

变分推断

- 1.变分推断思路
- 2.分解概率分布推导过程
- 3.一元高斯分布实例
- 4.指数族分布
- 5.指数族分布+VI 的推导过程

一、变分推断思路

1.1 变分推断 (VI, Variational Inference) 背景:

贝叶斯公式: $p(\theta|X) = \frac{p(X|\theta)p(\theta)}{p(X)}$, 我们通过似然函数, 参数的先验分布和观测的数据求出后验概率。其分母部分通常是在数据中求积分得到, 需要得到这部分数据很难。我们的问题集中在如何求后验概率上, 关于如何得到后验概率一般分为两类方法: 精确推断和近似推断。近似推断又分为确定性近似和随机性近似, 下面介绍的变分推断则是确定性近似的一种, 而随机性近似的方法有 Gibbs Sampling、Gibbs Sampling 等。

现在围绕如何根据可观测数据变量 X 来计算潜在变量 Z 的后验概率分布 $p(Z|X)$ 展开讨论。假设有一个纯粹的贝叶斯模型, 其中每个参数都有一个先验概率分布, 这个模型也可以有隐含 (潜在) 变量及其参数, 我们将所有的潜在变量和参数组成一个集合, 记为 Z , 将观测变量的集合记为 X , 例如: $X = \{x_1, x_2, \dots, x_N\}$, $Z = \{z_1, z_2, \dots, z_N\}$, 则可以确定联合概率分布 $P(X, Z)$, 我们的目标是要找到后验概率 $P(Z|X)$ 的近似以及 $P(X)$ 的近似。

1.2 数学框架

假设 Z 符合 q 分布, 即 $q(Z)$ 。在贝叶斯公式 $p(X) = \frac{p(X, Z)}{p(Z|X)}$ 分子分母上下同除以 $q(Z)$:

$$\log P(x) = \log P(x, z) - \log P(z|x) = \log \frac{P(x, z)}{Q(z; \lambda)} - \log \frac{P(z|x)}{Q(z; \lambda)}$$

对 Z 取期望可得:

$$\begin{aligned} \ln p(X) &= E_q[\ln \frac{p(X, Z)}{q(Z)} - \ln \frac{p(Z|X)}{q(Z)}] = E_q[\ln \frac{p(X, Z)}{q(Z)}] - E_q[\ln \frac{p(Z|X)}{q(Z)}] \\ &= \int q(Z) \ln \frac{p(X, Z)}{q(Z)} dZ - \int q(Z) \ln \frac{p(Z|X)}{q(Z)} \\ &= ELBO(q(Z)) + KL(q(Z)||p(Z|X)) \quad (1.1) \end{aligned}$$

这个推导过程与EM算法里面是一样的, $ELBO$ 是 **Evidence Lower Bound**,

为什么称之为 $ELBO$ 呢? $p(x)$ 一般被称之为 evidence, 又因为 $KL(q||p) \geq 0$, 所以 $p(x) \geq E_{q(z; \lambda)}[\log p(x, z) - \log q(z; \lambda)]$, 这就是为什么被称为 $ELBO$

不同的是:

(1) EM算法里面有参数 θ , 而VI里面整合到随机变量 Z 里面去了。比如 $\ln p(X)$ 在EM算法中是 $\ln p(X|\theta) = ELBO(q(Z, \theta)) + KL(q(Z)||p(Z|X, \theta))$ 。

(2) EM算法的目的是为了解一个特定的 θ 优化值, 使观测数据 X 的对数似然函数 $\ln p(X|\theta)$ 最大, 思路是不断提升 $ELBO$ 来提升 $\ln p(X|\theta)$: 在E-step中首先给定当前 θ^t , 即设定概率分布 $q(Z)^t = p(Z|X, \theta^t)$ 使上

式右边第二项KL散度为0，则 θ^t 使得下界ELOB达到了当前限定的值，为 $\ln p(X|\theta)$ ；在M-step中，将E-step获得的分布 $q(Z)^t$ 代入到ELBO中，然后通过调节模型参数 θ 来最大化ELBO，得到下一轮的 θ^{t+1} ，若算法没收敛，此时新的 $KL(q||p)>0$ ，因此 $\ln p(X)$ 的提升值大于ELBO提升值。通过不断迭代E-step和M-step来不断提升ELBO从而最大化 $\ln p(X|\theta)$ 。

而VI的思想是希望用概率分布 $q(Z)$ 来近似模拟 $p(Z|X)$ ，很简单的想法就是希望 $KL(q||p) = 0$ ，通过上式可知，当给定数据后， $\ln p(X)$ 是一个固定的上界，则期望通过最大化ELOB来最小化KL。VI的变分体现在调整输入到ELOB中的概率分布函数 $q(Z)$ 来使得ELOB最大化。

什么是泛函？什么是变分法？[参见PRML附录D:变分法]

