

Exercises: Multiple Linear Regression Part II

Valentina Zangirolami - valentina.zangirolami@unimib.it

December 14, 2023

(Referring to the theoretical parts: 9, 10, 12, 13, 14, 15, 16, 17, 18, 19)

1 Exercise 1

Among $n = 31$ cherry trees, data about *Volume*, *Diameter* and *Height* of each tree were collected.

Exercise 1.1

Specify an appropriate regression model for the response variable *Volume* and the related assumptions. Propose a transformation for all variables by thinking about the geometric relationship of our variables.

An appropriate regression model is the multiple regression model which can be defined as follows

$$\text{Volume}_i = \beta_1 + \beta_2 \text{Diameter}_i + \beta_3 \text{Height}_i + \varepsilon_i$$

with the following assumptions:

(i) $\varepsilon_i \stackrel{i.i.d.}{\sim} N(0, \sigma^2) \quad i=1, \dots, n$

(ii) linear independence among independent variables:

the vectors $x_j \in \mathbb{R}^n$, $j=1, \dots, p$, are linearly independent.

(iii) linearity

(i) and (iii) imply $y_i \stackrel{i.i.d.}{\sim} N(\mu_i, \sigma^2)$, where $\mu_i = \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}$

$i=1, \dots, n$

Answers

The geometry suggests the following model

$$\log(\text{Volume}) = \log(\pi/4) + 2 \log(\text{Diameter}) + \log(\text{Height}) + \varepsilon$$

Therefore, we can consider the following model

$$(a) \quad Y_i = \beta_1 + \beta_2 X_{i2} + \beta_3 X_{i3} + \varepsilon_i \quad , \quad \varepsilon_i \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$$

where

$$\cdot Y_i = \log(\text{Volume}_i)$$

$$\cdot X_{i2} = \log(\text{Diameter}_i)$$

$$\cdot X_{i3} = \log(\text{Height}_i)$$

Exercise 1.2

Knowing that

$$(X^T X)^{-1} = \begin{bmatrix} 96.572 & 3.139 & -24.165 \\ 3.139 & 0.849 & -1.227 \\ -24.165 & -1.227 & 6.310 \end{bmatrix} \quad X^T y = \begin{bmatrix} 101.455 \\ 263.056 \\ 439.896 \end{bmatrix}$$

Find the estimate for β and write the estimated model.

(Previous data refers to logarithmic case)

The model (a) can be re-written using matrix notation as

$$\underline{y} = X \underline{\beta} + \underline{\varepsilon} \quad \text{where } \underline{\varepsilon} \sim N(0, \sigma^2 I)$$

where

$$X = \begin{bmatrix} 1 & x_{12} & x_{13} \\ 1 & x_{22} & x_{23} \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \\ 1 & x_{31,2} & x_{31,3} \end{bmatrix} \quad \underline{\beta} = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix} \quad \underline{\varepsilon} = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_{31} \end{bmatrix}$$

The estimate of β corresponds to

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

and thus

$$\tilde{\underline{\beta}} = \begin{bmatrix} 96.572 & 3.139 & -24.165 \\ 3.139 & 0.849 & -1.227 \\ -24.165 & -1.227 & 6.310 \end{bmatrix}^{-1} \begin{bmatrix} 101.455 \\ 263.056 \\ 439.896 \end{bmatrix} = \begin{bmatrix} -6.6418 \\ 2.0494 \\ 1.314 \end{bmatrix}$$

The estimated model

$$\underline{y} = X \begin{bmatrix} -6.6418 \\ 2.0494 \\ 1.314 \end{bmatrix} + \underline{\varepsilon}$$

Using the notation of the equation (a), we have

$$\hat{\beta}_1 = -6.6438 \quad \hat{\beta}_2 = 2.0494 \quad \hat{\beta}_3 = 1.314$$

Exercise 1.3

Knowing the maximum likelihood estimate $\hat{\sigma}^2 = 0.00598$, find the unbiased estimate s^2 . Then, calculate the estimated variance-covariance matrix of $\hat{\beta}$.

