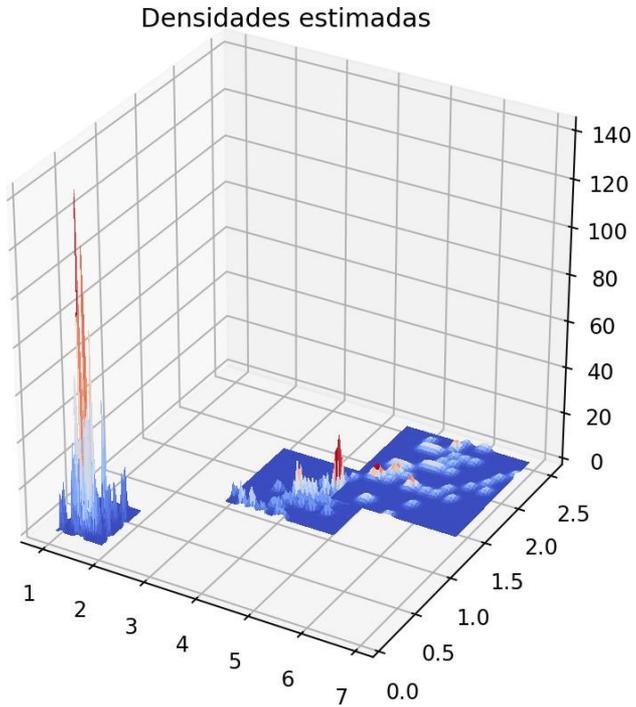
# Ejercicio 2: ventanas de Parzen

Efecto del ancho del kernel  $r = 1$



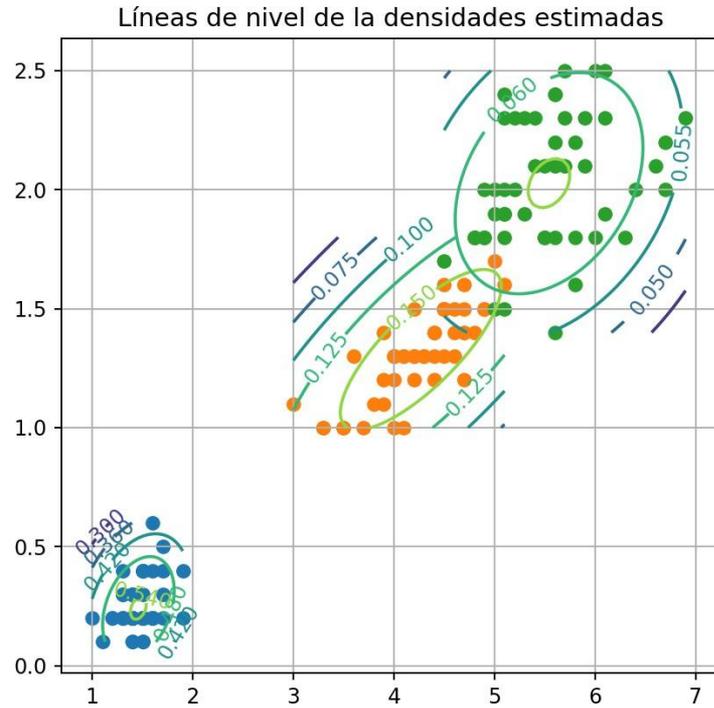
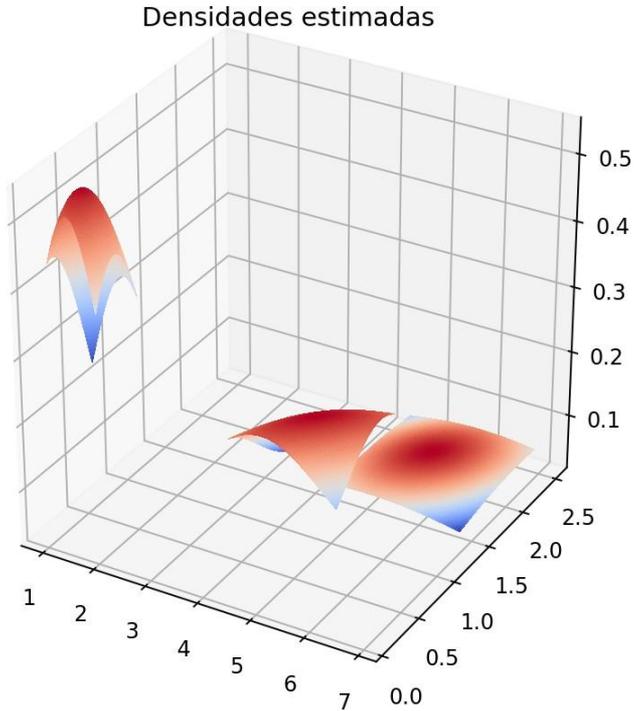
## Ejercicio 2: ventanas de Parzen

