

变分推断 Variational Inference

一、背景

频率角度 → 优化问题

$$\left\{ \begin{array}{l} \text{① } f(\omega) = \omega^T x + b, \text{ loss function} \\ \hat{\omega} = \arg \max L(\omega) \\ \left. \begin{array}{l} \text{① 解析解: } \hat{\omega} = (X^T X)^{-1} X^T Y \\ \text{② 数值解: 可用 GD, SGD 求解} \end{array} \right. \\ \text{SVM (分类): } f(\omega) = \text{sign}(\omega^T x + b), \\ \text{loss function: } \begin{cases} \min \sum_{i=1}^N w^T w \\ \text{s.t. } y_i(w^T x_i + b) \geq 1 \end{cases} \quad \text{有约束.} \\ \text{可用 Lagrange 乘子法或 QP 求解.} \\ EM: \hat{\theta} = \arg \max_{\theta} \log P(X|\theta) \\ \theta^{(t+1)} = \arg \max_{\theta} \int_z \log P(X, z|\theta) P(z|X, \theta^{(t)}) dz. \end{array} \right.$$

贝叶斯角度 → 积分问题.

$$P(\theta|X) = \frac{P(X|\theta)P(\theta)}{P(X)} \rightarrow \text{先验}$$

贝叶斯推断: 指在贝叶斯框架中求出后验概率.

贝叶斯决策: 可理解为预测.

$X \rightarrow N$ 个样本.

$$\begin{aligned} \text{对于新样本 } \tilde{x}, \text{ 有 } P(\tilde{x}|X) &= \int_{\theta} P(\tilde{x}, \theta|X) d\theta \\ &= \int_{\theta} P(\tilde{x}|\theta) P(\theta|X) d\theta \end{aligned}$$

Inference 求出后验, 就可以做决策了.

后验

$$= E_{\theta|X}[P(\tilde{x}|\theta)]$$

二、变分推断一公式推导.

X : 观测变量

Z : 隐变量 + 参数,

(X, Z) : 完整数据

将参数看作隐变量，这里就不写出来了。

$$\log \underline{P(X)} = \log P(X, Z) - \log P(Z|X)$$

$$= \log \frac{P(X, Z)}{q(Z)} - (\log \frac{P(Z|X)}{q(Z)})$$

等式两边对 Z 求期望：

$$\begin{aligned} \text{左边} &= \int_Z \log P(X) q(Z) dZ = \log P(X) \int_Z q(Z) dZ = \log P(X) \\ \text{右边} &= \underbrace{\int_Z q(Z) \cdot \log \frac{P(X, Z)}{q(Z)} dZ}_{\text{ELBO}} - \underbrace{\int_Z q(Z) \cdot \log \frac{P(Z|X)}{q(Z)} dZ}_{\text{KL}(q||p)} \\ &= J(q) + \text{KL}(q||p) \\ &\stackrel{\text{变分.}}{\downarrow} \quad \geq 0 \end{aligned}$$

由于 $P(X)$ 与 q 无关，因此等式左边： $\log P(X)$ 一定时， $J(q) + \text{KL}(q||p)$ 为固定值， q 无论如何取值，最大只能达到与 $\log P(X)$ 一样大。

我们想找一个 $q(Z)$ ，使 $q(Z)$ 近似等于 $P(Z|X)$ ，因此，我们希望找到一个 $q(Z)$ 使 $J(q)$ 能够达到最大。

$$\text{则 } \tilde{q}(Z) = \arg \max_{q(Z)} J(q) \Rightarrow \tilde{q}(Z) \approx P(Z|X)$$

假设 $q(Z)$ 可以划分为 M 个划分，这 M 个划分相互独立。

平均场理论结果

$$q(Z) = \prod_{i=1}^M q_i(z_i)$$

求解 q_i ，固定 $q_1, q_2, \dots, q_{j-1}, q_{j+1}, \dots, q_M$

$$J(q) = \underbrace{\int_Z q(Z) \log P(X, Z) dZ}_{①} - \underbrace{\int_Z q(Z) \log Q(Z) dZ}_{②}$$

