Chapter 9. Variable Selection and Model Building



We want: 1) bias errors to be small for prediction purposes

— include as many X's as possible.

2) variance of the prediction to be small

— i.e. $\sum_{i=1}^{n} Var \hat{y}_i/n = \frac{p}{n}\sigma^2$, small — p to be small.

Seek for a compromise between 1) and 2) - no unique procedure.

Usually iterative approach is employed including (1) a particular variable selection criterion

(2) diagnostic check for the resulting model.

Assume there are K candidate regressors, x_1, \ldots, x_K .

Full model :
$$y_i = \beta_0 + \sum_{j=1}^{K} \beta_j x_{ij} + \epsilon_i, \quad i = 1, \dots, n,$$

or

$$\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\epsilon}.$$

- **- - E** → - •

Let r be the number of regressors that are deleted from the full model. Rewrite the full model as

$$\mathbf{y} = X_p \beta_p + X_r \beta_r + \epsilon, \ X = (X_p, X_r), \beta = (\beta'_p, \beta'_r)'$$

such that the subset model is

$$\mathbf{y}=X_p\boldsymbol{\beta}_p+\boldsymbol{\epsilon}.$$

Then the LSE of
$$\beta$$
 is $\hat{\beta}^{\star} = (X'X)^{-1}X'\mathbf{y} = \begin{pmatrix} \hat{\beta}^{\star}_p \\ \hat{\beta}^{\star}_r \end{pmatrix}$ and we use

 $\hat{\sigma_{\star}^2}$ to denote the estimate of σ^2 ;

 \hat{y}_i^* to denote the fitted values in full model.

Also, let $\hat{\beta}_p = (X'_p X_p)^{-1} X'_p \mathbf{y}$ be the LSE of β_p in the subset model; $\hat{\sigma}^2$ and \hat{y}_i for estimate of σ^2 and fitted values. **<u>Results</u>**: 1. $\hat{\beta}_p$ is a *biased* estimate of β_p (unless $\beta_r = 0$ or $X'_{n}X_{r}=0$). 2. $Var(\hat{\beta}_{p}^{\star}) - Var(\hat{\beta}_{p})$ is p.s.d. (" \geq 0"). i.e. deleting variables never increases the variances of the estimates of the (remaining) parameters. 3. $Var(\hat{\beta}_{p}^{\star}) - MSE(\hat{\beta}_{p})$ " ≥ 0 ", if $Var(\hat{\beta}_{r}^{\star}) - \beta_{r}\beta_{r}'$ " > 0". ($|\beta_r| < s.e.(\hat{\beta}_r^*)$ in one-dimension) $\therefore \hat{\beta}_{p}$ has smaller MSE if $Var(\hat{\beta}_{r}^{\star}) - \beta_{r}\beta_{r}^{\prime}$ ">0".

Results: 4.
$$E(\hat{\sigma}_{\star}^2) = \sigma^2$$
 but $E(\hat{\sigma}^2) \ge \sigma^2$.
5. $Var(\hat{y}^{\star}) \ge MSE(\hat{y})$ if $Var(\hat{\beta}_r^{\star}) - \beta_r \beta_r' \text{ "} \ge 0$ ".

<u>Conclusions</u>: 1. Deleting variables may improve the precision, but potentially introduces bias.

 Retaining negligible variables may increase the variances of the estimates of the parameters and the predicted response. 1) R^2 . The R^2 in a *p*-term subset model including β_0 is

$$R_p^2 = \frac{SS_R(p)}{SS_T} = 1 - \frac{SS_{Res}(p)}{SS_T}.$$

Note: 1. $R_p^2 \uparrow$ in p and attains maximum when p = K + 1. 2. There are $\begin{pmatrix} K \\ p-1 \end{pmatrix}$ subset models of size p. 3. Look at the point where an additional variable is not useful. i.e. it provides only a small increase in R_p^2 .



▲□▶ ▲□▶ ▲目▶ ▲目▶ 三日 - 釣�?

| The Halu Cement Data. N = 4. | | | |
|------------------------------|--|---|---|
| I | П | III | IV |
| .675(4) | .979 (12)* | .98234(124)* | .98237(1234) |
| .666(2) | .972(14)* | .98228(123)* | |
| .534(1) | .935(34) | .98128(134)* | |
| .286(3) | .847(23) | .97282(234)* | |
| | .680(24) | | |
| | .548(13) | | |
| | .675(4) .666(2) .534(1) .286(3) | I II .675(4) .979 (12)* .666(2) .972(14)* .534(1) .935(34) .286(3) .847(23) .680(24) .548(13) | IIIIII.675(4).979 (12)*.98234(124)*.666(2).972(14)*.98228(123)*.534(1).935(34).98128(134)*.286(3).847(23).97282(234)*.680(24).548(13) |

<u>Ex.</u> The Hald Cement Data. K = 4.



