

2.1 Logistic回归

CSDN学院
2017年10月

► 大纲



- Logistic回归基本原理
- 多类Logistic回归
- Scikit learn 中的Logistic回归实现
- 分类模型的评价
- 模型选择与参数调优
- 案例分析



► 分类

- 给定训练数据 $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^N$, 分类任务学习一个从输入 \mathbf{x} 到输出 y 的映射 f :

$$\hat{y} = f(\mathbf{x}) = \arg \max_c p(y = c | \mathbf{x}, \mathcal{D})$$

- 其中 y 为离散值, 其取值范围称为标签空间: $\mathcal{Y} = \{1, 2, \dots, C\}$
- 当 $C=2$ 时, 为两类分类问题, 计算出 $p(y=1 | \mathbf{x})$ 即可。此时分布为Bernoulli分布:

$$p(y | \mathbf{x}) = \text{Ber}(y | \mu(\mathbf{x}))$$

– 其中 $\mu(\mathbf{x}) = \mathbb{E}(y | \mathbf{x}) = p(y=1 | \mathbf{x})$

► Logistic回归模型

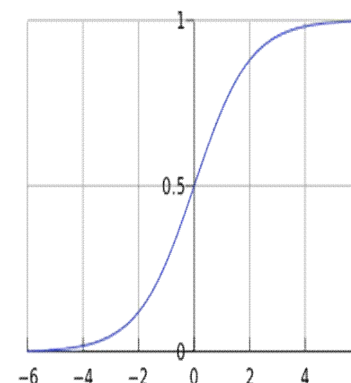
- Logistic回归模型同线性回归模型类似，也是一个线性模型，只是条件概率 $p(y|\mathbf{x})$ 的形式不同：

$$p(y|\mathbf{x}) = \text{Ber}(y|\mu(\mathbf{x})),$$

$$\mu(\mathbf{x}) = \text{sigm}(\mathbf{w}^T \mathbf{x})$$

- 其中sigmoid函数（S形函数）定义为

$$\text{sigm}(\eta) = \frac{1}{1 + \exp(-\eta)} = \frac{\exp(\eta)}{\exp(\eta) + 1}$$



- 亦被称为logistic函数或logit函数，将实数 η 变换到 $[0,1]$ 区间。
 - 因为概率取值在 $[0,1]$ 区间。
 - Logistic回归亦被称为logit回归

► 极大似然估计

$$\text{Ber}(x | \theta) = \theta^x (1 - \theta)^{1-x}$$

- Logistic回归： $p(y | \mathbf{x}, \mathbf{w}) = \text{Ber}(y | \mu(\mathbf{x}))$, $\mu(\mathbf{x}) = \text{sigm}(\mathbf{w}^T \mathbf{x})$
- 令 $\mu_i = \mu(\mathbf{x}_i)$, 则负log似然为

$$\begin{aligned} J(\mathbf{w}) = NLL(\mathbf{w}) &= -\sum_{i=1}^N \log \left[(\mu_i)^{y_i} \times (1 - \mu_i)^{(1-y_i)} \right] \\ &= \sum_{i=1}^N - \left[y_i \log(\mu_i) + (1 - y_i) \log(1 - \mu_i) \right] \end{aligned}$$

Logistic损失

极大似然估计 等价于 最小Logistic损失

优化求解：梯度下降 / 牛顿法

► 梯度

- 目标函数为

$$J(\mathbf{w}) = \sum_{i=1}^N -[y_i \log(\mu_i) + (1 - y_i) \log(1 - \mu_i)]$$

- 梯度为

$$g(\mathbf{w}) = \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = \frac{\partial}{\partial \mathbf{w}} \left[\sum_{i=1}^N -[y_i \log(\mu_i) + (1 - y_i) \log(1 - \mu_i)] \right]$$

$$J(\mathbf{w}) = -\sum_{i=1}^N \left[y_i \log(\mu_i) + (1 - y_i) \log(1 - \mu_i) \right]$$

