

EM 算法

EM 算法的推导

EM 算法的推广

EM算法的标准定义

输入：观测数据Y, 隐变量, Z参数 θ , 联合分布 $P(Y, Z | \theta)$, 条件分布 $P(Y | Z, \theta)$

输出： 模型参数 $\hat{\theta}$

(1) 选择参数的初始值 $\theta^{(0)}$, 开始迭代:

(2) E步: 以 $\theta^{(t)}$ 为每次迭代参数的估计值, 计算第t+1次迭代的E步, 计算

$$Q(\theta, \theta^{(t)}) = E_z[\log P(Y, Z | \theta) | Y, \theta^{(t)}] = \sum_z \log P(Y, Z | \theta) p(Z | Y, \theta^{(t)})$$

这里, $P(Z | Y, \theta^{(t)})$ 是在给定观测数据Y和当前的参数估计 $\theta^{(t)}$ 下隐变量数据Z的条件概率分布;

(3) M步: 求取 $Q(\theta, \theta^{(t)})$ 极大化的 θ , 确定第t+1次迭代的参数的估计值 $\theta^{(t+1)}$

$$\theta^{(t+1)} = \arg \max_{\theta} Q(\theta, \theta^{(t)})$$

(4) 重复第(2)步和第(3)步, 直到收敛:

$Q(\theta, \theta^{(t)}) = E_z[\log P(Y, Z | \theta) | Y, \theta^{(t)}]$
式 $= \sum_z \log P(Y, Z | \theta) p(Z | Y, \theta^{(t)})$ 可记为 $Q(\theta, \theta^{(t)})$ 是EM算法的核心, 称为Q函数 (Q function).

说明

下面关于EM算法有几点说明:

步骤 (1) 参数的初始值可以任意选择, 但要注意EM算法对初始值敏感数据。

步骤(2)E步 $Q(\theta, \theta^{(t)})$: Q函数式中Z表示观测数据, Y表示隐数据, 注意 $Q(\theta, \theta^{(t)})$ 对第t个Z表示条件极大化的参数, 第t个Z表示参数的当前估计值, 每次迭代时求取Q函数及其极大。

步骤(3)M步: 求 $Q(\theta, \theta^{(t)})$ 极大化, 得到 $\theta^{(t+1)}$, 完成一次迭代 $\theta^{(t)} \rightarrow \theta^{(t+1)}$, 需要保证每次迭代似然函数值或Q函数值增大。

步骤(4)给出停止迭代的条件, 一般是对较小的正数 ϵ , θ_k, ϵ 满足

$$|\theta^{(t)} - \theta^{(t-1)}| < \epsilon \quad \text{或} \quad |Q(\theta^{(t)}, \theta^{(t-1)}) - Q(\theta^{(t-1)}, \theta^{(t-1)})| < \epsilon$$

则停止迭代。

GMM使用若干个高斯分布的加权之和作为对观测数据集进行建模的基础分布, 而由中心极限定理我们知道, 大量独立同分布的随机变量的均值在做适当标准化之后会依分布收敛于高斯分布, 这使得高斯分布具有普适性的建模能力, 继而奠定了使用高斯分布作为主要构成部件的GMM进行数据建模的理论基础。

EM算法通过迭代地构造似然函数下界的方式不断地提升似然函数的取值, 从而完成对含有隐变量模型的参数估计, 其典型的应用包括GMM、HMM (Hidden Markov Model, 隐马尔可夫模型) 的参数估计

GMM模型使用 K 个高斯分布的加权之和作为其概率密度函数, 具体地

$$p(x; \Theta) = \sum_{k=1}^K \alpha_k \cdot \mathcal{N}(x; \bar{\mu}_k, \Sigma_k)$$

其中参数表示为:

$$\Theta = \{\alpha_k, \bar{\mu}_k, \Sigma_k\}_{k=1}^K$$

为保证概率密度为1, 应该满足

$$\int_{-\infty}^{+\infty} p(x) dx = \sum_{k=1}^K \left[\alpha_k \int_{-\infty}^{+\infty} \mathcal{N}(x; \bar{\mu}_k, \Sigma_k) dx \right] = \sum_{k=1}^K \alpha_k = 1 \quad (2)$$

