

LEC.VIII. T-test & Code

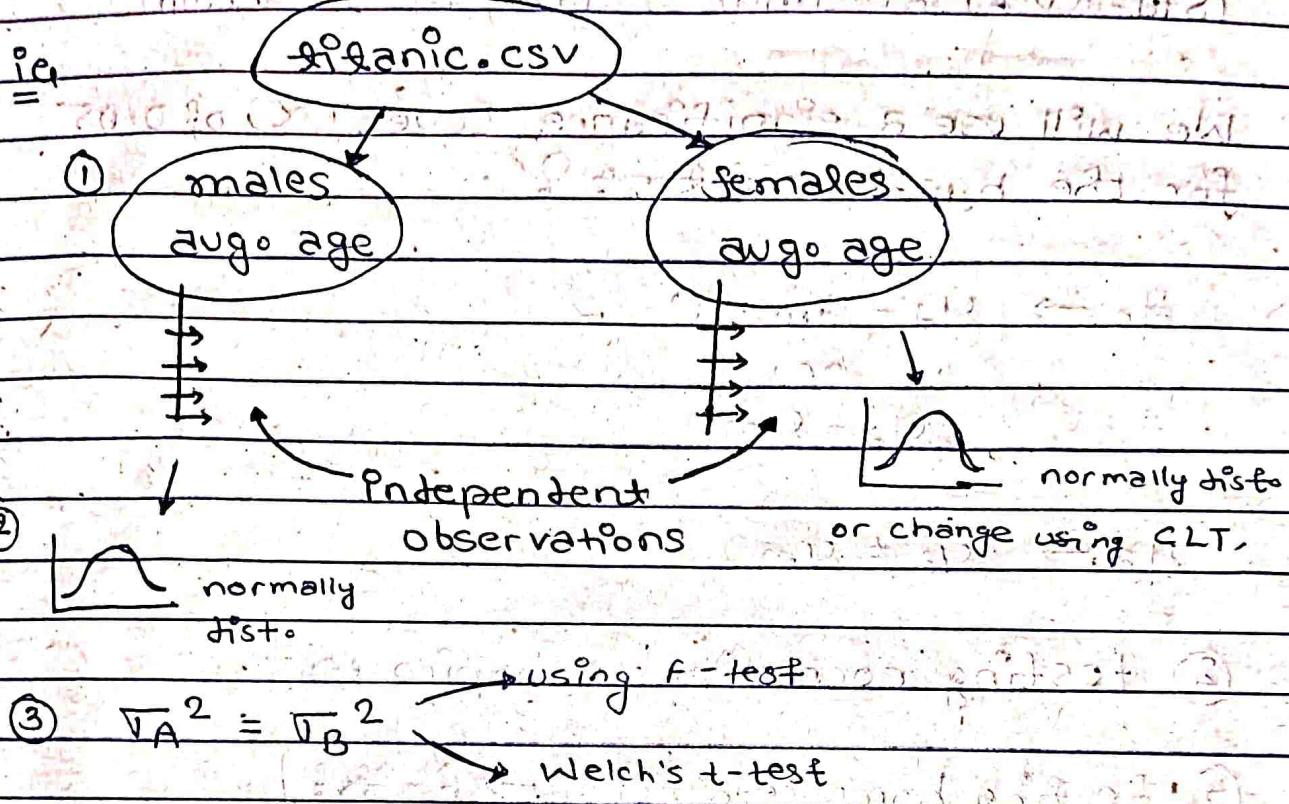
(B) Independent two sample test:

Used to compare the means of two-independent samples. The null hypothesis states that there's no significant difference b/w the means of two samples while the alternate hypothesis states that there's a significant difference.

○ Assumptions for t-test:

- ① Independence of two samples, The two sample must be independent, meaning there's no relation b/w the observations in one group & the observations in other group. The subjects in the two groups should be selected randomly & independently.
- ② Normality, The data in each of the two groups should be approximately normal dist. The t-test is considered robust to mild violations of normality, especially when the sample sizes are large ($n \geq 30$) and the sample sizes of the two groups are similar.
- ③ Equal Variances (Homoscedasticity), The variance of the two populations should be approximately equal. These assumptions can be checked using F-test, or if this assumption not met you can use Welch's t-test.
- ④ Random Sampling, The data should be collected using a random sampling method from the respective populations. This ensures that the sample is representative of population & reduces the risk of selection bias.