$$\begin{aligned} ① &= \int_{Z_j} \prod_{i=1}^M q_i(z_i) \log P(X, Z) dZ_1 dZ_2 \dots dZ_M \\ &= \int_{Z_j} q_j(z_j) \left(\int_{Z_1 \dots Z_{j-1} \dots Z_{j+1} \dots Z_M} \prod_{i=1}^{j-1} q_i(z_i) \log P(X, Z) dZ_1 \dots dZ_{j-1} \dots dZ_{j+1} \dots dZ_M \right) dZ_j \\ &= \int_{Z_j} q_j(z_j) \left(\int_{Z_1 \dots Z_{j-1} \dots Z_{j+1} \dots Z_M} \log P(X, Z) \prod_{i=1}^{j-1} q_i(z_i) dZ_1 \dots dZ_{j-1} \dots dZ_{j+1} \dots dZ_M \right) dZ_j \\ &\quad \underbrace{\quad \quad \quad}_{E \prod_{i=1}^{j-1} q_i(z_i)} [\log P(X, Z)] \end{aligned}$$

$$= \int_{Z_j} q_j(z_j) \cdot E_{\prod_{i \neq j} q_i(z_i)} [\log P(x, z)] \cdot dz_j$$

$$\begin{aligned} \textcircled{2} &= \int_Z q(z) \log q(z) dz \\ &= \int_Z \prod_{i=1}^M q_i(z_i) \sum_{j=1}^M \log q_j(z_j) dz \\ &= \underbrace{\int_Z \prod_{i=1}^M q_i(z_i) [\log q_1(z_1) + \log q_2(z_2) + \dots + \log q_M(z_M)] dz}_{\text{写为 } M \text{ 个积的和}}. \end{aligned}$$

其中一项：

$$\begin{aligned} \int_Z \prod_{i=1}^M q_i(z_i) \log q_i(z_i) dz &= \int_Z q_1 q_2 \dots q_M \log q_1 dz \\ &= \int_{Z_1 Z_2 \dots Z_M} q_1 q_2 \dots q_M \log q_1 dz_1 dz_2 \dots dz_M \\ &= \int_{Z_1} q_1 \log q_1 dz_1 \cdot \int_{Z_2} q_2(z_2) dz_2 \dots \int_{Z_M} q_M(z_M) dz_M \\ &= \int_{Z_1} q_1 \log q_1 dz_1 \\ \Rightarrow &= \sum_{i=1}^M \int_{Z_i} q_i(z_i) \log q_i(z_i) dz_i \\ &= \int_{Z_j} q_j(z) \log q_j(z) dz_j + C \quad (\text{因为现在只关心 } q_j) \end{aligned}$$

$$\textcircled{1} \rightarrow : \int_{Z_j} q_j(z_j) \cdot \underbrace{E_{\prod_{i \neq j} q_i(z_i)} [\log P(x, z)]}_{\text{写为 log似然式, 设为 } \hat{P}(x, z_j)} \cdot dz_j$$

$$\begin{aligned} \text{则 } L(q) &= \textcircled{1} - \textcircled{2} = \int_{Z_j} q_j(z_j) \log \frac{\hat{P}(x, z_j)}{q_j(z_j)} dz_j \\ &= \boxed{-KL(q_j || \hat{P}(x, z_j))} \end{aligned}$$

$$\therefore q_j(z_j) = \hat{P}(x, z_j) \text{ 或 } \log q_j(z_j) = \log \hat{P}(x, z_j)$$

三、变分推断再回首

对上一节进行澄清：

X : 观测变量, $X = \{x^{(i)}\}_{i=1}^N$

Z : 隐变量, $Z = \{z^{(i)}\}_{i=1}^N$

上一节中: $\log P(X) = \underbrace{\text{ELBO}}_{\geq 0} + \underbrace{\text{KL}(q||p)}_{\geq 0} \geq \mathcal{L}(q)$

对于此项, 由于是似然, $\log P(X) = \log \prod_i P_0(x^{(i)}) = \sum_i \log P_0(x^{(i)})$, 对于 MLE 由于样本独立同分布, 因此 MLE 使 $\log P(X)$ 达到最大, 实际上是使每一项都达到最大。

改为: $\log P_0(x^{(i)}) = \text{ELBO} + \text{KL}(q||p) \geq \mathcal{L}(q)$

目标函数: $\hat{q} = \arg \min_q \text{KL}(q||p) = \arg \max_q \mathcal{L}(q)$

由上一节:

$$\log q_i(z_i) = E_{\prod_{j \neq i} q_j(z_j)} [\log P(x^{(i)}, z)] dz_j + \text{Const}$$

$$= \int_{z_2} \int_{z_3} \cdots \int_{z_{i-1}} \int_{z_{i+1}} \cdots \int_{z_m} q_1 \cdots q_{i-1} q_{i+1} \cdots q_m \cdot \log P(x^{(i)}, z) dz_2 \cdots dz_{i-1} dz_{i+1} \cdots dz_m$$

