

Regression and ANOVA

Justin Patterson

2021-07-05

Contents

Preface	5
1 Introduction	7
1.1 The Setting	7
1.2 The Process	7
2 ANOVA Fundamentals	11
2.1 Law of Total Variance	11
2.2 Partitioning the SS	12
2.3 The ANOVA Table	14
2.4 Step 1: Make up Data	14
2.5 Checking the Assumptions	14
2.6 Regression and Categorical Variables	16
3 Regression Fundamentals	25

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 [13], which was built on top of R Markdown and **knitr** [14].

Chapter 1

Introduction

1.1 The Setting

A statistical analysis plan has been specified and a corresponding experiment has been performed. You are now ready to analyze the data according to plan. Questions you may have include: Are the results due to chance? Do some of the treatments explain more of the variation in the response than the other treatments? What are simultaneous confidence intervals for the mean response under each treatment?

1.2 The Process

This booklet provides the guidance needed to turn raw experimental data into results!

1.2.1 Models Provide the Foundation for Statistical Inference

Statistical models (whether they be parametric or non-parametric) provide the foundation for statistical estimation and null hypothesis significance testing (NHST) [11]. Therefore, analyzing data from an experiment requires some kind of statistical model of how such data was generated. Analysis of Variance (ANOVA) is a collection of statistical models and their associated estimation procedures used to analyze differences among means [1]. It is by far the most common paradigm for analyzing experimental data. Statistical models used in ANOVA can be broadly classified into one of three categories:

1. Fixed-Effects Models

2. Random-Effects Models

3. Mixed-Effects Models

We will mainly focus on fixed-effects models.

1.2.2 What Does Regression Have to do with it?

It should be noted that the central NHST of an ANOVA is the F -Test. Interestingly, the hypothesis that a proposed regression model fits the data well and the hypothesis that a data set in a regression analysis follows the simpler of two proposed linear models that are nested within each other all are tested using F -Tests [5].

Regression analysis is a collection of statistical models and their associated estimation procedures used to analyze the relationship between one or more response variables and one or more explanatory variables [10]. If we think of the explanatory variables as the levels of the factors in the experiment, then regression analysis can be used to analyze the relationship between the factors and the response variable just like ANOVA.

Indeed, regression analyses often include an ANOVA table (just like ANOVA) that also is based on partitioning sums of squares (measures of variability).

To do this we will

1. Give brief insights into the finer points of the theory.
 - Clarify key relationships.
 - Provide synonyms for common terminology.
2. Go through guided walk-throughs with simulated data.
 - Check assumptions.
 - Compare to results from permutation tests and bootstrapping.

We will give some theory and general guidelines. We get practice with checking the assumptions. For several simulated datasets, we will perform the classic, omnibus ANOVA F -Test and will estimate confidence intervals for linear combinations of means. We will compare our results to those from permutation tests and bootstrapping.

We will compare the results of various NHSTs of contrasts of means of various classic ANOVA null hand compare the results to permutation tests. is the multiple linear model (which also shows up in multiple linear regression). We will attempt to explain how the multiple linear model ties together ANOVA and multiple linear regression.

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

We will cover the theory behind ANOVA and its relationship to regression.

You can label chapter and section titles using `{#label}` after them, e.g., we can reference Chapter 1.

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

Reference a figure by its code chunk label with the `fig:` prefix. 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

ANOVA Fundamentals

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

Analysis of variance (ANOVA) is a collection of statistical models and their associated estimation procedures used to compare means. It may be difficult for the neophyte to understand, but ANOVA compares means by analyzing variances. ANOVA is based on the *Law of Total Variance* (2.1), where the observed variance in a particular variable is partitioned into components attributable to different sources of variation.[1]

2.1 Law of Total Variance

The law of total variance, also known as EVE's law [7][8], is very important for understanding how ANOVA works. The law says that for random variables X and Y on the same probability space, we have

$$\text{Var}(Y) = \text{E}[\text{Var}(Y|X)] + \text{Var}[\text{E}(Y|X)]. \quad (2.1)$$

