



# ***Aprendizaje no supervisado y análisis de agrupamientos***

Luis Vázquez

GTI - IIE

Facultad de Ingeniería

Universidad de la República

# Aprendizaje no supervisado

- ⑥ Partimos de conjunto de entrenamiento  $\{x_j; j = 1, 2, \dots, N\}$  tq. no conocemos sus etiquetas de clase  $\gamma_i$ .
- ⑥ Cuando:
  - △ no disponemos del conocimiento de un experto
  - △ el etiquetado de cada muestra individual es impracticable
- ⑥ **Ejemplo:** Aplicaciones con sensores remotos, como imágenes satelitales de terrenos.

Costoso o imposible recoger información real del tipo de suelo sensado en cada punto.

# Heurística del aprendizaje

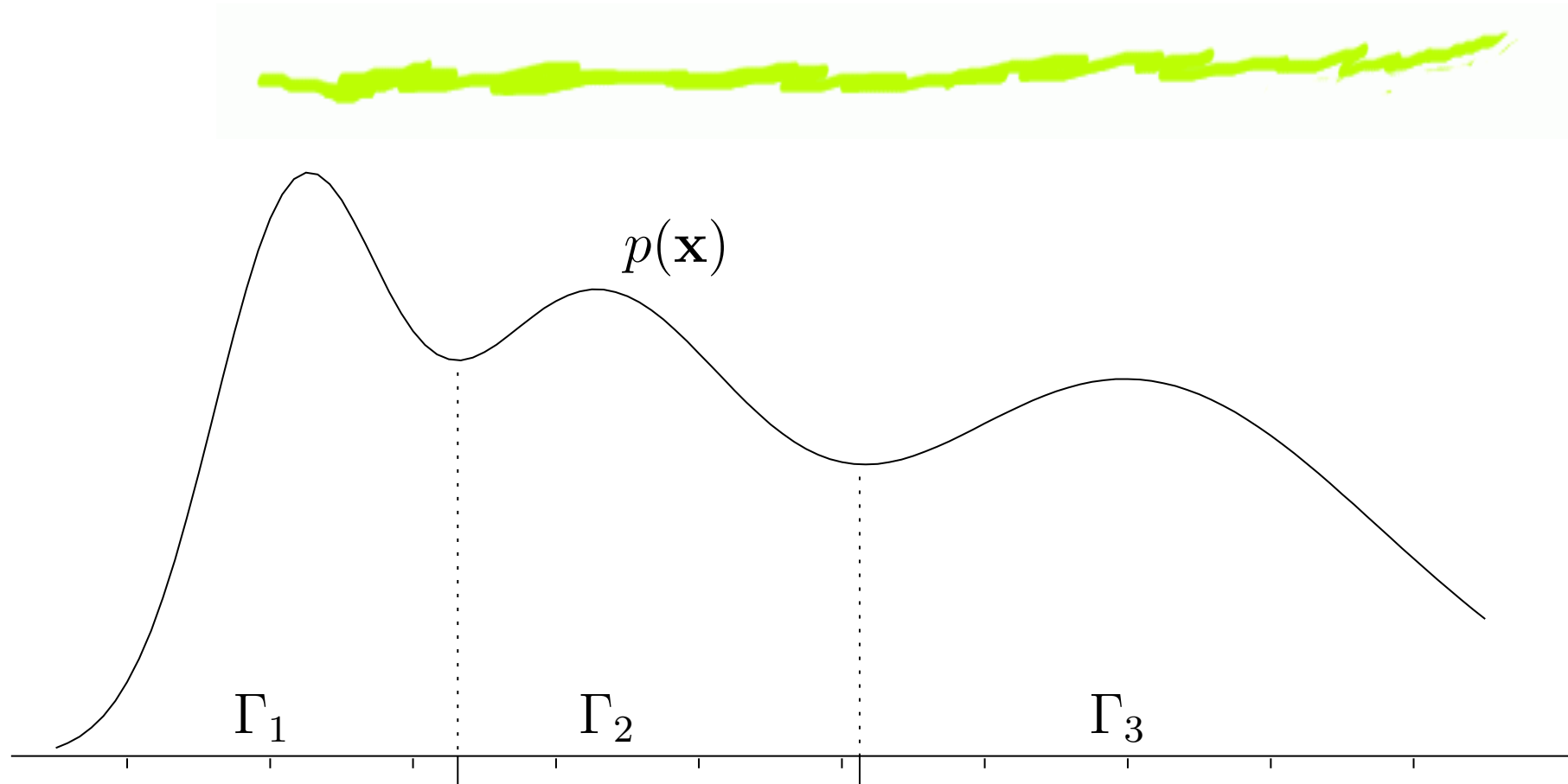
- ⑥ Dado conjunto de entrenamiento suficientemente grande podríamos inferir

