

Variational Inference

频率派 \rightarrow 优化问题

① 模型: $f(w) = w^T \cdot x$, ② loss function

策略: $L(w) = \sum_{i=1}^N \|w^T x_i - y_i\|^2$
 $\hat{w} = \arg \min L(w)$ ↑ 求解

SVM (分类问题)

① $f(w) = \text{sign}(w^T x + b)$ ③ 算法: 解析解 $\frac{\partial L(w)}{\partial w} = 0 \Rightarrow w^* = (x^T x)^{-1} x^T y$
 ② loss function: $\min \frac{1}{2} w^T w$
 $s.t. y_i (w^T x_i + b) \geq 1, i=1, \dots, N$ 数值解: Gradient Descent / SGD

③ QP, 拉格朗日, 对偶.

EM: $\hat{\theta} = \arg \max_{\theta} \log p(x|\theta)$
 $\theta^{(t+1)} = \arg \max_{\theta} \int_z \log p(x, z|\theta) \cdot p(z|x, \theta^{(t)}) dz$

贝叶斯角度 - 积分问题

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

后验 \uparrow 似然 \uparrow 先验

积分 $= \int_{\theta} p(x|\theta) \cdot p(\theta) d\theta$

贝叶斯 inference: 求后验 \rightarrow { 精确推断, 近似推断 }

贝叶斯 决策 (预测) x 确定性近似 VI
随机性近似 MCMC / MH / Gibbs.

给定样本 \hat{x} , 求 $p(\hat{x}|x)$

$p(\hat{x}|x) = \int_{\theta} p(\hat{x}, \theta|x) \cdot d\theta = \int_{\theta} p(\hat{x}|\theta) \cdot \underbrace{p(\theta|x)}_{\text{后验}} d\theta = \text{可理解为关于后验求期望 (积分)}$

\rightarrow 引入 θ = $E_{\theta|x} [p(\hat{x}|\theta)]$

$\hat{x} \leftarrow x$
 θ



x : observed data

z : latent variable + parameter
[隐变量 + 参数]

(x, z) : complete data

也看作 z , 合并在 z 里面。

做一个 EM.

$$\log P(x) = \log P(x, z) - \log P(z|x)$$

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

对 $q(z)$ 求期望。

$$\text{左边} = \int_z \log P(x) \cdot q(z) dz = \log P(x) \cdot \int_z q(z) dz$$

$$\therefore \int_z q(z) dz = \log P(x)$$

$$\text{右边} = \int_z q(z) \cdot \log \frac{P(x, z)}{q(z)} dz - \int_z q(z) \log \frac{P(z|x)}{q(z)} dz$$

$$\begin{aligned} &\xrightarrow{\text{ELBO}} \mathcal{L}(q) + \text{KL}(q||P) \\ &\xrightarrow{q \rightarrow 0} \end{aligned}$$

\therefore 找一个 $q(z)$, 让其接近 $\log P(x)$. 那么 KL 就接近 0 $\Rightarrow q(z) \approx P(z|x)$

$$\therefore \tilde{q}(z) = \arg \max_{q(z)} \mathcal{L}(q) \Rightarrow \tilde{q}(z) \approx P(z|x)$$

假设 $q(z)$ 可以合成 m 个, 且其相互独立 \Rightarrow 平均场理论
mean theory

$$q(z) = \prod_{i=1}^m q_i(z_i)$$

$$\text{右边 } \mathcal{L}(q) = \int_z q(z) \log P(x, z) dz - \int_z q(z) \log q(z) dz$$

则 $q(z)$ 的 $\mathcal{L}(q)$

$$\text{拆成 } z_1, z_2, \dots, z_m$$

$$\text{①} = \int_z \prod_{i=1}^m q_i(z_i) \cdot \log P(x, z) dz_1 dz_2 \dots dz_m$$