Let us consider the case of One-way ANOVA using Fixed Effects. We observe the random variable Y_{ij} , the j^{th} response for the i^{th} level of the single factor. Let A_i be the event where the i^{th} level of this factor is observed. Then 2.1 can be interpreted as

$$\text{Var}(Y_{ij}) = \underbrace{\text{E}[\text{Var}(Y_{ij}|A_i)]}_{\text{variance of } Y \text{ within groups}} + \underbrace{\text{Var}[\text{E}(Y_{ij}|A_i)]}_{\text{variance of } Y \text{ between groups}} \quad (2.2)$$

Can 2.2 be simplified into something computationally useful for our ANOVA?

$$E(Y_{ij}|A_i) = \sum_{j=1}^{n_i} y_{ij} \frac{1}{n_i} = \bar{y}_i. \quad (2.3)$$

$$\text{Var}(Y_{ij}|A_i) = \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2 \frac{1}{n_i} \quad (2.4)$$

$$\sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{\bullet\bullet})^2 \frac{1}{N} = \sum_{i=1}^k \left(\sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{i\bullet})^2 \frac{1}{n_i} \right) \frac{n_i}{N} + \sum_{i=1}^k (\bar{y}_{i\bullet} - \bar{y}_{\bullet\bullet})^2 \frac{n_i}{N} \quad (2.5)$$

After multiplying both sides by N , we have this partition of the SS

$$\sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{\bullet\bullet})^2 = \sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{i\bullet})^2 + \sum_{i=1}^k n_i (\bar{y}_{i\bullet} - \bar{y}_{\bullet\bullet})^2 \quad (2.6)$$

Recall [4] the definition of the mean of a population. Let X be a random variable with a finite number of finite outcomes

Recall [6][3] the definition of the expected value of a continuous random variable conditioned on an event.

Recall [12] that the population variance of a finite population of size N with values y_i and mean μ is

$$\sigma^2 = \frac{1}{N} \underbrace{\sum_{i=1}^N (y_i - \mu)^2}_{\text{SS}} \quad (2.7)$$

where the sum is known as a sum of squares (SS).

2.2 Partitioning the SS

```
source(file.path("src", "get-SS.R"))
```

Let's look at this simulated experiment to see the SS Decomposition in action. How can we explain the variation that we see in the response?

```
source(file.path("src", "partitioning-SS.R"))
```

```
##
## Attaching package: 'dplyr'

## The following objects are masked from 'package:stats':
##
##   filter, lag

## The following objects are masked from 'package:base':
##
##   intersect, setdiff, setequal, union
```

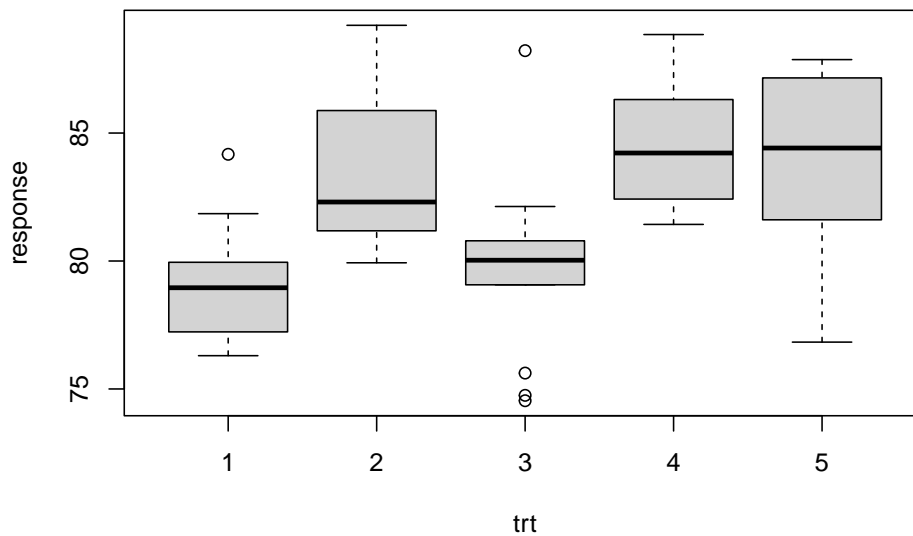


Figure 2.1: Responses were simulated for five different levels of a single factor. The respective sample sizes for each treatment group were 10, 14, 13, 9, and 8.