$$p(\mathbf{x}) = \sum_{i=1}^m P(\omega_i) p(\mathbf{x}|\omega_i)$$

- ⑥ Si la densidad conjunta es multimodal cada modo debería corresponder a cada clase presente.
- ⑥ Si cada clase fuera normal se podría recuperar parámetros de cada distribución y luego seguir con el diseño del clasificador ...
- ⑥ **Problema práctico:** estimar y analizar densidad conjunta en espacio de dimensión  $n$ :

Impracticable por su complejidad computacional.

# *Distribución conjunta multimodal*



**Figure 1: Regiones asociadas a cada clase**

# Heurística del aprendizaje II

- ⑥ Opción alternativa: clasificar los patrones en el conjunto de entrenamiento y usar estas etiquetas para un aprendizaje supervisado
- ⑥ Solo necesitamos un método indirecto que permita etiquetar automáticamente los patrones de entrenamiento.
- ⑥ Queremos forma de particionar el conjunto en clases con misma etiqueta

**métodos de agrupamiento o *clustering***

# Análisis de Agrupamientos

- ⑥ Nos interesan las modas en la función de densidad conjunta  $p(\mathbf{x})$
- ⑥ **Intuitivamente:** estarán asociadas a regiones con alta densidad en el espacio de observación.
- ⑥ Propósito de las técnicas de agrupamiento:
  - △ detectar y agrupar **enjambres de puntos**
  - △ analizar y extraer la estructura presente en el conjunto de muestras de entrenamiento.
- ⑥ **Conjunto de datos bien estructurado:** contiene varias regiones de alta densidad separadas por otras relativamente vacías o con poca densidad.

