



# EM应用例子-GMM参数估计

---



# EM算法（期望最大算法）

- 设一个完整数据集
- 存在完整数据集到观测数据（不完整集）的映射
- 以完整数据集的条件数学期望代替对数似然函数
- 为使积分进行，待估计参数使用猜测值，构成迭代

$$\mathbf{x} = [x_1, x_2, \dots, x_M]^T$$

$$\mathbf{y} = g(\mathbf{x})$$

$$\begin{aligned} E_{x|y} [\ln p_x(\mathbf{x}|\theta)] \\ = \int \ln p_x(\mathbf{x}|\theta) p(\mathbf{x}|\mathbf{y}, \theta) d\mathbf{x} \end{aligned}$$

$$\theta^{(0)} \rightarrow \theta^{(1)}, \dots, \rightarrow \theta^{(k)}$$



## EM算法描述

---

第 1 步：初始化，选择  $\theta$  的初始猜测值，令  $m = 0$ ，给出  $\theta^{(m)}$ ；

第 2 步：由观测数据  $y$  和  $\theta$  的猜测值  $\theta^{(m)}$ ，

得到完整数据集  $x$  的条件概率  $p(x|y, \theta^{(m)})$ ；

第 3 步：计算完整数据集下对数似然函数  $l(\theta|x) = \log p_x(x|\theta)$  的条件期望

$$Q(\theta|\theta^{(m)}) = E_{x|y, \theta^{(m)}} [\log p_x(x|\theta)] = \int \log p_x(x|\theta) p(x|y, \theta^{(m)}) dx$$



续

---

第 4 步：求  $\theta = \theta^{(m+1)}$  使得  $Q(\theta|\theta^{(m)})$  最大，即<sub>↵</sub>

$$\theta^{(m+1)} = \arg \max_{\theta \in \Omega} \{Q(\theta|\theta^{(m)})\}_{\substack{\uparrow \\ \downarrow}}$$

第 5 步：若满足停止条件，则  $\hat{\theta} = \theta^{(m+1)}$  为所得 MLE，若不满足停止条件，

令  $m := m + 1$  返回第 2 步。<sub>↵</sub>



## 特例**1**：数据缺失 或隐变量情况

---

观测数据集  $\mathbf{y}$

完整数据集  $\mathbf{x} = \{\mathbf{y}, \mathbf{z}\}$

隐藏数据集  $\mathbf{z}$

则有：

$$\begin{aligned} Q(\boldsymbol{\theta} | \boldsymbol{\theta}^{(m)}) &= E_{\mathbf{z} | \mathbf{y}, \boldsymbol{\theta}^{(m)}} [\log p_x(\mathbf{y}, \mathbf{z} | \boldsymbol{\theta})] \\ &= \int_{\mathfrak{Z}} \log p_x(\mathbf{y}, \mathbf{z} | \boldsymbol{\theta}) p(\mathbf{z} | \mathbf{y}, \boldsymbol{\theta}^{(m)}) d\mathbf{z} \end{aligned}$$



## 样本独立同分布情况

设  $p(\mathbf{x}|\boldsymbol{\theta}) = \prod_{i=0}^{N-1} p(\mathbf{x}(n_i)|\boldsymbol{\theta}),$

并且  $\mathbf{y}(n_i)$  仅与  $\mathbf{x}(n_i)$  有关

则有 
$$Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(m)}) = \sum_{i=0}^{N-1} Q_i(\boldsymbol{\theta}|\boldsymbol{\theta}^{(m)})$$

$$Q_i(\boldsymbol{\theta}|\boldsymbol{\theta}^{(m)}) = E_{\mathbf{x}(n_i)|\mathbf{y}(n_i), \boldsymbol{\theta}^{(m)}} [\log p_x(\mathbf{x}(n_i)|\boldsymbol{\theta})]$$

