# 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 $D = \left\{ \left( \vec{x_1}, y_1 \right), \left( \vec{x_2}, y_2 \right), ..., \left( \vec{x_n}, y_n \right) \right\}$，$y_i \in \{-1, +1\}$

**预测函数：**

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

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

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

**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：
$$P_f(D) = P\left( \vec{x_1} \right) f\left( \vec{x_1} \right) P\left( \vec{x_2} \right) f\left( \vec{x_2} \right) ... P\left( \vec{x_n} \right) f\left( \vec{x_n} \right)$$

预测函数 $h(x)$ 下，产生 $D$ 的概率：
$$P_h(D) = P\left( \vec{x_1} \right) h\left( \vec{x_1} \right) P\left( \vec{x_2} \right) h\left( \vec{x_2} \right) ... P\left( \vec{x_n} \right) h\left( \vec{x_n} \right)$$

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 $max\ likelihood(h)$，从而得到最佳的预测函数 $\underset{h}{arg\ max}\ P_h(D).$

$$P_h(D) = P\left( \vec{x_1} \right) P\left( y_1|\vec{x_1} \right) P\left( \vec{x_2} \right) P\left( y_2|\vec{x_2} \right) ... P\left( \vec{x_n} \right) P\left( y_n|\vec{x_n} \right)$$

$$\because P\left( y_i \mid \vec{x_i} \right) = \begin{cases} h\left( \vec{x_i} \right), & for\ y_i = +1 \\ 1 - h\left( \vec{x_i} \right), & for\ y_i = -1 \end{cases}$$

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

$$\therefore P\left(y_i \mid \overrightarrow{x_i}\right) = \begin{cases} h(x), & \text{for } y_i = +1 \\ h(-x), & \text{for } y_i = -1 \end{cases} = h(y_i x_i)$$

$$\therefore P_h(D) = P\left(\overrightarrow{x_1}\right)h\left(y_1\overrightarrow{x_1}\right) P\left(\overrightarrow{x_2}\right)h\left(y_2\overrightarrow{x_2}\right)... P\left(\overrightarrow{x_n}\right)h\left(y_n\mid \overrightarrow{x_n}\right)$$

$$\because P\left(\overrightarrow{x_i}\right)\text{固定，所以上式} \propto \prod h\left(y_i \overrightarrow{x_i}\right)$$

$$\therefore \max_h P_h(D) \propto \prod h\left(y_i \overrightarrow{x_i}\right)$$

$$\because \text{每个}h \text{ 对应一个参数}\overrightarrow{w}, \quad \therefore \max_{\overrightarrow{w}} likelihood\left(\overrightarrow{w}\right) \propto \prod s\left(y_i \overrightarrow{w}^T\overrightarrow{x_i}\right)$$

$$\because \text{连乘不容易求解，所以取对数变成连加：} \max_{\overrightarrow{w}} ln \prod s\left(y_i\overrightarrow{w}^T \overrightarrow{x_i}\right) = \max_h \sum ln\, s\left(y_i\overrightarrow{w}\,\overrightarrow{x_i}\right)$$

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

$$\min_{\overrightarrow{w}} \frac{1}{N}\sum - ln\, s\left(y_i\overrightarrow{w}\,\overrightarrow{x_i}\right)$$

$$\therefore s(z) = \frac{1}{1 + exp(-z)} \text{ 代入上式得到：} \min_{\overrightarrow{w}} \frac{1}{N}\sum ln\left(1 + exp\left(-y_i\, \overrightarrow{w}^T\, \overrightarrow{x_i}\right)\right)$$

$$\text{其中，定义交叉熵错误}cross\ entropy\ error : err\left(\overrightarrow{w}, \overrightarrow{x_i}, y_i\right) = ln(1 + exp\left(-y_i\overrightarrow{w}^T\, \overrightarrow{x_i}\right))$$

$$\text{综上，}error\ function: \frac{1}{N}\sum ln\left(1 + exp\left(-y_i\overrightarrow{w}^T\overrightarrow{x_i}\right)\right) = \frac{1}{N}\sum err\left(\overrightarrow{w}, \overrightarrow{x_i}, y_i\right)$$

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

$$E\left(\overrightarrow{w}\right) = \frac{1}{N}\sum ln\left(1 + exp\left(-y_i\overrightarrow{w}^T\overrightarrow{x_i}\right)\right)$$

$$\nabla E\left(\overrightarrow{w}\right) = \frac{\partial E\left(\overrightarrow{w}\right)}{\partial w_i} = \frac{1}{N}\sum\left(\frac{1}{\square}\right)(exp(\triangle))(-y_i\, \overrightarrow{x_i})$$

$$= \frac{1}{N}\sum\left(\frac{exp(\triangle)}{1 + exp(\triangle)}\right)\left(-y_i\, \overrightarrow{x_i}\right)$$

$$= \frac{1}{N}\sum s(\triangle)\left(-y_i\, \overrightarrow{x_i}\right)$$

$$= \frac{1}{N}\sum s\left(-y_i\overrightarrow{w}^T\overrightarrow{x_i}\right)\left(-y_i\, \overrightarrow{x_i}\right)$$

因为$E\left(\overrightarrow{w}\right)$是连续，可微的凸函数，所以令$\nabla E\left(\overrightarrow{w}\right) = 0$，可解的最佳$\overrightarrow{w}$。

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$\overrightarrow{w_0}$，设迭代次数为$t$，$t = 0, 1, 2,...;$

2. 计算梯度$\nabla E(\overrightarrow{w_t}) = \dfrac{1}{N}\sum s(-y_i\overrightarrow{w_t}^T\overrightarrow{x_i})(-y_i\overrightarrow{x_i});$

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

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