

# Regression and ANOVA

Justin Patterson

2021-06-11



# Contents

<b>Preface</b>	<b>5</b>
<b>1 Introduction</b>	<b>7</b>
<b>2 Literature</b>	<b>9</b>
<b>3 Methods</b>	<b>11</b>
<b>4 Applications</b>	<b>13</b>
4.1 Example one . . . . .	13
4.2 Example two . . . . .	13
<b>5 Final Words</b>	<b>15</b>
<b>6 ANOVA Fundamentals</b>	<b>17</b>
6.1 Law of Total Variance . . . . .	17
6.2 Step 1: Make up Data . . . . .	18
6.3 Checking the Assumptions . . . . .	18
6.4 Regression and Categorical Variables . . . . .	20



# Preface

This project is meant to be a personal guide to ANOVA and regression. The scope of this project does not include time series analysis. Also, the focus will not be on designing experiments, but rather on analyzing the data from experiments which have already been conducted. To accomplish this, we will use simulated experimental data.

## **Disclaimer**

This project was not written by an expert. That being said, I would appreciate any comments.

**Note** This book was constructed with the **bookdown** package [7], which was built on top of R Markdown and **knitr** [8].



# Chapter 1

## Introduction

You can label chapter and section titles using `{#label}` after them, e.g., we can reference Chapter 1. If you do not manually label them, there will be automatic labels anyway, e.g., Chapter 3.

Figures and tables with captions will be placed in `figure` and `table` environments, respectively.

```
par(mar = c(4, 4, .1, .1))  
plot(pressure, type = 'b', pch = 19)
```

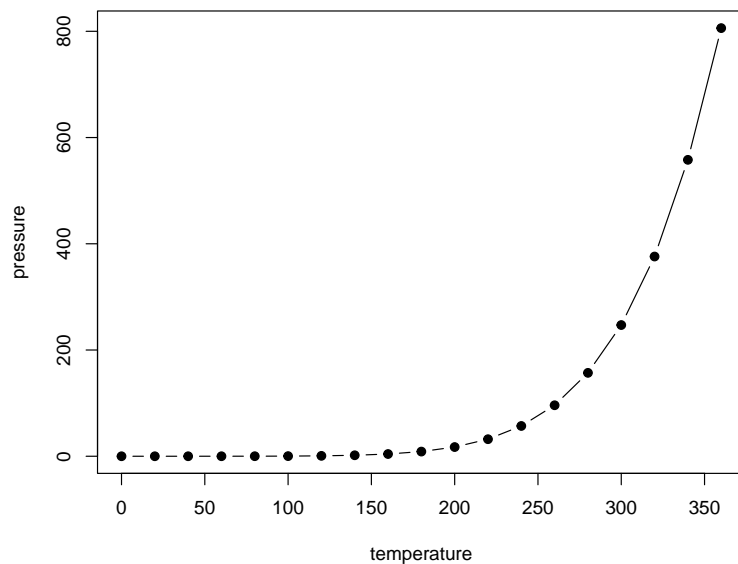


Figure 1.1: Here is a nice figure!

Table 1.1: Here is a nice table!

Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species
5.1	3.5	1.4	0.2	setosa
4.9	3.0	1.4	0.2	setosa
4.7	3.2	1.3	0.2	setosa
4.6	3.1	1.5	0.2	setosa
5.0	3.6	1.4	0.2	setosa
5.4	3.9	1.7	0.4	setosa
4.6	3.4	1.4	0.3	setosa
5.0	3.4	1.5	0.2	setosa
4.4	2.9	1.4	0.2	setosa
4.9	3.1	1.5	0.1	setosa
5.4	3.7	1.5	0.2	setosa
4.8	3.4	1.6	0.2	setosa
4.8	3.0	1.4	0.1	setosa
4.3	3.0	1.1	0.1	setosa
5.8	4.0	1.2	0.2	setosa
5.7	4.4	1.5	0.4	setosa
5.4	3.9	1.3	0.4	setosa
5.1	3.5	1.4	0.3	setosa
5.7	3.8	1.7	0.3	setosa
5.1	3.8	1.5	0.3	setosa

Reference a figure by its code chunk label with the `fig:` prefix, e.g., see Figure 1.1. Similarly, you can reference tables generated from `knitr::kable()`, e.g., see Table 1.1.

```
knitr::kable(
  head(iris, 20), caption = 'Here is a nice table!',
  booktabs = TRUE
)
```



## Chapter 2

# Literature

Here is a review of existing methods.