# Estructura de datos

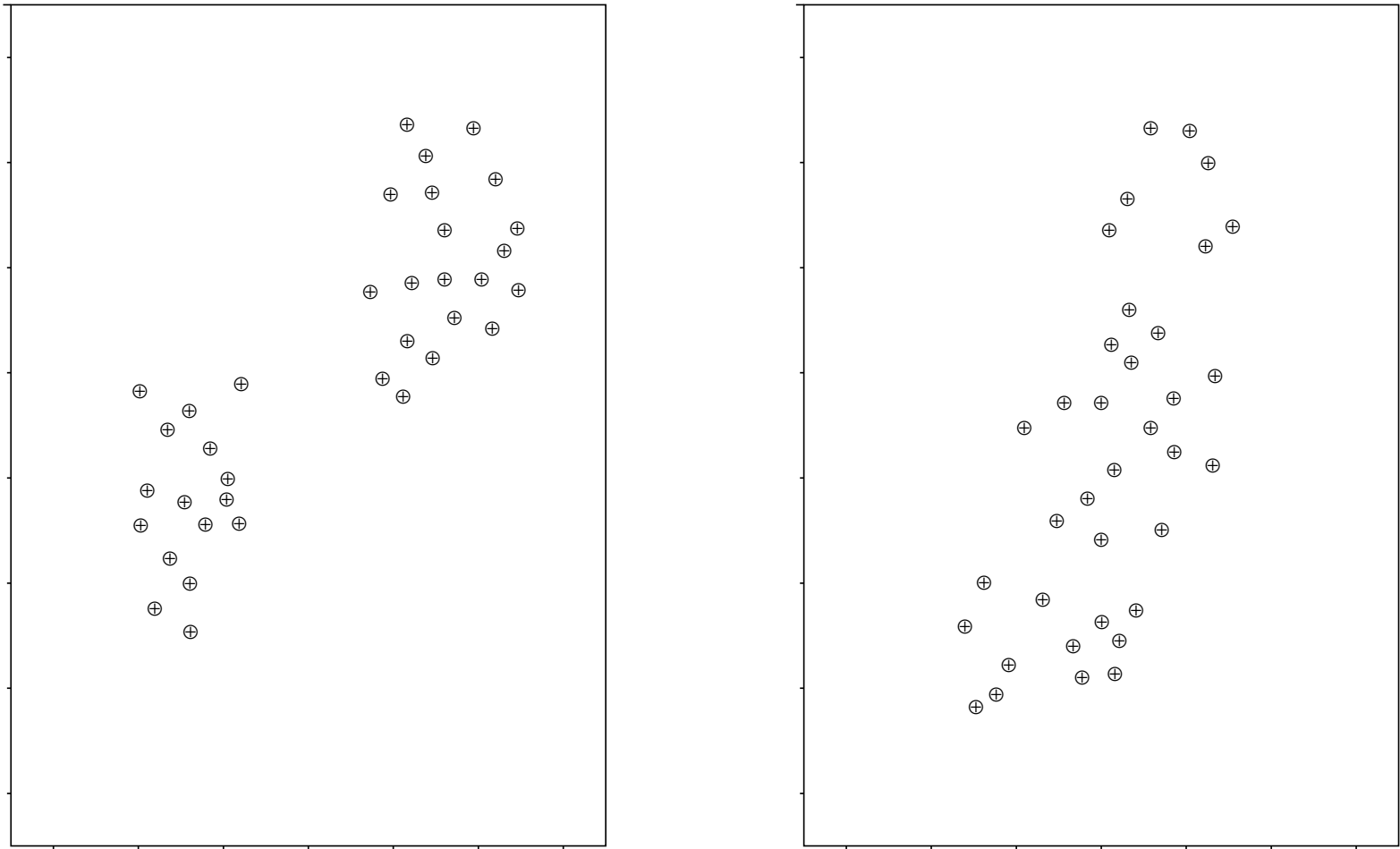


Figure 2: Datos estructurados y no estructurados

# Medidas de Similitud

- ⑥ Para decidir si  $\mathbf{x}$  pertenece a un agrupamiento necesitamos **medida de proximidad o similitud**.
- ⑥ Se han sugerido gran número de tales medidas pero las más usadas son medidas de distancia, en particular la distancia Euclídeana.
- ⑥ Medir afinidad de un punto a un agrupamiento:
  - △ similitud con otros puntos en el agrupamiento
  - △ similitud con modelo definido para el agrupamiento
- ⑥ **Ejemplo:** representar grupo  $i$  por su media  $\mu_i$ . La afinidad entre un punto y el grupo podría ser:

$$d(\mathbf{x}, \mu_i) = [(\mathbf{x} - \mu_i)^T (\mathbf{x} - \mu_i)]$$



# Criterios de Agrupamiento

- ⑥ Problema: particionar un conjunto de puntos en agrupamientos *óptimos*
- ⑥ Necesitamos algún **criterio de agrupamiento**
- ⑥ Permite definir *cuantitativamente* cuando una partición es mejor que otra.
- ⑥ El criterio de agrupamiento que definamos y el algoritmo de agrupamiento asociado estarán relacionados con la medida de similitud usada.

# Algoritmos de Agrupamiento

- Conocidos medida de similitud y criterios encontrar particion *óptima*  $\Rightarrow$  **problema de optimización!**
- Problema práctico: considerar todas las particiones

$$\frac{k^M}{k!} \quad \text{siendo} \quad \begin{cases} k & \text{nro. de grupos} \\ M & \text{nro. de patrones} \end{cases}$$

**Inviabile en la práctica!!**

$$M = 100, k = 5 \quad \Rightarrow \quad 10^{67}!!$$

- Se deben buscar métodos no exhaustivos, aunque no siempre garanticen óptimos globales.

# Clasificación de Algoritmos

- ⑥ Por estrategia de agrupamiento
  - △ Aglomerativos o incrementales: **bottom up**
  - △ Divisivos o decrementales: **top down**
  - △ Mixtos: **creación y mezcla**
- ⑥ Por criterio de búsqueda del óptimo
  - △ Heurísticos: no se optimiza la función criterio
  - △ Por Optimización: se usa una función criterio
- ⑥ Por información a priori
  - △ Número de clases conocido
  - △ Número de clases desconocido

# Medidas Objetivas Importantes

## ⑥ *Para mezclar*

- △ Distancias entre los centros de los agrupamientos

## ⑥ *Para dividir*

- △ Varianza interna en un agrupamiento
- △ Distancia entre puntos extremos de un agrupamiento

# Algoritmo de $k$ -medias

- ⑥ El conjunto de datos  $X$  contiene  $k$  agrupamientos  $X_i$  que pueden representarse adecuadamente con su valor medio  $\mu_i$ .
  - △ *Medida de similitud*: distancia Euclídeana a  $\mu_i$
  - △ *Criterio de agrupamiento*: suma total de la distancia cuadrática de cada punto al vector medio de su agrupamiento.
  
