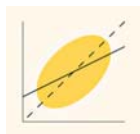


## 第八讲 线性回归模型



### 回归函数

$y$ : 响应变量

$\mathbf{x}$ : 自变量, 1维或多维

以 $\mathbf{x}$ 的某个函数 $f(\mathbf{x})$ 逼近 $y$ , 误差为 $E(y - f(\mathbf{x}))^2$   
最优逼近是 $f(\mathbf{x}) = E(y | \mathbf{x})$ , 这是因为  
$$E(y - f(\mathbf{x}))^2 \geq E(y - E(y | \mathbf{x}))^2$$

$E(y | \mathbf{x})$ 称为回归函数. 令  $\varepsilon = y - E(y | \mathbf{x})$ , 则容易验证 $\varepsilon$ 与 $\mathbf{x}$ 不相关。  
所以响应变量 $y$ 可以表示为:

$$y = E(y | \mathbf{x}) + \varepsilon.$$

对回归函数假设结构, 即得到各种回归模型。

### 一般线性回归模型

模型表示1:

$y$ : 响应变量,  $\mathbf{x}$ : 自变量(向量)。线性回归模型假设:

- (i) 线性回归函数:  $E(y | \mathbf{x}) = a + \mathbf{b}'\mathbf{x}$ ,
- (ii) 方差常数/齐性:  $\text{var}(y | \mathbf{x}) = \sigma^2$ .

其中 $a, \mathbf{b}, \sigma^2$ 是未知参数。

- 自变量 $\mathbf{x}$ 一维时, 称为简单线性回归模型 (simple linear regression model)
- 自变量 $\mathbf{x}$ 多维时, 称为多重线性回归模型 (multiple linear regression model)

令  $\varepsilon = y - E(y | \mathbf{x}) = y - (a + \mathbf{b}'\mathbf{x})$ ,

若要求 $\varepsilon$ 与 $\mathbf{x}$ 独立, 则有"等价的"/稍强的模型表示:

模型表示2:

$$y = a + \mathbf{b}'\mathbf{x} + \varepsilon,$$

其中

- (1)  $E(\varepsilon) = 0$
- (2)  $\text{var}(\varepsilon) = \sigma^2$  (方差齐性, Homoscedasticity)
- (3)  $\varepsilon$  与  $\mathbf{x}$  独立 (外生性, Exogeneity)

条件(1),(2)加上样本独立性, 称为Gauss - Markov假设

注:

(1) 误差0均值(均值常数,不依赖于x):  
保证可识别性(Identifiability),截距参数 $a$ 可识别

(2) 方差齐性(方差常数,不依赖于x):  
保证参数估计的最优性(Gauss - Markov定理)

(3) 外生性:保证参数估计的无偏性/因果

### 为什么线性?

1. 如果 $x, y$ 服从联合正态分布:  $\begin{pmatrix} y \\ x \end{pmatrix} \sim N\left(\begin{pmatrix} \mu_y \\ \mu_x \end{pmatrix}, \begin{pmatrix} \Sigma_{yy} & \Sigma_{yx} \\ \Sigma_{xy} & \Sigma_{xx} \end{pmatrix}\right)$ , 则

$$y|x \sim N(\mu_y + \Sigma_{yx}\Sigma_{xx}^{-1}(x - \mu_x), \Sigma_{yy} - \Sigma_{yx}\Sigma_{xx}^{-1}\Sigma_{xy}),$$

$$(a) E(y|x) = \mu_y + \Sigma_{yx}\Sigma_{xx}^{-1}(x - \mu_x) \triangleq a + b'x,$$

$$(b) \text{var}(y|x) = \Sigma_{yy} - \Sigma_{yx}\Sigma_{xx}^{-1}\Sigma_{xy} \triangleq \sigma^2.$$

或等价地(令  $\varepsilon = y - E(y|x) = y - (a + b'x)$ )

$$y = a + b'x + \varepsilon,$$

其中  $\varepsilon \sim N(0, \sigma^2)$ , 且与 $x$ 独立.

2. 如果 $x(=0,1)$ 是因子变量,代表随机化分组,  
而每一组内 $y$ 服从正态分布:

$$y|_{x=1} \sim N(\mu_1, \sigma^2), \quad y|_{x=0} \sim N(\mu_0, \sigma^2)$$

则可以表示为:

$$y|x \sim N(\mu_0 + (\mu_1 - \mu_0)x, \sigma^2) \triangleq N(a + bx, \sigma^2), x = 0, 1$$

或  $y = a + bx + \varepsilon$ ,  $\varepsilon \sim N(0, \sigma^2)$ , 且与 $x$ 独立.

### 模型的数据形式

数据:  $(y_i, \mathbf{x}_i)$ ,  $i = 1, 2, \dots, n$  独立, 其中

向量 $\mathbf{x}_i$ 为自变量的第 $i$ 个观察值(第一个分量为1, 对应于截距项);  
 $y_i$ 为响应变量,  $i = 1, 2, \dots, n$ .

