

ONE-WAY ANOVA (Analysis of variance)

In the cuckoo exercise we had 2 groups of observations and we wanted to test whether the means of the two groups were equal (assuming normality and homoscedasticity). In particular, we showed the equivalence between the two-sample t-test and a test of significance on the regression parameter of a simple em.

Let's generalize the setting and notation

Suppose we are testing the effectiveness of a treatment, and we measure the survival time Y on subjects divided into $G=3$ groups

The question of interest is whether the mean survival times of the three groups are equal or different. If they are different, then the treatments have different effectiveness.

We can use the same notation as in the cuckoo exercise, so that

$$\underline{y} = [y_1, \dots, y_{n_1}, \underbrace{y_{n_1+1}, \dots, y_{n_1+n_2}}_{\text{group 2 with } n_2 \text{ units}}, \underbrace{y_{n_1+n_2+1}, \dots, y_{n_1+n_2+n_3}}_{\text{group 3 with } n_3 \text{ units}}]^T$$

or we can use an alternative (equivalent) notation using 2 indices:

- index of the unit in each group $i=1, \dots, n_g$ \rightarrow number of units in group g
- index of the group $g=1, \dots, G$ \rightarrow number of groups

$$\Rightarrow Y_{ig} \sim N(\mu_g, \sigma^2) \quad \text{independent}$$

survival time of individual i from group g group-specific mean of group g common variance for all groups

With 3 groups, for example, we get

- group 1: n_1 individuals $\rightarrow \underline{Y}_1 = [Y_{11}, \dots, Y_{1n_1}]^T$
- group 2: n_2 individuals $\rightarrow \underline{Y}_2 = [Y_{21}, \dots, Y_{2n_2}]^T$
- group 3: n_3 individuals $\rightarrow \underline{Y}_3 = [Y_{31}, \dots, Y_{3n_3}]^T$

Let us denote with μ_g the mean survival time for group g ($g=1, \dots, G$)

The estimates are

$$\hat{\mu}_g = \bar{y}_g = \frac{1}{n_g} \sum_{i=1}^{n_g} Y_{ig}$$

If we want to test equality of the treatments, we test:

$$\begin{cases} H_0: \mu_1 = \mu_2 = \mu_3 \\ H_a: \text{at least one of them is different} \end{cases}$$

How do we formulate a linear model for this problem?

First, we define the vector of the response by concatenating each group-specific vector \underline{Y}_g

$$\underline{Y} = [\underline{Y}_1^T \underline{Y}_2^T \underline{Y}_3^T]^T = [\underbrace{Y_{11}, Y_{21}, \dots, Y_{31}}_{\text{group 1}}, \underbrace{Y_{12}, \dots, Y_{22}, \dots, Y_{32}}_{\text{group 2}}, \underbrace{Y_{13}, \dots, Y_{23}, \dots, Y_{33}}_{\text{group 3}}]^T$$

vector with $N = n_1 + n_2 + n_3$ elements.

Then, we define the matrix X of the covariates

We use DUMMY VARIABLES where

$$x_{ig} = \begin{cases} 1 & \text{if individual } i \text{ belongs to group } g \\ 0 & \text{otherwise} \end{cases}$$

for $i=1, \dots, n_g$ and $g=1, \dots, G$.

Remark:

consider $G=3$. If we define the matrix X as

$$X = [\underline{x}_1 \underline{x}_2 \underline{x}_3] = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & 0 & 1 \end{bmatrix} \rightarrow \underline{x} = \underline{x}_1 + \underline{x}_2 + \underline{x}_3 \quad \text{multicollinearity!} \quad \text{rank}(X) = 3 < 4 (\text{h.columns})$$

intercept x_0 indicator of group 1 indicator of group 2 indicator of group 3

To encode G groups, if we keep the intercept, we only need $G-1$ dummy variables.