## Chapter 3

# Methods

We describe our methods in this chapter.



## Chapter 4

# Applications

Some *significant* applications are demonstrated in this chapter.

### 4.1 Example one

### 4.2 Example two



## Chapter 5

# Final Words

We have finished a nice book.





## Chapter 6

# ANOVA Fundamentals

Analysis of variance (ANOVA) is a collection of statistical models and their associated estimation procedures used to analyze the differences among means. ANOVA is based on the *Law of Total Variance* (6.1), where the observed variance in a particular variable is partitioned into components attributable to different sources of variation.[1]

### 6.1 Law of Total Variance

The law of total variance, also known as EVE's law [2][3], is very important for understanding how ANOVA works.

$$\text{Var}(Y) = \text{E}[\text{Var}(Y|\mathbf{X})] + \text{Var}[\text{E}(Y|\mathbf{X})] \quad (6.1)$$

In the context of ANOVA with a response variable  $Y$  and a covariate vector  $\mathbf{X}$ , 6.1 can be interpreted as

$$\text{Var}(Y) = \underbrace{\text{E}[\text{Var}(Y|\mathbf{X})]}_{\text{variance of } Y \text{ within } X} + \underbrace{\text{Var}[\text{E}(Y|\mathbf{X})]}_{\text{variance of } Y \text{ between } X} \quad (6.2)$$

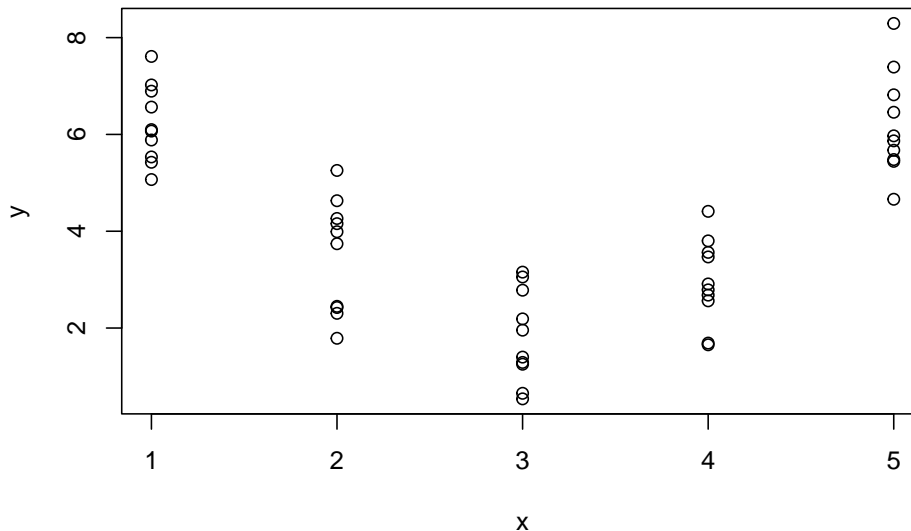
Recall [6] that the population variance of a finite population of size  $N$  is

$$\sigma^2 = \frac{1}{N} \underbrace{\sum_{i=1}^N (x_i - \mu)^2}_{\text{SS}} \quad (6.3)$$

where the sum is known as a sum of squares (SS). It is easy to see that both sides of 6.2 can be multiplied by  $N$  to give a relationship among sums of squares. The resulting relationship is known as a partition of the sum of squares. [5]

```
source(file.path("src", "law-of-total-variance-sim.R"))
```

```
## [1] "If we group Y by X, then we might understand Y's variability better."
```



```
## [1] "The population variance of Y is 4.01785074652354"
```

```
## [1] "The population variance of Y within X is 0.867547839457478"
```

```
## [1] "The population variance of Y between X is 3.15030290706606"
```

```
## [1] "The variance of Y = the variance of Y within X + the variance of Y between X"
```

