

One-way analysis of variance

Date: 2017-11-21

What it is

One-way analysis of variance, or ANOVA, compares means across 2 or more samples. The concept is similar to that of the independent samples t test, but generalizes that procedure to any number of groups.

When to use it

ANOVA can be used to compare means across two or more samples provided the following conditions are true:

- The samples are mutually independent (not paired or otherwise correlated with one another)
- Each sample is either:
 - Approximately normally distributed, OR
 - Large (using the rule of thumb $n \geq 30$)
- The population variances are equal across all groups¹

How to use it

An ANOVA test begins with hypotheses. If K is the number of groups, the hypotheses are:

H_0 : All population means are equal ($\mu_1 = \mu_2 = \dots = \mu_K$).

H_A : Not all the population means are equal.

Note that the alternative hypothesis is quite vague. There are a lot of ways that “not all” the means can be equal.² In the event that the null hypothesis is eventually rejected, a further analysis should be done to determine where any differences exist.

Suppose we want to compare average BMI across racial/ethnic groups. Then the hypotheses would be:

H_0 : All racial/ethnic groups have the same average BMI ($\mu_{white} = \mu_{black} = \mu_{hispanic} = \mu_{other}$). H_A : Not all racial/ethnic groups have the same average BMI.

Once the hypotheses are determined, there are a number of ways to run a one-way ANOVA in R. For a quick-and-dirty ANOVA, use the `oneway.test()` function. Provide the function with a formula structured as `outcome ~ groups`, the data set, and set `var.equal=TRUE` (assuming the equal variance rule of thumb is met). The result will

¹ A simple rule of thumb is to check the ratio of the largest sample standard deviation to the smallest. If $\frac{s_{max}}{s_{min}} \leq 2$ then this condition can be said to be reasonably well met.

² For example, μ_1 could be different from all the other means, or μ_2 could be different from all others. Or maybe μ_1 and μ_2 are the same, but different from the other means. Or maybe all K means are unique. The possibilities here are numerous...

include an test statistic (F), numerator and denominator degrees of freedom, and a p-value:

```
oneway.test(bmi ~ race, data = d, var.equal = TRUE)

##
## One-way analysis of means
##
## data:  bmi and race
## F = 7.9616, num df = 3, denom df = 496, p-value = 3.412e-05
```

For more detail, you first need to create a linear model with the `lm()` function, then run `anova()` on the result. Within the `lm()` function, You should use the same formula and data setup as in `oneway.test()`. The `anova()` results will include the same information as `oneway.test()`, but with a full ANOVA table, complete with sums of squares and mean squares:

```
bmi_model = lm(bmi ~ race, data = d)
anova(bmi_model)

## Analysis of Variance Table
##
## Response: bmi
##           Df Sum Sq Mean Sq F value    Pr(>F)
## race        3   985.8   328.61   7.9616 3.412e-05 ***
## Residuals 496 20472.5    41.28
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Why it works

The ANOVA procedure is highly computational, so we will stick with conceptual understanding here. One-way ANOVA compares the differences between groups (denoted by the race row in the table above) to the differences within groups (denoted by the Residuals row).³ If the differences between groups are small, then the F statistic will also be small; conversely, if the differences between groups are large, then the F statistic will be large. Therefore, large values of F provide strong evidence against the null hypothesis.

³ Explicitly, the F statistic is $F = \frac{MSB}{MSE}$, or mean square between groups divided by the mean square within groups.