

8 Sept

Chaptr 1 Simple & multiple linear regression

$$\underbrace{Y}_{m \times 1} \quad \underbrace{X}_{p \times 1} \quad (x_1, \dots, x_p)$$

Model $\underbrace{Y}_{m \times 1}$ on $\underbrace{X}_{p \times 1}$

- response
- outcome
- dependent variable
- predictor
- covariate
- independent variable

$$\underbrace{Y}_{m \times 1} \begin{cases} m=1 & \text{MATH 3423} \\ m>1 & \text{MATH 4424} \end{cases}$$

ultivariate \wedge Analysis
statistical

$$\underbrace{X}_{p \times 1} \begin{cases} p=1 & \text{simple} \\ p>1 & \text{multiple} \end{cases}$$

linear $f(x) = \alpha + \beta x \quad \text{OR} \quad \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$

\uparrow linear in parameter, no linear in x

e.g. $f(x) = \alpha + \beta \cdot \underbrace{x^2}_{x'}$

$f(x) = \alpha * \beta^x \quad \leftarrow$ non-linear

$\log f(x) = \log \alpha + x \log \beta$

linear regression $y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p + \underbrace{e}_{\text{r.v.}}$

- \Rightarrow ① y & (x_1, \dots, x_p) not perfectly linear ~~reg~~ related
- ② y is observed with error

①

$$y_{obs} = y_{true} + e_y \leftarrow \text{classical measurement error model}$$

$$y_{true} = \beta_0 + \beta_1 X + \epsilon_y$$

(y_{true} & X) are not perfectly linear related

Linear regression model

$$y_i = \beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{ip} + \epsilon_i \quad i=1, \dots, n$$

Assume

Purpose

- do prediction

- find the ~~important~~ important factors among X
(variable selection)

X - fixed \checkmark ~~fixed~~ $y|X \sim N(\beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{ip}, \sigma^2)$

X - random e.g. $X \sim N(\mu_x, \sigma_x^2)$

$$y = \beta_0 + \beta_1 X + \epsilon \Rightarrow y \sim N(\beta_0 + \beta_1 X, \sigma^2)$$

$\epsilon \sim N(0, \sigma^2)$

$\beta_0, \beta_1, \sigma^2 \Rightarrow$ dist. of (y, X)

\Rightarrow bivariate normal

$$\begin{pmatrix} y \\ X \end{pmatrix} \sim N \left(\begin{pmatrix} \beta_0 + \beta_1 \mu_x \\ \mu_x \end{pmatrix}, \begin{pmatrix} \beta_1^2 \sigma_x^2 + \sigma^2 & \beta_1 \sigma_x^2 \\ \beta_1 \sigma_x^2 & \sigma_x^2 \end{pmatrix} \right)$$

$\text{Var}(y)$ $\text{Cov}(X, Y)$ $\text{Var}(X)$

$$\Rightarrow y|X \sim N \left(\beta_0 + \beta_1 \mu_x + \frac{\beta_1 \sigma_x^2}{\sigma_x^2} (x - \mu_x), \beta_1^2 \sigma_x^2 + \sigma^2 - \frac{\beta_1^2 \sigma_x^4}{\sigma_x^2} \right)$$

$\sim N(\beta_0 + \beta_1 x, \sigma^2)$

Except If X - random & measured with error

e.g. $X_{obs} = X + \epsilon_x$

~~Write down the joint distribution~~

~~of y & $X \Rightarrow y|X_{obs}$~~

$$f(y|x) = y|x \sim N(\beta_0 + \beta_1 x, \sigma^2)$$

$$x \sim N(\mu_x, \sigma_x^2)$$

$$X_{obs} = X + \Sigma_x$$

$$\Rightarrow y | X_{obs} \Rightarrow \text{likelihood factor} \\ \Rightarrow \text{m.l.e.}$$