$$SS_{\text{total}} = \sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{\bullet\bullet})^2 = 760.608$$

$$SS_{\text{within}} = \sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{i\bullet})^2 = 502.377$$

$$SS_{\text{between}} = \sum_{i=1}^k n_i (\bar{y}_{i\bullet} - \bar{y}_{\bullet\bullet})^2 = 258.231$$

Notice that $SS_{\text{total}} = SS_{\text{within}} + SS_{\text{between}}$. Also, $\frac{SS_{\text{total}}}{n-1} = s^2 = 14.351$. We get these same numbers from an ANOVA as well.

```
anova(aov(response ~ trt, data = fabricated_data))

## Analysis of Variance Table
##
## Response: response
##           Df Sum Sq Mean Sq F value    Pr(>F)
## trt         4 258.23   64.558   6.2967 0.0003598 ***
## Residuals 49 502.38   10.253
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

2.3 The ANOVA Table

The ANOVA table summarizes the variability that we see in the response. For more information on mean squares (MS), which are calculated by dividing the SS by their corresponding degrees of freedom (DF), please see [2].

Table 2.1: One-way ANOVA table

Source of Variation	DF	SS	MS	F
Between treatment groups	$k - 1$	$\sum_{i=1}^k n_i (\bar{y}_{i\bullet} - \bar{y}_{\bullet\bullet})^2$	$\frac{SS_{\text{between}}}{k - 1}$	$\frac{MS_{\text{between}}}{MS_{\text{within}}}$
Within treatment groups	$n - k$	$\sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{i\bullet})^2$	$\frac{SS_{\text{within}}}{n - k}$	

2.4 Step 1: Make up Data