The total sum of squares divided by the total degrees of freedom is the total variance in the response variable.

```
library(cellWise)
```

```
library(knitr)
```

```
opts_chunk$set(tidy.opts=list(width.cutoff=50),tidy=TRUE)
```

## 6.2 Step 1: Make up Data

```
# dataset1
```

## 6.3 Checking the Assumptions

After running your ANOVA, check that the assumptions about the errors are met so that you can do statistical inference. Those assumptions are:

1.  $E(\epsilon_{ij}) = 0$ ,  $\text{Var}(\epsilon_{ij}) = \sigma_i^2 < \infty$ , for all  $i, j$ .
2. The  $\epsilon_{ij}$  are mutually independent and normally distributed.
3.  $\sigma_i^2 = \sigma^2$  for all  $i$ .

### 6.3.1 Checking Assumption 1

### 6.3.2 Assumption 1 was violated.

### 6.3.3 Checking Assumption 2

### 6.3.4 Assumption 2 was violated.

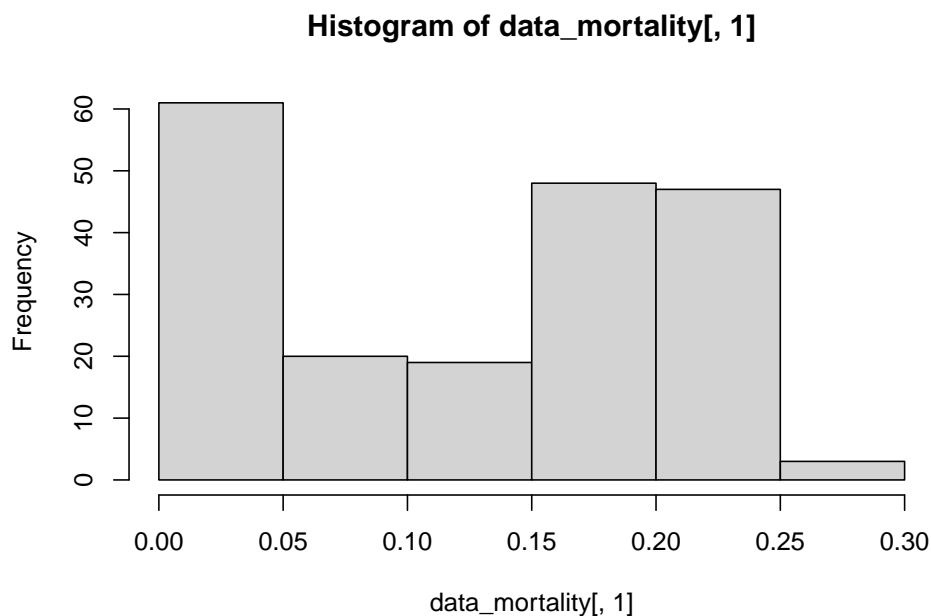
### 6.3.5 Checking Assumption 3

### 6.3.6 Assumption 3 was violated.

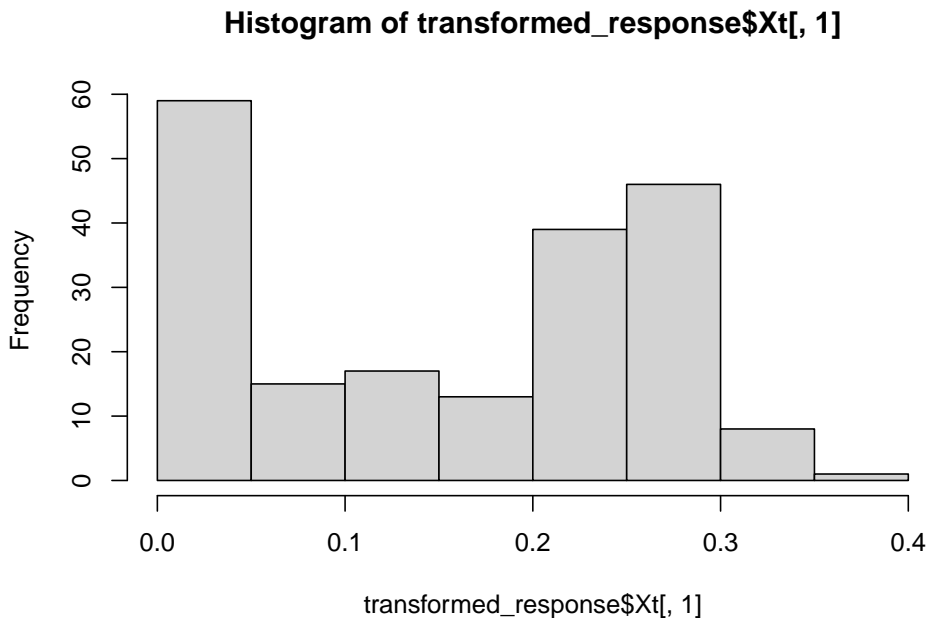
A variance-stabilizing transformation of the response variable may help.

```
data("data_mortality")
transformed_response = transfo(data_mortality, prestandardize = FALSE)
```

```
##
## The input data has 198 rows and 91 columns.
hist(data_mortality[, 1])
```



```
hist(transformed_response$Xt[, 1])
```



```
shapiro.test(data_mortality[, 1])
```

```
##
##  Shapiro-Wilk normality test
##
## data:  data_mortality[, 1]
## W = 0.86877, p-value = 4.552e-12
```

```
shapiro.test(transformed_response$Xt[, 1])
```

```
##
##  Shapiro-Wilk normality test
##
## data:  transformed_response$Xt[, 1]
## W = 0.88041, p-value = 1.968e-11
```

## 6.4 Regression and Categorical Variables

```
library(tidymodels)
```

```
## Registered S3 method overwritten by 'tune':
## method                from
## required_pkgs.model_spec parsnip

## -- Attaching packages ----- tidymodels 0.1.3 --
## v broom                0.7.6      v recipes                0.1.16
```

```
