

EM Algorithm

EM(Expectation maximization)算法，也即期望最大化算法，作为“隐变量”（属性变量不可知）估计的利器在自然语言处理（如HMM中的Baum-Welch算法）、高斯混合聚类、心理学、定量遗传学等含有隐变量的概率模型参数极大似然估计中有着十分广泛的应用。EM算法于1977年由Arthur Dempster, Nan Laird和Donald Rubin总结提出，其主要通过E步（expectation），M步（maximization）反复迭代直至似然函数收敛至局部最优解。由于其方法简洁、操作有效，EM算法曾入选“数据挖掘十大算法”，可谓是机器学习经典算法之一。

Introduction

EM算法推导一

对于概率模型，当模型中的变量均为观测变量时，我们可以直接使用给定数据通过最大似然估计（频率学派）或贝叶斯估计（贝叶斯学派）这两种方法求解。然而当我们的模型中存在隐变量时，我们将无法使用最大似然估计直接求解，这时即导出EM算法。

假设一个概率模型中同时存在隐变量 Z 和可观测变量 Y ，我们学习的目标是极大化观测变量 Y 关于模型参数 θ 的对数似然，即：

$$L(\theta) = \log P(Y|\theta) = \log \sum_Z P(Y, Z|\theta) = \log \left(\sum_Z P(Y|Z, \theta) P(Z|\theta) \right) \quad (1)$$

式（1）中我们假设直接优化 $P(Y|\theta)$ 是很困难的，但是优化完整数据的似然函数 $P(Y, Z|\theta)$ 相对容易，同时利用概率乘法公式将 $P(Y, Z|\theta)$ 展开。然而由于未观测变量 Z 的存在，上式仍求解困难，因此我们通过迭代逐步最大对数似然 $L(\theta)$ ，这里假设第 i 次迭代后 θ 的估计值为 θ^i 。根据要求，我们希望新估计的参数 θ 使 $L(\theta)$ 增加，即 $L(\theta) > L(\theta^i)$ ，且逐步使 $L(\theta)$ 达到最大，因此考虑两者之差：

$$L(\theta) - L(\theta^i) = \log \left(\sum_Z P(Y|Z, \theta) P(Z|\theta) \right) - \log(P(Y|\theta^i)) = \log \left(\sum_Z P(Z|Y, \theta^i) \frac{P(Y|Z, \theta) P(Z|\theta)}{P(Z|Y, \theta^i)} P(Z|Y, \theta^i) \right) - \log(P(Y|\theta^i)) \quad (2)$$

这里我们根据Jensen（琴生）不等式： $\log \sum_j \lambda_j y_j \geq \sum_j \lambda_j \log y_j$ ，其中 $\lambda_j \geq 0, \sum_j \lambda_j = 1$ ，有：

$$\log \left(\sum_Z P(Z|Y, \theta^i) \frac{P(Y|Z, \theta) P(Z|\theta)}{P(Z|Y, \theta^i)} P(Z|Y, \theta^i) \right) - \log(P(Y|\theta^i)) \geq \sum_Z P(Z|Y, \theta^i) \log \left(\frac{P(Y|Z, \theta) P(Z|\theta)}{P(Z|Y, \theta^i)} P(Z|Y, \theta^i) \right) - \log(P(Y|\theta^i)) \quad (3)$$

同时由于 $\sum_Z P(Z|Y, \theta^i) = 1$ ，式（3）可进一步写为：

$$\sum_Z P(Z|Y, \theta^i) \log \left(\frac{P(Y|Z, \theta) P(Z|\theta)}{P(Z|Y, \theta^i)} P(Z|Y, \theta^i) \right) - \log(P(Y|\theta^i)) = \sum_Z P(Z|Y, \theta^i) \log \left(\frac{P(Y|Z, \theta) P(Z|\theta)}{P(Z|Y, \theta^i) P(Y|\theta^i)} \right) \quad (4)$$

因此有：

$$L(\theta) \geq L(\theta^i) + \sum_Z P(Z|Y, \theta^i) \log \left(\frac{P(Y|Z, \theta) P(Z|\theta)}{P(Z|Y, \theta^i) P(Y|\theta^i)} \right) = B(\theta, \theta^i) \quad (5)$$

因此, $B(\theta, \theta^i)$ 即为 $L(\theta)$ 的下界。故当 $B(\theta, \theta^i)$ 增大时 $L(\theta)$ 也将同时增加, 为使 $L(\theta)$ 取得最大, 则我们必须在 $i+1$ 次迭代时选择的 θ^{i+1} 为使第 i 次迭代 $B(\theta, \theta^i)$ 取得最大的 θ^i , 即:

$$\begin{aligned} \theta^{i+1} &= \arg\max_{\theta} B(\theta, \theta^i) = \arg\max_{\theta} \{L(\theta^i) + \sum_Z P(Z|Y, \theta^i) \log \left(\frac{P(Y|Z, \theta) P(Z|\theta)}{P(Z|Y, \theta^i) P(Y|\theta^i)} \right)\} \\ &= \arg\max_{\theta} \{ \sum_Z P(Z|Y, \theta^i) \log(P(Y|Z, \theta) P(Z|\theta)) \} = \arg\max_{\theta} \{ \sum_Z P(Z|Y, \theta^i) \log P(Y, Z|\theta) \} \\ &= \arg\max_{\theta} Q(\theta, \theta^i) \end{aligned}$$

在上式的求解中我们略去了对求解 θ 极大化而言的常数项 $L(\theta^i)$ 和 $P(Z|Y, \theta^i) P(Y|\theta^i)$ 。

因此在EM算法的每一迭代中, 我们均需求解使得 $Q(\theta, \theta^i)$ 取得最大值的 θ , 使得下一迭代的 $\theta^{i+1} = \theta$, 这样如此反复提高最大似然 $L(\theta)$ 的下界, 直至逼近 $L(\theta)$ 的最优解 (最大值)。