拆解 $q_j(z_j)$, 其它固定

$$= \int_z q_j(z_j) \left(\prod_{i \neq j} q_i(z_i) \log P(x, z) dz_1 \dots dz_m \right) dz_j$$

$$\downarrow$$

$$\int_{z_1, \dots, z_m} P(x, z) \prod_{i \neq j} q_i(z_i) dz_1 \dots dz_m$$

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

全称为 $\log P(x, z_j)$

其中某一项 $\int_z \prod_{i=1}^m q_i(z_i) \cdot \log q_i(z_i) dz$

$$\begin{aligned} \text{②} &= \int_z q(z) \log q(z) dz = \int_z \prod_{i=1}^m q_i(z_i) \cdot \log \prod_{i=1}^m q_i(z_i) dz \\ &= \int_z \prod_{i=1}^m q_i(z_i) \cdot \sum_{i=1}^m \log q_i(z_i) dz \\ &= \int_z \prod_{i=1}^m q_i(z_i) [\log q_1(z_1) + \log q_2(z_2) + \dots + \log q_m(z_m)] dz \\ &= \text{有独立项} = \sum_{i=1}^m \int_{z_i} q_i(z_i) \cdot \log q_i(z_i) dz_i = \int_{z_1} q_1 \log q_1 dz_1 \end{aligned}$$

$$\text{把 } q_j(z_j) \text{ 之外看做已知 (固定)} \triangleq \int_{z_j} q_j(z_j) \log q_j(z_j) dz_j + C \rightarrow \text{常数}$$

$$\text{那么 } \text{①} - \text{②} = \int_{z_j} q_j(z_j) \cdot \log \frac{\hat{p}(x, z_j)}{q_j(z_j)} dz_j \xrightarrow{\text{常数略掉}}$$

$$= -\text{KL}(q_j || \hat{p}(x, z_j)) \leq 0 \quad \text{取等号时} \quad q_j(z_j) = \hat{p}(x, z_j)$$



变分推断 (VI)

x : observed variable $\rightarrow x = \{x_i\}_{i=1}^N = \{x^{(i)}\}_{i=1}^N$ 可视为已知, 或者是一个常数, $\Rightarrow L(\theta)$
 z : (latent variable) $\rightarrow z = \{z_i\}_{i=1}^N = \{z^{(i)}\}_{i=1}^N$ $\log p(x|\theta) = \underbrace{ELBO}_{L(\theta)} + \underbrace{KL(q(z)||p)}_{\geq 0} \geq L(\theta)$

(x, z) : complete data

θ : model parameter

$$ELBO = E_{q(z)} \left[\log \frac{p(x, z|\theta)}{q(z)} \right]$$

$$= E_{q(z)} [\log p(x, z|\theta)] + H[q(z)]$$

$$KL(q||p) = \int q(z) \cdot \log \frac{q(z)}{p(z|\theta)} \cdot dz$$

目标函数: $\hat{\theta} = \argmin_{\theta} KL(q||p) = \argmax_{\theta} L(\theta)$

假设已知 $\rightarrow \log \log p_0(x^{(i)}) = \underbrace{ELBO}_{L(\theta)} + \underbrace{KL(q||p)}_{\geq 0} \geq L(\theta)$

视频在此进行了变分推断的直观描述, 将本号入 $x^{(i)}$ 右端

$p_0(x^{(i)}, z) = p(x, z|\theta)$

VI:

平均场理论 \blacktriangle
 前提是假设 Assumption: $q(z) = \prod_{i=1}^M q_i(z_i)$

$$\log q_j(z_j) = E_{\pi} q_i(z_i) [\log p(x, z|\theta)] + c$$

$$= \int q_1 \dots \int q_{j-1} \int q_{j+1} \dots \int q_M q_1 \dots q_{j-1} q_{j+1} \dots q_M [\log p_0(x^{(i)}, z)]$$

$$dq_1, dq_2, \dots, dq_{j-1}, dq_{j+1}, \dots, dq_M$$

迭代优化过程 $\hat{q}_1(z_1) = \int q_2 \dots \int q_M q_2 \dots q_M [\log p_0(x^{(i)}, z)] dq_2 \dots dq_M$
 $\hat{q}_2(z_2) = \int q_1 \int q_3 \dots \int q_M q_3 \dots q_M [\log p_0(x^{(i)}, z)] dq_1 \dots dq_M$

$$\hat{q}_m(z_m) = \int q_1 \dots \int q_{m-1} \hat{q}_1 \dots \hat{q}_{m-1} [\log p_0(x^{(i)}, z)] dq_1 \dots dq_{m-1}$$

\blacktriangle 坐标下降法 coordinate Ascend.
 判定问题是: $L^{(t+1)} \leq L^{(t)}$ 时停止.

