

Introducción a la Inteligencia Artificial
Clase 7



Clase 7

1. Motivación
 - a. Aprendizaje No supervisado
 - b. Aplicaciones
2. Gaussian Mixture Models (Además ya vimos otros solo que no le pusimos etiqueta)
 - a. Aplicaciones
 - b. Formulación

Aprendizaje no supervisado

Aca puedo definir
un concepto de
error (error de
inferencia)

clusterización - K means
red. dim - PCA

Machine Learning Supervisado	Machine Learning no Supervisado
Proceso aleatorio \bar{X}, y	Proceso aleatorio \bar{X}
$\hat{f}_{y/\bar{x}}(y \bar{x})?$ \longrightarrow Bayes y M.V.	$\hat{f}_{\bar{x}}(\bar{x})?$ \longrightarrow Bayes y M.V.
Inferencias, predicciones	Clusterización, Reducción Dimensionalidad

\bar{X} datos

y label/regresada

$y \in \mathbb{R} \rightarrow$ regresión
 $y \in \mathbb{K} \rightarrow$ clasif.

$L(y, \hat{y})$

segmentación

tenemos métricas de relajación/equilibrio,
métricas de desigualdad, etc

en Kmeans $U = \sigma_{is} / \sigma_{eg}$

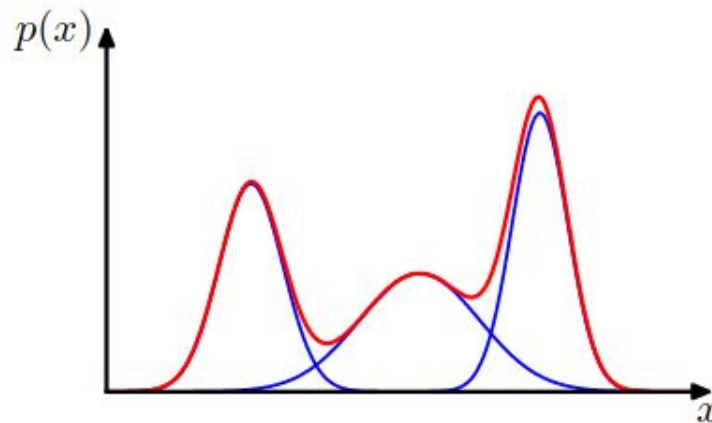
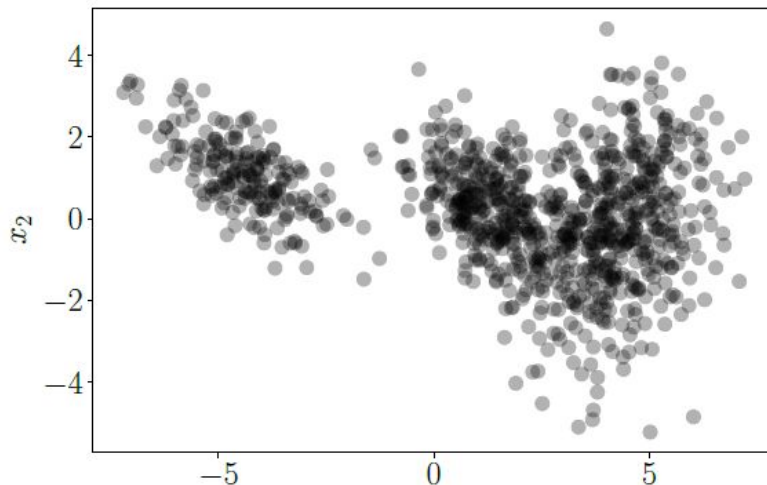
$U_1, U_2, U_3, \dots, U_n$

Aplicaciones Generales

- Data Mining
- Pattern Recognition
- Statistical Analysis

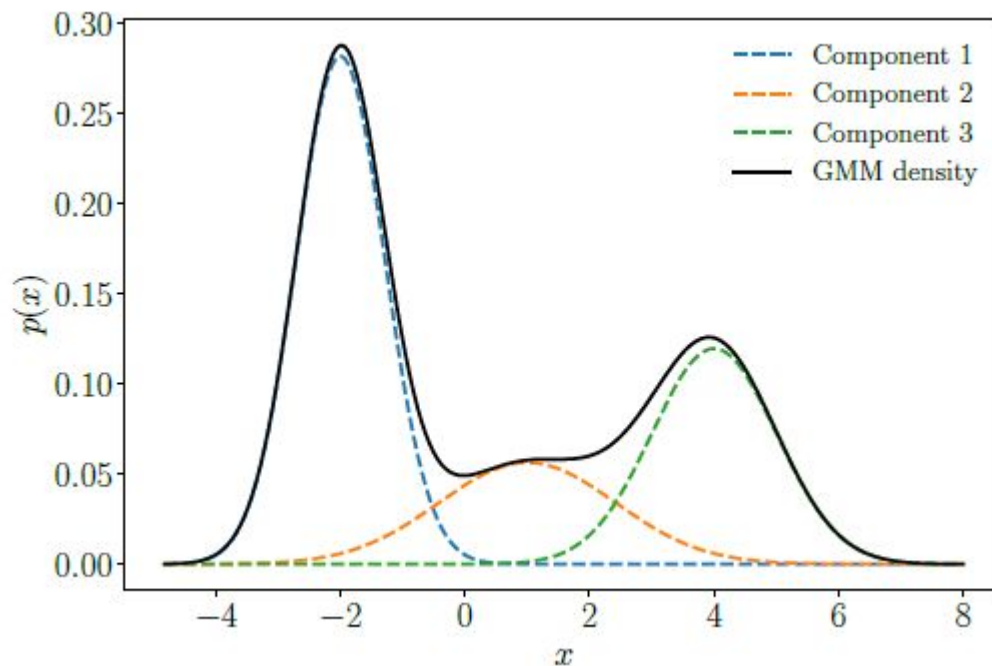
Aplicaciones Específicas

- Density Estimation
- Clustering
- Anomaly Detection
- Object Tracking
- Speech Feature Extraction