它实际上是一个迭代式,

$$\hat{q}_1(z_1) = \int_{q_2} \cdots \int_{q_m} q_2 \cdots q_m [\log P_0(x^{(1)}, z)] dz_2 \cdots dz_m$$

$$\hat{q}_2(z_2) = \int_{\hat{q}_1} \int_{q_3} \cdots \int_{q_m} \hat{q}_1 \cdots q_m [\log P_0(x^{(2)}, z)] dz_1 dz_3 \cdots dz_m$$

⋮

$$\hat{q}_m(z_m) = \int_{\hat{q}_1} \int_{\hat{q}_2} \cdots \int_{\hat{q}_{m-1}} \hat{q}_1 \cdots \hat{q}_{m-1} [\log P_0(x^{(m)}, z)] dz_1 dz_2 \cdots dz_{m-1}$$

坐标上升的思想。

Classical VI 存在的问题:

① 假设太强 (可分成 M 份推独立的变量)

② 有时候也会无法计算出来

四、随机梯度变分推断 (SGVI)

前面用坐标上升的方法求 $q(\mathbf{z})$, 很容易想到是否可以用梯度的方法求 $q(\mathbf{z})$

设 $q(\mathbf{z}) = q_{\phi}(\mathbf{z})$, ϕ 是 $q(\mathbf{z})$ 分布的参数.

$$\hat{\phi} = \arg \max_{\phi} L(\phi) \quad \text{应用梯度上升求 } L(\phi) \text{ 对 } \phi \text{ 的梯度}$$

于是求解 q 变为求解 ϕ .

$$\arg \max_{\phi} L(\phi) \rightarrow \arg \max_{\phi} L(\phi)$$

$$\hat{\phi} = \arg \max_{\phi} L(\phi)$$

$$\text{其中 } ELBO = E_{q_{\phi}(\mathbf{z})} [\log \frac{P_0(x^{(i)}, \mathbf{z})}{q_{\phi}(\mathbf{z})}] = E_{q_{\phi}(\mathbf{z})} [\log P_0(x^{(i)}, \mathbf{z}) - \log q_{\phi}(\mathbf{z})] = L(\phi)$$

求解:

$$\begin{aligned} \nabla_{\phi} L(\phi) &= \nabla_{\phi} E_{q_{\phi}} [\log P_0(x^{(i)}, \mathbf{z}) - \log q_{\phi}] \\ &= \nabla_{\phi} \int q_{\phi} [\log P_0(x^{(i)}, \mathbf{z}) - \log q_{\phi}] d\mathbf{z} \\ &= \underbrace{\int \nabla_{\phi} q_{\phi} [\log P_0(x^{(i)}, \mathbf{z}) - \log q_{\phi}] d\mathbf{z}}_{①} + \underbrace{\int q_{\phi} \nabla_{\phi} [\log P_0(x^{(i)}, \mathbf{z}) - \log q_{\phi}] d\mathbf{z}}_{②} \end{aligned}$$

$$\begin{aligned} ② &= \int q_{\phi} \nabla_{\phi} [\log P_0(x^{(i)}, \mathbf{z}) - \log q_{\phi}] d\mathbf{z} \\ &= - \int q_{\phi} \nabla_{\phi} \log q_{\phi} d\mathbf{z} \\ &= - \int q_{\phi} \cdot \frac{1}{q_{\phi}} \nabla q_{\phi} d\mathbf{z} \\ &= - \int \nabla_{\phi} q_{\phi} d\mathbf{z} \\ &= - \nabla_{\phi} \int q_{\phi} d\mathbf{z} \\ &= 0 \end{aligned}$$

$$\text{因此, } \nabla_{\phi} L(\phi) = ① = \int \nabla_{\phi} q_{\phi} [\log P_0(x^{(i)}, \mathbf{z}) - \log q_{\phi}] d\mathbf{z}$$

$$\begin{aligned} &= \int q_{\phi} \nabla_{\phi} \log q_{\phi} \cdot [\log P_0(x^{(i)}, \mathbf{z}) - \log q_{\phi}] d\mathbf{z} \\ &= E_{q_{\phi}} [\nabla_{\phi} \log q_{\phi} \cdot (\log P_0(x^{(i)}, \mathbf{z}) - \log q_{\phi})] \end{aligned}$$

