

You have **2** free member-only stories left this month.

[Sign up for Medium and get an extra one](#)



Wei-  
Meng  
Lee



Oct 18,  
2021

·  
10 min  
read



 [Listen](#)

# Statistics in Python — Using ANOVA for Feature Selection

Understand how to use  
ANOVA for comparing  
between a categorical and  
numerical variable



Photo by [Gabriel Gurrola](#) on [Unsplash](#)

In my previous article, I talked about using the chi-square statistics to select features from a dataset for machine learning. The chi-square test is used when both your independent and dependent variables are all *categorical* variables. However, what if your independent variable is *categorical* and your dependent variable is *numerical*? In this case, you have to use another statistic test known as ANOVA — **Analysis of Variance**.

And so in this article, our discussion will revolve around ANOVA and how you use it in machine learning for feature selection. Like all my previous articles, I will use a concrete example to explain the concept.

Before we get started, it is useful to summarize the different methods for feature selection that we have discussed so far :

Variable 1	Variable 2	Method to test for correlation (dependence)
Numerical	Numerical	Pearson correlation
Categorical (Ordinal)	Numerical	Spearman's rank correlation
Categorical (Nominal)	Categorical (Nominal)	Chi-Square
Categorical (Ordinal/Nominal)	Numerical	ANOVA

Image by author

If you need a refresher on **Pearson correlation**, **Spearman's rank correlation**, and **Chi-Square**, I suggest you go and check them out now (see the links below) and come back to this article once you are done. Some of the concepts discussed in this article is similar to that of the chi-square test, and so I recommend you check that out.

### Statistics in Pyth...

In my previous...

towardsdata  
science.com

## Statistics in Pyth...

Understa  
nd how t...

towardsdata  
science.com

## Statistics in Pyth...

Understa  
nd the...

towardsdata  
science.com

## What is ANOVA?

ANOVA is used for testing two variables, where:

- one is a *categorical* variable
- another is a *numerical* variable

ANOVA is used when the categorical variable has *at least 3 groups* (i.e three different unique values).

*If you want to compare just two groups, use the t-test. I will cover t-test in another article.*

ANOVA lets you know if your numerical variable changes according to the level of the categorical variable.

*ANOVA uses the  $f$ -tests to statistically test the equality of means.  $F$ -tests are named after its test statistic,  $F$ , which was named in honor of **Sir Ronald Fisher**.*

Here are some examples that makes it easier to understand when you can use ANOVA.

- You have a dataset containing information of a group of people pertaining to their social media usage and the number of hours they sleep:

Social Media Usage	Hours of sleep
Low	8
Medium	7
High	6

Image by author

You want to find out if the amount of social media usage

(categorical variable) has a direct impact on the number of hours of sleep (numerical variable).

- You have a dataset containing three different brands of medication and the number of days for the medication to take effect:

Brand	Number of Days to Take Effect
BrandX	1
BrandY	3
BrandZ	2

Image by author

You want to find out if there is a direct relationship between a specific brand and its effectiveness.

*ANOVA checks whether there is equal variance between groups of categorical feature with respect to the numerical response.*

If there is equal variance between groups, it means this feature has no impact on the response and hence it (the

categorical variable) cannot be considered for model training.

## Performing ANOVA by hand

The best way to understand ANOVA is to use an example. In the following example, I use a fictitious dataset where I recorded the reaction time of a group of people when they are given a specific type of drink.

### Sample Dataset

I have a sample dataset named **drinks.csv** containing the following content:

```
team,drink_type,reaction_
time
1,water,14
2,water,25
3,water,23
4,water,27
5,water,28
6,water,21
7,water,26
8,water,30
9,water,31
10,water,34
1,coke,25
2,coke,26
3,coke,27
4,coke,29
5,coke,25
6,coke,23
```

```
7, coke, 22
8, coke, 27
9, coke, 29
10, coke, 21
1, coffee, 8
2, coffee, 20
3, coffee, 26
4, coffee, 36
5, coffee, 39
6, coffee, 23
7, coffee, 25
8, coffee, 28
9, coffee, 27
10, coffee, 25
```

There are 10 teams in all — each team comprises of 3 persons. Each person in the team is given three different types of drinks — water, coke, and coffee. After consuming the drink, they were asked to perform some activities and their reaction time recorded. The aim of this experiment is to determine if the drinks have any effect on a person's reaction time.