EM算法推导二

这里我们采用变分的方法, 假设隐变量服从任一分布为 $q(Z)$, 则 $\sum_Z q(Z) = 1$ 。故对于 $L(\theta)$ 同样有:

$$\begin{aligned} L(\theta) &= \log P(Y|\theta) = \sum_Z q(Z) \log P(Y|\theta) = \sum_Z q(Z) \log \frac{P(Z|Y, \theta) P(Y|\theta)}{P(Z|Y, \theta)} \\ &= \sum_Z q(Z) \log \frac{P(Y, Z|\theta)}{P(Z|Y, \theta)} = \sum_Z q(Z) \log \frac{P(Y, Z|\theta)}{q(Z)} - \sum_Z q(Z) \log \frac{1}{q(Z)} \\ &= \underbrace{\sum_Z q(Z) \log \frac{P(Y, Z|\theta)}{q(Z)}}_{(1)} - \underbrace{\sum_Z q(Z) \log \frac{1}{q(Z)}}_{(2)} \end{aligned}$$

记 (1) 为 $\frac{1}{n} \sum_Z q(Z) \log \frac{P(Y, Z|\theta)}{q(Z)}$, (2) 为 $KL(q||p) = -\sum_Z q(Z) \log \frac{1}{q(Z)}$ 。其中 $KL(q||p)$ 即为KL散度 (相对熵), 主要反映变量 q 、 p 分布的相似性, 可以看出KL散度=交叉熵-信息熵, 故交叉熵在某种意义上与KL散度等价。有:

$$L(\theta) = \frac{1}{n} \sum_Z q(Z) \log \frac{P(Y, Z|\theta)}{q(Z)} + KL(q||p)$$

由于 $KL(q||p) \geq 0$, 因此 $\frac{1}{n} \sum_Z q(Z) \log \frac{P(Y, Z|\theta)}{q(Z)}$ 即为对数似然函数 $L(\theta)$ 的下界。同理在每一次迭代中我们均需要最大化下界 $\frac{1}{n} \sum_Z q(Z) \log \frac{P(Y, Z|\theta)}{q(Z)}$, 则在第 i 次迭代中即有:

$$\begin{aligned} q(Z) &= P(Z|Y, \theta^i) \quad \theta^{i+1} = \arg\max_{\theta} \frac{1}{n} \sum_Z q(Z) \log \frac{P(Y, Z|\theta)}{q(Z)} \\ &= \arg\max_{\theta} \{ \sum_Z P(Z|Y, \theta^i) \log \frac{P(Y, Z|\theta)}{P(Z|Y, \theta^i)} \} = \arg\max_{\theta} \{ \sum_Z P(Z|Y, \theta^i) (\log P(Y, Z|\theta) - \log P(Z|Y, \theta^i)) \} \\ &= \arg\max_{\theta} Q(\theta, \theta^i) + \text{const} \end{aligned}$$

式 (9) 中 $-\sum_Z P(Z|Y, \theta^i) \log P(Z|Y, \theta^i)$ 为一常数 const , 故式 (8) 与式 (6) 等价。因此, 综上所述, EM算法可描述为:

对于观测变量数据 Y 和隐变量数据 Z , 其联合分布为 $P(Y, Z|\theta)$, 条件分布为 $P(Z|Y, \theta)$:

- Step1. 参数初始化 θ^0 , 开始迭代。
- Step2. E步: 记 θ^i 为第 i 次迭代的参数 θ 的估计值, 则在第 $i+1$ 次迭代的E步中, 有:

$$Q(\theta, \theta^i) = E_Z [\log P(Y, Z|\theta) | Y, \theta^i] = \sum_Z P(Z|Y, \theta^i) \log P(Y, Z|\theta)$$

上式中, $P(Z|Y, \theta^i)$ 即为给定观测数据 Y 和当前估计参数 θ^i 下隐变量数据 Z 的条件概率分布。 Q 函数为对数似然函数 $\log P(Y, Z|\theta)$ 关于在给定观测数据 Y 和当前模型参数 θ^i 下对未观测

数据 Z 的条件概率分布 $P(Z|Y, \theta^i)$ 的期望。

- Step3. M步：计算使 $Q(\theta, \theta^i)$ 取得极大值的 θ ，确定第 $i+1$ 次迭代的参数估计值 θ^{i+1} ，有：

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

- Step4. 迭代Step2，Step3直至收敛。其收敛条件一般为给定较小的正数 ϵ_1, ϵ_2 ，若满足：

$$\|\theta^{i+1} - \theta^i\| < \epsilon_1 \text{ 或 } \|Q(\theta^{i+1}, \theta^i) - Q(\theta^i, \theta^i)\| < \epsilon_2$$

