

# 贪心学院推荐系统

## 机器学习基础

### 逻辑回归

简单的事情往往异乎寻常。——保罗·柯艾略

#### 线性模型

线性模型是一个通过属性的线性组合来进行预测的函数，即

$$f(\mathbf{x}) = w_1x_1 + w_2x_2 + \cdots + w_dx_d + b$$

一般用向量形式写成：

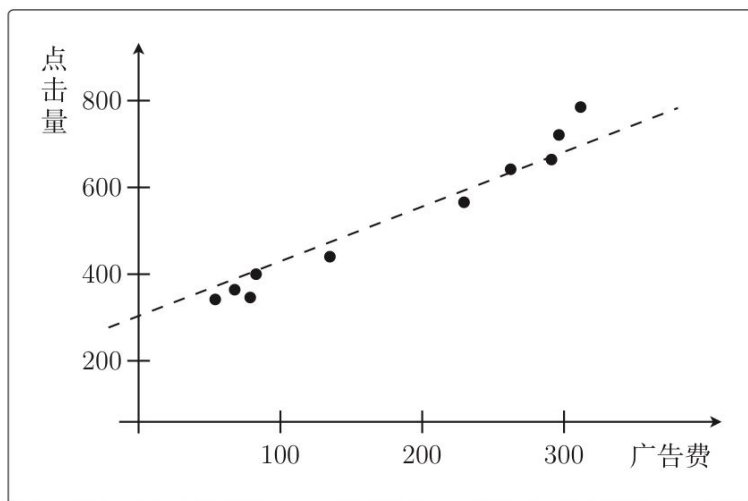
$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$$

其中的 $\mathbf{w}$ 和 $b$ 学得之后，模型就得以确定。

#### 线性回归

回归：regression，倒推的意思，意即由现有的数据倒推得（回归到）原来真实的函数。

给定数据集 $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$ ，其中 $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{id}), y_i \in \mathbb{R}$ 。"线性回归"试图学得一个线性模型尽可能准确地预测实值输出标记。



线性回归是利用最小二乘函数对一个或多个自变量之间关系进行建模的方法。

均方误差有非常好的几何意义，它对应了欧氏距离。基于均方误差最小化进行模型求解的方法成为“最小二乘法”。在线性回归中，最小二乘法就是试图找到一条直线，使所有样本到直线上的欧氏距离之和最小。

#### 最小二乘法

二乘其实是指平方的意思，为什么用平方呢？因为平方可以消除误差正负方向上的差异，单纯的只比较长度。

另一种通俗的说法叫距离（学术一点叫欧氏距离），距离不分上下、左右，只有大小，所以可以用来衡量目标与估计的所有方向偏差累积。

假设我们定义，估计在目标正方向上为正，负方向上为负。计算时只要加一个if-else判断就行了，并不一定非要用平方来代表距离，还可以用绝对值、三次方等。但是平方是一种很好用的数学技巧，从一个方面说来：比起绝对值，平方的微分更加简单。负负得正这个简单的计算规则让二乘法（即两个数相乘）这个名字得以发扬光大，成为一切拟和法的鼻祖。

$$E(\boldsymbol{\theta}) = \frac{1}{2} \sum_{i=1}^n (y^{(i)} - f_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}))^2$$

其中 $f_{\boldsymbol{\theta}}(\mathbf{x})$ 就是我们假设的以 $\boldsymbol{\theta}$ 为参数的线性模型。

- 举例：

假设 $f_{\boldsymbol{\theta}}(x) = \theta_0 + \theta_1 x$

广告费 $x$	点击量 $y$	$\theta_0=1$ 、 $\theta_1=2$ 时的 $f_{\boldsymbol{\theta}}(x)$
58	374	117
70	385	141
81	375	163
84	401	169

$$\begin{aligned}
 E(\boldsymbol{\theta}) &= \frac{1}{2} \sum_{i=1}^4 \left( y^{(i)} - f_{\boldsymbol{\theta}}(x^{(i)}) \right)^2 \\
 &= \frac{1}{2} \times ((374 - 117)^2 + (385 - 141)^2 + (375 - 163)^2 + (401 - 169)^2) \\
 &= \frac{1}{2} \times (66049 + 59536 + 44944 + 53824) \\
 &= 112176.5
 \end{aligned}$$

这个值本身没什么意义，我们用它来指导改变参数 $\boldsymbol{\theta}$ ，使这个值变得越来越小，即让误差变小。这种做法就叫做最小二乘法。

#### 通用数学公式解

误差方程为：

$$E(\mathbf{w}|\mathbf{X}, \mathbf{y}) = (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y})$$

其最优解为：

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

其中 $\mathbf{X}$ 为 $m \times n$ 的样本输入矩阵：

$$\begin{bmatrix}
 \mathbf{x}_{11} & \mathbf{x}_{12} & \dots & \mathbf{x}_{1n} \\
 \mathbf{x}_{21} & \mathbf{x}_{22} & \dots & \mathbf{x}_{2n} \\
 \dots & \dots & \dots & \dots \\
 \mathbf{x}_{m1} & \mathbf{x}_{m2} & \dots & \mathbf{x}_{mn}
 \end{bmatrix}$$

$\mathbf{y}$ 为 $m \times 1$ 的列向量，一般称为 $labels$ ：

$$\begin{bmatrix}
 \mathbf{y}_1 \\
 \mathbf{y}_2 \\
 \dots \\
 \mathbf{y}_m
 \end{bmatrix}$$

$\mathbf{w}$ 为 $n \times 1$ 列向量，就是待求的拟和权重参数：

$$\begin{bmatrix}
 \mathbf{w}_1 \\
 \mathbf{w}_2 \\
 \dots \\
 \mathbf{w}_m
 \end{bmatrix}$$