Let's first load the dataset into a Pandas DataFrame:

```
import pandas as pd
df =
pd.read_csv('drinks.csv')
```



Record the *observation size*,  
which we will make use of  
later:

```
observation_size =  
df.shape[0] # number of  
observations
```

**observation\_size**

	team	drink_type	reaction_time
0	1	water	14
1	2	water	25
2	3	water	23
3	4	water	27
4	5	water	28
5	6	water	21
6	7	water	26
7	8	water	30
8	9	water	31
9	10	water	34
10	1	coke	25
11	2	coke	26
12	3	coke	27
13	4	coke	29
14	5	coke	25
15	6	coke	23
16	7	coke	22
17	8	coke	27
18	9	coke	29
19	10	coke	21
20	1	coffee	8
21	2	coffee	20
22	3	coffee	26
23	4	coffee	36
24	5	coffee	39
25	6	coffee	23
26	7	coffee	25
27	8	coffee	28
28	9	coffee	27
29	10	coffee	25

Image by author

## Visualizing the dataset

It is useful to visualize the  
distribution of the data using a

Boxplot:

```
df =  
df.boxplot('reaction_time'  
, by='drink_type')
```

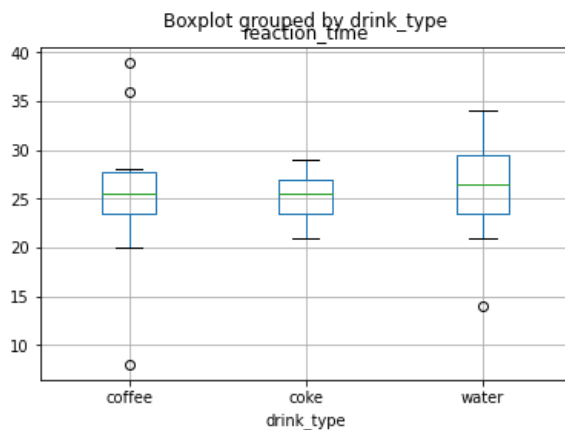


Image by author

You can see that the three types of drinks have about the same median reaction time.

## Pivoting the dataframe

To facilitate the calculation for ANOVA, we need to pivot the dataframe:

```
df =  
df.pivot(columns='drink_t'  
ype', index='team')  
display(df)
```

reaction_time			
drink_type	coffee	coke	water
team			
1	8	25	14
2	20	26	25
3	26	27	23
4	36	29	27
5	39	25	28
6	23	23	21
7	25	22	26
8	28	27	30
9	27	29	31
10	25	21	34

Image by author

The columns represent the three different types of drinks and the rows represents the 10 teams. We will also use this chance to record the *number of items in each group*, as well as the *number of groups*, which we will make use of later:

```
n = df.shape[0]    # 10;
number of items in each
group
k = df.shape[1]    # 3;
number of groups
```

n	Number of items in each group	reaction_time			
		drink_type	coffee	coke	water
		team			
		1	8	25	14
		2	20	26	25
		3	26	27	23
		4	36	29	27
		5	39	25	28
		6	23	23	21
		7	25	22	26
		8	28	27	30
		9	27	29	31
		10	25	21	34
			</		

## Defining the Hypotheses

- **H<sub>0</sub>** (Null hypothesis) — that there is no difference among group means.
- **H<sub>1</sub>** (Alternate hypothesis) — that at least one group differs significantly from the overall mean of the dependent variable.

## Step 1 — Calculating the means for all groups

We are now ready to begin our calculations for ANOVA. First, let's find the mean for each group:

```
df.loc['Group Means'] =  
df.mean()  
df
```

drink_type	reaction_time		
	coffee	coke	water
team			
1	8.0	25.0	14.0
2	20.0	26.0	25.0
3	26.0	27.0	23.0
4	36.0	29.0	27.0
5	39.0	25.0	28.0
6	23.0	23.0	21.0
7	25.0	22.0	26.0
8	28.0	27.0	30.0
9	27.0	29.0	31.0
10	25.0	21.0	34.0
Group Means	25.7	25.4	25.9

Image by author

From here, you can now calculate the **overall mean**:

Get started

[Sign In](#)



**Wei-Meng Lee**

489 Followers

ACLP Certified Trainer |  
Blockchain, Smart Contract, Data  
Analytics, Machine Learning, Deep  
Learning, and all things tech  
(<http://calendar.learn2develop.net>).

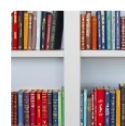
Follow

### More from Medium



Wei-... in Toward...

