

What is Anova Test?

The ANOVA test, or analysis of variance, is a statistical test used to compare the means of three or more groups to determine if there are any statistically significant differences between them. It allows us to determine if the variation between group means is greater than the variation within the groups.

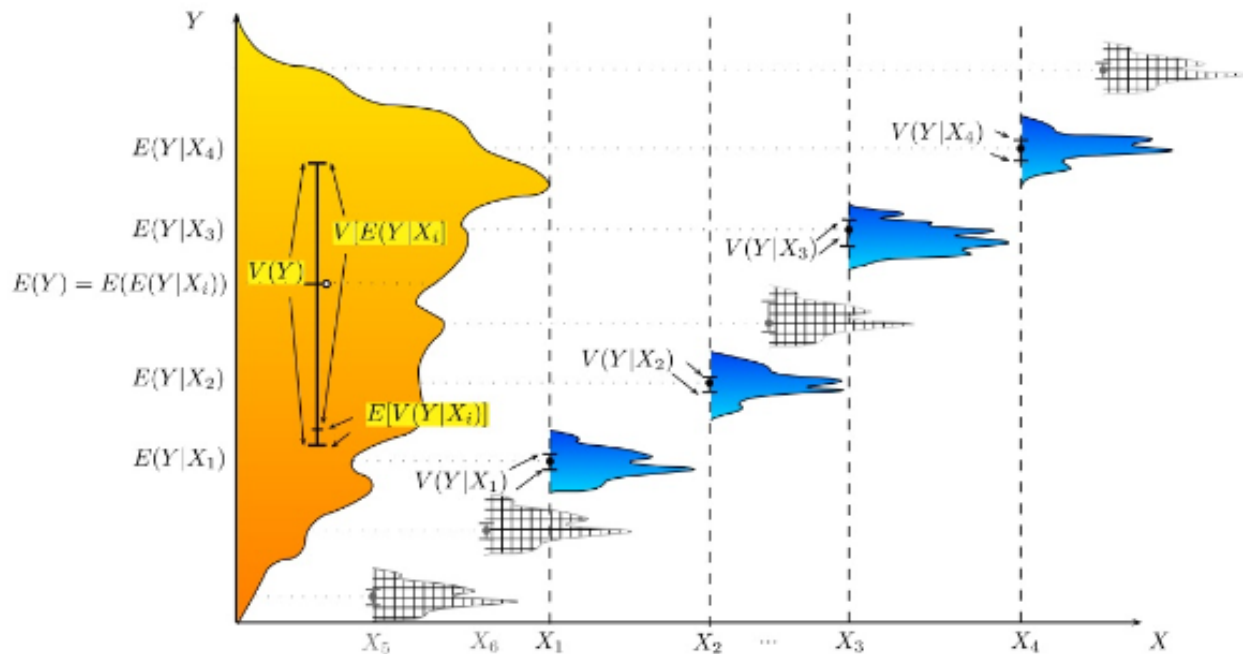


Figure 3: ANOVA : very good fit

* One-Way Anova

One-way ANOVA (**Analysis of Variance**) is a statistical method used to compare the means of three or more independent groups to determine if there are any significant differences between them. It is an extension of the t-test, which is used for comparing the means of two independent groups. The term "one-way" refers to the fact that there is only **one independent variable (factor) with multiple levels (groups) in this analysis**.

The primary purpose of one-way ANOVA is to test

- the null hypothesis that all the group means are equal.
- The alternative hypothesis is that at least one group mean is significantly different from the others.

Steps:

- Define the null and alternative hypotheses.

- Calculate the overall mean (grand mean) of all the groups combined and mean of all the groups individually.
- Calculate the "between-group" and "within-group" sum of squares (SS).
- Find the between group and within group degree of freedoms•
- Calculate the "between-group" and "within-group" mean squares (MS) by dividing their respective sum of squares by their degrees of freedom.
- Calculate the F-statistic by dividing the "between-group" mean square by the "within group"

Analysis of Variance(ANOVA)

