

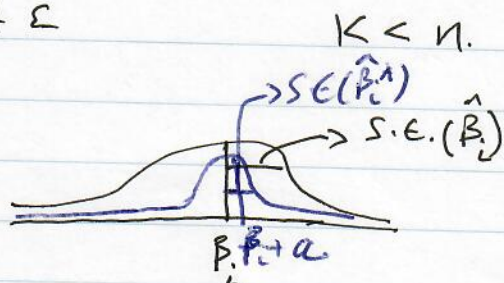
Notes

(1)

$$y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n + \varepsilon$$

$$y = \beta_0^* + \beta_1^* X_1 + \beta_2^* X_2 + \dots + \beta_k^* X_k + \varepsilon$$

$$S.E.(\hat{\beta}_k^*) \leq S.E.(\hat{\beta}_i)$$

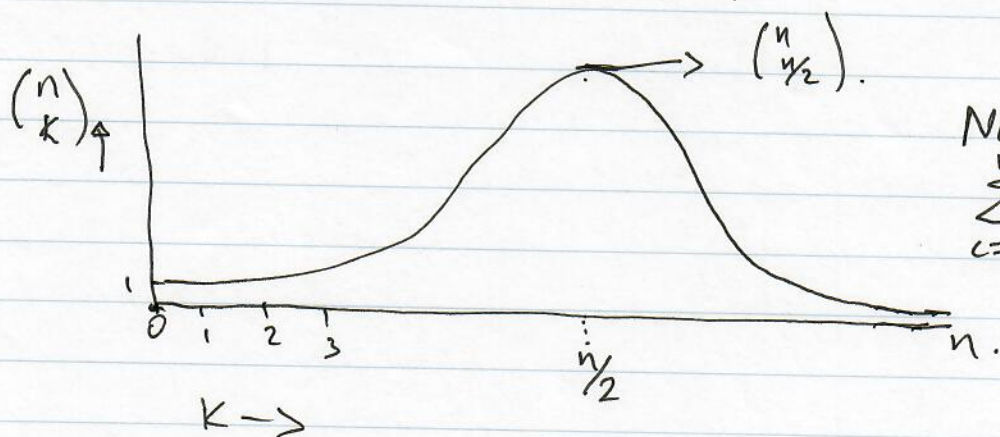


Finding all subset models with k variables it requires to assess $\binom{n}{k} = \frac{n!}{k!(n-k)!}$ subsets.

E.g.

$$k=1 \quad \binom{n}{k} = n \quad \text{and} \quad k=n-1 \quad \binom{n}{n-1} = n.$$

$$k=n \quad \binom{n}{n} = 1, \quad k=0 \quad \binom{n}{0} = 1. \quad (\text{note } 0! = 1)$$



Note:

$$\sum_{i=0}^n \binom{n}{i} = 2^n.$$

There are more combinations for $k = n/2$.

Comparing subsets with the same number of variables:

we can use R^2 or equivalently the RSS (Residual Sum of Squares).

C_p is used to compare subsets of different size.

FORWARD

Consider that there are n variables x_1, x_2, \dots, x_n .

First fit all single variable models, i.e.

$$y = \beta_0 + \beta_1 x_1 + \varepsilon \quad p^{(1)}$$

$$y = \beta_0 + \beta_1 x_2 + \varepsilon \quad p^{(2)}$$

$$\vdots$$

$$y = \beta_0 + \beta_1 x_n + \varepsilon \quad p^{(n)}$$

Find the p -value of β_1 in each model, say $p^{(i)}$ for subset $y = \beta_0 + \beta_1 x_i + \varepsilon$

Find the smallest p -value, say $p^{(2)}$, i.e.
 $p^{(2)} = \min(p^{(1)}, p^{(2)}, \dots, p^{(n)})$.

Assume that we have $\alpha\%$ (α percentage point, eg 5% , 10% , 1% , etc).

