

期望最大值算法

回忆

如果 $f(x)$ 的二阶导数 > 0 ，则 f 是凸函数；
当 x 是向量时，函数对应的海森矩阵 (hessian) $H > 0$ ，是半正定矩阵。
如果 $f(x)$ 的二阶导数 > 0 ，则 f 是严格凸函数
当 x 是向量时，函数对应的海森矩阵 (hessian) $H > 0$ ，是正定矩阵。

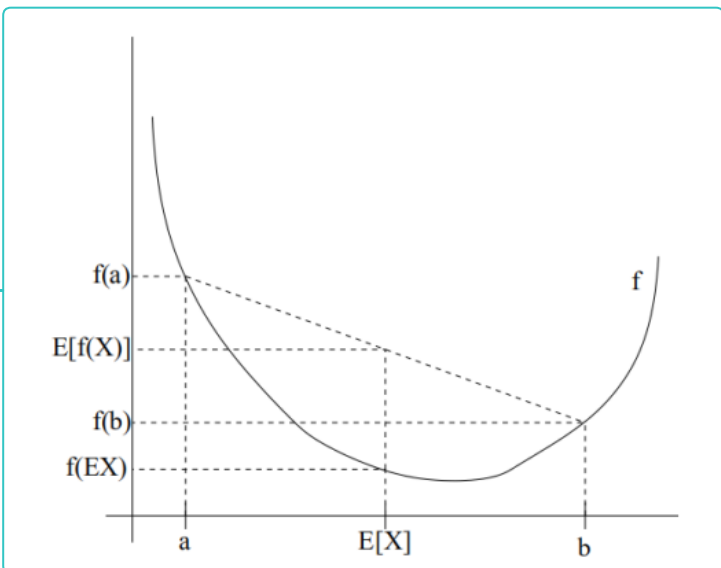
Theorem. Let f be a convex function, and let X be a random variable.
Then:

$$E[f(X)] \geq f(EX).$$

Moreover, if f is strictly convex, then $E[f(X)] = f(EX)$ holds true if and only if $X = E[X]$ with probability 1 (i.e., if X is a constant).

简而言之就是函数的期望 $>$ 期望的函数值

如果 f 是严格凸函数，则 $E[f(x)] = f(EX)$ 当且仅当 $X = E[X]$ 时成立，比如 X 取常数



更特殊的，可以画一个 $f = x^2$ 的函数，一看便知

注：当 $-f$ 是严格凸函数，也就是 f 是严格凹函数的时候，不等式变号

jensen's inequality

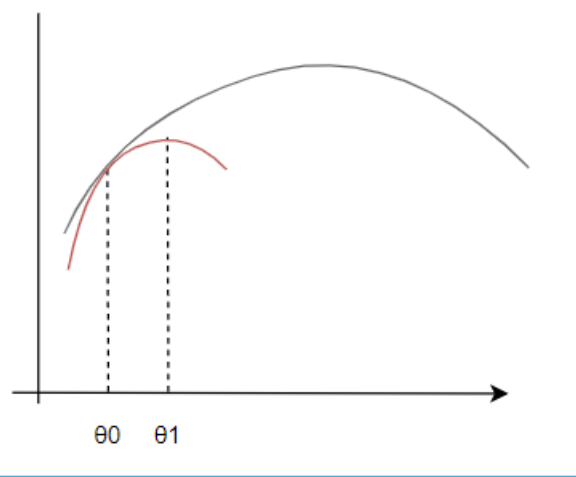
假如有一个估计问题 (estimation problem)，训练样本集 $\{x(1), \dots, x(m)\}$ 包含了 m 个独立样本。用模型 $p(x, z)$ 对数据建模，拟合参数 (parameters)，其中的似然函数

$$\begin{aligned} \ell(\theta) &= \sum_{i=1}^m \log p(x; \theta) \\ &= \sum_{i=1}^m \log \sum_z p(x, z; \theta) \end{aligned}$$

(Likelihood) 为

因为 z_i 是隐含随机变量，是未知的，所以解得 θ 很难，我们的策略是使用一种替代，找一个似然函数的下界函数。在 E 步找到一个下界函数，在 M 步对下界函数进行优化

图解



该算法的过程就是：初始化 θ_0 ，找到一个下界函数，再找到该下界函数的最大值点，设为 θ_1 ，再根据 θ_1 找到新的下界函数，再找该下界函数的最大值点，设为 θ_2 ，循环反复，最终收敛到 $l(\theta)$ 的局部最优值

