

逻辑斯蒂回归与最大熵模型

W.J.Z

1 逻辑斯蒂回归

1.1 逻辑斯蒂分布

设 X 为连续随机变量, X 服从逻辑斯蒂分布是指 X 具有以下的分布函数和密度函数:

$$F(x) = P(X \leq x) = \frac{1}{1 + e^{\frac{-(x-\mu)}{\gamma}}} \quad (1)$$

$$f(x) = F'(x) = \frac{e^{\frac{-(x-\mu)}{\gamma}}}{\gamma \left(1 + e^{\frac{-(x-\mu)}{\gamma}}\right)^2} \quad (2)$$

μ 为位置参数, $\gamma > 0$ 为形状参数。

1.2 逻辑斯蒂回归模型

设 Y 的取值为 0 或者 1, $w = (w^1, w^2, \dots, w^n, b)^T, x = (x^1, x^2, \dots, x^n, 1)^T$, 则二项逻辑斯蒂回归模型为如下的条件概率分布:

$$P(Y = 1|x) = \frac{\exp(w \cdot x)}{1 + \exp(w \cdot x)} \quad (3)$$

$$P(Y = 0|x) = \frac{1}{1 + \exp(w \cdot x)} \quad (4)$$

一件事情发生的几率是指改事件发生的概率与该事件不发生概率的比值, 对于逻辑斯蒂回归而言, 改事件的对数几率为:

$$\log \frac{P(Y = 1|x)}{1 - P(Y = 1|x)} = w \cdot x \quad (5)$$

输出 $Y = 1$ 的对数几率是输入 x 的线性函数表示的模型，线性函数的值越接近正无穷，概率值就越接近 1；线性函数的值越接近负无穷，概率值就越接近 0。

假设离散型随机变量 Y 的取值集合为 $1, 2, 3, \dots, K$ ，多项逻辑斯蒂回归模型是：

$$P(Y = k|x) = \frac{\exp(w_k \cdot x)}{1 + \sum_{k=1}^{K-1} \exp(w_k \cdot x)} \quad (6)$$

$$P(Y = K|x) = \frac{1}{1 + \sum_{k=1}^{K-1} \exp(w_k \cdot x)} \quad (7)$$

1.3 模型参数估计

设 $y_i \in \{0, 1\}$ 应用极大似然估计法评估模型参数，设：

$$P(Y = 1|x) = \pi(x), P(Y = 0|x) = 1 - \pi(x), \quad (8)$$

似然函数为：

$$\prod_{i=1}^N [\pi(x_i)]^{y_i} [1 - \pi(x_i)]^{1-y_i} \quad (9)$$

对数似然函数为：

$$L(w) = \sum_{i=1}^N [y_i \log \pi(x_i) + (1 - y_i) \log (1 - \pi(x_i))] \quad (10)$$

$$= \sum_{i=1}^N \left[y_i \log \frac{\pi(x_i)}{1 - \pi(x_i)} + \log (1 - \pi(x_i)) \right] \quad (11)$$

$$= \sum_{i=1}^N [y_i (w \cdot x_i) - \log (1 + \exp(w \cdot x_i))] \quad (12)$$

使用梯度下降法求极大值。

2 最大熵模型

2.1 最大熵原理

最大熵原理认为要选择的概率模型首先必须满足已有的事实，即约束条件，在没有更多信息的情况下，那些不确定的部分都是“等可能的”。

2.2 最大熵模型

设联合分布的经验分布和边缘分布的经验分布为：

$$P(X = x, Y = y) = \frac{v(X = x, Y = y)}{N} \quad (13)$$

$$P(X = x) = \frac{v(X = x)}{N} \quad (14)$$

$v(X = x, Y = y)$ 表示训练数据中样本 (x, y) 出现的频数， $v(X = x)$ 表示训练数据中输入 x 的频数， N 表示训练样本容量。

用特征函数表示输入 x 与输出 y 之间的某一个事实，其定义为：

$$f(x, y) = \begin{cases} 1, & x \text{ 与 } y \text{ 满足某一事实} \\ 0, & \text{其他} \end{cases} \quad (15)$$

特征函数 $f(x, y)$ 关于联合经验分布的期望值为：

$$E_{p(x, y)} = \sum_{x, y} P(x, y) f(x, y) \quad (16)$$

特征函数关于模型 $P(Y|X)$ 与边缘经验分布的期望值为：

$$E_{p(y|x)} = \sum_{x, y} P(x) P(y|x) f(x, y) \quad (17)$$

假设两个期望值相等，则：

$$E_{p(x, y)} = E_{p(y|x)} \quad (18)$$

公式 18 将作为模型学习的约束条件。定义在条件概率分布的条件熵为：

$$H(P) = - \sum_{x, y} P(x) P(y|x) \log P(y|x) \quad (19)$$

公式 19 为最大熵模型，式中的对数为自然对数。

2.3 最大熵模型的学习

最大熵模型的学习等价于约束最优化问题。

$$\max H(P) = - \sum_{x, y} P(x) P(y|x) \log P(y|x)$$

$$E_{p(x,y)} = E_{p(y|x)}$$

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

按照最优化的习惯，将求最大值问题改写为等价的求最小值问题。

$$\min -H(P) = \sum_{x,y} P(x) P(y|x) \log P(y|x)$$

$$E_{p(x,y)} - E_{p(y|x)} = 0$$

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

将约束最优化的原始问题转换为无约束最优化的对偶问题，引入拉格朗日乘子，得：

$$L(P, w) = -H(P) + w_0 \left(1 - \sum_y P(y|x) \right) + \sum_{i=1}^n w_i (E_{p(x,y)} - E_{p(y|x)})$$

$$= \sum_{x,y} P(x) P(y|x) \log P(y|x) + w_0 \left(1 - \sum_y P(y|x) \right) + \sum_{i=1}^n w_i (E_{p(x,y)} - E_{p(y|x)})$$

对 $L(P, w)$ 对 $P(y|x)$ 求偏导得：

$$P(y|x) = \frac{\exp(\sum_{i=1}^n w_i f_i(x, y))}{\exp(1 - w_0)} \quad (20)$$

将 $P(y|x)$ 带入 $L(P, w)$ 对 w 求偏导，得出 w_i 的值进而得出 $P(y|x)$ 的值。