To find the unbiased estimate for δ^2 , we need to compute

$$s^2 = \frac{1}{n-p} \underline{e}^\top \underline{e} = \frac{m}{m-p} \frac{1}{m} \underline{e}^\top \underline{e} = \frac{m}{m-p} \hat{\delta}^2$$

where p is the number of coefficients (the intercept is included) and

$$\underline{e} = \begin{bmatrix} y_1 - \hat{y}_1 \\ y_2 - \hat{y}_2 \\ \vdots \\ y_{31} - \hat{y}_{31} \end{bmatrix} = \underline{y} - \hat{\underline{y}} = (I - P)\underline{y}$$

$P = X(X^\top X)^{-1}X^\top \underline{y}$ is the projection matrix and I

is the identity matrix

$$I = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & & \vdots \\ \vdots & & \ddots & \\ 0 & \cdots & \cdots & 1 \end{bmatrix}$$

31×31

Knowing $\hat{\delta}^2 = 0.00598$, we can find s^2 easily

$$s^2 = \frac{31}{28} 0.00598 = 0.0066207$$

We also know the distribution of the estimator $\hat{\beta}$:

$$\hat{\beta} \sim N_p (\beta, \delta^2 (X^\top X)^{-1})$$

where each estimator $\hat{\beta}_i : \hat{\beta}_i \sim N(\beta_i, \delta^2 (X^\top X)^{-1})$

and the variance/covariance matrix is $\text{Var}(\hat{\beta}) = \delta^2 (X^\top X)^{-1}$

Its estimate corresponds to

$$\widehat{\text{Var}}(\hat{\beta}) = S^2 (X^T X)^{-1} = 0.0066207 \begin{bmatrix} 96.572 & 3.139 & -24.165 \\ 3.139 & 0.849 & -1.227 \\ -24.165 & -1.227 & 6.310 \end{bmatrix}^{-1}$$
$$= \begin{bmatrix} 0.639 & 0.0208 & -0.16 \\ 0.0208 & 0.0056 & -0.0081 \\ -0.16 & -0.0081 & 0.042 \end{bmatrix}$$

Exercise 1.4

Perform a statistical test to evaluate if the regression coefficient related to the *Diameter* is significant. Specify the system of hypothesis, the test statistic and the p-value. Then, perform another statistical test with the following null hypothesis $H_0 : \beta_2 = \beta_3 = 0$.
 (Restricted Model $\text{SS}_E = 8.309$)

Inference on $\hat{\beta}_2$

System of hypothesis

$$\left\{ \begin{array}{l} H_0: \beta_2 = 0 \\ H_1: \beta_2 \neq 0 \end{array} \right.$$

Test statistic:

$$t_2^{\text{obs}} = \frac{\hat{\beta}_2 - 0}{\sqrt{\text{Var}(\hat{\beta}_2)}} = \frac{2.0494}{\sqrt{0.0056}} = 24.38626$$

$$\alpha^{\text{obs}} = 2 \min \{ P(t_{28} \leq 24.38626), P(t_{28} \geq 24.38626) \} \approx 0$$

We reject H_0 .

$$\left\{ \begin{array}{l} H_0: \beta_2 = \beta_3 = 0 \\ H_1: \exists j \in \{2, 3\} \text{ s.t. } \beta_j \neq 0 \end{array} \right.$$

$$f^{\text{obs}} = \frac{\frac{8.309 - 0.1855}{3-1}}{\frac{0.1855}{31-3}} = 613.094$$

$$\alpha^{\text{obs}} = P(f_{2, 28} > 613.094) \approx 0 \quad \text{We reject } H_0.$$

2 Exercise 2

Among $n=32$ births, the following three variables were collected:

- *Weight*: birth weight in grams of baby
- *Smoking*: Smoking status of mother (yes or no)
- *Gest*: length of gestation in weeks

Describe the equation of the multiple regression model specifying the nature of each variable (without including the interaction term). Then, express the model in matrix form: $\underline{Y} = \underline{X}\underline{\beta} + \underline{\varepsilon}$, explicitly stating \underline{Y} , \underline{X} , $\underline{\beta}$, and $\underline{\varepsilon}$ (and its distribution).