由于目标函数为非凸函数，因此EM算法并不能保证收敛至全局最小值，即EM算法的求解结果与初值的选择有较大关系，该算法对初值敏感。

上述推导过程可由下图表示：

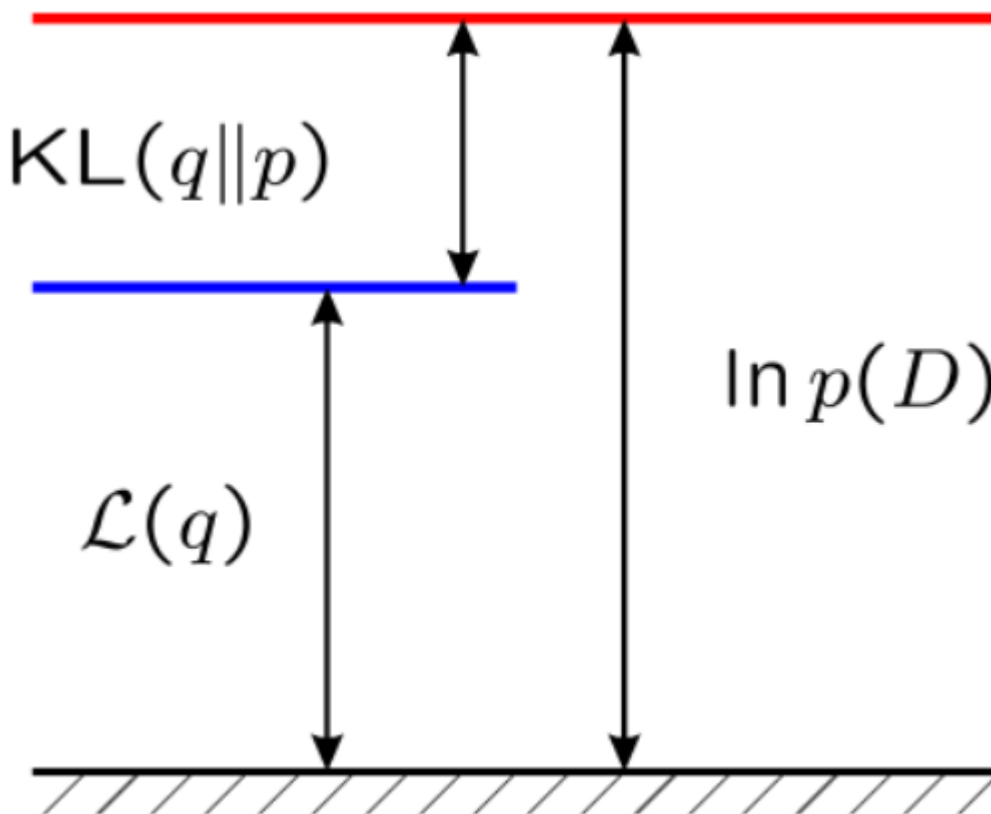


图1. EM算法

图1即对应式（8），可以看出 $\log P(Y|\theta)$ 由 $\frac{1}{n} \mathcal{L}(q, \theta)$ 和 $KL(q||p)$ 两部分组成。其中 $\frac{1}{n} \mathcal{L}(q, \theta)$ 即为 $\log P(Y|\theta)$ 的下界。

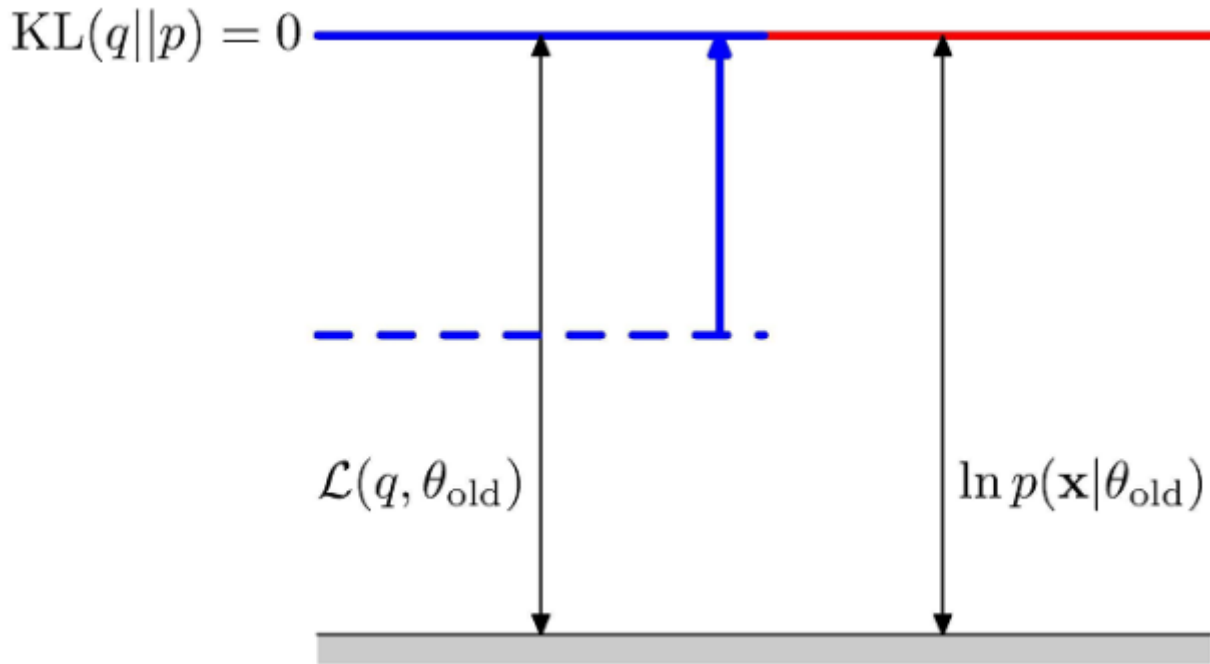


图2. EM算法

在M步中我们总期望最大化 $\frac{1}{n} \sum \log p(y_i|\theta)$ ，即使得 $\log P(Y|\theta)$ 的下界取得最大，也即最大化对数似然。故此时 $KL(q||p)$ 取得最小值为0。求解最大化 Q 函数，得到 $i+1$ 次迭代变量的估计 θ 。

