# EM算法(Expectation-Maximum)

#### 原理篇

- 本质:概率模型,去估计一个密度函数,最大化对数似然函数去估计参数。无法直接用极大化对数似然函数得到模型分布的参数,用启发式跌代法,用于求解含有隐变量的最大似然估计、最大后延概率估计问题。
  - 1. 隐变量:对概率模型有一定的影响,但无法观测;

观测数据X+隐变量 = 完全数据

- 2. 无法用极大似然的两种情况:数据缺失、含有未知隐变量
- 作用:经常存在数据缺失或者不可用的的问题,这时候直接处理数据比较困难,而数据添加办法有很多种,常用的有神经网络拟合、添补法、卡尔曼滤波法等等,但是EM算法之所以能迅速普及主要源于它算法简单,稳定上升的步骤能非常可靠地找到"最优的收敛值"。
- 目标:是使包含隐变量的数据集的后验概率或似然函数最大化,进而得到最优的参数估计
- 思想:我们可以发现我们的算法里已知的是观察数据,未知的是隐含数据和模型参数,在E步,我们所做的事情是固定模型参数的值,优化隐含数据的分布,而在M步,我们所做的事情是固定隐含数据分布,优化模型参数的值。

#### 主要步骤:

已知:概率分布、随机抽取的样本;

未知:分类、模型参数

• E-step: 猜想隐含数据,更新隐含数据和模型参数

• M-step: 基于观察数据和猜测的隐含数据求极大对数似然, 求解模型K

• E-step:基于前面得到的模型K,继续猜测隐含数据,继续极大化对数似然

• M-step:求模型参数

• 直至模型分布无明显变化, 算法收敛

EM算法步骤如下:

输入:观测变量数据X,隐变量Z,联合分布 $p(X,Z|\Theta)$ 

输出:模型参数Θ

(1) 选择初始参数Θ\_0

(2)E步:记 $\Theta$ i为第i次迭代参数 $\Theta$ 的估计值,在第i+1次迭代的E步,计算 $Q(\Theta,\Theta g)$ ;

(3)M步:确定第i+1次迭代的参数的估计值Θ(i+1),即为:

 $\Theta^{(i+1)} = argmax_{\Theta}Q(\Theta, \Theta^{(i)})$ 

(4)重复(2)和(3)步,直至收敛

## EM算法的引入篇

三枚硬币的模型:

$$egin{split} p(y| heta) &= \sum_{z} P(y,z| heta) = \sum_{z} p(z| heta) p(y|z, heta) \ &= \pi p^y (1-p)^{1-y} + (1-\pi) q^y (1-q)^{1-y} \end{split}$$

N枚硬币的模型:

$$egin{aligned} P(Y|\Theta) &= \sum_{z} P(Z|\Theta) P(Y|Z,\Theta) \ &= \prod_{j=1}^{n} [\pi p^{y_j}] (1-p)^{1-y_j} + (1-\pi) q^{y_j} (1-q)^{q-y_j} \end{aligned}$$

求模型的参数 Θ=(π,p,q)的极大似然估计, **即目标函数**为:

$$\Theta = argmaxlogp(Y|\Theta) = argmaxlog\sum_{\tilde{z}} P(Y|Z,\Theta)p(Z,\Theta)$$

面对含隐变量的概率模型,目标是**极大化观测数据(不完全数据)Y关于参数Θ的对数似然函数**,即极大化(取对数方便计算):

$$egin{aligned} L(\Theta) &= log P(Y|\Theta) = log \sum_{z} P(Y,Z|\Theta) \\ &= log (\sum_{z} P(Y|Z,\Theta) P(Z|\Theta)) \end{aligned}$$

EM算法通过迭代逐步近似极大化 $L(\theta)$ 的 + Jensen不等式

$$\begin{split} L(\Theta) - L(\Theta^i) &= log(\sum_Z P(Y|Z,\Theta)P(Z|\Theta)) - logP(Y|\Theta^{(i)}) \\ &= log[\sum_Z P(Z|Y,\Theta^{(i)}) \frac{P(Y|Z)P(Z|\Theta)}{P(Z|Y,\Theta^i)}] - logP(Y|\Theta^{(i)}) \\ & \quad \text{$\pm Jensen$} \\ & \quad \pm \sum_Z P(Z|Y,\Theta^{(i)}) = 1 \\ &\geq \sum_Z P(Z|Y,\Theta^{(i)}) log \frac{P(Y|Z,\Theta)P(Z|\Theta)}{P(Z|Y,\Theta^{(i)})} - \frac{logP(Y|\Theta^{(i)}) * \sum_Z P(Z|Y,\Theta^{(i)})}{\sum_Z P(Z|Y,\Theta^{(i)})} \\ &= \sum_Z P(Z|Y,\Theta^{(i)}) * log \frac{P(Y|Z,\Theta)P(Z|\Theta)}{P(Z|Y,\Theta^{(i)})P(Y|\Theta^{(i)})} \end{split}$$

