

Analysis of Variance (ANOVA)

XDASI Fall 2021

November 11, 2021

Contents

Suggested reading	1
Introduction	1
Assumptions	2
Partitioning the variation	2
Mean sums of squares	3
The χ^2 distribution	4
χ^2 with different sample sizes	4
The F -statistic	5
F-distribution for samples from the same population	5
F-distribution for populations with different means	5
<p>p-values</p>	6

Suggested reading

- Whitlock & Schluter, Ch 15 and online lab
- Online tutorials
 - Antoine Soetewey (UCLouvain, Belgium) - Stats and R blog: ANOVA (2020-10-02)
 - (Steven Doogue, 2019-07-09) - Chapter 7.1: One-way ANOVA

Introduction

We've previously used t -tests and non-parametric methods to compare two samples. What if we need to compare more than two samples? The problem with just performing multiple t -tests is that each test has a certain Type I error rate, and when we perform multiple tests this error simply compounds. This issue can be addressed by ***correction for multiple hypothesis testing***, which we will cover later.

If doing all pairwise t -tests between groups is not a good approach, what can we do instead? ANOVA to the rescue! ANOVA extends t -tests to more than two groups by allowing the comparison of the means between multiple groups. However, instead of simply comparing the groups using the difference in their means and their SD, ANOVA compares the ***overall variation between groups*** to the ***variation within each group***. If the overall variation is significantly greater than the individual variation, then we consider the groups to be different.

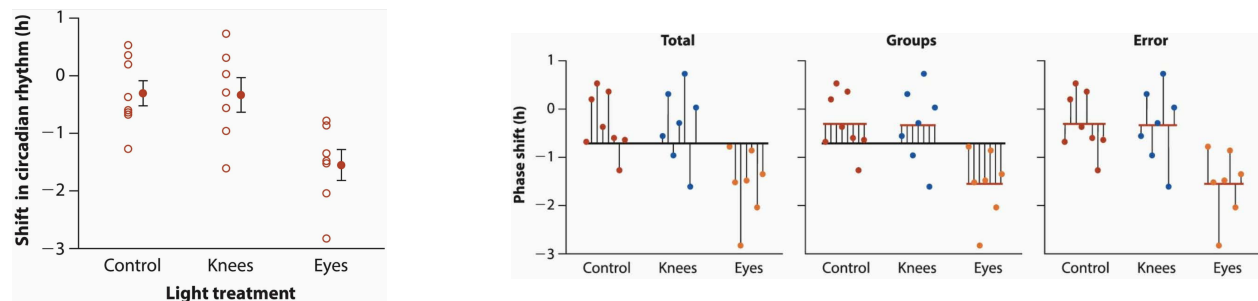
The simplest form of ANOVA is ***one-way ANOVA***, which we will discuss here.

Assumptions

Like t -tests, ANOVA assumes that the distributions of the samples are relatively normal, and also that their variances are similar. When these do not hold, a non-parametric analog ANOVA can be performed called the ***Kruskal-Wallis test***.

Partitioning the variation

The basic idea behind ANOVA is illustrated in Figures 15.1-1 and 15.1-2:



To determine whether there is a significant difference between groups, the variation is decomposed into three parts:

- Total variation with respect to the grand mean
- Variation between group means and the grand mean
- Variation within each group w.r.t. its group mean

These differences are calculated using the ***sum of squares*** of the difference between each data point and the mean used for comparison:

- Total sum of squares (SST): all data points vs. the grand mean

$$SST = \sum_{i=1}^n \sum_{j=1}^m (x_{ij} - \bar{X})^2$$

- Between-groups sum of squares (SSB): the means of the groups vs. the grand mean (note that we multiply each term by the number of points in each group)

$$SSB = \sum_{j=1}^m n_j (\bar{X}_j - \bar{X})^2$$

- Within-groups sum of squares (SSW):

$$SSW = \sum_{i=1}^n \sum_{j=1}^m (x_{ij} - \bar{X}_j)^2$$

where n is the number of data points in each group and m is the total number of groups.

