

# Variable Selection

Consider:  $Y = X\beta + \varepsilon$ .  $\varepsilon \sim N_n(0, \sigma^2 I_n)$  where

$$Y \in \mathbb{R}^n, \beta \in \mathbb{R}^p \text{ (unknown)}, X \in \mathbb{R}^{n \times p}, r(X) = p.$$

$H_0: (0 \ I_{p-1})\beta = 0$  is test for whether all the variables influence  $Y$  (data).  $\beta = (\beta_0 \ \beta_1 \ \dots \ \beta_{p-1})$

$X = (J_n \ X_1 \ \dots \ X_{p-1})$ .  $\beta_0$  is fixed.

$H_0: \beta_i = 0$  is test for single one.

## (1) Consequence of Selection:

We want to reduce the number of variables but retain the accuracy of estimation at the same time.

Suppose  $X_2 \in \mathbb{R}^{n \times 2}$ .  $\beta_2 \in \mathbb{R}^2$ .  $X = (X_2 \ X_0)$ .

$\beta = \begin{pmatrix} \beta_2 \\ \beta_0 \end{pmatrix}$ . Rewrite the model in:

$$Y = X_2 \beta_2 + X_0 \beta_0 + \varepsilon$$

Choose  $q-1$  variables  $X_2$  from  $p-1$  variables

$X$  then generate an alternative model:

$$Y = X_2 \beta_2 + \varepsilon. \text{ correspond } H_0: (0 \ I_{q-1})\beta = 0 \quad (\Lambda = (0, I_{q-1}))$$

$$R = \frac{\sum (\hat{\Lambda}^T \hat{\beta})^T (\hat{\Lambda}^T (X^T X)^{-1} \hat{\Lambda})^{-1} \hat{\Lambda}^T \hat{\beta} / r(\Lambda)}{MSE} = \frac{\hat{\beta}_0^T D^{-1} \hat{\beta}_0}{r MSE} \geq F_{\alpha}^{(q-1, n-q)}$$

where  $D = (X_t^T (I - P_{X_2}) X_t)^{-1}$ ,  $P_{X_2} = X_2 (X_2^T X_2)^{-1} X_2^T$

$$B_2 = \begin{pmatrix} I_1 & 0 \\ -X_t^T X_2 (X_2^T X_2)^{-1} & I_t \end{pmatrix} X^T X \begin{pmatrix} I_1 & -(X_2^T X_2)^{-1} X_2^T X_t \\ 0 & I_t \end{pmatrix}$$

$$= \begin{pmatrix} X_2^T X_2 & 0 \\ 0 & D^{-1} \end{pmatrix}. \text{ Denote } A = (X_2^T X_2)^{-1} X_2^T X_t$$

$$\Rightarrow (X^T X)^{-1} = \begin{pmatrix} (X_2^T X_2)^{-1} + A D A^T & -A D \\ -D A^T & D \end{pmatrix}$$

Denote: i)  $\hat{\beta} = (X^T X)^{-1} X^T Y$ ,  $\hat{\sigma}^2 = \frac{Y^T (I - M) Y}{n - r(X)}$ ,  $\hat{\eta} = X^T \hat{\beta} \in \mathbb{R}^n$

ii)  $\tilde{\beta}_2 = (X_2^T X_2)^{-1} X_2^T Y$ ,  $\tilde{\sigma}_2^2 = \frac{Y^T (I - M_2) Y}{n - r(X_2)}$ ,  $\tilde{\eta} = X_2^T \tilde{\beta}_2 \in \mathbb{R}^n$

where  $\tilde{x} = (\tilde{x}_2 \tilde{x}_t)$ ,  $\eta$  is estimate

$$\text{iii) } \text{MSE}(\hat{\theta}) = E[(\hat{\theta} - \theta)(\hat{\theta} - \theta)^T]$$

Thm. (Influence on estimation)

i)  $E(\hat{\beta}) = \beta$ . if the full model is correct.