## 推导过程

误差方程展开：

$$\begin{aligned}(\mathbf{X}\mathbf{w} - \mathbf{y})^T(\mathbf{X}\mathbf{w} - \mathbf{y}) &= ((\mathbf{X}\mathbf{w})^T - \mathbf{y}^T)(\mathbf{X}\mathbf{w} - \mathbf{y}) \\&= \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - (\mathbf{X}\mathbf{w})^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \mathbf{w} + \mathbf{y}^T \mathbf{y} \\&= \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - 2(\mathbf{X}\mathbf{w})^T \mathbf{y} + \mathbf{y}^T \mathbf{y}\end{aligned}$$

其中用到了 $\alpha^T \beta = \beta^T \alpha$ 这一等式。

化简的最后结果的极值（小）在对 $\mathbf{w}$ 求导为零处，所以有

$$2\mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{X}^T \mathbf{y} = 0$$

整理得：

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

---

## 线性回归

### 一元线性回归

我们将求解 $w$ 和 $b$ 使得 $E(w, b)$ （下边提到）最小化的过程称为线性回归模型的最小二乘参数估计。

- 求解偏置 $b$ 的公式推导

- 由最小二乘法导出损失函数 $E(w, b)$ ：

$$\begin{aligned}E(w, b) &= \sum_{i=1}^m (y_i - f(x_i))^2 \\&= \sum_{i=1}^m (y_i - (wx_i + b))^2 \\&= \sum_{i=1}^m (y_i - wx_i - b)^2\end{aligned}$$

- 求解：

- 定理一：二元函数判断凹凸性：

设 $f(x, y)$ 在区域 $D$ 上具有二阶连续偏导数，记 $A = f''_{xx}(x, y)$ ,  $B = f''_{xy}(x, y)$ ,  $C = f''_{yy}(x, y)$ 则：

- (1) 在 $D$ 上恒有 $A > 0$ ，且 $AC - B^2 \geq 0$ 时， $f(x, y)$ 在区域 $D$ 上是凸函数；
- (2) 在 $D$ 上恒有 $A < 0$ ，且 $AC - B^2 \geq 0$ 时， $f(x, y)$ 在区域 $D$ 上是凹函数。

- 定理二：二元凹凸函数求最值：

设 $f(x, y)$ 是在开区域 $D$ 内具有连续偏导数的凸（凹）函数， $(x_0, y_0) \in D$ 且

$f'_x(x_0, y_0) = 0, f'_y(x_0, y_0) = 0$ ，则 $f(x_0, y_0)$ 必为 $f(x, y)$ 在 $D$ 内的最小值（或最大值）。

---

接下来用以上两个定理来求解：

- 证明损失函数 $E(w, b)$ 是关于 $w$ 和 $b$ 的凸函数：

- 求 $A = f''_{xx}(x, y)$ ：

$$\begin{aligned}\frac{\partial E(w, b)}{\partial w} &= \frac{\partial \left[ \sum_{i=1}^m (y_i - wx_i - b)^2 \right]}{\partial w} \\&= \sum_{i=1}^m \frac{\partial (y_i - wx_i - b)^2}{\partial w} \\&= \sum_{i=1}^m 2 \cdot (y_i - wx_i - b) \cdot (-x_i) \\&= 2 \left( w \sum_{i=1}^m x_i^2 - \sum_{i=1}^m (y_i - b)x_i \right)\end{aligned}$$

$$\begin{aligned}
\frac{\partial^2 E(w, b)}{\partial w^2} &= \frac{\partial\left(\frac{\partial E(w, b)}{\partial w}\right)}{\partial w} \\
&= \frac{\partial\left[2\left(w\sum_{i=1}^m x_i^2 - \sum_{i=1}^m (y_i - b)x_i\right)\right]}{\partial w} \\
&= 2\sum_{i=1}^m x_i^2
\end{aligned}$$

■ 求  $B = f''_{xy}(x, y)$ :

$$\begin{aligned}
\frac{\partial^2 E(w, b)}{\partial w \partial b} &= \frac{\partial\left(\frac{\partial E(w, b)}{\partial w}\right)}{\partial b} \\
&= \frac{\partial\left[2\left(w\sum_{i=1}^m x_i^2 - \sum_{i=1}^m (y_i - b)x_i\right)\right]}{\partial b} \\
&= 2\sum_{i=1}^m x_i
\end{aligned}$$

■ 求  $C = f''_{yy}(x, y)$ :

$$\begin{aligned}
\frac{\partial E(w, b)}{\partial b} &= \frac{\partial\left[\sum_{i=1}^m (y_i - wx_i - b)^2\right]}{\partial b} \\
&= \sum_{i=1}^m \frac{\partial(y_i - wx_i - b)^2}{\partial b} \\
&= \sum_{i=1}^m 2 \cdot (y_i - wx_i - b) \cdot (-1) \\
&= 2\left(mb - \sum_{i=1}^m (y_i - wx_i)\right) \\
\frac{\partial^2 E(w, b)}{\partial b^2} &= \frac{\partial\left(\frac{\partial E(w, b)}{\partial b}\right)}{\partial b} \\
&= \frac{\partial\left[2\left(mb - \sum_{i=1}^m (y_i - wx_i)\right)\right]}{\partial b} \\
&= 2m
\end{aligned}$$

■ 以上:  $A = 2\sum_{i=1}^m x_i^2$ ,  $B = 2\sum_{i=1}^m x_i$ ,  $C = 2m$ .