```
# dataset1
```

2.5 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 ϵ_{ij} are mutually independent and normally distributed.
3. $\sigma_i^2 = \sigma^2$ for all i .

2.5.1 Checking Assumption 1

2.5.2 Assumption 1 was violated.

2.5.3 Checking Assumption 2

2.5.4 Assumption 2 was violated.

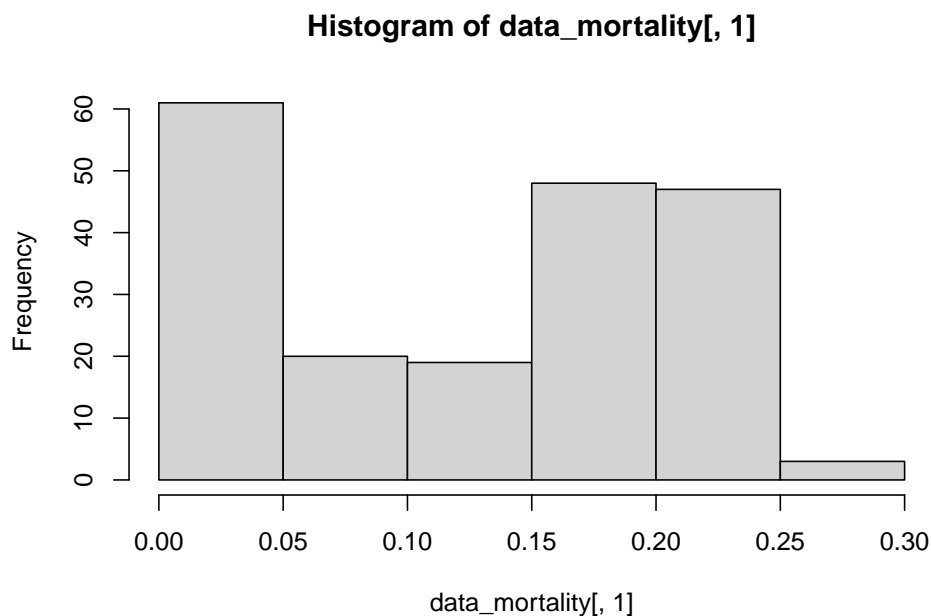
2.5.5 Checking Assumption 3

2.5.6 Assumption 3 was violated.

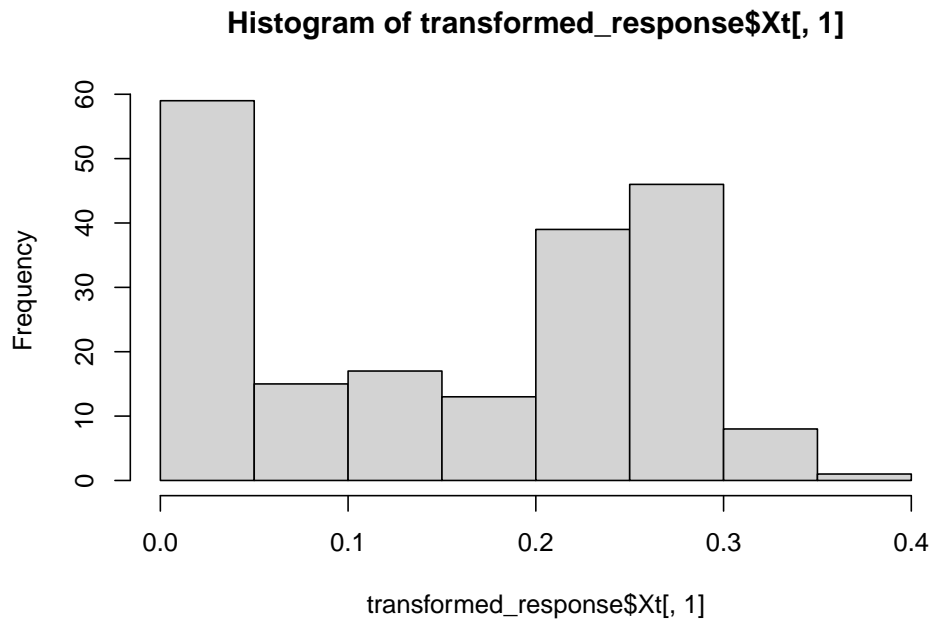
A variance-stabilizing transformation of the response variable may help.

```
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
```

2.6 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 rsample                0.1.0
```



```
## v dials          0.0.9      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
## v recipes        0.1.16

## -- 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 [9] 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 ϵ_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 β 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.

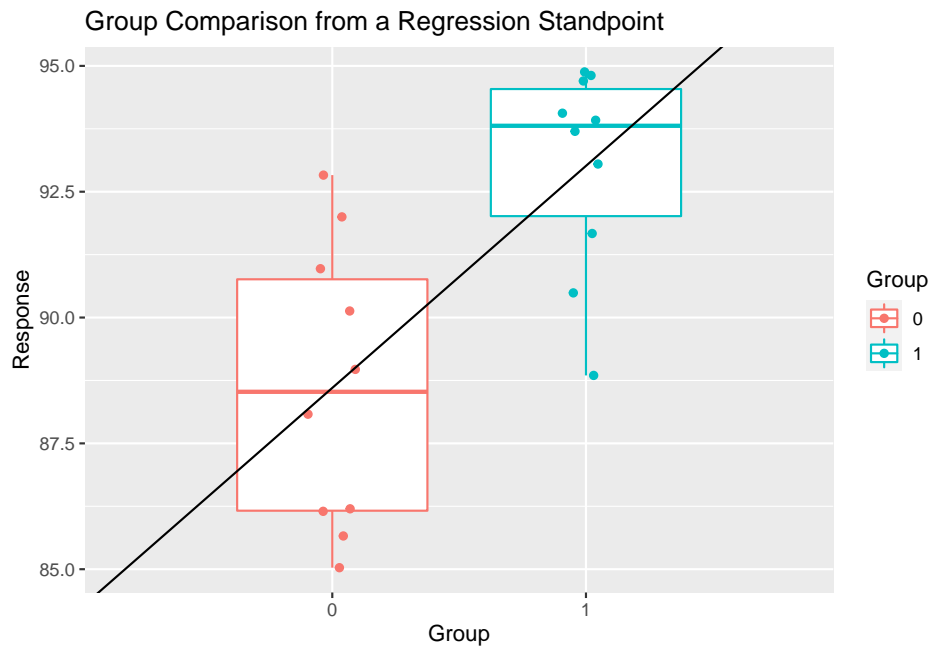
Let's make up some data and try it!