所以我们设置 Q_i 为 z 的某个分布，并且满足 $Q_i(z) > 0$ ， $\sum_z Q_i(z) = 1$
(如果 z 是连续的，则 Q_i 是概率密度函数，需要将求和符号换成积分符号)

则有

$$\sum_i \log p(x^{(i)}; \theta) = \sum_i \log \sum_{z^{(i)}} p(x^{(i)}, z^{(i)}; \theta) \quad (1)$$

$$= \sum_i \log \sum_{z^{(i)}} Q_i(z^{(i)}) \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \quad (2)$$

$$\geq \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \quad (3)$$

公式第三步解释

期望定义： $E[g(z)] = \sum p(z)g(z)$

$\left[\frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right]$ 基于 z_i 的期望为

$$\sum_{z^{(i)}} Q_i(z^{(i)}) \left[\frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right]$$

$f(x) = \log x$ 是凹函数

利用 Jensen 不等式，就得到了

$$f\left(E_{z^{(i)} \sim Q_i} \left[\frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right] \right) \geq E_{z^{(i)} \sim Q_i} \left[f\left(\frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right) \right]$$

其中上面的角标 “ $z^{(i)} \sim Q_i$ ” 就表明这个期望是对于依据分布 Q_i 来确定的 $z^{(i)}$ 的。

也就得到了第三步的式子

那么 Q_i 该如何选择呢？如果我们对 θ 有某种当前估计，那么将下界函数和 θ 联系起来是很自然的事。为了让下界和 θ 的特定值联系密切，我们有必要让上面的不等式的等号成立。Jensen 不等式的时候我们说过， $x = EX$ 则等号成立，而 x 取常数可以使 $x = EX$

$$\frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} = c$$

所以我们让

$$Q_i(z^{(i)}) \propto p(x^{(i)}, z^{(i)}; \theta)$$

要实现这一条件，只需

$$\sum_z Q_i(z^{(i)}) = 1$$

由于我们已经知道 Q_i 是一个概率分布，所以

$$\begin{aligned} Q_i(z^{(i)}) &= \frac{p(x^{(i)}, z^{(i)}; \theta)}{\sum_z p(x^{(i)}, z; \theta)} \\ &= \frac{p(x^{(i)}, z^{(i)}; \theta)}{p(x^{(i)}; \theta)} \\ &= p(z^{(i)} | x^{(i)}; \theta) \end{aligned}$$

所以

因此，在给定 $x^{(i)}$ 和参数 θ 的设置下，我们可以简单地把 Q_i 设置为 $z^{(i)}$ 的后验分布 (posterior distribution)

至此， Q_i 已经选择完毕，也就是 E 步骤完毕，在接下来的 M 步骤中，需要最大化等式 (3)，然后得到参数 θ ，循环往复

Repeat until convergence {

(E-step) For each i , set

$$Q_i(z^{(i)}) := p(z^{(i)} | x^{(i)}; \theta).$$

(M-step) Set

$$\theta := \arg \max_{\theta} \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}$$

设 $\theta^{(t)}$ 和 $\theta^{(t+1)}$ 是上面 EM 迭代过程中的某两个参数

接下来我们只需证明 $\ell(\theta^{(t)}) \leq \ell(\theta^{(t+1)})$ ，这就表明 EM 迭代过程总是让似然函数对数 (log-likelihood) 单调递增

$$\text{参数起点设置为 } \theta^{(t)}, \text{ 我们选择 } Q_i^{(t)}(z^{(i)}) := p(z^{(i)} | x^{(i)}; \theta^{(t)})$$

$$\ell(\theta^{(t)}) = \sum_i \sum_{z^{(i)}} Q_i^{(t)}(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta^{(t)})}{Q_i^{(t)}(z^{(i)})}$$