(14) or (12) is the best, but (4) is the best when p = 2. Read Page 271 in text.

Note: 4. Large R_p^2 is preferred. Choose subsets of regressors producing an R^2 greater than

$$\begin{split} R_0^2 &= 1 - (1 - R_{K+1}^2)(1 + d_{\alpha,n,K}), \text{ where} \\ d_{\alpha,n,K} &= K \; F_{K,n-K-1;\alpha}/(n-K-1), \text{ called } R_{adequate}^2(\alpha) \\ \text{subsets.} \end{split}$$

<u>Ex.</u> (cont'd)

$$R_0^2 = 1 - (1 - R_5^2)(1 + \frac{5F_{4,8;.05}}{8}) = .94885.$$

▲ロ ▶ ▲ 聞 ▶ ▲ 臣 ▶ ▲ 臣 ▶ ▲ 回 ▶ ▲ 回 ▶



æ



ヘロト ヘヨト ヘヨト ヘヨト

- **<u>Note</u>**: 1. For each p, there are $\begin{pmatrix} K \\ p-1 \end{pmatrix}$ possible subset models. 1) Average over all possibilities for each p.
 - 2) Find the smallest p that 'begins' to get close to $MS_{Res}(full).$
 - 3) Choose the subset of size p that yields

 $MS_{Res} \approx MS_{Res}$ (full).

Note: 2.
$$R_{Adj,p}^2 = 1 - \frac{n-1}{n-p}(1-R_p^2) = 1 - \frac{n-1}{n-p}(\frac{SS_{Res}(p)}{SS_T})$$

= $1 - \frac{n-1}{SS_T}MS_{Res}(p)$.
 $\therefore \min MS_{Res}(p)$ and $\max R_{Adj,p}^2$ are equivalent.
Ex.(Cont'd) Again (14) or (12). See P. 274.

|▲□▶ ▲圖▶ ▲≣▶ ▲≣▶ | 差||||の Q @

3) Mallow's C_p statistics.

Recall that $E(SS_{Res}(p)) = (n-p)\sigma^2 + \sum_{i=1}^n B_i^2$.

Define

$$C_p = SS_{Res}(p)/\hat{\sigma^2} - (n-2p), \text{ where } \hat{\sigma^2} = MS_{Res}(\mathsf{full}).$$

<u>Note</u>: 1. If the model is adequate (unbiased), $\sum_{i=1}^{n} B_i^2 = 0$ $\implies C_p \approx (n-p)\sigma^2/\sigma^2 - (n-2p) = p$, ideally. 2. Special case: p = K + 1, full model (unbiased).

$$C_p \equiv (n - K - 1) - (n - 2(K + 1)) = K + 1 = p.$$

3. If the model is biased, $\sum_{i=1}^{n} B_i^2 > 0$.

$$C_{p} \approx \frac{(n-p)\sigma^{2} + \sum B_{i}^{2}}{\sigma^{2}} - (n-2p) = \frac{\text{bias}^{2}}{\sigma^{2}} + \frac{p\sigma^{2}}{\sigma^{2}}$$
$$= (\text{bias}^{2} + \sum_{i=1}^{n} Var\hat{y}_{i})/\sigma^{2}$$
$$= MSE_{p}/\sigma^{2}.$$

... we want to choose a model with $C_p \approx p$ (such that bias ≈ 0) and p as small as possible (such that predicted variance is small).





Theoretically $E(C_p) \ge p$, but due to random variation C_p may

э

_ ▶ <

4) $PRESS_p$.

Recall

$$PRESS_p = \sum_{i=1}^{n} (y_i - \hat{y}_{(i)})^2 = \sum_{i=1}^{n} \left(\frac{e_i}{1 - h_{ii}}\right)^2$$
, in a *p*-term model.

The model with smaller $PRESS_p$ values are preferrable.

Computational Techniques

I Best Subset Regressions.

1) Specify M = number of best subsets considered.