Formulación

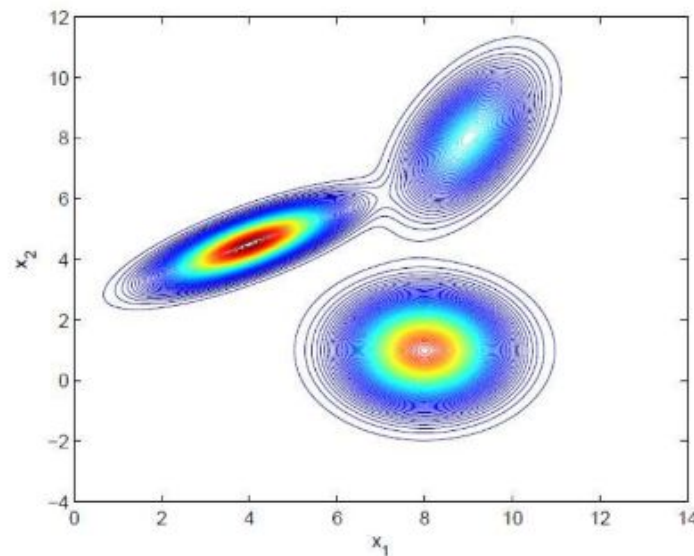
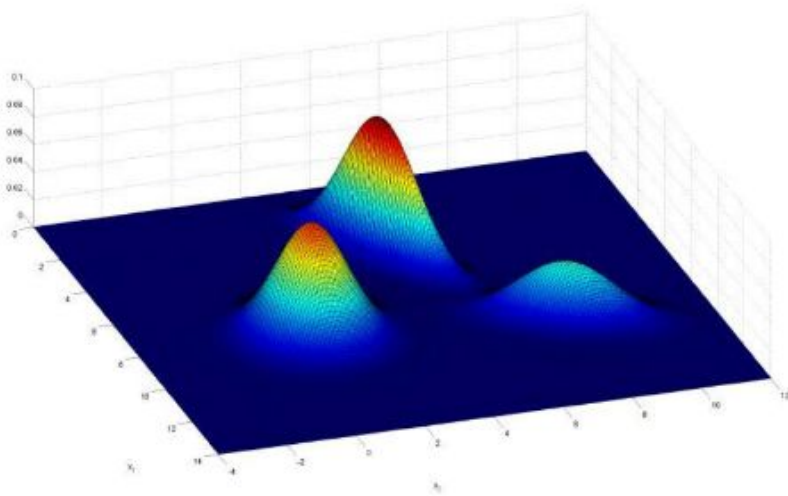
$$p(x|\theta, \hat{\sigma}) = \sum_i \lambda_i \mathcal{N}(x|\mu_i, \sigma_i^2)$$



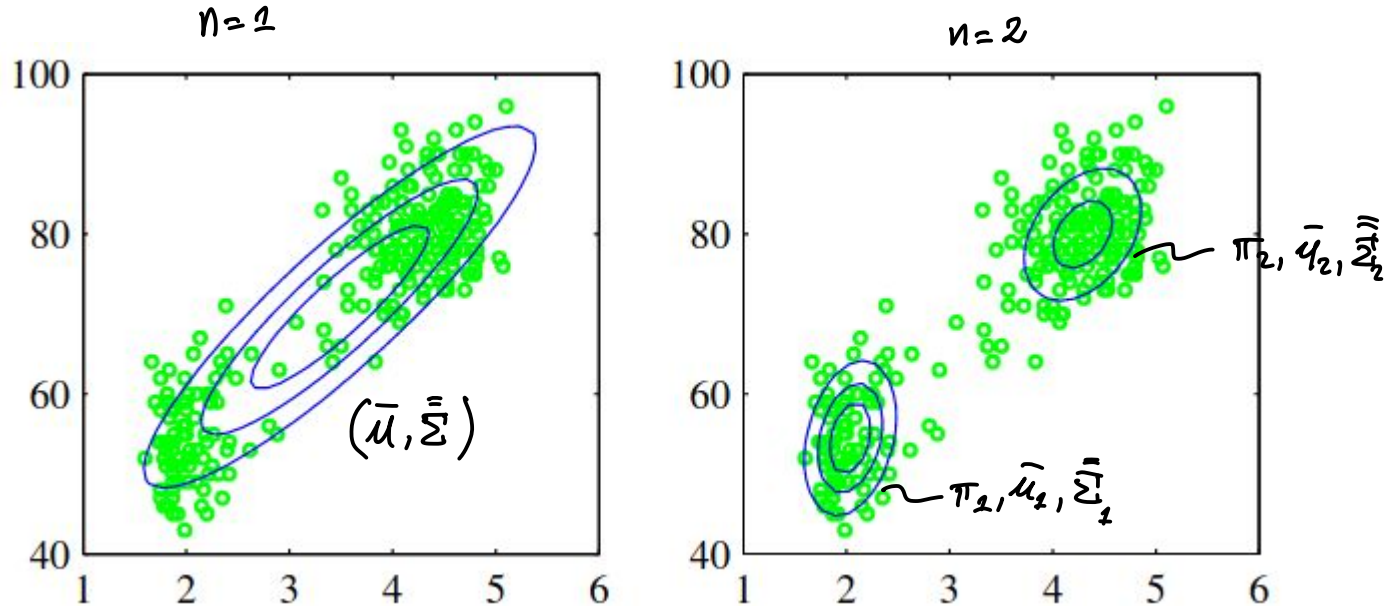
$$p(x|\theta) = 0.5\mathcal{N}(x|-2, \frac{1}{2}) + 0.2\mathcal{N}(x|1, 2) + 0.3\mathcal{N}(x|4, 1)$$

