

## logistic 回归

判别模型, 拟合的是  $P(Y|x)$ ,  $x$  没有随机性.

## 二项 logistic 回归

利用一个线性模型去拟合 logit 变换

输出  $Y \in \{0, 1\}$

$$P(Y=1|x) = \pi(x)$$

$$\text{logit } \pi(x) = \log \frac{P(Y=1|x)}{P(Y=0|x)} = \log \frac{\pi(x)}{1-\pi(x)} = w \cdot x$$

$$\pi(x) = \frac{\exp(w \cdot x)}{1 + \exp(w \cdot x)}$$

过拟合: 当  $w \cdot x$  趋于  $+\infty$  时,  $\pi(x)$  趋近于 1, 导致模型“过于自信”

## 多项 logistic 回归

同样是利用一个线性模型去拟合 logit 变换

$Y \in \{1, 2, 3, \dots, K\}$

$$P(Y=k|x) = \pi_k(x)$$

$$\text{logit } \pi_k(x) = \log \frac{P(Y=k|x)}{P(Y=1|x)} = w_k \cdot x$$

$$\begin{aligned} \sum_{k=1}^K P(Y=k|x) &= \sum_{k=1}^{K-1} e^{w_k \cdot x} P(Y=1|x) + P(Y=K|x) \\ &= P(Y=1|x) \left( \sum_{k=1}^{K-1} e^{w_k \cdot x} + 1 \right) \end{aligned}$$

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

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

## 学习算法

极大似然估计  $\Rightarrow$  估计  $W_k$  (因为  $W$  是参数)

共有  $N$  个样本, 以二分类为例  $P(Y=1|X)=z(X)$

$$L(W) = \log \prod_{i=1}^N P(Y=y_i|X=x_i)$$

$$= \log \prod_{i=1}^N [z(x_i)]^{y_i} [1-z(x_i)]^{1-y_i}$$

$$= \sum_{i=1}^N y_i \log z(x_i) + (1-y_i) \log (1-z(x_i))$$

$$= \sum_{i=1}^N y_i \log \frac{z(x_i)}{1-z(x_i)} + \log (1-z(x_i))$$

$$= \sum_{i=1}^N y_i (W \cdot x_i) - \log (1 + e^{W \cdot x_i}) \quad \leftarrow \text{梯度下降法求解该函数的极值}$$

## 最大熵模型

最大熵原理: 在满足约束条件的模型集合中选择熵最大的模型

有输入 条件熵  
无输入 熵

最大熵模型:

给定训练集  $T = \{(x_1, y_1), \dots, (x_n, y_n)\}$ , 以及特征函数  $f_i(x, y), i=1, 2, \dots, n$ .

学习最大熵模型等价于求解如下约束最优化问题:

$$\max_{P \in C} H(P) = - \sum_{x,y} \tilde{p}(x,y) P(y|x) \log P(y|x) \quad \text{条件熵}$$

$\leftarrow$  属于  $P(x,y)$   $\tilde{p}(x,y)$  为  $x$  的经验分布

$$\text{s.t.} \quad E_P(f_i) = E_{\tilde{P}}(f_i)$$

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

$$E_P(f_i) = \sum_{x,y} \tilde{p}(x,y) P(y|x) f_i(x,y) \quad E_{\tilde{P}}(f_i) = \sum_{x,y} \tilde{p}(x,y) f_i(x,y)$$

$\leftarrow x$  与  $y$  的联合经验分布

上述问题可通过拉格朗日对偶性转化

有约束的最优化问题利用拉格朗日函数进行转换

$$\min L(P, W) = -H(P) + W_0 (1 - \sum_y P(y|x)) + \sum_{i=1}^n W_i (E_P(f_i) - E_{\tilde{P}}(f_i))$$

$\leftarrow$  极小极大问题

$$\text{等价于 } \min_{P \in C} \max_W -H(P) + W_0 (1 - \sum_y P(y|x)) + \sum_{i=1}^n W_i (E_P(f_i) - E_{\tilde{P}}(f_i))$$

对偶问题

等价于  $\max_w \min_{p \in \mathcal{C}} -l(p) + w_0(1 - \sum_i p(y|x)) + \sum_{i=1}^n w_i (E_p[f_i] - E_{\tilde{p}}[f_i])$

$$\varphi(w) = \min_{p \in \mathcal{C}} -l(p) + w_0(1 - \sum_i p(y|x)) + \sum_{i=1}^n w_i (E_p[f_i] - E_{\tilde{p}}[f_i])$$

$$= \min_{p \in \mathcal{C}} \sum_{x,y} \tilde{p}(x) \log p(y|x) p(y|x) + w_0(1 - \sum_i p(y|x)) + \sum_{i=1}^n w_i (\sum_{x,y} \tilde{p}(x,y) f_i(x,y) - \sum_{x,y} \tilde{p}(x) p(y|x) f_i(x,y))$$

$$= \min_{p \in \mathcal{C}} L(p, w)$$

$$\frac{\partial L(p, w)}{\partial p(y|x)} = \sum_{x,y} \tilde{p}(x) (1 + \log p(y|x)) - \sum_i w_i - \sum_{i=1}^n w_i \sum_{x,y} \tilde{p}(x) f_i(x,y)$$

$$= \sum_{x,y} \tilde{p}(x) (1 + \log p(y|x)) - \sum_{x,y} \tilde{p}(x) w_0 - \sum_{x,y} \tilde{p}(x) \sum_{i=1}^n w_i f_i(x,y)$$

$$= \sum_{x,y} \tilde{p}(x) (1 + \log p(y|x)) - w_0 - \sum_{i=1}^n w_i f_i(x,y)$$

$$\downarrow = 0$$

$$1 + \log p(y|x) - w_0 - \sum_{i=1}^n w_i f_i(x,y) = 0$$

$$\log p(y|x) = w_0 + \sum_{i=1}^n w_i f_i(x,y) - 1$$

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

$$\sum_i p(y|x) = 1 \Rightarrow \sum_i \exp(\sum_{i=1}^n w_i f_i(x,y)) = \exp(1 - w_0)$$

$$p_w(y|x) = \frac{\exp(\sum_{i=1}^n w_i f_i(x,y))}{\sum_i \exp(\sum_{i=1}^n w_i f_i(x,y))} = \frac{\exp(\sum_{i=1}^n w_i f_i(x,y))}{Z_w(x)}$$