Source of Variation	Sum of Squares	Degrees of Freedom	Mean Squares (MS)	F
Within	$SS_w = \sum_{j=1}^k \sum_{i=1}^l (X - \bar{X}_j)^2$	$df_w = k - 1$	$MS_w = \frac{SS_w}{df_w}$	$F = \frac{MS_b}{MS_w}$
Between	$SS_b = \sum_{j=1}^k (\bar{X}_j - \bar{X})^2$	$df_b = n - k$	$MS_b = \frac{SS_b}{df_b}$	
Total	$SS_t = \sum_{j=1}^n (\bar{X}_j - \bar{X})^2$	$df_t = n - 1$		

- Calculate the p-value associated with the calculated F-statistic using the F-distribution and the appropriate degrees of freedom. The p-value represents the probability of obtaining an F-statistic as extreme or more extreme than the calculated value, assuming the null hypothesis is true.
 - Choose a significance level (alpha), typically 0.05.
 - Compare the calculated p-value with the chosen significance level (alpha).
- If the p-value is less than or equal to alpha, reject the null hypothesis in favour of the alternative hypothesis, concluding that there is a significant difference between at least one pair of group means.
 - If the p-value is greater than alpha, fail to reject the null hypothesis, concluding that there is not enough evidence to suggest a significant difference between the group means.

It's important to note that one-way ANOVA only determines if there is a significant difference between the group means; it does not identify which specific groups have significant differences. To determine which pairs of groups are significantly different, post-hoc tests, such as Tukey's HSD or Bonferroni, are conducted after a significant ANOVA result.

number category

A	B	C
3	1	8
6	8	6
3	9	10
\bar{X}_A	\bar{X}_B	\bar{X}_C

ANOVA

$H_0: \mu_A = \mu_B = \mu_C$

H_1 : atleast 1 of them is significant

$n = 9$ $K = 3$

$\bar{X} = 6$ $\bar{X}_A = 4$ $\bar{X}_B = 6$ $\bar{X}_C = 8$

[SST] \rightarrow Sum of square Total $\rightarrow df = n - 1 = 9 - 1 = 8$

$$(6-3)^2 + (6-6)^2 + (6-3)^2 + (6-1)^2 + (6-8)^2 + (6-9)^2 + (6-8)^2 + (6-6)^2 + (6-10)^2$$

$$9 + 0 + 9 + 25 + 4 + 9 + 4 + 0 + 16 = 76$$

[SSW] \rightarrow Sum of squares within $df = 6$ $nc = 9 - 3 = 6$

A	B	C
3	1	8
6	8	6
3	9	10
\bar{X}_A	\bar{X}_B	\bar{X}_C

$\bar{X}_A = 4$ $\bar{X}_B = 6$ $\bar{X}_C = 8$

$(4-3)^2 + (4-6)^2 + (4-3)^2 + (6-1)^2 + (6-8)^2 + (6-9)^2 + (8-8)^2 + (8-9)^2 + (8-10)^2$

