

# Logistic Regression

用于分类，主要是得到预测结果的概率值：

$$f(x) = P(y|x) \in [0, 1]$$

给定数据集  $D = \left\{ (\vec{x}_1, y_1), (\vec{x}_2, y_2), \dots, (\vec{x}_n, y_n) \right\}, y_i \in \{-1, +1\}$

预测函数：

$$s(z) = \frac{1}{1 + \exp(-z)}$$

$$z = \sum w_i x_i = \vec{w}^T \vec{x}$$

$$h(x) = \frac{1}{1 + \exp(-\vec{w}^T \vec{x})}$$

**Error function** 通过极大似然估计得到：

目标函数  $f(x)$  下，产生  $D$  的概率：

$$P_f(D) = P(\vec{x}_1) f(\vec{x}_1) P(\vec{x}_2) f(\vec{x}_2) \dots P(\vec{x}_n) f(\vec{x}_n)$$

预测函数  $h(x)$  下，产生  $D$  的概率：

$$P_h(D) = P(\vec{x}_1) h(\vec{x}_1) P(\vec{x}_2) h(\vec{x}_2) \dots P(\vec{x}_n) h(\vec{x}_n)$$

若  $h$  与  $f$  很接近，则  $P_h(D)$  与  $P_f(D)$  很接近，称  $P_h(D)$  为似然。

$f$  既然产生了  $D$ ，可以想成“ $f$  产生  $D$  的概率极大”（所以  $f$  产生了  $D$ ）；

因此，最佳的  $h$ ，与  $f$  最接近的  $h$ ，其产生  $D$  的概率也该是极大的。

所以，要求最大似然  $\max_h \text{likelihood}(h)$ ，从而得到最佳的预测函数  $\arg \max_h P_h(D)$ 。

$$P_h(D) = P(\vec{x}_1) P(y_1 | \vec{x}_1) P(\vec{x}_2) P(y_2 | \vec{x}_2) \dots P(\vec{x}_n) P(y_n | \vec{x}_n)$$

$$\because P(y_i | \vec{x}_i) = \begin{cases} h(\vec{x}_i), & \text{for } y_i = +1 \\ 1 - h(\vec{x}_i), & \text{for } y_i = -1 \end{cases}$$

又  $\because$  logistic 函数有性质：  $1 - h(x) = h(-x)$

$$\therefore P(y_i | \vec{x}_i) = \begin{cases} h(x), & \text{for } y_i = +1 \\ h(-x), & \text{for } y_i = -1 \end{cases} = h(y_i x_i)$$

$$\therefore P_h(D) = P(\vec{x}_1)h(y_1 \vec{x}_1) P(\vec{x}_2)h(y_2 \vec{x}_2) \dots P(\vec{x}_n)h(y_n \vec{x}_n)$$

$$\therefore P(\vec{x}_i) \text{ 固定, 所以上式 } \propto \prod h(y_i \vec{x}_i)$$

$$\therefore \max_h P_h(D) \propto \prod h(y_i \vec{x}_i)$$

$$\therefore \text{每个 } h \text{ 对应一个参数 } \vec{w}, \therefore \max_{\vec{w}} \text{likelihood}(\vec{w}) \propto \prod s(y_i \vec{w}^T \vec{x}_i)$$

$$\therefore \text{连乘不容易求解, 所以取对数变成连加: } \max_{\vec{w}} \ln \prod s(y_i \vec{w}^T \vec{x}_i) = \max_h \sum \ln s(y_i \vec{w}^T \vec{x}_i)$$

$\therefore$  似然最大, 即误差最小, 目的也是求误差函数, 所以  $\max$  变为  $\min$  加一个负号:

$$\min_{\vec{w}} \frac{1}{N} \sum -\ln s(y_i \vec{w}^T \vec{x}_i)$$

$$\therefore s(z) = \frac{1}{1 + \exp(-z)} \text{ 代入上式得到: } \min_{\vec{w}} \frac{1}{N} \sum \ln(1 + \exp(-y_i \vec{w}^T \vec{x}_i))$$

其中, 定义交叉熵错误 *cross entropy error*:  $\text{err}(\vec{w}, \vec{x}_i, y_i) = \ln(1 + \exp(-y_i \vec{w}^T \vec{x}_i))$

$$\text{综上, error function: } \frac{1}{N} \sum \ln(1 + \exp(-y_i \vec{w}^T \vec{x}_i)) = \frac{1}{N} \sum \text{err}(\vec{w}, \vec{x}_i, y_i)$$

**如何解  $\min$  error function 得到最佳  $\vec{w}$ : 梯度及梯度下降**

$$E(\vec{w}) = \frac{1}{N} \sum \ln(1 + \exp(-y_i \vec{w}^T \vec{x}_i))$$

$$\begin{aligned} \nabla E(\vec{w}) &= \frac{\partial E(\vec{w})}{\partial w_i} = \frac{1}{N} \sum \left( \frac{1}{\square} \right) (\exp(\Delta)) (-y_i \vec{x}_i) \\ &= \frac{1}{N} \sum \left( \frac{\exp(\Delta)}{1 + \exp(\Delta)} \right) (-y_i \vec{x}_i) \\ &= \frac{1}{N} \sum s(\Delta) (-y_i \vec{x}_i) \\ &= \frac{1}{N} \sum s(-y_i \vec{w}^T \vec{x}_i) (-y_i \vec{x}_i) \end{aligned}$$

因为  $E(\vec{w})$  是连续, 可微的凸函数, 所以令  $\nabla E(\vec{w}) = 0$ , 可解得最佳  $\vec{w}$ 。

数据线性可分，即将数据正正好分离的线存在，才有求这条线的可能，才能用上面的方法令 $\nabla = 0$ 来得到解析解；若非线性可分，这条线不存在，则 $\nabla = 0$ 无解。现实任务中，数据通常非线性可分。

所以，另外有迭代优化方法来得到最佳 $\vec{w}$ :

1. 设置权值向量 $\vec{w}$ 初始值为 $\vec{w}_0$ ，设迭代次数为 $t$ ， $t = 0, 1, 2, \dots$ ;

2. 计算梯度 $\nabla E(\vec{w}_t) = \frac{1}{N} \sum s(-y_i \vec{w}_t^T \vec{x}_i) (-y_i \vec{x}_i)$ ;

3. 更新 $\vec{w}$ :  $\vec{w}_{t+1} = \vec{w}_t - \eta \nabla E(\vec{w}_t)$

4. 直到 $\nabla E(\vec{w}_t) \approx 0$  或迭代次数足够多