



## 10.1. The active terms

- Variable selection

- Aim: Identify the correct model
  - select the useful predictor
  - Ignore the non-informative terms
- Y v.s.  $X_1, X_2, \dots, X_{999}$ 
  - Divide  $X=(X_1, X_2, \dots, X_{999})$  into two sets,  $X_A$ , and  $X_I$ ,
  - so that  $E(Y|X) = E(Y|X_A) = X_A\beta_A$

$$X = (X_A, X_I)$$

$X_A$  = active terms

$X_I$  = inactive terms

given  $X_A$

given  
all  $X_1 \dots X_{999}$

$X = (X_1, X_2, \dots, X_5)$  aliased terms :  $X_3 \equiv c_0 + c_1 X_1 + \dots + c_5 X_5$

## 10.1. Active terms and multicollinearity

- Multicollinearity

- some terms can be approximated by linear combination of the other terms.
  - e.g.  $X_3 \approx c_0 + c_1 X_1 + c_2 X_2 + c_4 X_4 + c_5 X_5$
- In this case,  $X'X$  is close to singular ( $\det=0$ ),
  - $\hat{\beta} = (X'X)^{-1}X'Y$  and  $\text{Var}(\hat{\beta}) = \sigma^2(X'X)^{-1}$  can be huge.
- We should avoid including all variables with multicollinearity in the regression model
  - e.g. set  $X_A = (X_1, X_2, X_4, X_5)$ ,
  - $X_I = (X_3)$ , since  $X_3$  can be explained by  $X_1, X_2, X_4, X_5$

$\downarrow$   
 $(X'X)^{-1}$   
 does  
not  
exist

$\downarrow$   
 $\hat{\beta} = (X'X)^{-1}X'Y$   
 does  
not  
exist

$$SYY = RSS + SS_{reg}$$

## 10.1. Active terms and multicollinearity

- How to detect multicollinearity?
  - Check  $(X'X)^{-1}$ ?
    - problem: don't know how large is large.
  - A better method:  $R_j^2$ , the coefficient of determination for the regression

$$R^2 = \frac{SS_{reg}}{SYY}$$

1 = good fit  
 0 = bad fit

$$X_j = c_0 + c_1 X_1 + \dots + c_{j-1} X_{j-1} + c_{j+1} X_{j+1} + \dots + c_p X_p + e$$

- $R_j^2 \approx 1 \rightarrow$  multicollinearity, i.e. some terms can be approximated by linear combination of the other terms.

e.g.  $(X_1, X_2, X_3, X_4, X_5)$

If all  $R_j^2$  are small  $\rightarrow$  no

$$X_1 = C_0 + C_1 X_2 + C_2 X_3 + C_3 X_4 + C_4 X_5 + e$$

$$X_2 = \dots (X_1, X_3, X_4, X_5) \rightarrow R_2^2$$

$$X_3 = (X_1, X_2, X_4, X_5) \rightarrow R_3^2$$

$$X_4 = \dots \rightarrow R_4^2$$

$$X_5 = \dots \rightarrow R_5^2$$

## 10.1. Active terms and multicollinearity

- Relationship between  $\text{Var}(\hat{\beta}_j)$  and  $R_j^2$

- Using the idea of Added Variable Plot

Let  $X_O = (1 \ X_1 \ X_2 \dots X_{j-1} \ X_{j+2} \dots X_p)$ ,  $H_O = X_O(X_O'X_O)^{-1}X_O'$

For  $Y = X_O\beta_O + X_j\beta_j + e$ ,

$\hat{\beta}_j$  = Regression coefficient between  $(I - H_O)Y$  and  $(I - H_O)X_j$

$$= (X_j'(I - H_O)X_j)^{-1}X_j'(I - H_O)Y \quad \leftarrow (X'X)^{-1}X'$$

$$\text{Var}(\hat{\beta}_j) = \sigma^2 (X_j'(I - H_O)X_j)^{-1}$$

$$Y = x\beta + e$$

$\hat{e}_Y \propto \hat{e}_{x_j}$

$$\leftarrow \text{Var}(AX) = A \text{Var}(X) A'$$

- For the regression

$$X_j = c_0 + c_1X_1 + \dots + c_{j-1}X_{j-1} + c_{j+1}X_{j+1} \dots + c_pX_p + e, \quad \text{or}$$

$$X_j = X_Oc_O + e,$$

we have  $R_j^2 = 1 - \frac{RSS}{TSS} = 1 - \frac{X_j'(I - H_O)X_j}{X_O'X_O}$

$$\leftarrow Y'(I - H)Y$$

$$\frac{SSE}{SYY} = \frac{SSE}{TSS} = \frac{TSS - RSS}{TSS} = 1 - \frac{RSS}{TSS}$$



