# Chapter 10

Variable Selection

## 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<sub>1</sub>, X<sub>2</sub>, ... X<sub>999</sub>
    - Divide X=(X<sub>1</sub>, X<sub>2</sub>, ... X<sub>999</sub>) into two sets, X<sub>A</sub>, and X<sub>I</sub>,
    - so that  $E(Y|X) = E(Y|X_A) = X_A \beta_A$ 
      - X<sub>A</sub> = active terms
      - X<sub>I</sub> =inactive terms

#### 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  $\operatorname{Var}(\hat{\beta}) = \sigma^2 (X'X)^{-1}$

can be huge.

- We should avoid including all variables with multicollinearity in the regression model
  - e.g. set  $X_A = (X_1, X_2, X_4, X_5)$ ,
  - $X_1 = (X_3)$ , since  $X_3$  can be explained by  $X_1, X_2, X_4, X_5$

#### Remarks

- Multicollinearity
  - some terms can be <u>approximated</u> by linear combination of the other terms.
  - (X'X) is close to non-singular. Inverse exists but unstable

• e.g. 
$$X_3 \approx c_0 + c_1 X_1 + c_2 X_2 + c_4 X_4 + c_5 X_5$$

- Perfect/Exact multicollinearity or Aliased
  - some term is <u>exactly expressed</u> by linear combination of the other terms.
  - (X'X) is singular. Inverse does not exist.
  - e.g.  $X_3 = c_0 + c_1 X_1 + c_2 X_2 + c_4 X_4 + c_5 X_5$

• How to detect multicollinearity?

- Check (X'X)<sup>-1</sup>?
  - problem: don't know how large is large.
- A better method: R<sup>2</sup><sub>j</sub>, the coefficient of determination for the regression

$$X_{j} = c_{0} + c_{1}X_{1} + \dots + c_{j-1}X_{j-1} + c_{j+1}X_{j+1} + \dots + c_{p}X_{p} + e$$

R<sup>2</sup><sub>j</sub>≈1 → multicollinearity, i.e. some terms can be <u>approximated</u> by linear combination of the other terms.

• Relationship between  $Var(\hat{\beta}_i)$  and  $R_i^2$ 

- 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  $Y = X_0 \beta_0 + X_j \beta_j + e$ ,

 $\hat{\beta}_{j} = Regression \ coefficient \ between (I - H_{o})Y \ and (I - H_{o})X_{j}$ 

$$= \left(X'_{j}(I - H_{o})X_{j}\right)^{-1}X'_{j}(I - H_{o})Y$$
$$Var(\hat{\beta}_{j}) = \sigma^{2}\left(X'_{j}(I - H_{o})X_{j}\right)^{-1}$$

For the regression

$$\begin{split} X_{j} &= c_{0} + c_{1}X_{1} + \ldots + c_{j-1}X_{j-1} + c_{j+1}X_{j+1} \ldots + c_{p}X_{p} + e, \qquad on \\ X_{j} &= X_{O}c_{O} + e, \end{split}$$

we have  $R_{j}^{2} = 1 - \frac{RSS}{TSS} = 1 - \frac{X'_{j}(I - H_{o})X_{j}}{SX_{j}X_{j}}$ 

## 10.1. Variance Inflation Factor (VIF)



- In practice, R<sup>2</sup> > 0.7 is regarded as strong correlation
  - Becareful if  $R_i^2 > 0.7$  or VIF > 1/0.3 = 3.33