上述模型就是“要求的最大熵模型”，可以看到还有要优化的参数  $w_i$ 。

这些参数就是要通过最优化算法进一步求解的

1) 梯度下降, 迭代尺度法, 拟牛顿法

首先,我们需要定义这个最优化问题,两种方法:极大似然估计或者是上述的对偶问题

① 极大似然估计

$$L(P) = \log \prod_{x,y} P(y|x)^{\tilde{p}(x,y)}$$

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

$$= \sum_{x,y} \tilde{p}(x,y) \left( \log \exp \sum_{i=1}^n w_i f_i - \log Z_{w(x)} \right)$$

$$= \sum_{x,y} \tilde{p}(x,y) \sum_{i=1}^n w_i f_i - \sum_{x,y} \tilde{p}(x,y) \log Z_{w(x)}$$

$$= \sum_{x,y} \tilde{p}(x,y) \sum_{i=1}^n w_i f_i - \sum_x \log Z_{w(x)} \sum_y \tilde{p}(x,y)$$

$$L(w) = \sum_{x,y} \tilde{p}(x,y) \sum_{i=1}^n w_i f_i - \sum_x \log Z_{w(x)} \tilde{p}(x)$$

$$W_i^* = \max_{W_i} \sum_{x,y} \tilde{p}(x,y) \sum_{i=1}^n w_i f_i - \sum_x \log Z_{w(x)} \tilde{p}(x)$$

迭代尺度法推导

$$w \rightarrow w + \delta$$

$$L(w + \delta) - L(w) = \sum_{x,y} \tilde{p}(x,y) \left[ \sum_{i=1}^n (w_i + \delta_i) f_i - w_i f_i \right]$$

$$- \sum_x \tilde{p}(x) \left[ \log Z_{w+\delta}(x) - \log Z_w(x) \right]$$

$$= \sum_{x,y} \tilde{p}(x,y) \sum_{i=1}^n \delta_i f_i - \sum_x \tilde{p}(x) \log \frac{Z_{w+\delta}(x)}{Z_w(x)}$$

$$- \log 2, 1 - 2$$

$$- \sum_x \tilde{p}(x) \log \frac{Z_{w+\delta}(x)}{Z_w(x)} \geq \sum_x \tilde{p}(x) \left( 1 - \frac{Z_{w+\delta}(x)}{Z_w(x)} \right)$$

$$\geq \sum_x \tilde{p}(x) - \sum_x \tilde{p}(x) \frac{Z_{w+\delta}(x)}{Z_w(x)}$$

$$\geq 1 - \sum_x \tilde{p}(x) \frac{Z_{w+\delta}(x)}{Z_w(x)}$$

$$\frac{Z_{wt+\delta}(x)}{Z_{wt}(x)} = \frac{\sum_i \exp(\sum_{j=1}^n (wt+\delta)_j f_{ij})}{Z_{wt}(x)}$$

$$= \frac{\sum_i \exp(\sum_{j=1}^n w_j f_{ij}) \cdot \exp(\sum_{j=1}^n \delta_j f_{ij})}{Z_{wt}(x)}$$

$$P_{wt}(x) \leftarrow \sum_i P_{wt}(y|x) \exp(\sum_{j=1}^n \delta_j f_{ij})$$

$$L(wt+\delta) - L(w) \geq \sum_{x,y} \tilde{p}(x,y) \sum_{j=1}^n \delta_j f_{ij} + 1 - \sum_x \tilde{p}(x) \sum_y P_{wt}(y|x) \exp(\sum_{j=1}^n \delta_j f_{ij})$$

Jeson 不等式:  $\varphi$  为凸函数,  $\sum_i \alpha_i = 1$ .

则有  $\varphi(\sum_i \alpha_i x_i) \leq \sum_i \alpha_i \varphi(x_i)$

$$\exp(\sum_{j=1}^n \delta_j f_{ij}) = \exp(\sum_{j=1}^n \frac{f_{ij}}{f^{\#}(x,y)} f^{\#}(x,y) \delta_j) \quad (f^{\#}(x,y) = \sum_{j=1}^n f_{ij})$$

$$\leq \sum_{j=1}^n \frac{f_{ij}}{f^{\#}(x,y)} e^{f^{\#}(x,y) \delta_j}$$

$$L(wt+\delta) - L(w) \geq \sum_{x,y} \tilde{p}(x,y) \sum_{j=1}^n \delta_j f_{ij} + 1 - \sum_x \tilde{p}(x) \sum_y P_{wt}(y|x) \sum_{j=1}^n \frac{f_{ij}}{f^{\#}(x,y)} e^{f^{\#}(x,y) \delta_j}$$

下界最大化时对  $\delta_j$  求偏导

$$\frac{\partial (L(wt+\delta) - L(w))}{\partial \delta_j} = \sum_{x,y} \tilde{p}(x,y) f_{ij} - \sum_x \tilde{p}(x) \sum_y P_{wt}(y|x) f_{ij} e^{f^{\#}(x,y) \delta_j} = 0$$