

T- Test :- A t-test is a type of inferential statistic which is used to determine if there is a significant difference between the means of two groups which may be related in certain features.

```

1 from scipy.stats import ttest_1samp
2 import numpy as np
3 import pandas as pd
4 ages = pd.read_csv('/content/ages.csv')
5 print(ages)
6 ages_mean = np.mean(ages)
7 print(ages_mean)
8 tset, pval = ttest_1samp(ages, 30)
9 print("p-values",pval)
10 if pval.any() < 0.05:    # alpha value is 0.05 or 5%
11     print("we are rejecting null hypothesis")
12 else:
13     print("we are accepting null hypothesis")

```

```

      n  age
0     1   20
1     2   99
2     3   42
3     4   79
4     5   66
5     6   55
6     7   74
7     8   39
8     9   33
9    10   81
10   11   36
11   12   27
12   13   95
13   14   50
14   15   15
15   16   25
16   17   84
17   18   56
18   19   50
19   20   98
20   21   20
21   22   59
22   23   96
23   24   72
24   25   67
25   26   39
26   27   34
27   28   11
28   29   43
29   30   28
n      15.5
age    53.1
dtype: float64

```

```
p-values [6.46799481e-10 5.02652426e-05]  
we are accepting null hypothesis
```

Two sampled T-test :-The Independent Samples t Test or 2-sample t-test compares the means of two independent groups in order to determine whether there is statistical evidence that the associated population means are significantly different. The Independent Samples t Test is a parametric test. This test is also known as: Independent t Test.

```
1 from scipy.stats import ttest_ind  
2 import numpy as np  
3 import pandas as pd  
4 from pandas import DataFrame  
5  
6 sales = pd.read_csv('/content/weeks.csv')  
7 df=DataFrame(data=sales, columns=['week1', 'week2'])  
8 df
```

	week1	week2
0	144	2663
1	102	2461
2	1781	1046
3	2749	759
4	3222	4362
5	1907	2422
6	4968	2480
7	1996	1576
8	1963	602
9	1971	2049
10	2207	4775
11	3377	776



```

1 week1_mean = np.mean(sales['week1'])
2 week2_mean = np.mean(sales['week2'])
3 print("week1 mean value:", week1_mean)
4 print("week2 mean value:", week2_mean)
5 week1_std = np.std(sales['week1'])
6 week2_std = np.std(sales['week2'])
7 print("week1 std value:", week1_std)
8 print("week2 std value:", week2_std)
9 ttest, pval = ttest_ind(sales['week1'], sales['week2'])
10 print("p-value", pval)
11 if pval < 0.05:
12     print("we reject null hypothesis")
13 else:
14     print("we accept null hypothesis")

```

```

week1 mean value: 2221.633333333333
week2 mean value: 2633.9
week1 std value: 1486.1584366711677
week2 std value: 1412.0187758430598
p-value 0.2832938455413314
we accept null hypothesis

```

Paired sampled t-test: The paired sample t-test is also called dependent sample t-test. It's an univariate test designed to identify a significant difference between 2 related variables.

```

21 2007 4059

```

```

1 import pandas as pd
2 from scipy import stats
3 df = pd.read_csv("/content/blood_pressure.csv", encoding='utf-8')

```

```

4 df[['bp_before', 'bp_after']].describe()
5 ttest,pval = stats.ttest_rel(df['bp_before'], df['bp_after'])
6 print(pval)
7 if pval<0.05:
8     print("reject null hypothesis")
9 else:
10    print("accept null hypothesis")
    0.0011297914644840823
    reject null hypothesis

```

Z Test

Several different types of tests are used in statistics (i.e. f test, chi square test, t test). You would use a Z test if:

- Your sample size is greater than 30. Otherwise, use a t test.
- Data points should be independent from each other. In other words, one data point is not related or does not affect another data point.
- Your data should be normally distributed. However, for large sample sizes (over 30) this doesn't always matter.
- Your data should be randomly selected from a population, where each item has an equal chance of being selected.
- Sample sizes should be equal if at all possible.

```

1 import pandas as pd
2 from scipy import stats
3 from statsmodels.stats import weightstats as stests
4 ztest ,pval = stests.ztest(df['bp_before'], x2=None, value=156)
5 print(float(pval))
6 if pval<0.05:
7     print("reject null hypothesis")
8 else:
9     print("accept null hypothesis")

0.6651614730255063
accept null hypothesis
/usr/local/lib/python3.7/dist-packages/statsmodels/tools/_testing.py:19: FutureWarning
import pandas.util.testing as tm

```

Two-sample Z test

In two sample z-test, similar to t-test, we are checking two independent data groups and deciding whether the sample mean of two group is equal or not.