Formulación

$$p(x) = \underbrace{0.3}_{\pi_1} \mathcal{N}\left(x \mid \underbrace{\begin{pmatrix} 4 \\ 4.5 \end{pmatrix}}_{\mu_1}, \underbrace{\begin{pmatrix} 1.2 & 0.6 \\ 0.6 & 0.5 \end{pmatrix}}_{\Sigma_1}\right) + \underbrace{0.5}_{\pi_2} \mathcal{N}\left(x \mid \underbrace{\begin{pmatrix} 8 \\ 1 \end{pmatrix}}_{\mu_2}, \underbrace{\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}}_{\Sigma_2}\right) + \underbrace{0.2}_{\pi_3} \mathcal{N}\left(x \mid \underbrace{\begin{pmatrix} 9 \\ 8 \end{pmatrix}}_{\mu_3}, \underbrace{\begin{pmatrix} 0.6 & 0.5 \\ 0.5 & 1.5 \end{pmatrix}}_{\Sigma_3}\right)$$

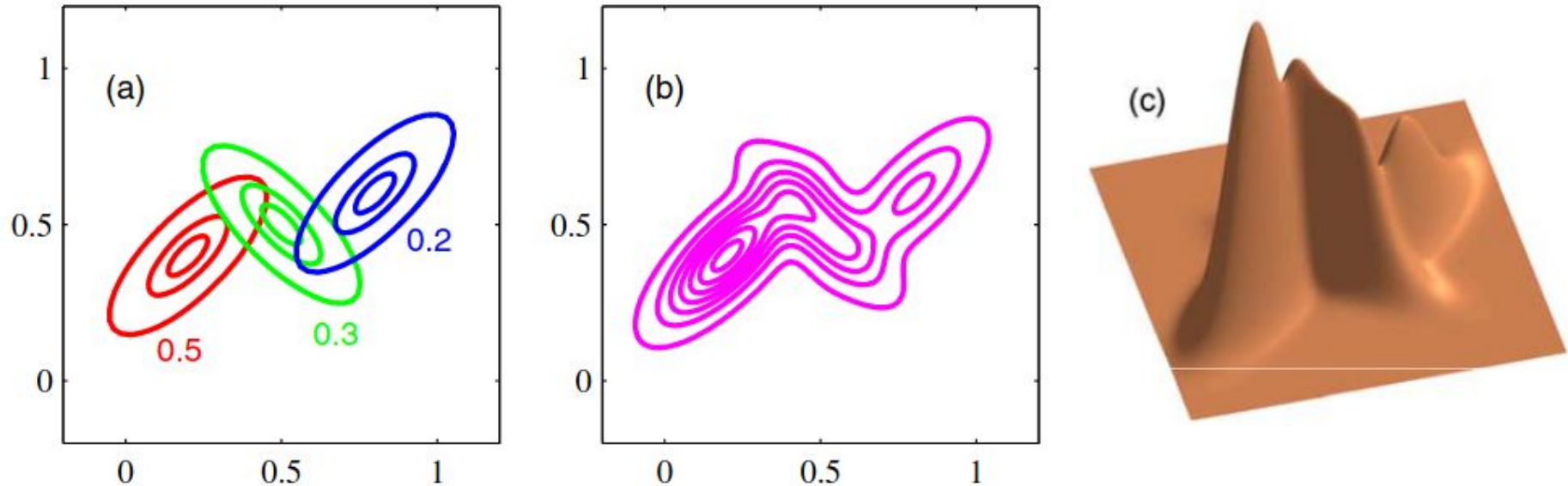


Gaussian Mixture Models: Estudio de fenómenos naturales



“Old Faithful” dataset. 272 mediciones de erupciones del “Old Faithful” geyser en el Parque Nacional Yellowstone. El eje horizontal representa la duración de una erupción (medida en minutos) y el vertical el tiempo hasta la próxima erupción.

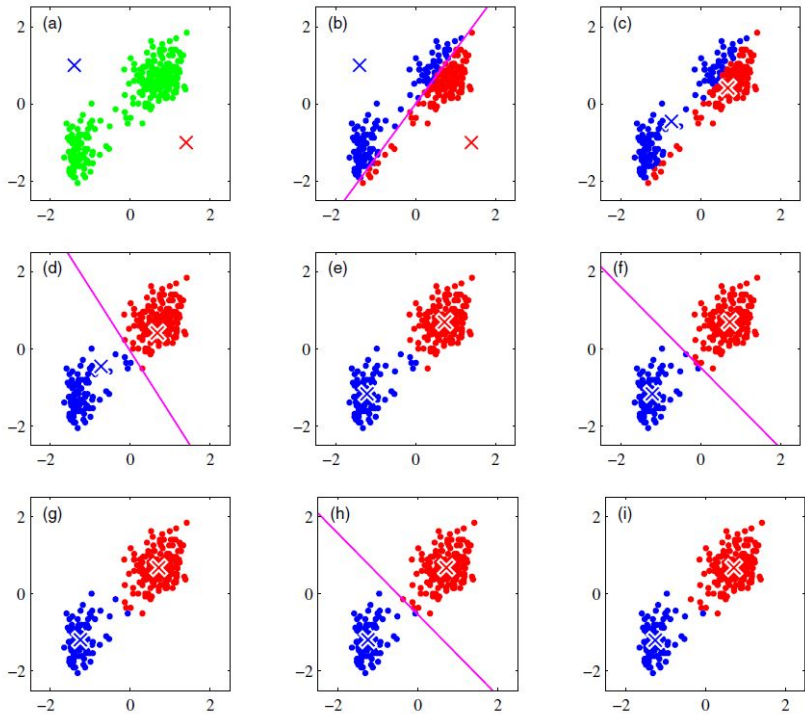
Gaussian Mixture Models: Clustering



Ejemplo de Gaussian Mixture. En la imagen (a) se muestran las tres distribuciones subyacentes indicando con colores sus variables latentes. En la imagen (b) las curvas de nivel de la distribución conjunta y en la (c) la densidad.