注意泛函与复合函数不同

1.3 思考

为了计算方便，却又能更充分、更灵活地提供对 $P(Z|X)$ 的近似概率分布 $q(Z)$ ，我们通常会选择一种限定类型的概率分布。核心思想包含两步：假设分布 $q(Z; \lambda)$ ，通过改变分布的参数 λ ，使 $q(z; \lambda)$ 靠近 $p(z|x)$ 总结成一句话就是，**为真实的后验分布引入了一个参数化的模型**。即：用一个简单的分布 $q(Z_i; \lambda_i)$ 拟合复杂的分布 $p(Z|X)$ 。这种策略将计算 $p(z|x)$ 的问题转化成优化问题了

$$\lambda^* = \arg \min_{\lambda} \text{divergence}(p(z|x), q(z; \lambda))$$

收敛后，就可以用 $q(z; \lambda)$ 来代替 $p(z|x)$ 了。

也可以用基于平均随机场的VI方法，对 $q(Z)$ 做了进一步的限定和简化，设定 $q(Z) = \prod_{i=1}^M q_i(Z_i; \lambda_i)$ ，然后对所有单个概率分布进行自由形式的变分最优化。

二、分解概率分布推导过程

平均场理论：将数量巨大的互相作用的多体问题转化成每一个粒子处在一种弱周期场中的单体问题，即：对某个独立的小个体，所有其他个体对它产生的作用可以用一个平均的量给出，如此，简化后的模型成为一个单体问题。

总而言之，平均场理论是用于简化复杂模型的理论，譬如对于一个概率模型： $P(x_1, x_2, x_3, \dots, x_n)$ ，利用平均场理论，找到另一个模型 $Q(x_1, x_2, x_3, \dots, x_n) = Q(x_1)Q(x_2)Q(x_3)\dots Q(x_n)$ ，使得 Q 尽量和 P 一致，并可以来近似代替 $p(x_1, x_2, x_3, \dots, x_n)$

基于平均场理论，将 Z 的元素划分为若干个互不相交的组，记为 Z_i 或 Z_j ，其中 $i, j = 1, 2, 3 \dots M$ 。假设 q 的分布关于这些分组可以分解，则 $q(Z) = \prod_{i=1}^M q_i(Z_i)$ ，将 $q(Z) = \prod_{i=1}^M q_i(Z_i)$ 代入ELBO式子：

$$\begin{aligned} L(q) &= \int q(Z) \ln p(X, Z) dZ - \int q(Z) \ln q(Z) dZ \\ &= \int \left(\prod_{i=1}^M q_i(Z_i) \right) \ln p(X, Z) dZ - \int \left(\prod_{i=1}^M q_i(Z_i) \right) \ln \prod_{i=1}^M q_i(Z_i) dZ \\ &= \int \left(\prod_{i=1}^M q_i(Z_i) \right) \ln p(X, Z) dZ - \int \left(\prod_{i=1}^M q_i(Z_i) \right) \sum_{i=1}^M \ln q_i(Z_i) dZ \\ &= L_1(q) - L_2(q) \quad (2.1) \end{aligned}$$

先处理第一项 $L_1(q)$ ：

$$L_1(q) = \int \left(\prod_{i=1}^M q_i(Z_i) \right) \ln p(X, Z) dZ$$

$$\begin{aligned}
&= \int_{Z_j} q_j(Z_j) \int_{\bar{Z}_j} \left(\prod_{i=1, i \neq j}^M q_i(Z_i) \right) \ln p(X, Z) d\bar{Z}_j dZ_j \\
&= \int_{Z_j} q_j(Z_j) \left[E_{q(\bar{Z}_j)} \ln p(X, Z) \right] dZ_j
\end{aligned}$$

其中: $E_{q(\bar{Z}_j)} \ln p(X, Z) = \int_{\bar{Z}_j} \left(\prod_{i=1, i \neq j}^M q_i(Z_i) \right) \ln p(X, Z) d\bar{Z}_j$

为了方便计算, 我们定义函数 $\ln \check{p}(X, Z_j) = E_{i \neq j} [\ln p(X, Z)]$, 这个式子将要刻画关于定义在所有 \bar{z}_j 上的 q 的概率分布。而函数 $\ln \check{p}(X, Z_j)$ 是没有归一化的, 其概率之和可能不等于1, 因此我们需要归一化 $\ln \check{p}(X, Z_j)$, 得到概率分布 $q_j^*(Z_j)$:

$$q_j^*(Z_j) = \frac{\check{p}(X, Z_j)}{\int_{Z_j} \check{p}(X, Z_j) dZ_j} = \frac{e^{E_{i \neq j} [\ln p(X, Z)]}}{\int_{Z_j} e^{E_{i \neq j} [\ln p(X, Z)]} dZ_j} = \frac{\check{p}(X, Z_j)}{C_1}$$

其中 $C_1 = \int_{Z_j} \check{p}(X, Z_j) dZ_j = \int_{Z_j} e^{E_{i \neq j} [\ln p(X, Z)]} dZ_j$, 则:

$$\ln \check{p}(X, Z_j) = \ln q_j^*(Z_j) + \ln C_1$$

因此, 第一项为:

$$\begin{aligned}
L_1(q) &= \int_{Z_j} q_j(Z_j) * (\ln q_j^*(Z_j) + \ln C_1) dZ_j \\
&= \int_{Z_j} q_j(Z_j) * \ln q_j^*(Z_j) dZ_j + \int_{Z_j} q_j(Z_j) * \ln C_1 dZ_j \\
&= \int_{Z_j} q_j(Z_j) * \ln q_j^*(Z_j) dZ_j + \ln C_1 \quad (2.2)
\end{aligned}$$

