

对数几率回归

对数几率回归

二分类问题

通过极大似然估计 w 和 b

Logistic Regression 比较好的翻译是周志华《机器学习》的“对数几率回归”，亦被翻译为“逻辑回归”或“逻辑斯蒂回归”。

虽然名字是回归，但对数几率回归是一种常用于分类(classification)的算法，可以解决二分类或多分类问题，因为对数几率回归的输出是离散的，线性回归的输出才是连续的。

二分类问题

现假设二分类问题的输出标记为 $y \in \{0, 1\}$ ，可以有下列 阶跃函数 (unit-step function)：

$$y = \begin{cases} 0, & z < 0 \\ 0.5, & z = 0 \\ 1, & z > 0 \end{cases} \quad (1)$$

单位阶跃函数不连续，我们希望得到其单调可微的替代函数，对数几率函数 (logistic function) 正是这样一个常用的替代函数。

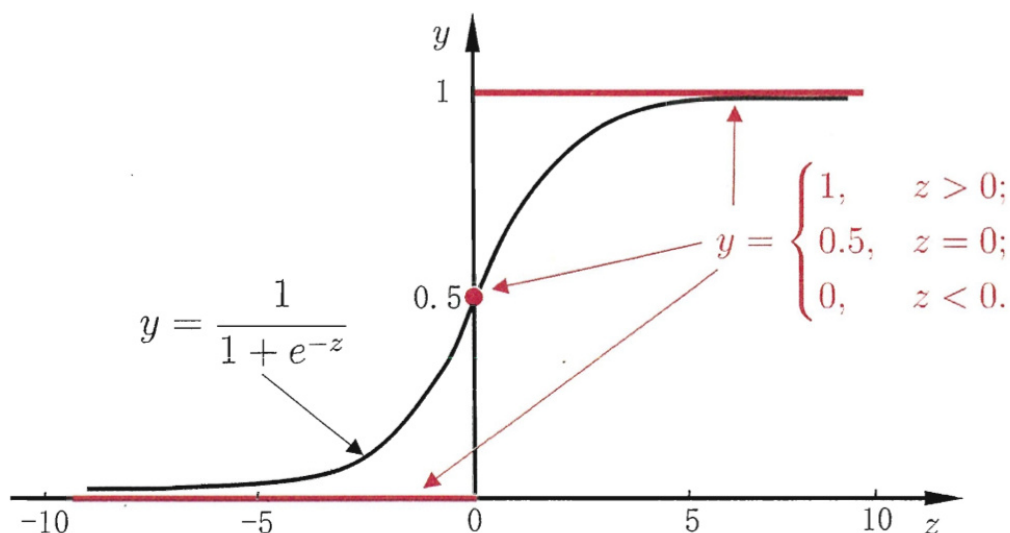


图 3.2 单位阶跃函数与对数几率函数

$$y = \frac{1}{1 + e^{-z}} \quad (2)$$

其实这个函数就是 **Sigmoid 函数**，将 z 用 $w^T x + b$ 表示：

$$y = \frac{1}{1 + e^{-(w^T x + b)}} \quad (3)$$

取对数，可转化为为：

$$\ln \frac{y}{1 - y} = w^T x + b \quad (4)$$

对数几率的名字怎么来的呢？若将 y 视为样本 x 作为正例的相对可能性，则 $1 - y$ 是其反例可能性，两者（正反例可能性）的比值即为**几率 odds**： $\frac{y}{1-y}$ ，又由于我们取了对数，因此 $\ln \frac{y}{1-y}$ 称为**对数几率（log odds，亦称为 logit）**

$$\ln \frac{y}{1 - y} = w^T x + b \quad (5)$$

这个式子实际上是用**线性回归模型** $w^T + b$ 预测的结果来逼近真实标记的**对数几率**，因此此模型也被称为**对数几率回归 Logistic Regression（或 Logit Regression）**，虽然名字是回归，但实际是一种**分类学习方法**。