The **total sum of squares (SST)** equals the **within-group SS (SSW)**, also called the **error SS (SSE)**, plus the **between-group SS (SSB or SSG)**:

$$SST = SSW + SSB$$

Notations differ between different sources. Table 10.1 from Ken Aho's book summarizes the equations involved as follows:

TABLE 10.1

General Form of One-Way ANOVA; α_i Represents the True i th Factor-Level Effect in Factor A

Variation Source	df	SS	MS	E(MS)	F*
A (among groups)	$a - 1$	$SS_A = \sum_{i=1}^a n_i (\bar{Y}_i - \bar{Y})^2$	$MS_A = \frac{SS_A}{a - 1}$	$\sigma^2 + \sum_{i=1}^a n_i \frac{\alpha_i^2}{a - 1}$	$\frac{MS_A}{MSE}$
Error (within groups)	$n - a$	$SSE = \sum_{i=1}^a \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_i)^2$	$MSE = \frac{SSE}{n - a}$	σ^2	
Total	$n - 1$	$SSTO = \sum_{i=1}^a \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y})^2$			

Mean sums of squares

To compare between-group and within-group differences, we will compare their **mean sums of squares**.

As usual, we need to take into account the **degrees of freedom**, which is just the number of data points in each equation minus one (the mean used for comparison).

$$MSB = \frac{SSB}{m - 1}$$

$$MSW = \frac{SSW}{N - m}$$

F-statistic

If there is no difference between the groups, we would expect the total variation between the groups to be the same as the variation within the groups, so the ratio of these should be 1. When this is not true, then we know our samples are different, and the ratio should be >1 . This ratio is called the F -statistic:

$$F = \frac{MSB}{MSW}$$

Below we will see that the F -statistic follows an expected distribution under H_o , so we can use this to find a p -value for the observed differences in the mean sums of squares.