$$\begin{aligned} g(\mathbf{w}) &= \frac{\partial}{\partial \mathbf{w}} J(\mathbf{w}) = \sum_{i=1}^N \left[-y_i \times \frac{1}{\mu(\mathbf{x}_i)} \frac{\partial}{\partial \mathbf{w}} \mu(\mathbf{x}_i) + (1 - y_i) \times \frac{1}{1 - \mu(\mathbf{x}_i)} \frac{\partial}{\partial \mathbf{w}} \mu(\mathbf{x}_i) \right] \\ &= \sum_{i=1}^N \left[-y_i \times \frac{1}{\mu(\mathbf{x}_i)} + (1 - y_i) \times \frac{1}{1 - \mu(\mathbf{x}_i)} \right] \frac{\partial}{\partial \mathbf{w}} \mu(\mathbf{x}_i) \leftarrow \\ &= \sum_{i=1}^N \left[-y_i \times \frac{1}{\mu(\mathbf{x}_i)} + (1 - y_i) \times \frac{1}{1 - \mu(\mathbf{x}_i)} \right] \mu(\mathbf{x}_i) (1 - \mu(\mathbf{x}_i)) \mathbf{x}_i \\ &= \sum_{i=1}^N \left[-y_i \times [1 - \mu(\mathbf{x}_i)] + (1 - y_i) \mu(\mathbf{x}_i) \right] \mathbf{x}_i \\ &= \sum_{i=1}^N \left[-y_i + \mu(\mathbf{x}_i) \right] \mathbf{x}_i \\ &= \sum_{i=1}^N \left[\mu(\mathbf{x}_i) - y_i \right] \mathbf{x}_i \end{aligned}$$

$$\frac{\partial}{\partial \mathbf{w}} \mu(\mathbf{x}) = \mu(\mathbf{x}) (1 - \mu(\mathbf{x})) \mathbf{x}$$

$$\mu(\mathbf{x}) = \frac{\exp(\mathbf{w}^T \mathbf{x})}{\exp(\mathbf{w}^T \mathbf{x}) + 1}$$

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

$$\begin{aligned} \frac{\partial}{\partial \mathbf{w}} \mu(\mathbf{x}) &= \frac{\frac{\partial}{\partial \mathbf{w}} [\exp(\mathbf{w}^T \mathbf{x})] (\exp(\mathbf{w}^T \mathbf{x}) + 1) - \exp(\mathbf{w}^T \mathbf{x}) \frac{\partial}{\partial \mathbf{w}} [\exp(\mathbf{w}^T \mathbf{x}) + 1]}{[\exp(\mathbf{w}^T \mathbf{x}) + 1]^2} \\ &= \frac{\exp(\mathbf{w}^T \mathbf{x}) (\exp(\mathbf{w}^T \mathbf{x}) + 1) \frac{\partial}{\partial \mathbf{w}} (\mathbf{w}^T \mathbf{x}) - \exp(\mathbf{w}^T \mathbf{x}) \exp(\mathbf{w}^T \mathbf{x}) \frac{\partial}{\partial \mathbf{w}} (\mathbf{w}^T \mathbf{x})}{[\exp(\mathbf{w}^T \mathbf{x}) + 1]^2} \\ &= \frac{\exp(\mathbf{w}^T \mathbf{x})}{[\exp(\mathbf{w}^T \mathbf{x}) + 1]^2} \mathbf{x} = \mu(\mathbf{x}) (1 - \mu(\mathbf{x})) \mathbf{x} \end{aligned}$$

$$\frac{\partial}{\partial \mathbf{w}} (\mathbf{w}^T \mathbf{x}) = \frac{\partial}{\partial \mathbf{w}} (\mathbf{x}^T \mathbf{w}) = \mathbf{x}$$

$$\mu(\mathbf{x}) = \frac{\exp(\mathbf{w}^T \mathbf{x})}{\exp(\mathbf{w}^T \mathbf{x}) + 1}$$

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

► 梯度

- 目标函数为

$$J(\mathbf{w}) = \sum_{i=1}^N -[y_i \log(\mu_i) + (1 - y_i) \log(1 - \mu_i)]$$

- 梯度为

$$g(\mathbf{w}) = \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = \sum_{i=1}^N (\mu(\mathbf{x}_i) - y_i) \mathbf{x}_i$$

算法与线性回归 $g(\mathbf{w}) = \sum_{i=1}^N (f(\mathbf{x}_i) - y_i) \mathbf{x}_i$ 看起来一样！
当然 $f(\mathbf{x})$ 不同（线性回归中 $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ ）
事实上所有的线性模型的梯度都是如此

- 二阶Hessian矩阵为

$$\mathbf{H}(\mathbf{w}) = \frac{\partial}{\partial \mathbf{w}} [\mathbf{g}(\mathbf{w})^T] = \sum_{i=1}^N \left(\frac{\partial}{\partial \mathbf{w}} \mu_i \right) \mathbf{x}_i^T$$

$$= \mu_i (1 - \mu_i) \mathbf{x}_i \mathbf{x}_i^T = \mathbf{X}^T \text{diag}(\mu_i (1 - \mu_i)) \mathbf{X}$$