Efecto del ancho del kernel  $r = 0.1$



## Ejercicio 2: ventanas de Parzen

Efecto del ancho del kernel  $r = 4$



## Ejercicio 3: K-means

- Técnica de agrupamiento de datos.
- Objetivo: Dividir los datos de entrada  $x_1, \dots, x_N$  en  $K$  subconjuntos  $S_1, \dots, S_K$  y elegir centros  $\mu_1, \dots, \mu_K$  para cada cluster.
- Se busca minimizar

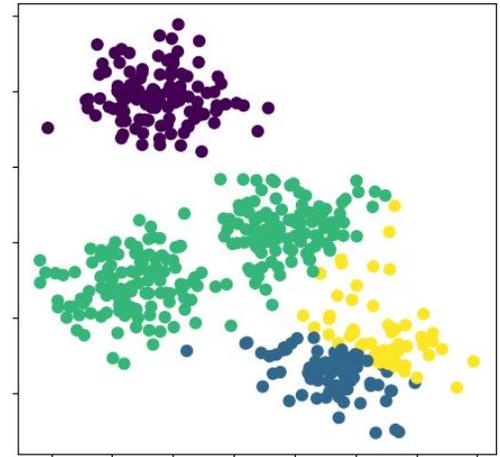
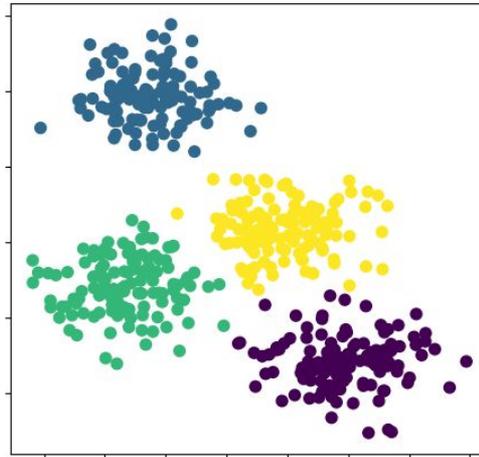
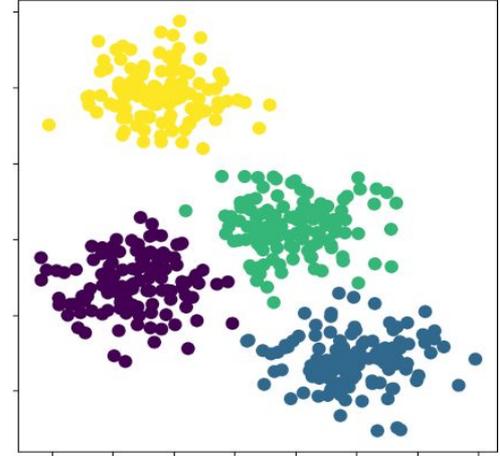
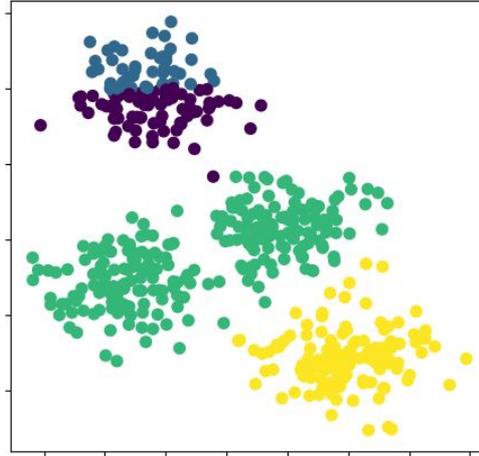
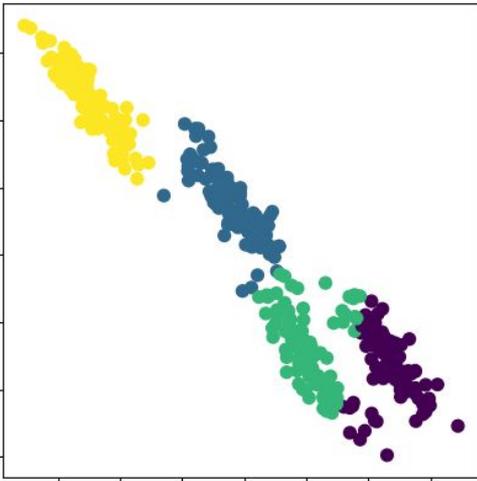
$$E_{\text{in}}(S_1, \dots, S_K, \mu_1, \dots, \mu_K) = \sum_{k=1}^K \sum_{x_n \in S_k} \|x_n - \mu_k\|^2$$