再处理第二项, $L_2(q) = \int \left(\prod_{i=1}^M q_i(Z_i) \right) \sum_{i=1}^M \ln q_i(Z_i) dZ$, 将其展开 (注意每一个函数 q_i 只涉及到一个变量 Z_i , 函数 q_i 是泛函变量):

$$\begin{aligned}
L_2(q) &= \int_{Z_1} \int_{Z_2} \dots \int_{Z_M} \left[q_1(Z_1) q_2(Z_2) \dots q_M(Z_M) \right] * \left[\ln q_1(Z_1) + \ln q_2(Z_2) \dots + \ln q_M(Z_M) \right] dZ_1 \dots dZ_{M-1} dZ_M \\
&= \sum_{i=1}^M \left(\int_{Z_1} \int_{Z_2} \dots \int_{Z_M} \ln q_i(Z_i) * \left[q_1(Z_1) q_2(Z_2) \dots q_M(Z_M) \right] dZ_1 \dots dZ_{M-1} dZ_M \right)
\end{aligned}$$

对求和内的第 i 项来说:

$$\begin{aligned}
&\int_{Z_1} \int_{Z_2} \dots \int_{Z_M} \ln q_i(Z_i) * \left[q_1(Z_1) q_2(Z_2) \dots q_M(Z_M) \right] dZ_1 \dots dZ_{M-i} dZ_M \\
&= \int_{Z_i} \left[\int_{Z_{j \neq i}} \dots \int \ln q_i(Z_i) * \left[q_1(Z_1) q_2(Z_2) \dots q_M(Z_M) \right] \left(\prod_{j \neq i}^M dZ_j \right) \right] dZ_i \\
&= \int_{Z_i} \ln q_i(Z_i) * \left[\int_{Z_{j \neq i}} \dots \int \left[q_1(Z_1) q_2(Z_2) \dots q_M(Z_M) \right] \left(\prod_{j \neq i}^M dZ_j \right) \right] dZ_i
\end{aligned}$$

上式中, 中括号内的部分, 对不包括 i 内的 $q(Z)$ 的所有其他维度积分, 相当于计算 $q(Z) = \prod_{i=1}^M q_i(Z_i)$ 的边缘概率, 即: 中括号内的值直接等于 $q_i(Z_i)$ 。则上式转化为:

$$\int_{Z_i} \ln q_i(Z_i) * q_i(Z_i) dZ_i$$

因此我们将每部分的结果代入 $L_2(q)$ 中得到第二部分：

$$\begin{aligned} L_2(q) &= \sum_{i=1}^M \left(\int_{Z_i} \ln q_i(Z_i) * q_i(Z_i) dZ_i \right) = \int_{Z_j} \ln q_j(Z_j) * q_j(Z_j) dZ_j + \sum_{i \neq j} \left(\int_{Z_i} \ln q_i(Z_i) * q_i(Z_i) dZ_i \right) \\ &= \int_{Z_j} \ln q_j(Z_j) * q_j(Z_j) dZ_j + \ln C_2 \quad (2.3) \end{aligned}$$

因为我们将 $q(Z) = \prod_{i=1}^M q_i(Z_i)$ 的每一项分开处理，现在只考虑第 j 项，则将上式中不涉及 Z_j 的项视为常数项 $\ln C_2$ 。

将两部分结果 (2.2) (2.3) 代入 (2.1) 中，得到 $L(q)$ ：

$$\begin{aligned} L(q) &= L(q_j) = \int_{Z_j} q_j(Z_j) * \ln q_j^*(Z_j) dZ_j - \int_{Z_j} \ln q_j(Z_j) * q_j(Z_j) dZ_j + \ln C_1 - C_2 \\ &= \int_{Z_j} q_j(Z_j) * \ln \frac{\ln q_j^*(Z_j)}{q_j(Z_j)} dZ_j + C = -KL(q_j || q_j^*) + C \end{aligned}$$

其中， $q_j^*(Z_j)$ 是归一化的 $\check{p}(X, Z_j)$ ，表示关于定义在所有 $i \neq j$ 的 Z_i 上的 q 概率分布的期望。

我们假设让所有的 $\{q_{i \neq j}\}$ 固定，分析上式，若要最大化 $L(q)$ ，相当于最大化 $KL(q_j || q_j^*)$ 的相反数，即最小化 $KL(q_j || q_j^*)$ ，而当 $q_j(Z_j) = q_j^*(Z_j)$ 时KL散度最小。因此，最优解的一般表达式为

$$\ln q_j^*(Z_j) = E_{q(\bar{Z}_j)} \ln p(X, Z) - \ln C_1 \quad (2.4)$$

这个解表示为了得到因子 q_j 最优解的对数，我们只需要考虑所有隐含变量 Z 和可见变量 x 上的联合概率分布的对数，然后关于所有其他因子 $\{q_{i \neq j}\}$ 取期望即可。

