$$C_p = 2p' - n + \frac{Res\,S.S.|p'}{\hat{\sigma}^2\,\text{full model}}$$

1. Dec

3. (15 marks)   An experiment was conducted to model $Y$ with five explanatory variables $X_1$, $X_2$, $X_3$, $X_4$ and $X_5$. We desire an equation of relating $Y$ to the other variables. The goal is to find variables that should be further studied with the eventual goal of developing a prediction equation. The following table gives RSS for all possible regressions. Total sum of squares is equal to 5.0634 and the number of observations is equal to 20.

| No. of parameters in the model | RSS | Model |
|---|---|---|
| 2 | 2.0338 | $X_1$ |
| 2 | 5.0219 | $X_2$ |
| 2 | 1.5370 | $X_3$ |
| 2 | 2.5044 | $X_4$ |
| 2 | 1.5563 | $X_5$ |
| 3 | 1.5921 | $X_1, X_2$ |
| 3 | 1.4397 | $X_1, X_3$ |
| 3 | 1.7462 | $X_1, X_4$ |
| 3 | 1.4963 | $X_1, X_5$ |
| 3 | 1.4707 | $X_2, X_3$ |
| 3 | 2.4381 | $X_2, X_4$ |
| 3 | 1.4388 | $X_2, X_5$ |
| 3 | 1.4590 | $X_3, X_4$ |
| 3 | 1.0850 | $X_3, X_5$ |
| 3 | 1.3287 | $X_4, X_5$ |
| 4 | 1.2582 | $X_1, X_2, X_3$ |
| 4 | 1.4257 | $X_1, X_2, X_4$ |
| 4 | 1.2764 | $X_1, X_2, X_5$ |
| 4 | 1.3894 | $X_1, X_3, X_4$ |
| 4 | 1.0644 | $X_1, X_3, X_5$ |
| 4 | 1.3204 | $X_1, X_4, X_5$ |
| 4 | 1.3900 | $X_2, X_3, X_4$ |
| 4 | 0.9871 | $X_2, X_3, X_5$ |
| 4 | 1.2178 | $X_2, X_4, X_5$ |
| 4 | 1.0634 | $X_3, X_4, X_5$ |
| 5 | 1.2199 | $X_1, X_2, X_3, X_4$ |
| 5 | 0.9871 | $X_1, X_2, X_3, X_5$ |
| 5 | 1.1565 | $X_1, X_2, X_4, X_5$ |
| 5 | 1.0388 | $X_1, X_3, X_4, X_5$ |
| 5 | 0.9653 | $X_2, X_3, X_4, X_5$ |
| 6 | 0.9652 | $X_1, X_2, X_3, X_4, X_5$ |

$$\Rightarrow C_p = 2*2 - 20 + \frac{1.5370}{0.9652/14} = 6.294$$

$$\Rightarrow C_p = 2*3 - 20 + \frac{1.085}{0.9652/14} = \boxed{1.738}$$

Best model $= X_3, X_5$

$$\Rightarrow C_p = 2*4 - 20 + \frac{0.9871}{0.9652/14} = 2.318$$

$$\Rightarrow C_p = 2*5 - 20 + \frac{0.9653}{0.9652/14} = 4.001$$
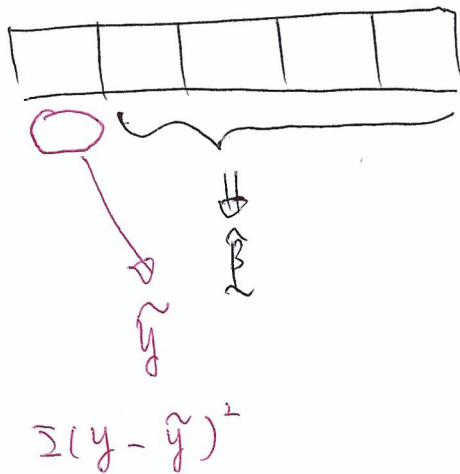
$$\Rightarrow C_p = 2p' - n + \frac{0.9652}{0.9652/(n-p)} = p' = 6$$