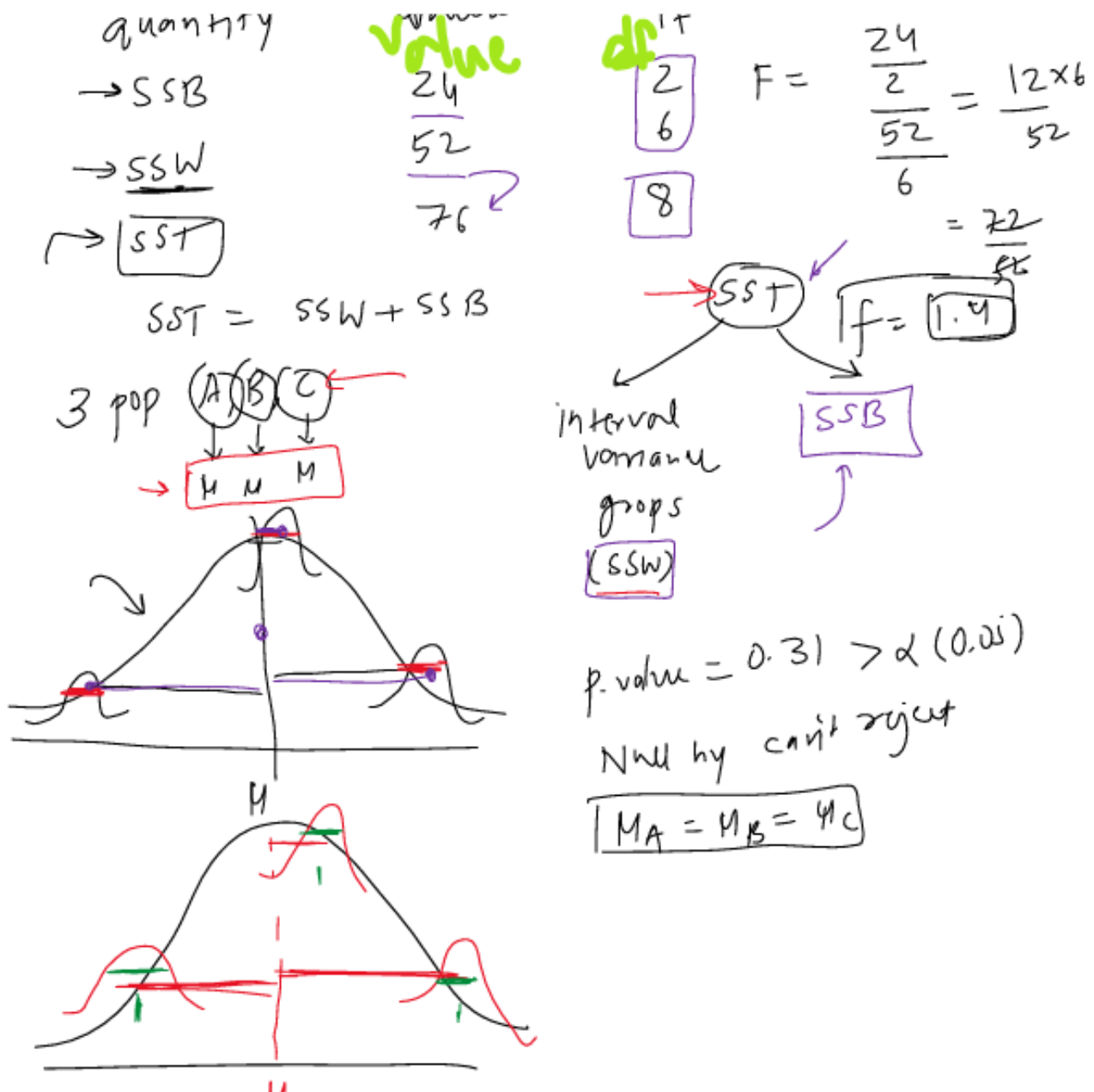
$1 + 4 + 1 + 25 + 4 + 9 + 0 + 4 + 4$

[SSW = 52]

$SSB \rightarrow df = 2$

$$3 \times (6-4)^2 + 3 \times (6-1)^2 + 3 \times (6-8)^2$$

$$3 \times 4 + 0 + 3 \times 4 = 24 = SSB$$



In [3]:

```
import scipy.stats as stats
```

```
f_statistic = 1.4 # The F-statistic value you've calculated
df1 = 2          # Degrees of freedom for the numerator (between groups)
df2 = 6          # Degrees of freedom for the denominator (within groups)
```

