

EM

与 Kmeans 不同，高斯混合聚类采用概率模型表达聚类。先介绍 EM：

考虑 $p(X, Z | \theta) = \prod_{n=1}^N p(x_n, z_n)$ ，其中：

observed: $X = \{x_n\}_{n=1}^N$; latent: $Z = \{z_n\}_{n=1}^N$

Goal: Estimate θ via MLE (or MAP).

$$\hat{\theta} = \operatorname{argmax}_{\theta} \log p(X | \theta) = \operatorname{argmax}_{\theta} \log \sum_z p(X, Z | \theta) \quad (\text{discrete})$$

$$= \operatorname{argmax}_{\theta} \log \int_Z p(X, Z | \theta) dZ. \quad (\text{continuous})$$

$$\log p(X | \theta) = \log \sum_z p(X, Z | \theta) = \log \sum_z q_n(z) \frac{p(X, Z | \theta)}{q_n(z)}$$

$$\geq \sum_z q_n(z) \log \frac{p(X, Z | \theta)}{q_n(z)} \quad (\text{since } \log(\cdot) \text{ is concave; Jensen})$$

$$= \sum_z q_n(z) \log p(X, Z | \theta) - \underbrace{\sum_z q_n(z) \log q_n(z)}_{\text{constant! since irrelevant with } \theta}$$

(if $q_n(z) = p(Z | X, \theta)$ ，则为等号)

故令 $q_n(z) = p(Z | X, \theta)$ (独立): $\log p(X | \theta) = [\mathbb{E}[\log p(X, Z | \theta)]] + \text{const}$

* lower-bound factor! EM 就是最大化它！

EM 流程: E: 根据当前 θ^t 推后验 $P(Z | X, \theta^t)$, 计算: 后验 先验

$$Q(\theta | \theta^t) = \mathbb{E}_{Z|X, \theta^t} [\log p(Z, X | \theta)] = \sum_z p(z | X, \theta^t) \log p(X, z | \theta)$$

$$M: \theta^{t+1} = \operatorname{argmax}_{\theta} \{ Q(\theta | \theta^t) + [\log p(\theta)] \} \quad \text{若加: MAP; 不加: MLE}$$

latent variable; cluster 中, $|z|=k$

GMM: 假设生成 N 个点: $x_n, n = 1, 2, \dots, N$

$$\sum_{k=1}^K \pi_k = 1.$$

首先从 K 个高斯中选一个。此选择基于混合成分先验概率 π 的。

$$z_n \sim \text{Multinomial}(z_n | \pi) \quad (\text{多项式分布})$$

然后生成点服从 $N(x_n | \mu_k, \Sigma_k)$, 其中 μ_k 是第 k 个高斯的值,

Σ_k 是第 k 个高斯的协方差 (suppose $z_n = k$)

说白了就是在 $\{\pi_i\}_{i=1}^K$ 的概率下, 抽一个 z_i



上述描述了在有 π_i, μ_i, Σ_i 下，如何抽 n 个点出来。

那么学 GMM 呢？考虑：

observed: $X = \{x_1, \dots, x_N\}$ latent: $Z = \{z_1, \dots, z_N\}$

$z_i \in \{1, \dots, K\} \Rightarrow N$ 点分 K 类

则 complete-data log-likelihood:

$$\log p(x, z | \theta) = \sum_{i=1}^N \log \pi_{zi} + \log N(x_i | \mu_{zi}, \Sigma_{zi})$$

则 E-step: posterior: $\gamma_{ik} = p(z_i=k | x_i, \theta^{old})$

$$\text{则 } \gamma_{ik} = \frac{p(x_i | z_i=k, \theta^{old})}{\sum_{j=1}^K p(x_i | \theta^{old})} \frac{p(z_i=k | \theta^{old})}{p(x_i | \theta^{old})}$$

$$= \frac{\pi_k^{old} N(x_i | \mu_k^{old}, \Sigma_k^{old})}{\sum_{j=1}^K \pi_j^{old} N(x_i | \mu_j^{old}, \Sigma_j^{old})}$$

Bring in Gaussian

M-step: $Q(\theta, \theta^{old}) = \mathbb{E}_{z|x, \theta^{old}} [\log p(x, z | \theta)]$

$$= \sum_{i=1}^N \sum_{k=1}^K \underbrace{\gamma_{ik}}_{\text{先验}} \underbrace{[\log \pi_k + \log N(x_i | \mu_k, \Sigma_k)]}_{\text{后验}}, \text{ constraint: } \sum_k \pi_k = 1$$

$$L = Q(\theta, \theta^{old}) + \lambda (1 - \sum_{k=1}^K \pi_k)$$

$$\frac{\partial L}{\partial \pi_k} = \sum_{i=1}^N \frac{\gamma_{ik}}{\pi_k} - \lambda = 0 \Rightarrow \pi_k \propto \sum_{i=1}^N \gamma_{ik}$$

$$\text{又 } \sum_k \pi_k = 1, \text{ 故 } \sum_{k=1}^K \sum_{i=1}^N \gamma_{ik} = 1$$

$$\left\{ \begin{array}{l} \pi_k^{new} = \frac{1}{N} \sum_{i=1}^N \gamma_{ik}; \text{ 同理: } \frac{\partial Q}{\partial \mu_k} = \sum_{i=1}^N \gamma_{ik} \sum_k (x_i - \mu_k)^{-1} \geq 0 \\ \mu_k^{new} = \frac{\sum_{i=1}^N \gamma_{ik} x_i}{\sum_{i=1}^N \gamma_{ik}} \\ \Sigma_k^{new} = \frac{\sum_{i=1}^N \gamma_{ik} (x_i - \mu_k^{new})(x_i - \mu_k^{new})^T}{\sum_{i=1}^N \gamma_{ik}} \end{array} \right.$$



附：PCA：主成分分析 $\rightarrow \mathbb{R}^{d'}$

Input: $D = \{x_1, x_2, \dots, x_m\}$, 低维空间维数 d'

Algorithm: 中心化: $x_i \leftarrow x_i - \frac{1}{m} \sum_{i=1}^m x_i$

计算协方差矩阵: XX^\top

特征值分解

取最大 d' 个特征值对应的 eigenvector $w_1, w_2, \dots, w_{d'}$

\Rightarrow Output $w^* = (w_1, w_2, \dots, w_{d'}) \in \mathbb{R}^{d' \times d'}$

$x_i \cdot w^* \top$ 后便进入低维空间

PCA 是最常用的降维方法！

