

12.1: One-way ANOVA

A one-way ANOVA (ANalysis Of VAriance) is a generalization of the independent samples t -test to compare more than 2 groups. (Actually an independent samples t -test and an ANOVA with two groups are the same thing). The hypotheses to be tested, in comparing the mean of k groups, with a one-way ANOVA are :

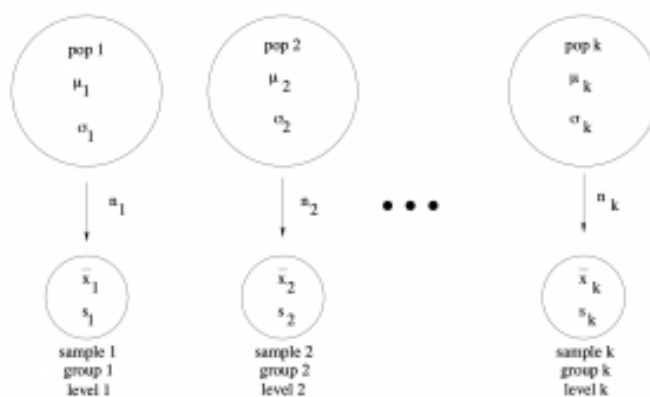
$$H_0 : \mu_1 = \mu_2 = \dots = \mu_k$$

H_1 : At least one of the means is different from the others.

The following assumptions must be met for ANOVA (the version we have here) to be valid :

1. Normally distributed populations (although ANOVA is robust to violations of this condition).
2. Independent samples (between subjects).
3. Homoscedasticity : $\sigma_1^2 = \sigma_2^2 = \dots = \sigma_k^2$. (ANOVA is robust to violations of this too, especially for larger sample sizes.)

The concept of ANOVA is simple but we need to learn some terminology so we can understand how other people talk about ANOVA. Each sample set from each population is referred to as a *group* or each population is called a group.



There will be k groups with sample sizes n_1, n_2, \dots, n_k with the total number of data points being $N = \sum_{i=1}^k n_i$. For an ANOVA, the concept of independent variable (IV) and dependent variable (DV) become important (the IV in a single sample or a paired t -test is trivially a number like k or 0). The groups comprise different values of one IV. The IV is discrete with k values or *levels*.

In raw form, the test statistic for a one-way ANOVA is

$$F_{\text{test}} = F_{\nu_1, \nu_2} = \frac{s_B^2}{s_W^2} \quad (12.1.1)$$

where

$$\nu_1 = k - 1 \text{ (d.f.N.)} \quad \nu_2 = N - k \text{ (d.f.D.)} \quad (12.1.2)$$

are the degrees of freedom you use when looking up F_{crit} in the **F Distribution Table** and where

$$s_B^2 = \frac{\sum_{i=1}^k n_i (\bar{x}_i - \bar{x}_{\text{GM}})^2}{k - 1} \quad (12.1.3)$$

is the variance between groups, and

$$s_W^2 = \frac{\sum_{i=1}^k (n_i - 1) s_i^2}{\sum_{i=1}^k (n_i - 1)} \quad (12.1.4)$$

is the variance within groups. Here n_i , \bar{x}_i and s_i are the sample size, mean and standard deviation for sample i and \bar{x}_{GM} is the grand mean:

$$\bar{x}_{\text{GM}} = \frac{\sum_{i=1}^k \sum_{j=1}^{n_i} x_{ij}}{N} = \frac{\sum_{i=1}^k n_i \bar{x}_i}{N} \quad (12.1.5)$$

where x_{ij} is data point j in group i .

So you can see that ANOVA, the analysis of variance, is about comparing two variances. The within variance s_W^2 is the variance of all the data lumped together, just as the grand mean \bar{x}_{GM} is the mean of all the data lumped together. It is the noise. You can see that the within variance is the weighted mean (weighted by $n_i - 1$) of the group sample variances — a little algebra shows that this is the variance of all the data lumped together. The between variance s_B^2 a variance of the sample means \bar{x}_i . It is the signal. If the sample means were all exactly the same then the between variance s_B^2 would be zero. So the higher F_{test} the more likely the means are different. F_{test} is a signal-to-noise ratio. If the means were all the same in the population then s_B^2 would follow a χ_{k-1}^2 distribution and s_W^2 (whether the population means were the same or not) would follow a χ_{N-k}^2 distribution. Thus if the population means were all the same (H_0) then the F test statistic follows a F_{ν_1, ν_2} distribution which has an expected value^[1] (mean) of about 1. F_{test} must be sufficiently bigger than 1 to reject H_0 .

The analysis of the variances can be broken down further, to sums of squares, with the following definitions^[2]:

$$s_B^2 = \text{MS}_B \quad \text{between groups mean square}$$

and

$$s_W^2 = \text{MS}_W \quad \text{within groups mean square.}$$

Next we note that $\nu_1 = k - 1$ and $\nu_2 = N - k = \sum_{i=1}^k (n_i - 1)$ so

$$\text{MS}_B = \frac{\text{SS}_B}{\nu_1}$$

and

$$\text{MS}_W = \frac{\text{SS}_W}{\nu_2}$$

where

$$\text{SS}_B = \sum_{i=1}^k n_i (\bar{x}_i - \bar{x}_{\text{GM}})^2 \quad \text{sum of squares between groups}$$

and

$$\text{SS}_W = \sum_{i=1}^k (n_i - 1) s_i^2 \quad \text{sum of squares within groups}$$

so that

$$F_{\text{test}} = \frac{\text{MS}_B}{\text{MS}_W}$$

Why are sums of squares so prominent in statistics? (They will show up in linear regression too.) Because squares are the essence of variance. Look at the formula for the normal distribution, [Equation 5.1](#). The exponent is a square. Mean and variance are all you need to completely characterize a normal distribution. Means are easy to understand, so sums of square focus our attention to the variance of normal distributions. If we make an assumption that all random noise has a normal distribution (which can be justified on general principles) then the sums of squares will tell the whole statistical story. Sums of squares also tightly links statistics to linear algebra (see [Chapter 17](#)) because the Pythagorus Theorem, which gives distances in ordinary geometrical spaces, is about sums of squares.