平均场方法通过恰当的分割独立变量子集 $\{Z_i\}_{i=1}^M$ ，使得相应的 $q_i(Z_i)$ 用闭合解。由于 q_i 的分布依赖于其他变量子集的分布，因此平均场是一个不断迭代的方法以保证各子集分布收敛到局部最优解(下界相对各子集分布是凸函数)。

至此我们已经从理论上找到了变分贝叶斯推断的通用公式求法：对于每一个 $q_j(Z_j)$ ：令 $q_j(Z_j) = q_j^*(Z_j)$ 循环，直到收敛。

虽然从理论上推导了变分推断的框架算法，但是在实践中，对于不同模型，最重要的是考虑如何对隐变量进行拆解，以及假设各变量子集服从何种分布。然后手动推导 $q_j^*(Z_j)$ ，简要说来，**推导变分贝叶斯模型一般分为四个步骤：**

- 1.确定好研究模型各个参数的共轭先验分布(如果想做full bayes model)
- 2.写出研究模型的联合分布 $P(Z, X)$
- 3.根据联合分布确定变分分布的形式 $q(Z)$
- 4.对于每个变分因子 $q_j(Z_j)$ 求出 $P(Z, X)$ 关于不包含变量 Z_j 的数学期望，再规整化为概率分布

当然这个过程并不简单，对于实际模型，其推导一般比较繁冗复杂，很容易出错，所以后来便有学者研究出更加一般更加自动化的基于概率图模型的算法框架——**VMP**(Variational Message Passing)。如果模型是指数族的模型，都可以套用VMP自动得到算法求解。

三、实例

假设给定一个数据集 $D = \{x_1, \dots, x_N\}$ ，其中每一个观测的单变量数据 x_n 均独立的抽样于**高斯分布**。那么似然函数为：

$$p(D|\mu, \tau) = \left(\frac{\tau}{2\pi}\right)^{N/2} \exp\left\{-\frac{\tau}{2} \sum_{i=1}^N (x_i - \mu)^2\right\}$$

其中 $\tau = \frac{1}{\sigma^2}$, 表示**精度**。我们现在的目标是估计参数 μ 和 τ 。我们引入两个参数的共轭先验分布

$$p(\mu|\tau) = N(\mu|\mu_0, (\lambda_0\tau)^{-1})$$

$$p(\tau) = Gam(\tau|a_0, b_0)$$

实质上, 对于这个简单问题的后验概率是能精确计算出来的, 由于共轭的性质, 参数的后验分布也是其先验分布。这里, 我们采用变分推断的思想进行后验概率的计算, 首先我们假设:

$$q(\mu, \tau) = q_\mu(\mu)q_\tau(\tau)$$

则对于 $q_\mu(\mu)$

$$\ln q_\mu^*(\mu) = E_\tau[\ln p(D|\mu, \tau) + \ln p(\mu|\tau)] + const = -\frac{E[\tau]}{2} \{\lambda_0(\mu - \mu_0)^2 + \sum_n (x_n - \mu)^2\} + const$$

注意我们只关注 μ , 其它无关项全融入常数项。由上式可以清楚推断 $q_\mu(\mu)$ 服从高斯分布 $N(\mu|\mu_N, \lambda_N^{-1})$, 注意我们并没有事先预设其 $q_\mu(\mu)$ 的分布形式, 完全推断与似然函数和预设的先验分布。其中

$$\mu_N = \frac{\lambda_0\mu_0 + N\bar{x}}{\lambda_0 + N}$$

$$\lambda_N = (\lambda_0 + N)E[\tau]$$

综上, 由 μ 的后验分布 $q_\mu^*(\mu)$, 可得知其期望为 μ_N 。注意当 $N \rightarrow \infty$, $\mu_N = \bar{x}$ 。同理, τ 的估计采用类似的方法。

四、指数族分布

4.1 指数族分布形式

如果一类分布可以写成如下形式, 则它可以叫做 **指数族分布**(exponential family distributions)。

$$P(X|\eta) = h(X)\exp\{\eta^T\Phi(X) - A(\eta)\} \quad (4.1)$$

η 是自然参数 (也叫 canonical parameter), 其中 X 是观测数据, $\Phi(X)$ 是观测数据的充分统计量(sufficient statistic), $A(\eta)$ 是对数归一化项 (也叫对数分配函数, log partition function, log normalization)。

其中: $h(X)$ 只含有 X , 而 $A(\eta)$ 只含有自然参数 η , $\eta^T\Phi(X)$ 相当于数据和参数的一个内积, 这样把参数和数据分开使得很多计算变得简单很多, 而且也产生了一定的规律和结论, 后面会详细介绍。

大多数概率分布都是指数族分布, 都可以表示成上面公式给出的形式:

- 1) 伯努利分布: 对 0、1 问题进行建模
- 2) 多项式分布: 对 K 个离散结果的事件建模
- 3) 泊松分布: 对计数过程进行建模
- 4) 伽马分布与指数分布: 对间隔的正数进行建模
- 5) Beta 分布: 对小数进行建模
- 6) Dirichlet 分布: 对小数进行建模

- 7) Wishart分布：对协方差进行建模
- 8) 高斯分布

4.2 一元高斯分布函数化成指数族分布的形式

一元高斯分布函数形式如下：

$$P(X|\theta) = P(X|\mu, \delta) = \frac{1}{\sqrt{2\pi\delta^2}} e^{-\frac{(x-\mu)^2}{2\delta^2}}$$

目的是将此式子化为 (4.1) 式的样式。

观察 (4.1) 式，其实 η 也是参数的一个函数，这个地方目的是要将参数 θ 映射到对应的 η 中，并且找到对应的对数归一化项 $A(\eta)$ 和前面的系数 $h(X)$ 。

推导过程如下：

$$\begin{aligned} P(X|\theta) &= \frac{1}{\sqrt{2\pi\delta^2}} e^{-\frac{(x-\mu)^2}{2\delta^2}} = e^{\log \frac{1}{\sqrt{2\pi\delta^2}}} * e^{-\frac{(x-\mu)^2}{2\delta^2}} \\ &= \exp\left\{\log \frac{1}{\sqrt{2\pi\delta^2}} - \frac{(x^2 - \mu x + \mu^2)}{2\delta^2}\right\} \\ &= \exp\left\{\begin{bmatrix} \frac{\mu}{\delta^2} & \frac{-1}{2\delta^2} \end{bmatrix} \begin{bmatrix} x \\ x^2 \end{bmatrix} + (\log 2\pi\delta^2)^{-1/2} - \frac{\mu^2}{\delta^2}\right\} \quad (4.2) \end{aligned}$$

可以看出式 (4.1) 中 $h(X) = 1$ 且：将 η 命为 $[\eta_1 \quad \eta_2]^T = \begin{bmatrix} \frac{\mu}{\delta^2} & \frac{-1}{2\delta^2} \end{bmatrix}^T$ ，将 $\Phi(X)$ 命为 $\begin{bmatrix} x \\ x^2 \end{bmatrix}$ ，反过来解得：

$$\mu = -\frac{\eta_1}{2\eta_2}; \delta^2 = -\frac{1}{2\eta_2} \quad (4.3)$$

将上式代入到 (4.2) 的后半部分 $(\log 2\pi\delta^2)^{-1/2} - \frac{\mu^2}{\delta^2}$ 中，可以得到这部分关于 η 的表达式，如下：

$$A(\eta) = -\frac{\eta_1^2}{4\eta_2} + \frac{1}{2} \log \frac{-\pi}{\eta} \quad (4.4)$$

总结：根据上面的推导过程，将一元高斯分布概率密度函数转化为了标准的指数族分布函数，其中 $\eta = [\eta_1 \quad \eta_2]^T$ ， $\Phi(X) = \begin{bmatrix} x \\ x^2 \end{bmatrix}$ ， $A(\eta)$ 为 (4.4)， $h(X) = 1$ 。

4.3 关于 $A(\eta)$

根据 (4.1)，将指数部分的 $A(\eta)$ 分开，式子转化为：

$$P(X|\eta) = \exp\{A(\eta)\}^{-1} h(X) \exp\{\eta^T \Phi(X)\} \quad (4.5)$$

上式左右两边对 x 积分，可以得到：

$$1 = \exp\{A(\eta)\}^{-1} \int h(X) \exp\{\eta^T \Phi(X)\} dx$$

因此：

$$\exp\{A(\eta)\} = \int h(X) \exp\{\eta^T \Phi(X)\} dx \quad (4.6)$$

这也是为什么将 $A(\eta)$ 称为对数归一化项的原因。式子(4.6)中，左右两边都对 η 求导可得：

$$\begin{aligned} A'(\eta) \exp\{A(\eta)\} &= \frac{\partial}{\partial \eta} \int h(X) \exp\{\eta^T \Phi(X)\} dx \\ A'(\eta) \exp\{A(\eta)\} &= \int h(X) \exp\{\eta^T \Phi(X)\} \Phi(X) dx \\ A'(\eta) &= \frac{\int h(X) \exp\{\eta^T \Phi(X)\} \Phi(X) dx}{\exp\{A(\eta)\}} = \int h(X) e^{\eta^T \Phi(X) - A(\eta)} \Phi(X) dx \\ &= \int P(X|\eta) \Phi(X) dx = E_{P(X|\eta)}[\Phi(X)] \end{aligned}$$

因此可以得到结论：

$$A'(\eta) = E_{P(X|\eta)}[\Phi(X)] \quad (4.7)$$

即：对数归一化项 $A(\eta)$ 的一阶导数等于充分统计量的期望。

这里也可以理解为什么 $\Phi(X)$ 称为充分统计量，因为我们只需要知道 $\Phi(X)$ 就可以刻画这个分布了，比如一

元高斯分布中的 $\Phi(X) = \begin{bmatrix} x \\ x^2 \end{bmatrix}$ ，只需要知道 $E_{x \sim P(X|\eta)}[\Phi(x)]$ 和 $E_{x \sim P(X|\eta)}[\Phi(x^2)]$ 就可明确这个分布了。