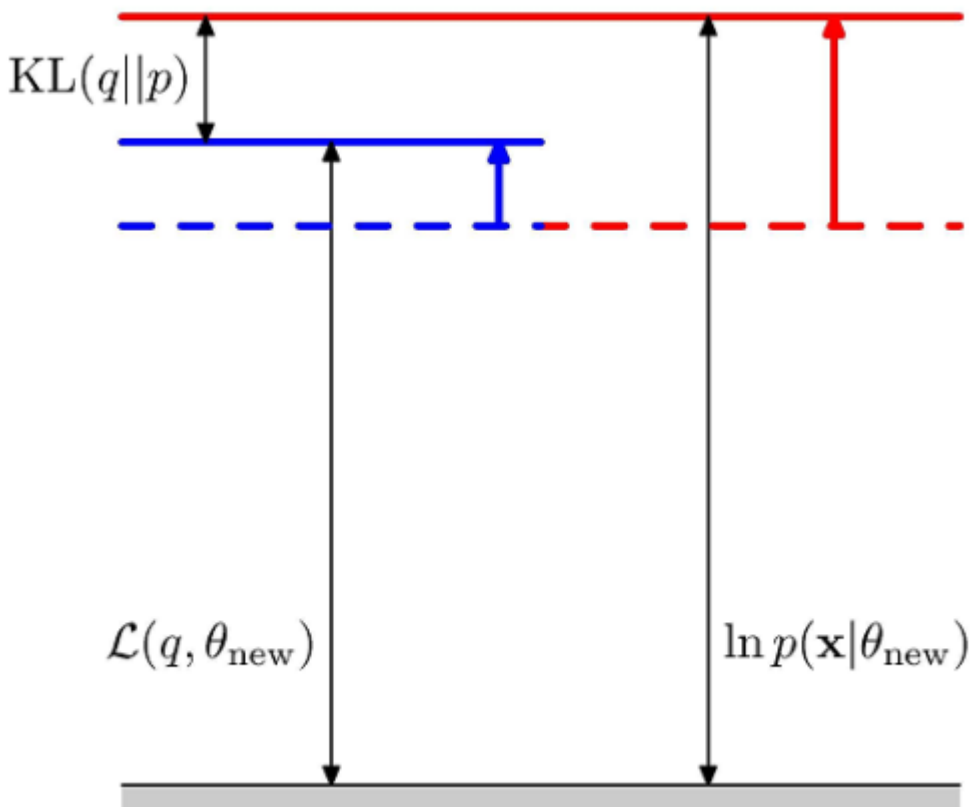


图3. EM算法

从图3中可以明显看出在 θ^i 更新后，对数似然 $L(\theta)$ 的下界 $\frac{1}{n} \sum \log p(y_i|\theta)$ 和 $KL(q||p)$ 均得到提高。此时在继续求解 $\theta^{i+1} = \arg\max_{\theta} Q(\theta, \theta^i)$ 。如此反复迭代，通过不断提高的 $L(\theta)$ 下界，使得其取得局部最大值。