对数几率回归的优点：

1. 直接对分类可能性建模，无需实现假设数据分布，避免了假设分布不准确带来的问题
2. 不仅预测出类别，还得到近似概率预测
3. 对数几率函数是任意阶可导的凸函数，有很好的数学性质，可利用很多数值优化算法求最优解

通过极大似然估计w和b

还是上面那个公式

$$\ln \frac{y}{1-y} = \mathbf{w}^T \mathbf{x} + b \quad (6)$$

我们把 y 视为类后验概率，上式可重写为：

$$\ln \frac{p(y=1|\mathbf{x})}{p(y=0|\mathbf{x})} = \mathbf{w}^T \mathbf{x} + b \quad (7)$$

$$p(y=1|\mathbf{x}) = \frac{e^{\mathbf{w}^T \mathbf{x} + b}}{1 + e^{\mathbf{w}^T \mathbf{x} + b}} \quad (8)$$

$$p(y=0|\mathbf{x}) = \frac{1}{1 + e^{\mathbf{w}^T \mathbf{x} + b}} \quad (9)$$

我们通过极大似然法（maximum likelihood method）估计 w 和 b ，给定数据集 $\{(\mathbf{x}_i, y_i)\}_{i=1}^m$ ，对率回归模型最大化对数似然 $\ell(w, b)$

$$\ell(\mathbf{w}, b) = \sum_{i=1}^m \ln p(y_i | \mathbf{x}_i; \mathbf{w}, b) \quad (10)$$

为了方便，把 w 和 b 堆叠一起标记为 $\beta = (\mathbf{w}; b)$ ，令 $\hat{\mathbf{x}} = (\mathbf{x}; 1)$ 则 $\mathbf{w}^T \mathbf{x} + b = \beta^T \hat{\mathbf{x}}$ 。

对数似然项 $p(y_i | \mathbf{x}_i; \mathbf{w}, b)$ 写成：

$$p(y_i | \mathbf{x}_i; \mathbf{w}, b) = y_i p_1(\hat{\mathbf{x}}_i; \beta) + (1 - y_i) p_0(\hat{\mathbf{x}}_i; \beta) \quad (11)$$

又由于

$$p(y=1|\mathbf{x}) = \frac{e^{\mathbf{w}^T \mathbf{x} + b}}{1 + e^{\mathbf{w}^T \mathbf{x} + b}} \quad (12)$$

$$p(y=0|\mathbf{x}) = \frac{1}{1 + e^{\mathbf{w}^T \mathbf{x} + b}} \quad (13)$$

最大化 $\ell(\mathbf{w}, b) = \sum_{i=1}^m \ln p(y_i | \mathbf{x}_i; \mathbf{w}, b)$ 等价于最小化 $\ell(\mathbf{w}, b) = \sum_{i=1}^m -\ln p(y_i | \mathbf{x}_i; \mathbf{w}, b)$ ，即负对数

由于 y_i 只能取 0, 1 值，因此 $\ln p_0, \ln p_1$ 只能有一个生效， $-\ln p_0 = 1 + e^{\beta^T \hat{\mathbf{x}}_i}$ ，而 $-\ln p_1 = 1 + e^{\beta^T \hat{\mathbf{x}}_i} + 1 + \beta^T \hat{\mathbf{x}}_i$ ，那么应该最小化下列 $\ell(\beta)$

$$\ell(\beta) = \sum_{i=1}^m \left(-y_i \beta^T \hat{\mathbf{x}}_i + \ln(1 + e^{\beta^T \hat{\mathbf{x}}_i}) \right) \quad (14)$$

$\beta = (\mathbf{w}; \mathbf{b})$ ，这是关于 β 的高阶可导连续凸函数，根据凸优化理论，可以通过经典数值优化算法：**梯度下降法 gradient descent method**、**牛顿法 Newton method** 得到最优解，于是得到

$$\beta^* = \arg \min_{\beta} \ell(\beta) \quad (15)$$