2) According to the above criteria, give a list of the selected best subsets for each p.

- <u>Note</u>: 1. Not all 2^K equations are considered, so may leave out sensible ones.
 - Unless K is small, otherwise, it is impossible to make a detailed examination of *all* possible regressions.

Read pp. 270-277 in text.

II Stepwise-Type Procedures

- (1) Stepwise regression:
 - i) Fit a simple linear regression for each of the *K* variables. Test: slope=0. (i.e. $y = \beta_0 + \epsilon$ versus $y = \beta_0 + \beta_k x_k + \epsilon$.) Compute $F_k = \frac{SS(\beta_k | \beta_0)}{SS_{Res}(\beta_k, \beta_0)/(n-2)}, \quad k = 1, 2, ..., K$.

The regressor with the largest F value is the *candidate* for first addition. Without loss of generality, assume $F_1 = \max_k F_k = F^*$. If $F^* > F_{1,n-2;\alpha} \equiv F_{IN}$, add into the model; otherwise, *stop* with no variable entering the model. ii) Suppose now $y = \beta_0 + \beta_1 x_1 + \epsilon$.

Compute the partial-*F* statistics for each $k \neq 1$.

$$F_{k|1} = \frac{SS(\beta_k|\beta_1,\beta_1)}{MS_{Res}(\beta_k,\beta_1,\beta_0)} = \left(\frac{\hat{\beta}_k}{s.e.(\hat{\beta}_k)}\right)^2 = t_k^2, k \neq 1.$$

If $F_{2|1} = \max_k F_{k|}$, then test

$$H_0: \left(egin{array}{c} eta_1\ eta_2\end{array}
ight) = 0 \ \ (\ ext{versus full}: \ y = eta_0 + eta_1 x_1 + eta_2 x_2 + \epsilon).$$

 x_2 enters the model if the test is significant; otherwise, stop.

iii) Suppose now $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$.

Examine whether *any* of the other variables already in the model should be *dropped*.

Compute $F_{2|1}$ and $F_{1|2}$.

Let $F^{\star} = \min(F_{2|1}, F_{1|2})$. If $F^{\star} < F_{OUT}$, insignificant, then

the corresponding variable is removed from the model.

iv) Examine the next candidate for addition, then examine *any* of the variables already in the model should be deleted.

- <u>Note</u>: 1. The stepwise regression allows an *x*-variable brought into the model at an early stage to be dropped subsequently if it is no longer helpful in conjunction with variables added at the later stages.
 - 2. The partial *F* test can be replaced by *t*-test (: $F_{1,\nu} = t_{\nu}^2$).
 - The test values for the partial F's are often called "F to enter" and "F to remove".

<u>Ex.</u> Hald Cement Data (Cont'd). K = 4, n = 13.

1)
$$F_{1|-} = 12.6, F_{2|-} = 21.96, F_{3|-} = 4.4, F_{4|-}^{\star} = 22.8$$
. First to
enter? x_4 ? $y = \beta_0 + \beta_4 x_4 + \epsilon$?
 $F^{\star} = 22.8 > F_{1,,11;.05} = 4.84, \therefore x_4$ enters.

2) Next to enter?
$$F_{.|4} = ?$$

 $F_{1|4}^{\star} = 108.22, F_{2|4} = 0.17, F_{3|4} = 40.29. x_1?$
(ls $y = \beta_0 + \beta_4 x_4 + \beta_1 x_1 + \epsilon$ significant?)
Test $H_0 = \begin{pmatrix} \beta_4 \\ \beta_1 \end{pmatrix} = 0$ in $y = \beta_0 + \beta_4 x_4 + \beta_1 x_1 + \epsilon$.
 $F^{\star} = 176.63$, significant! $\therefore x_1$ enters.

- 3) Next to exit? $y = \beta_0 + \beta_1 x_1 + \beta_4 x_4 + \epsilon$, x_1 or x_4 ? $F_{4|1} = 159.30, F_{1|4}^{\star} = 108.22 > F_{1,10;.05} = 4.96$. Both are significant. \therefore no exit.
- 4) Next to enter? $F_{.|14} = ?$

$$F_{2|14}^{\star} = 5.03, F_{3|14} = 4.24, x_2?$$

Test significance for $y = \beta_0 + \beta_1 x_1 + \beta_4 x_4 + \beta_2 x_2 + \epsilon$. Overall

F = 166.83, significant! x_2 enters.