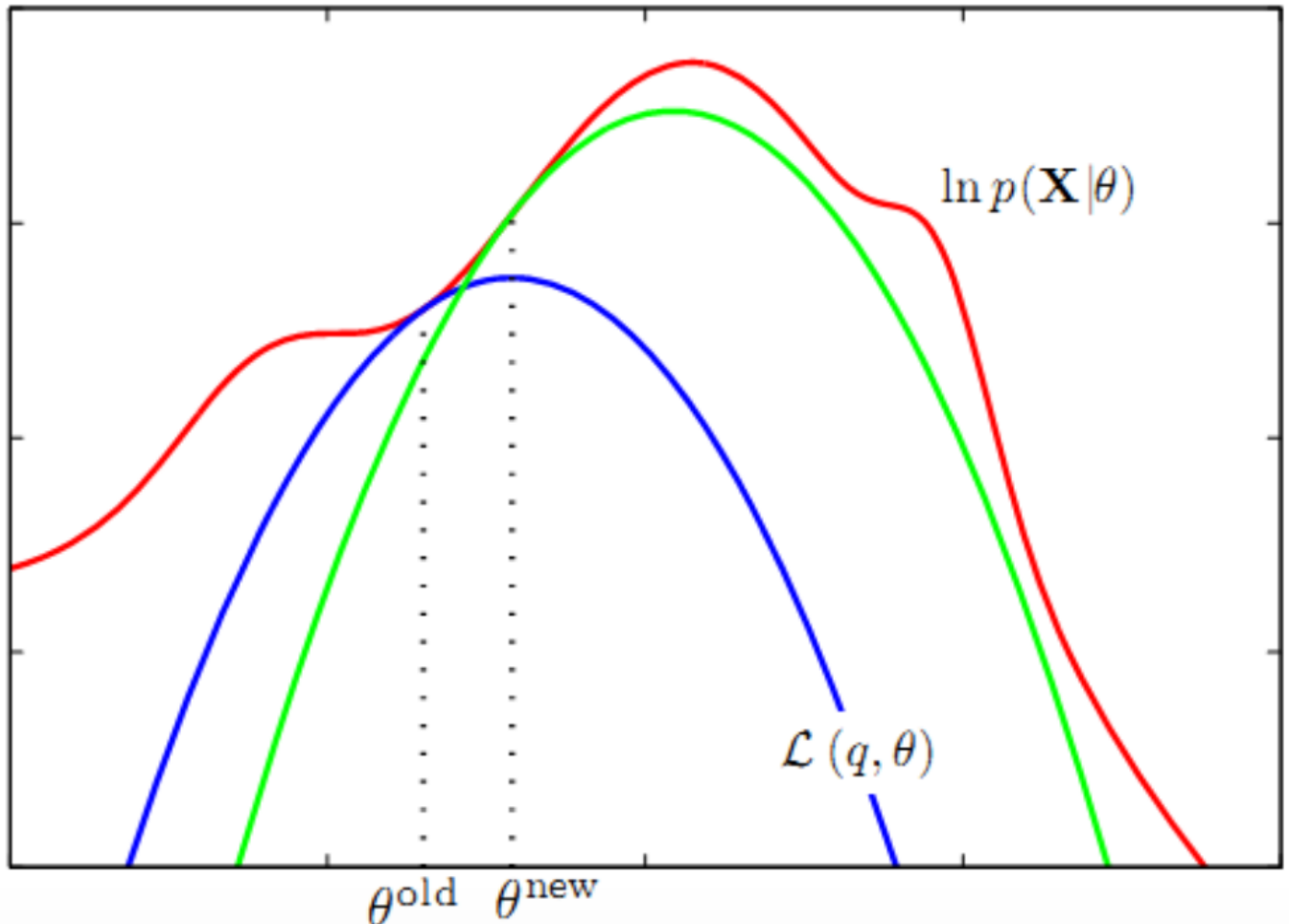


图4. EM算法迭代过程

从图4中我们也能明显看出，通过 θ^i 的反复迭代，我们不断提高对数似然的下界 $\frac{1}{n} \log L(q, \theta)$ 使之最后收敛于对数似然的局部最大解。

由上文讨论我们已经知道，通过EM反复迭代，其对数似然的下界将不断提高，然而我们却还是要问对于这种启发式的方法，即下界不断增大的这种方法，其等价于求解对数似然的最大值吗？或者说通过不断优化下界，算法就一定会收敛到似然函数的最大值吗？我们对此能否给出理论证明呢？

EM算法收敛性的理论证明

这里我们分别给出两种方法的理论证明。

收敛性证明方法一

这里主要利用变分思想，参照式（7）有：

$$\begin{aligned} \log(P(Y|\theta)) &= \int_{\mathcal{Z}} q(Z) \log \frac{P(Y, Z|\theta)}{P(Z|Y, \theta)q(Z)} \\ &= \int_{\mathcal{Z}} q(Z) \log \frac{P(Y, Z|\theta)}{q(Z)} - \int_{\mathcal{Z}} q(Z) \log P(Z|Y, \theta) \\ &= \int_{\mathcal{Z}} q(Z) \log \frac{P(Y, Z|\theta)}{q(Z)} - \text{KL}(q(Z) \| P(Z|Y, \theta)) \end{aligned}$$

由于 $\text{KL}(q(Z) \| P(Z|Y, \theta)) \geq 0$ 恒成立，且我们根据第 i 次迭代参数 θ^i 估计 $q(Z)$ 有 $q(Z) = P(Z|Y, \theta^i)$ 。故式（10）即为：

$$\log P(Y|\theta^{i+1}) \geq \int_{\mathcal{Z}} q(Z) \log \frac{P(Y, Z|\theta^{i+1})}{P(Z|Y, \theta^i)q(Z)} = \log P(Y|\theta^i) \quad (11)$$

故对于每一次的迭代均能保证 $\log P(Y|\theta^{i+1}) \geq \log P(Y|\theta^i)$ ，即可将EM算法理解为在变量坐标空间内，利用坐标下降法最大化对数似然下界的过程，故算法最终能够收敛至局部极小值点。

收敛性证明方法二

这里我们使用Jensen不等式进行证明，即对于凸函数 ϕ ，有 $\phi(E[f(x)]) \leq E[\phi(f(x))]$ ，故 $-\log E[f(x)] \leq E[-\log(f(x))] \implies E[\log(f(x))] \geq \log E[f(x)]$ （其中 $-\log(x)$ 为凸函数）。因此有：

$$\log(P(Y|\theta)) = \log \int_{\mathcal{Z}} q(Z) P(Y|\theta) \geq \int_{\mathcal{Z}} q(Z) \log P(Y|\theta) = \int_{\mathcal{Z}} q(Z) \log \frac{P(Y, Z|\theta)}{P(Z|Y, \theta)} \tag{12}$$

同样根据第 i 次迭代参数 θ^i 估计 $q(Z)$ ，则式(12)为：

$$\log P(Y|\theta^{i+1}) \geq \int_{\mathcal{Z}} q \log \frac{P(Y, Z|\theta^i)}{P(Z|Y, \theta^i)} P(Z|Y, \theta^i) = \log P(Y|\theta^i) \tag{13}$$

