

# 期望最大算法

## Expectation Maximization Algorithm (EM)

① 算法导出

② 收敛性证明

③ 广义EM 以及EM的变种

# 算法导出

模型无隐变量  $\Rightarrow$  MLE . MAP

模型含隐变量  $\Rightarrow$  EM

$X$ : 观测变量,  $z$ : 隐变量  $\Rightarrow (X, z)$  完全数据

$\theta$ : 参数

同样, 我们的目标是最大似然, (这里是观测变量的似然)

$$\max \log P(X|\theta)$$

下面从两个角度出发, 导出EM的公式

ELBO + KL Divergence

$$\begin{aligned}\log P(X|\theta) &= \log P(X, z|\theta) - \log P(z|X, \theta) \\ &= \log \frac{P(X, z|\theta)}{q(z)} - \log \frac{P(z|X, \theta)}{q(z)}\end{aligned}$$

两边同时对  $z$  求  $q(z)$  的期望

$$\text{左边} = \int_z q(z) \log P(X|\theta) dz = \log P(X|\theta) \int_z q(z) dz = \log P(X|\theta)$$

$$\text{右边} = \underbrace{\int_z q(z) \log \frac{P(X, z|\theta)}{q(z)} dz}_{\text{ELBO}} - \underbrace{\int_z q(z) \log \frac{P(z|X, \theta)}{q(z)} dz}_{\text{KL Divergence}}$$

ELBO  
(Evidence Lower Bound)

KL Divergence

## KL Divergence

用于衡量两个分布相似度的指标

定义两个分布：真实分布  $p(x)$ ，预测分布  $q(x)$ 。

$$\sim | D_{KL}(P \parallel q) = \int p(x) \ln \frac{p(x)}{q(x)} dx.$$

性质：① 非负性， $D_{KL} \geq 0$

② 信息变换不变性， $y=ax+b$   $D_{KL}(p(x) \parallel q(x)) = D_{KL}(p(y) \parallel q(y))$

③ 非对称性， $D_{KL}(p \parallel q) \neq D_{KL}(q \parallel p)$

④ 值域  $D_{KL}$  在一定条件下可趋近于无穷。

$$\text{例 } \log P(X|\theta) = ELBO + D_{KL}(q(z) \parallel P(z|X, \theta))$$

$$\log P(X|\theta) \geq ELBO \quad (\text{当且仅当 } q(z) = p(z|X, \theta) \text{ 时, 等号成立})$$

对于  $\log P(X|\theta)$  来说，其内包含有未观测之变量，所以往往直接  $\max \log P(X|\theta)$  是不可行的，得不到解析解。自然地想到通过迭代的方式求解。这一部分，我们不直接对  $\max \log P(X|\theta)$  进行迭代，而是通过 ELBO 进行迭代。

此时，就有了一种想法，我们不去  $\max \log P(X|\theta)$ ，而是去  $\max ELBO$  从而使得  $\log P(X|\theta)$  取得更大的值。

当  $q(z) = p(z|X, \theta)$  时，两者相等。故取  $q(z) = p(z|X, \theta)$ 。

在上述推导中， $q(z)$  是我们添加进来的分布，其无论为什么都不影响推导过程，故我们将  $q(z)$  作为迭代的媒介。令  $q(z) = p(z|X, \theta^{(i)})$ ，

$\theta^{(i)}$  表示第  $i$  步迭代得到的  $\theta$ 。（这里为什么选用  $q(z)$  作为媒介是我自己分析的）

$$k: \theta^{(i+1)} = \arg \max_{\theta} \text{ELBO}$$

$$= \arg \max_{\theta} \int_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}, \theta^{(i)}) \log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{p(\mathbf{z}|\mathbf{x}, \theta^{(i)})} d\mathbf{z}.$$

$$= \arg \max_{\theta} \int_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}, \theta^{(i)}) (\log p(\mathbf{x}, \mathbf{z}|\theta) - \log p(\mathbf{z}|\mathbf{x}, \theta^{(i)})) d\mathbf{z}.$$

与  $\theta$  无关.