继续推导还可以得出：

$$A''(\eta) = \text{Var}[\Phi(X)] \quad (4.8)$$

其中： $A'(\eta)$ 是凸函数。

对于式子4.7的理解：

(1) 在极大似然估计中，我们的目标是求参数 η 的值使得对数似然函数 $P(X|\eta)$ 最大，此时 $P(X|\eta)$ 对参数 η 求偏导然后令导数等于0，则可以求出 η_{MLE} 的值。对于指数族分布，我们可以证明：

$$\frac{\partial}{\partial \eta} \sum_{i=1}^n [\log P(X|\eta)] = \sum_{i=1}^n [\Phi(X_i) - A'(\eta)]$$

令上式为0，可以得到：

$$A'(\eta) = E_{P(X|\eta)}[\Phi(X)]$$

与(4.7)的结论相同，因此我们可以得到求解极大似然估计的简单方法：

对于指数族分布，先将分布的对数归一化项对参数求导得到 $A'(\eta)$ （等式左边），然后将充分统计量对分布 $P(X|\eta)$ 求均值（等式右边），让两者相等，得到参数值 η_{MLE} 。

(2) 在贝叶斯推断中，对于不同的参数值 η ， η 可以取 $\eta_1, \eta_2, \dots, \eta_N$ 。不同的 η_i 决定不同的概率密度函数 $P(X|\eta_i)$ ，所求出来的期望（等号右边） $E_{P(X|\eta)}[\Phi(X)]$ 也不同，即：等号右边的值是由参数 η 决定的。因此得到一个求期望的新方法：将 $A'(\eta)$ 对 η 求导，然后将 η_i 的值代入，则可以得到等式右边的期望。这将应用在指数族分布的变分推断中。

4.4 指数族分布的共轭先验

根据贝叶斯公式可知：

$$P(\theta|X) \propto P(X|\theta) * P(\theta)$$

共轭表示：似然函数 $P(X|\theta)$ 与先验分布 $P(\theta)$ 共轭，即后验分布 $P(\theta|X)$ 与先验分布 $P(\theta)$ 的形式相同，只是参数不同。若 $P(X|\theta)$ 服从指数族分布，则

$$P(X|\theta) = h(X) \exp\{\theta^T \Phi(X) - A(\theta)\}$$

结论：理论上来说，任何一个指数族似然函数，都必定有一个指数族的先验分布与其共轭，使得后验分布与先验分布有相同的形式。而此时，只需要另先验分布的充分统计量 $\Phi(X)$ 的第二部分与似然函数的对数归一化项 $A_1(\theta)$ 相等即可。下证此结论。

假设先验分布 $P(\theta)$ 和似然函数都服从指数族分布，并为先验分布引入参数 α ，则先验分布为：

$$P(\theta|\alpha) = h(\theta) \exp\{\alpha^T \Phi(\theta) - A_1(\alpha)\} \quad (4.9)$$

似然函数为：

$$P(X|\theta) = h(X) \exp\{\theta^T \Phi(X) - A_2(\theta)\} \quad (4.10)$$

假设

$$\Phi(\theta) = \begin{bmatrix} \theta \\ -g(\theta) \end{bmatrix}, \quad \alpha = \begin{bmatrix} \alpha_1 \\ \alpha_2 \end{bmatrix}$$

其中 θ 和 α_1 是标量，而 $g(\theta)$ 和 α_2 是矢量且长度相同。我们的目的是让（4.9）和（4.10）相乘的结果有着标准的指数族分布的形式。（4.9）和（4.10）相乘(观测数据 X 视为常量, $A_1(\alpha)$ 视为常量)：

$$\begin{aligned} P(\theta|\alpha)P(X|\theta) &= h(\theta) \exp\{\alpha^T \Phi(\theta) - A_1(\alpha)\} * h(X) \exp\{\theta^T \Phi(X) - A_2(\theta)\} \\ &\propto h(\theta) \exp\{\alpha^T \Phi(\theta) - A_2(\theta) + \theta^T \Phi(X)\} \\ &\propto h(\theta) \exp\{\alpha_1 \theta - \alpha_2 g(\theta) - A_2(\theta) + \theta \Phi(X)\} \\ &= h(\theta) \exp\{(\alpha_1 + \Phi(X))\theta - \alpha_2 g(\theta) - A_2(\theta)\} \end{aligned}$$

若令 $g(\theta) = A_2(\theta)$ ，则上式可以简化为：

$$P(X|\theta)P(\theta|\alpha) \propto h(\theta) \exp\{(\alpha_1 + \Phi(X))\theta - (\alpha_2 + 1)A_2(\theta)\} \quad (4.11)$$

