

Clustering GMM

GaussianMixtureModel

PhD(e). Jonnatan Arias Garcia – jonnatan.arias@utp.edu.co –
jariasg@uniquindio.edu.co

PhD. David Cardenas peña - dcardenas@utp.edu.co

PhD. Hernán Felipe Garcia - hernanf.garcia@udea.edu.co

Agrupamiento probabilístico

- Mezcla de funciones de probabilidad
- Algoritmo EM

Mezcla de funciones de probabilidad

- Una forma de aproximar funciones de probabilidad multimodales es a través de una mezcla de funciones de probabilidad.
- De las mezclas de funciones de probabilidad, la mezcla de Gaussianas es una de las más conocidas,

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

donde K es el número de componentes de la mezcla, y los parámetros π_k son probabilidades que satisfacen

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

Variable latente \mathbf{z}

- ❑ Se introduce una variable aleatoria binaria de K dimensiones \mathbf{z} con representación 1 de K .
- ❑ El vector \mathbf{z} puede tomar uno de K estados, de acuerdo a cuál de los elementos es diferente de cero.
- ❑ La distribución marginal sobre \mathbf{z} se especifica como

$$p(z_k = 1) = \pi_k.$$

- ❑ De forma compacta, esta distribución se escribe como

$$p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k}.$$

Distribución condicional de \mathbf{x} dado \mathbf{z}

- La distribución condicional de \mathbf{x} dado un valor particular de \mathbf{z} , es una Gaussiana

$$p(\mathbf{x}|z_k = 1) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- En forma compacta,

$$p(\mathbf{x}|\mathbf{z}) = \prod_{k=1}^K \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k}$$

Distribución marginal de \mathbf{x}

- ❑ La probabilidad conjunta está dada por $p(\mathbf{z})p(\mathbf{x}|\mathbf{z})$.
- ❑ La distribución marginal de \mathbf{x} se obtiene como,

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

- ❑ Si existen varios datos observados, a cada dato observado \mathbf{x}_n le corresponde una variable latente \mathbf{z}_n .
- ❑ Este es una nueva formulación de la distribución de mezclas usando variables latentes, lo que permite trabajar con la distribución conjunta $p(\mathbf{x}, \mathbf{z})$.

Distribución condicional de \mathbf{z} dado \mathbf{x}

- ❑ Otra cantidad que juega un papel importante es la probabilidad condicional de \mathbf{z} dado \mathbf{x} .
- ❑ Esta probabilidad está dada como

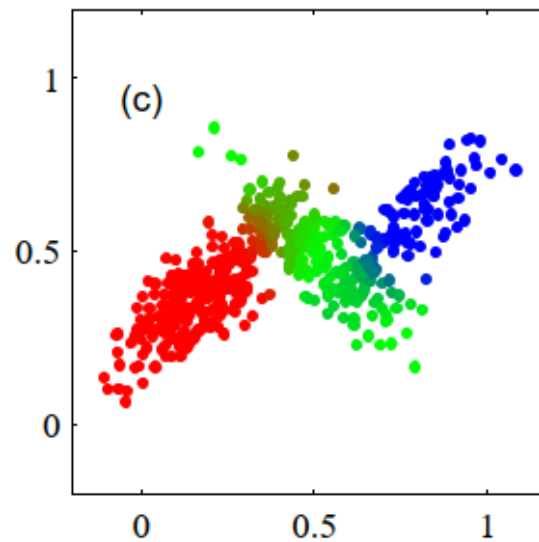
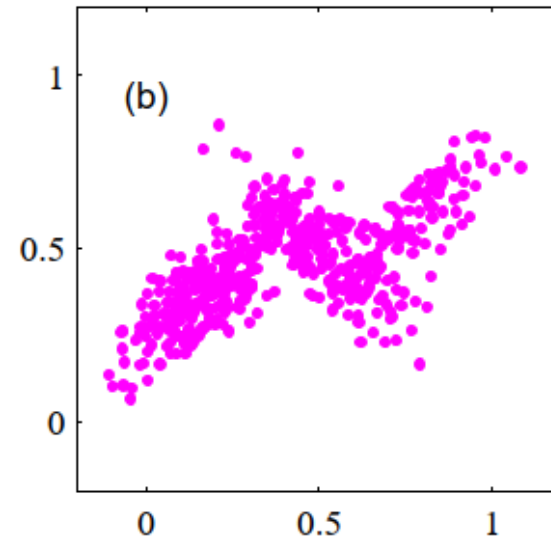
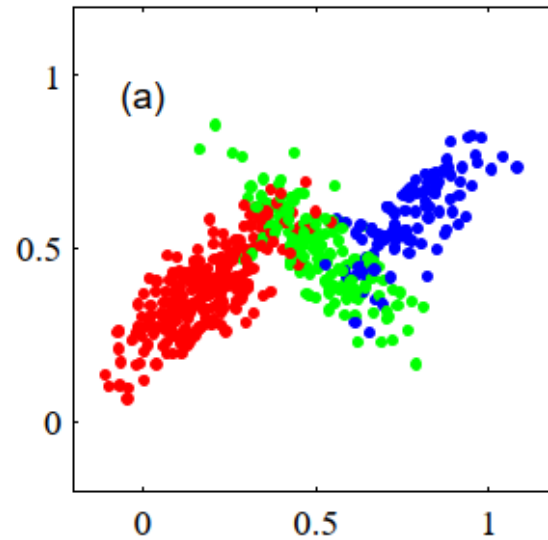
$$\begin{aligned}\gamma(z_k) \equiv p(z_k = 1 | \mathbf{x}) &= \frac{p(z_k = 1)p(\mathbf{x} | z_k = 1)}{\sum_{j=1}^K p(z_j = 1)p(\mathbf{x} | z_j = 1)} \\ &= \frac{\pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}\end{aligned}$$