这里

$$= \int \log p_x(\mathbf{x}(n_i)|\boldsymbol{\theta}) p(\mathbf{x}(n_i)|\mathbf{y}(n_i), \boldsymbol{\theta}^{(m)}) d\mathbf{x}$$



# EM算法解高斯混合模型

---

观测向量集表示为  $\mathbf{y} = \{\mathbf{y}_i | 0 \leq i < N\}$  是i.i.d.的

GMM模型

$$p(\mathbf{y}_i) = \sum_{k=1}^K w_k N(\mathbf{y}_i | \boldsymbol{\mu}_k, \mathbf{C}_k)$$

约束条件

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

由i.i.d.样本集估计参数集

$$\boldsymbol{\theta} = \{w_k, \boldsymbol{\mu}_k, \mathbf{C}_k | k = 1, 2, \dots, K\}$$

定义隐变量  $z_i \in \{1, 2, \dots, K\}$

隐变量向量  $\mathbf{z} = [z_0, z_1, \dots, z_{N-1}]^T$

完整数据集  $\mathbf{x} = \{\mathbf{y}, \mathbf{z}\}$

由

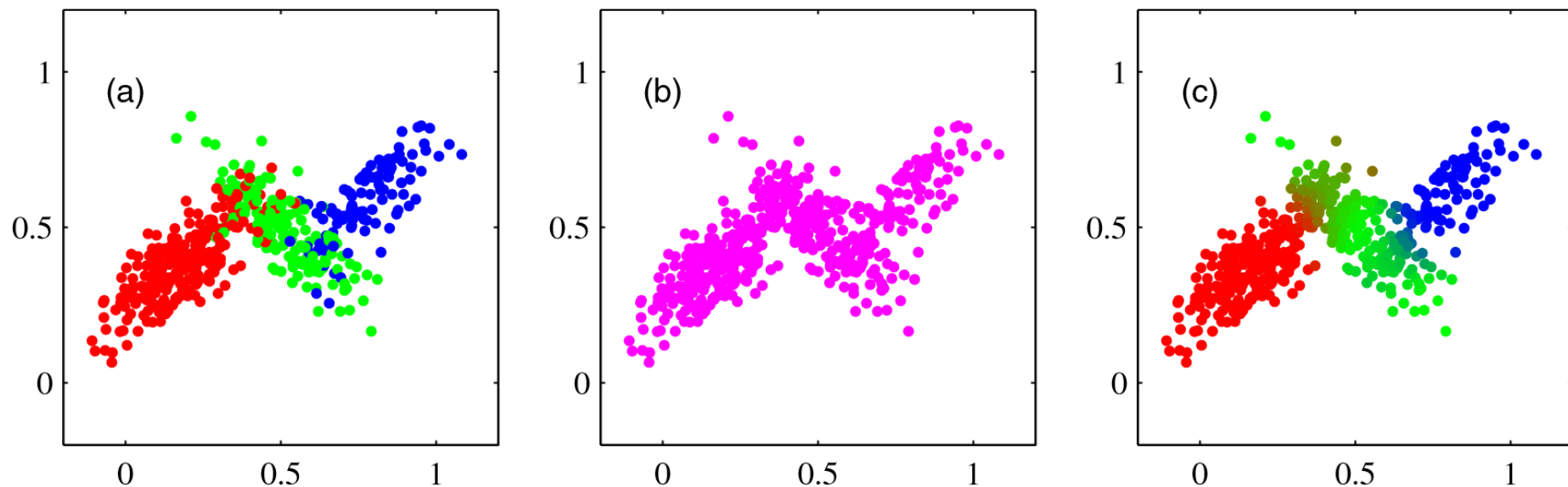
$$p_x(\mathbf{y}_i, z_i = k | \boldsymbol{\theta}) = w_k N(\mathbf{y}_i | \boldsymbol{\mu}_k, \mathbf{C}_k)$$