Find the best model by $C_p$, forward selection, backward selection and stepwise selection. Write down how to get the best model on details. Choose critical values for both ENTRY and STAY to be 2. Comment the results.
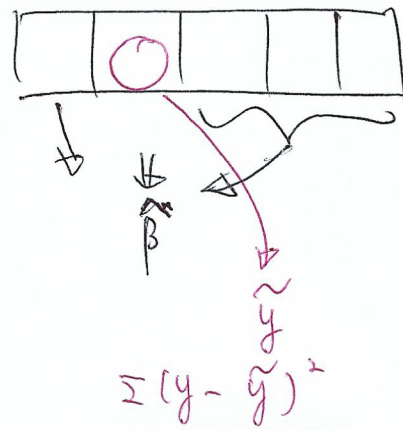
$$\hat{\sigma}^2\,\text{full model} = \frac{0.9652}{20 - 6^{P.3}} = \frac{0.9652}{14} = 0.06894$$

①

(6) Cross-validation

① 

$$\Sigma(y-\tilde{y})^2$$

② 

$$\Sigma(y-\tilde{y})^2$$

Calculate $\Sigma\Sigma(y-\tilde{y})^2$

CV PRESS

PRESS statistic

↑ Predicted

all obs. without $i$th obs.

$n_1 = n-1 \rightarrow$ fit a model

$\Rightarrow \hat{\beta}$

$n_2 = 1 \rightarrow \tilde{y}_{i(-i)}$   $i = 1, \cdots, n$

↑ $i$th obs

$$\left( y_i - \tilde{y}_{i(-i)} \right)$$

PRESS residual

$$PRESS = \sum_{i=1}^{n} \left( y_i - \tilde{y}_{i(-i)} \right)^2$$

$$= \sum_{i=1}^{n} \frac{\hat{e}_i^2}{(1-h_{ii})^2}$$

$\hat{e}_i = y_i - \hat{y}_i$ — residual

$h_{ii} = (i,i)^{th}$ element in $H$

where $H = X(X^TX)^{-1}X^T$

$h_{ii} = x_i^T (X^TX)^{-1} x_i$

↗ leverage

$x_i^T$ — $i$th row in $X$

②

Simple linear regression

$$y_i = \beta_0 + \beta_1 x_i + e_i \qquad i = 1, \dots, n$$

Centred model

$$\rightarrow X^T X = \begin{pmatrix} n & 0 \\ 0 & S_{xx} \end{pmatrix} \qquad (X^T X)^{-1} = \begin{pmatrix} \frac{1}{n} & 0 \\ 0 & \frac{1}{S_{xx}} \end{pmatrix}$$

$$h_{ii} = \begin{pmatrix} 1 & x_i - \bar{x} \end{pmatrix} \begin{pmatrix} \frac{1}{n} & 0 \\ 0 & \frac{1}{S_{xx}} \end{pmatrix} \begin{pmatrix} 1 \\ x_i - \bar{x} \end{pmatrix}$$

$$= \frac{1}{n} + \boxed{\frac{(x_i - \bar{x})^2}{S_{xx}}} \longleftarrow \text{diff between } x_i \text{ & } \bar{x}$$

$$\boxed{1} \geq h_{ii} \geq \boxed{\frac{1}{n}} \longleftarrow H - \frac{1}{n} J \text{ is positive semi definite}$$

$$\uparrow$$

$$I - H \text{ is positive semi definite (Chapter 1)}$$

$$\Rightarrow \text{all } \cancel{\text{diagon}} \text{ diagonal elements are non-negative}$$

$$\Rightarrow 1 - h_{ii} \geq 0$$

$$\Rightarrow 1 \geq h_{ii}$$

Define $H_c = X_c (X_c^T X_c)^{-1} X_c^T$

$$\text{Model} = y = \alpha 1 + X_c \beta + e$$

$$\hat{y} = \hat{\alpha} 1 + X_c \hat{\beta}$$

$$= \hat{\alpha} 1 + X_c (X_c^T X_c)^{-1} X_c^T y$$

$$= [\frac{1}{n} J + H_c] y$$

$$\uparrow$$
matrix with $\cancel{\text{eq}}$ all elements equal to 1

$$= H y$$

(3)

$$\Rightarrow \underset{\sim}{H} = \frac{1}{n} \underset{\sim}{J} + \underset{\sim}{H_c}$$

$$\Rightarrow \underset{\sim}{H} - \frac{1}{n} \underset{\sim}{J} = \underset{\sim}{H_c}$$