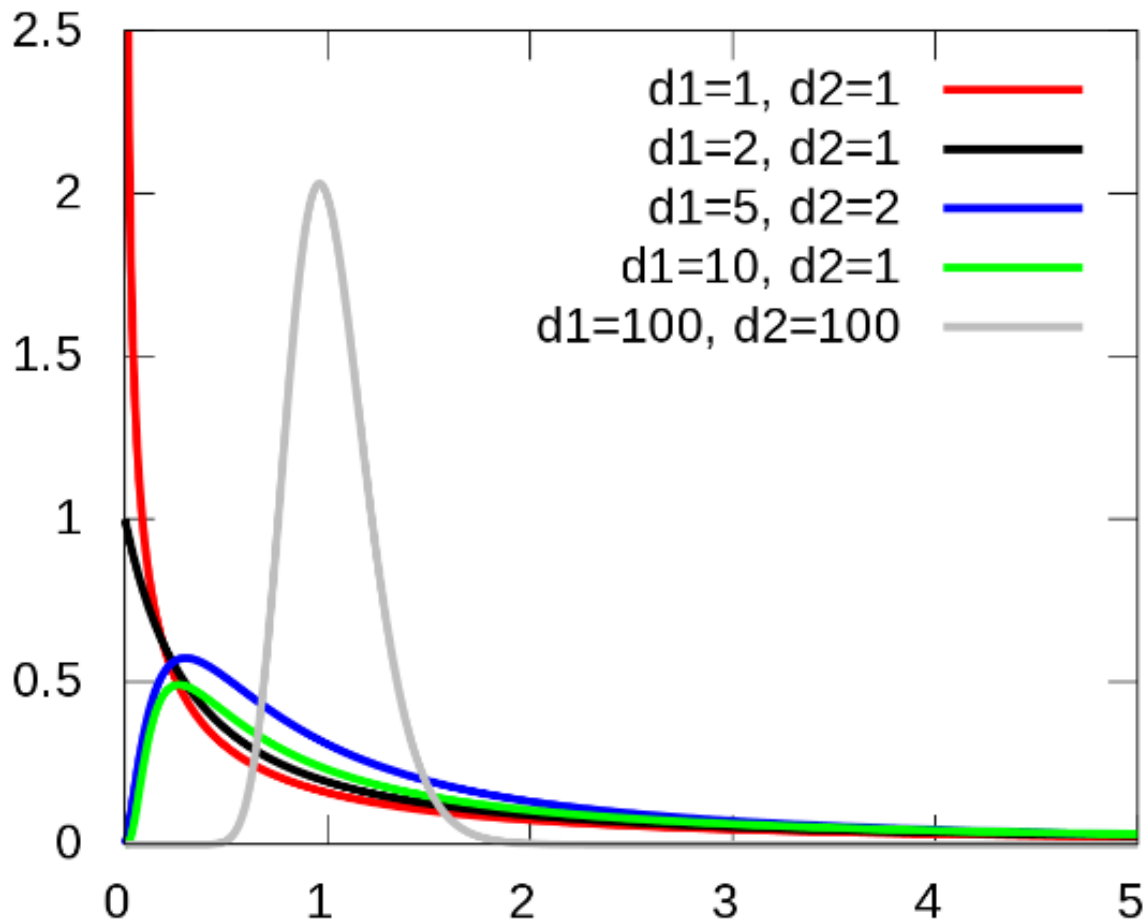
```
p_value = stats.f.sf(f_statistic, df1, df2)
print("P-value:", p_value)
```

P-value: 0.31696093163035305

F - Distribution

The F-distribution, also known as the Fisher-Snedecor distribution, is a probability distribution that arises in statistical inference and is used in various statistical tests, including the analysis of variance (ANOVA) test.

- The F-distribution has two parameters: the numerator degrees of freedom (df_1) and the denominator degrees of freedom (df_2). The shape of the F-distribution depends on these degrees of freedom.



1. **Continuous probability distribution:** The F-distribution is a continuous probability distribution used in statistical hypothesis testing and analysis of variance (ANOVA).
2. **Fisher-Snedecor distribution:** It is also known as the Fisher-Snedecor distribution, named after Ronald Fisher and George Snedecor, two prominent statisticians.
3. **Degrees of freedom:** The F-distribution is defined by two parameters - the degrees of freedom for the numerator (df_1) and the degrees of freedom for the denominator (df_2).
4. **Positively skewed and bounded:** The shape of the F-distribution is positively skewed, with its left bound at zero. The distribution's shape depends on the values of the degrees of freedom.
5. **Testing equality of variances:** The F-distribution is commonly used to test hypotheses about the equality of two variances in different samples or populations.

6. **Comparing statistical models:** The F-distribution is also used to compare the fit of different statistical models, particularly in the context of ANOVA.
7. **F-statistic:** The F-statistic is calculated by dividing the ratio of two sample variances or mean squares from an ANOVA table. This value is then compared to critical values from the F-distribution to determine statistical significance.
8. **Applications:** The F-distribution is widely used in various fields of research, including psychology, education, economics, and the natural and social sciences, for hypothesis testing and model comparison.