故算法最终能够收敛至局部极小值点。

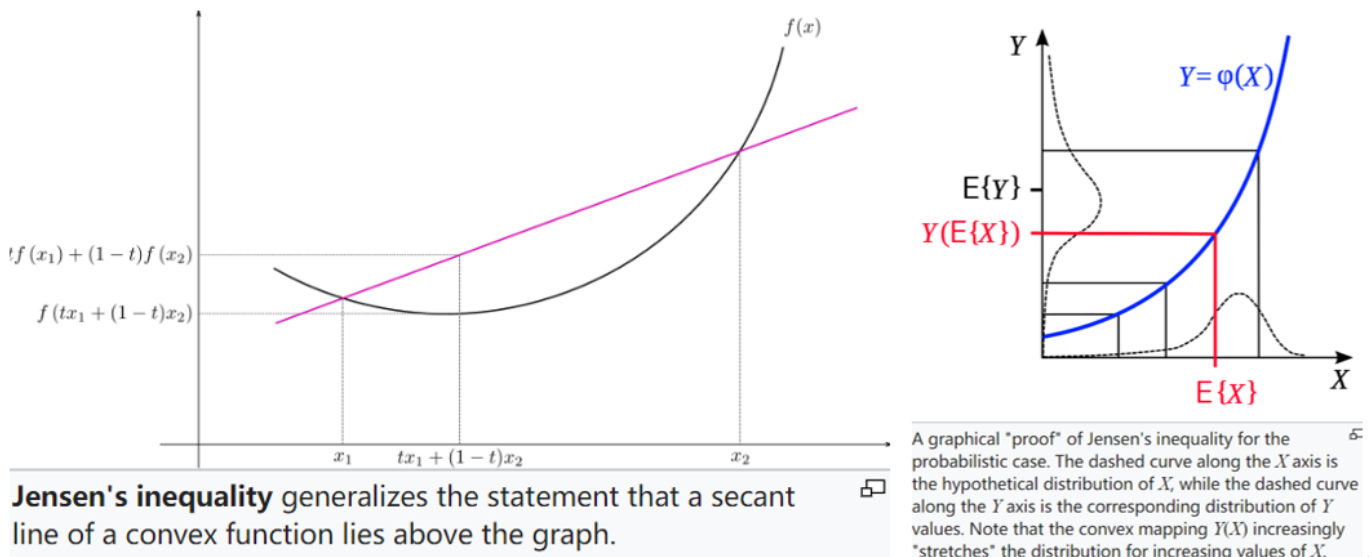


图5. Jensen不等式

GMM

与K-Means聚类不同，高斯混合聚类采用概率模型来刻画每个样本的簇类，即为一种“软划分”方法。这里我们首先回忆多元高斯模型。

对于 n 维样本空间 χ 中的随机向量 x ， x 服从高斯分布 $x \sim N(\mu, \Sigma)$ ，其概率密度函数如下：

$$p(x) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} e^{-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)} \tag{14}$$

其中 μ 为 n 维均值向量， Σ 为 $n \times n$ 协方差矩阵。如下图所示：

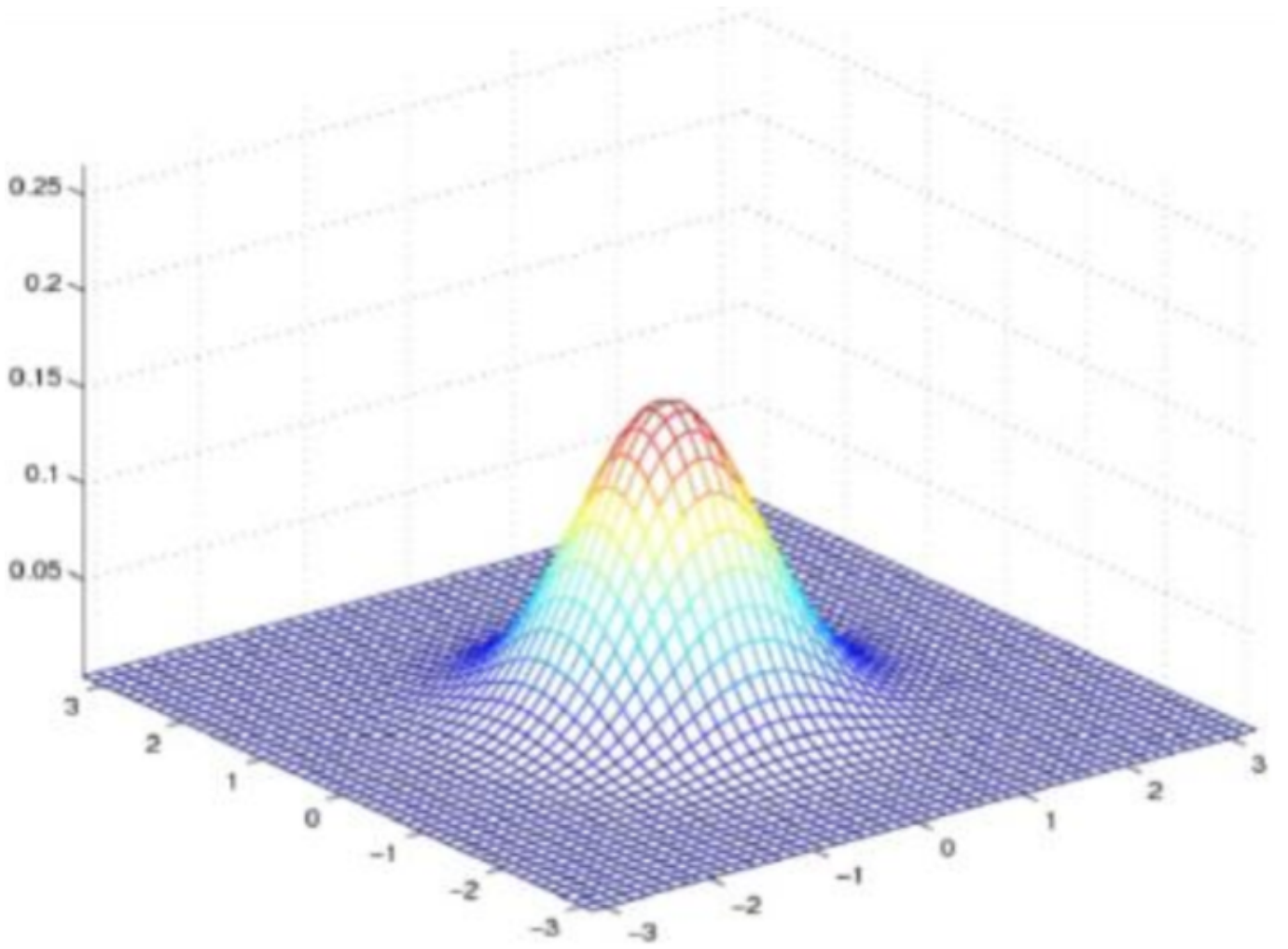


图6. 多元高斯变量

这里假设存在 k 个簇，且每一个簇均服从高斯分布。我们以概率 π_k 随机选择一个簇 k ，并从该簇的分布中采样样本点，如此得到观测数据 X ，则其似然函数为：

$$P(X|\theta) = P(X|\pi, \mu, \sum) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k N(x_n | \mu_k, \sum_k) \right\} \quad (15)$$

其中 $\sum_{k=1}^K \pi_k = 1, 0 \leq \pi_k \leq 1$

观察式（15）发现函数 $P(X|\theta)$ 由于 \log 中有求和运算，所有参数均耦合在一起，故求导困难，因而梯度下降优化较为困难。因此我们有必要采用一种新的优化算法。