$$\beta_t = 0 \text{ or } X_2^T X_t = 0 \Leftrightarrow E(\tilde{\beta}_2) = \beta_2$$

ii)  $\text{Var}(\hat{\beta})_2 - \text{Var}(\tilde{\beta}_2) \geq 0$

iii)  $\text{Var}(\hat{\beta})_t - \beta_t \beta_t^T \geq 0 \Rightarrow$

$$\text{Var}(\hat{\beta})_2 - E[(\tilde{\beta}_2 - \beta_2)(\tilde{\beta}_2 - \beta_2)^T] \geq 0$$

iv)  $E(\tilde{\sigma}_2^2) \geq E(\hat{\sigma}^2) = \sigma^2$ . "=" holds if and only if  $\beta_t = 0$ .

Rmk: i) means  $\tilde{\beta}_2$  isn't a unbiased estimator generally

ii) means if we used the reduced model but the origin model is true. then Variance of estimator will reduce.

iii) means the variables  $\beta_+$  discarded exactly effects the estimation. Since  $\beta_+$  can't be estimated and have large Variance ( $\text{Var}(\hat{\beta}_+) \geq \beta_+ \beta_+^T$ ). Delete  $\beta_+$  can reduce Variance.

iv) means  $\tilde{\sigma}_2^2$  isn't unbiased estimator for  $\sigma^2$  if the origin model is true.

Pf: i)  $E(\tilde{\beta}) = (X^T X)^{-1} X^T X \beta = \beta$

$$E(\tilde{\beta}_2) = (X_2^T X_2)^{-1} X_2^T X \beta = (X_2^T X_2)^{-1} X_2^T (X_2 \ X_+) \begin{pmatrix} \beta_2 \\ \beta_+ \end{pmatrix}$$

$$= \beta_2 + (X_2^T X_2)^{-1} X_2^T X_+ \beta_+ = \beta_2 + A \beta_+$$

ii)  $\text{Var}(\tilde{\beta}) = (X^T X)^{-1} X^T \text{Var}(Y) X (X^T X)^{-1}$

$$= \sigma^2 (X^T X)^{-1} = \sigma^2 \begin{pmatrix} (X_2^T X_2)^{-1} + ADA^T & -AD \\ -DA^T & D \end{pmatrix}$$

$$\therefore \text{Var}(\hat{\beta})_2 = \sigma^2 ((X_2^T X_2)^{-1} + ADA^T)$$

with  $\text{Var}(\tilde{\beta}_2) = \sigma^2 (X_2^T X_2)^{-1}$

$$\Rightarrow \text{Var}(\hat{\beta})_2 - \text{Var}(\tilde{\beta}_2) = \sigma^2 ADA^T \geq 0$$

$$\begin{aligned} \text{MSE}(\tilde{\beta}_2) &= \text{Var}(\tilde{\beta}_2) + (E(\tilde{\beta}_2) - \beta_2)(E(\tilde{\beta}_2) - \beta_2)^T \\ &= \text{Var}(\tilde{\beta}_2) + A\beta_2\beta_2^T A^T \end{aligned}$$

$$\therefore \text{Var}(\hat{\beta}_2) - \text{MSE}(\tilde{\beta}_2) = A(\sigma^2 D - \beta_2\beta_2^T)A^T$$

$$\text{iv) } E(\tilde{\sigma}_n^2) = \frac{1}{n-r(X_2)} E(Y^T(I - M_2)Y)$$

$$= \frac{1}{n-r(X_2)} \text{tr}((I - M_2)E(Y Y^T))$$

$$= \frac{1}{n-r(X_2)} \text{tr}((I - M_2)(\sigma^2 I + X P P^T X^T))$$

$$= \sigma^2 + \frac{1}{n-r(X_2)} P^T X^T (I - M_2) X P$$

$$= \sigma^2 + \frac{1}{n-r(X_2)} \beta_2^T X_2^T (I - M_2) X_2 \beta_2$$

Thm. (Influence on estimate)

i)  $E(\tilde{\eta}) = X^T \beta$  if the origin model is true

$$\beta_2 = 0 \Leftrightarrow E(\tilde{\eta}) = X^T \beta$$

ii)  $\text{Var}(\eta - X^T \hat{\beta}) \geq \text{Var}(\eta - X_2^T \tilde{\beta}_2)$

iii)  $\text{Var}(\hat{\beta})_t - \beta_2 \beta_2^T \geq 0 \Rightarrow \text{Var}(\eta - X^T \hat{\beta}) \geq E(\eta - X_2^T \tilde{\beta}_2)$

Rmk: It means if origin model is true.

then  $\tilde{\beta}_2$  isn't unbiased. But its MSE

will reduce if  $\text{Var}(\hat{\beta})_t \geq \beta_2 \beta_2^T$ , i.e. its variance is large.