令 $\hat{\alpha}_1 = (\alpha_1 + \Phi(X)), \hat{\alpha}_2 = (\alpha_2 + 1)$ 则：

$$P(X|\theta)P(\theta|\alpha) \propto h(\theta) \exp\{\hat{\alpha}_1 \theta - \hat{\alpha}_2 A_2(\theta)\} = h(\theta) \exp\{\hat{\alpha} \Phi(\theta)\} \quad (4.11)$$

式子（4.11）符合指数族分布的标准形式，即：若令 $g(\theta) = A_1(\alpha)$ 两者相乘后得到的也是一个指数族分布，反过来也成立。证毕。

结论是：让参数的先验分布的第二部分 $g(\theta)$ ，取似然函数的充分统计量取似然函数的对数归一化项 $A_2(\theta)$ ，则先验分布与似然函数共轭，得到的后验概率分布的参数为 $\hat{\alpha}_1 = (\alpha_1 + \Phi(X)), \hat{\alpha}_2 = (\alpha_2 + 1)$ 。

五、指数族分布的变分推断

变分推断的目的是找到一个分布 $q(Z)$ 去近似于后验概率 $p(Z|X)$ 。即：将式子（1.1）中的 KL 散度最小化，而数据是固定的，即等价于将 $ELBO$ 最大化：

$$\begin{aligned} ELBO : L(q) &= \int q(Z) \log p(X, Z) dZ - \int q(Z) \log q(Z) dZ \\ &= E_{q(Z)} \log p(X, Z) - E_{q(Z)} \log q(Z) \quad (5.1) \end{aligned}$$

在第二部分推导过程中，将模型中的变量分为了观测变量和隐含变量，现在将**隐含变量**拆开，比如拆解成两部分，例如： $Z = \{Z, \beta\}$ ，其中我们选择的两个分布 Z, β 是独立的，类似于第二部分所述的 $q(Z) = \prod_{i=1}^M q_i(Z_i)$ 。

对于 $\{Z, \beta\}$ ，其后验概率有如下推导：

$$P(\beta, Z|X) = P(\beta|Z, X)P(Z|X) = P(Z|\beta, X)P(\beta|X) \quad (5.2)$$

对于 (5.2) 中的 $P(\beta|Z, X)$ 和 $P(Z|\beta, X)$ ，现假设他们为指数族分布，将其写成指数族分布的标准形式：

$$P(\beta|\eta(Z, X)) = h(\beta) \exp\{\eta(Z, X)^T \Phi(\beta) - A_1(\eta(Z, X))\} \quad (5.3)$$

$$P(Z|\eta(\beta, X)) = h(Z) \exp\{\eta(\beta, X)^T \Phi(Z) - A_2(\eta(\beta, X))\} \quad (5.4)$$

我们的目标变换为：想用个变分函数 $q(\beta)$ 近似于 $P(\beta|Z, X)$ ，其中 $q(\beta)$ 也是指数族分布，**引入参数 λ** ， $q(\beta|\lambda)$ 的形式如下：

$$q(\beta|\lambda) = h(\beta) \exp\{\lambda^T \Phi(\beta) - A(\lambda)\} \quad (5.5)$$

不断优化 λ 的值，使式 (5.5) 趋近于式 (5.3)

同理，想用个变分函数想 $q(Z)$ 近似于 $P(Z|\beta, X)$ ，其中 $q(Z)$ 也是指数族分布，**引入参数 γ** ， $q(Z|\gamma)$ 的形式如下：

$$q(Z|\gamma) = h(Z) \exp\{\gamma^T \Phi(Z) - A(\gamma)\} \quad (5.6)$$

不断优化 γ 的值，使式 (5.6) 趋近于式 (5.4)

在式 (5.1) 中， $L(q)$ 中的 q 为隐含变量 $\{Z, \beta\}$ 的函数，而 Z 依赖于参数 λ ， β 依赖于参数 γ ，即：可以将ELBO记为： $L(q(\lambda, \gamma))$ 。对式子 $L(q(\lambda, \gamma))$ ，若想优化参数，可以使用常用的迭代优化，即：先固定 γ ，优化 λ ，然后固定 λ ，优化 γ ，不停迭代。（这里是两个参数，同理，若有多个参数则需要固定其他的所有参数来优化其中一个参数）

式 (5.1) 将 Z 拆分，且是独立的，即 $q(Z, \beta) = q(Z)q(\beta)$ ，并引入参数 λ, γ ，变换为：

$$\begin{aligned} L(q(\lambda, \gamma)) &= E_{q(Z, \beta)} \log p(X, Z, \beta) - E_{q(Z, \beta)} \log q(Z, \beta) \\ &= E_{q(Z, \beta)} \log [P(\beta|Z, X)P(Z|X)P(X)] - E_{q(Z, \beta)} \log [q(\beta)q(Z)] \\ &= E_{q(Z, \beta)} \log P(\beta|Z, X) + E_{q(Z, \beta)} \log P(Z|X) + E_{q(Z, \beta)} \log P(X) - E_{q(Z, \beta)} \log q(\beta) - E_{q(Z, \beta)} \log q(Z) \quad (5.7) \end{aligned}$$