If $p^{(2)} < \alpha\%$ then x_2 is significant and we select the variable x_2 to enter the model.

If $p^{(2)} \geq \alpha\%$ then x_2 is insignificant, we stop and select the previous model which is $y = \beta_0 + \varepsilon$. (the empty model)

Assume that x_2 is significant. Then x_2 is retained for ever and we fit all remaining variables (one each time) with x_2 . That is $n-1$ regressions (subsets) are evaluated:

(3)

$$\begin{array}{l}
 y = \beta_0 + \beta_1 x_2 + \beta_2 x_1 + \varepsilon \quad p^{(1)} \\
 y = \beta_0 + \beta_1 x_2 + \beta_2 x_3 + \varepsilon \quad p^{(3)} \\
 \vdots \\
 y = \beta_0 + \beta_1 x_2 + \beta_2 x_n + \varepsilon \quad p^{(n)}
 \end{array}
 \left. \vphantom{\begin{array}{l} \\ \\ \\ \end{array}} \right\} \begin{array}{l} \text{P-values of } \beta_2 \text{ in} \\ \text{each of the } n-1 \\ \text{regressions.} \end{array}$$

Let $p^{(1)} = \min(p^{(1)}, p^{(3)}, \dots, p^{(n)})$, i.e. $p^{(1)}$ is the smallest p-value.

If $p^{(1)} \geq \alpha\%$ then x_1 is insignificant, we stop and select the previous model which is $y = \beta_0 + \beta_1 x_2 + \varepsilon$.

If $p^{(1)} < \alpha\%$, then x_1 is significant, thus, x_1 enters the model (for ever), and we consider all remaining variables (one at a time) in addition to (x_2, x_1) . That is, we consider the $n-2$ regressions

$$\begin{array}{l}
 y = \beta_0 + \beta_1 x_2 + \beta_2 x_1 + \beta_3 x_3 + \varepsilon \quad p^{(3)} \\
 y = \beta_0 + \beta_1 x_2 + \beta_2 x_1 + \beta_3 x_4 + \varepsilon \quad p^{(4)} \\
 \vdots \\
 y = \beta_0 + \beta_1 x_2 + \beta_2 x_1 + \beta_3 x_n + \varepsilon \quad p^{(n)}
 \end{array}
 \left. \vphantom{\begin{array}{l} \\ \\ \\ \end{array}} \right\} \begin{array}{l} n-2 \\ \text{P-values} \end{array}$$

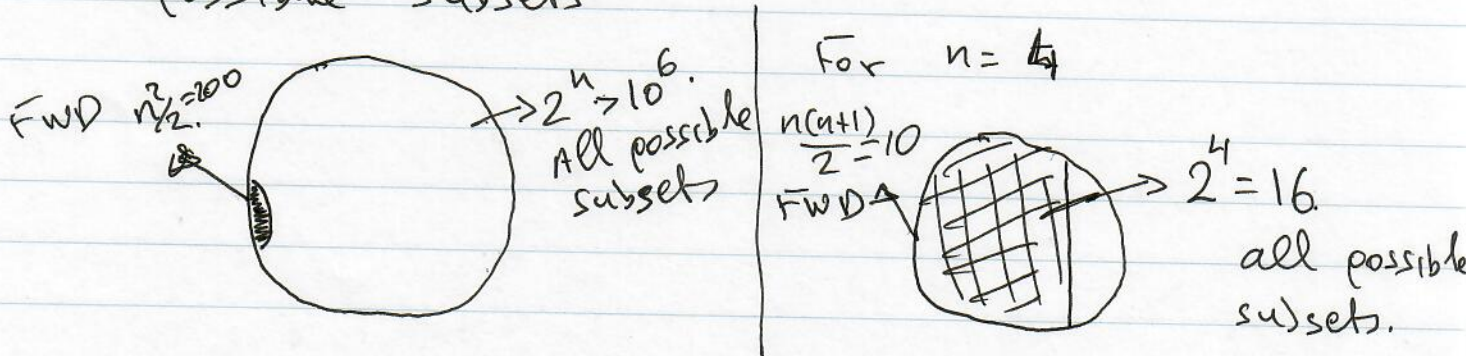
We repeat the same process until a p-value is found to be insignificant, or until all variables have been selected.

