

最大熵模型

Max Entropy Model

对某一训练数据集, $T = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$

$$\tilde{p}(X=x, Y=y) = \frac{v(X=x, Y=y)}{N}$$

$$\tilde{p}(X=x) = \frac{v(X=x)}{N} \quad v \text{ 为频数}$$

feature function. $f(x, y) = \begin{cases} 1, & x \text{ 与 } y \text{ 满足某事实} \\ 0, & \text{else} \end{cases}$

$$E_{\tilde{p}}(f) = \sum_{x,y} \tilde{p}(x,y) f(x,y)$$

$$E_p(f) = \sum_{x,y} \underbrace{\tilde{p}(x)}_{\prod?} p(y|x) f(x,y)$$

of $E_{\tilde{p}}(f) = E_p(f) \Rightarrow$ 约束

(2) 模型能否获取训练数据中的信息?

$$C \equiv \{p \in \mathcal{P} \mid E_p(f_i) = E_{\tilde{p}}(f_i), i=1, \dots, n\}$$

Cond-Ent

$$H(p) = - \sum_{x,y} \tilde{p}(x) p(y|x) (\log p(y|x))$$

$$H(Y|X) = \sum_{x \in X} p(x) H(Y|X=x)$$

$$\begin{aligned}
 &= - \sum_{x \in X} p(x) \sum_{y \in Y} p(y|x) \log p(y|x) \\
 &= - \sum_{x,y} p(x) p(y|x) \log p(y|x)
 \end{aligned}$$

2. 最大熵模型的定义

Input: $T = \{(x_1, y_1), \dots, (x_n, y_n)\}$,
 n 个特征函数 $f_i(x, y)$, $i = 1, \dots, n$,

$$\max_{p \in \mathcal{C}} H(p) = - \sum_{x,y} \tilde{p}(x) p(y|x) \log p(y|x)$$

$$\text{s.t. } \mathbb{E}_p(f_i) = \mathbb{E}_{\tilde{p}}(f_i), \quad i = 1, 2, \dots, n$$

$$\sum_y p(y|x) = 1$$



$$\begin{aligned}
 \min_{p \in \mathcal{C}} -H(p) &= \sum_{x,y} \tilde{p}(x) p(y|x) \log p(y|x) \\
 \text{s.t. } \mathbb{E}_p(f_i) - \mathbb{E}_{\tilde{p}}(f_i) &= 0, \quad i = 1, 2, \dots, n
 \end{aligned}$$

$$\sum_y p(y|x) = 1$$

为规范化, 引入拉格朗日乘因子 $\omega_0, \omega_1, \dots, \omega_n$

$$\begin{aligned}
 L(p, \omega) &= -H(p) + \omega_0 (1 - \sum_y p(y|x)) \\
 &\quad + \sum_{i=1}^n \omega_i (\mathbb{E}_{\tilde{p}}(f_i) - \mathbb{E}_p(f_i))
 \end{aligned}$$

$$= \sum_{x,y} \tilde{p}(x) p(y|x) \log p(y|x) + \omega_0 (1 - \sum_y p(y|x))$$

$$L(p, \omega) = \sum_{i=1}^n \omega_i (E_{\tilde{p}}(f_i) - E_p(f_i))$$

$$p_{\omega} = \arg \min_{p \in \mathcal{Q}} L(p, \omega) = p_{\omega}(y|x) \Rightarrow \text{最优解}$$

$$\begin{aligned} \frac{\partial L(p, \omega)}{\partial p(y|x)} &= \sum_{x,y} \tilde{p}(x) (1 + \log p(y|x)) - \sum_y \omega_0 \\ &\quad - \sum_{x,y} (\tilde{p}(x) \sum_{i=1}^n \omega_i f_i(x,y)) \end{aligned}$$

$$= \sum_{x,y} \tilde{p}(x) \left(1 + \log p(y|x) - \sum_{i=1}^n \omega_i f_i(x,y) \right)$$

$$\because \sum_y \omega_0 = \sum_y \omega_0 \left(\sum_x \tilde{p}(x) \right)$$

$$\therefore \rightarrow = \sum_{x,y} \tilde{p}(x) \left(\log p(y|x) + 1 - \omega_0 - \sum_{i=1}^n \omega_i f_i(x,y) \right)$$

$$\frac{\partial L(p, \omega)}{\partial p(y|x)} = 0 \Rightarrow p(y|x) = \frac{\exp \left(\sum_{i=1}^n \omega_i f_i(x,y) \right)}{\exp(1 - \omega_0)}$$

$$\because \sum_y p(y|x) = 1$$

$$\Rightarrow \sum_y \frac{\exp \left(\sum_{i=1}^n \omega_i f_i(x,y) \right)}{\exp(1 - \omega_0)} = 1$$

$$\Rightarrow \exp(1 - \omega_0) = \sum_y \exp \left(\sum_{i=1}^n \omega_i f_i(x,y) \right)$$

$$\Rightarrow \tilde{p}_{\omega}(y|x) = \frac{1}{Z_{\omega}(x)} \exp \left(\sum_{i=1}^n \omega_i f_i(x,y) \right)$$

$$\text{with } Z_{\omega}(x) = \sum_y \exp \left(\sum_{i=1}^n \omega_i f_i(x,y) \right)$$