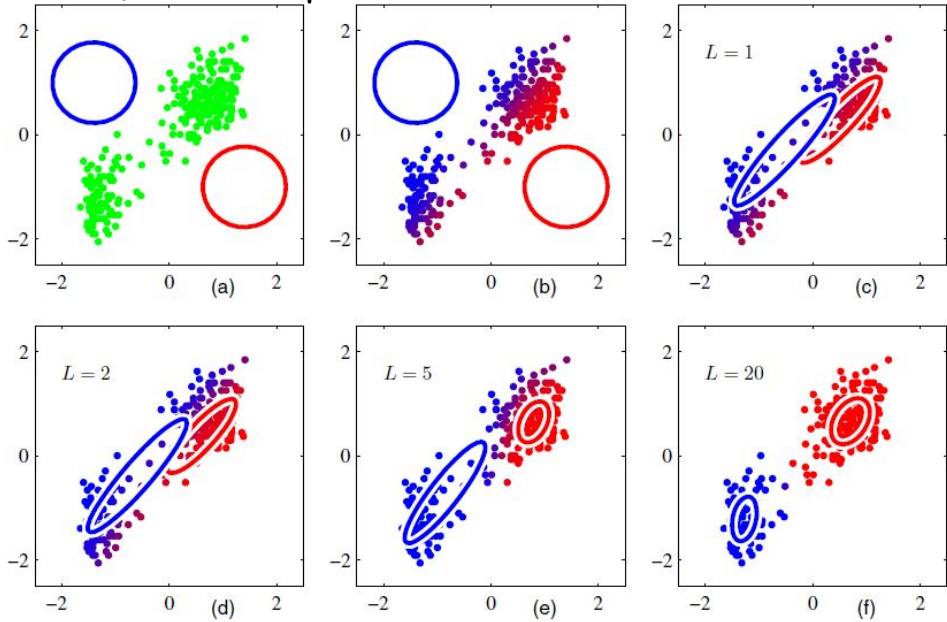
Gaussian Mixture Models: Clustering

μ_1, μ_2



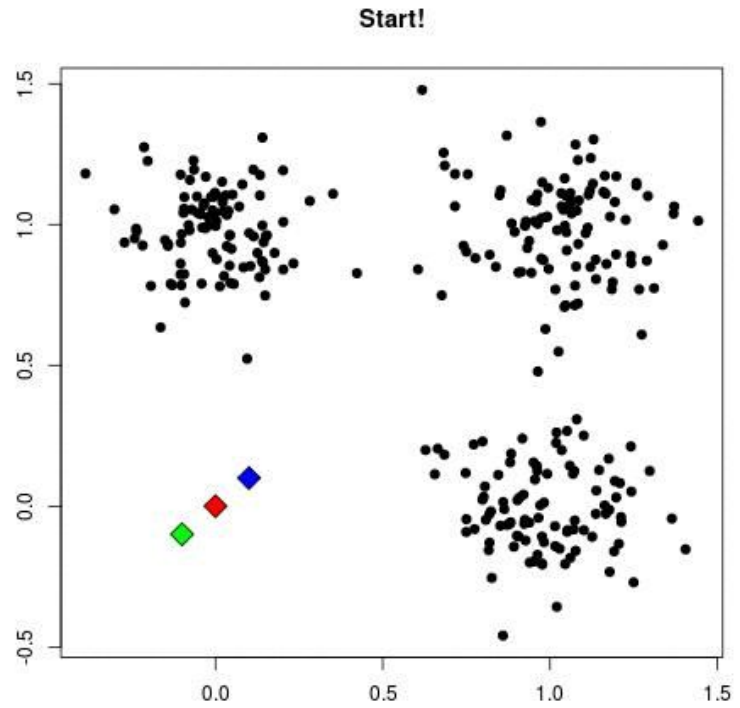
KMeans

- $\mu_1^1, \sigma_1^1 (\sigma_1^1, \sigma_2^1)$
- $\mu_1^2, \sigma_1^2 (\sigma_1^2, \sigma_2^2)$
- \vdots