Data contain three variables:

- Weight is a quantitative variable
- Smoking is a qualitative variable
- Gest is a quantitative variable

The multiple regression model can be represented by

$$\text{Weight}_i = \beta_1 + \beta_2 \text{Gest}_i + \beta_3 D_i + \varepsilon_i \quad \varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2) \quad i=1, \dots, 32$$

where $D_i = \begin{cases} 1, & \text{if Smoking}_i = Y \\ 0, & \text{otherwise} \end{cases}$

and using matrix notation:

$$\underbrace{\text{Weight}}_{32 \times 1} = \underbrace{X}_{32 \times 3} \underbrace{\underline{\beta}}_{3 \times 1} + \underbrace{\underline{\varepsilon}}_{32 \times 1} \quad \underline{\varepsilon} \sim N_3(0, \sigma^2 I)$$

$$\Rightarrow \text{Weight} \sim N_3(X\underline{\beta}, \sigma^2 I) \quad \text{where } I \text{ is an identity matrix}$$

$$X = \begin{bmatrix} 1 & \text{Gest}_1 & \text{Smoking}_1 \\ 1 & \text{Gest}_2 & \text{Smoking}_2 \\ 1 & \text{Gest}_3 & \text{Smoking}_3 \\ \vdots & \vdots & \vdots \\ 1 & \text{Gest}_{32} & \text{Smoking}_{32} \end{bmatrix} \quad \underline{\beta} = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix} \quad \underline{\varepsilon} = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_{32} \end{bmatrix}$$

Exercise 2.1

Knowing the following values

$$\widehat{Weight}^{(n)} = -2390 + 143 \text{ Gest} \quad (\text{Smoking} = 0)$$

$$\widehat{Weight}^{(y)} = -2635 + 143 \text{ Gest} \quad (\text{Smoking} = 1)$$

What is the estimated value of the regression coefficient associated to the dummy variable? Interpret the results and explain, theoretically, why the two equations are different.

$$\hat{\beta}_3 = -2635 + 2390 = -245$$

Repeating to the previous notation (ex.2), the two equations can be represented by

$$\widehat{Weight}^{(n)} = \hat{\beta}_1 + \hat{\beta}_2 \text{ Gest} \quad (D_i = 0)$$

$$\widehat{Weight}^{(y)} = (\hat{\beta}_1 + \hat{\beta}_3) + \hat{\beta}_2 \text{ Gest} \quad (D_i = 1)$$

In the second equation, the intercept increases by $\hat{\beta}_3$ due to $D_i = 1$ (Smoking = Yes).

Exercise 2.2

After adding an interaction term, we obtained the following estimated regression model

$$\widehat{Weight} = -2546.138 + 147.207 \text{ Gest} + 71.574 \text{ Smoke} - 8.178 \text{ Gest} \times \text{Smoke}$$

Provide the estimated regression model for each group (Smoke: yes or no) and interpret the results.

$$(\text{Smoking} = \text{No}) \quad \widehat{Weight}^{(n)} = -2546.138 + 147.207 \text{ Gest}$$

$$(\text{Smoking} = \text{Yes}) \quad \widehat{Weight}^{(y)} = -2474.564 + 139.029 \text{ Gest}$$

Considering the estimated regression model for $\text{Smoking} = \text{No}$,
on average the birth weight of baby increases by
147.207 grams as Gest increases by one week.

Instead, considering smokers:

On average the birth weight of baby increases by 139.029 grams
as Gest increases by one week.