↓ 即为最大熵模型

⇒ 等价于极大似然估计

3. 最优化算法

常用方法: 改进的迭代法, SGD, 牛顿法或拟牛顿法

(1) 牛顿法 = $\min_{x \in \mathbb{R}^n} f(x) \Rightarrow x^*$ (极小值解)

迭代过程中, 第 k 次迭代值为 $x^{(k)}$, 则

$$f(x) = f(x^{(k)}) + g_k^T (x - x^{(k)}) + \frac{1}{2} (x - x^{(k)})^T H(x^{(k)}) (x - x^{(k)})$$

Hesse matrix $H(x^{(k)}) = \left[\frac{\partial^2 f}{\partial x_i \partial x_j} \right]_{n \times n}, x \in \mathbb{R}^n$

$$\text{if } f(x) \text{ 取得极小值} \Rightarrow \nabla f(x^{(k+1)}) = 0$$

$$\Rightarrow g_k^T + H(x^{(k)}) (x^{(k+1)} - x^{(k)}) = 0$$

$$\Rightarrow x^{(k+1)} = x^{(k)} - H_k^{-1} g_k$$

$$\text{OR } x^{(k+1)} = x^{(k)} + p_k, H_k p_k = -g_k$$

且该算法为:

Input: $f(x), x \in \mathbb{R}^n, g(x) = \nabla f(x), H(x), \varepsilon$

Output: x^*

① 取初始点 $x^{(0)}, k=0$

② 计算 $g_k = g(x^{(k)})$

③ if $\|g_k\| < \varepsilon$, then $x^* = x^{(k)}$, return

④ 计算 $H_k = H(x^{(k)})$, $p_k = -H_k^{-1}g_k$

⑤ $x^{(k+1)} = x^{(k)} + p_k, k = k+1$, 转 ②

(2) 牛顿法

① $n \times n$ $G_k = G(x^{(k)})$ 非奇异 $H_k^{-1} = (H^{-1}(x^{(k)}))$

② H_k 是下对称阵

$$\begin{aligned} \nabla f(x^{(k+1)}) &= g_k + H_k(x^{(k+1)} - x^{(k)}) \\ \Rightarrow \underbrace{g_{k+1} - g_k}_{y_k} &= H_k \underbrace{(x^{(k+1)} - x^{(k)})}_{\delta_k} \end{aligned}$$

$$\Rightarrow H_k^{-1}y_k = \delta_k$$

$$\text{③ 牛顿法 } G_k y_k = \delta_k$$

$$\boxed{\text{BFGS 算法}}: B_{k+1} \delta_k = y_k$$

$$13 \quad B_{k+1} = B_k + P_k + Q_k$$

$$\Rightarrow B_k \delta_k = B_k \delta_k + P_k \delta_k + Q_k \delta_k$$

$$\begin{cases} P_k \delta_k = y_k \\ Q_k \delta_k = -B_k \delta_k \end{cases}$$

$$\Rightarrow B_{k+1} = B_k + \frac{y_k y_k^T}{y_k^T \delta_k} - \frac{B_k \delta_k \delta_k^T B_k}{\delta_k^T B_k \delta_k} \quad (*)$$

Input: $f(x)$, $g(x) = \nabla f(x)$, ε ;

Output: 极小值 x^* ;

Identity

① 选取 $x^{(0)}$, B_0 为 positive-definite matrix
 $k=0$,

② 计算 $g_k = g(x^{(k)})$, if $\|g_k\| < \varepsilon$,
then $x^* = x^{(k)}$, return

③ 由 $B_k P_k = -g_k \Rightarrow P_k$

④ 一维搜索: 求 λ_k , 使得

$$f(x^{(k)} + \lambda_k P_k) = \min_{\lambda > 0} f(x^{(k)} + \lambda P_k)$$

$$⑤ \quad x^{(k+1)} = x^{(k)} + \lambda_k P_k$$

⑥ $g_{k+1} = g(x^{(k+1)})$, if $\|g_{k+1}\| < \varepsilon$,
then $x^* = x^{(k+1)}$

else 按 (*) 求 B_{k+1}

⑦ $k=k+1$, 转至③

对于最大熵模型而言:

$$\min_{w \in \mathbb{R}^n} f(w) = \sum_x \tilde{p}(x) \log \sum_y \exp \left(\sum_{i=1}^n w_i f_i(x, y) \right) \\ - \sum_{x, y} \tilde{p}(x, y) \sum_{i=1}^n w_i f_i(x, y)$$

$$g(w) = \left(\frac{\partial f(w)}{\partial w_1}, \frac{\partial f(w)}{\partial w_2}, \dots, \frac{\partial f(w)}{\partial w_n} \right)$$

$$\text{with } \frac{\partial f(w)}{\partial w_i} = \sum_{x, y} \tilde{p}(x) P_w(y|x) f_i(x, y) - E_p(f_i)$$

$$P_w(y|x) = \frac{\exp \left(\sum_{i=1}^n w_i f_i(x, y) \right)}{\sum_y \exp \left(\sum_{i=1}^n w_i f_i(x, y) \right)}$$

上次修改: 23:06