Consider removing \underline{x}_1 . Then X becomes

$$X = [\underline{x}_2 \underline{x}_3] = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ \vdots & \vdots & \vdots \\ 1 & 0 & 0 \\ 1 & 1 & 0 \\ \vdots & \vdots & \vdots \\ 1 & 1 & 0 \\ 1 & 0 & 1 \\ \vdots & \vdots & \vdots \\ 1 & 0 & 1 \end{bmatrix} \quad \text{with } (N \times G) \text{ matrix}$$

$$\text{and } \underline{\epsilon} \sim N_N(0, \sigma^2 I_n)$$

or, equivalently

$$Y_i = \beta_1 + \beta_2 x_{i2} + \beta_3 x_{i3} + \epsilon_i$$

with

$$x_{i2} = \begin{cases} 1 & \text{for } i=n_1+1, \dots, n_1+n_2 \\ 0 & \text{otherwise} \end{cases}$$

$$x_{i3} = \begin{cases} 1 & \text{for } i=n_1+n_2+1, \dots, n_1+n_2+n_3 \\ 0 & \text{otherwise} \end{cases}$$

Let's study the expected value of observations in each group according to the model

$$E[Y_{i1}] = \beta_1 \quad \text{for } i=1, \dots, n_1$$

$$E[Y_{i2}] = \beta_1 + \beta_2 \quad \text{for } i=n_1+1, \dots, n_1+n_2$$

$$E[Y_{i3}] = \beta_1 + \beta_3 \quad \text{for } i=n_1+n_2+1, \dots, n_1+n_2+n_3$$

INTERPRETATION:

• INTERCEPT: β_1 is the mean of Y_{i1} when $g=1$ (when all dummy variables are equal to zero) (mean of the group for which we removed the dummy variable)

This group is said to be the REFERENCE GROUP: it is the BASELINE

A classical example is the control group (i.e. the "no treatment") in clinical trials.

$$\Rightarrow \beta_1 = E[Y_{i1}]$$

The other groups are described in terms of DEVIATION FROM THE BASELINE.

• β_2 is the difference in the mean of Y_{i2} w.r.t. the mean of Y_{i1} from the model we have

$$E[Y_{i2}] = \beta_1 + \beta_2$$

$$\Rightarrow \beta_2 = E[Y_{i2}] - E[Y_{i1}]$$

$$= \mu_2 - \mu_1$$

• β_3 is the difference in the mean of Y_{i3} w.r.t. the mean of Y_{i1}

$$E[Y_{i3}] = \beta_1 + \beta_3$$

$$\Rightarrow \beta_3 = E[Y_{i3}] - E[Y_{i1}]$$

$$= \mu_3 - \mu_1$$

Remark: we automatically get the estimates of the regression parameters:

Reparametrization

$$\begin{cases} \mu_1 = \beta_1 \\ \mu_2 = \beta_1 + \beta_2 \\ \mu_3 = \beta_1 + \beta_3 \end{cases} \Leftrightarrow \begin{cases} \beta_1 = \mu_1 \\ \beta_2 = \mu_2 - \mu_1 \\ \beta_3 = \mu_3 - \mu_1 \end{cases}$$

Invariance of the MLR w.r.t. reparametrizations

$$\begin{cases} \hat{\beta}_1 = \hat{\beta}_1 \\ \hat{\beta}_2 = \hat{\beta}_2 - \hat{\beta}_1 \\ \hat{\beta}_3 = \hat{\beta}_3 - \hat{\beta}_1 \end{cases} \Rightarrow \begin{cases} \hat{\beta}_1 = \bar{y}_1 \\ \hat{\beta}_2 = \bar{y}_2 - \bar{y}_1 \\ \hat{\beta}_3 = \bar{y}_3 - \bar{y}_1 \end{cases}$$

We can easily compute the predicted values \hat{y}_{ig}

$$\hat{y}_{i1} = \hat{\beta}_1 = \bar{y}_1 \quad i=1, \dots, n_1$$

$$\hat{y}_{i2} = \hat{\beta}_1 + \hat{\beta}_2 = \bar{y}_1 + \bar{y}_2 - \bar{y}_1 = \bar{y}_2 \quad i=n_1+1, \dots, n_1+n_2$$

$$\hat{y}_{i3} = \hat{\beta}_1 + \hat{\beta}_3 = \bar{y}_1 + \bar{y}_3 - \bar{y}_1 = \bar{y}_3 \quad i=n_1+n_2+1, \dots, n_1+n_2+n_3$$

\Rightarrow The predicted values are the group-specific means.

Finally, the test about equality of the group-specific means becomes

$$\begin{cases} H_0: \beta_2 = \beta_3 = 0 \\ H_a: \text{at least one is } \neq 0 \end{cases}$$

\hookrightarrow test about the overall significance of the model