$\alpha_k, k = 1, \dots, K$ 即为隐变量 Z 所服从的分布 Q 的初始值, 即 $\alpha_k = Q(Z = k), k \in \{1, \dots, K\}$. 由(1)可以推出 GMM模型对不完全数据分布的似然估计通过求似然函数 $p(x; \Theta) = \sum_{k=1}^K p(x; \Theta) p(Z = k | \Theta)$ 得到。

样本 x_i 生成的概率可表示为

$$p(x_i; \Theta) = \sum_{k=1}^K p(x_i, z_i; \Theta) = \sum_{k=1}^K \mathcal{N}(x_i | Z = k; \bar{\mu}_k, \Sigma_k) Q(Z = k) = \sum_{k=1}^K \alpha_k \cdot \mathcal{N}(x_i | Z = k; \bar{\mu}_k, \Sigma_k)$$

但是在实际问题中, 常常存在隐变量, 上式不能完全使用, 变换后为:

$$\begin{aligned} \mathcal{N} \mathcal{L}(\Theta) &= -\log \prod_{i=1}^N p(x_i; \Theta) \\ &= -\sum_{i=1}^N \log p(x_i; \Theta) \\ &= -\sum_{i=1}^N \log \left[\sum_{k=1}^K p(x_i, z_i; \Theta) \right] \\ &= -\sum_{i=1}^N \log \left[\sum_{k=1}^K \alpha_k \cdot \mathcal{N}(x_i | z_i; \Theta) \right] \end{aligned}$$

$\min \mathcal{N} \mathcal{L}(\Theta)$
s.t. $\sum_{k=1}^K \alpha_k = 1$

受Jensen's inequality的启发, 如果我们能在对数函数的连加操作里面构造出一个关于 Z 的期望, 那我们就能将连加操作移至对数函数的外面了, 具体地, 对于对数似然函数我们有

$$\begin{aligned} \mathcal{L}(\Theta) &= \sum_{i=1}^N \log \left[\sum_{k=1}^K p(x_i, z_i; \Theta) \right] \\ &= \sum_{i=1}^N \log \left[\sum_{k=1}^K p(z_i | x_i; \Theta_{i-1}) \cdot \frac{p(x_i, z_i; \Theta)}{p(z_i | x_i; \Theta_{i-1})} \right] \\ &\geq \sum_{i=1}^N \sum_{k=1}^K p(z_i | x_i; \Theta_{i-1}) \cdot \log \left[\frac{p(x_i, z_i; \Theta)}{p(z_i | x_i; \Theta_{i-1})} \right] \end{aligned}$$

需要注意的是, 由于对数函数的凹凸性正好相反, 所以此处 Jensen's inequality中的符号应该取反., 我们借助在第 t-1次迭代中得到的参数估计值 Θ_{i-1} 来获得 x_i 关于 z_i 的后验分布

$$p(z_i | x_i; \Theta_{i-1}) = \frac{p(x_i, z_i; \Theta_{i-1}) \cdot p(z_i; \Theta_{i-1})}{\sum_{k=1}^K p(x_i, z_i; \Theta_{i-1}) \cdot p(z_i; \Theta_{i-1})}$$

再使用这个关于 z_i 的后验分布构造期望

$$\mathbf{E}_{z_i} \left[\frac{p(x_i, z_i; \Theta)}{p(z_i | x_i; \Theta_{i-1})} \right] = \sum_{k=1}^K p(z_i | x_i; \Theta_{i-1}) \cdot \frac{p(x_i, z_i; \Theta)}{p(z_i | x_i; \Theta_{i-1})}$$