设已得到参数向量的猜测值  $\boldsymbol{\theta}^{(m)}$  则

$$\begin{aligned} p(z_i = k | \mathbf{y}_i, \boldsymbol{\theta}^{(m)}) &= \frac{p_x(\mathbf{y}_i, z_i = k | \boldsymbol{\theta}^{(m)})}{p(\mathbf{y}_i | \boldsymbol{\theta}^{(m)})} \\ &= \frac{w_k^{(m)} N(\mathbf{y}_i | \boldsymbol{\mu}_k^{(m)}, \mathbf{C}_k^{(m)})}{\sum_{l=1}^K w_l^{(m)} N(\mathbf{y}_i | \boldsymbol{\mu}_l^{(m)}, \mathbf{C}_l^{(m)})} \end{aligned}$$



# GMM: 隐变量观点示例



左：已知  $\pi_k$  和参数  $\mu_k, \Sigma_k$ ，仿真产生若干数据，并记下  $z_i = k$ ，用红绿蓝表示样本点产生那个  $k$  分量，故：左图是联合分布

中：去掉颜色，即去掉隐变量信息，实际样本是不知隐变量

右：已知  $\pi_k, \mu_k, \Sigma_k$ ，用中图的各样本点坐标估计：

$$\gamma_{ik} = p(z_i = k | y_i)$$

实际中，只有中图的样本点集：估计  $\pi_k, \mu_k, \Sigma_k$

简记

$$\gamma_{ik}^{(m)} = p(z_i = k | \mathbf{y}_i, \boldsymbol{\theta}^{(m)})$$

$$\sum_{k=1}^K \gamma_{ik}^{(m)} = 1$$

由数据缺失和**i.i.d.**情况，得

$$Q_i(\boldsymbol{\theta} | \boldsymbol{\theta}^{(m)}) = E_{z_i | y_i, \boldsymbol{\theta}^{(m)}} [\log p_x(\mathbf{y}_i, z_i | \boldsymbol{\theta})]$$

$$= \sum_{k=1}^K p(z_i = k | \mathbf{y}_i, \boldsymbol{\theta}^{(m)}) \log p_x(\mathbf{y}_i, z_i = k | \boldsymbol{\theta})$$

$$= \sum_{k=1}^K \gamma_{ik}^{(m)} \log(w_k N(\mathbf{y}_i | \boldsymbol{\mu}_k, \mathbf{C}_k))$$

$$= \sum_{k=1}^K \gamma_{ik}^{(m)} \left[ \log w_k - \frac{1}{2} \log |\mathbf{C}_k| - \frac{1}{2} (\mathbf{y}_i - \boldsymbol{\mu}_k)^T \mathbf{C}_k^{-1} (\mathbf{y}_i - \boldsymbol{\mu}_k) \right] + C$$



E步

$$\underline{\text{令}} \quad n_k^{(m)} = \sum_{i=0}^{N-1} \gamma_{ik}^{(m)}$$

$$Q(\boldsymbol{\theta} | \boldsymbol{\theta}^{(m)}) = \sum_{i=0}^{N-1} Q_i(\boldsymbol{\theta} | \boldsymbol{\theta}^{(m)}) =$$

$$\sum_{i=0}^{N-1} \sum_{k=1}^K \gamma_{ik}^{(m)} \left[ \log w_k - \frac{1}{2} \log |\mathbf{C}_{\mathbf{k}}| - \frac{1}{2} (\mathbf{y}_i - \boldsymbol{\mu}_k)^T \mathbf{C}_{\mathbf{k}}^{-1} (\mathbf{y}_i - \boldsymbol{\mu}_k) \right]$$

$$= \sum_{k=1}^K n_k^{(m)} \log w_k - \frac{1}{2} \sum_{k=1}^K n_k^{(m)} \log |\mathbf{C}_{\mathbf{k}}|$$

$$- \frac{1}{2} \sum_{i=0}^{N-1} \sum_{k=1}^K \gamma_{ik}^{(m)} \left[ (\mathbf{y}_i - \boldsymbol{\mu}_k)^T \mathbf{C}_{\mathbf{k}}^{-1} (\mathbf{y}_i - \boldsymbol{\mu}_k) \right]$$

