

TEST about a SUBSET OF β

Consider the model (for $i=1, \dots, n$)

$$Y_i = \beta_1 + \beta_2 x_{i2} + \dots + \beta_{p_0} x_{i p_0} + \beta_{p_0+1} x_{i p_0+1} + \dots + \beta_p x_{ip} + \epsilon_i \quad \epsilon_i \sim N(0, \sigma^2)$$

we want to test jointly, $(\beta_{p_0+1}, \dots, \beta_p) = 0$

$$\begin{cases} H_0: \beta_{p_0+1} = \dots = \beta_p = 0 \\ H_1: H_0 \text{ at least one of them is } \neq 0 \quad (\exists \epsilon \in \{\beta_{p_0+1}, \dots, \beta_p\} : \beta_r \neq 0) \end{cases}$$

Preliminary considerations:

• Under H_1

We have p covariates

$$Y_i = \beta_1 + \beta_2 x_{i2} + \dots + \beta_{p_0} x_{i p_0} + \beta_{p_0+1} x_{i p_0+1} + \dots + \beta_p x_{ip} + \epsilon_i$$

We call it the "full model".

When we estimate the model, we obtain:

- estimate $\hat{\beta}$ (p -dim. vector)
- residuals $\underline{e} = \underline{y} - X\hat{\beta}$
- sum of squared residuals $\underline{e}^T \underline{e}$
- estimate of σ^2 , $\hat{\sigma}^2 = \frac{1}{n} \underline{e}^T \underline{e}$. Distribution of the estimator $\frac{n\hat{\sigma}^2}{\sigma^2} \sim \chi^2_{n-p}$

• Under H_0

We have a model with $p_0 < p$ covariates

$$Y_i = \beta_1 + \beta_2 x_{i2} + \dots + \beta_{p_0} x_{i p_0} + \epsilon_i$$

We call it the "restricted model".

We are constraining the coefficients $(\beta_{p_0+1}, \dots, \beta_p)$ to be equal to zero.

When we estimate the model, we obtain:

- estimate $\hat{\beta}^0$ (p_0 -dim. vector)
- residuals $\tilde{\underline{e}} = \underline{y} - X\hat{\beta}^0$
- sum of squared residuals $\tilde{\underline{e}}^T \tilde{\underline{e}}$
- estimate of σ^2 , $\tilde{\sigma}^2 = \frac{1}{n} \tilde{\underline{e}}^T \tilde{\underline{e}}$. Distribution of the estimator $\frac{n\tilde{\sigma}^2}{\sigma^2} \sim \chi^2_{n-p_0}$

Remark:

The test about a subset of parameters is a test for comparing two models.

Notice that the two models are NESTED, meaning that the model under H_0 is included into the model under H_1 (it can be obtained from the full model using a set of constraints).

If the models are not nested you can not use this test to compare them.

How we test the hypothesis:

It is useful to write the model in a way to highlight the separation between the unconstrained parameters and the ones we are testing.

First, we formulate the model so that the parameters to test are the last p_0 (simply sort the covariates)

Then, we write

$$\underline{\beta} = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_{p_0} \\ \beta_{p_0+1} \\ \vdots \\ \beta_p \end{bmatrix} = \begin{bmatrix} \underline{\beta}^{(0)} \\ \underline{\beta}^{(1)} \end{bmatrix} \quad \begin{matrix} \underline{\beta}^{(0)} \in \mathbb{R}^{p_0} \\ \underline{\beta}^{(1)} \in \mathbb{R}^{p-p_0} \end{matrix} \quad \rightarrow \text{the system of hypothesis becomes}$$

$$\begin{cases} H_0: \underline{\beta}^{(1)} = 0 \\ H_1: \underline{\beta}^{(1)} \neq 0 \end{cases}$$

Similarly, we write the matrix X as the juxtaposition of two submatrices

$$X = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1p_0} & x_{1,p_0+1} & \dots & x_{1p} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np_0} & x_{n,p_0+1} & \dots & x_{np} \end{bmatrix} = \begin{bmatrix} X^{(0)} & X^{(1)} \end{bmatrix}$$

$n \times p_0 \quad n \times (p-p_0)$

Hence we obtain

FULL MODEL (H_1)

$$\underline{y} \sim N_n(X\underline{\beta}, \sigma^2 I)$$

$$\underline{y} = X\underline{\beta} + \underline{\epsilon} = \begin{bmatrix} X^{(0)} & X^{(1)} \end{bmatrix} \begin{bmatrix} \underline{\beta}^{(0)} \\ \underline{\beta}^{(1)} \end{bmatrix} + \underline{\epsilon}$$

$$= X^{(0)}\underline{\beta}^{(0)} + X^{(1)}\underline{\beta}^{(1)} + \underline{\epsilon}$$

$$\hat{\underline{\beta}} = \begin{bmatrix} \hat{\underline{\beta}}^{(0)} \\ \hat{\underline{\beta}}^{(1)} \end{bmatrix} = (X^T X)^{-1} X^T \underline{y}$$

RESTRICTED MODEL (H_0)

$$\underline{y} \sim N_n(X^{(0)}\underline{\beta}^{(0)}, \sigma^2 I)$$

$$\underline{y} = X^{(0)}\underline{\beta}^{(0)} + \underline{\epsilon}$$

$$\tilde{\underline{\beta}}^{(0)} = (X^{(0)T} X^{(0)})^{-1} X^{(0)T} \underline{y}$$

We know that $\tilde{\underline{e}}^T \tilde{\underline{e}} \geq \underline{e}^T \underline{e}$, since the model under H_0 is a constrained version of the full model. In particular, the difference between the two will be large if the coefficients that I have forced to zero are actually relevant for the analysis.

