

X : observed data $\rightarrow X = (x_1, x_2, \dots, x_N)$
 (X, Z) : complete data

θ : parameter $\rightarrow \theta = \{p_1, p_2, \dots, p_K, \mu_1, \mu_2, \dots, \mu_K, \Sigma_1, \Sigma_2, \dots, \Sigma_K\}$

$$\hat{\theta}_{MLE} = \arg \max_{\theta} \log P(X) = \arg \max_{\theta} \log \prod_{i=1}^N P(x_i | z_i, \theta)$$

$$= \arg \max_{\theta} \sum_{i=1}^N \log P(x_i | z_i, \theta)$$

$$= \arg \max_{\theta} \sum_{i=1}^N \log \left[\sum_{k=1}^K p_k N(x_i | \mu_k, \Sigma_k) \right]$$

混合高斯分布

无法直接得到解析解，所以用 EM 算法求近似解。

使用 EM 算法。

$$\theta^{t+1} = \arg \max_{\theta} \underbrace{E_{Z|X, \theta^t} [\log P(X, Z | \theta)]}_{Q(\theta, \theta^t)}$$

$$Q(\theta, \theta^t) = \int_Z \log P(X, Z | \theta) P(Z | X, \theta^t) dZ$$

$$= \sum_{i=1}^N \log \prod_{z_i=1}^N P(x_i, z_i | \theta) \prod_{i=1}^N P(z_i | x_i, \theta^t)$$

$$= \sum_{z_1, z_2, \dots, z_N} \sum_{i=1}^N \log P(x_i, z_i | \theta) \cdot \prod_{i=1}^N P(z_i | x_i, \theta^t)$$

z_1 : 第一个样本属于各个分布的概率

$$\sum_{z_1, z_2, \dots, z_N} \left[\log P(x_1, z_1 | \theta) + \log P(x_2, z_2 | \theta) + \dots + \log P(x_N, z_N | \theta) \right] \prod_{i=1}^N P(z_i | x_i, \theta^t)$$

$$\sum_{z_1, z_2, \dots, z_N} \log P(x_1, z_1 | \theta) \prod_{i=1}^N P(z_i | x_i, \theta^t)$$

$$= \sum_{z_1, z_2, \dots, z_N} \log P(x_1, z_1 | \theta) P(z_1 | x_1, \theta^t) \prod_{i=2}^N P(z_i | x_i, \theta^t)$$

$$= \sum_{z_1, z_2, \dots, z_N} \log P(x_1, z_1 | \theta) P(z_1 | x_1, \theta^t) \sum_{z_2, z_3, \dots, z_N} \prod_{i=2}^N P(z_i | x_i, \theta^t)$$

$$= \sum_{z_1, z_2, \dots, z_N} \log P(x_1, z_1 | \theta) P(z_1 | x_1, \theta^t) \underbrace{\sum_{z_2, z_3, \dots, z_N} \prod_{i=2}^N P(z_i | x_i, \theta^t)}_{=1}$$

$$= \sum_{z_1, z_2, \dots, z_N} \log P(x_1, z_1 | \theta) \cdot P(z_1 | x_1, \theta^t) + \dots +$$

$$\sum_{z_1, z_2, \dots, z_N} \log P(x_N, z_N | \theta) \cdot P(z_N | x_N, \theta^t)$$

$$= \sum_{i=1}^N \sum_{z_i} \log P(x_i, z_i | \theta) \cdot P(z_i | x_i, \theta^t)$$