**Ensemble  
Learning in  
sklearn**



Frank Ne... in C...





```
overall_mean =  
df.iloc[-1].mean()  
overall_mean #  
25.666666666666668
```

7	25.0	22.0	26.0
8	28.0	27.0	30.0
9	27.0	29.0	31.0
10	25.0	21.0	34.0
Group Means	25.7	25.4	25.9

$$\frac{25.7 + 25.4 + 25.9}{3} = 25.667$$

Image by author

## Step 2 — Calculate the Sum of Squares

Now that we have calculated the *overall mean*, we can proceed to calculate the following:

- Sum of squares of all observation — **SS\_total**
- Sum of squares within — **SS\_within**
- Sum of squares between — **SS\_between**

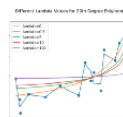
**Sum of squares of all observation — SS\_total**

## How to Set X and y in Pandas



Bran... in Toward...

## Complete Guide to Regression Analysis Using...



Ruks... in Towar...

## Encoding Categorical Variables: One-h...

one-hot encoding	d1
	1
	0
	0

[Help](#) [Status](#) [Writers](#) [Blog](#) [Careers](#)  
[Privacy](#) [Terms](#) [About](#) [Knowable](#)

The *sum of squares of all observation* is calculated by deducting each observation from the *overall mean*, and then summing all the squares of the differences:

$$SS_{total} = \Sigma(X_{ij} - \bar{X})^2$$

Image by author

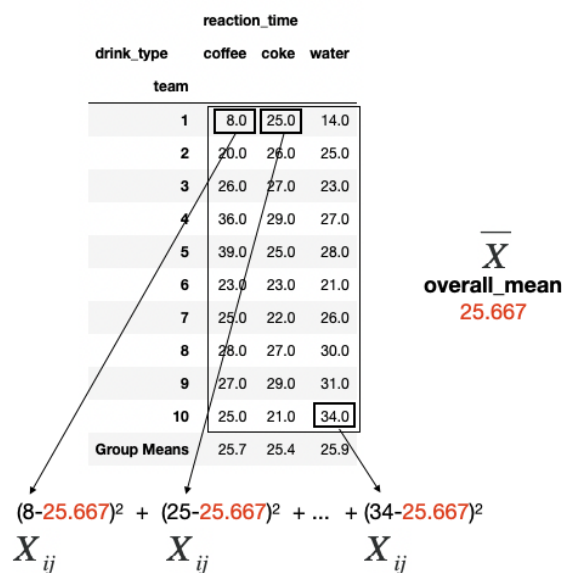


Image by author

Programmatically, **SS\_total** is computed as:

```
SS_total =
(((df.iloc[:-1] -
overall_mean)**2).sum()).
sum()
SS_total    #
```

## Sum of squares within — **SS\_within**

The *sum of squares within* is the sum of squared deviations of scores around their group's mean:

		reaction_time		
		coffee	coke	water
team				
1	8.0	25.0	14.0	
2	20.0	26.0	25.0	
3	26.0	27.0	23.0	
4	36.0	29.0	27.0	
5	39.0	25.0	28.0	
6	23.0	23.0	21.0	
7	25.0	22.0	26.0	
8	28.0	27.0	30.0	
9	27.0	29.0	31.0	
10	25.0	21.0	34.0	
Group Means		25.7	25.4	25.9

$(8-25.7)^2 + (20-25.7)^2 + \dots + (25-25.7)^2$

$(25-25.4)^2 + (26-25.4)^2 + \dots + (21-25.4)^2$

$(14-25.9)^2 + (25-25.9)^2 + \dots + (34-25.9)^2$

Image by author

Programmatically, **SS\_within** is computed as:

```
SS_within =
((df.iloc[:-1] -
df.iloc[-1])**2).sum()).s
um()
```



## Sum of Squares between — SS\_between

Next we calculate the sum of squares of the group means from the overall mean:

$$SS_{between} = n \sum (X_j - \bar{X})^2$$

Image by author

reaction_time			
drink_type	coffee	coke	water
team			
1	8.0	25.0	14.0
2	20.0	26.0	25.0
3	26.0	27.0	23.0
4	36.0	29.0	27.0
5	39.0	25.0	28.0
6	23.0	23.0	21.0
7	25.0	22.0	26.0
8	28.0	27.0	30.0
9	27.0	29.0	31.0
10	25.0	21.0	34.0
Group Means	25.7	25.4	25.9