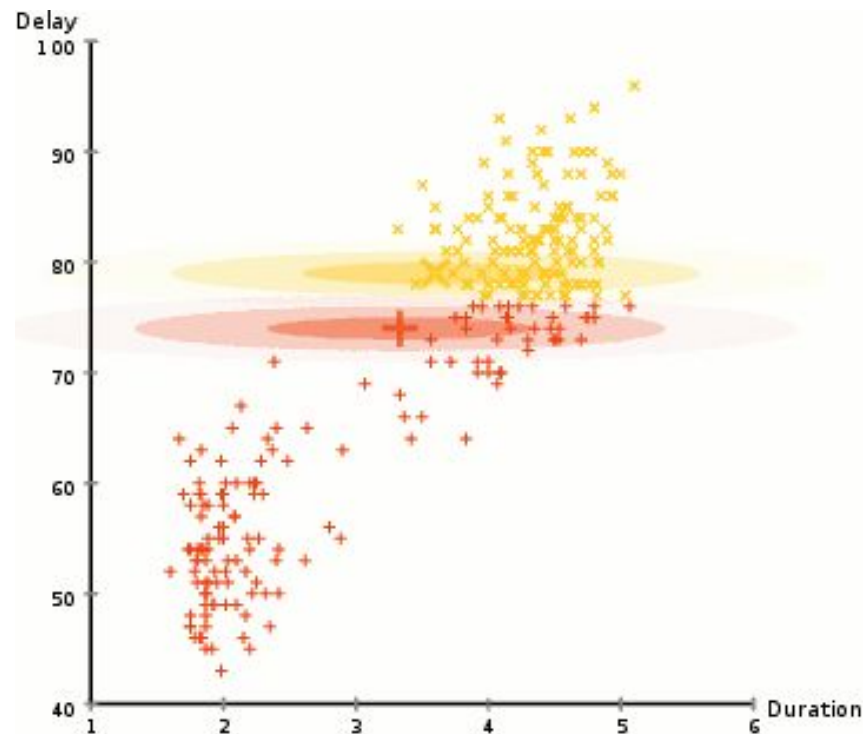


GMM

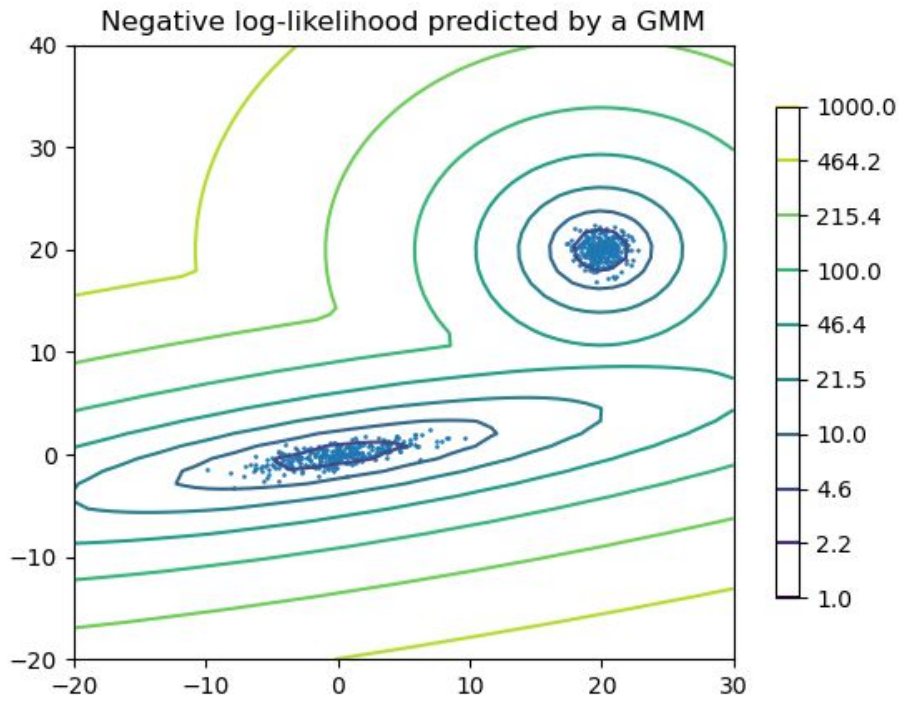
Gaussian Mixture Models - kMeans



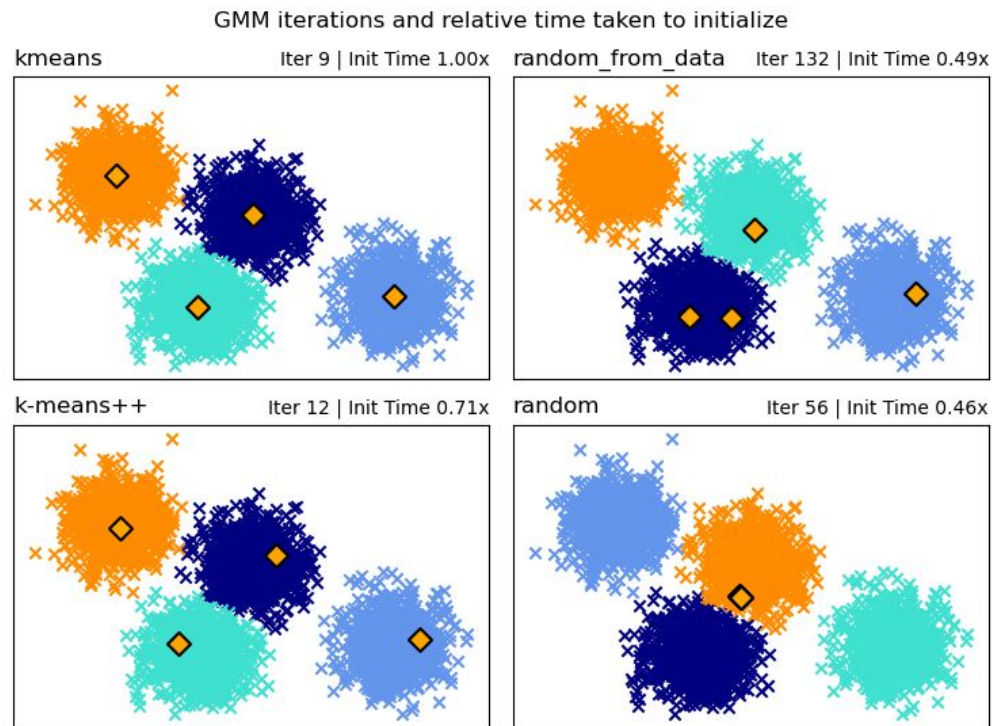
Gaussian Mixture Models: Clustering



Gaussian Mixture Models: Detección de anomalías



Gaussian Mixture Models: Inicialización



Formulación

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k p_k(\mathbf{x})$$

$$0 \leq \pi_k \leq 1, \quad \sum_{k=1}^K \pi_k = 1,$$

Mixture Models - General

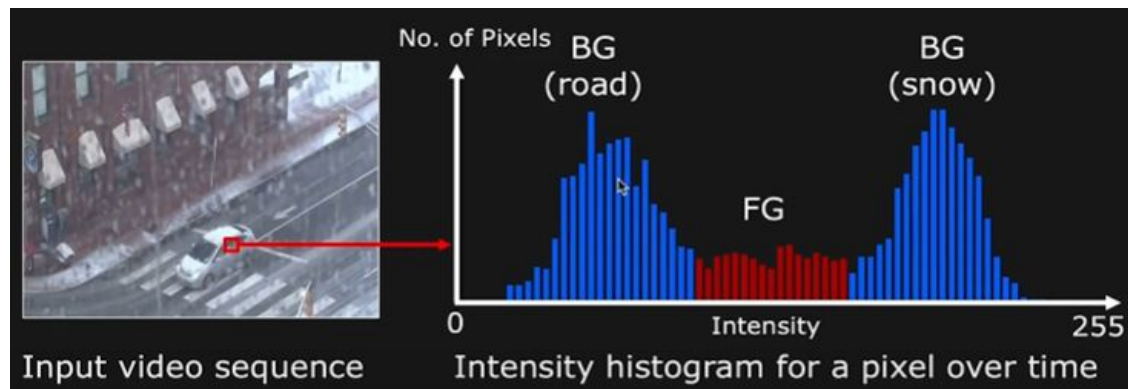
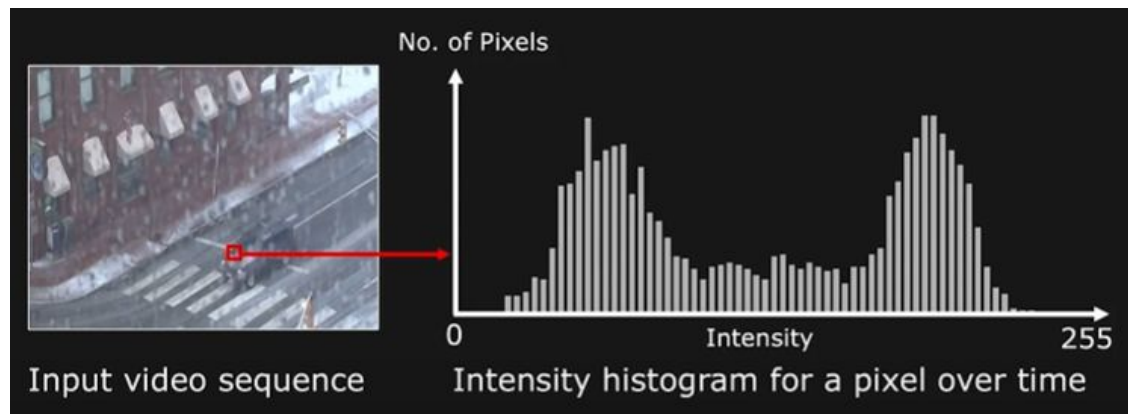