Pf: i)  $E(\hat{\eta}) = X^T E(\hat{\beta}) = X^T \beta$ .

$$E(\hat{\eta}_2) = X_2^T (\beta_2 + A\beta_1) = X_2^T \beta_2 + X_2^T A\beta_1$$

ii)  $\text{Var}(\eta - X^T \hat{\beta}) = \text{Var}(\eta) + X^T \text{Var}(\hat{\beta}) X$   
 $= \sigma^2 (I + X^T (X^T X)^{-1} X)$

$$\text{Var}(\eta - X_2^T \hat{\beta}_2) = \text{Var}(\eta) + X_2^T \text{Var}(\hat{\beta}_2) X_2$$

$$= \sigma^2 (I + X_2^T (X_2^T X_2)^{-1} X_2)$$

$$\Rightarrow \text{Var}(\eta - X^T \hat{\beta}) - \text{Var}(\eta - X_2^T \hat{\beta}_2) =$$

$$\sigma^2 (A^T X_2 - X_1)^T D (A^T X_2 - X_1) \geq 0$$

iii)  $E(\eta - X_2^T \hat{\beta}_2)^2 = \text{Var}(\eta - X_2^T \hat{\beta}_2) + E(\eta - X_2^T \hat{\beta}_2)^2$   
 $= \text{Var}(\eta - X_2^T \hat{\beta}_2) + (X^T \beta - X_2^T (\beta_2 + A\beta_1))^2$   
 $= \text{Var}(\eta - X_2^T \hat{\beta}_2) + (A^T X_2 - X_1)^T \beta_1 \beta_1^T (A^T X_2 - X_1)$

$$\Rightarrow \text{Var}(\eta - X^T \hat{\beta}) - E(\eta - X_2^T \hat{\beta}_2)^2 =$$

$$\sigma^2 (A^T X_2 - X_1)^T (\sigma^2 D - \beta_1 \beta_1^T) (A^T X_2 - X_1) \geq 0$$

with  $\text{Var}(\hat{\beta})_1 - \beta_1 \beta_1^T = \sigma^2 D - \beta_1 \beta_1^T$ .

Rmk: Back to the hypothesis testing: Replace para.

in  $\sigma^2 D - \beta_1 \beta_1^T \geq 0$  with estimator, then,

$$\hat{\sigma}^2 D - \hat{\beta}_1 \hat{\beta}_1^T \geq 0 \Leftrightarrow \hat{\beta}_1^T D^{-1} \hat{\beta}_1 / t \hat{\sigma}^2 \leq \frac{1}{t}$$

i.e. if  $\frac{1}{t} < F_{\alpha}(t, n-p)$ , then we accept

$H_0: \beta_1 = 0$ . (set  $\gamma = D^{-\frac{1}{2}} \hat{\beta}_1$ , prove:  $\hat{\sigma}^2 I \geq \alpha \alpha^T$

$\Leftrightarrow \hat{\sigma}^2 \geq \alpha^T \alpha$ , it's easy to see)

## (2) Principles:

If we have  $p$  variables to be selected.