• 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)
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
    R_{j1}=summary(Im(x1~x2+x3+x4)) (IF1=1/(1-R_{j1}))
    Rj2=summary(Im(x2~x1+x3+x4))$r.squared; VIF2=1/(1-Rj2)
    R_{j3}=summary(Im(x_3 x_1 + x_2 + x_4)) r.squared; VIF3=1/(1-R_{j3})
    R_{j4}=summary(Im(x4~x1+x2+x3)) squared; VIF_{4}=1/(1-R_{j4})
    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(Im(y \sim x1 + x2 + x4))
    summary(Im(y \sim x2 + x3 + x4))
    summary(lm(y \sim x1 + x3 + x4))
```

## 10.2. Automatic Variable Selection procedure

- For all possible candidate models, we compute
  - Akaike Information Criteria (AIC)  $n \log\left(\frac{RSS}{n}\right) + 2p_{c}$ Number of parameters in the model e.g. p+1 in regression
  - Bayesian Information Criteria (BIC)

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

Mallow's C<sub>p</sub> Statistics

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

Predicted residual sum of Square (PRESS)

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

## 10.2. Automatic Variable Selection procedure



## 10.2. Automatic Variable Selection procedure



- When you have Y and  $X_1, X_2, \dots X_p$ 
  - For each possible model  $y=\beta_0+\beta_1x_1, y=\beta_0+\beta_1x_1+...+\beta_px_p, y=\beta_{p-1}x_{p-1}+\beta_px_p$  etc)
    - Find AIC, BIC, C<sub>p</sub>, PRESS
    - Report the best model (smallest value) w.r.t each criteria
    - How many possible models ? (how many combinations)
      - 2<sup>p</sup>
      - If p=20, 2<sup>20</sup>=1048576
      - If a regression takes 1s, how long does it take for model selection?

(e.g

#### Solutions

- Forward Selection
- Backward Selection

#### **Forward Selection**

- Set-up
  - K=Total number of terms may be added
  - i=1, L=1 (current number of terms, first start with intercept only)
  - V<sub>i</sub> (current criterion values)
- Step i
  - Each of the time, add one term beyond the current model
  - Obtain K+1-L criterion values
  - Stop if
    - All terms are included
    - All K+1-L criterion values are greater than V<sub>i</sub>, i.e. additional terms does not improve the fitting.
  - Set V<sub>i+1</sub> be the minimum of the K+1-L criterion values. Add the corresponding term to the current model, go to Step i+1

## Forward Selection

• Each time add 1 variable, until no improvement

#### Example

#Data

- set.seed(1);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); x4=c(3,1,4,2,7,3,4,5,6,3)
- y=3+x1+2\*x2+2\*x4+rnorm(10,0,0.5)
- n=10; V<sub>0</sub>= n\*log(sum(lm(y~1)\$residuals^2)/n)+2\*1=35.49
- # Step 1: (current model  $y = \beta_0$ ,  $V_0 = 35.49$ )
  - AIC1a=n\*log(sum(lm(y~x1)\$residuals^2)/n)+2\*(1+1)
  - AIC1b=n\*log(sum(lm(y~x2)\$residuals^2)/n)+2\*(1+1)
  - AIC1c=n\*log(sum(lm(y~x3)\$residuals^2)/n)+2\*(1+1)  $n \log \frac{n \log n}{n \log n}$
  - AIC1d=n\*log(sum(lm(y~x4)\$residuals^2)/n)+2\*(1+1)
  - print(c(AIC1a, AIC1b, AIC1c, AIC1d))
  - V<sub>1</sub>=29.30, add x2,



n

#### **Forward Selection**

Each time add 1 variable, until no improvement

- # Step 2: (current model  $y = \beta_0 + \beta_2 X_2$ ,  $V_1 = 29.30$ )
  - AIC2a=n\*log(sum(Im(y~x2+x1)\$residuals^2)/n)+2\*(2+1)
  - AIC2b=n\*log(sum(lm(y~x2+x3)\$residuals^2)/n)+2\*(2+1)