- ⑥ **Objetivo del algoritmo:**

Encontrar entre todas las particiones de  $X$  en  $k$  conjuntos  $\{X_i ; i = 1, 2, \dots, k\}$  aquella que minimiza el criterio de agrupamiento.

# Algoritmo de $k$ -medias:

## Planteo Formal

- ⑥ Encontrar los agrupamientos  $\{\mathbf{X}_i\}$  que minimizan

$$J = \sum_{i=1}^k J_i = \sum_{i=1}^k \sum_{j=1}^{N_i} d(\mathbf{x}_{ij}, \mu_i) \quad \mathbf{x}_{ij} \in \mathbf{X}_i, \quad N_i = \#\mathbf{X}_i$$

entre todas las posibles  $k$ -particiones de  $\mathbf{X}$ .

- ⑥ Considero el efecto de un cambio atómico en la configuración de agrupamientos:

*sacar  $\mathbf{x}$  de  $\mathbf{X}_l$  y pasarlo a otro grupo  $\mathbf{X}_r$ .*

- ⑥ Esta reasignación solo afecta los grupos  $l$  y  $r$  con:

$$\bar{\mu}_l = \mu_l + \frac{1}{N_l - 1}(\mu_l - \mathbf{x}) \quad \bar{\mu}_r = \mu_r - \frac{1}{N_r + 1}(\mu_r - \mathbf{x})$$

# Deducción

- Calculamos el valor medio de  $X_i$  antes y después de la reasignación

$$\mu_l = \frac{1}{N_l} \sum_{j=1}^{N_l} \mathbf{x}_j \quad \bar{\mu}_l = \frac{1}{N_l - 1} \left( \sum_{j=1}^{N_l} \mathbf{x}_j - \mathbf{x} \right)$$

- De aquí resulta

$$\begin{aligned} (N_l - 1)\bar{\mu}_l &= N_l \mu_l - \mathbf{x} \Rightarrow \bar{\mu}_l = \frac{N_l}{N_l - 1} \mu_l - \frac{1}{N_l - 1} \mathbf{x} \\ &\Rightarrow \bar{\mu}_l = \mu_l + \frac{1}{N_l - 1} (\mu_l - \mathbf{x}) \end{aligned}$$

- Análogamente se verifica la segunda identidad.

# Deducción de cambio atómico global

- Para hallar el cambio global en  $J$  basta calcular los cambios en  $J_l$  y  $J_r$
- Para el nuevo agrupamiento  $l$ -ésimo tengo

$$\begin{aligned}\bar{J}_l &= \sum_{j=1}^{N_l-1} d(\mathbf{x}_j, \bar{\boldsymbol{\mu}}_l) = \sum_{j=1}^{N_l-1} (\mathbf{x}_j - \bar{\boldsymbol{\mu}}_l)^T (\mathbf{x}_j - \bar{\boldsymbol{\mu}}_l) = \\ &= \sum_{j=1}^{N_l} \left( \mathbf{x}_j - \boldsymbol{\mu}_l + \frac{\boldsymbol{\mu}_l - \mathbf{x}}{N_l - 1} \right)^T \left( \mathbf{x}_j - \boldsymbol{\mu}_l + \frac{\boldsymbol{\mu}_l - \mathbf{x}}{N_l - 1} \right) - \left\| \mathbf{x} - \boldsymbol{\mu}_l + \frac{\boldsymbol{\mu}_l - \mathbf{x}}{N_l - 1} \right\|^2 \\ &= J_l - \frac{2}{N_l - 1} (\boldsymbol{\mu}_l - \mathbf{x}) \underbrace{\sum_{j=1}^{N_l} (\mathbf{x}_j - \boldsymbol{\mu}_l)}_0 + \frac{N_l - N_l^2}{(N_l - 1)^2} (\boldsymbol{\mu}_l - \mathbf{x})^T (\boldsymbol{\mu}_l - \mathbf{x})\end{aligned}$$