Chapter10.pdf - Adobe Reader

File Edit View Window Help

Tools Sign Comment

## 10.1. Variance Inflation Factor (VIF)

- Relationship between  $\text{Var}(\hat{\beta}_j)$  and  $R_j^2$ 
  - $\text{Var}(\hat{\beta}_j) = \sigma^2 (X'_j (I - H_O) X_j)^{-1}$
  - $R_j^2 = 1 - \frac{\text{RSS}}{\text{TSS}} = 1 - \frac{X'_j (I - H_O) X_j}{S X_j X_j}$
  - Therefore,  $\text{Var}(\hat{\beta}_j) = \frac{\sigma^2}{(1 - R_j^2) S X_j X_j}$
  - Variance Inflation Factor (VIF)

$y = \beta_0 + \beta_j x_j + e$  (Simple linear)

$\text{Compare to } \text{Var}(\hat{\beta}_j) = \frac{\sigma^2}{S X_j X_j}$

$\text{IF } R_j^2 \approx 1 \Rightarrow \frac{1}{1 - R_j^2} \text{ very large}$

$\text{IF } R_j^2 \approx 0 \Rightarrow \frac{1}{1 - R_j^2} \approx 1$

$\frac{1}{1 - R_j^2} \rightarrow$

11.69 x 8.27 in

Page 7 of 23

NEXT PC MODE POWERPOINT BALL PEN HIGH LIGHT SIZE COLOR LINE RECT. CIRCLE SELECT ERASE ALL ERASE REC RECORDING 00:00:00 SAVE OPEN HELP EXIT

## 10.1. Active terms and multicollinearity

- Example

```
x1=c(1,3,2,4,5,2,3,1,0,5)
x2=c(8,9,7,2,5,9,6,4,4,1)
x3=2*x1-5*x2+rnorm(10,0,0.1) x3 ≈ 2x1 - 5x2, multicollinearity
x4=c(3,1,4,2,7,3,4,5,6,3)
y=3+x1+2*x2+2*x4+rnorm(10,0,0.5)
summary(lm(y~x1+x2+x3))
#Find VIF
Rj1=summary(lm(x1~x2+x3+x4))$r.squared; VIF1=1/(1-Rj1)
Rj2=summary(lm(x2~x1+x3+x4))$r.squared; VIF2=1/(1-Rj2)
Rj3=summary(lm(x3~x1+x2+x4))$r.squared; VIF3=1/(1-Rj3)
Rj4=summary(lm(x4~x1+x2+x3))$r.squared; VIF4=1/(1-Rj4)
print(rbind(c("Rj",Rj1,Rj2,Rj3,Rj4),c("VIF",VIF1,VIF2,VIF3,VIF4)))
#Modified fitting by deleting either one of x1,x2,x3
summary(lm(y~x1+x2+x4))
summary(lm(y~x2+x3+x4))
summary(lm(y~x1+x3+x4))
```

$y$  vs  $x_1, x_2, \dots, x_{999}$

Model 1 :  $y = \beta_0 + \beta_1 x_1 \rightarrow \text{RSS}_1$ ,  $\text{RSS} = \sum (y_i - \hat{y}_i)^2$  measures

Model 2 :  $y = \beta_0 + \beta_2 x_2 \rightarrow \text{RSS}_2$ , the goodness of fit of  
a regression

; ; ;

Model k  $y = \beta_0 + \beta_1 x_1 + \dots + \beta_{10} x_{10} + \dots \rightarrow \text{RSS}_k$

Q: Can we use RSS to select a good model?