```
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
## -4.1630 -2.4145  0.5275  1.7145  4.2280
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  88.6020     0.7761 114.162 < 2e-16 ***
## group        4.4110     1.0976   4.019 0.000805 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.454 on 18 degrees of freedom
## Multiple R-squared:  0.4729, Adjusted R-squared:  0.4436
## F-statistic: 16.15 on 1 and 18 DF,  p-value: 0.0008053
```

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  97.28    97.28   16.15 0.000805 ***
## Residuals     18 108.42     6.02
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

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 = 4.0188, df = 18, p-value = 0.0008053
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##  2.105066 6.716934
## sample estimates:
## mean of x mean of y
##   93.013   88.602
```

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 2.105066 6.716934
data1_t_test_independent_samples$conf.int

## [1] 2.105066 6.716934
## 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] 0.0008053323
summary(data1_ANOVA_independent_samples)[[1]][[1, 5]]

## [1] 0.0008053323
data1_t_test_independent_samples$p.value

## [1] 0.0008053323

```

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
## 37   Qc3    Quebec   chilled   175   21.0
##  6   Qn1    Quebec nonchilled   675   39.2
## 84   Mc3 Mississippi   chilled  1000   19.9
## 51   Mn2 Mississippi nonchilled   175   22.0
## 76   Mc2 Mississippi   chilled   675   13.7

```

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

Chapter 3

Regression Fundamentals

Bibliography

- [1] *Analysis of variance*. Wikipedia. URL: https://en.wikipedia.org/wiki/Analysis_of_variance (visited on 06/09/2021).
- [2] *ANOVA: What is estimated by the Mean Squares?* CrossValidated. URL: <https://stats.stackexchange.com/questions/229649/anova-what-is-estimated-by-the-mean-squares> (visited on 06/18/2021).
- [3] George Casella and Roger L. Berger. *Statistical Inference*. 2nd ed. Brooks/Cole Cengage Learning, p. 150.
- [4] *Expected value*. Wikipedia. URL: https://en.wikipedia.org/wiki/Expected_value#Finite_case (visited on 06/14/2021).
- [5] *F-test*. Wikipedia. URL: <https://en.wikipedia.org/wiki/F-test> (visited on 06/19/2021).
- [6] *L09.2 Conditioning A Continuous Random Variable on an Event*. MIT OpenCourseWare. URL: https://www.youtube.com/watch?v=mHj4A1gh_ws (visited on 06/14/2021).
- [7] *Law of total variance*. Wikipedia. URL: https://en.wikipedia.org/wiki/Law_of_total_variance (visited on 06/09/2021).
- [8] *Law of total variance intuition*. Mathematics StackExchange. URL: <https://math.stackexchange.com/a/3377007> (visited on 06/10/2021).
- [9] *Linear model*. Wikipedia. URL: https://en.wikipedia.org/wiki/Linear_model (visited on 05/31/2021).
- [10] *Regression analysis*. Wikipedia. URL: https://en.wikipedia.org/wiki/Regression_analysis (visited on 06/19/2021).
- [11] *Statistical Model*. Wikipedia. URL: https://en.wikipedia.org/wiki/Statistical_model (visited on 06/19/2021).
- [12] *Variance*. Wikipedia. URL: https://en.wikipedia.org/wiki/Variance#Population_variance (visited on 06/11/2021).
- [13] 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>.
- [14] 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/>.