- H_0 : mean of two group is 0
- H_1 : mean of two group is not 0

```

1 ztest ,pval1 = stats.ztest(df['bp_before'], x2=df['bp_after'], value=0,alternati
2 print(float(pval1))
3 if pval<0.05:
4     print("reject null hypothesis")
5 else:
6     print("accept null hypothesis")

0.002162306611369422
accept null hypothesis

```

ANOVA (F-TEST) :

- The t-test works well when dealing with two groups, but sometimes we want to compare more than two groups at the same time. For example, if we wanted to test whether voter age differs based on some categorical variable like race, we have to compare the means of each level or group the variable. We could carry out a separate t-test for each pair of groups, but when you conduct many tests you increase the chances of false positives. The analysis of variance or ANOVA is a statistical inference test that lets you compare multiple groups at the same time.
- $F = \text{Between group variability} / \text{Within group variability}$

One Way F-test (ANOVA) :

- It determines whether two or more groups are similar or not based on their mean similarity and f-score.

```

1 df_anova = pd.read_csv('/content/plant_growth.csv')
2 df_anova = df_anova[['weight', 'group']]
3 grps = pd.unique(df_anova.group.values)
4 d_data = {grp:df_anova['weight'][df_anova.group == grp] for grp in grps}
5
6 F, p = stats.f_oneway(d_data['ctrl'], d_data['trt1'], d_data['trt2'])
7 print("p-value for significance is: ", p)
8 if p<0.05:
9     print("reject null hypothesis")
10 else:
11     print("accept null hypothesis")

p-value for significance is: 0.0159099583256229
reject null hypothesis

```

Two Way F-test:

- Two way F-test is an extension of one-way F-test.
- It is used when we have two independent variables and 2+ groups.


- A two-way F-test does not tell which variable is dominant.
- If we need to check individual significance, then Post-hoc testing needs to be performed.

```

1 import statsmodels.api as sm
2 from statsmodels.formula.api import ols
3 df_anova2 = pd.read_csv("https://raw.githubusercontent.com/Opensourcefordatascier
4 model = ols('Yield ~ C(Fert)*C(Water)', df_anova2).fit()
5 print(f"Overall model F({model.df_model: .0f},{model.df_resid: .0f}) = {model.fv
6 res = sm.stats.anova_lm(model, typ= 2)
7 res

```

Overall model F(3, 16) = 4.112, p = 0.0243

	sum_sq	df	F	PR(>F)	
C(Fert)	69.192	1.0	5.766000	0.028847	
C(Water)	63.368	1.0	5.280667	0.035386	
C(Fert):C(Water)	15.488	1.0	1.290667	0.272656	
Residual	192.000	16.0	NaN	NaN	

Chi-Square Test:

- The test is applied when you have two categorical variables from a single population.
- It is used to determine whether there is a significant association between the two variables.
- For example, in an election survey, voters might be classified by gender (male or female) and voting preference (Democrat, Republican, or Independent).
- We could use a chi-square test for independence to determine whether gender is related to voting preference.

```

1 df_chi = pd.read_csv('/content/chi_square_test.csv')
2 df_chi

```

Gender Like Shopping?**0 Male No****1 Female Yes**

```

1 contingency_table=pd.crosstab(df_chi["Gender"],df_chi["Like Shopping?"])
2 print('contingency_table :-\n',contingency_table)
3 #Observed Values
4 Observed_Values = contingency_table.values
5 print("Observed Values :-\n",Observed_Values)
6 b=stats.chi2_contingency(contingency_table)
7 Expected_Values = b[3]
8 print("Expected Values :-\n",Expected_Values)
9 no_of_rows=len(contingency_table.iloc[0:2,0])
10 no_of_columns=len(contingency_table.iloc[0,0:2])
11 ddof=(no_of_rows-1)*(no_of_columns-1)
12 print("Degree of Freedom:-",ddof)
13 alpha = 0.05
14 from scipy.stats import chi2
15 chi_square=sum([(o-e)**2./e for o,e in zip(Observed_Values,Expected_Values)])
16 chi_square_statistic=chi_square[0]+chi_square[1]
17 print("chi-square statistic:-",chi_square_statistic)
18 critical_value=chi2.ppf(q=1-alpha,df=ddof)
19 print('critical_value:',critical_value)
20 #p-value
21 p_value=1-chi2.cdf(x=chi_square_statistic,df=ddof)
22 print('p-value:',p_value)
23 print('Significance level: ',alpha)
24 print('Degree of Freedom: ',ddof)
25 print('chi-square statistic:',chi_square_statistic)
26 print('critical_value:',critical_value)
27 print('p-value:',p_value)
28 if chi_square_statistic>=critical_value:
29     print("Reject H0,There is a relationship between 2 categorical variables")
30 else:
31     print("Retain H0,There is no relationship between 2 categorical variables")
32
33 if p_value<=alpha:
34     print("Reject H0,There is a relationship between 2 categorical variables")
35 else:
36     print("Retain H0,There is no relationship between 2 categorical variables")

```

contingency_table :-

Like Shopping? No Yes

Gender

Female 2 3

Male 2 2

Observed Values :-

[[2 3]

[2 2]]

Expected Values :-

[[2.22222222 2.77777778]

```
[1.77777778 2.22222222]]
Degree of Freedom:- 1
chi-square statistic:- 0.090000000000000008
critical_value: 3.841458820694124
p-value: 0.7641771556220945
Significance level: 0.05
Degree of Freedom: 1
chi-square statistic: 0.090000000000000008
critical_value: 3.841458820694124
p-value: 0.7641771556220945
Retain H0, There is no relationship between 2 categorical variables
Retain H0, There is no relationship between 2 categorical variables
```

✓ 0s completed at 4:17 PM