而后由Jensen's inequality即可得到

$$\log \mathbf{E}_{z_i} \left[\frac{p(x_i, z_i; \Theta)}{p(z_i | x_i; \Theta_{i-1})} \right] \geq \mathbf{E}_{z_i} \left[\log \frac{p(x_i, z_i; \Theta)}{p(z_i | x_i; \Theta_{i-1})} \right] \quad (16)$$
$$\log \left[\sum_{k=1}^K p(z_i | x_i; \Theta_{i-1}) \cdot \frac{p(x_i, z_i; \Theta)}{p(z_i | x_i; \Theta_{i-1})} \right] \geq \sum_{k=1}^K p(z_i | x_i; \Theta_{i-1}) \cdot \log \left[\frac{p(x_i, z_i; \Theta)}{p(z_i | x_i; \Theta_{i-1})} \right] \quad (17)$$
$$\mathcal{L}(\Theta_k) \geq \mathcal{L}(\Theta_{k-1}) \geq \mathcal{L}(\Theta_{k-2}) \geq \mathcal{L}(\Theta_{k-3}) \geq \dots \geq \mathcal{L}(\Theta_0) \quad (18)$$

我们记下届函数为:

$$B(\Theta_k, \Theta_{k-1}) = \sum_{i=1}^N \mathbf{E}_{z_i} \left[\log \frac{p(x_i, z_i; \Theta)}{p(z_i | x_i; \Theta_{k-1})} \right] = \sum_{i=1}^N \sum_{k=1}^K p(z_i | x_i; \Theta_{k-1}) \cdot \log \left[\frac{p(x_i, z_i; \Theta)}{p(z_i | x_i; \Theta_{k-1})} \right]$$

即对数似然函数与下届函数有

$$\mathcal{L}(\Theta) \geq B(\Theta_k, \Theta_{k-1})$$

我们记下届函数为:

$$B(\Theta_k, \Theta_{k-1}) = \sum_{i=1}^N \mathbf{E}_{z_i} \left[\log \frac{p(x_i, z_i; \Theta)}{p(z_i | x_i; \Theta_{k-1})} \right] = \sum_{i=1}^N \sum_{k=1}^K p(z_i | x_i; \Theta_{k-1}) \cdot \log \left[\frac{p(x_i, z_i; \Theta)}{p(z_i | x_i; \Theta_{k-1})} \right]$$

即对数似然函数与下届函数有

$$\mathcal{L}(\Theta) \geq B(\Theta_k, \Theta_{k-1})$$

算法

输入：观测数据 y_1, y_2, \dots, y_N ，高斯混合模型；
输出：高斯混合模型参数。

(1) 初始化参数 $\theta^{(0)}$, 开始迭代

(2) E步: 以 $\theta^{(t)}$ 为每次迭代参数的估计值, 计算第t+1次迭代的E步, 计算

$$Q(\theta, \theta^{(t)}) = E_z[\log P(Y, Z | \theta) | Y, \theta^{(t)}] = \sum_z \log P(Y, Z | \theta) p(Z | Y, \theta^{(t)})$$

(3) M步: 求 $Q(\theta, \theta^{(t)})$ 极大化的 θ , 确定第t+1次迭代的参数的估计值 $\theta^{(t+1)}$

$$\theta^{(t+1)} = \arg \max_{\theta} Q(\theta, \theta^{(t)})$$

(4) 重复第(2)步和第(3)步, 直到收敛。

1

输入：观测数据 Y, 隐变量, Z参数 θ , 联合分布 $P(Y, Z | \theta)$, 条件分布 $P(Y | Z, \theta)$

输出： 模型参数 $\hat{\theta}$

(1) 选择参数的初始值 $\theta^{(0)}$, 开始迭代:

(2) E步: 以 $\theta^{(t)}$ 为每次迭代参数的估计值, 计算第t+1次迭代的E步, 计算

$$Q(\theta, \theta^{(t)}) = E_z[\log P(Y, Z | \theta) | Y, \theta^{(t)}] = \sum_z \log P(Y, Z | \theta) p(Z | Y, \theta^{(t)})$$