(Independent Two Sample or Unpaired T-test) is a statistical method used to compare means of two independent groups to determine if there's a significant difference b/w them.



Q) Suppose a website owner claims that there's no difference in the avg. time spent on their website b/w desktop & mobile users. To test this claim, we collect data from 30 desktop users and 30 mobile users regarding the time spent on the website in minutes. The sample statics are as follows —

desktop users = { 12, 15, 18, 16, 20, 17, 14, 22, 19, 21, 23, 18, 25, 17, 16, 24, 20, 19, 22, 18, 15, 14, 23, 16, 12, 21, 19, 17, 20, 14 }

mobile-users = { 10, 12, 14, 13, 16, 15, 11, 17, 14, 16, 18, 14, 20, 15, 14, 19, 16, 15, 17, 14, 12, 11, 18, 15, 10, 16, 15, 13, 16, 11 }

Desktop user

mobile user

n1 : 30

n2 : 30

mean 1 : 10.5 min

mean 2 : 14.3 min

(std-dev 1) : 3.5 min

(std-dev 2) : 2.7 min

We will use a significance level (α) of 0.05 for the hypothesis test?

$$H_0 \rightarrow \mu_d - \mu_m = 0$$

or

$$\mu_d = \mu_m$$

$$H_a \rightarrow \mu_d \neq \mu_m$$

② testing normality \rightarrow Shapiro test

③ Testing homoscedasticity \rightarrow F-test /

Welch's t-test

(discussed earlier)

desktop

mobile

statistic $\rightarrow 0.97031141$ 0.97143560

p-value $\rightarrow 0.779096$ 0.57916086

variances \rightarrow using Levene

if p-value < 0.05

statistic: 2.9439

$\sqrt{A^2} \neq \sqrt{B^2}$ \rightarrow p-value ≈ 0.9153

p-value > 0.05

$$\sqrt{A^2} = \sqrt{B^2}$$

$0.9153 > 0.05$

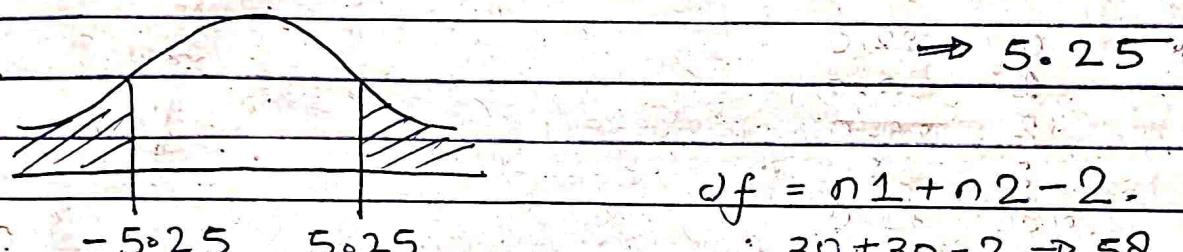
$$\sqrt{A^2} = \sqrt{B^2}$$

Both assumptions are true.

$$t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

$\bar{x}_1 = \frac{\bar{x}_1 - \bar{x}_2}{s_1 \sqrt{n_1}}$ replaces

$$\Rightarrow \frac{18.5 - 14.3}{\sqrt{\frac{(3.5)^2}{30} + \frac{(2.7)^2}{30}}} \Rightarrow 4.2$$



$$df = n_1 + n_2 - 2$$

$$30 + 30 - 2 \Rightarrow 58$$

p-value $\rightarrow 0.00000 - 256$

reject Null Hypothesis

Single Sample t-test code:

```
import seaborn as sns
import pandas as pd
import numpy as np

df = sns.load_dataset('titanic')
df.head()

pop = df['age'].dropna()
pop
# 0    22.0
# 2    30.0
# 2    26.0
# 2    32.0
# 2    32.0
# Name: age, length: 74, dtype: float64

sample_age = pop.sample(25).values
print('sample-age:', sample_age)

t_statistic, p_value = stats.ttest_1samp(sample_age, pop.mean())
print('t-statistic:', t_statistic)
print('p-value:', p_value)

alpha = 0.05

if p_value < alpha:
    print("Reject Null Hypo")
else:
    print("Fail to reject Null Hypo")

# Ho: The mean Age is 35
# Ha: The mean Age is less than 35.

# As the sample size is less than 30 (i.e. 25)
# Applying Shapiro test to check whether
# its Normally Distributed
```