假设 $(y_i, \mathbf{x}_i)$ ,  $i = 1, 2, \dots, n$ , 独立, 满足线性回归模型:

$$y = \mathbf{x}'\beta + \varepsilon, \quad \varepsilon \sim (0, \sigma^2), \quad \varepsilon \text{ 与 } \mathbf{x} \text{ 独立}$$

其中 $\beta$ 为  $p \times 1$  回归系数列向量(第一个分量是截距项). 即

$$y_i = \mathbf{x}_i'\beta + \varepsilon_i, i = 1, 2, \dots, n$$

其中 $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$  iid  $\sim (0, \sigma^2)$ , 且 $\varepsilon_i$ 与 $\mathbf{x}_i$ 独立.

记

$$\mathbf{Y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} \mathbf{x}_1' \\ \mathbf{x}_2' \\ \vdots \\ \mathbf{x}_n' \end{pmatrix}, \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{pmatrix}, \quad \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix},$$

其中  $\mathbf{X}$  称为设计阵(第  $i$  行为  $\mathbf{x}_i'$ )。

模型等价地以矩阵-向量形式表示为:

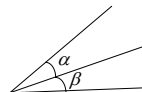
$$\mathbf{Y}_{n \times 1} = \mathbf{X}_{n \times p} \boldsymbol{\beta}_{p \times 1} + \boldsymbol{\varepsilon}_{n \times 1}, \quad \boldsymbol{\varepsilon} \sim (\mathbf{0}, \sigma^2 \mathbf{I}_n), \quad \boldsymbol{\varepsilon} \text{ 与 } \mathbf{X} \text{ 独立}.$$

或以前两阶矩表示为:

$$(a) \quad E(\mathbf{Y}|\mathbf{X}) = \mathbf{X}\boldsymbol{\beta},$$

$$(b) \quad \text{var}(\mathbf{Y}|\mathbf{X}) = \sigma^2 \mathbf{I}_n$$

例1. 一个航海员使用六分仪测量下图中的角度  $\alpha, \beta$ , 他共测量了3次, 第一、二次分别测量  $\alpha, \beta$ , 得到测量值  $y_1, y_2$  第三、四次测量角度  $\alpha + \beta$  得到测量值  $y_3, y_4$ , 假设四次测量独立, (可加)误差均值为0, 方差为  $\sigma^2$ , 写出线性模型



$$y_1 = \alpha + \varepsilon_1$$

$$y_2 = \beta + \varepsilon_2$$

$$y_3 = \alpha + \beta + \varepsilon_3$$

$$y_4 = \alpha + \beta + \varepsilon_4$$

$$\varepsilon_i \text{ 独立 } \sim (0, \sigma^2)$$

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad E\boldsymbol{\varepsilon} = \mathbf{0}, \quad \text{var}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I}_4$$

$$\mathbf{Y} = \begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \\ 1 & 1 \end{pmatrix}, \quad \boldsymbol{\beta} = \begin{pmatrix} \alpha \\ \beta \end{pmatrix}, \quad \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \varepsilon_3 \\ \varepsilon_4 \end{pmatrix}.$$

例2(单因素方差分析模型).  $n$  个个体被随机地分配于  $K$  组分别接受某种处理(*treatment*), 并测量某指标。假设第  $k$  组有  $n_k$  个个体, 指标测量为  $y_{k1}, \dots, y_{kn_k} \text{ iid } \sim N(\mu_k, \sigma^2), k = 1, \dots, K$

记  $\varepsilon_{kj} = y_{kj} - \mu_k, j = 1, \dots, n_k \text{ iid } \sim N(0, \sigma^2), k = 1, \dots, K$ , 即

$$y_{kj} = \mu_k + \varepsilon_{kj},$$

模型可表示为:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, E\boldsymbol{\varepsilon} = \mathbf{0}, \text{var}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I}_n$ , 其中

$$\mathbf{Y} = \begin{pmatrix} y_{11} \\ \vdots \\ y_{1n_1} \\ \vdots \\ y_{K1} \\ \vdots \\ y_{Kn_K} \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} 1 & 0 & \cdots & 0 \\ \vdots & \vdots & \cdots & \vdots \\ 1 & 0 & \cdots & 0 \\ \vdots & \vdots & \cdots & \vdots \\ 0 & 0 & \cdots & 1 \\ \vdots & \vdots & \cdots & \vdots \\ 0 & 0 & \cdots & 1 \end{pmatrix}, \quad \boldsymbol{\beta} = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_K \end{pmatrix}, \quad \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_{11} \\ \vdots \\ \varepsilon_{1n_1} \\ \vdots \\ \varepsilon_{K1} \\ \vdots \\ \varepsilon_{Kn_K} \end{pmatrix}$$

为了检验  $K$  个均值相同, 通常重新参数化, 比如令

$$\beta_1 = \mu_1, \quad \beta_2 = \mu_2 - \mu_1, \quad \dots, \quad \beta_K = \mu_K - \mu_1,$$

对于上述新参数  $\tilde{\boldsymbol{\beta}} = (\beta_1, \dots, \beta_K)'$ , 写出模型

$$\mathbf{Y} = \mathbf{X}\tilde{\boldsymbol{\beta}} + \boldsymbol{\varepsilon}, \quad E\boldsymbol{\varepsilon} = \mathbf{0}, \quad \text{var}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I}_n$$

中的设计阵  $\mathbf{X}$