## v dials          0.0.9      v rsample      0.1.0
## v dplyr          1.0.6      v tibble     3.1.2
## v ggplot2        3.3.3      v tidyr     1.1.3
## v infer          0.5.4      v tune      0.1.5
## v modeldata      0.1.0      v workflows 0.2.2
## v parsnip        0.1.6      v workflowsets 0.0.2
## v purrr          0.3.4      v yardstick 0.0.8

## -- Conflicts ----- tidymodels_conflicts() --
## x purrr::discard() masks scales::discard()
## x dplyr::filter()  masks stats::filter()
## x dplyr::lag()     masks stats::lag()
## x recipes::step()  masks stats::step()
## * Use tidymodels_prefer() to resolve common conflicts.

library(ggplot2)
```

There is a profound connection between linear regression and ANOVA. In order to see this, you have to understand that the categorical variables of an ANOVA can be coded with numbers, which allows them to be used in a linear regression model. Let us recall [4] the multiple linear regression model.

Given a random sample of  $n$  observations  $(Y_i, X_{i1}, \dots, X_{ip})$ ,  $i = 1, \dots, n$ , the basic multiple linear regression model is

$$Y_i = \beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{ip} + \epsilon_i, \quad i = 1, \dots, n$$

where each  $\epsilon_i$  is a random variable with a mean of 0. In matrix form, this can be written as

$$\begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} 1 & X_{1,1} & X_{1,2} & \dots & X_{1,p} \\ 1 & X_{2,1} & X_{2,2} & \dots & X_{2,p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & X_{n,1} & X_{n,2} & \dots & X_{n,p} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix} + \begin{bmatrix} \epsilon_0 \\ \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

Here, the  $X_{i,j}$  represent our coded categorical variables. These categorical variables are coded according to the hypotheses of interest. In many cases, the coding is done so that the newly coded variables are contrasts of the old categorical variables.

