

高斯混合模型

Gaussian Mixture Model

- ① 介绍
- ② 极大似然估计
- ③ EM 算法

介绍 (Introduction)

K均值聚类:

构建的目标函数: $J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2$

其中 r_{nk} 表示 x_n 属于哪一类, μ_k 表示该类的聚集中心.

显然, 目标就是最小化 J .

方法也很简单. 先固定 μ_k , 更新 r_{nk} ; 再固定 r_{nk} , 更新 μ_k . 类似 EM.

有时候 $\|x_n - \mu_k\|^2$ 这种度量会失效. 故将其看作一类函数 $V(x_n, \mu_k)$

则引中为 $J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} V(x_n, \mu_k)$.

同样 现在我们将 GMM 看作是简单的 线性叠加:

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

下面, 我们从概率角度看 GMM

我们引入一个隐变量 z . z 有 K 个分量, 将其看作离散随机变量.

分布如下:

每个分量取值为 0 或 1, 且 $\sum_{k=1}^K z^{(k)} = 1$.

z	$z^{(1)}$	$z^{(2)}$	$z^{(3)}$	\dots	$z^{(K)}$
p	π_1	π_2	π_3	\dots	π_K

$$\text{即 } p(z^{(k)}=1) = \pi_k$$

当 $z^{(k)}=1$ 时, 表明当前的数据点属于第 k 个高斯

$$k. | p(z) = \prod_{k=1}^K \pi_k^{z^{(k)}}$$

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

分布. 和 K 均值聚类中的 r_{nk} 类似.

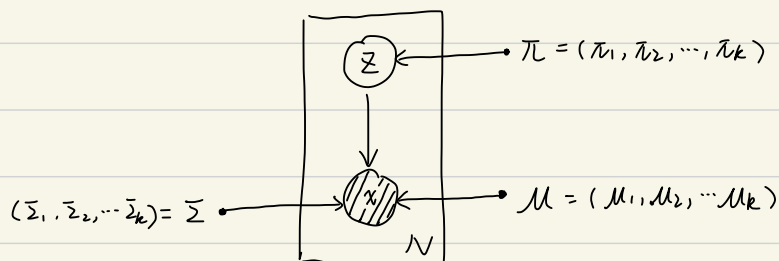
$$\text{即 } p(x | z) = \prod_{k=1}^K N(x | \mu_k, \Sigma_k)^{z^{(k)}}$$

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

则生成过程可以表示为先从 $p(z)$ 中采样得到 $z = z^{(k)}$.

再从 $p(x|z = z^{(k)})$ 中采样得到 x .

一组 N 个独立同分布数据点, 其概率图模型表示方法如下.



阴影表示观测变量
空心点表示参数.

同时, 由联合概率可以求得 x 的边缘概率.

$$p(x) = \sum_z p(x, z) = \sum_{k=1}^K p(x, z = z^{(k)})$$

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

和线性叠加的表达式一样.

极大似然估计 (MLE)

为了估计参数 π, μ, Σ , 我们用 MLE.

X : observed data

(X, Z) : complete data

$\theta: \pi, \mu, \Sigma$

$$\log P(X|\theta) = \sum_{i=1}^N \log \sum_{k=1}^K \pi_k N(x_i | \mu_k, \Sigma_k)$$

我将从下面 3 个方面来说明直接应用 MLE 的不可行性.

① 奇异值 (singularities)

假设 K 个高斯分布中, 某个高斯分布的 $\Sigma_k = G_k^2 I$, (一般的协方差矩阵率也会出现下述现象, 但为了方便计算与理解, 这里取 $\Sigma_k = G_k^2 I$)