$\uparrow$ positive $\overset{semi}{\text{definite}}$

$\Rightarrow$ all diagonal elements of $\underset{\sim}{H} - \frac{1}{n} \underset{\sim}{J}$ are non-negative

$$\Rightarrow h_{ii} - \frac{1}{n} \geq 0$$

$$\Rightarrow h_{ii} \geq \frac{1}{n}$$

$$PRESS = \sum_{i=1}^{n} \frac{\hat{e}_i^2}{(1-h_{ii})^2} \quad - \text{smallest } PRESS$$

(c) largest $R^2 = 1 - \frac{Res\, S.S.}{total\, SS.}$

$\Leftrightarrow$ smallest $Res\, S.S.$

$R^2 \uparrow$ when # of indep. variables $\uparrow$

$\Rightarrow$ Best model = full model

(d) largest $R^2_{adj} = 1 - \dfrac{MSE}{total\, MS.} \overset{\longleftarrow}{\underset{\longleftarrow}{\begin{array}{l} Res\,S.S./(n-p') \\ total\, S.S./(n-1) \end{array}}}$

$$= 1 - \frac{Res\, S.S.}{total\, S.S.} * \frac{n-1}{n-p'}$$

$R^2_{adj}$ may be negative

(f) $AIC$ = Akaike Information Criteria — smallest

$$= -2 \log L + \underbrace{2 p'}_{penalty}$$

(g) $BIC$ = Bayesian Information Criteria — smallest

$$= -2 \log L + \underbrace{\ln n * p'}_{penalty} \qquad \log_e(n) > 2$$

$\textcircled{4}$

# Linear Regression

$$L(\hat{\beta}, \hat{\sigma}) = \frac{1}{(2\pi)^{n/2} (\hat{\sigma}^2)^{n/2}} \exp\left\{ - \frac{\sum_{i=1}^{n}(y_i - \underset{\sim}{x_i^T} \hat{\beta})^2}{2\hat{\sigma}^2} \right\}$$

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^{n}(y_i - \underset{\sim}{x_i^T} \hat{\beta})^2}{\boxed{n}} \quad \text{m.l.e}$$

$$\Rightarrow \log L = -\frac{n}{2}\ln(2\pi) - \frac{n}{2}\ln(\hat{\sigma}^2) - \frac{n}{2}$$

$$\Rightarrow -2\log L = n\ln(2\pi) + n\ln(\hat{\sigma}^2) + n$$

$$\hat{\sigma}^2 = \frac{Res.S.S.}{n-p'} \quad \text{to est.}$$

$$= n\ln\left(\frac{Res.S.S.}{n}\right) + \text{constant}$$

$$\sigma^2 \ \text{(unbiased est.)}$$

$$\Rightarrow AIC = n\ln\left(\frac{Res.S.S.}{n}\right) + 2p'$$

$$BIC = n\ln\left(\frac{Res.S.S.}{n}\right) + \ln(n) * p'$$

# Chapter 5   Residual Analysis

**Assume**  $- e_i \sim N(0, \sigma^2)$

   $- y \ \& \ \underset{\sim}{x}$  are linear related

**Residual**   $\hat{e}_i = y_i - \hat{y}_i$    $E(\hat{e}_i) = 0$

$$Var(\hat{e}_i) = \sigma^2(1 - h_{ii}) \quad \left.\begin{array}{c}\\ \\ \end{array}\right\} \leftarrow \text{Chapter 1}$$

$$\neq \text{constant}$$

**Standardized residual**

① $\quad r_i = \dfrac{\hat{e}_i}{\sqrt{\hat{\sigma}^2(1-h_{ii})}}$  — (internally studentized residual)

② $\quad t_i = \dfrac{\hat{e}_i}{\boxed{\hat{\sigma}_{(-i)}}\sqrt{1-h_{ii}}}$  — (externally studentized residual)

$$\underset{\sim}{x}_1, \ \text{----} \ \underset{\sim}{x}_{i-1}, \ \boxed{\underset{\sim}{x}_i}, \ \underset{\sim}{x}_{i+1} \ \text{----}, \ \underset{\sim}{x}_n$$

$$y_1, \ \text{----} \ y_{i-1}, \ \boxed{y_i}, \ y_{i+1} \ \text{----}, \ y_n$$

<span style="color:green">↑ ignore the ith obs</span>

$\Rightarrow (n-1)$ obs $\Rightarrow$ Fit a model of $y$ on $\underset{\sim}{x}$

$$\hat{\beta}_{0(-i)}, \ \hat{\beta}_{1(-i)}, \ \text{----}, \ \hat{\beta}_{p(-i)}, \ \hat{\sigma}^2_{(-i)}$$