- ❑ Se puede ver a π_k como la probabilidad a priori de que $p(z_k = 1)$ y a $\gamma(z_k)$ como la probabilidad a posteriori correspondiente una vez se ha observado \mathbf{x} .
- ❑ Esta cantidad puede verse como la responsabilidad que la componente k asume para explicar la observación \mathbf{x} .

Definición del tipo de datos

- Al conjunto $\{\mathbf{X}, \mathbf{Z}\}$ se le conoce como el conjunto completo de datos.
- Al conjunto de datos observados \mathbf{X} se le conoce como los datos incompletos.
- Del conjunto $\{\mathbf{X}, \mathbf{Z}\}$ sólo se conoce \mathbf{X} . La única información sobre \mathbf{Z} está en la función de probabilidad $p(\mathbf{Z}|\mathbf{X}, \theta)$.

Datos incompletos y completos



Algoritmo EM

Función de verosimilitud logarítmica

- Se parte de un conjunto de datos $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ que se desean modelar con una mezcla de Gaussianas.
- Este conjunto de datos se representan con una matriz \mathbf{X} de dimensiones $N \times D$ y filas \mathbf{x}_n^\top .
- Similarmente, las variables latentes correspondientes se denotan por una matriz \mathbf{Z} con filas \mathbf{z}_n^\top y de dimensiones $N \times K$.
- La función de verosimilitud logarítmica está dada por

$$\ln p(\mathbf{X}|\pi, \mu, \Sigma) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k) \right\}$$

- Encontrar los parámetros $\theta^{\text{old}} = \{\{\pi_k\}_{k=1}^K, \{\mu_k\}_{k=1}^K, \{\Sigma_k\}_{k=1}^K\}$, que maximicen la función de verosimilitud de los datos incompletos.

Algoritmo EM

Dada una distribución conjunta $p(\mathbf{X}, \mathbf{Z}|\theta)$, el objetivo es maximizar la función de verosimilitud $p(\mathbf{X}|\theta)$ con respecto a los parámetros θ .

1. Escoger un valor inicial para los parámetros θ^{old} .
2. **Paso E.** Evaluar $p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}})$.
3. **Paso M.** Evaluar θ^{new} dada por

$$\theta^{\text{new}} = \arg \max_{\theta} Q(\theta, \theta^{\text{old}}),$$

donde

$$Q(\theta, \theta^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\theta).$$

4. Verificar la convergencia de la función de verosimilitud o de los parámetros. Si no se satisface el criterio de convergencia, luego $\theta^{\text{old}} \leftarrow \theta^{\text{new}}$ y volver al paso 2.