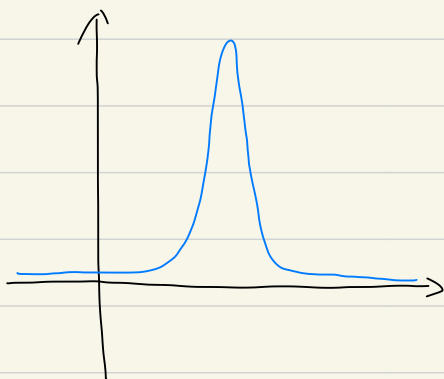
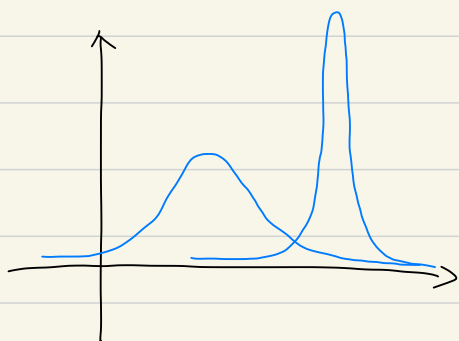
当出现某个数据点, $x_n = \mu_k$, 则高斯分布会退化为

$$N(x_n | \mu_k, \Sigma_k) = \frac{1}{(2\pi)^{\frac{D}{2}}} \cdot \frac{1}{G_j^D}, \quad D \text{ 为高斯分布的维度.}$$

我们注意到, 对于该高斯分布, 当 $G_j \rightarrow 0$ 时, $N(x_n | \mu_k, \Sigma_k) \rightarrow \infty$. 这就导致了 $\log\text{-likelihood} \rightarrow \infty$

从而, 在 K 个分量中, 其中一个分量趋近于无穷, 其余的是一个有限的值, 导致求解错误.

那为何在单一的高斯分布上不会出现该情况呢？



如上，两幅图中，左边是 GMM，右边是单一高斯分布。
可以看出，在 GMM 中，两个分量之间是加权求和的操作，
也就是说，如果一个数据点导致一个分量为无穷，其他数
据点仍然会有一个有限的值，这就导致了总的 likelihood
趋近于无穷。

在单一的高斯分布中，若出现 x_n ，使得当 $\sigma \rightarrow 0$ 时，当前
数据点的似然趋近于无穷，但注意 $p(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\{-\frac{1}{2\sigma^2}(x-\mu)^T(x-\mu)\}$ 。
对于其他数据点，已的指数项中 $x-\mu \neq 0$ ，也就导致了，当 $\sigma \rightarrow 0$ 时，
其他数据点的似然会以指数速度趋向于 0，从而使得整体
的 likelihood 不会趋近于无穷。

② Identifiability

由于总的概率是由 K 个分量混合而成的，这就导致了，对于任意给定的最大似然解（有 K 个分量），我们可以对这 K 个分量进行全排列，则有 $K!$ 种组合方式，也就是说有 $K!$ 种解都能得到该似然，这也就违反了统计中 identifiability 的概念。这个概念在表示模型参数时，是非常重要的。（PRML 中介绍的，但我不理解的是 $K!$ 种解只是换了个位置，为何会对模型参数的表示产生问题）。

③ 解析解

我们注意到 $\log P(X|\theta)$ 中， \log 中有对 k 的求和，从而 \log 不直接作用于高斯分布，当令其导数为 0 时，无法得到解析解。

（视频和 PRML 中都下了此结论，但我实际还没算过）

EM算法

$$EM: \theta^{(t+1)} = \arg \max_{\theta} \mathbb{E}_{z \sim p(z|x, \theta^{(t)})} [\log p(x, z | \theta)]$$

注意，K均值算法常被用于EM之前的初始化高斯混合模型的参数。

E-step:

$$\begin{aligned} Q(\theta, \theta^{(t)}) &= \mathbb{E}_{z \sim p(z|x, \theta^{(t)})} [\log p(x, z | \theta)] \\ &= \sum_{k=1}^K \left[\sum_{i=1}^N p(z_i^{(k)} = 1 | x_i, \theta^{(t)}) \log p(x_i, z_i^{(k)} = 1 | \theta) \right] \\ &= \sum_{k=1}^K \sum_{i=1}^N \frac{\pi_k N(x_i | \mu_k, \Sigma_k^{(t)})}{\sum_{j=1}^K \pi_j N(x_i | \mu_j^{(t)}, \Sigma_j^{(t)})} \log \pi_k N(x_i | \mu_k, \Sigma_k) \end{aligned}$$

