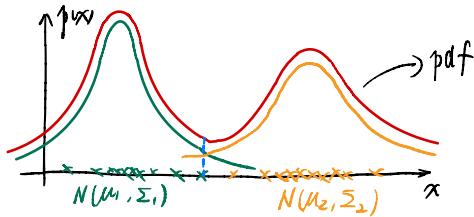


# 高斯混合模型 (GMM)

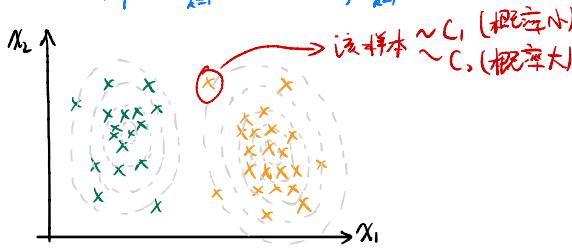
## 一、模型介绍

一般来说，将数据假设为高斯分布是合理的。



从几何角度来看：加权平均值 → (多个高斯分布叠加而成)

$$\text{则 } p(x) = \sum_{k=1}^K \alpha_k N(x|\mu_k, \Sigma_k), \quad \sum_{k=1}^K \alpha_k = 1, \quad \alpha_k \text{ 为权重.}$$



从混合模型角度来看：

1: 观测变量

2: 隐变量

→ 对应的样本  $x$  是属于哪一个高斯分布，因此它是 **离散的随机变量**。

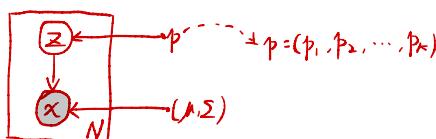
$$\begin{array}{c|cccc} z & C_1 & C_2 & \dots & C_K \\ \hline p(z) & p_1 & p_2 & \dots & p_K \end{array}, \quad \sum_{k=1}^K p_k = 1$$

GMM 为生成模型，其生成过程为：

① 假设有一个骰子，其每个面都有涂料，但涂料的重量不同，即指这个骰子每个面出现的概率不同，每个面出现的概率为  $p_k$ 。

② 接下来根据出现概率为  $p_k$  的某个高斯分布，在这个高斯分布上采样。

其概率图模型可表示为：



混合模型角度  $P(X)$  为：

$$\begin{aligned} P(X) &= \sum_{\mathbf{Z}} P(X, \mathbf{Z}) \\ &= \sum_{k=1}^K P(X, Z=c_k) \\ &= \sum_{k=1}^K P(Z=c_k) \cdot P(X|Z=c_k) \\ &= \sum_{k=1}^K p_k \cdot N(\mathbf{x} | \mu_k, \Sigma_k) \end{aligned}$$

$\Delta$  根率值.

求参数：

假设：  $X$ : 观测数据，  $X = (X_1, X_2, \dots, X_N)$

$(X, Z)$ : 完备数据

$\theta$ : 参数，  $\theta = \{\pi_1, \pi_2, \dots, \pi_K, \mu_1, \mu_2, \dots, \mu_K, \Sigma_1, \Sigma_2, \dots, \Sigma_K\}$

则：

$$\begin{aligned} \hat{\theta}_{MLE} &= \arg \max \log P(X) = \arg \max \log \prod_{i=1}^N P(X_i) = \arg \max \sum_{i=1}^N \log P(X_i) \\ &= \arg \max \sum_{i=1}^N \log \sum_{k=1}^K p_k \cdot N(x_i | \mu_k, \Sigma_k) \end{aligned}$$

由于  $\log(\Delta + \Delta + \Delta + \dots + \Delta)$  中是连加符号，因此，MLE 不能求出其解析解。  
1. 用 MLE 求解 GMM，无法得出解析解。

应用 EM 算法求解： ( $EM: \theta^{(t+1)} = \arg \max_{\theta} E_{Q(\theta, \theta^{(t)})} [\log P(X, Z | \theta)]$ )  $\theta^{(t)}$  是常数。

E-step:

$$\begin{aligned} Q(\theta, \theta^{(t)}) &= \int_Z \log P(X, Z | \theta) \cdot P(Z | X, \theta^{(t)}) dz \\ &= \sum_{z_1}^N \log \prod_{i=1}^N P(X_i, z_i | \theta) \cdot \prod_{i=1}^N P(z_i | x_i, \theta^{(t)}) \quad \text{Zi 在 GMM 中由等式} \\ &= \sum_{z_1, z_2, \dots, z_N} \sum_{i=1}^N \log P(X_i, z_i | \theta) \cdot \prod_{i=1}^N P(z_i | x_i, \theta^{(t)}) \\ &= \sum_{z_1, z_2, \dots, z_N} [\log P(X_1, z_1 | \theta) + \log P(X_2, z_2 | \theta) + \dots + \log P(X_N, z_N | \theta)] \prod_{i=1}^N P(z_i | x_i, \theta^{(t)}) \end{aligned}$$

