

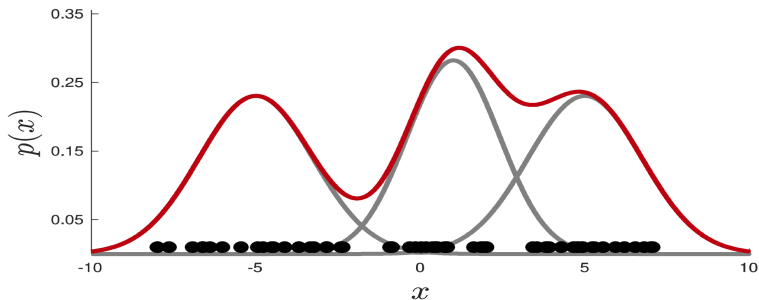
Clasificación Bayesiana

Estimación paramétrica de densidades desconocidas

FaMAF

2019

Modelos de mezclas de gaussianas



Es frecuente encontrar conjuntos de datos que requieren la combinación de varias funciones de densidad para modelarlos de manera más cercana a la distribución real. Se requiere una combinación de tres funciones de densidad para modelar el conjunto de datos representados con puntos negros.

Modelos de mezclas de gaussianas

- Modelar una distribución de probabilidad desconocida mediante la combinación lineal de m funciones de densidad (componentes o modelos):

$$p(\mathbf{x}|\Theta) = \sum_{j=1}^m p(\mathbf{x}|\boldsymbol{\theta}_j) p_j \quad (1)$$

donde $\Theta = [\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_m, p_1, \dots, p_m]^T$

- Y se satisface

$$\sum_{j=1}^m p_j = 1 \quad \text{y} \quad \int_{\mathbf{x}} p(\mathbf{x}|\boldsymbol{\theta}_j) d\mathbf{x} = 1 \quad (2)$$

Maximización de la expectativa

- Función de log-verosimilitud en función de Θ :

$$\mathcal{L}(\Theta) = \sum_{i=1}^{\ln} \sum_{j=1}^N p(\mathbf{x}_i | \theta_j) p_j \quad (3)$$

- Imposible generar una solución analítica (resolver $\nabla_{\Theta} \mathcal{L}(\Theta)$); no hay información sobre cuál componente generó la muestra \mathbf{x}_i .
- El algoritmo maximización de la expectativa (EM) introduce una variable indicadora $y_i = j$ si la muestra \mathbf{x}_i fue generada por la j -ésima componente.
- EM maximiza la expectativa de $\mathcal{L}(\Theta)$ condicionada a las muestras observadas \mathbf{X} y la iteración actual del estimado de Θ .

Maximización de la expectativa

Input: $X = \{\mathbf{x}_i | i = 1, \dots, n\}, m, \epsilon$

Output: $\hat{\Theta}$

1. $\hat{\Theta}^0 \leftarrow \left[\left(\hat{\theta}_1, \hat{p}_1 \right)^0, \dots, \left(\hat{\theta}_m, \hat{p}_m \right)^0 \right]$
2. $t \leftarrow 0$
3. **do**
4. $t \leftarrow t + 1$
5. **Paso-E:** $Q \left(\Theta; \hat{\Theta}^t \right) \leftarrow \mathbb{E} \left[\sum_{i=1}^{\infty} \ln \left(p \left(\mathbf{x}_i | j; \hat{\Theta}_j^t \right) p_j^t \right) \right]$
6. **Paso-M:** $\hat{\Theta}^{t+1} \leftarrow \arg \max_{\Theta} Q \left(\Theta; \hat{\Theta}^t \right)$
7. **until** $\left| Q \left(\Theta; \hat{\Theta}^t \right) - Q \left(\Theta; \hat{\Theta}^{t+1} \right) \right| < \epsilon$

Maximización de la expectativa

El algoritmo EM para mezclas Gaussianas:

1. Inicialización: definir los parámetros iniciales de m componentes

$$\text{Gaussianas: } \hat{\Theta}^0 = \left[\left(\hat{\mu}_1, \hat{\Sigma}_1, \hat{p}_1 \right)^0, \dots, \left(\hat{\mu}_m, \hat{\Sigma}_m, \hat{p}_m \right)^0 \right]$$