A contrast is a linear combination of variables such that the coefficients sum to 0.

$$\sum_i a_i \theta_i \quad \text{such that} \quad \sum_i a_i = 0$$

Unlike in ANOVA, in regression, it is best to use coding schemes based on orthogonal and fractional contrasts. Orthogonal contrasts are a set of contrasts

in which, for any distinct pair, the sum of the cross-products of the coefficients is 0.

$$\sum_i a_i b_i = 0$$

I believe that a fractional contrast is such that

$$\sum_i |a_i| = 2$$

Categorical variable coding schemes can be easily expressed in a matrix format. The convention is to have the old categorical variables as the row headers and the newly coded variables as the column headers. In such a matrix, the  $[c_{ij}]$  entry indicates the value of the  $j^{th}$  level of the new variable for the  $i^{th}$  level of the old variable. Here is an example of such a matrix constructed using orthogonal and fractional contrasts.

```
(contr_mat = matrix(data = c(1, 0, -1, 0.5, -1, 0.5),
  nrow = 3, ncol = 2))
```

```
##      [,1] [,2]
## [1,]    1  0.5
## [2,]    0 -1.0
## [3,]   -1  0.5
```

Interpreting this coding scheme in the context of our linear model, we see that

$$\begin{aligned} E(Y_i | X_{i1} = 1, X_{i2} = \tfrac{1}{2}) &= \beta_0 + \beta_1 + \tfrac{1}{2}\beta_2 = \mu_1 \\ E(Y_i | X_{i1} = 0, X_{i2} = -1) &= \beta_0 - \beta_2 = \mu_2 \\ E(Y_i | X_{i1} = -1, X_{i2} = \tfrac{1}{2}) &= \beta_0 - \beta_1 + \tfrac{1}{2}\beta_2 = \mu_3 \end{aligned}$$

or, in matrix format,

$$\begin{bmatrix} 1 & 1 & \frac{1}{2} \\ 1 & 0 & -1 \\ 1 & -1 & \frac{1}{2} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \end{bmatrix}$$

We can solve this for  $\beta$  for interpretation's sake.

```
solve(cbind(rep(1, nrow(contr_mat)), contr_mat))
```

```
##      [,1]      [,2]      [,3]
## [1,] 0.3333333 0.3333333 0.3333333
## [2,] 0.5000000 0.0000000 -0.5000000
## [3,] 0.3333333 -0.6666667 0.3333333
```

$$\begin{aligned}
\beta_0 &= \frac{\mu_1 + \mu_2 + \mu_3}{3} &= & \text{grand mean response} \\
2\beta_1 &= \mu_1 - \mu_3 &= & \text{difference in the mean response between levels 1} \\
& & & \text{and 3 of the old categorical variable} \\
\frac{3}{2}\beta_2 &= \frac{\mu_1 + \mu_3}{2} - \mu_2 &= & \text{difference in the mean response between level 2} \\
& & & \text{and the average of levels 1 and 3 of the old cate-} \\
& & & \text{gorical variable}
\end{aligned}$$

Let's look at another contrast matrix and see if we can interpret it.

```
contr.helmert(n = 3)
```

```
##      [,1] [,2]
## 1     -1  -1
## 2      1  -1
## 3      0   2
```

```
solve(cbind(rep(1, 3), contr.helmert(n = 3)))
```

```
##           1           2           3
## [1,] 0.3333333 0.3333333 0.3333333
## [2,] -0.5000000 0.5000000 0.0000000
## [3,] -0.1666667 -0.1666667 0.3333333
```

$$\begin{aligned}
\beta_0 &= \frac{\mu_1 + \mu_2 + \mu_3}{3} &= & \text{grand mean response} \\
2\beta_1 &= \mu_2 - \mu_1 &= & \text{difference in the mean response between levels 2} \\
& & & \text{\& 1 of the old categorical variable} \\
3\beta_2 &= \mu_3 - \frac{\mu_1 + \mu_2}{2} &= & \text{difference in the mean response between level 3} \\
& & & \text{and the average of levels 1 and 2 of the old cate-} \\
& & & \text{gorical variable}
\end{aligned}$$