$\bar{X}$   
**overall\_mean**  
 25.667

$10(25.7 - 25.667)^2 + 10(25.4 - 25.667)^2 + 10(25.9 - 25.667)^2$   
 $X_j \quad X_j \quad X_j$

Image by author

Programmatically, **SS\_between** is computed as:

```
SS_between = (n *
(df.iloc[-1] -
overall_mean)**2).sum()
SS_between #
```

1.2666666666666667

You can verify that:

$$SS_{total} = SS_{between} + SS_{within}$$

## Creating the ANOVA Table

With all the values computed, you can now complete the ANOVA table. Recall you have the following variables:

Variable	Description
observation_size	Number of observations
n	Number of items in each group (teams)
k	Number of groups (e.g. water, coke, coffee)

Image by author

You can compute the various *degrees of freedoms* as follows:

```
df_total =  
observation_size - 1  
# 29  
df_within =  
observation_size - k  
# 27  
df_between = k - 1  
# 2
```

From the above, compute the various *mean squared* values:

```
mean_sq_between =  
SS_between / (k - 1)  
# 0.6333333333333335  
mean_sq_within = \  
    SS_within /  
    (observation_size - k)  
# 37.08888888888889
```

Finally, you can calculate the **F-value**, which is the ratio of two variances:

```
F = mean_sq_between /  
    mean_sq_within      #  
0.017076093469143204
```

*Recall earlier that I mentioned ANOVA uses the f-tests to statistically test the equality of means.*

Once the F-value is obtained, you now have to refer to the *f-distribution table* (see [http://www.socr.ucla.edu/Applets.dir/F\\_Table.html](http://www.socr.ucla.edu/Applets.dir/F_Table.html) for one example) to obtain the **f-critical value**. The f-

distribution table is organized based on the  $\alpha$  value (usually 0.05). So you need to first locate the table based on  $\alpha=0.05$ :

F Table for $\alpha = 0.05$																			
$\alpha$	$df_1$	2	3	4	5	6	7	8	9	10	12	15	20	24	30	40	60	120	$\infty$
0.05	2	18.5128	19.0000	19.1643	19.2468	19.2964	19.3259	19.3512	19.3710	19.3848	19.3959	19.4125	19.4251	19.4358	19.4451	19.4524	19.4578	19.4617	19.4657
0.05	3	10.1280	10.5521	9.7666	9.1172	8.9135	8.6867	8.5452	8.4125	8.2855	8.1629	8.0299	7.8962	7.7625	7.6288	7.4951	7.3614	7.2277	7.0940
0.05	4	7.7086	8.0442	6.9112	6.2582	6.0503	5.8242	5.6817	5.5489	5.4162	5.2835	5.1508	5.0181	4.8854	4.7527	4.6200	4.4873	4.3546	4.2219
0.05	5	6.5979	6.8703	5.9059	5.2412	5.0263	4.7902	4.6477	4.5149	4.3822	4.2495	4.1168	3.9841	3.8514	3.7187	3.5860	3.4533	3.3206	3.1879
0.05	6	5.9874	6.2155	5.4059	4.7312	4.5063	4.2602	4.1177	3.9849	3.8522	3.7195	3.5868	3.4541	3.3214	3.1887	3.0560	2.9233	2.7906	2.6579
0.05	7	5.4914	5.6775	5.0559	4.3712	4.1363	3.8802	3.7377	3.6049	3.4722	3.3395	3.2068	3.0741	2.9414	2.8087	2.6760	2.5433	2.4106	2.2779
0.05	8	5.0774	5.2275	4.7039	4.0192	3.7743	3.5182	3.3757	3.2429	3.1102	2.9775	2.8448	2.7121	2.5794	2.4467	2.3140	2.1813	2.0486	1.9159
0.05	9	4.7274	4.8475	4.4239	3.7392	3.4843	3.2282	3.0857	2.9529	2.8202	2.6875	2.5548	2.4221	2.2894	2.1567	2.0240	1.8913	1.7586	1.6259
0.05	10	4.4374	4.5375	4.2139	3.5292	3.2643	2.9982	2.8557	2.7229	2.5902	2.4575	2.3248	2.1921	2.0594	1.9267	1.7940	1.6613	1.5286	1.3959
0.05	12	4.1774	4.2575	4.0339	3.3492	3.0743	2.8082	2.6657	2.5329	2.4002	2.2675	2.1348	2.0021	1.8694	1.7367	1.6040	1.4713	1.3386	1.2059
0.05	15	3.7454	3.8075	3.6839	3.0992	2.8143	2.5482	2.4057	2.2729	2.1402	2.0075	1.8748	1.7421	1.6094	1.4767	1.3440	1.2113	1.0786	0.9459
0.05	20	3.2934	3.3375	3.3139	2.7292	2.4343	2.1682	2.0257	1.8929	1.7602	1.6275	1.4948	1.3621	1.2294	1.0967	0.9640	0.8313	0.6986	0.5659
0.05	24	2.9774	3.0075	3.0039	2.4192	2.1143	1.8482	1.7057	1.5729	1.4402	1.3075	1.1748	1.0421	0.9094	0.7767	0.6440	0.5113	0.3786	0.2459
0.05	30	2.6234	2.6475	2.6539	2.0692	1.7643	1.4982	1.3557	1.2229	1.0902	0.9575	0.8248	0.6921	0.5594	0.4267	0.2940	0.1613	0.0286	-0.1041
0.05	40	2.2474	2.2675	2.2739	1.6892	1.3843	1.1182	0.9757	0.8429	0.7102	0.5775	0.4448	0.3121	0.1794	0.0467	-0.0860	-0.2187	-0.3514	-0.4841
0.05	60	1.8534	1.8675	1.8739	1.2892	0.9843	0.7182	0.5757	0.4429	0.3102	0.1775	0.0448	-0.0879	-0.2206	-0.3533	-0.4860	-0.6187	-0.7514	-0.8841
0.05	120	1.4574	1.4675	1.4739	1.0892	0.7843	0.5182	0.3757	0.2429	0.1102	-0.0225	-0.1548	-0.2871	-0.4194	-0.5517	-0.6840	-0.8163	-0.9486	-1.0809
0.05	$\infty$	1.2854	1.2975	1.3039	0.9192	0.6143	0.3482	0.2057	0.0729	-0.0598	-0.1921	-0.3244	-0.4567	-0.5890	-0.7213	-0.8536	-0.9859	-1.1182	-1.2505