The interpretation of intercept does not provide relevant information,
given that ~~does~~ makes sense focusing with the case of Gest=0

Exercise 2.3

Consider the following table where for the two models M_0 (the restricted model) and M_4 (corresponds to the model specified in the ex. 2.2) are expressed the residual sum of squares (SSE). Complete the table.

Model	D.o.f.	SSE
Resticted (M_0)	30	839951.03
Unconstrained (M_4)	28	384391.46

Further, perform a statistical test with the following system of hypothesis (where Smoking has two classes: Y and N)

$$\begin{cases} H_0: \mu_N = \mu_Y \\ H_1: \mu_N \neq \mu_Y \end{cases}$$

Compute the test statistic and the p-value.

Degree of freedom (from table)

$$H_0: m - p_0 = 32 - 2 = 30$$

$$H_4: m - p = 32 - 4 = 28$$

where p_0 are the number of regression coefficients for H_0 and p are the number of regression coefficients for H_4 .

Given the following system of hypothesis

$$\begin{cases} H_0: \mu_N = \mu_Y \\ H_1: \mu_N \neq \mu_Y \end{cases} \quad (4-2)$$

TEST STATISTIC

$$F^{\text{obs}} = \frac{839951.03 - 384391.46 / 2}{384391.46 / 32 - 4} = 16.59$$

P-VALUE

$$\alpha^{\text{obs}} = P(F_{2,28} > 16.59) = 1 - P(F_{2,28} < 16.59) \approx 0$$

We reject the null hypothesis. The coefficients β_3 and β_4 are significant.

The previous system of hypothesis can be rewritten as

$$\left\{ \begin{array}{l} H_0: \beta_3 = \beta_4 = 0 \\ H_1: \beta_3 \neq \beta_4 \neq 0 \end{array} \right.$$

Exercise 2.4

Now, we want just to test the regression coefficient related to the interaction. The SSE of the model without the interaction is equal to 387069.83. Perform a valid test and discuss about the results.

System of hypothesis

$$\begin{cases} H_0: \beta_4 = 0 \\ H_1: \beta_4 \neq 0 \end{cases}$$

TEST STATISTIC:

$$f^{\text{obs}} = \frac{387069.83 - 384391.6 / 1}{384391.6 / 28} = 0.20$$

P-VALUE

$$\alpha^{\text{obs}} = P(F_{2,28} \geq 0.20) = 0.6621$$

$$(F_{1,28; 0.3379} \approx 0.20)$$

We cannot reject H_0 .

Exercise 2.5 Knowing the $SSE = 3735789.2$ of the model $y = \beta_1 + \varepsilon$.

Compute R^2 and adjusted- R^2 related to M_4 . Explain the difference.

Let consider

$$M_0: Y = \beta_1 \mathbb{1} + \varepsilon \quad \text{where} \quad \tilde{\sigma}^2 = \frac{\underline{e}_0^\top \underline{e}_0}{m} = \frac{1}{m} \sum_{i=1}^m (y_i - \bar{y})^2$$

and

$$M_4: \text{Delete } \beta_2 \quad y_i = \beta_1 + \beta_2 \text{Gest}_i + \beta_3 \text{Smoke}_i + \beta_4 \text{Gest}_i \times \text{Smoke}_i + \varepsilon_i$$

with $\hat{\sigma}^2 = \frac{\underline{e}_4^\top \underline{e}_4}{m}$ where \underline{e}_4 refers to residuals of M_4

Then, we know

$$\frac{\tilde{\sigma}^2 - \hat{\sigma}^2}{\tilde{\sigma}^2} = \frac{\tilde{\sigma}^2}{\tilde{\sigma}^2} - 1 = \frac{\sum_{i=1}^m (y_i - \bar{y})^2}{\sum_{i=1}^m (\underline{e}_i)^2} - 1 = \frac{1}{1-R^2} - 1 = \frac{R^2}{1-R^2}$$