Perhaps you have heard of polynomial regression? Polynomial regression is just a special case of linear regression in a different basis. In polynomial regression, (just like multiple linear regression) if you use all of your explanatory variables, then you will likely get multi-collinearity problems.

```
contr.poly(n = 3)
```

```
##           .L           .Q
## [1,] -7.071068e-01  0.4082483
## [2,] -7.850462e-17 -0.8164966
## [3,]  7.071068e-01  0.4082483
```

```
(A = solve(cbind(rep(1, 3), contr.poly(n = 3))))
```

```
##           [,1]      [,2]      [,3]
##      0.3333333 0.3333333 0.3333333
## .L -0.7071068 0.0000000 0.7071068
## .Q  0.4082483 -0.8164966 0.4082483
```

The first matrix shows how to code the levels of your categorical variable and the second matrix is used for interpretation.

$$\begin{aligned}
\beta_0 &= \frac{\mu_1 + \mu_2 + \mu_3}{3} &&= \text{grand mean response} \\
\beta_1 &= -0.707\mu_1 + 0.707\mu_3 &&= \text{measure of a linear trend in the mean response} \\
\beta_2 &= 0.408\mu_3 - 0.816\mu_2 + 0.408\mu_1 &&= \text{measure of a quadratic trend in the mean response}
\end{aligned}$$

For example, we can test whether the difference between the means from two populations are equal by doing a linear regression or an ANOVA.

Let's make up some data and try it!

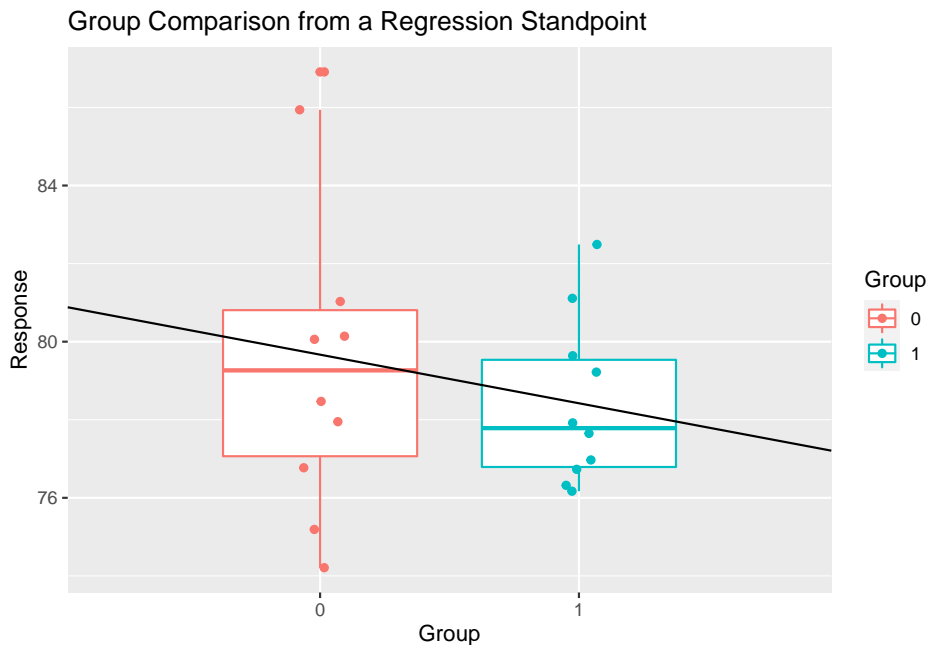
```
source(file.path("src", "fabricate.R"))
design = data.frame(group = c(0, 1), n = c(10, 10))
data1 = fabricate(flr = design)
```