$$p(\mathbf{x} | \boldsymbol{\theta}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

$$0 \leq \pi_k \leq 1, \quad \sum_{k=1}^K \pi_k = 1,$$

$$\boldsymbol{\theta} := \{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k : k = 1, \dots, K\}$$

Gaussian Mixture Models

Gaussian Mixture Models - Object Tracking



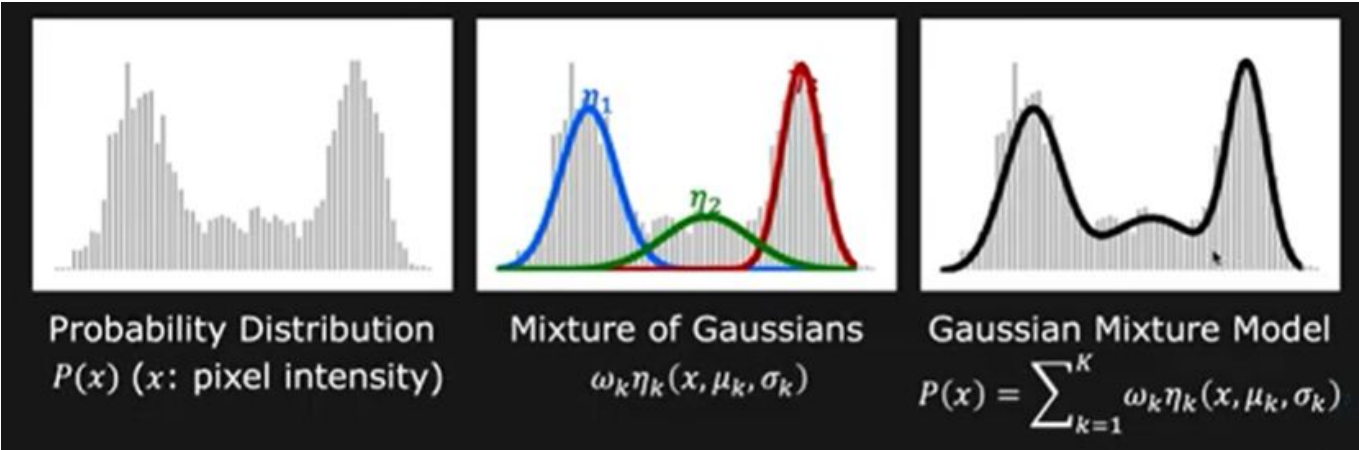
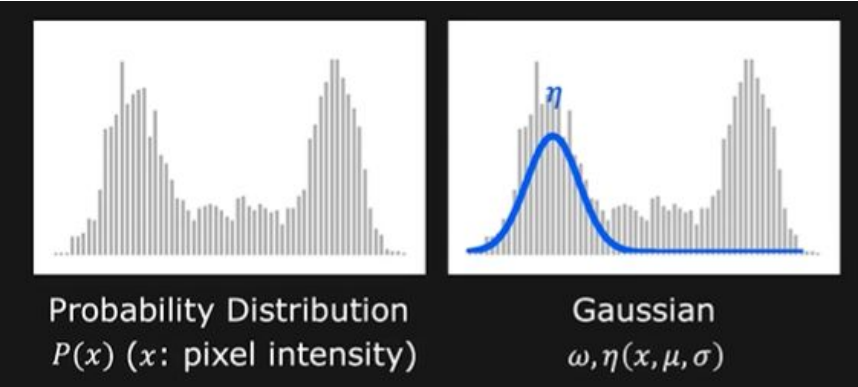
Gaussian Mixture Models - Object Tracking

$$P(\mathbf{X}) \cong \sum_{k=1}^K \omega_k \eta_k(\mathbf{X}, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad \text{such that} \quad \sum_{k=1}^K \omega_k = 1$$

