

$$D = \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(N)}, y^{(N)})\}$$

$$x^{(i)} = \begin{bmatrix} 1 \\ x_1^{(i)} \\ \vdots \\ x_d^{(i)} \end{bmatrix} \in \mathbb{R}^{d+1} \quad y^{(i)} \in \{0, 1\}$$

Logistic Function:

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

导数

$$\sigma'(z) = \sigma(z) \cdot (1 - \sigma(z))$$

$$\begin{aligned} \sigma'(z) &= \left(\frac{1}{1+e^z} \right)' \\ &= -\frac{1}{(1+e^z)^2} \cdot (e^z)' \\ &= -\frac{1}{(1+e^z)^2} \cdot e^z \cdot (1) \\ &= -\frac{1}{1+e^z} \cdot \frac{e^z}{1+e^z} \cdot (-1) \\ &= \sigma(z) \cdot (1 - \sigma(z)) \end{aligned}$$

模型参数

$$w = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix} \xrightarrow{\text{bias}}$$

定义概率

$$p(y=1|x;w) = \sigma(w^T x)$$

$$p(y=0|x;w) = 1 - \sigma(w^T x)$$

单个样本的统一概率表示:

$$p(y|x;w) = \sigma(w^T x)^y \cdot (1 - \sigma(w^T x))^{(1-y)}$$

数据集似然, (IID假设)

$$L(D) = \prod_{i=1}^N p(y^{(i)}|x^{(i)};w) = \prod_{i=1}^N \sigma(w^T x^{(i)})^{y^{(i)}} \cdot (1 - \sigma(w^T x^{(i)}))^{(1-y^{(i)})}$$

Log 似然

$$\begin{aligned} L(D) &= \sum_{i=1}^N \log p(y^{(i)}|x^{(i)};w) \\ &= \sum_{i=1}^N \log \{ \sigma(w^T x^{(i)})^{y^{(i)}} \cdot (1 - \sigma(w^T x^{(i)}))^{(1-y^{(i)})} \} \\ &= \sum_{i=1}^N \{ y^{(i)} \log \sigma(w^T x^{(i)}) + (1-y^{(i)}) \log (1 - \sigma(w^T x^{(i)})) \} \end{aligned}$$

损失函数 Negative log likelihood (NLL)

$$L(w) = - \sum_{i=1}^N \{ y^{(i)} \log \sigma(w^T x^{(i)}) + (1-y^{(i)}) \log (1 - \sigma(w^T x^{(i)})) \}$$

通过判断Hessian矩阵是否半正定验证Logistic Regression的损失函数为凸函数。

为简化Notation, 我们用偷懒的方式:

可以先验证单个数据的损失函数是凸函数, 然后根据凸函数加法的保凸性, 可知完整的损失函数也是凸函数。

单个样本损失函数:

$$l(w) = -y \log \sigma(w^T x) - (1-y) \log (1 - \sigma(w^T x))$$

一阶导数:

$$\begin{aligned} \frac{\partial l(w)}{\partial w} &= -y \frac{\partial \log \sigma(w^T x)}{\partial w} - (1-y) \cdot \frac{\partial \log (1 - \sigma(w^T x))}{\partial w} \\ &= -y \frac{1}{\sigma(w^T x)} \cdot \frac{\partial \sigma(w^T x)}{\partial w} - (1-y) \frac{1}{1 - \sigma(w^T x)} \cdot \frac{\partial (1 - \sigma(w^T x))}{\partial w} \\ &= -y \frac{1}{\sigma(w^T x)} \cdot \sigma(w^T x) \cdot (1 - \sigma(w^T x)) \cdot x \\ &\quad - (1-y) \cdot \frac{1}{1 - \sigma(w^T x)} \cdot [-\sigma(w^T x) (1 - \sigma(w^T x))] \cdot x \\ &= -y (1 - \sigma(w^T x)) x + (1-y) \sigma(w^T x) x \\ &= -yx + \sigma(w^T x) y \cdot x + \sigma(w^T x) \cdot x - y \sigma(w^T x) \cdot x \\ &= (\sigma(w^T x) - y) \cdot x \end{aligned}$$

$$\frac{\partial l(w)}{\partial w} = \underbrace{(\sigma(w^T x) - y)}_{\in \mathbb{R}} \cdot \underbrace{x}_{\in \mathbb{R}^{d+1}}$$

