

Test on significant impact of variables

In a multiple linear regression model we have

$$E[Y] = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$$

Thus β_i , $i=1, 2, \dots, n$ is the change in $E[Y]$ if x_i is changed with one unit and all the rest of the variables are kept unchanged.

We want to test. Does a variable have a significant impact on the response given that the other variables are in the model.

such a test for x_j , $j=1, 2, \dots, k$ is

$$H_0: \beta_j = 0 \text{ against } H_1: \beta_j \neq 0$$

We know $\hat{\beta}$ is independent of $\hat{\sigma}^2 = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n-k-1}$ if

$$\epsilon \sim N(0, \sigma^2 I)$$

Let c_{jj} be the $(j+1)$ th diagonal element in $(X^T X)^{-1}$

$$\text{Then } T = \frac{\hat{\beta}_j}{\hat{\sigma} \sqrt{c_{jj}}} = \frac{\frac{\hat{\beta}_j}{\sigma}}{\frac{\sqrt{c_{jj}}}{\sigma}} = \frac{\frac{\hat{\beta}_j}{\sigma}}{\sqrt{\frac{\hat{\sigma}^2 (n-k-1)}{(n-k-1) \sigma^2}}} \sim \frac{N(0, 1)}{\sqrt{\frac{\chi^2(n-k-1)}{n-k-1}}} \text{ if } H_0 \text{ is true}$$

and therefore t-distributed with $n-k-1$ degrees of freedom. We reject if $|T_{obs}| \geq t_{\frac{\alpha}{2}, n-k-1}$

For the test $H_0: \beta_j = \beta_{j0}$ against $H_1: \beta_j \neq \beta_{j0}$, we use

$T = \frac{\hat{\beta}_1 - \beta_{10}}{SE_{\hat{\beta}_1}}$ the same way except that

$\frac{d}{2}$ is substituted by d if the test is one sided.

choice of a fitted model

Explanatory variables with small or no influence on the response can give the model a bad prediction ability.

$$\text{We have: } \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

$$\text{or } SS_T = SS_E + SSR$$

In a good model $y_i - \hat{y}_i$ should be small and $\hat{y}_i - \bar{y} \approx y_i - \bar{y}$, $i = 1, 2, \dots, n$. A measure of how much variation in the data that can be explained by the model is the coefficient of multiple determination.

$$R^2 = \frac{SS_R}{SS_T} = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}, \quad 0 \leq R^2 \leq 1$$

$R^2 = 0,84$ tells us that 84% of the variation in the data can be explained by the model.

$$\text{Note that } R^2 = \frac{SS_R}{SS_T} = \frac{SS_T - SS_E}{SS_T} = 1 - \frac{SS_E}{SS_T}$$

A problem with R^2 is that it will always increase when we increase the number of explanatory variables.

To avoid that we have introduced R^2_{adjusted} , defined as

$$R^2_{\text{adj}} = 1 - \frac{\frac{SSE}{m-k-1}}{\frac{SS_T}{m-1}} = 1 - \frac{(m-1)S^2}{SS_T}$$

Maximizing $R^2_{\text{adj}} \Leftrightarrow$ minimizing $\hat{\sigma}_{\text{LS}}^2 = S^2$

Other criteria

def $\hat{\sigma}^2 = \frac{SSE}{n}$ i.e. the ML estimator
smallest value of.

$AIC \stackrel{a}{=} n \ln(\hat{\sigma}^2) + 2(k+2)$ now $\hat{\sigma}^2$ is also counted as parameter.

$$\hat{\sigma}^2 = \frac{SSE}{n} = \frac{SSE}{m-k-1} \cdot \frac{m-k-1}{m} = \beta^2 \frac{(m-k-1)}{m}$$

$$\text{Therefore } AIC = n \left(\ln \hat{\sigma}^2 - \ln \left(\frac{m}{m-k-1} \right) + \frac{2(k+2)}{m} \right)$$

Smallest value of:

$$BIC \stackrel{a}{=} n \ln(\hat{\sigma}^2) + \underbrace{\ln(n)(k+2)}_{\text{full model}}$$

$$\text{Mallows' } Cp \stackrel{a}{=} \frac{SSE}{\beta^2} - (n - 2(k+1))$$

$$\text{Note } E[Cp] \approx m - (k+1) - m + 2(k+1) = k+1$$

if the model is correct.

Variable selection methods

These algorithms are normally performed as partial F-tests, adding and removing one factor at the time

Forward selection

1. Start with only β_0 in the model

2. Find $\max_j R(\beta_j) = \max_j \{ SSR(\beta_0, \beta_j) - SSR(\beta_0) \}$

3. If $\max_j \frac{R(\beta_j)}{\frac{SSE}{n-2}} = \frac{R(\beta_m)}{\frac{SSE}{n-2}} < f_{d, 1, n-2}$ stop, no variable

is entered into the model

4. If $\frac{R(\beta_m)}{\frac{SSE}{n-2}} \geq f_{d, 1, n-2}$, add x_m to the model

Find $\max_{j \neq m} R(\beta_j | \beta_m) = \max_{j \neq m} \{ SSR(\beta_0, \beta_m, \beta_j) - SSR(\beta_0, \beta_m) \}$ and.

go to step 3. If $\max_{j \neq m} \frac{R(\beta_j | \beta_m)}{\frac{SSE}{n-3}} < f_{d, 1, n-3}$ no variable

is entered. Otherwise Proceed in the same way.

Note that the degrees of freedom in the partial F-test is reduced by one for each variable that is entered.

Backward elimination

Define $\underline{\beta} \setminus \beta_j = (\beta_0, \dots, \beta_{j-1}, \beta_{j+1}, \dots, \beta_n)$

1. Start with all variables in the model
2. Find $\min \int R(\beta_j | \underline{\beta} \setminus \beta_j) = \min \{ SSR(\beta_0, \dots, \beta_n) - SSR(\beta_0, \dots, \beta_{j-1}, \beta_{j+1}, \dots, \beta_n) \}$
3. If $\frac{\min \int R(\beta_j | \underline{\beta} \setminus \beta_j)}{\frac{SSE}{m-k-1}} \geq f_{d,1,m-k-1}$ stop.

No variable is removed from the model

4. If $\frac{R(\beta_m | \underline{\beta} \setminus \beta_m)}{\frac{SSE}{m-k-1}} < f_{d,1,m-k-1}$ remove X_m

Find $\min_{j \neq m} R(\beta_j | \underline{\beta} \setminus \{\beta_m, \beta_j\}) = \min_{j \neq m} \{ SSR(\beta_0, \dots, \beta_{m-1}, \beta_{m+1}, \dots, \beta_n) - SSR(\underline{\beta} \setminus \{\beta_m, \beta_j\}) \}$

- If $\frac{\min_{j \neq m} R(\beta_j | \underline{\beta} \setminus \{\beta_m, \beta_j\})}{\frac{SSE}{m-k}} \geq f_{d,1,m-k}$ stop.

Otherwise proceed in the same way