$$\begin{aligned}
AC - B^2 &= 2m \cdot 2 \sum_{i=1}^m x_i^2 - \left( 2 \sum_{i=1}^m x_i \right)^2 \\
&= 4m \sum_{i=1}^m x_i^2 - 4 \left( \sum_{i=1}^m x_i \right)^2 \\
&= 4m \sum_{i=1}^m x_i^2 - 4 \cdot m \cdot \frac{1}{m} \left( \sum_{i=1}^m x_i \right)^2 \\
&= 4m \sum_{i=1}^m x_i^2 - 4m \cdot \bar{x} \cdot \sum_{i=1}^m x_i \\
&= 4m \left( \sum_{i=1}^m x_i^2 - \sum_{i=1}^m x_i \bar{x} \right) \\
&= 4m \sum_{i=1}^m (x_i^2 - x_i \bar{x} - x_i \bar{x} + x_i \bar{x}) \\
&= 4m \sum_{i=1}^m (x_i^2 - x_i \bar{x} - x_i \bar{x} + \bar{x}^2) \\
&= 4m \sum_{i=1}^m (x_i - \bar{x})^2
\end{aligned}$$

我们可知  $AC - B^2 = 4m \sum_{i=1}^m (x_i - \bar{x})^2 \geq 0$ , 也即损失函数  $E(w, b)$  是关于  $w$  和  $b$  的凸函数。

■ 补充:  $\sum_{i=1}^m x_i \bar{x} = \bar{x} \sum_{i=1}^m x_i = \bar{x} \cdot m \cdot \frac{1}{m} \cdot \sum_{i=1}^m x_i = m \bar{x}^2 = \sum_{i=1}^m \bar{x}^2$ .

令一阶偏导数等于0解出  $b$ :

$$\begin{aligned}
\frac{\partial E_{(w,b)}}{\partial b} &= 2 \left( mb - \sum_{i=1}^m (y_i - wx_i) \right) = 0 \\
mb - \sum_{i=1}^m (y_i - wx_i) &= 0 \\
b &= \frac{\sum_{i=1}^m (y_i - wx_i)}{m} = \frac{1}{m} \sum_{i=1}^m y_i - w \cdot \frac{1}{m} \sum_{i=1}^m x_i = \bar{y} - w\bar{x}.
\end{aligned}$$

令一阶偏导数为0解出  $w$ :

$$\begin{aligned}
\frac{\partial E_{(w,b)}}{\partial w} &= 2 \left( w \sum_{i=1}^m x_i^2 - \sum_{i=1}^m (y_i - b)x_i \right) = 0 \\
w \sum_{i=1}^m x_i^2 - \sum_{i=1}^m (y_i - b)x_i &= 0 \\
w \sum_{i=1}^m x_i^2 &= \sum_{i=1}^m y_i x_i - \sum_{i=1}^m b x_i
\end{aligned}$$

将  $b = \bar{y} - w\bar{x}$  代入上式可得:

$$\begin{aligned}
w \sum_{i=1}^m x_i^2 &= \sum_{i=1}^m y_i x_i - \sum_{i=1}^m (\bar{y} - w\bar{x}) x_i \\
w \sum_{i=1}^m x_i^2 &= \sum_{i=1}^m y_i x_i - \bar{y} \sum_{i=1}^m x_i + w\bar{x} \sum_{i=1}^m x_i \\
w \sum_{i=1}^m x_i^2 - w\bar{x} \sum_{i=1}^m x_i &= \sum_{i=1}^m y_i x_i - \bar{y} \sum_{i=1}^m x_i \\
w \left( \sum_{i=1}^m x_i^2 - \bar{x} \sum_{i=1}^m x_i \right) &= \sum_{i=1}^m y_i x_i - \bar{y} \sum_{i=1}^m x_i \\
w &= \frac{\sum_{i=1}^m y_i x_i - \bar{y} \sum_{i=1}^m x_i}{\sum_{i=1}^m x_i^2 - \bar{x} \sum_{i=1}^m x_i} = \frac{\sum_{i=1}^m y_i x_i - \bar{x} \sum_{i=1}^m y_i}{\sum_{i=1}^m x_i^2 - \frac{1}{m} (\sum_{i=1}^m x_i)^2} = \frac{\sum_{i=1}^m y_i (x_i - \bar{x})}{\sum_{i=1}^m x_i^2 - \frac{1}{m} (\sum_{i=1}^m x_i)^2}
\end{aligned}$$

其中:

$$\begin{aligned}
\blacksquare \bar{y} \sum_{i=1}^m x_i &= \frac{1}{m} \sum_{i=1}^m y_i \sum_{i=1}^m x_i = \bar{x} \sum_{i=1}^m y_i \\
\blacksquare \bar{x} \sum_{i=1}^m x_i &= \frac{1}{m} \sum_{i=1}^m x_i \sum_{i=1}^m x_i = \frac{1}{m} (\sum_{i=1}^m x_i)^2
\end{aligned}$$

为提高运算速度, 将求解 $w$ 的过程向量化可得 (核心: 将累加的形式抽象成向量的点乘):

$$w = \frac{\sum_{i=1}^m y_i (x_i - \bar{x})}{\sum_{i=1}^m x_i^2 - \frac{1}{m} (\sum_{i=1}^m x_i)^2}$$

将 $\frac{1}{m} (\sum_{i=1}^m x_i^2) = \bar{x} \sum_{i=1}^m x_i = \sum_{i=1}^m x_i \bar{x}$ 代入分母可得:

$$w = \frac{\sum_{i=1}^m y_i (x_i - \bar{x})}{\sum_{i=1}^m x_i^2 - \sum_{i=1}^m x_i \bar{x}} = \frac{\sum_{i=1}^m (y_i x_i - y_i \bar{x})}{\sum_{i=1}^m (x_i^2 - x_i \bar{x})}$$