Then there're  $2^p - 1$  possible regression equation

Denote: i)  $SST = Y^T Y$ ,  $SSE = Y^T (I - M) Y$

$$SSE_2 = Y^T (I - M_2) Y.$$

$$\text{ii) } R^2 = 1 - \frac{SSE}{SST} \quad R_2^2 = 1 - \frac{SSE_2}{SST}$$

Rmk:  $\forall z \leq p$ ,  $SSE \leq SSE_2$ ,  $R^2 \geq R_2^2$

### ① For fitting the model:

i) Minimize the mean of  $SSE_2$ :

i.e. find  $z$ .  $\hat{\sigma}_2^2 = \min_r \frac{SSE_r}{n-r}$

Rmk:  $SSE_2 \uparrow$  as  $z \downarrow$ ,  $n-z \uparrow$  as  $z \downarrow$

if  $\{x_i\}_2^2$  effects  $\eta$  significantly

then  $SSE_2 \downarrow$  fast as  $z \uparrow$ .

if not, then  $SSE_2 \downarrow$  slowly as  $z \uparrow$

$\frac{1}{n-z} \uparrow$  as  $z \uparrow$ . is penalty of increase of number of variables.

ii) Maximize  $\bar{R}_2^2 = 1 - \frac{SSE_2 / (n-z)}{SST / (n-1)} = 1 - (1 - R_2^2) \frac{n-1}{n-z}$

Rmk:  $\beta_2 = 0 \Rightarrow \bar{R}_2^2 \geq \bar{R}_{2+1}^2$

We call  $\bar{R}_2^2$  adjustment complex decision coefficient.

② For prediction:

Denote:  $X = \begin{pmatrix} X_1^T \\ \vdots \\ X_n^T \end{pmatrix} = \begin{pmatrix} X_{1q}^T & X_{1t}^T \\ \vdots & \vdots \\ X_{nq}^T & X_{nt}^T \end{pmatrix}$ ,  $X_i = \begin{pmatrix} X_{iq} \\ X_{it} \end{pmatrix}$

$$\Rightarrow X^T X = \begin{pmatrix} X_q^T X_q & X_q^T X_t \\ X_q^T X_t & X_t^T X_t \end{pmatrix}$$

$$X_k^T X_j = \sum_{i=1}^n X_{ik} X_{ij} \quad k, j \in \{q, t\}$$

i) Minimize the estimator of  $JJ_t$ :

$$JJ_t = \sum_{i=1}^n \text{Var}(Y_i - X_{iq}^T \beta_2) = \sum_{i=1}^n (1 + X_{iq}^T (X_q^T X_q)^{-1} X_{iq}) \sigma^2$$

Note that  $RMS = n\sigma^2 + \text{tr}((X_q^T X_q)^{-1} \sum_{i=1}^n (X_{iq} X_{iq}^T) \sigma^2)$   
 $= (n+q) \sigma^2$

Replace  $\sigma^2$  by its estimator  $\hat{\sigma}_2^2$ .

$$\Rightarrow \text{minimize } \hat{JJ}_t = (n+q) \hat{\sigma}_2^2 = \frac{n+q}{n-2} SSE_2$$

ii) Minimize  $S_2 = \hat{\sigma}_2^2 / (n-q-1)$

It's from: consider  $\eta_j = \beta_0 + \sum_{i=1}^{q-1} \beta_i X_{ji} + \varepsilon_j = \bar{\beta}_0 + \sum \beta_i (X_{ji} - \bar{X}_j) + \varepsilon_j$

i.e.  $Y = P_n X \beta + (I - P_n) X \beta + \varepsilon$ . (Centralization),  $P_n = J_n J_n^T / n$

$$\hat{Y} = P_n Y + (I - P_n) X \hat{\beta}, \quad \hat{\beta} = (X^T (I - P_n) X)^{-1} X^T (I - P_n) Y$$

$\Rightarrow$  estimate  $\hat{\eta} = \bar{\eta} + (X - \bar{X})^T \hat{\beta}$ ,  $\text{Var}(\eta - \hat{\eta}) = (\frac{n+1}{n} + (X - \bar{X})^T S^{-1} (X - \bar{X})) \sigma^2$

where  $S = X^T (I - P_n) X = \sum (X_k - \bar{X})(X_k - \bar{X})^T$ ,  $x$  from  $X$ .  $\begin{pmatrix} Y \\ X \end{pmatrix} \sim N(n, \Sigma)$