正定矩阵，凸优化



► 牛顿法

- 亦称牛顿-拉夫逊 (Newton-Raphson) 方法
 - 牛顿在17世纪提出的一种近似求解方程的方法
 - 使用函数 $f(x)$ 的泰勒级数的前面几项来寻找方程 $f(x)=0$ 的根
- 在求极值问题中，求 $g(\mathbf{w}) = \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = 0$ 的根
 - 对应处 $J(\mathbf{w})$ 取极值

► 牛顿法

- 将导数 $\mathbf{g}(\mathbf{w})$ 在 \mathbf{w}^t 处进行Taylor展开：

$$0 = \mathbf{g}(\hat{\mathbf{w}}) = \mathbf{g}(\mathbf{w}^t) + (\hat{\mathbf{w}} - \mathbf{w}^t) \mathbf{H}(\mathbf{w}^t) + Op(\hat{\mathbf{w}} - \mathbf{w}^t)$$

- 去掉高阶无穷小 $Op(\hat{\mathbf{w}} - \mathbf{w}^t)$ ，从而得到

$$\mathbf{g}(\mathbf{w}^t) + (\hat{\mathbf{w}} - \mathbf{w}^t) \mathbf{H}(\mathbf{w}^t) = 0 \Rightarrow \hat{\mathbf{w}} = \mathbf{w}^t - \mathbf{H}^{-1}(\mathbf{w}^t) \mathbf{g}(\mathbf{w}^t)$$

- 因此迭代机制为：

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \mathbf{H}^{-1}(\mathbf{w}^t) \mathbf{g}(\mathbf{w}^t)$$

- 也被称为二阶梯度下降法，移动方向： $\mathbf{H}(\mathbf{w}^t) \mathbf{d} = -\mathbf{g}(\mathbf{w}^t)$
- Vs. 一阶梯度法，移动方向： $\mathbf{d} = -\mathbf{g}(\mathbf{w}^t)$ 移动

► Iteratively Reweighted Least Squares (IRLS)

不止于代码

引入记号：

$$\mathbf{g}^t(\mathbf{w}) = \mathbf{X}^T (\boldsymbol{\mu}^t - \mathbf{y}), \quad \mu_i^t = \text{sigm}\left(\left(\mathbf{w}^t\right)^T \mathbf{x}_i\right)$$

$$\mathbf{H}^t(\mathbf{w}) = \mathbf{X}^T \mathbf{S}^t \mathbf{X}, \quad \mathbf{S}^t := \text{diag}\left(\mu_1^t(1-\mu_1^t), \dots, \mu_N^t(1-\mu_N^t)\right)$$

根据牛顿法的结果：

$$\begin{aligned} \mathbf{w}^{t+1} &= \mathbf{w}^t - \left(\mathbf{H}^t\right)^{-1} \mathbf{g}^t \\ &= \mathbf{w}^t + \left(\mathbf{X}^T \mathbf{S}^t \mathbf{X}\right)^{-1} \mathbf{X}^T (\mathbf{y} - \boldsymbol{\mu}^t) \\ &= \left(\mathbf{X}^T \mathbf{S}^t \mathbf{X}\right)^{-1} \left[\left(\mathbf{X}^T \mathbf{S}^t \mathbf{X}\right) \mathbf{w}^t + \mathbf{X}^T (\mathbf{y} - \boldsymbol{\mu}^t) \right] \\ &= \left(\mathbf{X}^T \mathbf{S}^t \mathbf{X}\right)^{-1} \mathbf{X}^T \left[\mathbf{S}^t \mathbf{X} \mathbf{w}^t + \mathbf{y} - \boldsymbol{\mu}^t \right] \\ &= \left(\mathbf{X}^T \mathbf{S}^t \mathbf{X}\right)^{-1} \mathbf{X}^T \mathbf{S}^t \left[\mathbf{X} \mathbf{w}^t + \left(\mathbf{S}^t\right)^{-1} (\mathbf{y} - \boldsymbol{\mu}^t) \right] \quad \text{where } \mathbf{z}^t = \mathbf{X} \mathbf{w}^t + \left(\mathbf{S}^t\right)^{-1} (\mathbf{y} - \boldsymbol{\mu}^t) \\ &= \left(\mathbf{X}^T \mathbf{S}^t \mathbf{X}\right)^{-1} \mathbf{X}^T \mathbf{S}^t \mathbf{z} \end{aligned}$$

加权最小二乘问题的标准方程

$$\text{回忆最小二乘：} \hat{\mathbf{w}}_{OLS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

但权重矩阵 \mathbf{S} 不是常数，而是依赖参数向量 \mathbf{w} 。因此我们必须使用标准方程来迭代计算，每次使用新的权向量 \mathbf{w} 来修正权重矩阵 \mathbf{S} 。因此该算法被称为迭代再加权最小二乘（iterative reweighted least squares, IRLS）。

► IRLS (cont.)

