

Correspondence between the t-test for comparing the means of two independent Gaussian samples with equal variances and test on the regression coefficient of a Gaussian LM.

We can reformulate the test using a Gaussian linear model.

Write the full vector of the response as :

$$\underline{y} = \begin{bmatrix} \underline{y}^R \\ \underline{y}^W \end{bmatrix} = (\underbrace{y_1, \dots, y_n}_{\text{Robin}}, \underbrace{y_{n+1}, \dots, y_{n+m}}_{\text{Wren}})^T \quad (n+m)\text{-dimensional vector}$$

MODEL FORMULATION

$$Y_i = \mu_i + \varepsilon_i$$

$$= \beta_1 + \beta_2 x_i + \varepsilon_i \quad \varepsilon_i \sim N(0, \sigma^2) \text{ iid} \quad i = 1, \dots, n+m$$

$x_i$  is a DUMMY variable (indicator variable)

$$x_i = \begin{cases} 0 & \text{if the } i\text{-th egg is in a ROBIN's nest} \\ 1 & \text{if the } i\text{-th egg is in a WREN's nest} \end{cases} \rightarrow X = \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ \vdots & \vdots \\ 1 & 0 \\ 1 & 1 \\ 1 & 1 \\ \vdots & \vdots \\ 1 & 1 \end{bmatrix} = \begin{bmatrix} \text{Robin} \\ \text{Wren} \end{bmatrix}$$

Let's see what happens to  $Y_i$  depending on the bird species:

• if egg<sub>i</sub> is in a ROBIN's nest

$$x_i = 0 \Rightarrow \mu_i = \beta_1 + \beta_2 \cdot 0 = \beta_1 \Rightarrow Y_i \sim N(\beta_1, \sigma^2) \text{ for } i = 1, \dots, n$$

$$\text{This is the group of eggs from robins} \Rightarrow Y_i \sim N(\mu^R, \sigma^2) \Rightarrow \beta_1 = \mu^R$$

• if egg<sub>i</sub> is in a WREN's nest

$$x_i = 1 \Rightarrow \mu_i = \beta_1 + \beta_2 \cdot 1 = \beta_1 + \beta_2 \Rightarrow Y_i \sim N(\beta_1 + \beta_2, \sigma^2) \text{ for } i = n+1, \dots, n+m$$

$$\text{This is the group of eggs from wrens} \Rightarrow Y_i \sim N(\mu^W, \sigma^2) \Rightarrow \beta_1 + \beta_2 = \mu^W$$

Remark: this is a reparametrization :

a one-to-one correspondence between  $(\mu^R, \mu^W)$  and  $(\beta_1, \beta_2)$

$$\begin{cases} \mu^R = \beta_1 \\ \mu^W = \beta_1 + \beta_2 \end{cases} \Leftrightarrow \begin{cases} \beta_1 = \mu^R \\ \beta_2 = \mu^W - \mu^R \end{cases}$$

$$\text{The correspondence also holds for the ML estimates: } \begin{cases} \hat{\beta}_1 = \hat{\mu}^R \\ \hat{\beta}_2 = \hat{\mu}^W - \hat{\mu}^R \end{cases}$$

$$\text{So if we want to test } H_0: \mu^R = \mu^W \Leftrightarrow H_0: \mu^W - \mu^R = 0 \Leftrightarrow H_0: \beta_2 = 0$$

To test this hypothesis using the linear model

$H_0: \beta_2 = 0$  we have seen the test on individual coefficients  $\rightarrow$  t-test in particular

$$T = \frac{\hat{\beta}_2 - 0}{\sqrt{\frac{\hat{\sigma}^2}{\sum_{i=1}^{n+m} (x_i - \bar{x})^2}}} \stackrel{H_0}{\sim} t_{n+m-2}$$

$\rightarrow$  expression of  $\text{Var}(\hat{\beta}_2)$  in the simple LM.