Computer programs, like SPSS, will output an ANOVA table that breaks down all the sums of squares and other pieces of the F test statistic :

Source	SS	ν	MS	F	p (sig)
Between (signal)	SS_B	$\nu_1 = k - 1$	$\text{MS}_B = \text{SS}_B / \nu_1$	Rendered by QuickLaTeX.com	p
Within (error)	SS_W	$\nu_2 = N - k$	$\text{MS}_W = \text{SS}_W / \nu_2$		
Totals	SS_T	$\nu_T = N - 1$			

Here p is the p -value of F_{test} , reported by SPSS as “sig” for significance. F_{test} is significant (you can reject H_0) if $p < \alpha$. You should be able to reconstruct an ANOVA table given only the SS values. Notice that the total degrees of freedom of the ANOVA is $\nu_T = \nu_1 + \nu_2 = N - 1$. One degree of freedom is used up in computing the grand mean, the rest in computing the variances, very similar to how $n - 1$ is the degrees of freedom for sample standard deviation s . If you think of degrees of freedom as the amount of information in the data then the one-way ANOVA uses up all the information in the data. This point will come up again when we consider post hoc comparisons.

Example 12.1 : A state employee wishes to see if there is a significant difference in the number of employees at the interchanges of three state toll roads. At $\alpha = 0.05$ is there a difference in the average number of employees at each interchange between the toll roads?

The data are :

Road 1 (group 1)	Road 2 (group 2)	Road 3 (group 3)
7	10	1
14	1	12
32	1	1
19	0	9
10	11	1
11	1	11

Solution :

0. Data reduction.

Using your calculators, find

$$n_1 = 6 \quad \bar{x}_1 = 15.5 \quad s_1^2 = 81.9$$

$$n_2 = 6 \quad \bar{x}_2 = 4.0 \quad s_2^2 = 25.6$$

$$n_3 = 6 \quad \bar{x}_3 = 5.83 \quad s_3^2 = 29.0$$

$N = 18$.

1. Hypothesis.

$$H_0 : \mu_1 = \mu_2 = \mu_3$$

H_1 : At least one of the means is different from the others.

2. Critical statistic.

Use the **F Distribution Table** with $\alpha = 0.05$; do not divide the table α (right tail area) by 2 in this case, there are no left and right tail tests in ANOVA. The degrees of freedom needed are $\nu_1 = k - 1 = 3 - 1 = 2$ (d.f.N.) and $\nu_2 = N - k = 18 - 3 = 15$ (d.f.D.). With that information

$$F_{\text{crit}} = 3.68 \quad (12.1.6)$$

3. Test statistic.

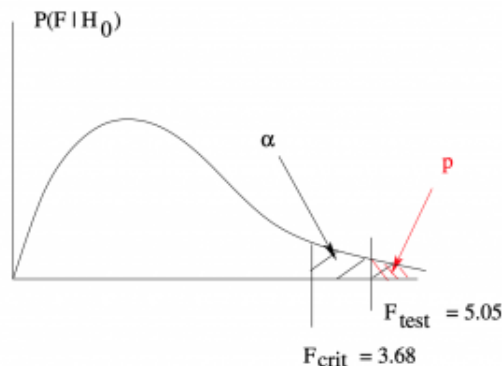
Compute, in turn :

$$\begin{aligned} \bar{x} &= \frac{n_1 \bar{x}_1 + n_2 \bar{x}_2 + n_3 \bar{x}_3}{N} = \frac{(6)(15.5) + (6)(4.0) + (6)(5.83)}{18} = 8.4 \\ s_B^2 &= \frac{\sum_{i=1}^k n_i (\bar{x}_i - \bar{x})^2}{k-1} = \frac{(6)(15.5 - 8.4)^2 + (6)(4.0 - 8.4)^2 + (6)(5.83 - 8.4)^2}{3-1} \\ &= \frac{459.18}{2} = 229.59 \\ s_W^2 &= \frac{\sum_{i=1}^k (n_i - 1) s_i^2}{\sum_{i=1}^k (n_i - 1)} = \frac{(6-1)(81.9) + (6-1)(25.6) + (6-1)(29.0)}{(18-3)} \\ &= \frac{682.5}{15} = 45.5 \end{aligned}$$

Note how we saved SS_B and SS_W for the ANOVA table. And finally

$$F_{\text{test}} = \frac{s_B^2}{s_W^2} = \frac{229.59}{45.5} = 5.05 \quad (12.1.7)$$

4. Decision.



Reject H_0 .

5. Interpretation.

Using one-way ANOVA at $\alpha = 0.05$ we found that at least one of the toll roads has a different average number of employees at their interchanges. The ANOVA table is :

Source	SS	ν	MS	F	p (sig)
Between (signal)	459.18	2	229.59	5.05	$p < 0.05$
Within (error)	682.5	15	45.5		
Totals	1141.68	17			

We did not compute p but a computer program like SPSS will.

□

1. The mean of the F_{ν_1, ν_2} distribution is $\mu_F = \frac{\nu_2}{\nu_2 - 2}$ if $\nu_2 > 2$. ↩
2. You might have heard of RMS for "root mean square". $\text{RMS} = \sqrt{\text{MS}} = \sqrt{s^2}$. RMS is standard deviation. ↩

This page titled [12.1: One-way ANOVA](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Gordon E. Sarty](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.