

高斯混合模型

$$p(x) = \sum_{k=1}^K \pi_k N(x | \mu_k, \Sigma_k)$$

$$\begin{cases} N(x | \mu_k, \Sigma_k) = \frac{1}{\sqrt{(2\pi)^d |\Sigma_k|}} \exp \left\{ -\frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) \right\} \\ \sum_{k=1}^K \pi_k = 1 \end{cases}$$

用极大似然法估计概率密度

输入: $\{x_i | i=1 \sim N\}$

最小化: $E(\{\pi_k, \mu_k, \Sigma_k | k=1 \sim K\})$

$$= - \sum_{i=1}^N \log [p(x_i)]$$

$$= - \sum_{i=1}^N \log \left[\sum_{k=1}^K \pi_k \cdot \frac{1}{\sqrt{(2\pi)^d |\Sigma_k|}} \exp \left\{ -\frac{1}{2} (x_i - \mu_k)^T \Sigma_k^{-1} (x_i - \mu_k) \right\} \right]$$

非凸问题, 无法求全局极值, 只能求局部极值

↓

① 梯度下降

② 启发式方法

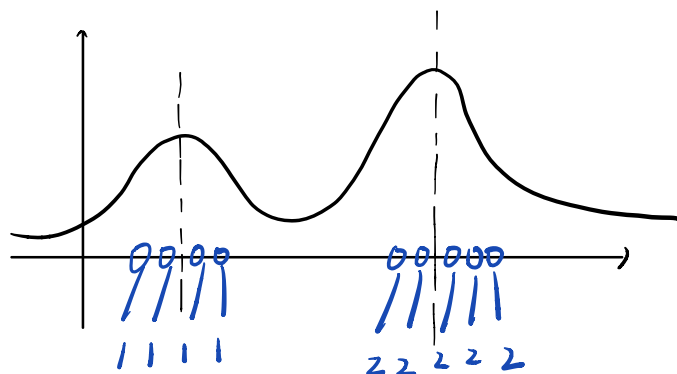
③ EM算法

1. 不需调参一定收敛

2. 编程简单

3. 理论优美

EM算法:



随机化所有 $\{\pi_k, \mu_k, \Sigma_k\} \quad k=1 \sim K$
 $K = \{1, 2\}$

高斯混合模型 EM (Expectation-Maximization)

① 随机化 $\{\pi_k, \mu_k, \Sigma_k\} \quad k=1 \sim K$

② E-Step:

$$V_{nk} = \frac{\pi_k \cdot N(X_n | \mu_k, \Sigma_k)}{\sum_{k=1}^K \pi_k N(X_n | \mu_k, \Sigma_k)} \quad (n=1 \sim N)$$

(第 n 个样本落在第 k 个高斯的概率)

③ 更新: M-Step

$$N_k = \sum_{n=1}^N V_{nk} \quad (\text{所有 } N \text{ 个样本中有多少属于第 } k \text{ 个样本})$$

$$\pi_k^{(new)} = \frac{N_k}{N} \quad (\pi_k: \text{第 } k \text{ 个高斯的概率})$$

$$\mu_k^{(new)} = \frac{1}{N_k} \sum_{n=1}^N V_{nk} \cdot X_n \quad (\mu_k: \text{第 } k \text{ 个高斯的均值})$$

$$\Sigma_k^{(new)} = \frac{1}{N_k} \sum_{n=1}^N V_{nk} (X_n - \mu_k^{(new)})(X_n - \mu_k^{(new)})^T$$

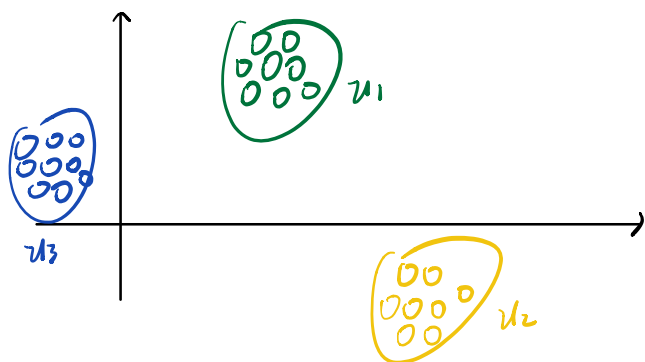
$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} (\mathbf{x}_n - \mathbf{u}_k) (\mathbf{x}_n - \mathbf{u}_k)^T$$

(Σ_k : 第k个高斯的协方差矩阵)

④ 回到②直到收敛

(应用: 高斯混合模型在说话人识别方面的应用)

K-均值聚类 (K-means clustering)



(如何让机器自动聚类?)

$$E(\{\mathbf{u}_k\}) = \sum_{k=1}^K \sum_{n=1}^N \|\mathbf{x}_i - \mathbf{u}_k\|^2$$

(第②③都使E变小, 而E有下界0
一定收敛)

问题: 输入N个样本 $\{\mathbf{x}_i\} \quad i=1 \sim N$

输出N个样本类别 $\{z_i\} \quad i=1 \sim N$

其中 $z_i = 1, 2, 3, \dots, K$

① 随机化 $\mathbf{u}_1, \dots, \mathbf{u}_K$

② E-Step:

$$z_i = \arg \min_k \|\mathbf{x}_i - \mathbf{u}_k\| \quad (\text{离谁近属于谁})$$

③ M-Step:

$$N_k = \sum_{i=1}^N \mathbb{I}(Z_i = k) \quad (\text{N个样本中有多少属于第k类})$$

$$\mu_k = \frac{1}{N_k} \sum_{\substack{i=1 \\ Z_i=k}}^N X_i \quad (\text{第k类样本的均值})$$

④ 回到② 直到收敛

(应用: 基于K-均值聚类的图像矢量化 / 压缩图形)