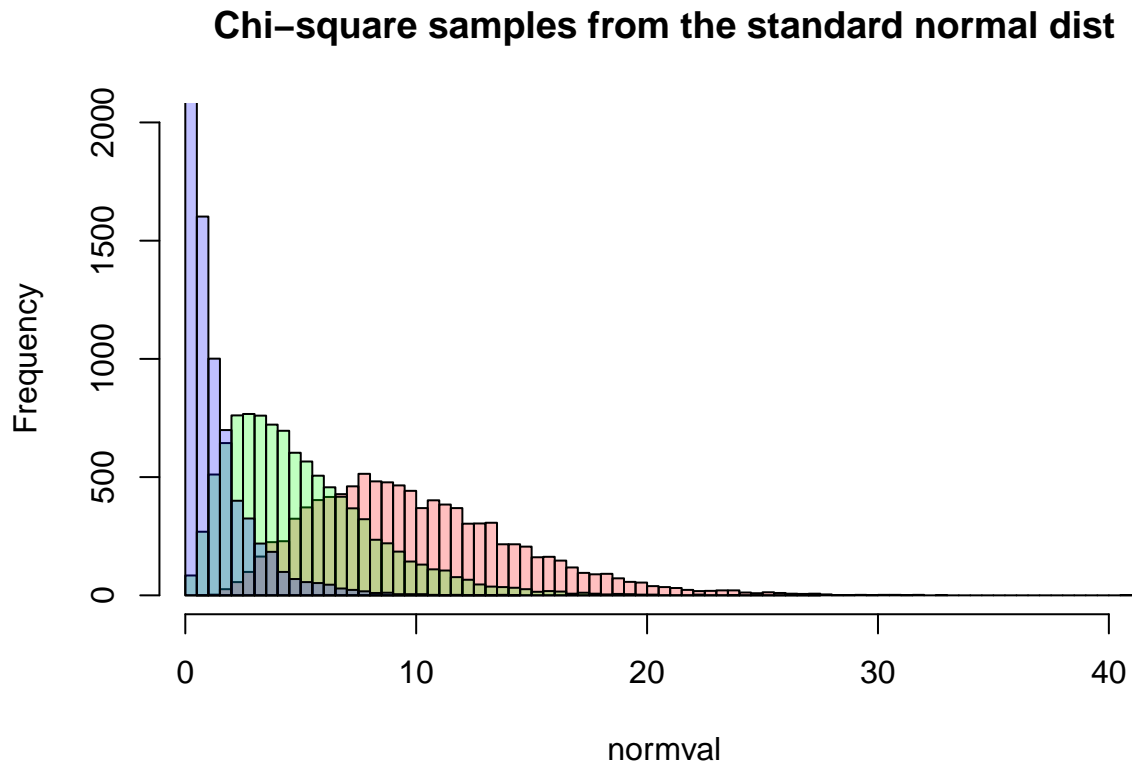
The χ^2 distribution

Since we are measuring differences using **sums of squares**, the differences will follow a χ^2 distribution, which represents the **sum of squared random values** selected from a **normal distribution**. The degrees of freedom, k , is simply the number of random values.

$$Q = \sum_{i=1}^k Z_i^2$$

χ^2 with different sample sizes

Let's simulate some data to see what it looks like when $k = 10, 5$, and 1 and the values are retrieved from a standard normal distribution. (We have also done this in a previous class.)



It is clear to see that as k increases, the distribution begins to look like a **normal distribution**.

This only happens when the **sample size is at least 10**, which is why it is not recommended to use the χ^2 test for small values of k (<10).

The F -statistic

To *compare two χ^2 distributions*, we can simply take a *ratio* of them (taking into account their respective degrees of freedom).

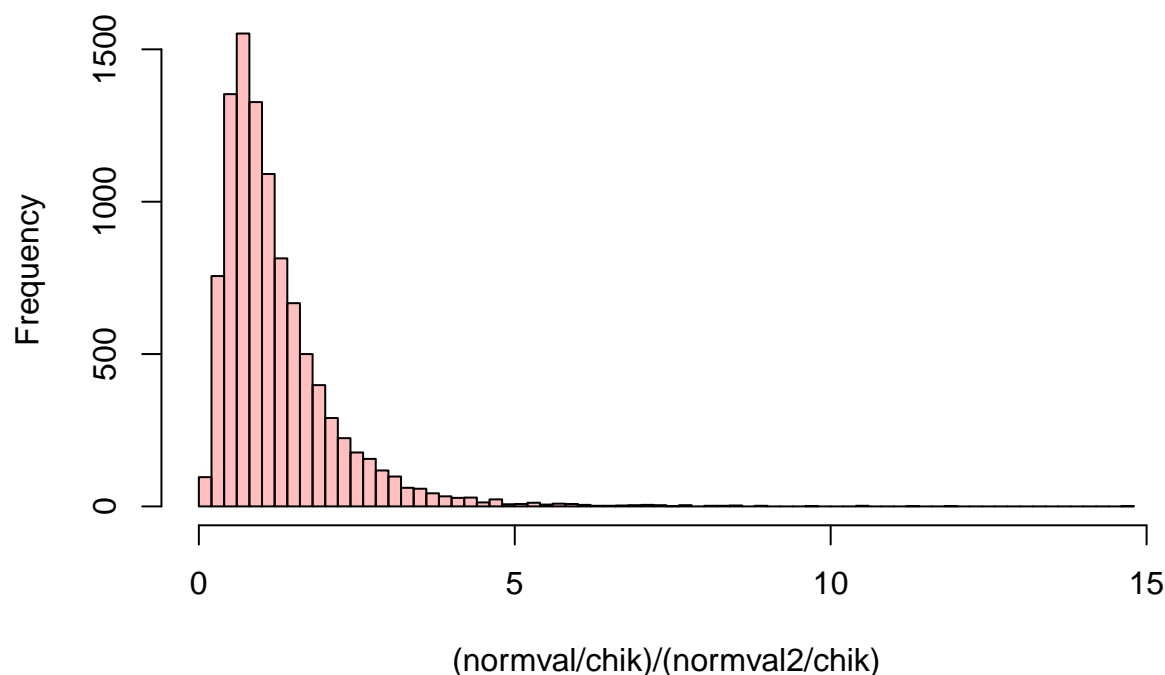
This distribution is called an F -**distribution** and the *observed ratio* is the F -*statistic*. (The F -distribution is named after the statistician Ron Fisher.)