这里首先我们令 $\frac{\partial P(X|\theta)}{\partial \mu_k} = 0$ ，则有：

$$-\sum_{n=1}^N \frac{\pi_k N(x_n | \mu_k, \sum_k)}{\sum_{j=1}^K \pi_j N(x_n | \mu_j, \sum_j)} \sum_k^{-1} (x_n - \mu_k) = 0 \quad (16)$$

这里我们记 $\gamma(z_{nk}) = \frac{\pi_k N(x_n | \mu_k, \sum_k)}{\sum_{j=1}^K \pi_j N(x_n | \mu_j, \sum_j)}$ ，则 $\gamma(z_{nk})$ 可以看为由参数 μ_k, \sum_k 对应的观测变量 x_n 的后验概率，即 x_n 从属于第 k 个簇的一种估计，或权值或“解释”。同时对式（16）左右两边同时乘以 \sum_k ，并进行移项操作，有：

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) x_n \quad N_k = \sum_{n=1}^N \gamma(z_{nk}) \quad (17)$$

同理我们令 $\frac{\partial P(X|\theta)}{\partial \sum_k} = 0$ ，有：

$$\sum_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (x_n - \mu_k)(x_n - \mu_k)^T \quad (18)$$

最后我们考虑混合系数即变量 π_k ，同理最大化对数似然 $\ln P(X|\pi, \mu, \sum)$ 。然而由式（15）知 π_k 需满足约束条件 $\sum_{k=1}^K \pi_k = 1$ ，故这里我们引入拉格朗日乘子法，即最大化下式：

$$\ln P(X|\pi, \mu, \sum) + \lambda (\sum_{k=1}^K \pi_k - 1) \quad (19)$$

式（19）对 π_k 求偏导为0有：

$$\sum_{n=1}^N \frac{\pi_k N(x_n|\mu_k, \sum_k)}{\sum_j \pi_j N(x_n|\mu_j, \sum_j)} + \lambda = 0 \quad (20)$$

上式两边同时乘以 π_k ，有：

$$\pi_k = -\frac{\sum_{n=1}^N \frac{\pi_k N(x_n|\mu_k, \sum_k)}{\sum_j \pi_j N(x_n|\mu_j, \sum_j)}}{\lambda} = \frac{-N_k}{\lambda} \quad (21)$$

这里我们将 $\sum_{k=1}^K \pi_k = 1$ 对 k 进行求和，则有 $\lambda = -N$ ，故：

$$\pi_k = \frac{N_k}{N} \quad (22)$$

这里需要注意的是由于 $\gamma(z_{nk}) = \frac{\pi_k N(x_n|\mu_k, \sum_k)}{\sum_j \pi_j N(x_n|\mu_j, \sum_j)}$ ，中仍存在隐变量 π_k ，并非为封闭解，故我们需要根据EM算法求解。具体如下：

- Step1. 初始化参数并计算对数似然；
- Step2. E步：依据当前模型参数，计算观测数据 x_i 属于簇 k 的概率（从属度）：

$$\gamma(z_{ik}) = \frac{\pi_k N(x_i|\mu_k, \sum_k)}{\sum_j \pi_j N(x_i|\mu_j, \sum_j)}$$

- Step3. M步：基于当前参数最大化对数似然函数，即重新求解新一轮迭代参数（第 $i+1$ 轮）：

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) x_n; \quad \sum_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (x_i - \mu_k)(x_i - \mu_k)^T; \quad \pi_k = \frac{\sum_{n=1}^N \gamma(z_{nk})}{N}$$

- Step4. 反复迭代直至收敛。

至此我们已经给出了EM算法求解GMM模型的具体方法。对比GMM与K-Means方法，我们可已看出由于概率的引入使得点到簇的从属关系为软分配，故其可以被用于非球形簇。