Put $\underset{\sim}{x}_i$ into the fitted model

$$\Rightarrow \hat{y}_{i(-i)}$$

$$\frac{\text{Res S.S.} (-i)}{(n-1) - p'}$$

$\underbrace{\qquad}_{\text{\# of obs.}}$ $\underbrace{\qquad}_{\substack{\text{\# of unknown} \\ \text{para. in} \\ \text{the model}}}$

$$t_i \sim t_{(n-1-p')} \ ?$$

Use $t_i$ to detect <u>outlier</u>  ① $\hat{e}_i, r_i, t_i$

<u>Residual Plot</u>

Check: linearity

constant variance

pattern $\Rightarrow$ transformation of $y$

& $\underset{\sim}{x}$

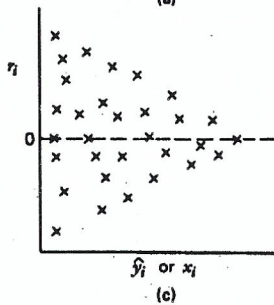

$\hat{e}_i, r_i, t_i$ ↑ ⟶ $\hat{y}_i$

※ If there is pattern

$\hat{e}_i, r_i, t_i$ ↑ ⟶ $x_{ij}$

$\hat{e}_i, r_i, t_i$

no pattern $\Rightarrow$ Assumptions are valid

variance (variation of residuals) increases as $\hat{y}_i$, or $x_i$ increases
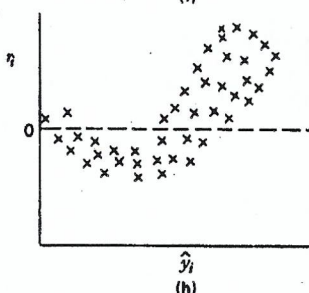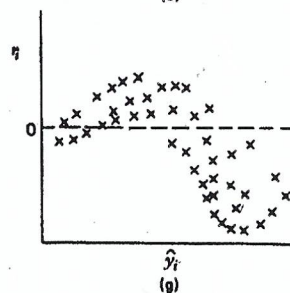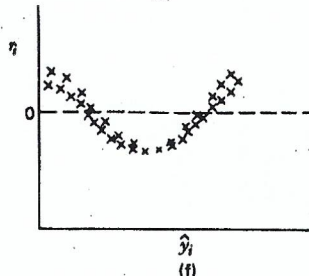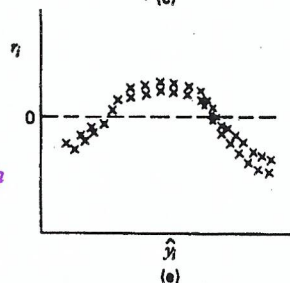
non-linear $\longrightarrow$

**Figure 6.3** Residual plots: (a) null plot; (b) right-opening megaphone; (c) left-opening megaphone; (d) double outward bow; (e) nonlinearity; (f) nonlinearity; (g) nonlinearity and nonconstant variance; (h) nonlinearity and nonconstant variance.

$\Rightarrow$ transformation of $y$ and/or $x$

e.g. count data $\Rightarrow \sim$ Poisson $(\mu)$    $E(y) = \mu$

$Var(y) = \mu$    smallest $-2 \log L$
largest likelihood

$-\sqrt{y}, \log(y)$

$y$ — area    $m^2$ $\boxed{\begin{array}{ccc} \sqrt{y} & \log(y) & y \\ x & \log(x) & x^2 \end{array}}$    $\boxed{\text{smallest Res S.S}}$
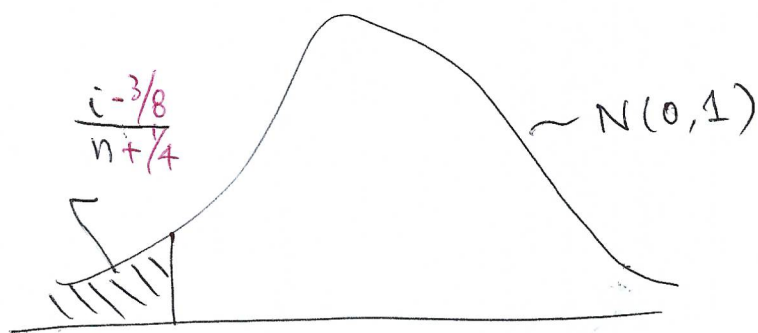
$x$ — perimeter    $m$    $\uparrow$ linear regression

