

# 读书报告

## 第三章 线性模型

刘精昌

September 27, 2016

# 引例

## 工资与教育水平关系

考查工人工资水平与其受教育关系：

- a 工资水平（每小时美元数）：用  $Y$  表示
- b 受教育程度（受教育年数）：用  $X$  表示
- c 非可观测因素，如工作经验、天生素质、工作时间等其他因素

# 引例

## 工资与教育水平关系

d 观测的数据:  $(x_i, y_i), i = 1, \dots, n$  (受访人数),

$i$	$x_i$	$y_i$	$i$	$x_i$	$y_i$
1	5.3	1.4	9	8.5	3.2
2	11.0	3.9	10	7.1	8.6
3	9	6.3	11	15	4
4	8.7	8.6	12	12.0	9.0
5	10	12	13	29	12
6	15.5	12	14	19.7	13.1
7	21	16	15	15.1	10
8	19	14.4	16	15.7	16

## 定义

$$y_i = \beta_0 + \mathbf{x}_{i1}\beta_1 + \cdots + \mathbf{x}_{i,p-1}\beta_{p-1} + \mathbf{e}_i, \quad i = 1, 2, \dots, n$$

$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} 1 & \mathbf{x}_{11} & \cdots & \mathbf{x}_{1,p-1} \\ 1 & \mathbf{x}_{21} & \cdots & \mathbf{x}_{2,p-1} \\ \vdots & \vdots & & \vdots \\ 1 & \mathbf{x}_{n1} & \cdots & \mathbf{x}_{n,p-1} \end{pmatrix} \begin{pmatrix} \omega_0 \\ \omega_1 \\ \vdots \\ \omega_{p-1} \end{pmatrix} + \begin{pmatrix} \mathbf{e}_1 \\ \mathbf{e}_2 \\ \vdots \\ \mathbf{e}_{p-1} \end{pmatrix}$$

$$\mathbf{y} = \mathbf{X}\beta + \mathbf{e}$$

$$y_i = \beta_0 + x_{i1}\beta_1 + \cdots + x_{i,p-1}\beta_{p-1} + \mathbf{e}_i, \quad i = 1, 2, \dots, n$$

(a) 拟合值:  $\hat{y}_i = \hat{\beta}_0 + x_{i1}\hat{\beta}_1 + \cdots + x_{i,p-1}\hat{\beta}_{p-1}$

(b) 残差 (residual) :  $\varepsilon_i = y_i - \hat{y}_i = y_i - \hat{\beta}_0 - x_{i1}\hat{\beta}_1 - \cdots - x_{i,p-1}\hat{\beta}_{p-1}$

(c) 残差平方和 (residual sum of squares RSS)

$$:RSS(\beta) = \sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \left( \hat{\mathbf{y}} - \mathbf{X}\hat{\beta} \right)^T \left( \hat{\mathbf{y}} - \mathbf{X}\hat{\beta} \right)$$

# 参数求解

最小二乘法：RSS 最小

$$1 \quad E(\hat{\beta}) = \left( y - X\hat{\beta} \right)^T \left( y - X\hat{\beta} \right) = y^T y - 2y^T X\hat{\beta} + \hat{\beta}^T X^T X \hat{\beta}$$

$$2 \quad \frac{\partial E(\hat{\beta})}{\partial \hat{\beta}} = 0$$

$$3 \quad X^T X \hat{\beta} = X^T y$$

$$4 \quad \hat{\beta} = (X^T X)^{-1} X^T y$$

## MLE

$$1 \quad y_i = \beta^T \mathbf{x}_i + \mathbf{e}_i, \mathbf{e}_i \sim N(0, \sigma^2)$$

$$2 \quad y_i \sim N(\beta^T \mathbf{x}_i, \sigma^2)$$

$$3 \quad \hat{\beta} \triangleq \arg \max_{\beta} \log p(D|\beta)$$

$$4 \quad l(\beta) \triangleq \log p(D|\beta) = \sum_{i=1}^n \log p(y_i|\beta) =$$

$$\sum_{i=1}^n \log \left[ \left( \frac{1}{2\pi\sigma^2} \right)^{\frac{1}{2}} \exp \left( -\frac{1}{2\sigma^2} (y_i - \beta^T \mathbf{x}_i)^2 \right) \right] =$$

$$-\frac{1}{2\sigma^2} \text{RSS} - \frac{N}{2} \log(2\pi\sigma^2)$$

# 平方和的定义

a、总平方和 (Total Sum of Squares TSS):

$$TSS \triangleq \sum_{i=1}^n (y_i - \bar{y})^2, \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

b、解释平方和 (Explained Sum of Squares ESS) :

$$ESS \triangleq \sum_{i=1}^n \left( \hat{y}_i - \bar{\hat{y}} \right)^2, \bar{\hat{y}} = \frac{1}{n} \sum_{i=1}^n \hat{y}_i$$