(3) M步: 求取 $Q(\theta, \theta^{(t)})$ 极大化的 θ , 确定第t+1次迭代的参数的估计值 $\theta^{(t+1)}$

$$\theta^{(t+1)} = \arg \max_{\theta} Q(\theta, \theta^{(t)})$$

(4) 重复第(2)步和第(3)步, 直到收敛。

2

输入：观测数据 Y, 隐变量, Z参数 θ , 联合分布 $P(Y, Z | \theta)$, 条件分布 $P(Y | Z, \theta)$

输出： 模型参数 $\hat{\theta}$

(1) 选择参数的初始值 $\theta^{(0)}$, 开始迭代:

(2) E步: 以 $\theta^{(t)}$ 为每次迭代参数的估计值, 计算第t+1次迭代的E步, 计算

$$Q(\theta, \theta^{(t)}) = E_z[\log P(Y, Z | \theta) | Y, \theta^{(t)}] = \sum_z \log P(Y, Z | \theta) p(Z | Y, \theta^{(t)})$$

(3) M步: 求取 $Q(\theta, \theta^{(t)})$ 极大化的 θ , 确定第t+1次迭代的参数的估计值 $\theta^{(t+1)}$

$$\theta^{(t+1)} = \arg \max_{\theta} Q(\theta, \theta^{(t)})$$

(4) 重复第(2)步和第(3)步, 直到收敛。

3

输入：观测数据 Y, 隐变量, Z参数 θ , 联合分布 $P(Y, Z | \theta)$, 条件分布 $P(Y | Z, \theta)$

输出： 模型参数 $\hat{\theta}$

(1) 选择参数的初始值 $\theta^{(0)}$, 开始迭代:

(2) E步: 以 $\theta^{(t)}$ 为每次迭代参数的估计值, 计算第t+1次迭代的E步, 计算

$$Q(\theta, \theta^{(t)}) = E_z[\log P(Y, Z | \theta) | Y, \theta^{(t)}] = \sum_z \log P(Y, Z | \theta) p(Z | Y, \theta^{(t)})$$

(3) M步: 求取 $Q(\theta, \theta^{(t)})$ 极大化的 θ , 确定第t+1次迭代的参数的估计值 $\theta^{(t+1)}$

$$\theta^{(t+1)} = \arg \max_{\theta} Q(\theta, \theta^{(t)})$$

(4) 重复第(2)步和第(3)步, 直到收敛。

EM算法简单理解

EM算法是求解这个问题的一种迭代算法, 它有3步:

E步

E步: 计算在当前迭代的模型参数下, 观测数据Y来自混合B的概率:

$$\mu^{(t+1)} = \frac{\pi^{(t)} (\mu^{(t)})^2 (1 - p^{(t)})^{-1}}{\pi^{(t)} (\mu^{(t)})^2 (1 - p^{(t)})^{-1} + (1 - \pi^{(t)}) (q^{(t)})^2 (1 - q^{(t)})^{-1}}$$

这个式子也是一目了然的, 分子代表选B并进行一次投票试验, 分母代表选B或C并进行一次投票试验, 两个一除就得到试验结果来自B的概率。

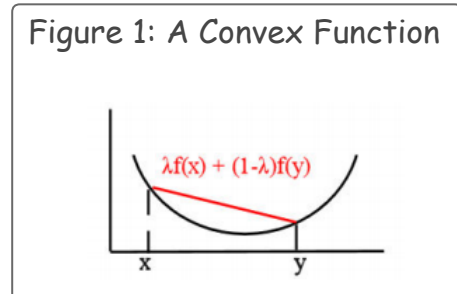
M步

M步: 估算下一个迭代的新的模型值:

$$\begin{aligned} \pi^{(t+1)} &= \frac{1}{n} \sum_{i=1}^n \mu_i^{(t+1)} \\ p^{(t+1)} &= \frac{\sum_{i=1}^n \mu_i^{(t+1)} y_i}{\sum_{i=1}^n \mu_i^{(t+1)}} \\ q^{(t+1)} &= \frac{\sum_{i=1}^n (1 - \mu_i^{(t+1)}) y_i}{\sum_{i=1}^n (1 - \mu_i^{(t+1)})} \end{aligned}$$

这个也好说, 把这n个(试验结果来自B的概率)求和得到期望, 平均后, 得到B出正面的似然估计, 同理有p和q。

重复以上步骤直到收敛



Jensen's inequality is a general result in convexity

$$\lambda f(x) + (1 - \lambda) f(y) \geq f(\lambda x + (1 - \lambda) y)$$

Jensen's Inequality

We can also generalize the result to expectation: $\mathbb{E}[f(X)] \geq f(\mathbb{E}[X])$.