Assumptions of Anova

1. **Independence:** The observations within and between groups should be independent of each other. This means that the outcome of one observation should not influence the outcome of another. Independence is typically achieved through random sampling or random assignment of subjects to groups.
2. **Normality:** The data within each group should be approximately normally distributed. While one-way ANOVA is considered to be robust to moderate violations of normality, severe deviations may affect the accuracy of the test results. If normality is in doubt, non-parametric alternatives like the Shapiro-wilk test can be considered.
3. **Homogeneity of variances:** The variances of the populations from which the samples are drawn should be equal, or at least approximately so. This assumption is known as homoscedasticity. If the variances are substantially different, the accuracy of the test results may be compromised. Levene's test or Bartlett's test can be used to assess the homogeneity of variances. If this assumption is violated, alternative tests such as Welch's ANOVA can be used.

Python case Study

In [4]:

```
import pandas as pd
import numpy as np
from scipy.stats import chisquare

# Load the Titanic dataset
url = "https://raw.githubusercontent.com/datasciencedojo/datasets/master/titanic.csv"
df = pd.read_csv(url).dropna(subset=['Age'])

df.head()
```

Out[4]:

PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450

In [5]:

```
df[df['Pclass'] == 3]['Age'].mean()
```

Out[5]:

25.14061971830986

In [6]:

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

# Fit the model
model = ols('Age ~ Pclass', data=df).fit()

# Perform ANOVA analysis
anova_result = sm.stats.anova_lm(model, typ=2)
print(anova_result)
```

	sum_sq	df	F	PR(>F)
Pclass	20511.429755	1.0	112.386893	1.756699e-24
Residual	129945.206190	712.0	NaN	NaN

Post-hoc Test

Post hoc tests, also known as post hoc pairwise comparisons or multiple comparison tests, are used in the context of ANOVA when the overall test indicates a significant difference among the group means. These tests are performed after the initial one-way ANOVA to determine which specific groups or pairs of groups have significantly different means.



The main purpose of post hoc tests is to control the family-wise error rate (FWER) and adjust the significance level for multiple comparisons to avoid inflated Type I errors. There are several post hoc tests available, each with different characteristics and assumptions.

Some common post hoc tests include:

1. **Bonferroni correction:** This method adjusts the significance level (α) by dividing it by the number of comparisons being made. It is a conservative method that can be applied when

making multiple comparisons, but it may have lower statistical power when a large number of comparisons are involved.

2. **Tukey's HSD (Honestly Significant Difference) test:** This test controls the FWER and is used when the sample sizes are equal and the variances are assumed to be equal across the groups. It is one of the most commonly used post hoc tests.

When performing post hoc tests, it is essential to choose a test that aligns with the assumptions of your data (e.g., equal variances, equal sample sizes) and provides an appropriate balance between controlling Type I errors and maintaining statistical power.

In [7]:

```
import scipy.stats as stats

# Perform t-test for three pairs of classes
for class1, class2 in [(1,2), (2,3), (3,1)]:
    print(class1, class2)
    print(stats.ttest_ind(df[df['Pclass'] == class1]['Age'], df[df['Pclass'] == class2]['Age'], nan_policy='omit'))
```

```
1 2
Ttest_indResult(statistic=5.485187676773201, pvalue=7.835568991415144e-08)
2 3
Ttest_indResult(statistic=3.927800191020872, pvalue=9.715078600777852e-05)
3 1
Ttest_indResult(statistic=-10.849122601201033, pvalue=6.134470007830625e-25)
```

In [9]:

```

from statsmodels.stats.multicomp import pairwise_tukeyhsd
import matplotlib.pyplot as plt
import pandas as pd