5) Next to exit? x_1 , x_2 or x_4 ? $F_{1|24} = 154.01, F_{2|14} = 5.03, F_{4|12}^{\star} = 1.86 < F_{1,9;.05} = 5.12,$ insignificant!

 \therefore x₄ should be removed from the model.

$$\therefore y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon.$$

 Next to enter? Only x₃? F_{3|12} = 1.83 < 5.12, insignificant! No entrance!

7) Exit?
$$x_1$$
 or x_2 ?
 $F_{1|2}^{\star} = 146.52, F_{2|1} = 208.58$, significant! No exit!
 $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$,

Overall F = 229.5, significant! Done!

$$\hat{y} = 52.58 + 1.468x_1 + 0.662x_2.$$

Note: x_4 entered the model in the earliest stage, but was out later!

(2) **Forward selection**. Only the most recent *entrance* is tested.

- i) $F_{1|-}^{\star} = \max_k F_{k|-}$, k = 1, say. Significant! x_1 is added.
- ii) Find $\max_{k\neq 1} F_{k|1}$, say k = 2. x_2 enters?

Test $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$. Overall $F(x_1, x_2)$ say. Significant!

Test H_0 : $\beta_2 = 0$ in $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$.

 $F_{2|1}$, significant. $\Longrightarrow x_2$ enters. (Partial-F.)

iii) Check $\max_{k \neq 1,2} F_{k|12}$, say $k = 3, x_3$?

Overall significant, and $F_{3|12}$ significant, enters!

iv) :

Until no entrance!

<u>**Ex.</u>(Cont'd)**.</u>

1)
$$F_{4|-}^{\star} = 22.8 = \max_{k} F_{k|-} > F_{1,11;.05} = 4.84$$
, significant!
 $y = \beta_0 + \beta_4 x_4 + \epsilon$.

2)
$$F_{1|4}^{\star} = 108.22$$
 is the largest. x_1 ?
(i) Test $H_0: (\beta_1, \beta_4) = (0, 0)$ vs. $H_1: \beta_1 \neq 0$ or $\beta_4 \neq 0$, significant!

(ii) Partial
$$F = F_{1|4}$$
, significant, $\therefore x_1$ enters.

3)
$$F_{2|14}^{\star} = 5.03 > F_{3|14} = 4.24$$
, x_2 ?
(i) Overall: significant! (ii) $F_{2|14} < F_{1,9;.05} = 5.12$,
insignificant. No entrance!

 \Box

æ –

Stop!
$$y = \beta_0 + \beta_1 x_1 + \beta_4 x_4 + \epsilon$$
.

- Note: 1. Only the last possible entrance is tested.
 - 2. Not recommended unless it is specially desired never to remove variables that were retained in the model.

(3) Backward elimination.

i) Include all variables in the model.

- ii) Find $F^* = \min_k F_{k|\text{all others}}$, say k = K. x_K removed? If $F^* < F_{OUT}$, insignificant, x_K is eliminated.
- iii) Model $y = \beta_0 + \beta_1 x_1 + \dots + \beta_{k-1} x_{k-1} + \epsilon$, overall significant!

Go to ii) until no variable is deleted.

<u>Ex.</u>(Cont'd).

▲□ > ▲□ > ▲ ≧ > ▲ ≧ > ▲ ≧ > りへぐ

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$$
, done!

<u>Note</u>: Once a variable has been removed, it is gone *forever*. Thus, all alternative models using the eliminated variables are not considered for possible examination.

Significance levels.

- If a large α is selected, more x-variables would be admitted (Large rejection region).
- 2 Most often, use the same α in testing for entrance and exit.
- One may like to set the 'exit α' > 'entry α' to 'protect' predictors already in the model.
- One may always use $F_{IN} = F_{OUT} = 4$.

- <u>Note</u>: 1. All the procedures discussed do not necessarily select the best model, but an acceptable one.
 - Run the stepwise regression procedure to determine the number of variables used, say q. Do all possible sets of size q, and choose the best one.
- **<u>Ex</u>**. Two competitive models for q = 2. x_1, x_2 or x_1, x_4 ??? Need apriori consideration or the experimenter's judgement. \Box

Homework 9: (Page 300) 9.1, 9.2, 9.7, 9.8, 9.12, 9.13, 9.17, 9.18, 9.19, 9.20.

Due: Jan. 7, 2009.