  - AIC2c=n\*log(sum(lm(y~x2+x4)\$residuals^2)/n)+2\*(2+1)
  - print(c(AIC2a, AIC2b, AIC2c))
  - V<sub>2</sub>=13.85, Add x4
- # Step 3: (current model  $y = \beta_0 + \beta_2 X_2 + \beta_4 X_4$ ,  $V_2 = 13.85$ )
  - AIC3a=n\*log(sum(lm(y~x2+x4+x1)\$residuals^2)/n)+2\*(3+1)
  - AIC3b=n\*log(sum(lm(y~x2+x4+x3)\$residuals^2)/n)+2\*(3+1)
  - print(c(AIC3a, AIC3b))
  - V<sub>3</sub>= -9.30, Add x1

#### **Forward Selection**

Each time add 1 variable, until no improvement

- # Step 4: (current model y=  $\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_4 X_4$ ,  $V_3 = -9.30$ )
  - AIC4a=n\*log(sum(lm(y~x2+x4+x1+x3)\$residuals^2)/n)+2\*(4+1)
  - print(c(AIC4a))
  - $V_4$  = -8.56 > -9.30 =  $V_3$ . STOP, without adding x3
  - Conclusion
    - The model  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_4 x_4$  is the best according to AIC using Forward Selection.

#### **Backward Selection**

- Set-up
  - K=Total number of terms may be added
  - i=1, L=K+1 (current number of terms, first fit all terms)
  - V<sub>i</sub> (current criterion values)
- Step i
  - Each of the time, remove one term from the current model
  - Obtain L-1 criterion values
  - Stop if
    - All terms are removed
    - All L-1 criterion values are greater than V<sub>i</sub>, i.e. removing terms worsen the fitting.
  - Set V<sub>i+1</sub> be the minimum of the L-1 criterion values. Remove the corresponding term from the current model, go to Step i+1

#### **Backward Selection**

Each time delete 1 variable, until no improvement

- #Data
- set.seed(2);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); x4=c(3,1,4,2,7,3,4,5,6,3)
- y=3+x1+2\*x2+2\*x4+rnorm(10,0,0.5)
- n=10;  $V_0 = n^{10}(sum(lm(y \sim x1 + x2 + x3 + x4))) = -4.48$
- # Step 1: (current model y=  $\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4$ , V<sub>0</sub>= -4.48)
  - AIC1a=n\*log(sum(lm(y~x2+x3+x4)\$residuals^2)/n)+2\*(3+1)
  - AIC1b=n\*log(sum(lm(y~x1+x3+x4)\$residuals^2)/n)+2\*(3+1)
  - AIC1c=n\*log(sum(lm(y~x1+x2+x4)\$residuals^2)/n)+2\*(3+1)
  - AIC1d=n\*log(sum(lm(y~x1+x2+x3)\$residuals^2)/n)+2\*(3+1)
  - print(c(AIC1a, AIC1b, AIC1c, AIC1d))
  - V<sub>1</sub>= -6.46, remove x1,

### **Backward Selection**

Each time delete 1 variable, until no improvement

- # Step 2: (current model  $y = \beta_0 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4$ ,  $V_0 = -6.46$ )
  - AIC2a=n\*log(sum(Im(y~x3+x4)\$residuals^2)/n)+2\*(2+1)
  - AIC2b=n\*log(sum(Im(y~x2+x4)\$residuals^2)/n)+2\*(2+1)
  - AIC2c=n\*log(sum(lm(y~x2+x3)\$residuals^2)/n)+2\*(2+1)
  - print(c(AIC2a, AIC2b, AIC2c))
  - $V_2$ =16.10> -6.46 = $V_1$ . STOP, without removing any variable
  - Conclusion
    - The model  $y = \beta_0 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4$  is the best according to AIC using Backward Selection.

## What you have learnt

- Regression: Relationship b/w variables
  - Y=Xβ+e
    - Least Sq Est  $\beta$ , Test if  $\beta$ =0, C.I. for  $\beta$ , Predict Y
    - Compare between models by F-test
  - Diagnostic check Residual/scatter/AV plot
  - Improvement
    - Non-constant variance -- WLS
    - Curve relationship -- 1) Poly Reg 2) Transform
    - Check Outlier 1) T-test 2) Cook's distance

• Variable selection: 1)AIC. 2)BIC. 3)C<sub>p</sub>. 4)Press