Note that  $\frac{n-2}{2-1} \frac{n}{n+1} (X - \bar{X})^T S^{-1} (X - \bar{X}) \sim F_{q-1, n-q-1}$

$$\Rightarrow E(\text{Var}(\eta - \hat{\eta} | X, X)) = \frac{n+1}{n} \frac{n-2}{n-2-1} (\sum Y - \sum Y X \Sigma_{XX}^{-1} \sum X Y)$$
, set  $S_2 = \frac{\hat{\sigma}_2^2}{n-2}$

iii)  $C_p$  Statistics:

Consider to minimize  $J_p = \frac{1}{\sigma^2} \sum_i (\hat{\eta}_i - E(y_i))^2$

$= \frac{1}{\sigma^2} \sum (x_i^T \hat{\beta}_p - x_i^T \beta)^2$ . Take its expectation:

$$E(J_p) = \frac{1}{\sigma^2} \sum_i \text{Var}(x_i^T \hat{\beta}_p) + (E(x_i^T \hat{\beta}_p) - x_i^T \beta)^2$$

$$= \sum x_i^T (X_2^T X_2)^{-1} x_i + \frac{1}{\sigma^2} \sum (x_i^T (\beta_0 + A\beta_1) - x_i^T \beta)^2$$

$$= q + \frac{1}{\sigma^2} \sum (x_i^T A\beta_1 - x_i^T \beta)^2$$

$$= q + \frac{1}{\sigma^2} \beta_1^T \left( \sum (A^T x_i - x_i) (A^T x_i - x_i)^T \right) \beta_1$$

$$= q + \frac{1}{\sigma^2} \beta_1^T D^{-1} \beta_1$$

$$= q + \frac{1}{\sigma^2} (n-q) (E(\hat{\sigma}_2^2) - \sigma^2)$$

$$\text{(since } E(\hat{\sigma}_2^2) = \sigma^2 + \frac{1}{n-q} \beta_1^T X_1^T (I - M_2) X_1 \beta_1)$$

$$\Rightarrow E(J_p) = \frac{E(SS E_2)}{\sigma^2} + 2q - n$$

$$\text{Choose } C_p = 2q - n + \frac{SS E_2}{\hat{\sigma}^2}$$

Remark:  $E(C_p) = 2 - t + \frac{n-p}{n-p-2} (t + \beta_1^T D^{-1} \beta_1 / \sigma^2)$

if  $n-p$  is large, when  $\beta_1 = 0$

then  $E(C_p) \approx 2$ .

$\Rightarrow$  Choose  $q$  st.  $C_p$  is small and

$|C_p - 2|$  is small.

iv) Minimize PRESS:

Denote:  $Y^{(-i)} = \begin{pmatrix} \eta_1 \\ \vdots \\ \eta_{i-1} \\ \eta_{i+1} \\ \vdots \\ \eta_n \end{pmatrix}$  ,  $X^{(-i)} = \begin{pmatrix} X_1^T \\ \vdots \\ X_{i-1}^T \\ X_{i+1}^T \\ \vdots \\ X_n^T \end{pmatrix}$  (remove  $i^{\text{th}}$  component)

Consider  $Y^{(-i)} = X^{(-i)} \beta + \varepsilon$

$\hat{\beta}^{(-i)} = (X^{(-i)T} X^{(-i)})^{-1} X^{(-i)T} Y^{(-i)}$  ,  $\hat{\varepsilon}^{(-i)} = \eta_i - X_i^T \hat{\beta}^{(-i)}$

Def: PRESS =  $\sum_{i=1}^n (\hat{\varepsilon}^{(-i)})^2$  , prediction of sum of square of error.