由于:

$$\begin{cases} \sum_{i=1}^m y_i \bar{x} = \bar{x} \sum_{i=1}^m y_i = \frac{1}{m} \sum_{i=1}^m x_i \sum_{i=1}^m y_i = \sum_{i=1}^m x_i \cdot \frac{1}{m} \cdot \sum_{i=1}^m y_i = \sum_{i=1}^m x_i \bar{y} \\ \sum_{i=1}^m y_i \bar{x} = \bar{x} \sum_{i=1}^m y_i = \bar{x} \cdot m \cdot \frac{1}{m} \sum_{i=1}^m y_i = m \bar{x} \bar{y} = \sum_{i=1}^m \bar{x} \bar{y} \\ \sum_{i=1}^m x_i \bar{x} = \bar{x} \sum_{i=1}^m x_i = \bar{x} \cdot m \cdot \frac{1}{m} \sum_{i=1}^m x_i = m \bar{x}^2 = \sum_{i=1}^m \bar{x}^2 \end{cases}$$

所以:

$$w = \frac{\sum_{i=1}^m (y_i x_i - y_i \bar{x})}{\sum_{i=1}^m (x_i^2 - x_i \bar{x})} = \frac{\sum_{i=1}^m (y_i x_i - y_i \bar{x} - y_i \bar{x} + y_i \bar{x})}{\sum_{i=1}^m (x_i^2 - x_i \bar{x} - x_i \bar{x} + x_i \bar{x})} = \frac{\sum_{i=1}^m (y_i x_i - y_i \bar{x} - x_i \bar{y} + \bar{x} \bar{y})}{\sum_{i=1}^m (x_i^2 - 2x_i \bar{x} + \bar{x}^2)} = \frac{\sum_{i=1}^m (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^m (x_i - \bar{x})^2}$$

令 $\mathbf{x} = (x_1, x_2, \dots, x_m)^T$ ,  $\mathbf{y} = (y_1, y_2, \dots, y_m)^T$

$\mathbf{x}_d = (x_1 - \bar{x}, x_2 - \bar{x}, \dots, x_m - \bar{x})^T$ ,  $\mathbf{y}_d = (y_1 - \bar{y}, y_2 - \bar{y}, \dots, y_m - \bar{y})^T$ .

则

$$w = \frac{\sum_{i=1}^m (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^m (x_i - \bar{x})^2} = \frac{\mathbf{x}_d^T \mathbf{y}_d}{\mathbf{x}_d^T \mathbf{x}_d}$$

- 值得注意的是：上述求 $w$ 的式子在高中选修1-2出现过。后续求 $b$ 的式子只用到了性质： $(\bar{x}, \bar{y})$ 一定在回归方程上。（我们上述求解是先求 $b$ 后求的 $w$ ，与课本相反）。

## 多元线性回归

多元线性回归的求解思路同一元线性回归的求解思路基本相同：



- 将 $w$ 和 $b$ 组合成 $\hat{w}$ ：

$$f(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i + b = \hat{\mathbf{w}}^T \hat{\mathbf{x}}_i$$

- 由最小二乘法导出损失函数 $E_{\hat{w}}$ ：

$$E_{\hat{w}} = \sum_{i=1}^m (y_i - f(\hat{\mathbf{x}}_i))^2 = \sum_{i=1}^m (y_i - \hat{\mathbf{w}}^T \hat{\mathbf{x}}_i)^2$$

- 进行向量化：

$$\text{令 } \mathbf{X} = (\hat{\mathbf{x}}_1^T, \hat{\mathbf{x}}_2^T, \dots, \hat{\mathbf{x}}_m^T)^T, \mathbf{y} = (y_1, y_2, \dots, y_m)^T$$

得：

$$E_{\hat{w}} = (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})^T (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})$$

求解：

- 证明损失函数是关于 $\hat{w}$ 的凸函数：

- 凸集：设集合 $D \in R^n$ ，如果对任意的点 $\mathbf{x}, \mathbf{y} \in D$ 与任意的 $a \in [0, 1]$ ，有 $a\mathbf{x} + (1-a)\mathbf{y} \in D$ ，则称集合 $D$ 是凸集。



- 梯度（多元实值函数的一阶导数）：设 $n$ 元函数 $f(\mathbf{x})$ 对自变量 $\mathbf{x} = (x_1, x_2, \dots, x_n)^T$ 的各分量 $x_i$ 的一阶偏导数 $\frac{\partial f(\mathbf{x})}{\partial x_i}$  ( $i = 1, 2, \dots, n$ )都存在，则称函数 $f(\mathbf{x})$ 在 $\mathbf{x}$ 处一阶可导，并称向量

$$\nabla f(\mathbf{x}) = \begin{pmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} \\ \frac{\partial f(\mathbf{x})}{\partial x_2} \\ \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_n} \end{pmatrix}$$

为函数 $f(\mathbf{x})$ 在 $\mathbf{x}$ 处的一阶导数或梯度，记为 $\nabla f(\mathbf{x})$ （列向量）。

- *Hessian* (海塞) 矩阵: 设  $n$  元函数  $f(\mathbf{x})$  对自变量  $\mathbf{x} = (x_1, x_2, \dots, x_n)^T$  的各分量  $x_i$  的二阶偏导数  $\frac{\partial^2 f(\mathbf{x})}{\partial x_i \partial x_j}$  ( $i = 1, 2, \dots, n; j = 1, 2, \dots, n$ ) 都存在, 则称函数  $f(\mathbf{x})$  在  $\mathbf{x}$  处二阶可导, 并称矩阵

$$\nabla^2 f(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2 f(\mathbf{x})}{\partial x_1^2} & \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_2} & \dots & \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_2 \partial x_1} & \frac{\partial^2 f(\mathbf{x})}{\partial x_2^2} & \dots & \frac{\partial^2 f(\mathbf{x})}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_n \partial x_1} & \frac{\partial^2 f(\mathbf{x})}{\partial x_n \partial x_2} & \dots & \frac{\partial^2 f(\mathbf{x})}{\partial x_n^2} \end{bmatrix}$$

