

# Math Basis of Logistic Regression

周陆延

1. 给出利用梯度下降法求解逻辑回归的前反向公式推导,在这一部分中,你只需要考虑最朴素的二分类逻辑回归模型

- 前向传播

线性回归

$$f(x) = \omega^\top x + b \quad (1)$$

令  $x' = (1, x)$  则  $f(x)$  可统一为

$$f(x) = \omega'^\top x \quad (2)$$

为处理二分类问题将线性回归  $(-\infty, +\infty)$  的值域转换为  $(0, 1)$  对概率进行预测, 嵌套  $\sigma(x) = \frac{1}{1+e^{-x}}$  函数

即逻辑回归的前向传播公式

$$h_\omega(x) = \frac{1}{1 + e^{-\omega^\top x}} \quad (3)$$

- 反向传播

根据逻辑回归的极大似然估计设置交叉熵损失函数

$$J(\omega) = -\frac{1}{m} \sum_{i=1}^m y \log \hat{y} + (1 - y) \log(1 - \hat{y}) \quad (4)$$

求梯度得

$$\begin{aligned} \frac{\delta J(\omega)}{\delta \omega} &= -\frac{1}{m} \sum_{i=1}^m y \frac{1}{\hat{y}} \hat{y}' + (1 - y) \frac{1}{1 - \hat{y}} \hat{y}' \\ &= -\frac{1}{m} \sum_{i=1}^m \left( \frac{y}{\hat{y}} + \frac{1 - y}{1 - \hat{y}} \right) \hat{y}' \\ &= -\frac{1}{m} \sum_{i=1}^m (\sigma(\omega^\top x^{(i)}) - y^{(i)}) x^{(i)} \end{aligned} \quad (5)$$

使用梯度下降法更新  $\omega := \omega - \alpha \frac{\delta J(\omega)}{\delta \omega}$  其中  $\alpha$  为学习率

2. 假设标签集不是  $\{0, 1\}$  而是  $\{1, -1\}$ , 将会有什么变化, 请给出推导  
似然函数变为

$$L(\omega) = \sqrt{\text{Pr}(x)^{1+y} (1 - \text{Pr}(x))^{1-y}} \quad (6)$$

取负对数得到损失函数

$$J(\omega) = -\frac{1}{2} ((1 + y) \log \hat{y} + (1 - y) \log(1 - \hat{y})) \quad (7)$$

对  $y = 1, y = -1$  分类讨论得

$$J(\omega) = \log(1 + e^{y^* \omega^T x}) \quad (8)$$

3. 在问题2.1的基础上，即标签集为  $\{0,1\}$  的情况，分别增加L1正则化和L2正则化，公式和模型效果分别会有什么变化，请给出推导

最基本的正则化方法是在原目标（代价）函数中添加惩罚项，对复杂度高的模型进行“惩罚”，即

$$\tilde{J}(x) = J(x) + \lambda \Omega(\omega) \quad (9)$$

正则化可理解为对原损失函数最优化过程添加约束

$$\begin{aligned} \min_{\omega} J(\omega) \\ s.t. \Omega(\omega) \leq C \end{aligned} \quad (10)$$

利用拉格朗日算子法，我们可将上述带约束条件的最优化问题转换为不带约束项的优化问题，构造拉格朗日函数

$$L(\omega, \lambda) = J(\omega) + \lambda(\Omega(\omega) - C) \quad (11)$$

设  $\lambda$  最优解为  $\lambda^*$  则对拉格朗日函数最小化等价于

$$\min_{\omega} J(\omega) + \lambda^* \Omega(\omega) \quad (12)$$

#### • L2 正则化

即使用L2范数作为惩罚

$$J(\omega) = -\frac{1}{m} \sum_{i=1}^m y \log \hat{y} + (1-y) \log(1-\hat{y}) + \frac{\lambda}{2m} \|\omega\|_2^2 \quad (13)$$

梯度则变为

$$\frac{\delta J(\omega)}{\delta \omega} = -\frac{1}{m} \left( \sum_{i=1}^m (\sigma(\omega^T x^{(i)}) - y^{(i)}) x^{(i)} \right) + \frac{\lambda}{m} \omega \quad (14)$$