二阶导数 (向量对向量求导, 结果为Hessian矩阵)

$$\begin{aligned} \frac{\partial^2 l(w)}{\partial w \partial w^T} &= \frac{\partial [(\sigma(w^T x) - y) \cdot x]}{\partial w^T} \\ &= \begin{bmatrix} \frac{\partial [(\sigma(w^T x) - y) \cdot x_0]}{\partial w^T} \\ \frac{\partial [(\sigma(w^T x) - y) \cdot x_1]}{\partial w^T} \\ \vdots \\ \frac{\partial [(\sigma(w^T x) - y) \cdot x_d]}{\partial w^T} \end{bmatrix} = \begin{bmatrix} \frac{\partial [(\sigma(w^T x) - y) \cdot x_0]}{\partial w_0} & \frac{\partial [(\sigma(w^T x) - y) \cdot x_0]}{\partial w_1} & \dots \\ \frac{\partial [(\sigma(w^T x) - y) \cdot x_1]}{\partial w_0} & \frac{\partial [(\sigma(w^T x) - y) \cdot x_1]}{\partial w_1} & \dots \\ \vdots & \vdots & \ddots \\ \frac{\partial [(\sigma(w^T x) - y) \cdot x_d]}{\partial w_0} & \frac{\partial [(\sigma(w^T x) - y) \cdot x_d]}{\partial w_1} & \dots \end{bmatrix} \\ &= \begin{bmatrix} \sigma(w^T x) (1 - \sigma(w^T x)) x_0^2, & \sigma(w^T x) (1 - \sigma(w^T x)) x_0 x_1, & \dots \\ \sigma(w^T x) (1 - \sigma(w^T x)) x_1 x_0, & \sigma(w^T x) (1 - \sigma(w^T x)) x_1^2, & \dots \\ \vdots & \vdots & \ddots \end{bmatrix} \end{aligned}$$

$$= \begin{bmatrix} \vdots & \vdots & \vdots \\ \sigma(w^T x)(1-\sigma(w^T x))x_d x_0, & \sigma(w^T x)(1-\sigma(w^T x))x_d x_1, & \vdots \end{bmatrix}$$

$$= \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_d \end{bmatrix} \cdot \underbrace{[\sigma(w^T x)(1-\sigma(w^T x))]}_{1 \times 1} \cdot [x_0, x_1, \dots, x_d]_{1 \times (d+1)}$$

$(d+1) \times 1$
 1×1
 $1 \times (d+1)$

$$= x \cdot \sigma(w^T x)(1-\sigma(w^T x)) \cdot x^T$$

$$\Rightarrow \frac{\partial^2 l(w)}{\partial w \partial w^T} = x \cdot \sigma(w^T x) \cdot (1-\sigma(w^T x)) \cdot x^T$$

令 $u \in \mathbb{R}^{d+1}$ 则

$$\begin{aligned} u^T \cdot \frac{\partial^2 l(w)}{\partial w \partial w^T} u &= u^T \cdot x \cdot \sigma(w^T x) \cdot (1-\sigma(w^T x)) \cdot x^T \cdot u \\ &= \underbrace{(x^T u)^T}_{\geq 0} \cdot \underbrace{(x^T u) \cdot \sigma(w^T x) \cdot (1-\sigma(w^T x))}_{> 0} \\ &\geq 0 \end{aligned}$$

$$\Rightarrow \frac{\partial^2 l(w)}{\partial w \partial w^T} \succcurlyeq 0$$

$\therefore l(w)$ 为凸函数

$\therefore L(w) = \sum_{i=1}^N l_i(w)$ 也为凸函数

不使用保凸性，也可由单个样本的损失得到整体损失函数的Hessian矩阵，直接对整体损失函数判断半正定属性。

$$\frac{\partial^2 L(w)}{\partial w \partial w^T} = \sum_{i=1}^N x^{(i)} \sigma(w^T x^{(i)}) (1 - \sigma(w^T x^{(i)})) x^{(i)T}$$

$$= \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(N)} \end{bmatrix} \begin{bmatrix} \sigma(w^T x^{(1)}) (1 - \sigma(w^T x^{(1)})) & & & \\ & \sigma(w^T x^{(2)}) (1 - \sigma(w^T x^{(2)})) & & \\ & & \ddots & \\ & & & \sigma(w^T x^{(N)}) (1 - \sigma(w^T x^{(N)})) \end{bmatrix} \begin{bmatrix} x^{(1)T} \\ x^{(2)T} \\ \vdots \\ x^{(N)T} \end{bmatrix}$$

$(D+1) \times N$
 $N \times N$
 $N \times (D+1)$

$$= X D X^T \quad (D_{ii} = \sigma(w^T x^{(i)}) (1 - \sigma(w^T x^{(i)})))$$

$$\forall u \in \mathbb{R}^{(D+1)}:$$

$$u^T X \cdot D \cdot X^T \cdot u = (X^T u)^T \cdot D \cdot (X^T u) \geq 0$$

$$\therefore \frac{\partial^2 L(w)}{\partial w \partial w^T} \succcurlyeq 0 \quad \therefore L(w) \text{ 为凸函数}$$