之前已经分析过了，这样选择 Q_i 可以使不等式的等号成立，也就是

第一个不等式来自

$$\ell(\theta) \geq \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}$$

第二个不等式是因为 $\theta^{(t+1)}$ 来自式子

$$\arg \max_{\theta} \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}$$

所以 (4) 式一定比 (5) 式大

第三个式子是因为我们前面就保证了 Jensen 不等式的等号成立了

因此，EM 算法能够导致似然函数的单调收敛

$$J(Q, \theta) = \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}$$

备注

如果我们定义

$$\ell(\theta) \geq J(Q, \theta)$$

在前面我们已经知道

EM 算法也可以看成是对 $J(Q, \theta)$ 的坐标上升，E 步骤在 Q 上对 J 进行了最大化，M 步骤则在 θ 上对 J 进行最大化。

MoG 混合高斯模型

之前的高斯混合模型要拟合的参数有 ϕ ， μ 和 Σ

E-step

$$w_j^{(i)} = Q_i(z^{(i)} = j) = P(z^{(i)} = j | x^{(i)}; \phi, \mu, \Sigma)$$

这里的 “ $Q_i(z^{(i)} = j)$ ” 表示的是在分布 Q_i 上 $z^{(i)}$ 取值 j 的概率

最大化关于参数 ϕ ， μ ， Σ 的似然函数的下界函数

$$\begin{aligned} \sum_{i=1}^m \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \phi, \mu, \Sigma)}{Q_i(z^{(i)})} \\ &= \sum_{i=1}^m \sum_{j=1}^k Q_i(z^{(i)} = j) \log \frac{p(x^{(i)} | z^{(i)} = j; \mu, \Sigma) p(z^{(i)} = j; \phi)}{Q_i(z^{(i)} = j)} \\ &= \sum_{i=1}^m \sum_{j=1}^k w_j^{(i)} \log \frac{1}{(2\pi)^{n/2} |\Sigma_j|^{1/2}} \exp \left(-\frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j) \right) \cdot \phi_j \end{aligned}$$

解 μ

求偏导

$$\begin{aligned} \nabla_{\mu_l} \sum_{i=1}^m \sum_{j=1}^k w_j^{(i)} \log \frac{1}{(2\pi)^{n/2} |\Sigma_j|^{1/2}} \exp \left(-\frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j) \right) \cdot \phi_j \\ &= -\nabla_{\mu_l} \sum_{i=1}^m \sum_{j=1}^k w_j^{(i)} \frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j) \\ &= \frac{1}{2} \sum_{i=1}^m w_l^{(i)} \nabla_{\mu_l} 2 \mu_l^T \Sigma_l^{-1} x^{(i)} - \mu_l^T \Sigma_l^{-1} \mu_l \\ &= \sum_{i=1}^m w_l^{(i)} (\Sigma_l^{-1} x^{(i)} - \Sigma_l^{-1} \mu_l) \end{aligned}$$

令上式为 0

$$\mu_l := \frac{\sum_{i=1}^m w_l^{(i)} x^{(i)}}{\sum_{i=1}^m w_l^{(i)}}$$

M-step

$$\sum_{i=1}^m \sum_{j=1}^k w_j^{(i)} \log \phi_j$$

发现只需要最大化

其余式子都可以当成常数看待，一求导就没了

然而还有附加约束 the ϕ_j 's sum to 1，因为 ϕ_j 表示的是 $\phi_j = p(z^{(i)} = j; \phi)$

$$\mathcal{L}(\phi) = \sum_{i=1}^m \sum_{j=1}^k w_j^{(i)} \log \phi_j + \beta \left(\sum_{j=1}^k \phi_j - 1 \right)$$

所以构造拉格朗日函数

$$\frac{\partial}{\partial \phi_j} \mathcal{L}(\phi) = \sum_{i=1}^m \frac{w_j^{(i)}}{\phi_j} + 1$$

对拉格朗日函数求导

令导数为 0

$$\phi_j = \frac{\sum_{i=1}^m w_j^{(i)}}{-\beta}$$

也就是

$$\phi_j \propto \sum_{i=1}^m w_j^{(i)}$$

对式子两边求和，得到

$$-\beta = \sum_{i=1}^m \sum_{j=1}^k w_j^{(i)} = \sum_{i=1}^m 1 = m$$

这里使用了条件 $w_j^{(i)} = Q_i(z^{(i)} = j)$

$$\phi_j := \frac{1}{m} \sum_{i=1}^m w_j^{(i)}$$

所以最终：

The EM algorithm

The EM algorithm

完整 EM 算法

Q_i 的选择 (E 步骤)

备注

怎么保证该算法收敛呢

对该式最大化可得 $\ell(\theta^{(t+1)})$ ：

$$\ell(\theta^{(t+1)}) \geq \sum_i \sum_{z^{(i)}} Q_i^{(t)}(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta^{(t+1)})}{Q_i^{(t)}(z^{(i)})} \quad (4)$$

$$\geq \sum_i \sum_{z^{(i)}} Q_i^{(t)}(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta^{(t)})}{Q_i^{(t)}(z^{(i)})} \quad (5)$$

$$= \ell(\theta^{(t)}) \quad (6)$$