If H_0 is true, removing $\underline{\beta}^{(1)}$ in the model will not make a big difference for predicting y .

Under H_0 , I expect $\frac{\tilde{\underline{e}}^T \tilde{\underline{e}}}{\underline{e}^T \underline{e}} \approx 1$

$$\Rightarrow \frac{\tilde{\sigma}^2}{\hat{\sigma}^2} \approx 1 \Rightarrow \frac{n\tilde{\sigma}^2}{n\hat{\sigma}^2} \approx 1 \Rightarrow \frac{SSE_{H_0}}{SSE_H} \approx 1 \Rightarrow \frac{SSE_{H_0}}{SSE_H} - 1 \approx 0$$

If H_0 is not true, removing $\underline{\beta}^{(1)}$ will lead to worse results (larger errors).

Under H_1 , I expect $\frac{\tilde{\underline{e}}^T \tilde{\underline{e}}}{\underline{e}^T \underline{e}} \gg 1$

$$\Rightarrow \frac{\tilde{\sigma}^2}{\hat{\sigma}^2} \gg 1 \Rightarrow \frac{SSE_{H_0}}{SSE_H} \gg 1 \Rightarrow \frac{SSE_{H_0}}{SSE_H} - 1 \gg 0$$

To perform the test, we are going to use again a function of $\frac{SSE_{H_0}}{SSE_H} - 1$

TEST STATISTIC and DISTRIBUTION

$$F = \frac{\frac{\tilde{\sigma}^2 - \hat{\sigma}^2}{p-p_0}}{\frac{\hat{\sigma}^2}{n-p}} \stackrel{H_0}{\sim} F_{p-p_0, n-p}$$

analogous formulations

$$F = \frac{\frac{\tilde{\sigma}^2 - \hat{\sigma}^2}{\hat{\sigma}^2} \cdot \frac{n-p}{p-p_0}}{\frac{\tilde{\underline{e}}^T \tilde{\underline{e}} - \underline{e}^T \underline{e}}{\underline{e}^T \underline{e}} \cdot \frac{n-p}{p-p_0}} = \frac{SSE_{H_0} - SSE_H}{SSE_H} \cdot \frac{n-p}{p-p_0} \stackrel{H_0}{\sim} F_{p-p_0, n-p}$$

Note to remember the degrees of freedom

- $\tilde{\sigma}^2 \sim \chi^2_{n-p_0}$
- $\hat{\sigma}^2 \sim \chi^2_{n-p}$

difference of the estimators = $\frac{\tilde{\sigma}^2 - \hat{\sigma}^2}{(n-p_0) - (n-p)} = \frac{\tilde{\sigma}^2 - \hat{\sigma}^2}{n-p_0-n+p} = \frac{\tilde{\sigma}^2 - \hat{\sigma}^2}{p-p_0}$

difference of the d.o.f. = $(n-p_0) - (n-p) = n-p_0-n+p = p-p_0$

its d.o.f. = $n-p$

$F = \frac{\frac{\tilde{\sigma}^2 - \hat{\sigma}^2}{p-p_0}}{\frac{\hat{\sigma}^2}{n-p}} \stackrel{H_0}{\sim} F_{p-p_0, n-p}$

With the data, we compute the observed value of the test, $f^{obs} = \frac{\tilde{\sigma}^2 - \hat{\sigma}^2}{\hat{\sigma}^2} \cdot \frac{n-p}{p-p_0}$

How do we define the reject and acceptance regions?

- acceptance region (values that suggest that the data support H_0)
under H_0 : $\tilde{\sigma}^2 \approx \hat{\sigma}^2 \Rightarrow f^{obs} \approx 0$
 - reject region (values that suggest that the data are against H_0)
under H_1 : $\tilde{\sigma}^2 \gg \hat{\sigma}^2 \Rightarrow f^{obs} \gg 0$ | reject H_0 for large values of f^{obs}
- hence $A = (0, k)$ and $R = (k, +\infty)$

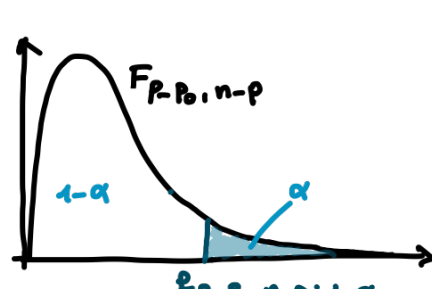
1) FIXED SIGNIFICANCE α

If we fix the significance α , k will be the quantile of level $(1-\alpha)$ of an $F_{p-p_0, n-p}$ distribution

$$R = (f_{p-p_0, n-p; 1-\alpha}, +\infty)$$

with the data: I can compute f^{obs}

- reject H_0 if $f^{obs} > f_{p-p_0, n-p; 1-\alpha}$



2) P-VALUE

Alternatively, the p-value is $\alpha^{obs} = P_{H_0}(F > f^{obs})$