# Perform Tukey's HSD test
tukey = pairwise_tukeyhsd(endog=df['Age'], groups=df['Pclass'], alpha=0.05)

# Plot simultaneous confidence intervals
tukey.plot_simultaneous()

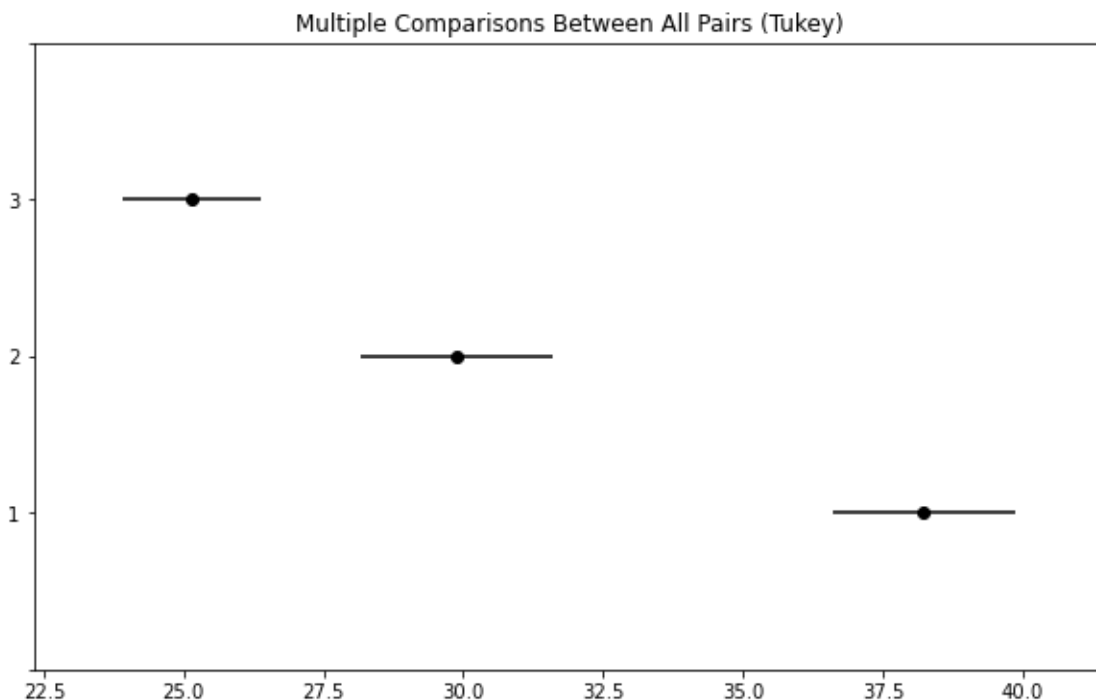
# Add a red vertical line at x=49.57
plt.vlines(x=49.57, ymin=-0.5, ymax=4.5, color="red")

```

C:\Users\user\anaconda3\lib\site-packages\statsmodels\sandbox\stats\multicomp.py:775: UserWarning: FixedFormatter should only be used together with FixedLocator
 ax1.set_yticklabels(np.insert(self.groupsunique.astype(str), 0, ''))

Out[9]:

<matplotlib.collections.LineCollection at 0x20b9cdf6670>



Why t-test is not used for more than 3 categories?

- 1. Increased Type I error:** When you perform multiple comparisons using individual t-tests, the probability of making a Type I error (false positive) increases. The more tests you perform, the higher the chance that you will incorrectly reject the null hypothesis in at least one of the tests, even if the null hypothesis is true for all groups.
- 2. Difficulty in interpreting results:** When comparing multiple groups using multiple t-tests, the interpretation of the results can become complicated. For example, if you have 4 groups and

you perform 6 pairwise t-tests, it can be challenging to interpret and summarize the overall pattern of differences among the groups.

3. **Inefficiency:** Using multiple t-tests is less efficient than using a single test that accounts for all groups, such as one-way ANOVA. One-way ANOVA uses the information from all the groups simultaneously to estimate the variability within and between the groups, which can lead to more accurate conclusions.

Applications in Machine Learning

1. **Hyperparameter tuning:** When selecting the best hyperparameters for a machine learning model, one-way ANOVA can be used to compare the performance of models with different hyperparameter settings. By treating each hyperparameter setting as a group, you can perform one-way ANOVA to determine if there are any significant differences in performance across the various settings.
2. **Feature selection:** One-way ANOVA can be used as a univariate feature selection method to identify features that are significantly associated with the target variable, especially when the target variable is categorical with more than two levels. In this context, the one-way ANOVA is performed for each feature, and features with low p-values are considered to be more relevant for prediction.
3. **Algorithm comparison:** When comparing the performance of different machine learning algorithms, one-way ANOVA can be used to determine if there are any significant differences in their performance metrics (e.g., accuracy, F1 score, etc.) across multiple runs or cross-validation folds. This can help you decide which algorithm is the most suitable for a specific problem.
4. **Model stability assessment:** One-way ANOVA can be used to assess the stability of a machine learning model by comparing its performance across different random seeds or initializations. If the model's performance varies significantly between different initializations, it may indicate that the model is unstable or highly sensitive to the choice of initial conditions.

- Prudhvi Vardhan

In []: