

Test for the significance of the regressor variables in a multiple regression

Suppose we have a set of k independent variables, x_1, \dots, x_k , and the dependent variable Y . We denote the observations as $(x_{1i}, x_{2i}, \dots, x_{ki}, y_i)$ for $i = 1(1)n$, the sample size. Now, a multiple linear regression equation

$$y_i = \alpha + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_k x_{ki} + \epsilon_i, \quad \dots \text{①}$$

for $i = 1(1)n$, is considered. In this model, ϵ_i is assumed to be normally distributed with mean 0 and variance σ^2_ϵ and ϵ_i 's are independent over $i = 1(1)n$, i.e. $\epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2_\epsilon)$.

Now, one may be interested to test whether all the k regressors, x_1, \dots, x_k , have significant presence in the regression model ①. To address the question, we do a statistical test for

$$H_0: \beta_j = 0 \quad \forall j = 1(1)k$$

which means that there is no dependence of Y on x_1, x_2, \dots, x_k .

Commonly, H_A can be considered as $H_A: \beta_j \neq 0$ for at least one j . For simplicity in the calculation of the least-square estimates of the parameters α and β_j 's, we can re-write the model ① as

$$y_i = \alpha' + \beta_1 (x_{1i} - \bar{x}_1) + \dots + \beta_k (x_{ki} - \bar{x}_k) + \epsilon_i \quad \dots \text{②}$$

$$\text{where } \bar{x}_j = \frac{1}{n} \sum_{i=1}^n x_{ji}.$$

$$\text{or, } y_i = \alpha' + \beta_1 x'_{1i} + \beta_2 x'_{2i} + \dots + \beta_k x'_{ki} + \epsilon_i, \quad \dots \text{②}$$

$$\text{where } x'_{ji} = x_{ji} - \bar{x}_j \quad \forall j = 1(1)k.$$

Using Least-square method, estimates of β_j 's are

$$\hat{\beta}_j = \frac{\sum_{i=1}^n (x'_{ji} y_i)}{\sum_{i=1}^n (x'_{ji})^2} = b_j, \text{ say.} \quad \dots \text{③}$$

Averaging the equation ② over i , we have $\bar{y} = \alpha' + \sum_{j=1}^k \beta_j \cdot 0$.

$$\Rightarrow \hat{\alpha}' = \bar{y}.$$

Now, the unrestricted residual sum of squares (SS) is

$$\begin{aligned}
S_1^2 &= \min_{\alpha', \beta_j} \sum_{i=1}^n \left(y_i - \hat{\alpha}' - \sum_{j=1}^k \beta_j \cdot (x'_{ji}) \right)^2 \\
&= \sum_{i=1}^n \left(y_i - \hat{\alpha}' - \sum_{j=1}^k b_j x'_{ji} \right) \quad [\text{replacing parameters by their least-square estimates}] \\
&= \sum_{i=1}^n \left[(y_i - \bar{y})^2 + \left\{ \sum_{j=1}^k b_j \cdot x'_{ji} \right\}^2 - 2(y_i - \bar{y}) \sum_{j=1}^k b_j \cdot x'_{ji} \right] \\
&= \sum_{i=1}^n (y_i - \bar{y})^2 + \sum_{i=1}^n \sum_{j=1}^k b_j^2 x'_{ji}^2 - 2 \sum_{j=1}^k b_j \sum_{i=1}^n y_i \cdot x'_{ji} \\
&\quad \left[\because 2 \sum_{j=1}^k b_j \sum_{i=1}^n x'_{ji} = 0 \text{ as } \sum_i x'_{ji} = 0 \right] \\
&= \sum_i (y_i - \bar{y})^2 + \sum_j b_j^2 \sum_i x'_{ji}^2 - 2 \sum_j b_j \sum_i x'_{ji} y_i \\
&= \sum_i (y_i - \bar{y})^2 - \sum_j b_j \sum_i x'_{ji} y_i \quad [\text{from (3), } b_j \sum_i x'_{ji}^2 = \sum_i x'_{ji} y_i] \\
&\equiv \sum_i (y_i - \bar{y})^2 - \sum_{j=1}^k b_j P_j \quad , \text{ where } P_j = \sum_{i=1}^n x'_{ji} y_i
\end{aligned}$$

Next, the restricted (i.e. under H_0) residual SS is

$$\begin{aligned}
S_2^2 &= \min_{\substack{\alpha', \beta_j \\ H_0}} \sum_{i=1}^n \left(y_i - \hat{\alpha}' - \sum_{j=1}^k \beta_j \cdot x'_{ji} \right)^2 \\
&= \sum_{i=1}^n (y_i - \bar{y})^2 , \text{ since under } H_0, \beta_j = 0 \forall j=1 \dots k
\end{aligned}$$

The d.f. of S_1^2 is $(n-k-1)$.

\because the first part of S_1^2 consists n observations on Y with one restriction to sample mean, \bar{y} , so it has $(n-1)$ d.f and second part consists k estimates b_1, b_2, \dots, b_k , that means k more restrictions]

The d.f. of S_2^2 is $(n-1)$.

$$\text{So, SS due to regression} = S_2^2 - S_1^2 = \sum_{j=1}^k b_j P_j = SSR, \text{ say.}$$

$$\text{SS due to error} = S_1^2 = SSE, \text{ say.}$$

$$\text{So, the d.f. of } SSR = \text{d.f.}(S_2^2) - \text{d.f.}(S_1^2) = n-1 - (n-k-1) = k$$

Under H_0 , SSR and SSE both follow χ^2 distribution with respective d.f.s. and they are independent.

$$\text{Hence, } F = \frac{SSR/k}{SSE/(n-k-1)} \sim F_{k, n-k-1}$$

Hence, if the observed F , say F_0 , is greater than $F_{\alpha; k, n-k-1}$. we reject H_0 at level α , otherwise we do not reject H_0 .