Exploiting this equivalence, we can find our R^2 \textcircled{K}

$$\frac{3735789.2}{384391.46} - 1 = 8.7187$$

$$\frac{R^2}{1-R^2} = 8.7187 \Rightarrow R^2 + 8.7187 R^2 = 8.7187 \Rightarrow R^2 = \frac{8.7187}{9.7187} = 0.8971$$

While the adjusted- R^2 (R^2_{adj})

$$R^2_{adj} = 1 - (1 - R^2) \frac{m-1}{m-p} = 1 - (1 - 0.8971) \frac{31}{28} = 0.886075$$

In the case of multiple regression model, the R^2_{adj} is a penalized R^2 with respect to the number of covariates ($p-1$) into the model.

R^2_{adj} penalizes model with many variables.

Therefore, however, R^2 and R^2_{adj} reach an higher value, near to 1, which means that our model is good.

\textcircled{K} You can also consider directly

$$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST} = 1 - \frac{384391.46}{3735789.2} = 0.8971$$

3 Exercise 3

Among $n = 15$ poisoned lab rats, data about survival time (Y) and antidote/treatment (T) were collected. The researchers consider three kind of treatment: A, B and C, where

Group	n_j	\bar{y}_j	s_j
A ($j = 1$)	4	0.347	0.2055
B ($j = 2$)	3	3.067	0.4618
C ($j = 3$)	8	1.764	0.5007

Compute the total deviance (total sum of squares - SST) and explain its components.

We are dealing with ANOVA one-way ANOVA.

In a sense, we know

$$\underbrace{\sum_{j=1}^3 \sum_{i=1}^{n_j} (y_{ij} - \bar{y})^2}_{SST} = \underbrace{\sum_{j=1}^3 (n_j - 1) s_j^2}_{SSE} + \underbrace{\sum_{j=1}^3 n_j (\bar{y}_j - \bar{y})^2}_{SSR}$$

where $j = \{1, 2, 3\}$ refers to each group and n_j is the number of observations for each group (see the table).

In this case, SSE corresponds to the deviance between groups while SSR corresponds to the deviance within groups.

We can obtain the average of survival time by

$$\bar{y} = \frac{1}{m} \sum_{j=1}^3 n_j \bar{y}_j = \frac{(4 \cdot 0.347 + 3 \cdot 3.067 + 8 \cdot 1.764)}{15} = 1.646733$$

Then, let's compute SSE and SSR

$$SSE = \underline{3 \cdot (0.2055)^2 + 2 \cdot (0.4618)^2 + 7 \cdot (0.5007)^2} = 2.308113$$

$$SSR = 4(0.347 - 1.646733)^2 + 3(3.067 - 1.646733)^2 + 8(1.764 - 1.646733)^2 = 12.91871$$

We are ready to find SST

$$SST = SSE + SSR = 2.308113 + 12.91871 = 15.22682$$

Exercise 3.1

Perform a statistical test to evaluate if the means (of each group) are homogeneous.

System of hypothesis

$$\begin{cases} H_0: \mu_1 = \mu_2 = \mu_3 \\ H_1: \exists j \text{ s.t. } \mu_j \neq \mu_i \quad (i \neq j, i, j \in \{1, 2, 3\}) \end{cases}$$

test statistic

$$f^{\text{obs}} = \frac{SSR / (J-1)}{SSE / (n-J)} \stackrel{H_0}{\sim} F_{J-1, n-J}$$

$$f^{\text{obs}} = \frac{12.92/2}{2.308/12} = 33.59$$

p-value

$$\alpha^{\text{obs}} = P(F_{2, 12} > 33.59) \approx 0 \quad \text{We reject null hypothesis}$$

Exercise 3.2

Specify an appropriate linear regression model. Given the previous data, find the estimates for each coefficient and interpret them.