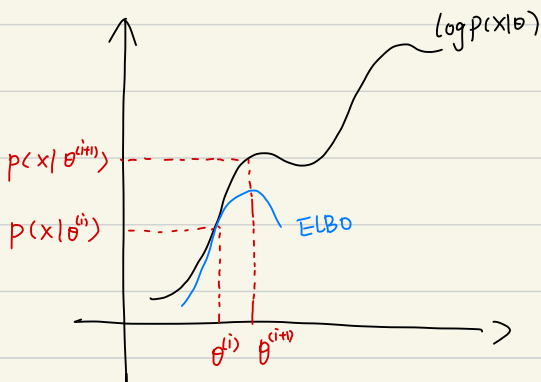
$$= \arg \max_{\theta} \int_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}, \theta^{(i)}) \log p(\mathbf{x}, \mathbf{z}|\theta) d\mathbf{z}$$

$$= \arg \max_{\theta} \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z}|\mathbf{x}, \theta^{(i)})} [\log p(\mathbf{x}, \mathbf{z}|\theta)]$$

这就是EM算法.

E步: 求  $\mathbb{E}_{\mathbf{z} \sim p(\mathbf{z}|\mathbf{x}, \theta^{(i)})} [\log p(\mathbf{x}, \mathbf{z}|\theta)]$

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



左图形象地表示了如何通过  $\max \text{ELBO}$  去  $\max \log p(\mathbf{x}|\theta)$ .

# Jensen Inequality

这一部分用 Jensen 不等式导出 EM 算法。

和之前 ELBO 部分一样。我们需要  $\max \log P(X|\theta)$ ，这里我们直接  $\max \log P(X|\theta)$ ，而不通过 ELBO。

## Jensen 不等式.

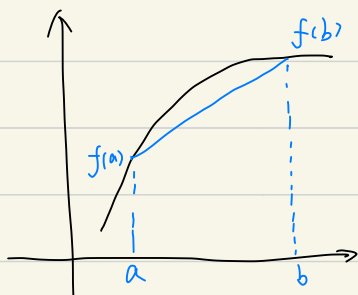
函数  $f(x)$  是 concave function.

则  $\forall t \in [0, 1]$ ,

$$f(ta + (1-t)b) \geq tf(a) + (1-t)f(b)$$

当  $t = \frac{1}{2}$  时.

$$f\left(\frac{a+b}{2}\right) \geq \frac{f(a)+f(b)}{2} \Rightarrow f(E) \geq E[f]$$



$$2|. \log P(X|\theta) = \log \int_z P(X, z|\theta) dz$$

$$= \log \int_z q(z) \frac{P(X, z|\theta)}{q(z)} dz$$

$$= \log E_{z \sim q(z)} \left[ \frac{P(X, z|\theta)}{q(z)} \right]$$

$$\geq E_{z \sim q(z)} \left[ \log \frac{P(X, z|\theta)}{q(z)} \right] \quad (\text{当且仅当 } \frac{P(X, z|\theta)}{q(z)} = c \text{ 时, 等式成立})$$

求  $q(z)$  的表达式:

$$\frac{p(x, z | \theta)}{q(z)} = c$$

$$q(z) = \frac{1}{c} p(x, z | \theta)$$

两边对  $z$  积分:

$$\int_z q(z) dz = \frac{1}{c} \int_z p(x, z | \theta) dz$$

$$1 = \frac{1}{c} p(x | \theta)$$

$$c = p(x | \theta)$$

$$\text{故} \quad q(z) = \frac{p(x, z | \theta)}{p(x | \theta)} = p(z | x, \theta).$$

$$\text{故} \quad \theta^{(i+1)} = \arg \max_{\theta} E_{z \sim p(z | x, \theta^{(i)})} \left[ \log \frac{p(x, z | \theta)}{p(z | x, \theta^{(i)})} \right]$$

与  $\theta$  无关.

$$= \arg \max_{\theta} E_{z \sim p(z | x, \theta^{(i)})} [\log p(x, z | \theta)]$$

李航的《统计学习方法》第九章对EM的导出也是用 Jensen 不等式, 但切入点在这里有点不同, 可以详细地读第九章全部一遍.