**Theorem**  $TSS = ESS + RSS$



# 拟合优度

## 判定系数的定义

$$R^2 = \frac{ESS}{TSS}$$

## 判定系数的性质

- $0 \leq R^2 \leq 1$
- $R^2$  越接近 1, 拟合效果越好
- $R^2$  越接近 0, 拟合效果越差

# 几何解释

# 增加惩罚项

## ridge regression

- $J(\beta) = \frac{1}{n}RSS + \lambda \|\beta\|_2^2$
- $\hat{\beta}_{ridge} = \min_{\beta} \arg J(\beta)$
- $\hat{\beta}_{ridge} = (\lambda I + X^T X)^{-1} X^T y$

## LASSO

- $J(\beta) = \frac{1}{n}RSS + \lambda \|\beta\|_1$

## 定义

$$p(y|x, w) = \text{Ber}(y | \text{sigm}(w^T x)) = \begin{cases} \text{sigm}(w^T x), y = 1 \\ 1 - \text{sigm}(w^T x), y = 0 \end{cases}$$

$$\text{sigm}(w^T x) = \frac{1}{1 + e^{w^T x}}$$

## MLE

$$\mu_i = \text{sigm}(w^T x_i) = \frac{1}{1 + e^{w^T x_i}}$$

$$\begin{aligned} NLL(w) &= - \sum_{i=1}^N \log \left( \mu_i^{I(y_i=1)} \times (1 - \mu_i)^{I(y_i=0)} \right) \\ &= - \sum_{i=1}^N [y_i \log \mu_i + (1 - y_i) \log (1 - \mu_i)] \end{aligned}$$

## MLE

$$\frac{\partial \mu_i}{\partial \mathbf{w}} = -\frac{\mathbf{x}_i \mathbf{e}^{\mathbf{w}^T \mathbf{x}_i}}{(1 + \mathbf{e}^{\mathbf{w}^T \mathbf{x}_i})^2} = -\mathbf{x}_i \mu_i (1 - \mu_i)$$

$$\begin{aligned} \frac{\partial NLL(\mathbf{w})}{\partial \mathbf{w}} &= \sum_{i=1}^N \frac{\partial NLL(\mathbf{w})}{\partial \mu_i} \frac{\partial \mu_i}{\partial \mathbf{w}} \\ &= \sum_{i=1}^N \left( \frac{y_i}{\mu_i} - \frac{1 - y_i}{1 - \mu_i} \right) \mathbf{x}_i \mu_i (1 - \mu_i) \\ &= \sum_{i=1}^N (\mathbf{x}_i y_i - \mathbf{x}_i \mu_i) \\ &= \mathbf{X}^T (\boldsymbol{\mu} - \mathbf{y}) \end{aligned}$$

$$\mathbf{H} = \frac{\partial^2 NLL(\mathbf{w})}{\partial \mathbf{w}^2} = \mathbf{X}^T \mathbf{S} \mathbf{X}, \mathbf{S} \triangleq \text{diag}(\mu_i (1 - \mu_i))$$

# Optimization

## Steepest descent

$$\theta_{k+1} = \theta_k - \eta_k \mathbf{g}_k$$

## Newton's method

$$\theta_{k+1} = \theta_k - \eta_k \mathbf{H}_k^{-1} \mathbf{g}_k$$

## softmax

$$\left\{ \left( \mathbf{x}^{(1)}, y^{(1)} \right), \left( \mathbf{x}^{(2)}, y^{(2)} \right), \dots, \left( \mathbf{x}^{(n)}, y^{(n)} \right) \right\}, y^{(i)} \in \{1, 2, \dots, k\}$$

$$\begin{bmatrix} p(y^{(i)} = 1 | \mathbf{x}^{(i)}; \theta) \\ p(y^{(i)} = 2 | \mathbf{x}^{(i)}; \theta) \\ \vdots \\ p(y^{(i)} = k | \mathbf{x}^{(i)}; \theta) \end{bmatrix} = \frac{1}{\sum_{j=1}^k \mathbf{e}^{w_j^T \mathbf{x}^{(i)}}} \begin{bmatrix} \mathbf{e}^{w_1^T \mathbf{x}^{(i)}} \\ \mathbf{e}^{w_2^T \mathbf{x}^{(i)}} \\ \vdots \\ \mathbf{e}^{w_k^T \mathbf{x}^{(i)}} \end{bmatrix}$$



# LDA

$$J(w) = \frac{w^T S_B w}{w^T S_W w}$$

## ■ between-class scatter matrix

$$S_B = (\mu_2 - \mu_1)(\mu_2 - \mu_1)^T$$

## ■ within-class scatter matrix

$$S_W = \sum_{i:y_i=1} (x_i - \mu_1)(x_i - \mu_1)^T + \sum_{i:y_i=2} (x_i - \mu_2)(x_i - \mu_2)^T$$

## 求解

$$1 \quad J'(w) = w^T S_B w - \lambda w^T S_W w, \lambda > 0$$

$$2 \quad \frac{dJ'(w)}{dw} = 0$$

$$3 \quad \lambda S_W w = S_B w$$

$$4 \quad S_B w = (\mu_2 - \mu_1)(\mu_2 - \mu_1)^T w = (\mu_2 - \mu_1)(m_2 - m_1)$$

$$5 \quad w \propto S_W^{-1}(\mu_2 - \mu_1)$$

## Q & A