其中，项： $\sum_{z_1, z_2, \dots, z_N} \log P(X_i, z_i | \theta) \stackrel{P(Z_i | X_i, \theta^{(t)})}{=} P(Z_i | X_i, \theta^{(t)})$

$$\begin{aligned} &= \sum_{z_1} \log P(X_1, z_1 | \theta) P(z_1 | X_1, \theta^{(t)}) \sum_{z_2} \prod_{i=2}^N P(z_i | X_i, \theta^{(t)}) \rightarrow = \sum_{z_1} P(z_1 | X_1, \theta^{(t)}) \sum_{z_2} P(z_2 | X_2, \theta^{(t)}) \cdots \sum_{z_N} P(z_N | X_N, \theta^{(t)}) \\ &= \sum_{z_1} \log P(X_1, z_1 | \theta) P(z_1 | X_1, \theta^{(t)}) \end{aligned}$$

$$\begin{aligned} \text{因此：} &= \sum_{z_1} \log P(X_1, z_1 | \theta) P(z_1 | X_1, \theta^{(t)}) + \dots + \sum_{z_N} \log P(X_N, z_N | \theta) P(z_N | X_N, \theta^{(t)}) \\ &= \sum_{i=1}^N \sum_{z_i} \log P(X_i, z_i | \theta) P(z_i | X_i, \theta^{(t)}) \end{aligned}$$

$$\begin{aligned}\therefore P(X) &= \sum_{k=1}^K p_k \cdot N(x| \mu_k, \Sigma_k) \\ \therefore P(X, Z) &= P(Z) \cdot P(X|Z) \\ &= p_z \cdot N(x| \mu_z, \Sigma_z) \\ \therefore P(Z|X) &= \frac{P(X, Z)}{P(X)} = \frac{p_z \cdot N(x| \mu_z, \Sigma_z)}{\sum_{k=1}^K p_k N(x| \mu_k, \Sigma_k)}\end{aligned}$$

$$\text{则 } Q(\theta, \theta^{(t)}) = \sum_{i=1}^N \sum_{z_i} \log [p_{z_i} \cdot N(x_i | \mu_{z_i}, \Sigma_{z_i})] \cdot \frac{p_{z_i} \cdot N(x_i | \mu_{z_i}, \Sigma_{z_i})}{\sum_{k=1}^K p_k N(x_i | \mu_k, \Sigma_k)}$$

M-step:

由 E-step:

$$\begin{aligned}Q(\theta, \theta^{(t)}) &= \sum_{i=1}^N \sum_{z_i} \log [p_{z_i} \cdot N(x_i | \mu_{z_i}, \Sigma_{z_i})] \cdot P(z_i | x_i, \theta^{(t)}) \\ &= \sum_{z_i} \sum_{i=1}^N \log [p_{z_i} \cdot N(x_i | \mu_{z_i}, \Sigma_{z_i})] \cdot P(z_i | x_i, \theta^{(t)}) \\ &= \sum_{k=1}^K \sum_{i=1}^N \log [p_k \cdot N(x_i | \mu_k, \Sigma_k)] \cdot P(z_i = c_k | x_i, \theta^{(t)}) \\ &= \sum_{k=1}^K \sum_{i=1}^N \log [p_k + \log N(x_i | \mu_k, \Sigma_k)] \cdot P(z_i = c_k | x_i, \theta^{(t)})\end{aligned}$$

$$\begin{aligned}\theta^{(t+1)} &= \arg \max_{\theta} Q(\theta, \theta^{(t)}) \\ \text{求 } p_k^{(t+1)}: \\ \left\{ \begin{array}{l} p_k^{(t+1)} = \arg \max_{p_k} \sum_{i=1}^N \log p_k \cdot P(z_i = c_k | x_i, \theta^{(t)}) \\ \text{s.t. } \sum_{k=1}^K p_k = 1 \end{array} \right.\end{aligned}$$

应用拉格朗日乘子法:

$$L(p, \lambda) = \sum_{k=1}^K \sum_{i=1}^N \log p_k \cdot P(z_i = c_k | x_i, \theta^{(t)}) + \lambda (\sum_{k=1}^K p_k - 1)$$

$$\text{对 } p_k \text{ 求导} \rightarrow \frac{\partial L}{\partial p_k} = \sum_{i=1}^N \frac{1}{p_k} \cdot P(z_i = c_k | x_i, \theta^{(t)}) + \lambda \stackrel{=} 0$$

求导与中  $p_1, p_K$  等式两边同乘  $p_k$ .

$$\cdots p_1 \text{ 都无关.} \implies \sum_{i=1}^N P(z_i = c_k | x_i, \theta^{(t)}) + p_k \lambda = 0$$

$$\implies \underbrace{\sum_{i=1}^N \sum_{k=1}^K}_{\text{(对 } z_i \text{ 求和)}} P(z_i = c_k | x_i, \theta^{(t)}) + \underbrace{\sum_{k=1}^K p_k \lambda}_{\lambda} = 0$$

$$\implies N + \lambda = 0 \implies \lambda = -N.$$

$$\text{故 } p_k^{(t+1)} = \frac{1}{N} \sum_{i=1}^N P(z_i = c_k | x_i, \theta^{(t)})$$

$$p^{(t+1)} = (p_1^{(t+1)}, p_2^{(t+1)}, \dots, p_K^{(t+1)})$$