还有PRML中, EM在第九章混合模型中有介绍, 是在第4节. 一般形式的EM算法, 也可以看一下.

# 收敛性证明

$$\text{EM公式: } \theta^{(i+1)} = \arg \max_{\theta} \int_{\mathbf{z}} p(\mathbf{z} | \mathbf{x}, \theta^{(i)}) \log p(\mathbf{x}, \mathbf{z}) d\mathbf{z}.$$

该算法是迭代算法  $\theta^{(i)} \rightarrow \theta^{(i+1)}$

目标是  $\max \log P(\mathbf{x} | \theta)$ .

如果算法可行并收敛, 那么就得  $\log P(\mathbf{x} | \theta^{(i+1)}) \geq \log P(\mathbf{x} | \theta^{(i)})$   
且有上界.

显然,  $p(\mathbf{x} | \theta)$  是一个概率, 最大值为 1. 故  $\log p(\mathbf{x} | \theta)$  一定有上界.

所以只需证明  $\log P(\mathbf{x} | \theta^{(i+1)}) \geq \log P(\mathbf{x} | \theta^{(i)})$

$$\log P(\mathbf{x} | \theta) = \log P(\mathbf{x}, \mathbf{z} | \theta) - \log p(\mathbf{z} | \mathbf{x}, \theta)$$

两边同时对  $p(\mathbf{z} | \mathbf{x}, \theta^{(i)})$  求期望.

$$\text{左边} = \int_{\mathbf{z}} p(\mathbf{z} | \mathbf{x}, \theta^{(i)}) \log P(\mathbf{x} | \theta) d\mathbf{z} = \log P(\mathbf{x} | \theta)$$

$$\text{右边} = \underbrace{\int_{\mathbf{z}} p(\mathbf{z} | \mathbf{x}, \theta^{(i)}) \log P(\mathbf{x}, \mathbf{z} | \theta) d\mathbf{z}}_{Q(\theta, \theta^{(i)})} - \underbrace{\int_{\mathbf{z}} p(\mathbf{z} | \mathbf{x}, \theta^{(i)}) \log p(\mathbf{z} | \mathbf{x}, \theta) d\mathbf{z}}_{H(\theta, \theta^{(i)})}.$$

$$\text{由于 } \theta^{(i+1)} = \arg \max_{\theta} \int_{\mathbf{z}} p(\mathbf{z} | \mathbf{x}, \theta^{(i)}) \log p(\mathbf{x}, \mathbf{z} | \theta) d\mathbf{z}.$$

所以, 显然,  $Q(\theta^{(i+1)}, \theta^{(i)}) \geq Q(\theta^{(i)}, \theta^{(i)})$

所以, 只需要证  $H(\theta^{(i+1)}, \theta^{(i)}) \leq H(\theta^{(i)}, \theta^{(i)})$  就行了.

$$H(\theta^{(i+1)}, \theta^{(i)}) - H(\theta^{(i)}, \theta^{(i)}) \quad \text{这里也可以用 Jensen 不等式 (见《统计学习法》)} \\ = \int_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}, \theta^{(i)}) \log \frac{p(\mathbf{z}|\mathbf{x}, \theta^{(i+1)})}{p(\mathbf{z}|\mathbf{x}, \theta^{(i)})} d\mathbf{z}$$

$$= -D_{KL}(p(\mathbf{z}|\mathbf{x}, \theta^{(i)}) \parallel p(\mathbf{z}|\mathbf{x}, \theta^{(i+1)}))$$

$$\leq 0$$

$$\text{所以 } H(\theta^{(i+1)}, \theta^{(i)}) \leq H(\theta^{(i)}, \theta^{(i)})$$

$$\text{从而 } \log p(\mathbf{x}|\theta^{(i+1)}) \geq \log p(\mathbf{x}|\theta^{(i)})$$

算法收敛性得到证明。

其实还需要证明最终得到的  $\theta$  是  $L(\theta) = \log p(\mathbf{x}|\theta)$  的稳定点，但是这里李航的《统计学习方法》中只给出了文献，并没有证明。



# 广义EM以及EM的变种

EM主要是用来解决 概率生成模型 的参数估计问题



$x$ : observed variable

$z$ : latent variable

$(x, z)$ : complete data

$\theta$ : parameter.