Denote:  $h_{ii} = X_i^T (X^T X)^{-1} X_i$   $i^{\text{th}}$  diagonal element of  $X(X^T X)^{-1} X^T$

To calculate PRESS:

$$\begin{aligned} \hat{\varepsilon}^{(-i)} &= \eta_i - X_i^T (X^{(-i)T} X^{(-i)})^{-1} (X^{(-i)T} Y^{(-i)}) \\ &= \eta_i - X_i^T (X^T X - X_i X_i^T)^{-1} (X^T Y - X_i \eta_i) \end{aligned}$$

$$(X^T X - X_i X_i^T)^{-1} = (X^T X)^{-1} + \frac{(X^T X)^{-1} X_i X_i^T (X^T X)^{-1}}{1 - h_{ii}} \quad (\text{easy to check})$$

$$\therefore \hat{\varepsilon}^{(-i)} = \hat{\varepsilon}_i + h_{ii} \eta_i - \frac{h_{ii} X_i^T \hat{\beta}}{1 - h_{ii}} + \frac{h_{ii} \eta_i}{1 - h_{ii}} \quad (\hat{\varepsilon}_i = \eta_i - X_i^T \hat{\beta})$$

$$= \frac{\hat{\varepsilon}_i}{1 - h_{ii}}$$

RMK: For  $PRESS_2 = \sum_{i=1}^n (\hat{\varepsilon}_2^{(-i)})^2$  . We have:

$$\hat{\varepsilon}_2^{(-i)} = \frac{\hat{\varepsilon}_{i2}}{1 - h_{i2}} \quad h_{i2} \text{ is } i^{\text{th}} \text{ diagonal element}$$

$$\text{of } M_2 = X_2 (X_2^T X_2)^{-1} X_2^T$$

③ By MLE:

i) AIC: (Akaike information criteria)

If  $I \sim N(0, \sigma^2 I_n)$ . Then:

$$\ln L(\beta, \sigma^2 | Y) = -\frac{n}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} (Y - X\beta)^T (Y - X\beta)$$

$$\ln L_{\max} = \ln L((X^T X)^{-1} X^T Y, \frac{SSE}{n})$$

$$= -\frac{n}{2} \ln\left(\frac{2\pi}{n}\right) - \frac{n}{2} \ln(SSE) - \frac{r}{2}$$

Def:  $AIC = -2 \log L(\hat{\theta} | X) + 2p$ ,  $p = \dim \theta$ .

$$\Rightarrow AIC_q = n \ln SSE_q + 2q \text{ in this case}$$

Find  $q$  to minimize  $AIC_q$

ii) BIC (Bayesian information criteria)

$$BIC_q = n \ln SSE_q + 2q \ln n$$

(3) Selection:

By above, we have three common criteria:

i)  $R^2$  criteria:  $\max_q \bar{R}_q^2$

ii)  $C_p$  criteria:  $\min_q C_p = \min_q \left( \frac{SSE_q}{SSE_p / (n-p)} + 2q - n \right)$

iii)  $AIC$ :  $\min_q n \ln SSE_q + 2q$

① Global Selection:

check  $2^p - 1$  possible regression models

for optimal  $q$ . But it's low efficient

when  $p$  is large.

e.g. choose  $A(k) = \{X_{1k}, \dots, X_{pk}\} \subset \{X_j\}_1^n$ .

Calculate criterion on  $A(k)$ .

## ② Stepwise method:

### i) Test for Significance:

Consider  $H_0: \beta_1 = \beta_2 = \dots = \beta_p = 0$ . i.e. test

whether  $\{X_j\}_1^p$  influences  $y$  a lot.

Denote:  $SST = Y^T (I - P_n) Y$ .  $df_T = n-1$ .

$SSR = Y^T (M - P_n) Y$ .  $df_R = p$

$SSE = Y^T (I - M) Y$ .  $df_E = n-p-1$

$\Rightarrow SST = SSR + SSE$ .  $SSE$  indep with  $SSR$

$$F = \frac{SSR/p}{SSE/(n-p-1)} \sim F(p, n-p-1) \text{ under } H_0.$$

We obtain  $p$ -value:  $p = P(F(p, n-p-1) > F)$

Consider  $H_0: \beta_{i_0} = 0$ . i.e. test the influence of individual variable.