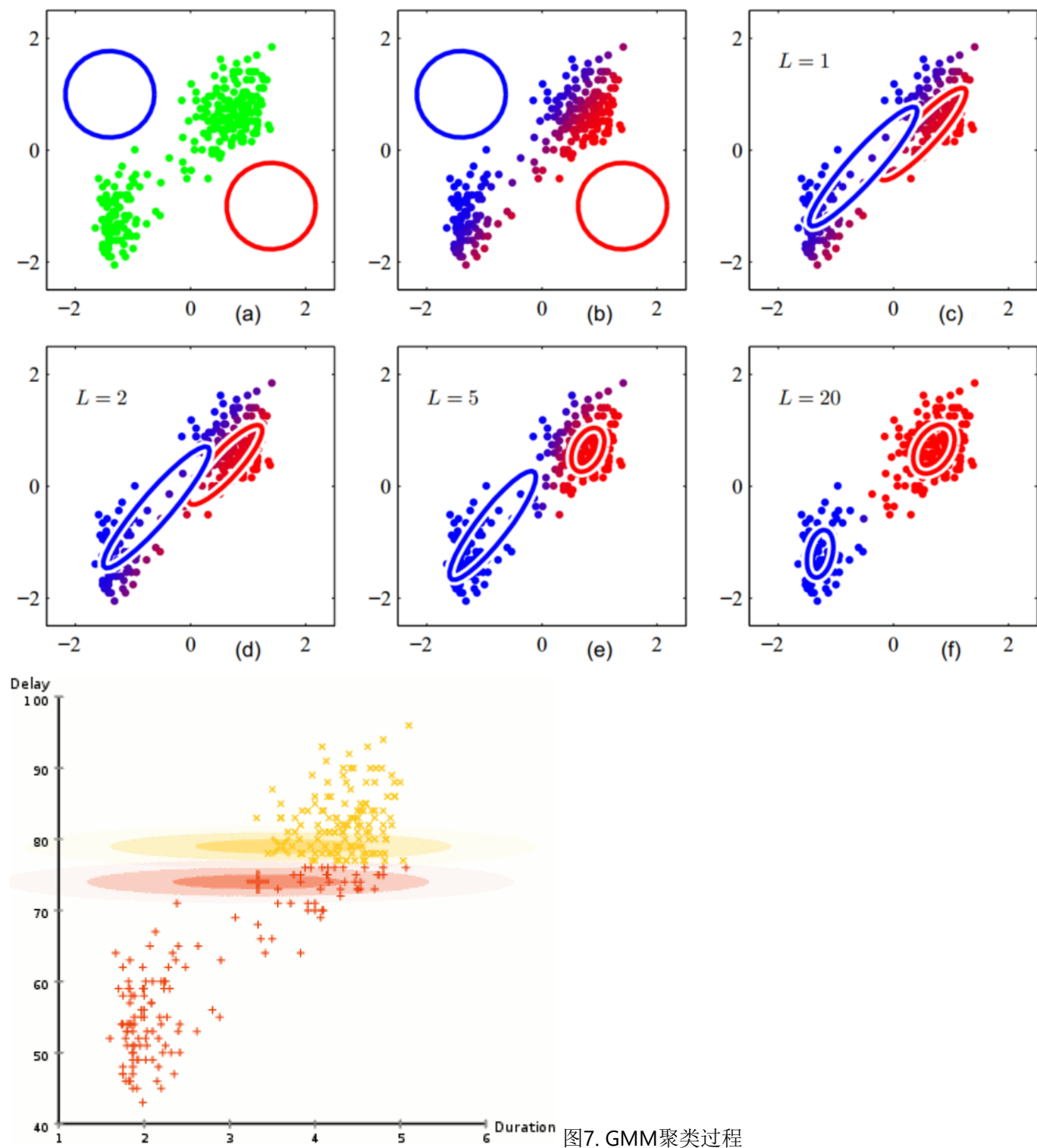


图7. GMM聚类过程

上图即为GMM算法的聚类过程。EM算法求解结果为局部最优解，其在隐变量的估计中应用广泛。

Different cluster analysis results on "mouse" data set:

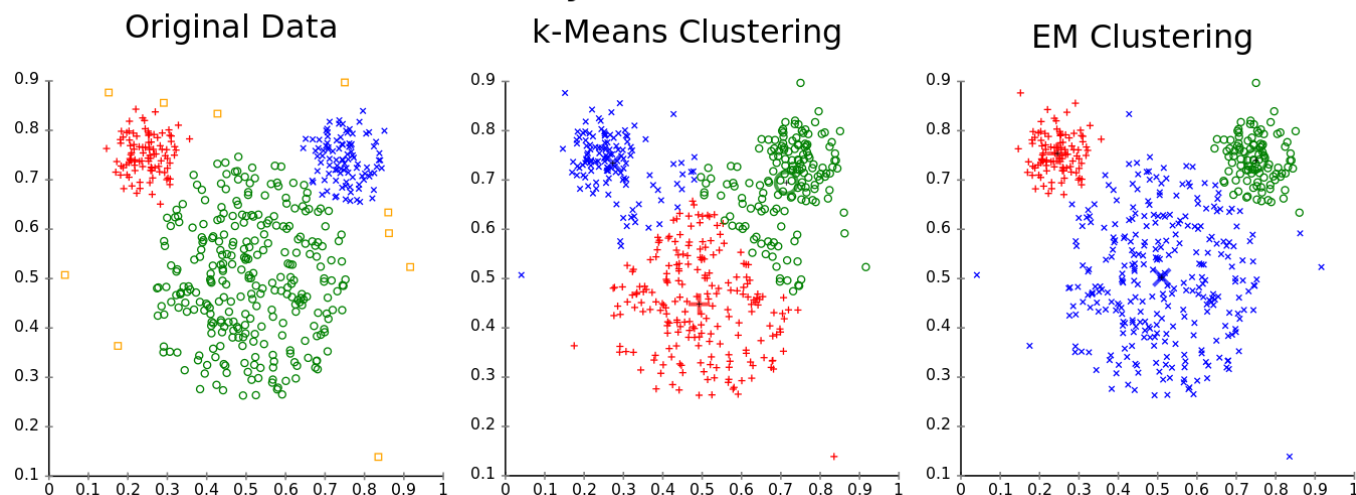


图8. GMM与K-Mean算法比较

由上图可以明显看出GMM相较于K-Means聚类有更佳的效果。

Variants

由于EM算法是只能收敛至局部极小值点，其对初值敏感。为克服这一缺陷，各种各样的启发式搜索算法如模拟退火法（其能较好的收敛至全局最优解）、随机重复爬山法等，通过多次随机的参数初始化或一定概率拒绝当前解从而使算法收敛至全局最优。此外卡尔曼滤波的思想同EM算法结合从而发展了Filtering and smoothing EM algorithms，以解决联合状态参数估计问题。共轭梯度与拟牛顿法也在EM中得到了应用。参数扩展期望最大化算法（PX-EM, parameter-expanded expectation maximization）通过协方差的调整引入额外的信息来修正M步中参数的估计以加速算法收敛。 α -EM由于不需要计算梯度或Hessi矩阵而使得算法收敛更快，同时也因而派生出了 α -HMM算法。

Reference

- [1] Dempster A P. Maximum likelihood estimation from incomplete data via the EM algorithm (with discussion[J]. Journal of the Royal Statistical Society, 1977, 39(1):1-38.
- [2] Jensen's inequality - Wikipedia
- [3] Bishop C M, 박원석. Pattern Recognition and Machine Learning, 2006[M]. Academic Press, 2006.
- [4] Expectation-maximization algorithm - Wikipedia