为  $f(\mathbf{x})$  在  $\mathbf{x}$  处的二阶导数或 *Hessian* 矩阵, 记为  $\nabla^2 f(\mathbf{x})$ , 若  $f(\mathbf{x})$  对  $\mathbf{x}$  各变元的所有二阶偏导数都连续, 则  $\frac{\partial^2 f(\mathbf{x})}{\partial x_i \partial x_j} = \frac{\partial^2 f(\mathbf{x})}{\partial x_j \partial x_i}$ , 此时易知  $\nabla^2 f(\mathbf{x})$  为对称矩阵。

- **定理一: 多元实值函数凹凸性判定定理:**

设  $D \subset R^n$  是非空开凸集。  $f: D \subset R^n \rightarrow R$  且  $f(\mathbf{x})$  在  $D$  上二阶连续可微。如果  $f(\mathbf{x})$  的 *Hessian* 矩阵  $\nabla^2 f(\mathbf{x})$  在  $D$  上是正定的, 则  $f(\mathbf{x})$  是  $D$  上的严格凸函数。

- **定理二: 凸充分性定理:**

若  $f: R^n \rightarrow R$  是凸函数, 且  $f(\mathbf{x})$  一阶连续可微, 则  $\mathbf{x}^*$  是全局解的充分必要条件是  $\nabla f(\mathbf{x}^*) = \mathbf{0}$ , 其中  $\nabla f(\mathbf{x})$  为  $f(\mathbf{x})$  关于  $\mathbf{x}$  的一阶导数 (也称梯度)。

接下来用以上两个定理来求解:

$$\begin{aligned} \frac{\partial E_{\hat{\mathbf{w}}}}{\partial \hat{\mathbf{w}}} &= \frac{\partial [(\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})^T(\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})]}{\partial \hat{\mathbf{w}}} \\ &= \frac{\partial [(\mathbf{y}^T - \hat{\mathbf{w}}^T \mathbf{X}^T)(\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})]}{\partial \hat{\mathbf{w}}} \\ &= \frac{\partial [\mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \hat{\mathbf{w}} - \hat{\mathbf{w}}^T \mathbf{X}^T \mathbf{y} + \hat{\mathbf{w}}^T \mathbf{X}^T \mathbf{X} \hat{\mathbf{w}}]}{\partial \hat{\mathbf{w}}} \\ &= \frac{\partial [-\mathbf{y}^T \mathbf{X} \hat{\mathbf{w}} - \hat{\mathbf{w}}^T \mathbf{X}^T \mathbf{y} + \hat{\mathbf{w}}^T \mathbf{X}^T \mathbf{X} \hat{\mathbf{w}}]}{\partial \hat{\mathbf{w}}} \\ &= \frac{\partial \mathbf{y}^T \mathbf{X} \hat{\mathbf{w}}}{\partial \hat{\mathbf{w}}} - \frac{\partial \hat{\mathbf{w}}^T \mathbf{X}^T \mathbf{y}}{\partial \hat{\mathbf{w}}} + \frac{\partial \hat{\mathbf{w}}^T \mathbf{X}^T \mathbf{X} \hat{\mathbf{w}}}{\partial \hat{\mathbf{w}}} \end{aligned}$$

由矩阵微分公式  $\frac{\partial \mathbf{x}^T \mathbf{a}}{\partial \mathbf{x}} = \frac{\partial \mathbf{a}^T \mathbf{x}}{\partial \mathbf{x}} = \mathbf{a}$ ,  $\frac{\partial \mathbf{x}^T \mathbf{B} \mathbf{x}}{\partial \mathbf{x}} = (\mathbf{B} + \mathbf{B}^T) \mathbf{x}$  可得:

$$\frac{\partial E_{\hat{\mathbf{w}}}}{\partial \hat{\mathbf{w}}} = -\mathbf{X}^T \mathbf{y} - \mathbf{X}^T \mathbf{y} + (\mathbf{X}^T \mathbf{X} + \mathbf{X}^T \mathbf{X}) \hat{\mathbf{w}} = 2\mathbf{X}^T (\mathbf{X} \hat{\mathbf{w}} - \mathbf{y})$$

则 *Hessian* 矩阵为:

$$\begin{aligned} \frac{\partial^2 E_{\hat{\mathbf{w}}}}{\partial \hat{\mathbf{w}} \partial \hat{\mathbf{w}}^T} &= \frac{\partial}{\partial \hat{\mathbf{w}}} \left( \frac{\partial E_{\hat{\mathbf{w}}}}{\partial \hat{\mathbf{w}}} \right) \\ &= \frac{\partial}{\partial \hat{\mathbf{w}}} [2\mathbf{X}^T (\mathbf{X} \hat{\mathbf{w}} - \mathbf{y})] \\ &= \frac{\partial}{\partial \hat{\mathbf{w}}} (2\mathbf{X}^T \mathbf{X} \hat{\mathbf{w}} - 2\mathbf{X}^T \mathbf{y}) \\ &= 2\mathbf{X}^T \mathbf{X} \end{aligned}$$

当  $\mathbf{X}^T \mathbf{X}$  是正定矩阵时, 损失函数  $E_{\hat{\mathbf{w}}}$  是关于  $\hat{\mathbf{w}}$  的凸函数。

$$\text{令 } \frac{\partial E_{\hat{\mathbf{w}}}}{\partial \hat{\mathbf{w}}} = 2\mathbf{X}^T (\mathbf{X} \hat{\mathbf{w}} - \mathbf{y}) = \mathbf{0}$$

解出:  $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$



有时分类问题也可以转化为回归问题，例如肺癌预测，我们可以用回归模型先预测出患肺癌的概率，然后再给定一个阈值，例如50%，概率值在50%以下的人划为没有肺癌，50%以上则认为患有肺癌。这种分类型问题的回归算法预测，最常用的就是逻辑回归。