Let's check out our data.

```
# Make a linear model
data1_lm_independent_samples = lm(response ~ group,
  data = data1)
# plot
ggplot(data = data1, aes(x = group, y = response, color = factor(group))) +
  geom_boxplot() + geom_jitter(height = 0, width = 0.1) +
  geom_abline(intercept = data1_lm_independent_samples$coefficients[1],
    slope = data1_lm_independent_samples$coefficients[2]) +
  labs(title = "Group Comparison from a Regression Standpoint",
    color = "Group", x = "Group", y = "Response") +
  scale_x_discrete(limits = c(0, 1))
```

```
## Warning: Continuous limits supplied to discrete scale.
## Did you mean `limits = factor(...)` or `scale_*_continuous()`?
```





The way you code your categorical variables in a linear model is extremely important. Different codings lead to different interpretations of the parameters (betas) in your model. For us, our model is

$$Y_i = \beta_0 + \beta_{i1}X_{i1} + \epsilon_i$$

From this, we have

$$E(Y_i|X_{i1} = 0) = \beta_0$$

$$E(Y_i|X_{i1} = 1) = \beta_0 + \beta_1$$

From which we can derive,

$$\beta_1 = E(Y_i|X_{i1} = 1) - E(Y_i|X_{i1} = 0)$$

So, our slope estimate is the estimated amount by which the mean of group1 is above that of the mean of group0.

Run linear regression

```
summary(data1_lm_independent_samples)
```

```
##
## Call:
## lm(formula = response ~ group, data = data1)
##
## Residuals:
```

```
##      Min      1Q  Median      3Q      Max
## -5.457 -1.813 -0.637   1.254   7.243
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    79.667     1.048   75.99  <2e-16 ***
## group         -1.245     1.483   -0.84   0.412
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.315 on 18 degrees of freedom
## Multiple R-squared:  0.0377, Adjusted R-squared:  -0.01577
## F-statistic: 0.7051 on 1 and 18 DF,  p-value: 0.4121
```

Run ANOVA

```
data1$group = as.factor(data1$group)
data1_ANOVA_independent_samples = aov(response ~ group,
  data = data1)
summary(data1_ANOVA_independent_samples)
```

```
##              Df Sum Sq Mean Sq F value Pr(>F)
## group          1   7.75    7.75    0.705  0.412
## Residuals     18 197.85   10.99
```

Run t-Test

```
(data1_t_test_independent_samples = t.test(x = data1[data1$group ==
  1, "response"], y = data1[data1$group == 0, "response"],
  paired = FALSE, var.equal = TRUE))
```

```
##
## Two Sample t-test
##
## data: data1[data1$group == 1, "response"] and data1[data1$group == 0, "response"]
## t = -0.8397, df = 18, p-value = 0.4121
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -4.359965  1.869965
## sample estimates:
## mean of x mean of y
## 78.422 79.667
```

Notice the similarities.

```
# Confidence interval for the difference in the
# means
confint(data1_lm_independent_samples, "group", level = 0.95)
```

```
##           2.5 %    97.5 %
## group -4.359965  1.869965

data1_t_test_independent_samples$conf.int

## [1] -4.359965  1.869965
## attr(,"conf.level")
## [1] 0.95

# p-values
with(summary(data1_lm_independent_samples), unname(pf(fstatistic[1],
  fstatistic[2], fstatistic[3], lower.tail = F)))

## [1] 0.412089

summary(data1_ANOVA_independent_samples)[[1]][[1, 5]]

## [1] 0.412089

data1_t_test_independent_samples$p.value

## [1] 0.412089
```

Now, let's look at something else. The CO2 data frame has 84 rows and 5 columns of data from an experiment on the cold tolerance of the grass species *Echinochloa crus-galli*.

```
data("CO2")
CO2[sample(nrow(CO2), size = 5), ]

##   Plant      Type Treatment conc uptake
## 15   Qn3      Quebec nonchilled   95   16.2
## 78   Mc3 Mississippi   chilled   95   10.6
## 46   Mn1 Mississippi nonchilled  350   30.0
## 10   Qn2      Quebec nonchilled  250   37.1
## 18   Qn3      Quebec nonchilled  350   42.1
```