Denote  $RSS^{(-i)}$  is  $RSS$  on  $\{X_k\}_{k \neq i}$ .

$$P_{i_0} = RSS - RSS^{(-i)} = \hat{\beta}_{i_0}^2 / l_{ii} \quad (*) \text{ (prove it later)}$$

where  $l_{ii}$  is the  $i$ th diagonal element of

$$L^{-1} = (X^T (I - P_n) X)^{-1}$$

$$F_{i_0} = \frac{P_{i_0}}{SSE/(n-p-1)} \sim F(1, n-p-1) \text{ under } H_0$$

Hint:  $SST = SSE + SSR = SSE^{(-i)} + SSR^{(-i)} \Rightarrow P_{i_0} = Q_{(-i)} - a$

## ii) Forward Selection:

Establish  $p$  regression equations with one variable. Then calculate each  $p$ -value ( $F \sim F_{(1, n-2)}$ )

$\Rightarrow$  Choose the variable with smallest  $p$ -value and introduce it into equation. (Denote  $x_1$ )

$\Rightarrow$  Establish  $p-1$  regression equations w.r.t  $\{x_1, x_i\}_{i=2}^p$ . Choose the pair with smallest  $p$ -value.

$\Rightarrow$  Repeat the process until we obtain the target number of variables.

Rmk: The later introduction of variables may reduce the significance of former variables.

## iii) Backward Selection:

Put all the variables into the equations

$\Rightarrow$  Calculate  $F_i$ . Discard the max one.

$\Rightarrow$  Put  $p-1$  variables into equations. then repeat the process before.

Rmk: It needs a lot of computations.

## iv) Stepwise Selection:

Step one: Discard variable. Suppose we have introduced  $\{x_{ik}\}_{k=1}^i$ .

Calculate  $P_{ik} = SSR_{\{ij\}_{j=1}^r} - SSR_{\{ij\}_{j \neq k}}$ . if

$X_0$  is the variable corresponds  $\min_{1 \leq k \leq r} P_{ik}$ .

then test:  $F_0 = \frac{P_0}{SSE/(n-r-1)}$ .  $P = P(F(1, n-r-1) > F_0 | \beta_0 = 0)$

if  $p \geq \alpha_{out} \Rightarrow$  discard  $X_0$

$p < \alpha_{out} \Rightarrow$  consider to introduce other variables

Step two: introduce new variables. suppose we have

$\{X_{jk}\}_{k=1}^{p-r}$  out of the equation.

$P_{jk} = SSR_{\{ij\}_{j=1}^r \cup \{jk\}} - SSR_{\{ij\}_{j=1}^r}$ . if

$X_{j_0}$  is the variable correspond  $\max_{1 \leq k \leq p-r} P_{jk}$

$F_{j_0} = \frac{P_{j_0}}{SSE_{\{ij\}_{j=1}^r \cup \{jk\}}/(n-r-2)}$   $\sim F(1, n-r-2)$

$P = P(F(1, n-r-2) > F_{j_0} | M_0)$ .  $M_0 = \beta_{j_0} = 0$ .

if  $p < \alpha_{out} \Rightarrow$  introduce  $X_{j_0}$  into equation

$p \geq \alpha_{out} \Rightarrow$  selection is over.

Pf of (\*): 
$$\begin{cases} Q = Y^T (I - X(X^T X)^{-1} X^T) Y \\ Q_{(i)} = Y^T (I - X_{(i)}(X_{(i)}^T X_{(i)})^{-1} X_{(i)}^T) Y \end{cases} \quad X = (X_1, X_i)$$

$$(X^T X)^{-1} = \begin{pmatrix} (X_1^T X_1)^{-1} + A D^{-1} A^T & -A D^{-1} \\ -D^{-1} A^T & D^{-1} \end{pmatrix}$$

$$A = (X_1^T X_1)^{-1} X_1^T X_i, \quad D = X_i^T (I - P_{(X_1)}) X_i = 1/c_{ii}$$

$$\Rightarrow P_i = Y^T (A^T X_1^T - X_i^T)^T D^{-1} (A^T X_1^T - X_i^T) Y$$

$$\hat{\beta}_i = (0 \ I) (X^T X)^{-1} X^T Y = (-A^T X_1^T + X_i^T) Y / D$$

$\Rightarrow P_i$  indept with  $Q$ , as well.