---

## 逻辑回归

---

如果面试官问你熟悉哪个机器学习模型，可以说 *SVM*，但千万别提 *LR*，因为细节真的太多了。

——阿泽

---

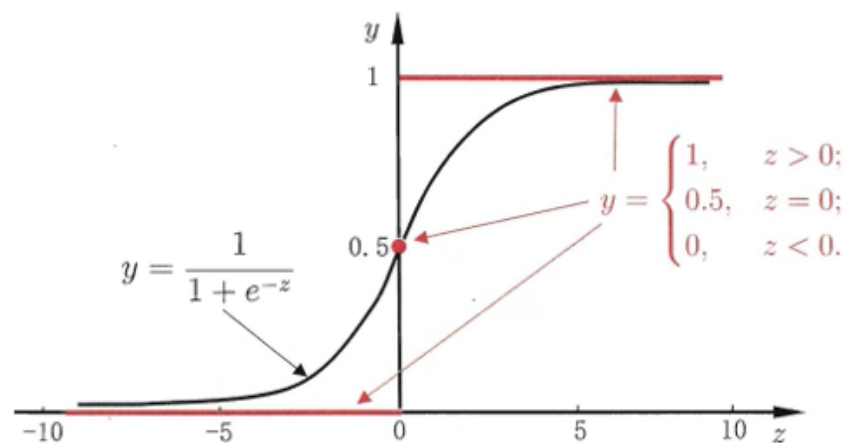
*Logistic*回归的本质是：假设数据服从这个分布，然后使用极大似然估计做参数的估计。

*Logistic*回归=线性回归+*sigmoid*函数

---

*logistic*函数

*logistic*函数中文名叫对数几率函数，因其形状像 *S*，是 *Sigmoid* 函数最重要的代表，故又称 *Sigmoid* 函数：



*Sigmoid*函数的性质：

- 将压缩到(0, 1)之间（可以由此联想到概率的取值）
- $\frac{1}{2}$ 处导数值最大
- $y(x)$ 的导数为 $y(x)(1 - y(x))$
- 两边梯度趋于饱和（作为激活函数在神经网络的弊端）
- 不以原点为中心（作为激活函数在神经网络的弊端）
- 单调递增

---

逻辑回归为何不叫逻辑分类？

回归来源于“线性回归”，使用线性回归去拟合逼近一个“界”（对数几率），使得按照这个“界”进行数据分类后得到的 $cost$ 最小。以概率0.5为分界线，将数据分为正例和反例。使得 $z > 0$ 对应于正例， $z < 0$ 对应于反例。因此是使用的回归思想去解决分类问题。

---

广义线性模型

---

指数族分布

看似无关的事物背后往往有着不可思议的联系。

——佚名

---

指数族分布是一类分布的总称，该类分布的分布律（或者概率密度）的一般形式如下：

$$p(y; \eta) = b(y)e^{\eta^T T(y) - a(\eta)}$$

其中,  $\eta$ 称为该分布的自然参数;  $T(y)$ 为充分统计量, 视具体的分布而定, 通常是等于随机变量 $y$ 本身;  $a(\eta)$ 为配分函数;  $b(y)$ 为关于随机变量 $y$ 的函数。(  $e^{-a(\eta)}$ 本质上是一个归一化常数, 保证 $\sum p(y; \eta) = 1$ )。也就是说 $T, a, b$ 确定了一种分布,  $\eta$ 是该分布的参数。

常见的伽玛分布、泊松分布、二项分布和正态分布均属于指数族分布。

举例说明**二项分布是指数族分布**:

已知二项分布的分布律为:

$$p(y) = \phi^y(1 - \phi)^{1-y}$$

其中 $y \in \{0, 1\}$ ,  $\phi$ 为 $y = 1$ 的概率。对上式恒等变形可得:

$$\begin{aligned} p(y) &= \phi^y(1 - \phi)^{1-y} \\ &= e^{(\ln(\phi^y(1-\phi)^{1-y}))} \\ &= e^{(\ln(\phi^y) + \ln(1-\phi)^{1-y})} \\ &= e^{(y\ln\phi + (1-y)\ln(1-\phi))} \\ &= e^{(y\ln\phi + \ln(1-\phi) - y\ln(1-\phi))} \\ &= e^{(y(\ln\phi - \ln(1-\phi)) + \ln(1-\phi))} \\ &= e^{y\ln(\frac{\phi}{1-\phi}) + \ln(1-\phi)} \end{aligned}$$

对比指数族分布的一般形式 $p(y; \eta) = b(y)e^{\eta^T T(y) - a(\eta)}$ 可知:

$$\begin{aligned} b(y) &= 1 \\ \eta &= \ln\left(\frac{\phi}{1-\phi}\right) \\ T(y) &= y \\ a(\eta) &= -\ln(1-\phi) = \ln(1 + e^\eta) \end{aligned}$$

故得证。

考虑一个分类或者回归问题, 我们就是想预测某个随机变量 $y$ ,  $y$ 是某些特征 $\mathbf{x}$ 的函数。为了推导出广义线性模型, 我们必须做出如下三条假设:

#### 广义线性模型的三条假设

- 在给定 $\mathbf{x}$ 的条件下, 假设随机变量 $y$ 服从某个**指数族分布**;
- 在给定 $\mathbf{x}$ 的条件下, 我们的目标是得到一个模型 $h(\mathbf{x})$ 能预测出 $T(y)$ 的**期望值**;
- 假设该指数族分布中的自然参数 $\eta$ 和 $\mathbf{x}$ 是线性相关的, 即 $\eta = \mathbf{w}^T \mathbf{x}$ 。