El algoritmo de Lloyd consta de los siguientes pasos:

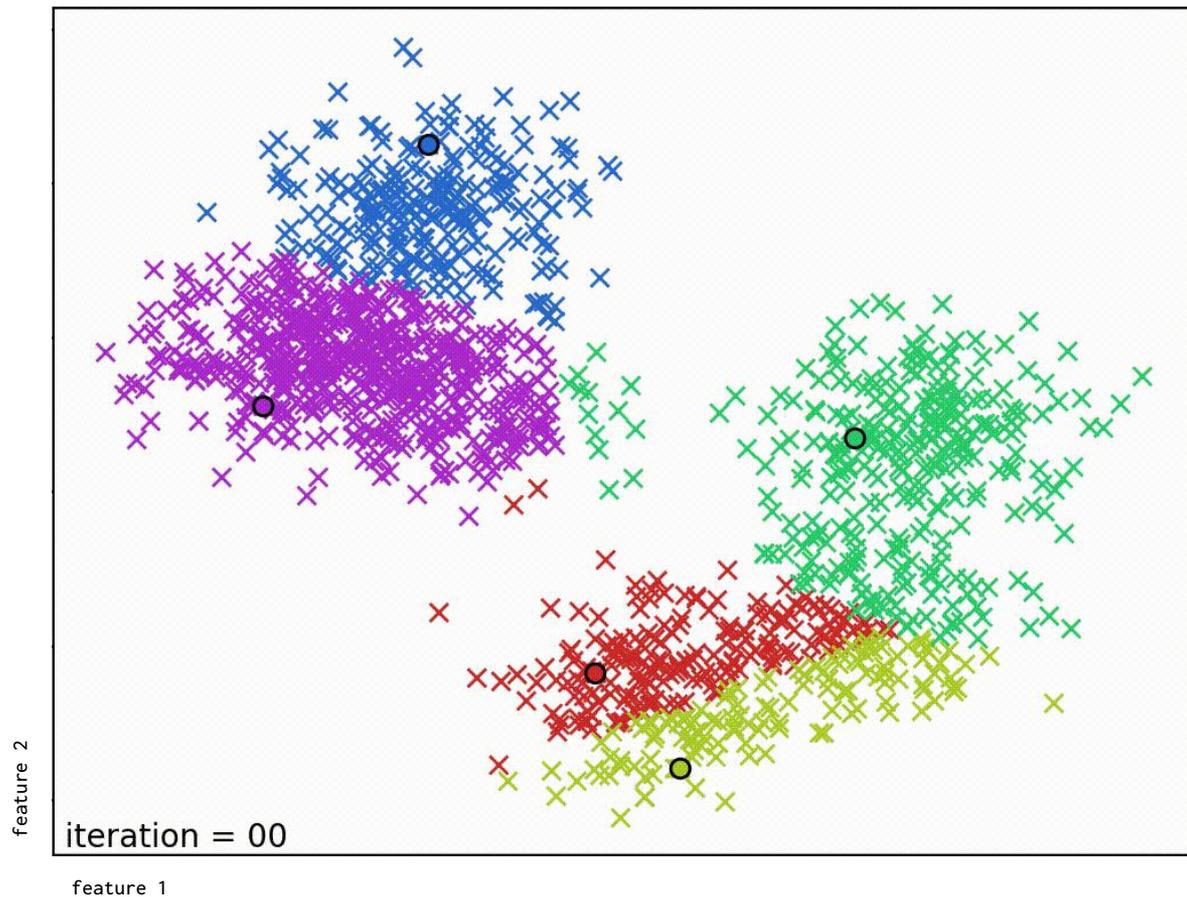
1. Se inicializan los centros  $\mu_j$  de los clusters
2. Se construye el cluster  $S_j$  con todos los puntos cuyo centro más cercano es  $\mu_j$
3. Se actualizan los centros  $\mu_j$  al centroide de los puntos pertenecientes a  $S_j$
4. Se repiten los pasos anteriores hasta que el algoritmo converge