即:

$$\Theta^{(i+1)} = argmax[L(\Theta^{(i)}) + \sum_{Z} P(Z|Y,\Theta^{(i)}) * log rac{P(Y|Z,\Theta)P(Z|\Theta)}{P(Z|Y,\Theta^{(i)})P(Y|\Theta^{(i)})}]$$

去掉与Θ无关的变量的式子:

$$egin{aligned} \Theta^{(i+1)} &= argmax[\sum_{Z} P(Z|Y,\Theta^{(i)}) * logP(Y|Z,\Theta)P(Z|\Theta)] \ &= argmaxQ(\Theta,\Theta^{(i)}) = E[logP(Z|Y,\Theta^{(i)})].\dots (1) \end{aligned}$$

观测数据 + 隐变量 = 完全数据

Jensen不等式:

$$\log(\sum lpha_i arphi(x_i)) >= \sum lpha_i log arphi(x_i) \ \sum lpha_i = 1$$
 if  $lpha_i >= 0$ 

### EM算法的收敛性

证明: $P(Y|\theta)$ 为观测数据的似然函数,且是递增的,即:

$$P(Y|\theta^{(i+1)}) \ge P(Y|\theta^{(i)})$$

证明如下:

$$P(Y| heta)=rac{P(Y,Z| heta)}{P(Z|Y, heta)}.\dots$$
. 有贝叶斯公式展开得 $log P(Y| heta)=log P(Y,Z| heta)-log P(Z|Y, heta).\dots$ 取对数

$$Q( heta, heta^{(i)}) = \sum_z log P(Y, Z| heta) P(Z|Y, heta^{(i)})$$
...由( $1$ )得

构造下式(因为取对数方便相减和相除,同时构造了贝叶斯公式):

在式(2)中分别取 $\theta$ 为 $\theta$ i 和 $\theta$ i+1, **并相减**:

$$log P(Y| heta^{(i+1)}) - log P(Y| heta^{(i)}) = \ [Q( heta^{(i+1)}, heta^{(i)}) - Q(( heta^{(i)}, heta^{(i)}))] - [H( heta^{(i+1)}, heta^{(i)}) - H(( heta^{(i)}, heta^{(i)}))]$$

其中对H:

$$\begin{split} [H(\theta^{(i+1)},\theta^{(i)})-H((\theta^{(i)},\theta^{(i)}))] &= \sum_z (logP(Z|Y,\theta^{(i+1})) = \\ &\sum_z (log\frac{P(Z|Y,\theta^{(i+1})}{P(Z|Y,\theta^{(i)})}P(Z|Y,\theta^{(i)}) \leq \\ &log(\sum_Z \frac{P(Z|Y,\theta^{(i+1})}{P(Z|Y,\theta^{(i)})}P(Z|Y,\theta^{(i)}))...Jensen$$
不等式得 
$$&log(\sum_Z P(Z|Y,\theta^{(i+1)})) = log1 = 0 \end{split}$$

对Q,由于Q的i+1项已经达到极大,所以有:

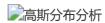
$$[Q(\theta^{(i+1)},\theta^{(i)})-Q((\theta^{(i)},\theta^{(i)}))]\geq 0$$

最后证明得到:

$$log P(Y| heta^{(i+1)}) \geq log P(Y| heta^{(i)})$$

# 应用篇--GMM

EM在GMM中的应用



图中可知:

- 1. 单个高斯拟合效果差,均值应该分布在数据密集处
  - 2. 混合高斯模型中的**隐变量**,同时,隐变量在概率模型中不能改变边缘分布,即:

$$p(x_i) = \int_{z_i} p(x_i|z_i) * p(z_i) dz_i = lpha_z = \sum_{z_i=1}^k lpha_{z_i} N(x_i|\mu_z, \Sigma_z)$$

每个数据都有一个隐变量,告诉你在哪个高斯模型中(由两个高斯扩展到n个高斯)

$$P(z_i=z_1|x_i,\Theta^g)=rac{a}{a+b}$$

3.

$$P(x) = \sum_{l=1}^k lpha_l * N(X|\mu_l, \Sigma_l) \qquad \quad \sum_{l=1}^k lpha_l = 1$$

$$\Theta = \{\alpha_1, \dots, \alpha_k, \mu_1, \dots, \mu_k, \Sigma_1, \dots, \Sigma_{k-1}\}$$

