

ANOVA

The a previous lab we introduced the two-group independent t -test as a method for comparing the means of two groups. In some settings, it is useful to compare the means across more than two groups. The methodology behind a two-group independent t -test can be generalized to a procedure called **analysis of variance (ANOVA)**. Assessing whether the means across several groups are equal by conducting a single hypothesis test rather than multiple two-sample tests is important for controlling the overall Type I error rate.

The material in this lab corresponds to Section 7.5 of *OpenIntro Statistics*.

FAMuSS: comparing change in non-dominant arm strength by *ACTN3* genotype

Is change in non-dominant arm strength after resistance training associated with genotype?

In the Functional polymorphisms Associated with Human Muscle Size and Strength study (FAMuSS), researchers examined the relationship between muscle strength and genotype at a particular location on the *ACTN3* gene. The `famuss` dataset loaded below contains a subset of data from the study.

The percent change in non-dominant arm strength, comparing strength after resistance training to before training, is stored as `ndrm.ch`. There are three possible genotypes (CC, CT, TT) at the $r577x$ position on the *ACTN3* gene; genotype is stored as `actn3.r577x`.

1. Load the data. Create a plot that shows the association between change in non-dominant arm strength and *ACTN3* genotype. Describe what you see.

```
#load the data
load(url('https://github.com/jbryer/DATA606/blob/master/data/famuss.rda?raw=true'))
#create plot
```

2. Assess whether the assumptions for conducting an ANOVA are reasonably satisfied: 1) observations are independent within and across groups, 2) the data within each group are nearly normal, and 3) the variability across the groups is about equal.

```
#check assumptions
```

3. Conduct a hypothesis test to address the question of interest. Let $\alpha = 0.05$.
 - a) Let the parameters μ_{CC} , μ_{CT} , and μ_{TT} represent the population mean change in non-dominant arm strength for individuals of the corresponding genotype. State the null and alternative hypotheses.
 - b) Use `summary(aov())` to compute the F -statistic and p -value. Interpret the p -value.

```
#use summary(aov())
```

- c) Complete the analysis using pairwise comparisons.

- i. What is the appropriate significance level α^* for the individual comparisons, as per the Bonferroni correction?

$$\alpha^* = \alpha/K, \text{ where } K = \frac{k(k-1)}{2} \text{ for } k \text{ groups}$$

```
#use R as a calculator (note you need to change eval = TRUE once you assign the variables below)
alpha =
k =
K = (k*(k-1))/2
alpha.star = alpha/k
alpha.star
```

- ii. Use `pairwise.t.test()` to conduct the pairwise two-sample *t*-tests.

```
#The pairwise.t.test() command uses a comma between the two variables x and y,
# instead of a tilde like the aov() command. pairwise.t.test(y, x, p.adj = "")
# NOTE: change eval = TRUE once you have this working

#use pairwise.t.test() without adjusting p-value
pairwise.t.test( , , p.adj = "none")
#alternatively, use pairwise.t.test() with bonferroni adjustment
pairwise.t.test( , , p.adj = "bonf")
```

- iii. Summarize the results.

NHANES: comparing BMI by educational level

Is body mass index (BMI) associated with educational attainment?

This section uses data from the National Health and Nutrition Examination Survey (NHANES), a survey conducted annually by the US Centers for Disease Control (CDC).¹ The dataset `nhanes.samp.adult.500` contains data for 500 participants ages 21 years or older that were randomly sampled from the complete NHANES dataset that contains 10,000 observations.

The variable `BMI` contains BMI information for the study participants. The variable `Education` records the highest level of education obtained: 8th grade, 9th - 11th grade, high school, some college, or college degree.

4. Load the data. Create a plot that shows the association between BMI and educational level. Describe what you see.

```
#load the data
load(url('https://github.com/jbryer/DATA606/blob/master/data/nhanes_samp_adult_500.rda?raw=true'))
#create a plot
```

5. Examine the normality and equal variance assumptions across the groups. Explain why it is advisable to restrict the analysis to participants who have completed at least 9th grade.

¹The dataset was first introduced in Chapter 1, Lab 1 (Introduction to Data).

```
#check assumptions
```

6. Conduct a hypothesis test to address the question of interest. Let $\alpha = 0.05$. Summarize the conclusions.

```
#conduct hypothesis test
```

Chick weights: comparing weight across feed supplements

Chicken farming is a multi-billion dollar industry, and any methods that increase the growth rate of young chicks can reduce consumer costs while increasing company profits. An experiment was conducted to measure and compare the effectiveness of various feed supplements on the growth rate of chicks. Newly hatched chicks were randomly allocated into groups, and each group was given a different feed supplement.

The `chickwts` dataset available in the `datasets` package contains the weight in grams of chicks at six weeks of age. For simplicity, this analysis will be limited to four types of feed supplements: linseed, meatmeal, soybean, and sunflower.

7. Run the following code to load the `chickwts` dataset and subset the data for the four feed supplements of interest.

```
#load the data
library(datasets)
data("chickwts")
#subset the four feed supplements
keep = (chickwts$feed == "linseed" | chickwts$feed == "meatmeal" |
  chickwts$feed == "soybean" | chickwts$feed == "sunflower")
chickwts = chickwts[keep, ]
#eliminate unused levels
chickwts$feed <- droplevels(chickwts$feed)
```

8. Analyze the data and report the results. Using language accessible to a non-statistician, discuss which feed supplement(s) is/are the most effective for increasing chick weight.

```
#check assumptions
#conduct analysis
```