Mezcla de Gaussianas: Aplicación del paso E

- Comenzando con un valor de θ^{old} se calcula la probabilidad a posteriori de las variables latentes \mathbf{Z} dados los datos \mathbf{X} y los parámetros θ^{old} .
- La función de probabilidad $p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}})$ tiene como elementos $\gamma(z_{n,k})$.
- Las probabilidades $\gamma(z_{n,k})$ están dadas como

$$\begin{aligned}\gamma(z_{n,k}) \equiv p(z_{n,k} = 1|\mathbf{x}_n) &= \frac{p(z_{n,k} = 1)p(\mathbf{x}_n|z_{n,k} = 1)}{\sum_{j=1}^K p(z_{n,j} = 1)p(\mathbf{x}_n|z_{n,j} = 1)} \\ &= \frac{\pi_k \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}\end{aligned}$$

- $p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}})$ es una tabla de dimensiones $N \times K$.

Mezcla de Gaussianas: Aplicación del paso M(I)

- ❑ Encontremos primero la función $Q(\theta, \theta^{\text{old}})$.

- ❑ La función $Q(\theta, \theta^{\text{old}})$ está dada como

$$Q(\theta, \theta^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\theta) = \mathbb{E}_{\mathbf{Z}}[\ln p(\mathbf{X}, \mathbf{Z}|\theta)]$$

- ❑ La función de verosimilitud de los datos completos para la mezcla de Gaussianas está dada como

$$p(\mathbf{X}, \mathbf{Z}|\pi, \mu, \Sigma) = \prod_{n=1}^N \prod_{k=1}^K \pi_k^{z_{nk}} \mathcal{N}(\mathbf{x}_n|\mu_k, \Sigma_k)^{z_{nk}}$$

- ❑ La verosimilitud logarítmica está como

$$\ln p(\mathbf{X}, \mathbf{Z}|\pi, \mu, \Sigma) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \{\ln \pi_k + \ln \mathcal{N}(\mathbf{x}_n|\mu_k, \Sigma_k)\}.$$

Mezcla de Gaussianas: Aplicación del paso M(II)

- La función $Q(\theta, \theta^{\text{old}})$ está dada como

$$\begin{aligned} Q(\theta, \theta^{\text{old}}) &= \mathbb{E}_{\mathbf{Z}}[\ln p(\mathbf{X}, \mathbf{Z} | \pi, \mu, \Sigma)] \\ &= \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \{ \ln \pi_k + \ln \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k) \}. \end{aligned}$$

- Nótese en la ecuación anterior, que $\mathbb{E}_{\mathbf{Z}}[z_{nk}]$ coincide con $\gamma(z_{nk})$,

$$\mathbb{E}_{\mathbf{Z}}[z_{nk}] = \sum_{z_{nk}} z_{nk} p(z_{nk} | \mathbf{X}, \theta^{\text{old}}) = p(z_{nk} = 1 | \mathbf{x}_n, \theta^{\text{old}}) = \gamma(z_{nk}).$$

- Dado $\gamma(z_{nk})$ la idea es ahora maximizar $Q(\theta, \theta^{\text{old}})$ con respecto a los parámetros $\theta = \{ \{ \pi_k \}_{k=1}^K, \{ \mu_k \}_{k=1}^K, \{ \Sigma_k \}_{k=1}^K \}$.

Mezcla de Gaussianas: Aplicación del paso M(III)

- La maximización de $Q(\theta, \theta^{\text{old}})$ con respecto a π_k conduce a

$$\pi_k^{\text{new}} = \frac{1}{N} \sum_{n=1}^N \gamma(z_{nk}) = \frac{N_k}{N},$$

donde $N_k = \sum_{n=1}^N \gamma(z_{nk})$.

- La maximización de $Q(\theta, \theta^{\text{old}})$ con respecto a μ_k conduce a

$$\mu_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n$$

- La maximización de $Q(\theta, \theta^{\text{old}})$ con respecto a Σ_k conduce a

$$\Sigma_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \mu_k^{\text{new}})(\mathbf{x}_n - \mu_k^{\text{new}})^{\top}.$$

Resumen

1. Se escoge un valor inicial para θ^{new} .
2. Paso **E**. Se calcula $\gamma(z_{nk})$, para $n = 1, \dots, N$ y $k = 1, \dots, K$.
3. Paso **M**. Se usan las fórmulas de actualización para π_k^{new} , μ_k^{new} y Σ_k^{new} , para $k = 1, \dots, K$.
4. Se verifica la convergencia de la función de verosimilitud o de los parámetros. Si no se satisface el criterio de convergencia, luego $\theta^{\text{old}} \leftarrow \theta^{\text{new}}$ y se repite desde el paso 2.

Ejemplo