Source:

[http://www.socr.ucla.edu/Applets.dir/F\\_Table.html](http://www.socr.ucla.edu/Applets.dir/F_Table.html)

Next, observe that the columns of the f-distribution table is based on **df1** while the rows are based on **df2**. You can get your **df1** and **df2** from the previous variables that we have created:

```
df1 = df_between    # 2
df2 = df_within     # 27
```

Using the values of **df1** and **df2**, you can now locate the **f-critical value** by locating the **df1** column and **df2** row:



	$df_1=1$	2	3
$df_2=1$	161.4476	199.5000	215.7073
2	18.5128	19.0000	19.1643
3	10.1280	9.5521	9.2766
4	7.7086	6.9443	6.5914
5	6.6079	5.7861	5.4095

26	4.2252	3.3690	2.9752
27	4.2100	<b>3.3541</b>	2.9604
28	4.1960	3.3404	2.9467
29	4.1830	3.3277	2.9340
30	4.1709	3.3158	2.9223

Table from

[http://www.socr.ucla.edu/Applets.dir/F\\_Table.html](http://www.socr.ucla.edu/Applets.dir/F_Table.html); annotations by author

From the above figure, you can see that the **f-critical value** is **3.3541**. Using this value, you can now decide if you will accept or reject the null hypothesis using the **F-distribution curve**:

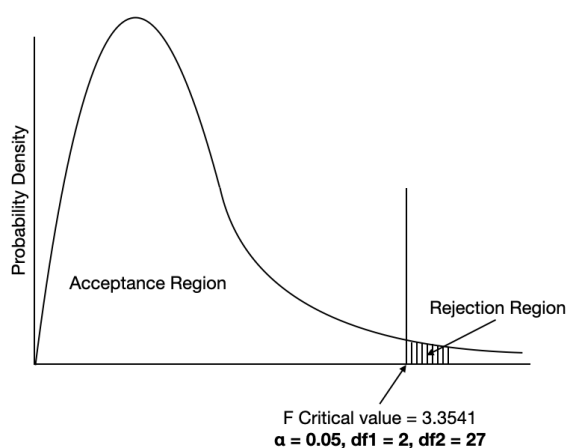


Image by author

Since the **f-value** (0.0171, which is what we can

calculated) is less than the f-critical value in the f-distribution table, we accept the null hypothesis — **this means there is no variance in different groups — all the means are the same.**

*For machine learning, this feature — drink\_type, should **not** be included for training as it seems the different types of drinks have no effect on the reaction time.*