Associated ANOVA table

Source of variation	d.f.	SS	MS	F_0
Due to multiple linear regression	k	$SSR = \sum_{j=1}^k b_j P_j$	$MSR = \frac{SSR}{k}$	$F_0 = \frac{MSR}{MSE}$
Error	$n-k-1$	$SSE = S_1^2$	$MSE = \frac{SSE}{n-k-1}$	
Total	$n-1$	$\sum_i (y_i - \bar{y})^2 = S_2^2$	-	-

■ Test for the significance of a subset of regressor variables in a given multiple ^{linear} regression model with k no. of regressor variables:

Suppose we have a multiple ^{linear} regression model of the variable Y on a set of k independent regressor variables x_1, x_2, \dots, x_k and the regression equation has the following form

$$y_i = \alpha + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_k x_{ki} + e_i, \quad \dots \quad (1)$$

for $i=1(1)n$, where $e_i \stackrel{iid}{\sim} N(0, \sigma_e^2)$ for $i=1(1)n$, by assumption.

Now, one may be interested to test whether a subset of s no. of regressor variables has significance in the multiple linear regression model (1), where $s < k$.

To address the question, we do a statistical hypothesis test for

$$H_0: \beta_j = 0 \quad \forall j=1(1)s.$$

against H_A : not all β_j equal to 0 for $j=1(1)s$,

i.e. $\beta_j \neq 0$ for at least one j .

For simplicity in the calculation of the least-square estimates of the parameters, we can re-write the model (1) as

$$y_i = \alpha' + \beta_1(x_{1i} - \bar{x}_1) + \dots + \beta_k(x_{ki} - \bar{x}_k) + e_i \\ = \alpha' + \beta_1 x'_{1i} + \dots + \beta_k x'_{ki} + e_i \quad \dots \quad (2)$$

, where $x'_{ji} = x_{ji} - \bar{x}_j$, $\bar{x}_j = \frac{1}{n} \sum_{i=1}^n x_{ji}$ & $j=1(1)k$.

Therefore, under H_0 , equation (2) reduces to

$$y_i = \alpha' + \beta_{s+1} x'_{s+1,i} + \dots + \beta_k x'_{ki} + e_i \quad \dots \quad (3)$$

Now the unrestricted residual SS from model (2) is

$$S_1^2 = \min_{\alpha', \beta_j's} \sum_{i=1}^n \left(y_i - \alpha' - \sum_{j=1}^k \beta_j \cdot x'_{ji} \right)^2 \\ = \sum_{i=1}^n (y_i - \bar{y})^2 - \sum_{j=1}^k b_j \cdot P_j, \text{ where } P_j = \sum_{i=1}^n x'_{ji} y_i$$

[Derivation of S_1^2 is already stated in earlier test]

Next, the restricted (i.e. under H_0) residual SS is

$$S_2^2 = \min_{\substack{\alpha', \beta_j's \\ H_0}} \sum_{i=1}^n \left(y_i - \alpha' - \sum_{j=s+1}^k \beta_j \cdot x'_{ji} \right)^2 \\ = \min_{\alpha', \beta_j's} \sum_{i=1}^n \left(y_i - \alpha' - \sum_{j=s+1}^k b_j^* P_j \right)^2 \\ = \sum_{i=1}^n (y_i - \bar{y})^2 - \sum_{j=s+1}^k b_j^* P_j,$$

[following the derivation similar to the derivation for S_1^2]

where b_j^* 's are the least square estimates of β_j 's, respectively, for the restricted model (3).

The degrees of freedom of S_1^2 and S_2^2 are $(n-k-1)$ and $(n-(k-s)-1)$, respectively.

Thus, SS due to regression involving s regressor variables

$$= S_2^2 - S_1^2 \\ = \sum_{j=1}^k b_j \cdot P_j - \sum_{j=s+1}^k b_j^* \cdot P_j \quad \text{with df.} = (n-k-1) - (n-k-s-1) \\ = S \\ = SSR_{s/k}, \text{ say.}$$

SS due to error = $S_1^2 = SSE$, say.

Under H_0 , $SSR_{S/k}$ and SSE both follow χ^2 -distribution with respective d.f.s. and they are independent.

$$\text{So, } F = \frac{(SSR_{S/k})/s}{SSE/(n-k-1)} \sim F_{s, n-k-1}.$$

If the observed F , say F_0 , is greater than $F_{\alpha; s, n-k-1}$, we reject H_0 at level α , otherwise we do not reject H_0 .

Associated ANOVA table.

Source of Variation	d.f	SS		F_0
Due to multiple linear regression of Y on x_1, \dots, x_s , after fitting x_1, \dots, x_k	s	$= \sum_{j=1}^k b_j p_j - \sum_{j=s+1}^k b_j^* p_j$	$MSR_{S/k}$	$F_0 = \frac{MSR_{S/k}}{MSE}$
Due to multiple linear regression of Y on x_{s+1}, \dots, x_k	$k-s$	$\sum_{j=s+1}^k b_j^* p_j$	MSR_{k-s}	
Due to multiple linear regression of Y on x_1, \dots, x_k	k	$\sum_{j=1}^k b_j \cdot p_j$	MSR_k	
Error	$n-k-1$	$\sum_{i=1}^n (y_i - \bar{y})^2 - \sum_{j=1}^k b_j \cdot p_j = S^2$	MSE	
Total	$n-1$	$\sum_{i=1}^n (y_i - \bar{y})^2$	—	—

Remark: If anyone wants to test whether the inclusion of variable x_i is at all necessary in the multiple linear regression of Y on x_1, x_2, \dots, x_k , then the above test procedure can be carried out with $s=1$.