考虑对模型的影响，令  $\omega^*$  为未正则化的目标函数的最优解，对  $J(\omega)$  作二阶泰勒展开近似（ $\omega^*$  为最优，无一阶导项；略去样本数量  $m$ ）

$$\hat{J}(\omega) = J(\omega^*) + \frac{1}{2}(\omega - \omega^*)^T H(\omega - \omega^*) \quad (15)$$

当  $\hat{J}(\omega)$  最小时，其梯度为

$$\frac{\delta \hat{J}(\omega)}{\delta \omega} = H(\omega - \omega^*) = 0 \quad (16)$$

加入惩罚项，记此时的最优解为  $\omega'$  得

$$\begin{aligned} H(\omega' - \omega^*) + \lambda \omega' &= 0 \\ \omega' &= (H + \lambda I)^{-1} H \omega^* \end{aligned} \quad (17)$$

由  $H$  实对称，将其合同到对角矩阵  $H = Q \Lambda Q^T$  带入上式得

$$\begin{aligned}
\omega' &= (Q\Lambda Q^\top + \lambda I)^{-1} Q\Lambda Q^\top \omega^* \\
&= (Q\Lambda Q^\top + Q(\lambda I)Q^\top)^{-1} Q\Lambda Q^\top \omega^* \\
&= (\Lambda + \lambda I)^{-1} \Lambda \omega^*
\end{aligned} \tag{18}$$

发现  $\omega'$  相比  $\omega^*$  是依据  $H$  的特征值在对应分量上做了缩放，在  $H$  特征值较大的方向影响较小，在特征值较小的方向影响较大，使对减少目标函数作用显著的参数被保留，作用微弱的参数被衰减

• L1 正则化

即使用L1范数作为惩罚

$$J(\omega) = -\frac{1}{m} \left( \sum_{i=1}^m y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}) \right) + \frac{\lambda}{m} \|\omega\|_1 \tag{19}$$

梯度则变为

$$\frac{\delta J(\omega)}{\delta \omega} = -\frac{1}{m} \left( \sum_{i=1}^m (\sigma(\omega^\top x^{(i)}) - y^{(i)}) x^{(i)} \right) + \frac{\lambda}{m} \text{sgn}(\omega) \tag{20}$$

与L2正则化类比，但是由于  $\text{sgn}(x)$  的特殊性，进一步假设  $H$  为对角阵（对数据进行主成分分析后成立），则

$$\begin{aligned}
\tilde{J}(\omega) &= J(\omega^*) + \frac{1}{2}(\omega - \omega^*)^\top H(\omega - \omega^*) + \lambda \|\omega\|_1 \\
&= J(\omega^*) + \sum_i \left( \frac{1}{2} H_{i,i} (\omega_i - \omega_i^*)^2 + \lambda |\omega_i| \right)
\end{aligned} \tag{21}$$

令

$$\tilde{J}'(\omega) = \sum_i H_{i,i} (\omega_i - \omega_i^*) + \lambda \text{sgn}(\omega_i) = 0 \tag{22}$$

得

$$\omega_i = \text{sgn}(\omega_i^*) \max \left\{ |\omega_i^*| - \frac{\lambda}{H_{i,i}}, 0 \right\} \tag{23}$$

可见当  $|\omega_i^*| \leq \frac{\lambda}{H_{i,i}}$  时会使得  $\omega_i$  变成 0，使得参数稀疏化

4. 给出核逻辑回归的对偶形式

为书写方便给出标签集为  $\{+1, -1\}$  的情况，其他情况可将标签集  $\{a, b\}$  进行映射  $y_i \rightarrow \frac{2}{b-a} \left( y_i - \frac{b-a}{2} \right)$

$$\begin{aligned}
\min_{\alpha} \frac{1}{2\lambda} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \kappa(x_i, x_j) + \sum_{i=1}^m [\alpha_i \log \alpha_i + (1 - \alpha_i) \log(1 - \alpha_i)] \\
\text{s.t. } 0 \leq \alpha_i < 1
\end{aligned} \tag{24}$$