M-step:

$$\theta^{(t+1)} = \arg \max_{\theta} \sum_{k=1}^K \sum_{i=1}^N \frac{\pi_k N(x_i | \mu_k^{(t)}, \Sigma_k^{(t)})}{\sum_{j=1}^K \pi_j N(x_i | \mu_j^{(t)}, \Sigma_j^{(t)})} \log \pi_k N(x_i | \mu_k, \Sigma_k)$$

以求解 π 为例：

$$\begin{cases} \pi^{(t+1)} = \arg \max_{\pi} \sum_{k=1}^K \sum_{i=1}^N \frac{\pi_k N(x_i | \mu_k^{(t)}, \Sigma_k^{(t)})}{\sum_{j=1}^K \pi_j N(x_i | \mu_j^{(t)}, \Sigma_j^{(t)})} \log \pi_k \\ \text{s.t. } \sum_{k=1}^K \pi_k = 1 \end{cases}$$

拉格朗日函数：
$$\mathcal{L}(\pi, \lambda) = \sum_{k=1}^K \sum_{i=1}^N \frac{\pi_k N(x_i | \mu_k^{(t)}, \Sigma_k^{(t)})}{\sum_{j=1}^K \pi_j N(x_i | \mu_j^{(t)}, \Sigma_j^{(t)})} \log \pi_k + \lambda (1 - \sum_{k=1}^K \pi_k)$$

$$\frac{\partial \mathcal{L}}{\partial \pi_k} = \frac{N}{\sum_{j=1}^K \pi_j} \frac{\pi_k N(x_i | \mu_k^{(t)}, \Sigma_k^{(t)})}{N(x_i | \mu_j^{(t)}, \Sigma_j^{(t)})} \cdot \frac{1}{\pi_k} - \lambda = 0$$

$$\Rightarrow \sum_{i=1}^N \frac{\pi_k N(x_i | \mu_k^{(t)}, \Sigma_k^{(t)})}{\sum_{j=1}^K \pi_j N(x_i | \mu_j^{(t)}, \Sigma_j^{(t)})} - \lambda \pi_k = 0$$

对 π , 到 π_k 都求偏导, 得到 K 个等式之后, 将其相加, 得

$$\sum_{i=1}^N \sum_{k=1}^K \frac{\pi_k N(x_i | \mu_k^{(t)}, \Sigma_k^{(t)})}{\sum_{j=1}^K \pi_j N(x_i | \mu_j^{(t)}, \Sigma_j^{(t)})} - \lambda \sum_{k=1}^K \pi_k = 0$$

~~~~~
1 (条件)

注意, 这个分数是  $p(z_i^{(t)} = 1 | x_i, \theta^{(t)})$ , 求和后为 1 (直接算也是 1, 利用概率的性质也是 1)

$$\Rightarrow N - \lambda = 0 \Rightarrow \lambda = N$$

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

另外, 两个参数结论如下:

$$\begin{cases} \mu_k^{(t+1)} = \frac{1}{N_k^{(t+1)}} \sum_{n=1}^N \gamma^{(t)}(z_n^{(k)}) x_n \\ \Sigma_k^{(t+1)} = \frac{1}{N_k^{(t+1)}} \sum_{n=1}^N \gamma^{(t)}(z_n^{(k)}) (x_n - \mu_k^{(t+1)})(x_n - \mu_k^{(t+1)})^T \end{cases}$$

$$\text{其中, } \gamma^{(t)}(z_n^{(k)}) = \frac{\pi_k N(x_n | \mu_k^{(t)}, \Sigma_k^{(t)})}{\sum_{j=1}^K \pi_j N(x_n | \mu_j^{(t)}, \Sigma_j^{(t)})}, \quad N_k^{(t+1)} = \sum_{i=1}^N \gamma^{(t)}(z_i^{(k)})$$

(结论出自于 PRML, 有空也可以看一下 PRML 上面有关这一节的介绍)