linear regression model

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + e_i \quad i=1, \dots, n$$

fixed

Assume

$$\begin{aligned} & \sim N(0, \sigma^2) \\ & \text{cov}(e_i, e_j) = 0 \quad i \neq j \\ & (\text{obs. are indep}) \end{aligned}$$

$$\begin{aligned} \underset{n \times 1}{\tilde{y}} &= \underset{n \times (p+1)}{\tilde{X}} \underset{(p+1) \times 1}{\tilde{\beta}} + \underset{n \times 1}{\tilde{e}} \\ \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} &= \begin{pmatrix} 1 & x_{11} & \dots & x_{1p} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} + \begin{pmatrix} e_1 \\ \vdots \\ e_n \end{pmatrix} \end{aligned}$$

intercept

regression coefficient of x_{ij} $j=1, \dots, p$

column with all elements equal to 1

$$\tilde{e} \sim MN(\underline{0}, \begin{pmatrix} \sigma^2 & & 0 \\ & \ddots & \\ 0 & & \sigma^2 \end{pmatrix})$$

$\sigma^2 I$

MATH 3423 $y_i \overset{\text{iid}}{\sim} N(\underline{\mu}, \sigma^2) \leftarrow$

$$E\left(\frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n-1}\right) = \sigma^2$$

MATH 3424 $y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + e_i$

\leftarrow diff. mean $\sim N(0, \sigma^2)$

$$y_i \approx N(\beta_0 + \beta_1 x_{i1}, \sigma^2) \leftarrow$$

Is $\frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n-1}$ unbiased est. of σ^2 ?

\Rightarrow Is $E\left(\frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n-1}\right)$ equal to σ^2 ? **NO** (3)

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \epsilon_i \sim N(0, \sigma^2)$$

$$E(y_i) = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} \quad E(\epsilon_i) = 0$$

$$\uparrow$$

$$\mu_i = \tilde{x}_i^T \beta$$

Define $g(\mu) = \tilde{x}^T \beta$

\uparrow link function

linear ① $g(\mu) = \mu$ — identity function

② $y \sim$ binary data Y/N

Assume $f(y|x) \sim \text{Bernoulli}(P)$

$$E(y) = P \quad \leftarrow \quad \mu = P$$

$$0 < P < 1$$

$$P = \frac{\exp(\tilde{x}^T \beta)}{1 + \exp(\tilde{x}^T \beta)}$$

$$\Rightarrow \frac{P}{1-P} = \exp(\tilde{x}^T \beta)$$

$$\Rightarrow \boxed{\ln\left(\frac{P}{1-P}\right)} = \tilde{x}^T \beta$$

$$\uparrow g(\mu) = \ln\left(\frac{\mu}{1-\mu}\right) \text{ — logit function}$$

③ y — count data e.g.

$$\sim \text{Poisson}(\mu)$$

$$E(y) = \mu$$

$$\mu = \exp(\tilde{x}^T \beta)$$

$$\Rightarrow \log_e(\mu) = \tilde{x}^T \beta$$

$$g(\mu) = \ln(\mu) \text{ — link function}$$

Section 2 Estimation

method of maximum likelihood

Unknown parameters: $\beta_0, \beta_1, \dots, \beta_p, \sigma^2$

Model = $y_i = \beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{ip} + e_i$

Assume

$N(0, \sigma^2)$

Method of least squares

$y_i = \beta_0 + \beta_1 X_{i1} + e_i \quad i=1, \dots, n$

~~unobservable~~

unobservable

$\hat{e}_i = y_i - (\hat{\beta}_0 + \hat{\beta}_1 X_{i1})$ — residual

observable

\hat{y}_i — fitted value of y_i

Method of least squares

Find $\hat{\beta}_0, \hat{\beta}_1$ such that

$\sum_{i=1}^n (y_i - \hat{y}_i)^2$

or $\sum_{i=1}^n (y_i - (\hat{\beta}_0 + \hat{\beta}_1 X_{i1}))^2$

is minimized

Residual ~~SS~~ Sum of Squares = Res. S. S.

