

8.1.3. 逻辑回归——不是回归的回归。

模型背景: $z = f(y)$, $y = w^T x + b$. $z \in \{0, 1\}$

从 y 的实值到 z 的 0-1 值 套用函数.

①. f : unit-step function: $f = \begin{cases} 0 & z < 0 \\ 0.5 & z = 0 \\ 1 & z > 0 \end{cases}$ 不可导原点

②. f : sigmoid function: $f = \frac{1}{1+e^x}$ 连续可导, 原点导值大变化快

\Rightarrow 模型 $z = \frac{1}{1+e^{-(w^T x + b)}}$ $\Leftrightarrow \ln \frac{z}{1-z} = w^T x + b$

2. 模型意义: $z=0$: 样本标记为负类 negative $z=1$: 样本标记为正类 positive.

所以 logistic regression 不是解决回归问题, 而是分类啦! (DIO 厨狂喜)

3. 模型求解的指标:

假使说 z 数值代表了一个右验概率 $p(z=1|x)$. 那么 $\ln \frac{p(y=1|x)}{p(y=0|x)} = w^T x + b$

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

那么我们试着用极大似然法求解: (极大似然函数是把 $y=1$ 的概率全乘再取对数)

$$L(w, b) = \sum_{i=1}^m \ln p(y_i | x_i). \quad \text{与线性回归一样, 记 } \beta = (w^T, b)^T \quad X = \begin{pmatrix} x_1^T \\ \vdots \\ x_m^T \end{pmatrix}$$

$$p(y_i | x_i) = y_i p(x_i, \beta) + (1 - y_i) p_0(x_i, \beta).$$

$$\Rightarrow L(w, b) = \sum_{i=1}^m \ln [y_i p(x_i, \beta) + (1 - y_i) p_0(x_i, \beta)].$$

解目标函数最小值

4. 求解模型:

求解方法有牛顿法和梯度下降法.

解 $L(\beta)$ 的最小值本质上也就是 $\frac{\partial L}{\partial \beta} = 0$ 方程. 牛顿法可用于方程根的数值求解

①. 牛顿法: $\beta^{t+1} = \beta^t - \left(\frac{\partial^2 L(\beta)}{\partial \beta \partial \beta^T} \right)^{-1} \frac{\partial L(\beta)}{\partial \beta}$

其中: $\frac{\partial L(\beta)}{\partial \beta} = - \sum_{i=1}^m \hat{x}_i (y_i - p_i(\hat{x}_i, \beta))$ $\frac{\partial^2 L(\beta)}{\partial \beta \partial \beta^T} = \sum_{i=1}^m \hat{x}_i \hat{x}_i^T p_i(\hat{x}_i, \beta) (1 - p_i(\hat{x}_i, \beta))$

②. 梯度下降法: $\beta_{in} = \beta_i - \eta \frac{\partial L}{\partial \beta}$ η 称为学习率.



补充: 梯度下降.

Gradient Descent is a common way to solve Optimization problem.

Target: $\min_{x \in \mathbb{R}^n} f(x)$. x^* is the best solution.

Actually, $f(x) = f(x^{(k)}) + \text{grad}_k^T (x - x^{(k)})$ ①.

$\text{grad}_k = \nabla f(x^{(k)})$, is called the gradient of $f(x)$.

\Rightarrow if determined a vector $P_k = \{P_1, P_2, \dots, P_n\}$ and a step λ .

① can be rewritten as: $x^{(k+1)} = x^{(k)} + \lambda P_k$.

Eh, we often say the step λ as learning rate, and use η instead:

and P_k , yes, as its name indicates, is associated with $\text{grad}(x^{(k)})$,

actually, $P_k = -\nabla f(x^{(k)})$. ($-\text{grad}(x^{(k)})$)

So we have the algorithm below:

(1). initialize $x^{(0)} \in \mathbb{R}^n$, $k=0$.

(2). compute $f(x^{(k)})$ & $\text{grad}(x^{(k)})$.

(3). if $\|\text{grad}(x^{(k)})\| \leq \varepsilon$, then break;

else: $\{P_k = -\text{grad}(x^{(k)})\}$, s.t. $\min_{\eta \geq 0} f(x^{(k)} + \eta P_k)$.

(4). calculate $x^{(k+1)} = x^{(k)} + \eta P_k$. then $f(x^{(k+1)})$.

(5). if $\|x^{(k+1)} - x^{(k)}\| < \varepsilon$ or $\|f(x^{(k+1)}) - f(x^{(k)})\| < \varepsilon$, break, and $x^* = x^{(k+1)}$.

else, $i++$, go to step (3) and recurrent.

Here, λ : learning rate; ε : error