An appropriate model can be

$$Y_i = \beta_1 + \beta_2 X_{i2} + \beta_3 X_{i3} + \varepsilon_i \quad \text{assuming } \varepsilon_i \sim N(0, \sigma^2), i=1, \dots, 15$$

and which implies $Y_i \sim N(\mu_i, \sigma^2)$, where

$$\mu_i = \beta_1 + \beta_2 X_{i2} + \beta_3 X_{i3}$$

Therefore, it also implies:

$$(a) \quad \begin{cases} Y_i \sim N(\beta_1, \sigma^2) & i=1, \dots, 4 \quad (\text{Group A}) \\ Y_i \sim N(\beta_1 + \beta_2, \sigma^2) & i=5, \dots, 7 \quad (\text{Group B}) \\ Y_i \sim N(\beta_1 + \beta_3, \sigma^2) & i=8, \dots, 15 \quad (\text{Group C}) \end{cases}$$

considering our covariates as follows

$$X_{i2} = \begin{cases} 1 & \text{if the antidote is B} \\ 0 & \text{otherwise} \end{cases}$$

$$X_{i3} = \begin{cases} 1 & \text{if the antidote is C} \\ 0 & \text{otherwise} \end{cases}$$

Basically, (a) guarantees

$$\mu_A = \beta_1; \mu_B = \beta_1 + \beta_2; \mu_C = \beta_1 + \beta_3$$

and if we want to find the maximum likelihood estimates, the property of unbiasedness we know that they satisfy the following conclusions:

$$\begin{aligned} \hat{\beta}_1 &= \bar{\mu}_A = \bar{Y}_A \\ \hat{\beta}_1 + \hat{\beta}_2 &= \bar{\mu}_B = \bar{Y}_B \quad (\text{for equivalence property}) \\ \hat{\beta}_1 + \hat{\beta}_3 &= \bar{\mu}_C = \bar{Y}_C \end{aligned}$$

Hence, we can find all estimates as follows

$$\hat{\beta}_1 = \bar{y}_A$$

$$\hat{\beta}_2 = \bar{y}_B - \bar{y}_A$$

$$\hat{\beta}_3 = \bar{y}_C - \bar{y}_A$$

Now, we are ready to compute our values

$$\hat{\beta}_1 = 0.347$$

$$\hat{\beta}_2 = 3.067 - 0.347 = 2.72$$

$$\hat{\beta}_3 = 1.464 - 0.347 = 1.417$$

The estimated regression model can be written as

$$\hat{y}_i = 0.347 + 2.72 x_{i2} + 1.417 x_{i3}$$

Interpretation of regression coefficients:

- $\hat{\beta}_1 = 0.347$ is the average of the survival time (y) when the group is A
- $\hat{\beta}_2 = 2.72$ is the average increase in y when the group is B
- $\hat{\beta}_3 = 1.417$ is the average increase in y when the group is C

Exercise 3.3

Propose an alternative test for ex. 3.1 specifying the system of hypothesis. Then, compute R^2 of our model.

To test if the means of each group are homogeneous, we can use equivalently the following system of hypothesis

$$\begin{cases} H_0: \beta_2 = \beta_3 = 0 \\ H_1: \beta_2 \neq 0 \vee \beta_3 \neq 0 \end{cases}$$

Therefore, the restricted model (or null model) is

$$H_0: y_i = \beta_1 + \varepsilon_i$$

and the unconstrained model (or full model) is

$$H_1: y_i = \beta_1 + \beta_2 x_{i2} + \beta_3 x_{i3} + \varepsilon_i$$

And we can point out: $\hat{\sigma}^2$ as the estimated variance of the null model while $\hat{\sigma}^2$ as the estimated variance of the full model, where

$$\hat{\sigma}^2 = \frac{1}{m} \sum_{i=1}^m (y_i - \bar{y})^2 = \frac{SST}{m} = \frac{15.23}{15} = 1.015$$

$$\hat{\sigma}^2 = \frac{1}{m} \sum_{i=1}^m (y_i - \hat{y})^2 = \frac{SSE}{m} = \frac{2.308}{15} = 0.1539$$

Hence, the test statistic corresponds to

$$f^{\text{obs}} = \frac{(\hat{\sigma}^2 - \hat{\sigma}^2) / (J-1)}{\hat{\sigma}^2 / (m-J)} = \frac{SSR / (J-1)}{SSE / (m-J)} = 33.59 \text{ as ex. 3.1}$$

We can find R^2 by calculating

$$R^2 = \frac{SSR}{SST} = \frac{15.23}{15.23} - \frac{2.308}{15.23} = 0.848457$$

The model explains 84.8% of the total variability of the response variable.