# M步

分别对各参数求导，并令导数为0，得

$$w_k^{(m+1)} = \frac{n_k^{(m)}}{\sum_{l=1}^K n_l^{(m)}} = \frac{n_k^{(m)}}{N}$$

$$\boldsymbol{\mu}_k^{(m+1)} = \frac{1}{n_k^{(m)}} \sum_{i=0}^{N-1} \gamma_{ik}^{(m)} \mathbf{y}_i$$

$$\mathbf{C}_k^{(m+1)} = \frac{1}{n_k^{(m)}} \sum_{i=0}^{N-1} \gamma_{ik}^{(m)} (\mathbf{y}_i - \boldsymbol{\mu}_k)(\mathbf{y}_i - \boldsymbol{\mu}_k)^T$$




## GMM参数的EM算法完整描述

S1 初始化  $m = 0$  , 设置停止门限  $\varepsilon$  ,

给出初值  $w_k^{(0)}$  ,  $\boldsymbol{\mu}_k^{(0)}$  ,  $\mathbf{C}_k^{(0)}$  ,  $k = 1, 2, \dots, K$  ,

并计算初始对数似然函数值

$$l(\boldsymbol{\theta}^{(m)} | \mathbf{y}) = \frac{1}{N} \sum_{i=0}^{N-1} \log \left( \sum_{k=1}^K w_k^{(m)} N(\mathbf{y}_i | \boldsymbol{\mu}_k^{(m)}, \mathbf{C}_k^{(m)}) \right)$$




---

S2 E-步, 对  $k = 1, 2, \dots, K$  计算<sup>4</sup>

$$\gamma_{ik}^{(m)} = \frac{w_k^{(m)} N(\mathbf{y}_i | \boldsymbol{\mu}_k^{(m)}, \mathbf{C}_k^{(m)})}{\sum_{l=1}^K w_l^{(m)} N(\mathbf{y}_i | \boldsymbol{\mu}_l^{(m)}, \mathbf{C}_l^{(m)})}, \quad i = 0, 1, \dots, N-1.$$

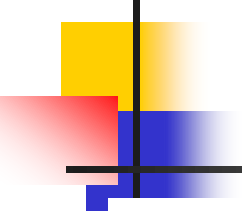
$$n_k^{(m)} = \sum_{i=0}^{N-1} \gamma_{ik}^{(m)} \quad ^4$$

S3 M 步, 对  $k = 1, 2, \dots, K$  计算。


$$w_k^{(m+1)} = \frac{n_k^{(m)}}{\sum_{l=1}^K n_l^{(m)}} = \frac{n_k^{(m)}}{N} \quad \leftarrow$$

$$\boldsymbol{\mu}_k^{(m+1)} = \frac{1}{n_k^{(m)}} \sum_{i=0}^{N-1} \gamma_{ik}^{(m)} \mathbf{y}_i \quad \leftarrow$$

$$\mathbf{C}_k^{(m+1)} = \frac{1}{n_k^{(m)}} \sum_{i=0}^{N-1} \gamma_{ik}^{(m)} (\mathbf{y}_i - \boldsymbol{\mu}_k)(\mathbf{y}_i - \boldsymbol{\mu}_k)^T \quad \leftarrow$$



---

S4 收敛性验证, 计算  $l(\boldsymbol{\theta}^{(m+1)}|\boldsymbol{y})$  并检查下式

$$\left| l(\boldsymbol{\theta}^{(m+1)}|\boldsymbol{y}) - l(\boldsymbol{\theta}^{(m)}|\boldsymbol{y}) \right| < \varepsilon \quad \leftarrow$$

若成立则停止, 否则  $m = m + 1$  转 S2。  $\leftarrow$



# GMM参数估计的EM算法实例

Old Faithful 数据集