What is a linear model? In the context of linear regression, a linear model is a relationship between the responses and the explanatory variables that is linear in the parameters.

```
CO2_recipe = recipe(uptake ~ ., data = CO2) %>%
  step_dummy(c("Type", "Treatment"))
# see contrasts() function
CO2_linear_model = linear_reg() %>%
  set_engine("lm", contrasts = list(Plant = "contr.poly"))
CO2_workflow = workflow() %>%
  add_model(CO2_linear_model) %>%
  add_recipe(CO2_recipe)
CO2_fit = CO2_workflow %>%
  fit(data = CO2)
```

```
CO2_fit %>%
  pull_workflow_fit() %>%
  tidy()
```

```
## # A tibble: 15 x 5
```

##	term	estimate	std.error	statistic	p.value
##	<chr>	<dbl>	<dbl>	<dbl>	<dbl>
##	1 (Intercept)	19.5	1.17	16.7	2.96e-26
##	2 Plant.L	-22.9	2.27	-10.1	2.17e-15
##	3 Plant.Q	-4.62	2.27	-2.03	4.57e- 2
##	4 Plant.C	4.67	2.27	2.06	4.34e- 2
##	5 Plant^4	2.34	2.27	1.03	3.06e- 1
##	6 Plant^5	4.31	2.27	1.90	6.13e- 2
##	7 Plant^6	-0.0390	2.27	-0.0172	9.86e- 1
##	8 Plant^7	-2.04	2.27	-0.897	3.73e- 1
##	9 Plant^8	-3.28	2.27	-1.44	1.53e- 1
##	10 Plant^9	-9.07	2.27	-4.00	1.56e- 4
##	11 Plant^10	0.546	2.27	0.241	8.10e- 1
##	12 Plant^11	1.91	2.27	0.843	4.02e- 1
##	13 conc	0.0177	0.00223	7.96	1.97e-11
##	14 Type_Mississippi	NA	NA	NA	NA
##	15 Treatment_chilled	NA	NA	NA	NA

# Bibliography

- [1] *Analysis of variance*. Wikipedia. URL: [https://en.wikipedia.org/wiki/Analysis\\_of\\_variance](https://en.wikipedia.org/wiki/Analysis_of_variance) (visited on 06/09/2021).
- [2] *Law of total variance*. Wikipedia. URL: [https://en.wikipedia.org/wiki/Law\\_of\\_total\\_variance](https://en.wikipedia.org/wiki/Law_of_total_variance) (visited on 06/09/2021).
- [3] *Law of total variance intuition*. Mathematics StackExchange. URL: <https://math.stackexchange.com/a/3377007> (visited on 06/10/2021).
- [4] *Linear model*. Wikipedia. URL: [https://en.wikipedia.org/wiki/Linear\\_model](https://en.wikipedia.org/wiki/Linear_model) (visited on 05/31/2021).
- [5] *Partition of sums of squares*. Wikipedia. URL: [https://en.wikipedia.org/wiki/Partition\\_of\\_sums\\_of\\_squares](https://en.wikipedia.org/wiki/Partition_of_sums_of_squares) (visited on 06/11/2021).
- [6] *Variance*. Wikipedia. URL: [https://en.wikipedia.org/wiki/Variance#Population\\_variance](https://en.wikipedia.org/wiki/Variance#Population_variance) (visited on 06/11/2021).
- [7] Yihui Xie. *bookdown: Authoring Books and Technical Documents with R Markdown*. R package version 0.22. 2021. URL: <https://CRAN.R-project.org/package=bookdown>.
- [8] Yihui Xie. *Dynamic Documents with R and knitr*. 2nd. ISBN 978-1498716963. Boca Raton, Florida: Chapman and Hall/CRC, 2015. URL: <http://yihui.name/knitr/>.