If all variables enter the model then the number of subsets (regressions) that have been evaluated are $n + (n-1) + (n-2) + \dots + (1)$

i.e. $\sum_{i=1}^n i = \frac{n(n+1)}{2} \approx \frac{n^2}{2}$

This implies that in the best case scenario the FORWARD method compares $n^2/2$ out of 2^n possible subsets.

E.g if $n = 20$, $n^2/2 = 200$ and $2^n > 10^6$ (1 million). Thus, the FORWARD checks 0.02% of all possible subsets.



BACKWARD

Fit all the variables, check the less significant variable and, (a) if the variable is significant then stop and select the model, or (b) if the variable is insignificant then delete it, and repeat the process. i.e.

$$y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n + \epsilon.$$

$p^{(1)}$ $p^{(2)}$ $p^{(n)}$

→ p-values of each variable

Let $p^{(2)} = \max(p^{(1)}, p^{(2)}, \dots, p^{(n)})$.

If $p^{(2)} < \alpha\%$, then keep X_2 and stop with the full model above.

If $p^{(2)} \geq \alpha\%$, then delete the variable X_2 (from ever) and continue without it, i.e. fit the model

$$y = \beta_0 + \underset{p^{(1)}}{\beta_1} X_1 + \underset{p^{(3)}}{\beta_2} X_3 + \dots + \underset{p^{(n)}}{\beta_n} X_n + \epsilon \rightarrow n-1 \text{ p-values.}$$

Repeat the process until you find an insignificant variable and STOP. Otherwise continue until all the variables are deleted.

If all variables are deleted then at most $n(n+1)/2$ or $n^2/2$ p-values have been assessed. The same as the FORWARD.

STEPWISE

Combines both FORWARD AND BACKWARD so that when a variable enters the model can be deleted at a latter stage, or when a variable ~~enters~~ ^{is deleted from} the model, then it can enter later

First apply the Forward to enter a variable and then a backward to delete all insignificant variables. This is repeated until the Forward does not enter a variable.

(6)

Example Stepwise

variables in model

candidate variables

step 0: (\emptyset)

(X_1, X_2, \dots, X_n)

step 1: FORWARD

$$y = \beta_0 + \beta_1 X_1 + \epsilon$$

$$\boxed{y = \beta_0 + \beta_1 X_2 + \epsilon}$$

$$y = \beta_0 + \beta_1 X_n + \epsilon$$

Suppose X_2 is the most significant and enter the model

model list (X_2)

candidate list

(X_1, X_3, \dots, X_n)

The Backward to $y = \beta_0 + \beta_1 X_2 + \epsilon$ will not result to any change.

Step 2:

FWD

$$\boxed{y = \beta_0 + \beta_1 X_2 + \beta_2 X_1 + \epsilon}$$

$$y = \beta_0 + \beta_1 X_2 + \beta_2 X_3 + \epsilon$$

$$y = \beta_0 + \beta_1 X_2 + \beta_2 X_n + \epsilon$$

X_1 has the smallest P-value and significant

BWD: $y = \beta_0 + \beta_1 X_2 + \beta_2 X_1 + \epsilon$. (nothing happens!)

Step 3:

FWD selects say X_n i.e. $y = \beta_0 + \beta_1 X_2 + \beta_2 X_1 + \beta_3 X_n + \epsilon$
BWD ~~may~~ ~~can~~ eliminate a variable
(only X_1 or X_2 but actually only X_2 ...).

The stepwise in the best case considers ~~n^2~~ n^2 subsets (double Pass FWD and BWD).

This is still small when n is not small.