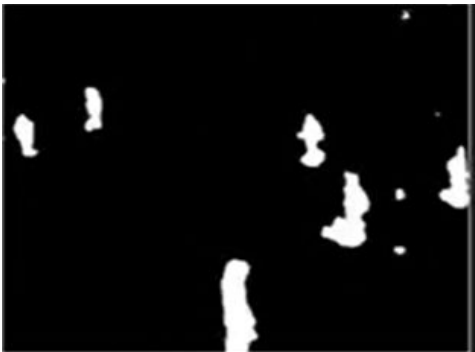
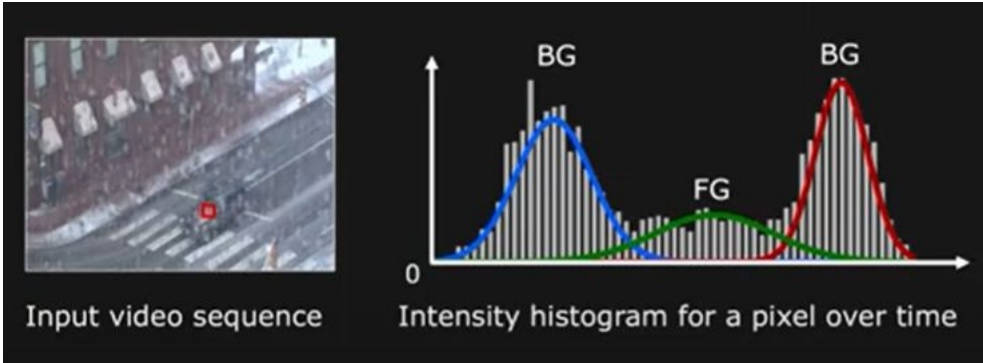
where: $\eta(\mathbf{X}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2} |\boldsymbol{\Sigma}|^{1/2}} e^{-\frac{1}{2}(\mathbf{X}-\boldsymbol{\mu})^T (\boldsymbol{\Sigma})^{-1} (\mathbf{X}-\boldsymbol{\mu})}$

Mean $\boldsymbol{\mu} = \begin{bmatrix} \mu_r \\ \mu_g \\ \mu_b \end{bmatrix}$ Covariance matrix $\boldsymbol{\Sigma} = \begin{bmatrix} \sigma^2 & 0 & 0 \\ 0 & \sigma^2 & 0 \\ 0 & 0 & \sigma^2 \end{bmatrix}$ (can be a full matrix)

Gaussian Mixture Models - Object Tracking



Gaussian Mixture Models - Object Tracking



\uparrow
 $\frac{\omega}{\sigma}$ Background

\downarrow
 $\frac{\omega}{\sigma}$ Foreground

GMM y EM - JAMBOARD

K - Means:

Partimos de un dataset $D = \{x_1, \dots, x_n\}$, $x_i \in \mathbb{R}^n$, con esto queremos encontrar K clusters que modelen a D .

Vamos a considerar μ_k centros con $k \in [1, \dots, K]$ llamados **centroides** que van a ser nuestros parámetros a estimar. ¿Cómo los estimamos?

Vamos a tomar un μ_k y asignar mi dato x_i de manera de minimizar una función de pérdida:

$$r_{nk} = \begin{cases} 1 & \text{si } x_n \in K \\ 0 & \text{o.w.} \end{cases}$$

Esto me
'labeliza' mis
datos

$$\underbrace{\begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{bmatrix}}_{\mu_k} \left\} x_i$$

Con esto definimos:

$$J = \sum_n \sum_k r_{nk} \|x_n - \mu_k\|^2$$

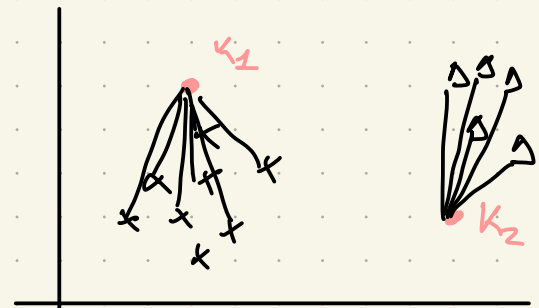
Poses a seguir:

(criterio)

1. inicializamos μ_k

2. Asignamos 'labels' según distancia $\Rightarrow \min_{r_{nk}} J(n, k)$ (μ_k fijo)

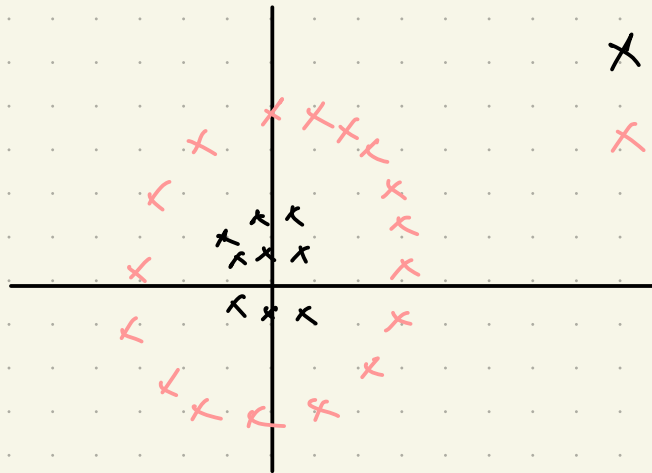
3. minimizar J respecto a μ_k , fijando r_{nk} (**Actualización de centroides**)



$$\min_{\mu_k} J(r_{nk}, \mu_k) \leadsto \mu_k = \frac{1}{\sum_n r_{nk}} \cdot \sum_n r_{nk} x_k = \sum_{n/x_n \in K} x_n \cdot \frac{1}{N_k}$$

$$K \text{ means} \rightarrow \text{dist euclidea} \quad \tilde{J} = \sum \sum r_{nk} N(\mu_k, x_n)$$

si yo vario esto \rightarrow K-proto
K-meoids



x label 1 (real)

x label 2 (real)

• caso atípico
en Kmeans.

GMM

$$p(x) = \sum_k \pi_k \mathcal{N}(x | \bar{\mu}_k, \bar{\Sigma}_k) ; \quad \sum_k \pi_k = 1$$

Consideramos z v.a. $\in \{0, 1\} \wedge \sum_k z_k = 1$.

Definimos $P(x, z) = P(x|z) P(z) \rightarrow P(z_k = 1) = \pi_k \rightarrow \text{prob del centro } k$

$P(x|z_k) \sim \mathcal{N}(\bar{\mu}_k, \bar{\Sigma}_k)$ (cada k tiene par. distintos)

con esto $p(x) = \sum \pi_k \cdot P(x|z_k)$

con esto tenemos de datos \rightarrow conjunta $P(x, z) = P(x|z) P(z)$

\rightarrow marginales $P(z) = [\pi_1, \dots, \pi_K]$

\rightarrow cond. $P(x|z) \sim \mathcal{N}$ ①

\rightarrow marg $p(x) = \sum_k (\dots)$

partimos de ①:

$$\gamma(z_k) = p(z_k=1|x) = \frac{p(z_k=1) p(x|z_k=1)}{\sum_{j=1}^K p(z_j=1) p(x|z_j=1)} = \frac{\pi_k \mathcal{N}(\bar{\mu}_k, \bar{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\bar{\mu}_j, \bar{\Sigma}_j)}$$

↑
distrib. a posteriori

$$\text{log likelihood} \leadsto \ln P(x|\bar{\pi}, \bar{\mu}, \bar{\Sigma}) = \sum_n \ln \left\{ \sum_k \pi_k \mathcal{N}(\bar{\mu}_k, \bar{\Sigma}_k) \right\}$$

$\max_{\pi, \mu, \Sigma} l \leadsto$ no tiene sol. cerrada
 \leadsto tenemos $K!$ soluciones.

para resolver esto vamos a usar **Expectation Maximization**.

consiste en derivar log likelihood e igualarlo a cero.

$$\partial_{\theta} l = 0$$

$$\hookrightarrow \partial_{\mu_k} l = 0 \rightarrow \mu_k^* = \frac{1}{N_k} \sum \gamma(z_{nk}) \cdot x_n$$

$$\partial_{\Sigma_k} l = 0 \rightarrow \Sigma_k^* = \frac{1}{N_k} \sum \gamma(z_{nk}) (x_n - \mu_k)(x_n - \mu_k)^t$$

$$\partial_{\pi_k} l = 0 \rightarrow \pi_k^* = N_k/N$$

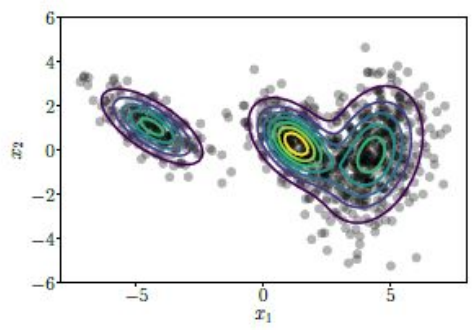
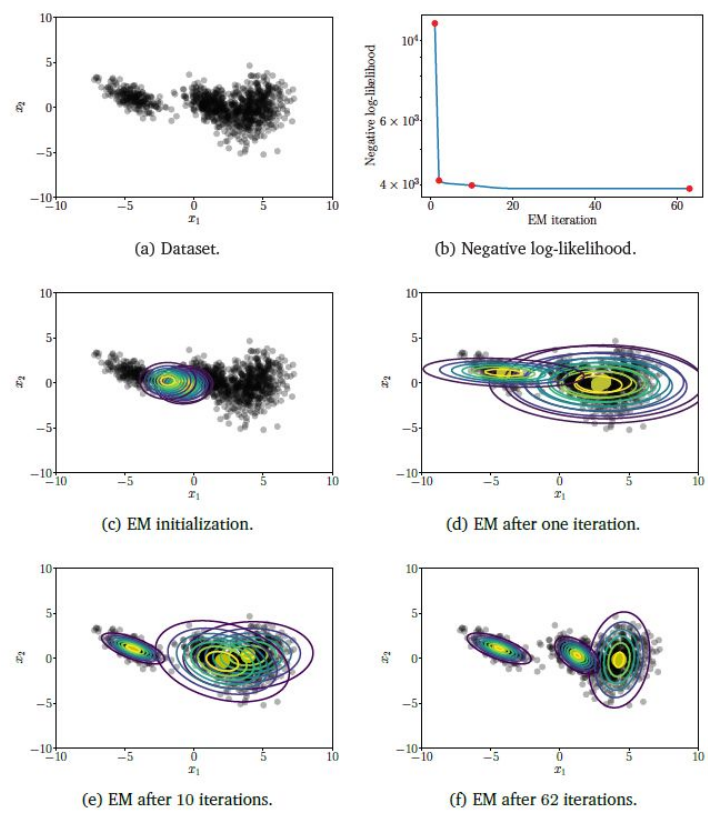
N_k : número efectivo de
ptos en la clase k

$$N_k = \sum_n \gamma(z_{nk})$$

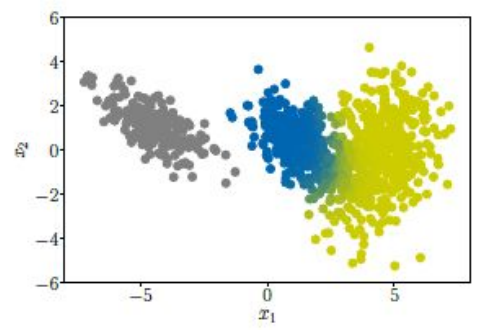
Pasos a seguir:

1. inicializamos μ_k, Σ_k, π_k
2. con estos param. vamos a evaluar $\gamma(z_{nk})$
3. Actualizamos los parametros $\mu_k', \Sigma_k', \pi_k'$
4. calcular el log likelihood nuevo.

Gaussian Mixture Models - Teoría



(a) GMM fit after 62 iterations.



(b) Dataset colored according to the responsibilities of the mixture components.

Notebooks

Bibliografía

- Mathematics for Machine Learning | Deisenroth, Faisal, Ong
- Pattern Recognition and Machine Learning | Bishop
- Gaussian Mixture Model | John McGonagle, Geoff Pilling, Andrei Dobre
- Expectation-Maximization Algorithms | Stanford CS229: Machine Learning
- First Principles of Computer Vision | Computer Science Department, School of Engineering and Applied Sciences, Columbia University