4 Exercise 4

Lab Precise got some measurements to see whether the tar contents (in milligrams) for three different brands of cigarettes are different. The measurements are showed in the following table.

Sample	Brand A	Brand B	Brand C
1	10.21	11.32	11.60
2	10.25	11.20	11.90
3	10.24	11.40	11.80
4	9.80	10.50	12.30
5	9.77	10.68	12.20
6	9.73	10.90	12.20

Write the equation of the linear regression model. Find and interpret the estimates of regression coefficients.

Equation of linear regression model :

$$y_i = \beta_1 + \beta_2 x_{i2} + \beta_3 x_{i3} + \varepsilon_i \quad \varepsilon_i \sim N(0, \sigma^2) \quad i=1, \dots, 18$$

where

$$x_{i2} = \begin{cases} 1, & \text{if Brand = B} \\ 0, & \text{otherwise} \end{cases} \quad x_{i3} = \begin{cases} 1, & \text{if Brand = C} \\ 0, & \text{otherwise} \end{cases}$$

As in the previous exercise, we can find the ML estimates as follows

$$\hat{\beta}_1 = \bar{y}_A, \quad \hat{\beta}_2 = \bar{y}_B - \bar{y}_A, \quad \hat{\beta}_3 = \bar{y}_C - \bar{y}_A$$

$$\hat{\beta}_1 = 10, \quad \hat{\beta}_2 = 11 - 10 = 1, \quad \hat{\beta}_3 = 12 - 10 = 2$$

- $\hat{\beta}_1 = 10$ represents the average of tar contents (in milligrams) when the Brand is A

- $\hat{\beta}_2$ is the average change in tar contents when the Brand is B

- $\hat{\beta}_3$ is the average change in tar contents when the Brand is C

Exercise 4.1

Complete the following table.

Sources of variation	D.o.f.	Deviance
Total	17	13.3748
Regression	2	12
Residual	15	1.3748

Then, perform the statistical test to evaluate if the means (of each group) are homogeneous.

We can compute SST, SSR and SSE as follows (you can also check the previous exercise)

$$SSR = \sum_{j=1}^3 m_j (\bar{y}_j - \bar{y})^2 = 16 \cdot (10-11)^2 + 6 \cdot (0)^2 + 6 \cdot (1)^2 = 12$$

$$SSE = \sum_{j=1}^3 \sum_{i=1}^{m_j} (y_{ij} - \bar{y}_j)^2 = 0.33 + 0.6648 + 0.38 = 1.3748$$

$$SST = SSR + SSE = 12 + 1.3748 = 13.3748$$

The degrees of freedom for SSR is equal to $J-1 = 3-1 = 2$

" " SSE is equal to $m-J = 18-3 = 15$

" " SST is equal to $m-1 = 18-1 = 17$

System of hypothesis

$$\left\{ H_0: \mu_A = \mu_B = \mu_C \right.$$

$$\left. H_1: \exists j \text{ st } \mu_j \neq \mu_i \quad (i \neq j, i, j \in \{A, B, C\}) \right)$$

test statistic:

$$f_{\text{obs}} = \frac{SSR / (J-1)}{SSE / (m-J)} = \frac{12/2}{1.3748/15} = 65.46404$$

P-value

$$d_{\text{obs}} = P(F_{2,15} > 65.46404) \approx \text{We reject } H_0$$