Object function

$$\hat{\theta} = \arg \max_{\theta} \sum_x p(x | \theta) p(x | z, \theta)$$
$$\hat{\theta} = \arg \max_{\theta} \log p(x, z | \theta) = \arg \max_{\theta} \log p(z | \theta) + \log p(x | z, \theta).$$

Let D, (x_1, \dots, x_N) be the observed data, and let Z, hidden random variables

where $q(z)$ represents an arbitrary distribution for the random variable Z.

$$\begin{aligned} \log p(x | \theta) &= \log \sum_z p(z | \theta) p(x | z, \theta) \frac{q(z)}{q(z)} \\ &= \log \mathbb{E}_q \left[\frac{p(z | \theta) p(x | z, \theta)}{q(z)} \right] \\ &\geq \mathbb{E}_q \left[\log \left[\frac{p(z | \theta) p(x | z, \theta)}{q(z)} \right] \right] \\ \mathcal{L}(\theta; q) &= \mathbb{E}_q[\log p(z | \theta)] + \mathbb{E}_q[\log p(x | z, \theta)] - \mathbb{E}_q[\log q(z)]. \end{aligned}$$

At the M-Step, we update the model parameter θ

$$\theta^{(t+1)} = \arg \max_{\theta} \mathcal{L}(q^{(t+1)}, \theta)$$

At the E-Step, we update the posterior value q

$$q^{(t+1)} = \arg \max_q \mathcal{L}(q, \theta^{(t)}) = p(z | x, \theta^{(t)}).$$

as we maximize the objective function, we are also maximizing the log likelihood function.

$$\begin{aligned} \mathcal{L}(p(z | x, \theta), \theta) &= \sum_z p(z | x, \theta) \log \frac{p(x, z | \theta)}{p(z | x, \theta)} \\ &= \sum_z p(z | x, \theta) \log \frac{p(x, z | \theta) p(z | \theta)}{p(z | x, \theta)} \\ &= \sum_z p(z | x, \theta) \log p(z | \theta) \\ &= \log p(x | \theta) \sum_z p(z | x, \theta) \\ &= \log p(x | \theta) \end{aligned}$$

The E step has the following simple form, which is the same for any mixture model:

$$\begin{aligned} Q(\theta, \theta^{(t)}) &= E_z[\log P(Y, Z | \theta) | Y, \theta^{(t)}] \\ &= \sum_z \log P(Y, Z | \theta) P(Z | Y, \theta^{(t)}) \\ \theta^{(t+1)} &= \arg \max_{\theta} Q(\theta, \theta^{(t)}) \end{aligned}$$

repeat

The lower bound is obtained via Jensen's inequality

$$\log \sum_i p_i f_i \geq \sum_i p_i \log f_i$$

the bound of EM

This shows that $\Theta(t+1)$ is indeed a better (or no worse) parameter than $\Theta(t)$ in terms of the marginal log likelihood function. By iterating, we arrive at a local maximum of function

It is easy to see that

$$Q(\Theta^{(t)}, \Theta^{(t)}) = \sum_{i=1}^n \log p(x_i | \Theta^{(t)}) = \ell(\Theta^{(t)}). \quad (18)$$

Since $\Theta^{(t+1)}$ maximizes Q, we have

$$Q(\Theta^{(t+1)}, \Theta^{(t)}) \geq Q(\Theta^{(t)}, \Theta^{(t)}) = \ell(\Theta^{(t)}). \quad (19)$$

On the other hand, Q lower bounds ℓ . Therefore

$$\ell(\Theta^{(t+1)}) \geq Q(\Theta^{(t+1)}, \Theta^{(t)}) \geq Q(\Theta^{(t)}, \Theta^{(t)}) = \ell(\Theta^{(t)}). \quad (20)$$