#### 对数几率回归

已知 $y$ 服从二项分布, 而二项分布属于指数族分布, 所以满足广义线性模型的**第一条假设**, 根据**第二条假设**我们可以推得模型 $h(\mathbf{x})$ 的表达式应该为:

$$h(\mathbf{x}) = E[T(y|\mathbf{x})]$$

由于二项分布的 $T(y|\mathbf{x}) = y|\mathbf{x}$ , 所以:

$$h(\mathbf{x}) = E[y|\mathbf{x}]$$

又因为 $E[y|\mathbf{x}] = 1 \times p(y = 1|\mathbf{x}) + 0 \times p(y = 0|\mathbf{x}) = \phi$ , 所以 $h(\mathbf{x}) = \phi$ 。

而由 $\eta = \ln(\frac{\phi}{1-\phi})$ 可推出 $\phi = \frac{1}{1 + e^{-\eta}}$ 。

所以可得:

$$h(\mathbf{x}) = \phi = \frac{1}{1 + e^{-\eta}}$$

再根据广义模型的**第三条假设**:  $\eta = \mathbf{w}^T \mathbf{x}$ ,  $h(\mathbf{x})$ 最终可化为:

$$h(\mathbf{x}) = \phi = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}} = p(y = 1 | \mathbf{x})$$

## 对数几率回归的极大似然估计

似然和概率是不同的。*likelihood*不是*probability*。

当把 $\mathbf{x}$ 固定而把 $f(\mathbf{x}, \boldsymbol{\theta})$ 看做 $\boldsymbol{\theta}$ 的函数时，它称为“似然函数”。这个名称的意义，可根据分析得到理解：这个函数对于不同的 $\boldsymbol{\theta}$ 的取值，反映了在观察结果 $\mathbf{x}$ 已知的条件下， $\boldsymbol{\theta}$ 的各种值的“似然程度”。注意，这里有些像贝叶斯公式中的推理：把观察值 $\mathbf{x}$ 看成结果，而把参数值 $\boldsymbol{\theta}$ 看成是导致这个结果的原因。现已有了结果，要反过来推算各种原因的概率。这里参数 $\boldsymbol{\theta}$ 是确定的值（虽然未知），并非随机变量，无概率可言，于是就改用“似然”这个词。——陈希孺《概率论与数理统计》

在统计学上，基于某些模型的参数（粗略地说，我们可以认为参数决定了模型），观测到某数据的概率称为概率；而已经观测到某数据，模型的参数取特定值的概率称为似然。

已知随机变量 $y$ 取1和0的概率分别为

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

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

将 $b$ 考虑进 $\mathbf{w}$ ，令 $\boldsymbol{\beta} = (\mathbf{w}; b)$ ， $\hat{\mathbf{x}} = (\mathbf{x}; 1)$ ，则 $\mathbf{w}^T \mathbf{x} + b$ 可以简写为 $\boldsymbol{\beta}^T \hat{\mathbf{x}}$ ，于是上式可化简为：

$$p(y = 1 | \mathbf{x}) = \frac{e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}}}{1 + e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}}}$$

$$p(y = 0 | \mathbf{x}) = \frac{1}{1 + e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}}}$$

为了简单表示，记：

$$p(y = 1 | \mathbf{x}) = \frac{e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}}}{1 + e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}}} = p_1(\hat{\mathbf{x}}; \boldsymbol{\beta})$$

$$p(y = 0 | \mathbf{x}) = \frac{1}{1 + e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}}} = p_0(\hat{\mathbf{x}}; \boldsymbol{\beta})$$

使用一个简单的技巧可得到随机变量 $y$ 的分布律表达式

$$p(y | \mathbf{x}; \mathbf{w}, b) = y \cdot p_1(\hat{\mathbf{x}}; \boldsymbol{\beta}) + (1 - y) \cdot p_0(\hat{\mathbf{x}}; \boldsymbol{\beta})$$

或者

$$p(y | \mathbf{x}; \mathbf{w}, b) = [p_1(\hat{\mathbf{x}}; \boldsymbol{\beta})]^y \times [p_0(\hat{\mathbf{x}}; \boldsymbol{\beta})]^{1-y}$$

根据对数似然函数的定义可得：

$$\ell_{(w,b)} = \ln L(w, b) = \sum_{i=1}^m \ln p(y_i | x_i; w, b)$$

将 $p(y | \mathbf{x}; \mathbf{w}, b) = y \cdot p_1(\hat{\mathbf{x}}; \boldsymbol{\beta}) + (1 - y) \cdot p_0(\hat{\mathbf{x}}; \boldsymbol{\beta})$ 代入可得：

$$\ell(\boldsymbol{\beta}) = \sum_{i=1}^m \ln (y_i p_1(\hat{\mathbf{x}}_i; \boldsymbol{\beta}) + (1 - y_i) \cdot p_0(\hat{\mathbf{x}}_i; \boldsymbol{\beta}))$$

将

$$\frac{e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i}}{1 + e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i}} = p_1(\hat{\mathbf{x}}_i; \boldsymbol{\beta})$$

$$\frac{1}{1 + e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i}} = p_0(\hat{\mathbf{x}}_i; \boldsymbol{\beta})$$

代入上式可得：

$$\begin{aligned}\ell(\beta) &= \sum_{i=1}^m \ln \left( \frac{y_i e^{\beta^T \hat{x}_i}}{1 + e^{\beta^T \hat{x}_i}} + \frac{1 - y_i}{1 + e^{\beta^T \hat{x}_i}} \right) \\ &= \sum_{i=1}^m \left( \ln(y_i e^{\beta^T \hat{x}_i} + 1 - y_i) - \ln(1 + e^{\beta^T \hat{x}_i}) \right)\end{aligned}$$

由于  $y_i \in \{0, 1\}$ , 所以可以用一个小技巧继续化简为:

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

损失函数一般转化为最小化问题, 加个负号即可:

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

- **补充:** 也可以用  $p(y|\mathbf{x}; \mathbf{w}, b) = [p_1(\hat{\mathbf{x}}; \beta)]^y \times [p_0(\hat{\mathbf{x}}; \beta)]^{1-y}$  来推导。

$$\begin{aligned}\ell(\beta) &= \sum_{i=1}^m \ln ([p_1(\hat{\mathbf{x}}_i; \beta)]^{y_i} \times [p_0(\hat{\mathbf{x}}_i; \beta)]^{1-y_i}) \\ &= \sum_{i=1}^m \ln ([p_1(\hat{\mathbf{x}}_i; \beta)]^{y_i}) + \ln ([p_0(\hat{\mathbf{x}}_i; \beta)]^{1-y_i}) \\ &= \sum_{i=1}^m [y_i \ln ([p_1(\hat{\mathbf{x}}_i; \beta)]) + (1 - y_i) \ln (p_0(\hat{\mathbf{x}}_i; \beta))] \\ &= \sum_{i=1}^m \{y_i [\ln(p_1(\hat{\mathbf{x}}_i; \beta)) - \ln(p_0(\hat{\mathbf{x}}_i; \beta))] + \ln(p_0(\hat{\mathbf{x}}_i; \beta))\} \\ &= \sum_{i=1}^m \left[ y_i \ln \left( \frac{p_1(\hat{\mathbf{x}}_i; \beta)}{p_0(\hat{\mathbf{x}}_i; \beta)} \right) + \ln(p_0(\hat{\mathbf{x}}_i; \beta)) \right]\end{aligned}$$

同样将

$$\begin{aligned}\frac{e^{\beta^T \hat{x}_i}}{1 + e^{\beta^T \hat{x}_i}} &= p_1(\hat{\mathbf{x}}_i; \beta) \\ \frac{1}{1 + e^{\beta^T \hat{x}_i}} &= p_0(\hat{\mathbf{x}}_i; \beta)\end{aligned}$$

代入可得:

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

但是  $\ell(\beta) = \sum_{i=1}^m \left( -y_i \beta^T \hat{x}_i + \ln(1 + e^{\beta^T \hat{x}_i}) \right)$  并不能像线性回归函数一样得到关于  $\beta$  的解析解。我们可以用经典的比如梯度下降法, 牛顿法来解。下面介绍梯度下降法求解。

**梯度下降法求解**

取逻辑回归的损失函数为:

$$\ell(\beta) = - \sum_{i=1}^m [y_i \ln ([p_1(\hat{\mathbf{x}}_i; \beta)]) + (1 - y_i) \ln (p_0(\hat{\mathbf{x}}_i; \beta))]$$

我们对其中的  $p_1(\hat{\mathbf{x}}_i; \beta)$  对  $\beta_j$  求导得到:

$$\frac{\partial p_1(\hat{\mathbf{x}}_i; \beta)}{\partial \beta_j} = p_1(\hat{\mathbf{x}}_i; \beta) \cdot (1 - p_1(\hat{\mathbf{x}}_i; \beta)) \cdot \hat{x}_{ij}$$

其中  $\mathbf{x}_i$  实际指的是第  $i$  个样本的特征向量, 即  $(x_{i1}, x_{i2}, \dots, x_{im}, 1)$ , 我们注意到只有  $x_{ij}$  会和  $\beta_j$  相乘, 因此一眼可看出求导结果。

$$\begin{aligned}
\frac{\partial \ell(\beta)}{\partial \beta_j} &= - \sum_{i=1}^m \left( y_i \frac{1}{p_1(\hat{\mathbf{x}}_i; \beta)} \frac{\partial p_1(\hat{\mathbf{x}}_i; \beta)}{\partial \beta_j} + (1 - y_i) \frac{1}{1 - p_1(\hat{\mathbf{x}}_i; \beta)} \frac{\partial p_1(\hat{\mathbf{x}}_i; \beta)}{\partial \beta_j} \right) \\
&= - \sum_{i=1}^m \left( \frac{y_i}{p_1(\hat{\mathbf{x}}_i; \beta)} - \frac{1 - y_i}{1 - p_1(\hat{\mathbf{x}}_i; \beta)} \right) \cdot \frac{\partial p_1(\hat{\mathbf{x}}_i; \beta)}{\partial \beta_j} \\
&= - \sum_{i=1}^m \left( \frac{y_i}{p_1(\hat{\mathbf{x}}_i; \beta)} - \frac{1 - y_i}{1 - p_1(\hat{\mathbf{x}}_i; \beta)} \right) \cdot p_1(\hat{\mathbf{x}}_i; \beta) \cdot (1 - p_1(\hat{\mathbf{x}}_i; \beta)) \cdot \hat{x}_{ij} \\
&= - \sum_{i=1}^m [y_i(1 - p_1(\hat{\mathbf{x}}_i; \beta)) - (1 - y_i)p_1(\hat{\mathbf{x}}_i; \beta)] \cdot x_{ij} \\
&= - \sum_{i=1}^m (y_i - p_1(\hat{\mathbf{x}}_i; \beta)) \cdot x_{ij} \\
&= \sum_{i=1}^m \left( \frac{e^{\beta^T \hat{\mathbf{x}}_i}}{1 + e^{\beta^T \mathbf{x}_i}} - y_i \right) \cdot x_{ij}
\end{aligned}$$

有了偏导，也就有了梯度 $\mathbf{G}$ (即偏导数组成的向量)。则梯度下降算法过程如下：

- 初始化向量 $\beta$ 的值为 $\theta_0$ ，将其代入 $G$ 得到当前位置的梯度；
- 用步长 $\alpha$ 乘以当前梯度，得到从当前位置下降的距离；
- 更新 $\theta_1$ ，其更新表达式为 $\theta_1 = \theta_0 - \alpha \mathbf{G}$ ；
- 重复以上步骤，直到更新至 $\theta_k$ 达到停止条件，此时 $\theta_k$ 就是我们所求的参数向量 $\beta$ 。

---

## 神经网络

## 正则化

## 常用优化算法

## 推荐系统基础

---