## Ejercicio 3: K-means

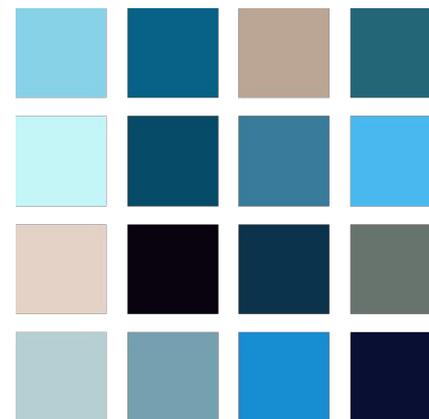
Resultados variando la inicialización.



# Ejercicio 3: K-means



↓ 16 color centroids



# Color quantization (~K-means)

Color Inspector 3D (v2.5) suarez-color.jpg

Color Space: HSL    Display Mode: Wu Quant    # Colors: 2    LUT

262144 Pixels, 2 Colors

Depth

Perspect...

Scale

H (x 1.0)    S (x 1.0)    L (x 1.0)    Brightness (+0)    Contrast (x 1.0)    Saturation (x 1.0)    Color Rotation (0°)

# EM para mezcla de gaussianas

El algoritmo EM para encontrar los parámetros del modelo GMM consta de los siguientes pasos:

1. Se inicializan  $\boldsymbol{\mu}_j$ ,  $\boldsymbol{\Sigma}_j$  y  $w$ .
2. Para cada muestra  $\mathbf{x}_n$  se calcula la probabilidad  $\gamma_{nj}$  de pertenencia a cada *cluster*  $j$ .

$$\gamma_{nj} = \frac{w_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\sum_{l=1}^L w_l \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_l, \boldsymbol{\Sigma}_l)}$$

3. Se actualizan los centros  $\boldsymbol{\mu}_j$ ,  $\boldsymbol{\Sigma}_j$  y  $w$ :

$$N_j = \sum_{n=1}^N \gamma_{nj} \quad \boldsymbol{\mu}_j^{new} = \frac{1}{N_j} \sum_{n=1}^N \gamma_{nj} \mathbf{x}_n$$

$$\boldsymbol{\Sigma}_j^{new} = \frac{1}{N_j} \sum_{n=1}^N \gamma_{nj} (\mathbf{x}_n - \boldsymbol{\mu}_j) (\mathbf{x}_n - \boldsymbol{\mu}_j)^T$$

$$w_j^{new} = \frac{N_j}{N}$$

4. Se repiten los pasos anteriores hasta que el algoritmo converge.

## Ejercicio 4: mezcla de Gaussianas

Resultados variando la inicialización.

