

Partial F-test on Subsets of Coefficients

Suppose we wish to test the hypothesis

$$H_0: \beta_2 = \beta_3 = 0 \quad H_1: \beta_2 \text{ and } \beta_3 \text{ are } \cancel{\text{not}} \text{ both } 0.$$

in a regression model with 3 explanatory variables

$$\text{i.e. } \overset{0}{y} = \overset{0}{\beta_0} + \overset{0}{\beta_1} x_1 + \overset{0}{\beta_2} x_2 + \overset{0}{\beta_3} x_3 = 0$$

We may then compute

$SS_R(\beta_1, \beta_2, \beta_3 | \beta_0) = SSR$ for the model with 3 regression variables and $SS_R(\beta_1 | \beta_0) = R(\beta_1)$ for the model with only x_1 as regression variable

$$\text{let } R(\beta_2, \beta_3 | \beta_1) = SSR - R(\beta_1)$$

$$\text{The F statistics is then } F = \frac{\frac{R(\beta_2, \beta_3 | \beta_1)}{2}}{\frac{SSE}{m-4}} \sim F_{2, m-4}$$

\rightarrow full model

We reject H_0 when $F_{obs} \geq f_{\alpha, 2, m-4}$

In general: $H_0: \beta_{i+1} = \beta_{i+2} = \dots = \beta_k = 0$, $H_1: \text{at least one of the coefficients } \neq 0$

$$R(\beta_{i+1}, \dots, \beta_k | \beta_1, \dots, \beta_i) = SSR(\beta_1, \dots, \beta_k | \beta_0) - SSR(\beta_1, \dots, \beta_i | \beta_0)$$

and we reject H_0 if

$$\frac{\frac{R(\beta_{i+1}, \dots, \beta_k | \beta_1, \dots, \beta_i)}{k-i}}{\frac{SSE}{m-k-1}} \geq f_{\alpha, k-i, m-k-1}$$

This test can be used to test if one variable can be taken out and if the regression is significant.

In case we test on one variable for instance

$$H_0: \beta_2 = 0 \quad H_1: \beta_2 \neq 0 \text{ in a model}$$

with 3 regression variables we compute

$$R(\beta_2 | \beta_1, \beta_3) = SSR - R(\beta_1, \beta_3)$$

The F statistic is

$$\frac{\frac{R(\beta_2 | \beta_1, \beta_3)}{1}}{\frac{SSE}{m-4}} \sim F_{1, m-4}$$

This is equivalent to a t-test.

12.9 Sequential Methods for Model Selection and Multicollinearity.

We have multicollinearity when two or more columns are strongly correlated (almost linear dependent). The degree of multicollinearity can be measured by the correlation coefficient

$$r_{ij} = \frac{\sum_{k=1}^m (x_{ik} - \bar{x}_i)(x_{jk} - \bar{x}_j)}{\left(\sum_{k=1}^m (x_{ik} - \bar{x}_i)^2 \sum_{k=1}^m (x_{jk} - \bar{x}_j)^2 \right)^{\frac{1}{2}}}.$$

N.B. Since the regression variables are not considered random, this expression does not estimate a true correlation.

Multicollinearity occurs often when the variation range of one or more explanatory variables are small.

Example If $x_i \in (0.95, 1.05)$, x_i will be strongly correlated by 1 and x_i^2 .

The best measure of multicollinearity is given by the VIF-factor.

$$VIF_j = \frac{1}{1 - R_j^2} \text{ where } R_j^2 \text{ is the } \cancel{\text{multiple}} \text{ multiple}$$

determination coefficient that is obtained by regressing x_j on all the other regression variables.

VIF should be below 10. If R_j^2 is greater than 0.99 Minitab will remove variable x_j .

What can be done if there exist strong multicollinearity

Remove column(s) columns

Collect more data

Respecify the model $(x_1, x_2, x_3) \rightarrow \begin{cases} \frac{x_1 + x_2}{x_3} \\ x_1 \cdot x_2 \cdot x_3 \end{cases}$

Try to center variables: $(\bar{x}_i - \bar{\bar{x}})$: polynomial models

center and scale variables.

$$\frac{(\bar{x}_i - \bar{\bar{x}})}{s_{x_i}}$$

PCA regression

PLS regression

RIDGE regression

12.9 Sequential methods for variable selection

An alternative to best subset regression or to try all regression models with 0, 1, 2, ..., k variables

Forward selection

1. Start with β_0 in the model
2. Find $\max_j R(\beta_j) = \max_j SSR(\beta_j | \beta_0) = \max_j \{ SSR(\beta_0, \beta_j) - SSR(\beta_0) \}$

3. If $\max_j \frac{R(\beta_j)}{SSE} \in f_{d, 1, m-2}$, stop no variables in the model.

4. If $\max_j \frac{R(\beta_j)}{SSE} \geq f_{d, 1, m-2}$ add x_j^m to the model

Find $\max_{i \neq m} R(\beta_i | \beta_m) = \max_{i \neq m} SSR(\beta_i | \beta_0, \beta_m) = \max_{i \neq m} \{ SSR(\beta_0, \beta_m, \beta_i) - SSR(\beta_0, \beta_m) \}$

and continue the same way.

Backward ~~selection~~ elimination

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

1. Start with all the variables in the model

2. Find $\min_j R(\beta_j | \underline{\beta} \setminus \beta_j) = \min_j \{ SSR(\beta_0, \dots, \beta_n) - SSR(\beta_0, \dots, \beta_{j-1}, \beta_{j+1}, \dots, \beta_n) \}$

3. If $\frac{\min_j R(\beta_j | \underline{\beta} \setminus \beta_j)}{SSE} \geq \alpha_{1, m-k-1}$ stop, no variable is removed.

4. If $\frac{\min_j R(\beta_j | \underline{\beta} \setminus \beta_j)}{SSE} = \frac{R(\beta_m | \underline{\beta} \setminus \beta_m)}{SSE} < \alpha_{1, m-k-1}$

remove x_m

Let $\underline{\beta} = \{\beta_1, \beta_2, \dots, \beta_k\} \setminus \beta_m$, $k = k-1$ and continue until all the variables are significant.

Stepwise regression

1. Start like forward selection

Assume x_1 and x_2 are chosen to enter the model in steps 1 and steps 2. Let $\underline{\beta} = \{\beta_1, \beta_2\}$

2. Find $\min_{j=1,2} R(\beta_j | \underline{\beta} \setminus \beta_j) = R(\beta_m | \underline{\beta} \setminus \beta_m)$

and investigate if x_m should be removed as for backward elimination.

3. Continue as for forward selection, but test

in each step if any of the variable chosen to enter can be removed.