F-distribution for samples from the same population

If we take two random samples of the same size from a normal distribution, square them, and then take the ratio of the mean sums of squares (taking into account the degrees of freedom), we will get an F -distribution.

Let's first see what it would look like if the means of the two populations that are being sampled are equal.

Histogram of F -distribution



Note that since we are taking random samples, the histogram will vary slightly each time we take new samples.

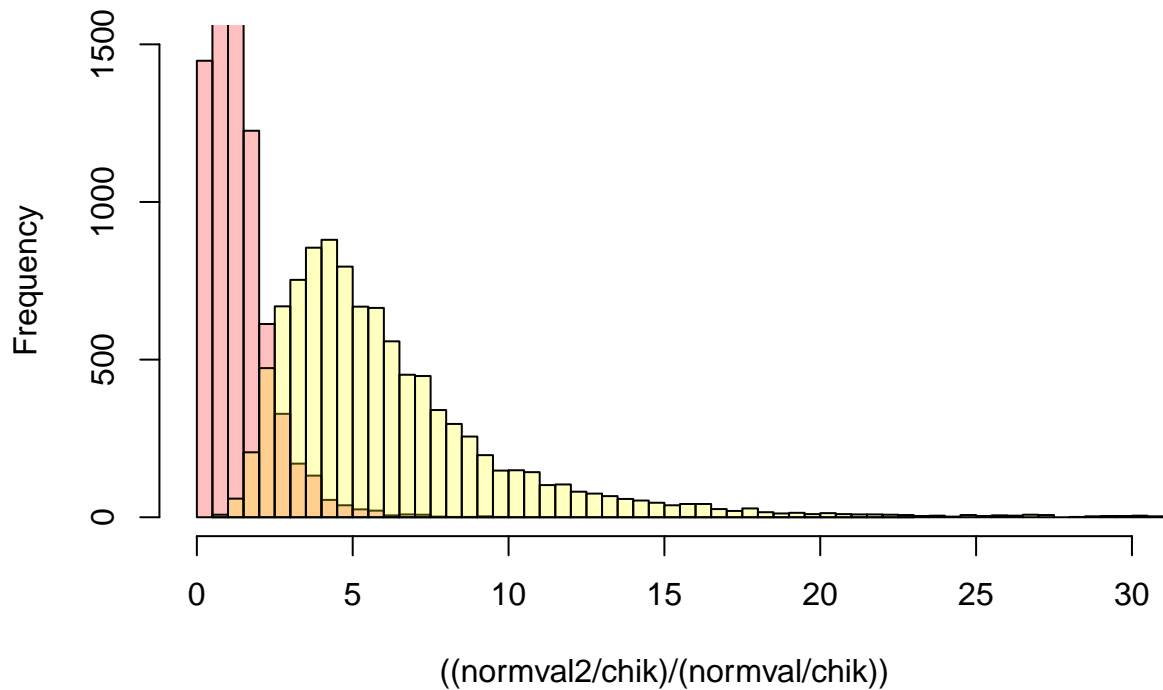
F-distribution for populations with different means

What if the means of our normal distributions are different?

We can make a second histogram showing the same ratio for data sampled from two normal distributions with different means: the standard normal and a normal distribution with mean = 2 and sd = 1.

Now, the ratio of the sums of the two samples will look quite different. Let's try this and superimpose the two histograms for comparison.

Histogram of F-distribution



Remember that *variance* is essentially a sum of squares as discussed above. So now we have the ability to compare two different variances and use a statistic to determine if they are significantly different.

p-values

We can use the value of F to find a p -value using the F -distribution for a particular sample size, given the degrees of freedom for each sample:

```
Fstat = sum(rnorm(chik, mean=2)**2)/sum(rnorm(chik)**2)
Fstat
```

```
## [1] 9.44284
```

```
pf(Fstat, chik, chik, lower.tail=F)
```

```
## [1] 0.0007284056
```