4. 目标函数为:

$$\Theta_{MLE} = argmax_{\Theta} * L(\Theta \mid X) = argmax_{\Theta}(\sum_{l=1}^{n} log * \sum_{l=1}^{k} lpha_{l} N(X \mid \mu_{l}, \Sigma_{l}))$$

该式子包含和(或积分)的对数,不能像单个高斯模型那样直接求导,再令导数为0来求解。这时我们需要利用 EM 算法通过迭代逐步近似极大化L(O)X)来求解。

#### EM算法在GMM中的应用:

高斯混合模型的概率分布模型如下:

$$\begin{split} P(Y|\theta) &= \sum_{K=1}^K \alpha_k \phi(Y|\theta) \\ \sharp \, & + : \sum_{K=1}^K \alpha_k = 1, \theta_k = (\mu_k, \sigma_k^2) \\ \theta &= (\alpha_1, \alpha_2, \dots, \alpha_k, \theta_1, \theta_2, \dots, \theta_k) \\ \phi(Y|\theta) &= \frac{1}{\sqrt{2\pi}\sigma_k} exp(-\frac{(y-\mu_k)^2}{2\sigma_k^2}) \end{split}$$

用 y 作为隐变量,去确定是哪个模型, y =1/0.

完整数据的似然函数为:

$$egin{aligned} P(y,\gamma| heta) &= \prod_{j=1}^N P(y_j,\gamma_{j1},\gamma_{j2},\ldots,\gamma_{jK}| heta) \ &= \prod_{K=1}^K \prod_{j=1}^N [lpha_k\phi(y_j| heta_k)]^{\gamma^{jK}} \end{aligned}$$

完全数据的对数似然为:

$$logP(y,\gamma| heta) = \sum_{K=1}^K [\sum_{j=1}^N loglpha_k + \sum_{j=1}^N \gamma_{jk} [log(rac{1}{\sqrt{2\pi}}) - log\sigma_k - rac{1}{2\sigma_K^2} (y_j - \mu_k)^2]]$$

EM算法中的E步:确定Q函数,对每一个隐变量求期望

根据当前模型参数,计算分模型k对观测数据y\_i的响应度

$$egin{aligned} Q( heta, heta^{(i)}) &= E[log P(y,\gamma| heta)|y, heta^{(i)}] \ &= E[\sum_{K=1}^K [\sum_{j=1}^N log lpha_k + \sum_{j=1}^N \gamma_{jk} [log (rac{1}{\sqrt{2\pi}}) - log \sigma_k - rac{1}{2\sigma_K^2} (y_j - \mu_k)^2]]] \ &= \sum_{K=1}^K [\sum_{j=1}^N (E_{\gamma_{jk}}) log lpha_K + \sum_{j=1}^N (E\gamma_{jk}) [log (rac{1}{\sqrt{2\pi}}) - log \sigma_K - ]rac{1}{2\sigma_k^2} (y_j - \mu_k)2] \end{aligned}$$

计算E

$$egin{aligned} \gamma_{jk} &= E(\gamma_{jk}|y, heta) \ &= P(\gamma_{jk} = 1|y, heta) * 1 + P(\gamma_{jk} = 0, heta) * 0 \ &= rac{P(\gamma_{jk} = 1,y_{j}| heta)}{P(y_{j}| heta)} = rac{P(\gamma_{jk} = 1,y_{j}| heta)}{\sum_{K=1}^{K} P(\gamma = 1,y_{j}| heta)} \ &= \ldots = rac{lpha_{k}\phi(y_{j}| heta_{k})}{\sum_{K=1}^{K} lpha_{k}\phi(y_{j}| heta_{k})} \end{aligned}$$

最终求得Q

$$Q(\theta, \theta^{(i)}) = \sum_{K=1}^{K} [\sum_{j=1}^{N} log \alpha_k + \sum_{j=1}^{N} \gamma_{jk} [log (\frac{1}{\sqrt{2\pi}}) - log \sigma_k - \frac{1}{2\sigma_k^2} (y_j - \mu_k)^2]]$$

#### M步,进行迭代模型的参数

$$egin{aligned} heta^{(i+1)} &= argmax_{ heta}Q( heta, heta^{(i)}) \ \mu_k &= rac{\sum_{j=1}^N \gamma_{jk}y_j}{\sum_{j=1}^N \gamma_{jk}} \ \sigma_k^2 &= rac{\sum_{j=1}^N \gamma_{jk}(y_j - \mu_k)^2}{\sum_{j=1}^N \gamma_{jk}} \ lpha_k &= rac{\sum_{j=1}^N \gamma_{jk}}{N} \end{aligned}$$

参考1: 知乎关于EM

参考2:GMM应用

参考3:EM九个境界

参考4:统计学习方法