这个期望可以通过MC采样来近似, 从而得到梯度, 然后利用梯度上升的方法来得到参数:

设采样样本为 $\mathbf{z}^{(l)} \sim q_{\phi}(\mathbf{z})$, $l = 1, 2, \dots, L$.

$$\text{则 } E_{q_{\phi}} [\nabla_{\phi} \log q_{\phi} \cdot (\log P_0(x^{(i)}, \mathbf{z}) - \log q_{\phi})] \approx \frac{1}{L} \sum_{l=1}^L \nabla_{\phi} \log q_{\phi}(\mathbf{z}^{(l)}) (\log P_0(x^{(i)}, \mathbf{z}^{(l)}) - \log q_{\phi}(\mathbf{z}^{(l)}))$$

但由于存在对数项，如果采样得到的 ϵ 很大， \log 的值会无穷大，才这么说。
 如果两个采样点差一点点， \log 的差距会非常大，会导致期望方程中分子的方差会非常大。
 就需要大量的样本才能得到较好的近似，L 如果过大，一般认为这是不可行的。
 因此采用重参数化技巧

重参数化技巧 (Reparameterization Trick)

由于 $\nabla_{\phi} \mathcal{L}(\phi) = \nabla_{\phi} E_{q_{\phi}} [\log P_{\theta}(x^{(i)}, z) - \log q_{\phi}]$ ，是对 q_{ϕ} 的分布求期望，
 如果我们可以对于一个确定的分布求期望，计算会变得很容易。

假设 $z = g_{\phi}(\epsilon, x^{(i)})$, $\epsilon \sim p(\epsilon)$

将 Z 的随机性转移到 \epsilon

$$z \sim q_{\phi}(z|x^{(i)})$$

$$\downarrow \quad \downarrow$$

$$\Rightarrow |q_{\phi}(z|x^{(i)}) dz| = |p(\epsilon) d\epsilon|$$

$$\epsilon \sim p(\epsilon)$$

$$\text{Q1} \quad \nabla_{\phi} \mathcal{L}(\phi) = \nabla_{\phi} E_{q_{\phi}} [\log P_{\theta}(x^{(i)}, z) - \log q_{\phi}]$$

$$= \nabla_{\phi} \int [\log P_{\theta}(x^{(i)}, z) - \log q_{\phi}] q_{\phi} dz$$

$$= \nabla_{\phi} \int [\log P_{\theta}(x^{(i)}, z) - \log q_{\phi}] \cdot p(\epsilon) d\epsilon$$

$$= \nabla_{\phi} E_{p(\epsilon)} [\log P_{\theta}(x^{(i)}, z) - \log q_{\phi}]$$

$$= E_{p(\epsilon)} [\nabla_{\phi} (\log P_{\theta}(x^{(i)}, z) - \log q_{\phi})]$$

$$= E_{p(\epsilon)} [\nabla_z (\log P_{\theta}(x^{(i)}, z) - \log q_{\phi}) \cdot \nabla_{\phi} z]$$

$$= E_{p(\epsilon)} [\nabla_z (\log P_{\theta}(x^{(i)}, z) - \log q_{\phi}) \cdot \nabla_{\phi} g_{\phi}(\epsilon, x^{(i)})]$$

链式法则 $\frac{\partial z}{\partial \phi} = \frac{\partial z}{\partial \epsilon} \cdot \frac{\partial \epsilon}{\partial \phi}$

进行 MC 采样，

$$\epsilon^{(t)} \sim p(\epsilon)$$

$$l = 1, 2, \dots, L$$

$$\nabla_{\phi} \mathcal{L}(\phi) \approx \frac{1}{L} \sum_{l=1}^L \nabla_z [\log P_{\theta}(x^{(i)}, z) - \log q_{\phi}(z|x^{(i)})] \cdot \nabla_{\phi} g_{\phi}(\epsilon^{(l)}, x^{(i)})$$

$\hookrightarrow g_{\phi}(\epsilon^{(l)}, x^{(i)})$

SGVI:

$$\phi^{(t+1)} \leftarrow \phi^{(t)} + \lambda^{(t)} \cdot \nabla_{\phi} \mathcal{L}(\phi)$$