We now compute the estimated regression model and show the equivalence with the previous procedure.

$$\text{We have a SIMPLE LINEAR MODEL } Y_i = \beta_1 + \beta_2 x_i + \varepsilon_i \quad \varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$$

From the previous lectures we know that the estimate of  $\beta_2$  in the simple lm is:

$$\hat{\beta}_2 = \frac{\sum_{i=1}^{n+m} (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^{n+m} (x_i - \bar{x})^2} = \frac{\sum_{i=1}^{n+m} x_i y_i - (n+m) \bar{x} \bar{y}}{\sum_{i=1}^{n+m} (x_i - \bar{x})^2}$$

we need to compute  $\bar{x}, \bar{y}, \sum_{i=1}^{n+m} x_i y_i, \sum_{i=1}^{n+m} (x_i - \bar{x})^2$

$$\bullet \bar{x} = \frac{1}{n+m} \sum_{i=1}^{n+m} x_i = \frac{m}{n+m}$$

$$\bullet \bar{y} = \frac{1}{n+m} \sum_{i=1}^{n+m} y_i = \frac{1}{n+m} \left( \sum_{i=1}^n y_i + \sum_{i=n+1}^{n+m} y_i \right) = \frac{1}{n+m} (n \bar{y}^R + m \bar{y}^W)$$

$$\bullet \sum_{i=1}^{n+m} x_i y_i = \sum_{i=n+1}^{n+m} y_i = m \bar{y}^W$$

$$\begin{aligned} \bullet \sum_{i=1}^{n+m} (x_i - \bar{x})^2 &= \sum_{i=1}^n (0 - \bar{x})^2 + \sum_{i=n+1}^{n+m} (1 - \bar{x})^2 = \sum_{i=1}^n \bar{x}^2 + \sum_{i=n+1}^{n+m} (1 - \bar{x})^2 = \\ &= n \cdot \left( \frac{m}{n+m} \right)^2 + \sum_{i=n+1}^{n+m} \left( 1 - \frac{m}{n+m} \right)^2 = \\ &= \frac{n m^2}{(n+m)^2} + m \cdot \frac{n^2}{(n+m)^2} = \frac{n m (n+m)}{(n+m)^2} = \frac{n m}{n+m} \end{aligned}$$

Hence

$$\begin{aligned} \hat{\beta}_2 &= \frac{m \bar{y}^W - (n+m) \cdot \frac{m}{n+m} \cdot \frac{1}{n+m} (n \bar{y}^R + m \bar{y}^W)}{\frac{n m}{n+m}} \\ &= \frac{\bar{y}^W - \frac{1}{n+m} (n \bar{y}^R + m \bar{y}^W)}{\frac{n}{n+m}} = \\ &= \frac{\frac{1}{n+m} (n \bar{y}^W + m \bar{y}^W - n \bar{y}^R - m \bar{y}^W)}{\frac{n}{n+m}} = \bar{y}^W - \bar{y}^R \end{aligned}$$

The estimate of  $\beta_1$  instead is  $\hat{\beta}_1 = \bar{y} - \hat{\beta}_2 \bar{x}$

in this case:

$$\begin{aligned} \hat{\beta}_1 &= \frac{1}{n+m} (n \bar{y}^R + m \bar{y}^W) - \frac{m}{n+m} (\bar{y}^W - \bar{y}^R) \\ &= \frac{1}{n+m} (n \bar{y}^R + m \bar{y}^W - m \bar{y}^W + m \bar{y}^R) \\ &= \frac{n+m}{n+m} \bar{y}^R = \bar{y}^R \end{aligned}$$