*You should only include a feature for training only if you reject the null hypothesis as this means that the values in the drink types have an effect on the reaction time.*

## **Using the Stats module to calculate f-score**

In the previous section, we manually calculated the f-value for our dataset. Actually, there is an easier way — use the **stats** module's **f\_oneway()** function

to calculate the f-value and p-value:

```
import scipy.stats as stats

fvalue, pvalue = stats.f_oneway(
    df.iloc[: -1, 0],
    df.iloc[: -1, 1],
    df.iloc[: -1, 2])

print(fvalue, pvalue)
# 0.0170760934691432
0.9830794846682348
```

The **f\_oneway()** function takes the groups as input and returns the ANOVA F and p-value:

reaction_time			
drink_type	coffee	coke	water
team			
1	8.0	25.0	14.0
2	20.0	26.0	25.0
3	26.0	27.0	23.0
4	36.0	29.0	27.0
5	39.0	25.0	28.0
6	23.0	23.0	21.0
7	25.0	22.0	26.0
8	28.0	27.0	30.0
9	27.0	29.0	31.0
10	25.0	21.0	34.0
Group Means	25.7	25.4	25.9

Pass these columns to **f\_oneway()**

In the above, the **f-value** is  
**0.0170760934691432**  
(identical to the one we  
calculated manually) and the **p-**  
**value** is  
**0.9830794846682348.**

Observe that the **f\_oneway()**  
function takes in a variable  
number of arguments:

```
fvalue, pvalue = stats.f_oneway(  
    df.iloc[: -1, 0], ←  
    df.iloc[: -1, 1], ← The values of each  
    df.iloc[: -1, 2]) ← column
```

Image by author

If you have many groups, it  
would be quite tedious to pass  
in the values of all the groups  
one by one. So, there is an  
easier way:

```
fvalue, pvalue =  
stats.f_oneway(  
  
    *df.iloc[: -1, 0:3].T.values  
    s  
    )
```

I will leave the above as an  
exercise for you to understand



how it works.

## Using the statsmodels module to calculate f-score

Another way to calculate the f-value is to use the **statsmodel** module. You first build the model using the **ols()** function, and then call the **fit()** function on the instance of the model. Finally, you call the **anova\_lm()** function on the fitted model and specify the type of ANOVA test to perform on it:

*There are 3 types of ANOVA tests to perform, but their discussion is beyond the scope of this article.*

```
import pandas as pd
import statsmodels.api as sm
from statsmodels.formula.api
import ols
```

```
df =
pd.read_csv('drinks.csv')
```

```
model =  
ols('reaction_time ~  
drink_type',  
data=df).fit()  
sm.stats.anova_lm(model,  
typ=2)
```

The above code snippet produces the following result, which is the same as the f-value that we calculated earlier (**0.017076**):

	sum_sq	df	F	PR(>F)
drink_type	1.266667	2.0	0.017076	0.983079
Residual	1001.400000	27.0	NaN	NaN

Image by author

The **anova\_lm()** function also returns the p-value (**0.983079**). You can make use of the following rules to determine if the categorical variable has any influence on the numerical variable:

- if  $p < 0.05$ , this means that the categorical variable has significant influence on the numerical variable
- if  $p > 0.05$ , this means that the categorical variable has

no significant influence on the numerical variable

Since the p-value is now 0.983079 ( $>0.05$ ), this means that the **drink\_type** has no significant influence on the **reaction\_time**.

## Summary

In this article, I have explained how ANOVA helps to determine if a categorical variable has influence on a numerical variable. So far the ANOVA test that we have discussed is known as the **one-way ANOVA** test. There are a few variations of ANOVA:

- **One-way ANOVA**— used to check how a numerical variable responds to the levels of *one* independent categorical variables
- **Two-way ANOVA** —used to check how a numerical variable responds to the levels of *two* independent categorical variables

- **Multi-way ANOVA** — used to check how a numerical variable responds to the levels of *multiple* independent categorical variables

Using a **two-way ANOVA** or **multi-way ANOVA**, you can investigate the combined impact of two (or more) independent categorical variables on one dependent numerical variable.

I hope you find this article useful. Stay tuned for the next article!

**Join  
Medium...**

As a  
Medium...

weimenglee  
.medium.c...



180

|



1

---

## Sign up for The Variable

By Towards Data Science

Every Thursday, the Variable delivers the very best of Towards Data Science: from hands-on tutorials and cutting-edge research to original features you don't want to miss. [Take a look.](#)

Get this  
newsletter