A: No: The fitting of Regression is based on

$$\min_{\beta_m} \text{RSS}(\beta_m)$$

More parameter (more x variables)  $\rightarrow$  small RSS  
 $\Rightarrow$  Choose the model with most number of variable



$$RSS = \sum (Y_i - \hat{Y}_i)^2$$

## 10.2. Automatic Variable Selection procedure

- For all possible candidate models, we compute
  - Akaike Information Criteria (AIC)  $-E_{Y|X} (\text{likelihood}(Y|X))$
  - Bayesian Information Criteria (BIC)  $- \text{posterior prob(model)}$
  - Mallow's C<sub>p</sub> Statistics  $\hat{\sigma}^2 \text{ using all } X \rightarrow \frac{RSS}{\hat{\sigma}^2} + 2p_c - n$
  - Predicted residual sum of Square (PRESS)

$$n \log\left(\frac{RSS}{n}\right) + 2p_c$$

Number of parameters in the model  
 e.g. p+1 in regression

$$n \log\left(\frac{RSS}{n}\right) + p_c \log(n)$$

$$\text{Mallow's C}_p \text{ Statistics} = \sum (Y_i - E(Y_i|X_i))^2$$

$\hat{Y}_i$  - from fitting using all data  
 $\hat{Y}_{i(i)}$  - is not affected by  $(X_i, Y_i)$   
 $\hat{Y}_{i(i)}$  - prediction of the  $Y_i$ , using reg fitting without  $(X_i, Y_i)$

## 10.2. Model Selection in Practice

- When you have  $Y$  and  $X_1, X_2, \dots, X_p$ 
  - For each possible model (e.g.  
 $y = \beta_0 + \beta_1 x_1, y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p, y = \beta_{p-1} x_{p-1} + \beta_p x_p$  etc)
    - Find AIC, BIC,  $C_p$ , PRESS
    - Report the best model (smallest value) w.r.t each criteria
    - How many possible models ? (how many combinations)
      - $2^p$
      - If  $p=20, 2^{20}=1048576 / 60 \times 60 \times 24$
      - If a regression takes 1s, how long does it take for model selection?
  - Solutions
    - Forward Selection
    - Backward Selection

## 10.2. Automatic Variable Selection procedure

lack of fit is large

- AIC

$$n \log\left(\frac{RSS}{n}\right) + 2p_c$$

- BIC

$$n \log\left(\frac{RSS}{n}\right) + p_c \log(n)$$

- $C_p$

$$\frac{RSS}{\hat{\sigma}^2} + 2p_c - n$$

- PRESS

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

- Smaller value implies a better model

- For AIC, BIC,  $C_p$ , they have a common structure:

**lack of fit + penalty for model complexity**

- Larger model → Smaller RSS, Bigger penalty
- Smaller model → Larger RSS, Smaller penalty
- When p is fixed, they yield the same result (min RSS)

- For PRESS, no penalty is need since different data are used for fitting and estimating errors

8.27 in

e.g.  $Y$   $\underline{\underline{X_1, X_2, X_3}}$

Model 1:  $Y = \beta_0 + \beta_1 X_1 + e \rightarrow AIC_1$

(use BIC/C<sub>p</sub>)

Model 2:  $Y = \beta_0 + \beta_2 X_2 + e \rightarrow AIC_2$

Press  
instead of  
 $AIC$ )

$\vdots \quad \vdots \quad \vdots$

Model K:  $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + e \rightarrow AIC_K$

$\vdots$

$\rightarrow AIC_{10}$

Model 10:

1 2 3    2 3  
1    int  
2  
3  
1 2  
1 3

e.g., If  $AIC_5 = \min_{K=1-10} AIC_K$

(2 3) = 8

$\Rightarrow$  select Model 5

# Exam

- Ch 1-8 & Ch 10
- 1 page of A4 size (double sided)
- 50 point from Ch 4-10 from MC game  
short question
- 30 points numerical calculation (iii)  $(\bar{X}, \bar{Y})^T$  are given  
eg  $\hat{\beta}$ , CI for  $\beta_1$ , CE for  $\beta_1$ , predict  $\hat{y}_0$
- 10 point  $\text{Cor}(\hat{e}_i, \hat{e}_j) = E(\text{Sum of } S_{ij})$

10 point "new"