⑦

## Normality assumption

Q-Q plot $\Rightarrow$ ~~normal~~ normality

Arrange $t_1, \ldots, t_n$

$\Rightarrow t_{(1)}, \ldots, t_{(n)}$



$\dfrac{i - 3/8}{n + 1/4}$     ~ $N(0,1)$

$t_{(i)}$
$\uparrow$
$z_{\frac{i}{n}} = \Phi^{-1}\left(\dfrac{i - 3/8}{n + 1/4}\right)$
$\uparrow$

$\Phi(z) = $ c.d.f. of $Z = N(0,1)$

---

## Transformation on $y$   Box-Cox transformation (Power transformation)

$$y_i^* = \begin{cases} \dfrac{y_i^\lambda - 1}{\lambda} & \lambda \neq 0 \\[2mm] \log_e(y_i) & \lambda = 0 \end{cases}$$

$\lambda$ — unknown
$\uparrow$
Find m.l.e. of $\lambda$

Assume $\underset{\sim}{Y}^* \sim N(X\underset{\sim}{\beta}, \sigma^2 \underset{\sim}{I})$

— Data: $\underset{\sim}{Y}, X$

— dist. of $Y$ is unknown
— dist. of $\underset{\sim}{Y}^* \sim N$
$\Rightarrow$ dist. of $Y$

Likelihood $(\underset{\sim}{\beta}, \sigma^2, \lambda \mid \underset{\sim}{Y}, X)$

$= \dfrac{1}{(2\pi\sigma^2)^{n/2}} \exp\left\{ -\dfrac{(\underset{\sim}{Y}^* - X\underset{\sim}{\beta})^2}{2\sigma^2} \right\} * J(\lambda, \underset{\sim}{Y})$

where $J(\lambda, \underset{\sim}{Y}) = \boxed{\prod_{i=1}^{n} y_i}^{\lambda - 1}$

$= GM(y)^{\lambda - 1}$
    $\underset{\sim}{\hspace{3mm}}$ geometric mean

⑧

linear regression
$$\log_e L = -\frac{n}{2}\log(2\pi) - \frac{n}{2}\log(\hat{\sigma}^2) - \frac{n}{2}$$

$$= -\frac{n}{2}\log\left(\frac{ResSS}{n}\right) + constant$$

log likelihood

$$likelihood(\lambda) = -\frac{n}{2}\log\left(\tilde{\sigma}^2(\lambda)\right)$$

$$\uparrow \qquad \tilde{\sigma}^2(\lambda) = \frac{\underset{\sim}{Y}^{*T}\underset{\sim}{Y}^* - \underset{\sim}{Y}^{*T}\underset{\sim}{X}(\underset{\sim}{X}^T\underset{\sim}{X})^{-1}\underset{\sim}{X}^T\underset{\sim}{Y}^*}{n}$$

profile likelihood

define $\underset{\sim}{Z}^\lambda = \underset{\sim}{Y}^* / J^{1/n}$

$\Rightarrow \log likelihood(\lambda) = -\frac{n}{2}\log\left(ResSS.\lambda(\underset{\sim}{Z}^\lambda)\right)$

where $\underset{\sim}{Z}^\lambda$ is $n \times 1$ vector with $Z_i^\lambda = \begin{cases} \frac{y^{\lambda}-1}{\lambda[GM(y)]^{\lambda-1}} & \lambda \neq 0 \\ GM(y)\ln(y_i) & \lambda = 0 \end{cases}$

For each $\lambda$, $\Rightarrow$ calculate $Z_i^\lambda$.

$\qquad\qquad$ fit $Z_i^\lambda$ on $\underset{\sim}{X}$

$\qquad\qquad \Rightarrow ResSS.\lambda(\underset{\sim}{Z})$



e.g. $\hat{\lambda} = 0.423$ ∴ 95% of $\hat{\lambda}$?

$\qquad\qquad\qquad$ Does it over 0.5 ?