- 回忆最小二乘问题：
 - 目标函数： $J(\mathbf{w}) = \sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x})^2 = (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w})$
 - 解： $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$
- 回忆加权最小二乘问题： $(\Sigma^{-1}$ ：权重矩阵)
 - 目标函数： $J(\mathbf{w}) = (\mathbf{y} - \mathbf{X}\mathbf{w})^T \Sigma^{-1} (\mathbf{y} - \mathbf{X}\mathbf{w})$
 - 解： $\hat{\mathbf{w}} = (\mathbf{X}^T \Sigma^{-1} \mathbf{X})^{-1} \mathbf{X}^T \Sigma^{-1} \mathbf{y}$
- IRLS中， $\mathbf{w}^{t+1} = (\mathbf{X}^T \mathbf{S}^t \mathbf{X})^{-1} \mathbf{X}^T \mathbf{S}^t \left[\mathbf{X}\mathbf{w}^t + (\mathbf{S}^t)^{-1} (\mathbf{y} - \boldsymbol{\mu}^t) \right]$
 - 相当于权重矩阵为 $\Sigma^{-1} = \mathbf{S}^t$
 - 由于 \mathbf{S}^t 是对角阵， \mathbf{S} 相当于给每个样本的权重为 $S_{ii}^t = \mu_i^t (1 - \mu_i^t)$ ，



$$\mathbf{z}_i^t = (\mathbf{w}^t)^T \mathbf{x}_i + \frac{y_i - \mu_i^t}{S_{ii}^t}$$

► Iteratively Reweighted Least Squares (cond)

- Iteratively reweighted least squares(IRLS)

```

1  $\mathbf{w} = \mathbf{0}_D$ 
2  $w_0 = \log(\bar{y} / (1 - \bar{y}))$ 
3 repeat
4    $\eta_i = w_0 + \mathbf{w}^T \mathbf{x}_i$ 
5    $\mu_i = \text{sigm}(\eta_i)$ 
6    $s_i = \mu_i(1 - \mu_i)$ 
7    $z_i = \eta_i + \frac{y_i - \mu_i}{s_i}$ 
8    $\mathbf{S} = \text{diag}(s_{1:N})$ 
9    $\mathbf{w} = (\mathbf{X}^T \mathbf{S} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{S} \mathbf{z}$ 
10 until converged

```

$$\mathbf{S} = \text{diag}(\mu_1(1 - \mu_1), \dots, \mu_N(1 - \mu_N))$$

$$\mathbf{z}_i = \mathbf{w}^T \mathbf{x}_i + \frac{y_i - \mu_i}{\mu_i(1 - \mu_i)}$$

Weighted least square

► 拟牛顿法

- 牛顿法比一般的梯度下降法收敛速度快，但是在高维情况下，计算目标函数的二阶偏导数的复杂度很大，而且有时候目标函数的海森矩阵无法保持正定，不存在逆矩阵，此时牛顿法将不再能使用。
- 因此，人们提出了**拟牛顿法**。其基本思想是：不用二阶偏导数而构造出可以近似Hessian矩阵(或Hessian矩阵的逆矩阵)的正定对称矩阵，进而再逐步优化目标函数。不同的构造方法就产生了不同的拟牛顿法（Quasi-Newton Methods）
 - BFGS / LBFGS / Newton-CG

► 正则化的Logistic回归

- 若损失函数取logistic损失，则Logistic回归的目标函数为

$$J(\mathbf{w}) = \sum_{i=1}^N -[y_i \log(\mu_i) + (1 - y_i) \log(1 - \mu_i)]$$

- 同线性回归类似，Logistic回归亦可加上L2正则

$$J(\mathbf{w}) = \sum_{i=1}^N -[y_i \log(\mu_i) + (1 - y_i) \log(1 - \mu_i)] + \lambda \|\mathbf{w}\|_2^2$$

- 或L1正则

$$J(\mathbf{w}) = \sum_{i=1}^N -[y_i \log(\mu_i) + (1 - y_i) \log(1 - \mu_i)] + \lambda |\mathbf{w}|$$

► 小结



- Logistic回归：
 - 损失函数：负log似然损失
 - 正则：L2/L1正则
 - 优化：梯度下降 / 牛顿法 / 拟牛顿法



THANK YOU



AI100