Finally

$$\begin{aligned} s^2 &= \frac{1}{n+m-2} \sum_{i=1}^{n+m} (y_i - \hat{y}_i)^2 = \frac{1}{n+m-2} \sum_{i=1}^{n+m} (y_i - \hat{\beta}_1 - \hat{\beta}_2 x_i)^2 = \\ &= \frac{1}{n+m-2} \sum_{i=1}^{n+m} (y_i - \bar{y}^R - (\bar{y}^W - \bar{y}^R) x_i)^2 = \\ &= \frac{1}{n+m-2} \left[ \sum_{i=1}^n (y_i - \bar{y}^R)^2 + \sum_{i=n+1}^{n+m} (y_i - \bar{y}^W)^2 \right] = \\ &= \frac{1}{n+m-2} \left[ \underbrace{\sum_{i=1}^n (y_i - \bar{y}^R)^2}_{(n-1) s_R^2} + \underbrace{\sum_{i=n+1}^{n+m} (y_i - \bar{y}^W)^2}_{(m-1) s_W^2} \right] = \frac{1}{n+m-2} [(n-1) s_R^2 + (m-1) s_W^2] \end{aligned}$$

Hence we obtain

$$\begin{aligned} \hat{\beta}_2 &= \bar{y}^W - \bar{y}^R \\ \hat{\beta}_1 &= \bar{y}^R \\ s^2 &= \frac{(n-1) s_R^2 + (m-1) s_W^2}{n+m-2} \\ \sum_{i=1}^{n+m} (x_i - \bar{x})^2 &= \frac{n m}{n+m} \end{aligned}$$

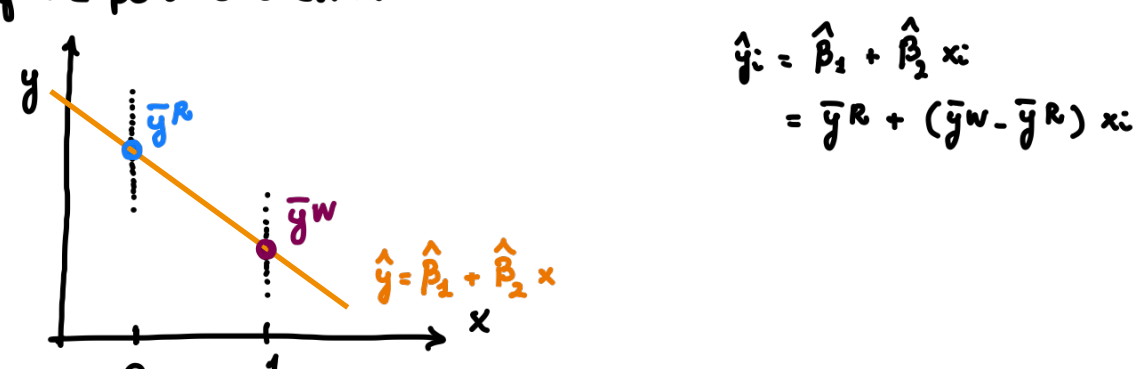
Going back to the test,

$$\begin{aligned} T &= \frac{\hat{\beta}_2}{\sqrt{\frac{\hat{\sigma}^2}{\sum_{i=1}^{n+m} (x_i - \bar{x})^2}}} = \frac{\bar{y}^W - \bar{y}^R}{\sqrt{\frac{(n-1) s_R^2 + (m-1) s_W^2}{n+m-2} \cdot \left( \frac{n m}{n+m} \right)^{-1}}} \\ &\Rightarrow T = \frac{\bar{y}^W - \bar{y}^R}{\sqrt{\frac{(n-1) s_R^2 + (m-1) s_W^2}{n+m-2} \cdot \frac{n+m}{n \cdot m}}} \stackrel{H_0}{\sim} t_{n+m-2} \end{aligned}$$

$\rightarrow$  it's the same expression as the two-sample t-test

Hence we have proven the correspondence of the two procedures.

if we plot the estimated model



Remark:

Notice that if we consider instead a covariate

$$z_i = \begin{cases} 1 & \text{if the bird is a robin} \\ 0 & \text{if the bird is a wren} \end{cases}$$

then  $\mu^W = \beta_1$  and  $\mu^R = \beta_1 + \beta_2$

is a different model but the result of inference is the same