生成模型就是假设  $(z) \rightarrow (x)$ . 给定  $z$  一个隐含变量  $p(z|\theta)$   
数据是由隐变量生成的.  $p(x, z|\theta) = p(z|\theta)p(x|z, \theta) \Rightarrow p(x|\theta) = \int_z p(x, z|\theta) dz$ .  
或  $p(x|\theta) = \frac{p(x, z|\theta)}{p(z|x, \theta)}$

从而  $\hat{\theta} = \arg\max p(x|\theta)$ .

## 广义EM

狭义EM的局限性:

由前面推导可知:  $\log P(x|\theta) = ELBO + D_{KL}(q||p) \xRightarrow{\Delta \text{后验概率}} \log P(x|\theta) \geq ELBO$

$$ELBO = \int_z q(z) \log \frac{p(x, z|\theta)}{q(z)} dz = \underbrace{\int_z q(z) \log p(x, z|\theta) dz}_{E_{q(z)}[\log p(x, z|\theta)]} - \underbrace{\int_z q(z) \log q(z) dz}_{\text{熵: } H[q(z)]}$$

$$D_{KL}(q||p) = - \int_z q(z) \log \frac{p(z|x, \theta)}{q(z)} dz$$

对于狭义的EM来说, 我们是令  $q(z) = p(z|x, \theta)$ .

但是大多数时候, 后验概率  $p(z|x, \theta)$  是 intractable 的, 求不出来. 显然我们不能简单地令  $q = p$ . 对  $q$ , 我们采用局部最优, 而不是全局最优

广义EM的导出.

$$\log p(x|\theta) = \text{ELBO} + D_{KL}(q||p).$$

我们通过上式可知: 当  $D_{KL}(q||p)$  越小时, ELBO 越大.

所以我们可以将  $q$  的选择也作为迭代的一部分. 令  $\text{ELBO} = L(q, \theta)$

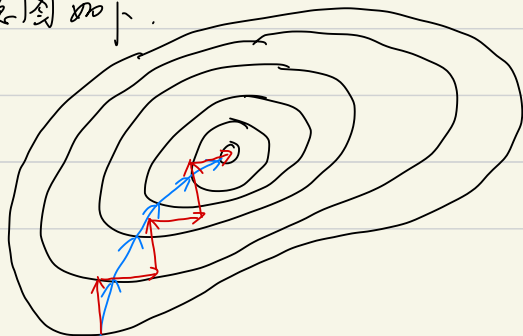
$$\begin{cases} \text{固定 } \theta, \hat{q} = \arg\min_q D_{KL}(q||p) = \arg\max_q L(q, \theta) \\ \text{固定 } \hat{q}, \theta = \arg\max_{\theta} L(\hat{q}, \theta) \end{cases}$$

则广义的EM为:

$$\begin{cases} \text{E-step: } q^{(i+1)} = \arg\max_q L(q, \theta^{(i)}) \\ \text{M-step: } \theta^{(i+1)} = \arg\max_{\theta} L(q^{(i+1)}, \theta) \end{cases}$$
 广义的EM有时也被称为F-MM

所以狭义的EM就是直接将E-step改为  $q^{(i+1)} = p(z|x, \theta^{(i)})$ . 是广义EM的一种特例.

由此可见, 广义的EM是对于两个目标  $q, \theta$ , 先固定  $\theta$  求  $q$ , 再固定  $q$  求  $\theta$ . 如此不断迭代. 类似于坐标上升法 (SNO算法 → 在svm中提了一嘴) 算法简单的示意图如下.



→: 梯度上升法

→: 坐标上升法.

EM 的变种

variational Inference / variational Bayes

VBEM / VEM : 用 VI / VB (变分推断) 来计算 E-step.

MCEM : 用 MC 采样方法计算 E-step.