2. Paso-E: Computar la esperanza de $\mathcal{L}(\Theta)$:

$$Q(\Theta; \hat{\Theta}^t) = \sum_{i=1}^n \sum_{j=1}^m \ln \left(p(\mathbf{x}_i | j; \hat{\theta}_j^t) \hat{p}_j^t \right) p(j | \mathbf{x}_i; \hat{\Theta}^t) \quad (4)$$

donde

$$p(j | \mathbf{x}; \hat{\Theta}^t) = \frac{p(\mathbf{x} | j; \hat{\theta}_j^t) \hat{p}_j^t}{\sum_{k=1}^m p_k(\mathbf{x} | \hat{\theta}_k^t) \hat{p}_k^t}$$

y

$$\ln \left(p(\mathbf{x} | j; \hat{\theta}_j^t) \hat{p}_j^t \right) = -\frac{1}{2}(\mathbf{x} - \hat{\mu})^T \hat{\Sigma}^{-1}(\mathbf{x} - \hat{\mu}) - \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\hat{\Sigma}| + \ln \hat{p}_j$$

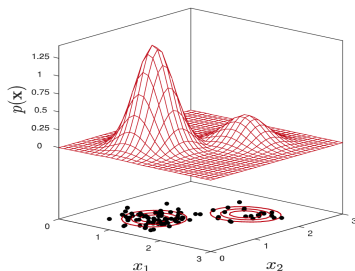
Maximización de la expectativa

3 Paso-M: Maximizar $Q(\Theta; \hat{\Theta}^t)$ con respecto de $\hat{\mu}_j$, $\hat{\Sigma}_j$ y \hat{p}_j :

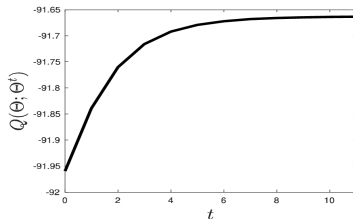
$$\begin{aligned}\hat{\mu}_j^{t+1} &= \frac{\sum_{i=1}^n p(j|\mathbf{x}_i; \hat{\Theta}^t) \mathbf{x}_i}{\sum_{i=1}^n p(j|\mathbf{x}_i; \hat{\Theta}^t)} \\ \hat{\Sigma}_j^{t+1} &= \frac{\sum_{i=1}^n p(j|\mathbf{x}_i; \hat{\Theta}^t) (\mathbf{x}_i - \hat{\mu}_j^{t+1}) (\mathbf{x}_i - \hat{\mu}_j^{t+1})^T}{\sum_{i=1}^n p(j|\mathbf{x}_i; \hat{\Theta}^t)} \\ \hat{p}_j^{t+1} &= \frac{1}{n} \sum_{i=1}^n p(j|\mathbf{x}_i; \hat{\Theta}^t)\end{aligned}\tag{5}$$

4 Finalizar si $|Q(\Theta; \hat{\Theta}^t) - Q(\Theta; \hat{\Theta}^{t+1})| \leq \epsilon$; en caso contrario regresar al Paso-E.

Maximización de la Expectativa



(a)



(b)

(a) Muestras bivariadas originadas a partir de dos funciones normales y estimación de la función de densidad mediante el algoritmo EM. (b) Valores de log-verosimilitud en función del número de iteraciones.

Inicialización de parámetros

1. Estimar las medias $\hat{\mu}_j$ para $j = 1, \dots, m$ componentes mediante el algoritmo k -means ($k = m$)
2. Generar variables indicadoras: $y_i = \arg \min_{j=1, \dots, m} \|\mathbf{x}_i - \hat{\mu}_j\|_2$
3. Estimar la probabilidad $\hat{p}_j = n_j/n$ para $j = 1, \dots, m$
4. Estimar la matriz de covarianza $\hat{\Sigma}_j$:

$$\hat{\Sigma}_j = \frac{1}{n_j - 1} \sum_{\mathbf{x}_i \in j} (\mathbf{x}_i - \hat{\mu}_j) (\mathbf{x}_i - \hat{\mu}_j)^T, \quad j = 1, \dots, m$$

5. Generar el arreglo $\hat{\Theta}^0 = \left[\left(\hat{\mu}_1, \hat{\Sigma}_1, \hat{p}_1 \right), \dots, \left(\hat{\mu}_m, \hat{\Sigma}_m, \hat{p}_m \right) \right]$

Número óptimo de componentes

Principio de descripción de longitud mínima (MDL): minimiza la suma de la complejidad del modelo $L(\mathbf{M})$ y la eficiencia de la descripción de los datos con respecto al modelo $L(\mathbf{X}, \mathbf{M})$:

$$L(\mathbf{X}, \mathbf{M}) = \min_{\mathbf{M} \in \mathcal{M}} [L(\mathbf{M}) + L(\mathbf{X}|\mathcal{M})] \quad (6)$$

Número óptimo de componentes

- Principio MDL para mezclas Gaussianas con m componentes:

$$L(\mathbf{X}, \hat{\Theta}_m) = \frac{\alpha_m}{2} \ln n - \sum_{i=1}^n \ln \sum_{j=1}^m p(\mathbf{x}_i | j; \hat{\theta}_j) \hat{p}_j \quad (7)$$

donde

$$\alpha_m = (m - 1) + md + \frac{md(d + 1)}{2}$$

- Número óptimo de componentes:

$$m^* = \arg \min_{k=1, \dots, m} [L(\mathbf{X}, \Theta_k)] \quad (8)$$

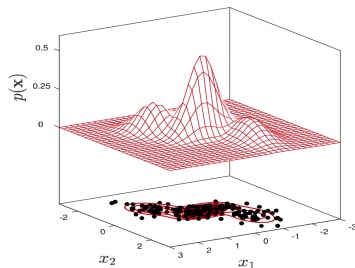
Número óptimo de componentes

Input: $\mathbf{X} = \{\mathbf{x}_i | i = 1, \dots, n\}, m_{\max}, \epsilon$

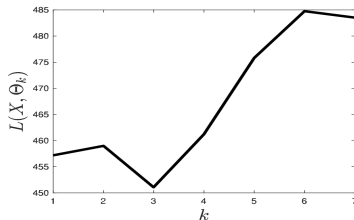
Output: m^*

1. $m^* \leftarrow 1$
2. for $(m = 1; m \leq m_{\max}; m++)$ do
3. $\hat{\Theta}_m \leftarrow \text{EM}(\mathbf{X}, m, \epsilon)$ // Algoritmo EM
4. $L(\mathbf{X}, \hat{\Theta}_m) \leftarrow \frac{\alpha_m}{2} \ln n - \sum_{i=1}^n \ln \sum_{j=1}^m p(\mathbf{x}_i | j; \hat{\theta}_j) \hat{p}_j$
5. if $(m > 1) \wedge (L(\mathbf{X}, \hat{\Theta}_m) < L(\mathbf{X}, \hat{\Theta}_{m-1}))$ then
6. $m^* \leftarrow m$
7. end
8. end

Número óptimo de componentes



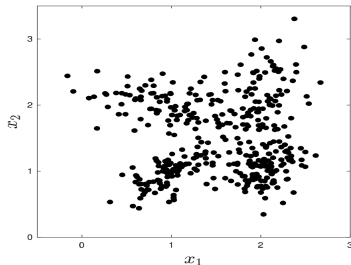
(a)



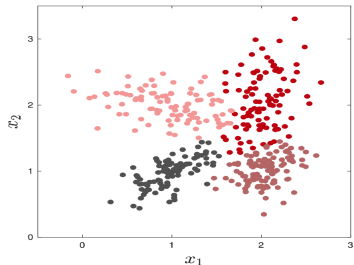
(b)

(a) Muestras bivariadas originadas a partir de tres funciones normales y estimación de la función de densidad mediante el algoritmo EM. El valor m^* se estimó mediante el principio MDL. (b) MDL en función del número de componentes donde $m^* = 3$.

Aprendizaje no supervisado



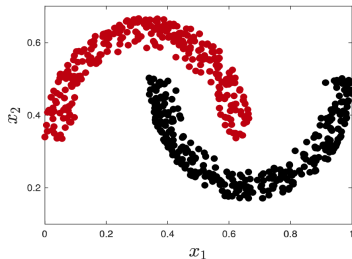
(a)



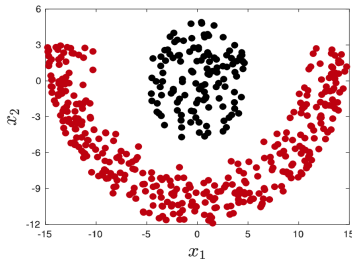
(b)

El método de mezclas Gaussianas es útil para agrupamiento de datos: (a) datos sin etiquetar y (b) grupos formados mediante mínima distancia Mahalanobis.

Aprendizaje supervisado



(a)



(b)

Para modelar distribuciones de clases con formas complejas, se puede utilizar un modelo de mezclas Gaussianas para aproximar las distribución real de cada clase.

Aprendizaje supervisado

Definir una colección de c modelos de mezclas Gaussianas:

$$\hat{\Theta}_{\Omega} = \left\{ \hat{\Theta}_{\omega_1}, \dots, \hat{\Theta}_{\omega_c} \right\}$$

donde el modelo para la i -ésima clase es

$$\hat{\Theta}_{\omega_i} = \{M_i, S_i, p_i\}$$

donde cada elemento en $\hat{\Theta}_{\omega_i}$ es

$$M_i = \left\{ \hat{\mu}_1, \dots, \hat{\mu}_{m_i^*} \right\}, \quad S_i = \left\{ \hat{\Sigma}_1, \dots, \hat{\Sigma}_{m_i^*} \right\}, \quad p_i = [\hat{p}_1, \dots, \hat{p}_{m_i^*}]$$

donde m_i^* es el número óptimo de modelos para la i -ésima clase.

Aprendizaje supervisado

Input: $X = \{(\mathbf{x}_i, y_i) \mid i = 1, \dots, n\}$ con $y_i \in \{\omega_1, \dots, \omega_c\}$ m_{\max}, ϵ

Output: Θ_Ω

1. $\hat{\Theta}_\Omega \leftarrow \emptyset$
2. for $(i = 1; i \leq c; i++)$ do
3. Identificar patrones de la clase $\omega_i : \mathbf{X}_{y=\omega_i} \in X$
4. $\hat{\Theta}_{\omega_i} \leftarrow \text{MDL}(\mathbf{X}_{y=\omega_i}, m_{\max}, \epsilon)$ // Algoritmo MDL
5. $\hat{\Theta}_\Omega \leftarrow \{\hat{\Theta}_\Omega + \hat{\Theta}_{\omega_i}\}$
6. end
7. end

Aprendizaje supervisado

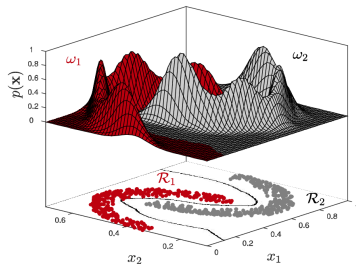
- Etapa de clasificación, un patrón de prueba \mathbf{x}_t se evalúa en el modelo entrenado de mezclas Gaussianas de la i -ésima clase, $\hat{\Theta}_{\omega_i}$:

$$g_i(\mathbf{x}_t) = \sum_{j=1}^{m_i^*} p(\mathbf{x} | \hat{\boldsymbol{\mu}}_j, \hat{\boldsymbol{\Sigma}}_j) \hat{p}_j, \quad \text{para } i = 1, \dots, c \quad (9)$$

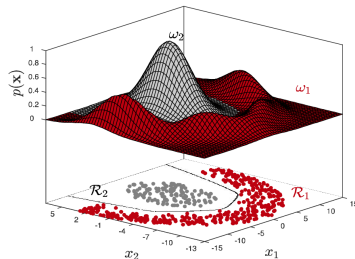
- La regla de decisión es

$$\text{Decidir } \omega_i \text{ si } g_i(\mathbf{x}_t) > g_j(\mathbf{x}_t) \text{ para todo } i \neq j \quad (10)$$

Aprendizaje supervisado



(a)



(b)

Estimación de densidad con mezclas Gaussianas para clases que no siguen una distribución normal. La frontera de decisión se define mediante (10)

Aprendizaje supervisado con datos faltantes

- Supongamos tener una muestra de entrenamiento D con datos faltantes D_b , esto es, $D = \{x_1, \dots, x_n\} = D_g \cup D_b$.
- El estimador de máxima verosimilitud es el vector $\hat{\theta} = (\hat{\theta}_1, \dots, \hat{\theta}_k)$ que maximiza la $\mathcal{L}(\theta) = \sum_{i=1}^n \ln p(\mathbf{x}_i | \theta)$, que no es calculable con datos faltantes.
- Una solución aproximada sería maximizar la siguiente función objetivo

$$\mathcal{L}(\theta) = E_{D_b} \left(\sum_{i=1}^n \ln p(\mathbf{x}_i | D_g, \theta) \right) \quad (11)$$

y utilizar el algoritmo de Maximización de la expectativa para calcular dicho máximo en forma iterativa.

Maximización de la expectativa

1. Initialize $\Theta^0, T, i = 0$
2. **do** $i \leftarrow i + 1$
3. **Paso-E:** $Q(\Theta; \hat{\Theta}^i)$
4. **Paso-M:** $\Theta^{i+1} \leftarrow \arg \max_{\Theta} Q(\Theta; \Theta^i)$

$$Q(\theta; \theta^i) = E_{D_b} \left(\sum_{i=1}^n \ln p(\mathbf{x}_i | D_g, \theta) \right) \quad (12)$$

5. **until** $|Q(\Theta; \Theta^i) - Q(\Theta; \Theta^{i+1})| < T$
6. return $\hat{\Theta} \leftarrow \Theta^{i+1}$

Ejemplo

Supongamos que tenemos dos distribuciones normales y los siguientes datos

$$D = \{x_1, x_2, x_3, x_4\} = \left\{ \begin{pmatrix} 0 \\ 2 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 2 \\ 2 \end{pmatrix}, \begin{pmatrix} * \\ 4 \end{pmatrix} \right\} = D_g \cup D_b \quad (13)$$

Entonces $D_b = x_{41}$ $\theta = (\mu_1, \mu_2, \sigma_1^2, \sigma_2^2)$. Sea $\theta^0 = (0, 0, 1, 1)$ el punto de inicio. La función objetivo es

$$Q(\theta; \theta^0) = E_{D_b} (\ln p(D_g, D_b) | D_g; \theta^0) \quad (14)$$

$$= \int_{-\infty}^{\infty} \left[\sum_{k=1}^3 \ln p(x_k | \theta) + \ln p(x_4 | \theta) \right] p(x_{41} | \theta^0; x_{42} = 4) dx_{41} \quad (15)$$

Ejemplo

$$= \sum_{k=1}^3 \ln p(x_k|\theta) + \int_{-\infty}^{\infty} \ln p\left(\begin{pmatrix} x_{41} \\ 4 \end{pmatrix} \middle| \theta\right) \cdot \frac{p\left(\begin{pmatrix} x_{41} \\ 4 \end{pmatrix} \middle| \theta^0\right)}{\int_{-\infty}^{\infty} p\left(\begin{pmatrix} x'_{41} \\ 4 \end{pmatrix} \middle| \theta^0\right) dx'_{41}} dx_{41}$$

La integral del denominador es constante y puede ser sacada fuera de la integral general. Sustituyendo $p\left(\begin{pmatrix} x_{41} \\ 4 \end{pmatrix} \middle| \theta^0\right)$ la densidad Gaussiana en el valor de inicio θ^0 el paso E resulta

Ejemplo

$$= \sum_{k=1}^3 \ln p(x_k | \theta) + \frac{1}{K} \int_{-\infty}^{\infty} \ln p \left(\begin{pmatrix} x_{41} \\ 4 \end{pmatrix} | \theta \right) \cdot \frac{\exp \left[-\frac{1}{2} (x_{41}^2 + 4^2) \right]}{2\pi \left| \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \right|} dx_{41} \quad (16)$$

$$= \sum_{k=1}^3 \ln p(x_k | \theta) - \frac{1 + \mu_1^2}{2\sigma_1^2} - \frac{(1 - \mu_2)^2}{2\sigma_2^2} - \ln(2\pi\sigma_1\sigma_2) \quad (17)$$

Esta función $Q(\theta; \theta^0)$ tiene un máximo en $\theta^1 = \begin{pmatrix} 0.75 \\ 2 \\ 0.938 \\ 2 \end{pmatrix}$

Iterando de la misma forma en tres pasos se obtiene $\theta = \begin{pmatrix} 1 \\ 2 \\ 0.667 \\ 2 \end{pmatrix}$

Ejemplo

Es importante destacar que el óptimo para la verosimilitud con todos los datos quizás sea otro valor. Por ejemplo, si el dato faltante es

$x_{14} = 1$, la solución es $\theta = \begin{pmatrix} 1 \\ 2 \\ 0.5 \\ 2 \end{pmatrix}$ y el cálculo de $Q(\theta; \theta^0)$ es

mucho más simple pues no involucra integrales.