以上就是经典的基于平均场理论的变分推断问题.

存在问题 ① 一些模型不存在(不满足)平均场理论, 如神经网络

② 此处求积分并不是都可以求出 intractable 例如 $p(z|x)$ 后验, 对于它, 不可求



随机梯度变分推断 (SGVI)
 $\theta^{(t+1)} \leftarrow \theta^{(t)} + \alpha \Delta$ 梯度方向, 求上梯度

$q(z)$ 其实是关于 x 的联合分布 $q(z|x)$.

$q(z)$ 是一个分布, 假定其参数为 ϕ , 最好的 ϕ 求出来, 那么最好的 q 也求出来了.

$$q(z) = q_\phi(z), \text{ 求 } \phi. \text{ 那么 } ELBO \text{ 写成相减形式} = E_{q_\phi(z)} [\log p_0(x^{(i)}, z) - \log q_\phi(z)]$$

$$\therefore \hat{\phi} = \argmax_{\phi} L(\phi) = L(\phi)$$

梯度 $\nabla_{\phi} L(\phi) = \nabla_{\phi} E_{q_\phi} [\log p_0(x^{(i)}, z) - \log q_\phi]$ (这里省略了 z)

$$= \nabla_{\phi} \int q_\phi [\log p_0(x^{(i)}, z) - \log q_\phi] dz$$



扫描全能王 创建

白板(12) 随机梯度变分推断 (SGVI)

承接上张笔记, 交换 $\int \nabla \phi \cdot \nabla \psi = \int \nabla \phi \cdot \nabla \psi - \log \psi \cdot \nabla \phi + \nabla \phi \cdot \log \psi$

$$= \int \nabla \phi \cdot \nabla \psi \cdot (\log p_0(x^{(i)}, z) - \log \psi) dz + \int \nabla \phi \cdot \nabla \psi \cdot \underbrace{(\log p_0(x^{(i)}, z) - \log \psi)}_{\text{与 } \phi \text{ 无关}} dz$$

= ① + ②

$$② = - \int \nabla \phi \cdot \nabla \psi \log \psi dz = - \int \nabla \phi \cdot \frac{1}{\psi} \nabla \psi dz = - \int \nabla \phi \cdot \nabla \log \psi dz = - \nabla \phi \cdot \underbrace{\int \nabla \log \psi dz}_{=1} = 0$$

对数积分梯度

$$① = \int \nabla \phi \cdot \nabla \psi \log \psi (\log p_0(x^{(i)}, z) - \log \psi) dz$$

$\hookrightarrow \int \nabla \phi \cdot \nabla \psi \log \psi$ 把 $\log \psi$ 增加出来, ② 使用过.

$$= E_{\psi} [\nabla \phi \cdot \nabla \psi \log \psi (\log p_0(x^{(i)}, z) - \log \psi)] \quad (\text{梯度期望表示})$$

$\nabla \phi$ 梯度可以表示成 $E_{\psi} \nabla \phi$ 的期望. 采用 mcmc, 蒙特卡罗方式估计

$$\begin{matrix} \psi(z) \\ \rightarrow \\ z^{(i)} \sim \psi(z) \\ i=1, 2, \dots, L \end{matrix}$$

$$\nabla \phi(\phi) = \nabla \phi E_{\psi} [\log p_0(x^{(i)}, z) - \log \psi] = E_{\psi} [\nabla \phi \log \psi (\log p_0(x^{(i)}, z) - \log \psi)]$$

$$z^{(i)} \sim \psi(z) \quad (i=1, 2, \dots, L) \sim \text{mcmc 采样} \approx \frac{1}{L} \sum_{i=1}^L \nabla \phi \log \psi(z^{(i)}) (\log p_0(x^{(i)}, z^{(i)}) - \log \psi(z^{(i)}))$$

上面的 $\log \psi$ 项, 可能会 $-\infty$, 很大.

导致 high variance. 方差大, 需要采样多, 但是 L 是固定的.

[还有部分没看]



扫描全能王 创建