我们**首先固定 γ** ，**优化 $q(\beta|\lambda)$** （此时由于固定了 γ ， $q(Z|\gamma)$ 也固定，将 γ 、 $q(Z|\gamma)$ 视为常数），式子 (5.7) 中的第二项与 λ （ $q(\beta)$ ）无关，第三项也无关，第五项也无关，剩下的只有第一项和第四项，其他的都视为常数，则 (5.7) 化简为：

$$E_{q(Z, \beta)} \log P(\beta|Z, X) - E_{q(Z, \beta)} \log q(\beta) + \text{constant} \quad (5.8)$$

想要使ELBO max,则优化式 (5.8)。对式 (5.8) 的指数族分布展开，即代入式 (5.3) (5.5)：

$$\begin{aligned} (5.8) &= E_{q(Z, \beta)} \log [h(\beta) \exp\{\eta(Z, X)^T \Phi(\beta) - A_1(\eta(Z, X))\}] - E_{q(Z, \beta)} \log [h(\beta) \exp\{\lambda^T \Phi(\beta) - A(\lambda)\}] \\ &= \left\{ E_{q(Z, \beta)} \log [h(\beta)] + E_{q(Z, \beta)} [\eta(Z, X)^T \Phi(\beta)] - E_{q(Z, \beta)} [A_1(\eta(Z, X))] \right\} - \end{aligned}$$

$$\left\{ E_{q(Z,\beta)} \log[h(\beta)] + E_{q(Z,\beta)} [\lambda^T \Phi(\beta)] - E_{q(Z,\beta)} [A(\lambda)] \right\} \quad (5.9)$$

上式中，第一项和第四项约掉，第三项只有 Z 没有 β ，视为常数，上式化简为：

$$E_{q(Z,\beta)} [\eta(Z, X)^T \Phi(\beta)] - E_{q(Z,\beta)} [\lambda^T \Phi(\beta)] + E_{q(Z,\beta)} [A(\lambda)] + \text{constant}$$

$$\text{将 } Z, \beta \text{ 分开} = E_{q(Z)} [\eta(Z, X)^T] * E_{q(\beta)} \log[\Phi(\beta)] - E_{q(\beta)} [\lambda^T \Phi(\beta)] + A(\lambda) \quad (5.10)$$

上式中的相乘项的因子 $E_{q(\beta)} \log[\Phi(\beta)]$ 。由上一节的公式 (4.7) 可知， $E_{q(\beta)} \log[\Phi(\beta)] = A'(\lambda)$ ， $E_{q(\beta)} [\lambda^T \Phi(\beta)] = \lambda^T * E_{q(\beta)} [\Phi(\beta)] = \lambda^T A'(\lambda)$ 。则，上式 (5.10) 化简为：

$$A'(\lambda) E_{q(Z)} [\eta(Z, X)^T] - \lambda^T A'(\lambda) + A(\lambda) + \text{constant} \quad (5.11)$$

若要最大化 (5.11)，只需要对 η 求导，即：

$$\frac{\partial L(q(\lambda, \gamma))}{\partial \lambda} = A''(\lambda) E_{q(Z)} [\eta(Z, X)^T] - \lambda^T A''(\lambda)$$

$$= A''(\lambda) (E_{q(Z)} [\eta(Z, X)^T] - \lambda^T)$$

令上式为0，解得：

$$\lambda = E_{q(Z)} [\eta(Z, X)^T] \quad (5.12)$$

其中， $A'(\lambda)$ 是凸函数。

这样**固定 γ 得到了 λ** ，同理，在 (5.7) 中通过**固定 λ 得到了 γ** 的结果为：

$$\gamma = E_{q(\beta)} [\eta(\beta, X)^T] \quad (5.13)$$

即：若要优化 $q(\beta|\lambda)$ ，则将 P 的参数分解为独立的参数，将剩下的所有参数的分布求期望则得到想要优化的 $q(\beta|\lambda)$ 的 λ 值。在第二部分中，我们将 $q(Z) = \prod_{i=1}^M q_i(Z_i)$ ，这里我们若也将 Z 分解为独立的 Z_i ，不同的 $q_i(Z_i)$ 由不同的参数决定，且与 $q_i(\eta)$ 耦合，其中 η 是除去 Z_i 后的所有参数。我们可以使用迭代优化，固定其他所有的分布来优化一个分布的参数，其结果等于计算 $E_{q_i(Z_i)} [\eta^T]$ ；然后再固定其他所有的分布优化另一个分布的参数，直至收敛，达到ELBO的最大值。

references:

- 1.PRML 10.1
- 2.PRML 2.4
- 3.PRML附录
- 4.https://blog.csdn.net/step_forward_ML/article/details/78077383
- 5.<https://www.scribd.com/document/388036022/Notes-on-Expectation-Maximization-and-Variational-Inference-pdf>
- 6.<https://www.scribd.com/document/388036017/Notes-on-Varational-Inference-pdf>
- 7.<https://www.scribd.com/document/388036021/VI简明推导-pdf>
- 8.<https://www.youtube.com/watch?v=arMoli91OZE&list=PLFze15KrfxbF0n1zTNoFlaDpxnSyfgNgc>