# Cambio global en reasignación

- 6 Luego de agrupar resulta:

$$\begin{cases} \bar{J}_l = J_l - \frac{N_l}{N_l - 1} (\boldsymbol{\mu}_l - \mathbf{x})^T (\boldsymbol{\mu}_l - \mathbf{x}) = J_l - \frac{N_l}{N_l - 1} d(\mathbf{x}, \boldsymbol{\mu}_l) \\ \bar{J}_r = J_r + \frac{N_r}{N_r - 1} (\boldsymbol{\mu}_r - \mathbf{x})^T (\boldsymbol{\mu}_r - \mathbf{x}) = J_r + \frac{N_r}{N_r - 1} d(\mathbf{x}, \boldsymbol{\mu}_r) \end{cases}$$

- 6 Obtenemos descenso en  $J \Leftrightarrow (J_l - \bar{J}_l) > (\bar{J}_r - J_r)$

$$\Leftrightarrow \frac{N_l}{N_l - 1} d(\mathbf{x}, \boldsymbol{\mu}_l) > \frac{N_r}{N_r - 1} d(\mathbf{x}, \boldsymbol{\mu}_r)$$

- 6 Si  $N_l$  y  $N_r$  son grandes  $\Rightarrow$  esta desigualdad se cumple  
 $\Leftrightarrow \mathbf{x}$  esta más cerca de  $\boldsymbol{\mu}_r$  que de  $\boldsymbol{\mu}_l$ .

# El algoritmo $k$ -mean

De lo anterior se deducen los pasos del algoritmo

- ⑥ **Paso 1:** Elegir una partición inicial en  $k$  grupos al azar
- ⑥ **Paso 2:** Calcular las medias  $\mu_i$  de los *clusters*
- ⑥ **Paso 3:** Seleccionar secuencialmente un punto  $\mathbf{x}$  del conjunto y, si corresponde, **reasignarlo al *cluster* que minimiza  $d(\mathbf{x}, \mu_i)$**
- ⑥ **Paso 4:** Si no hay reasignaciones en todo el conjunto terminar; sino volver al **Paso 2**.

# Algoritmos de agrupamiento dinámico

- El  $k$ -mean es caso particular de una familia de algoritmos

## *Dynamical Clustering Algorithms*

- El grupo  $\mathbf{X}_i$  se representa por un núcleo  $K_i(\mathbf{X}_i)$ :  
Puede ser un una función paramétrica, un conjunto de puntos u otro modelo
- Se define una medida de similitud  $\Delta(\mathbf{x}, K_i)$  y el criterio de agrupamiento

$$J = \sum_{i=1}^k \sum_{j=1}^{N_i} \Delta(\mathbf{x}_{ji}, K_i) \quad \mathbf{x}_{ji} \in \mathbf{X}_i$$

# Algoritmo de agrupamiento genérico

Por analogía el algoritmo de agrupamiento dinámico genérico sería

- ⑥ **Paso 1:** Elegir una partición inicial en  $k$  grupos al azar
- ⑥ **Paso 2:** Calcular el modelo  $K_i(\mathbf{X}_i)$  para cada *cluster*
- ⑥ **Paso 3:** Seleccionar uno a uno cada punto  $x$  del conjunto  $y$ , si corresponde, **reasignarlo al *cluster* con mayor afinidad**
- ⑥ **Paso 4:** Si no hay reasignaciones terminar el algoritmo; sino volver al **Paso 2**.

# Ejemplo de Modelo Gaussiano

- Un ejemplo importante de representación de un cluster es el *modelo gaussiano*

$$K_i(\mathbf{x}, \mathbf{X}_i) = \frac{1}{(2\pi)^{\frac{d}{2}} |\boldsymbol{\Sigma}_i|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i)\right) \cdot \frac{N_i}{N}$$

- La medida de similitud asociada es

$$\Delta(\mathbf{x}_{ji}, K_i) = (\mathbf{x} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) + \log |\boldsymbol{\Sigma}_i| - 2 \log N_i$$

- Este modelo se usa para recuperar las componentes de una mezcla de distribuciones normales.