from scipy.stats import shapiro

```
shapiro_age = shapiro(sample_age)
shapiro_age
#> ShapiroResult(statistic=0.955, pvalue=0.322)
```

pvalue > 0.05 --> Normally Distribute

Independent two-sample t-Test

import seaborn as sns

```
df = sns.load_dataset('titanic')  
df.head()
```

```
pop_male = df[df['sex'] == 'male']['age']  
dropna
```

```
pop_female = df[df['sex'] == 'female']['age']  
dropna()
```

```
sample_male = pop_male.sample(25)
```

```
sample_female = pop_female.sample(25)
```

alpha = 0.05

H0 : mean age of male & female are same

H1 : mean age of male is higher than female

from 'scipy.stats import levene'

```
Levens_test = levene(sample_male, sample  
- female)
```

print(Levens_test)

Output (statistic = 0.30, pvalue = 0.58)

0.58 > 0.05

```
import scipy.stats as stats  
t_statistic, p-value = stats.ttest_ind  
(sample_male, sample_female)  
print('t-statistic')  
print(p-value/2)
```

C Paired two sample test:

Used to compare the means of two related or dependent groups.

e.g. ① before-and-after studies — comparing the performance of ~~two~~ groups before and after an intervention or treatment.

② Matched or correlated groups — Comparing the performance of two groups that are correlated or matched in some way, such as siblings.

Assumptions: ① paired observations — The two set of observation must be related or paired in some way.

② Normality ③ Independence of pairs

Q) Let's assume that a fitness centre is evaluating the effectiveness of a new 8 week weight loss program. They enroll 15 participants in the program and measure their weight before & after program. The goal is to test whether the new weight loss program leads to a significant reduction in the participants weight?

significance level (α) = 0.05

name	wt before	wt after	difference (difference)
1	80	73	-7
2	92	90	-2
3	75	81	-6
4	68	67	-1
5	85	88	-3
6	78	76	-2
7	73	74	-1
8	90	91	-1
9	70	69	-1
10	88	88	0
11	76	77	-1
12	84	81	-3
13	82	80	-2
14	77	79	-2
15	91	80	-11

$$\sum d \Rightarrow -1 \quad \sum d^2 \Rightarrow 16$$

$$(\bar{d})^2 \Rightarrow 1$$

$$n = 15$$

$H_0: \mu_{\text{before}} = \mu_{\text{after}}$

$H_a: \mu_{\text{before}} > \mu_{\text{after}}$

check → difference → is normally distributed.

$$t = \frac{\sum d}{\sqrt{\frac{n(\sum d^2) - (\sum d)^2}{n-1}}} \Rightarrow \frac{-1}{\sqrt{\frac{15(165) - 1}{14}}} = -1$$

$$\Rightarrow \frac{-1}{\sqrt{\frac{2474}{14}}} \Rightarrow \frac{-1}{\sqrt{176.71}} \Rightarrow \frac{-1}{13.29} \Rightarrow -0.075$$

* Chi-Square Test

- A chi-square test (χ^2) is basically a comparison of two statistical data sets.
- The chi-square test is used to estimate how likely the observation that are made would be, by considering the assumption of the Null Hypo. as true.
- When we consider the null-speculation is true, the sampling distribution of the test-statistic is called chi-squared distribution.
- The chi-square test helps us to determine whether there's a notable difference b/w the normal freq. and observed frequencies.

↳ Applicable for categorical data such as men & women falling under categories of gender, age, height etc.

To calculate p-value, chi-square test is used

- $P \leq 0.05 \rightarrow$ rejected
- $P > 0.05 \rightarrow$ Accepted

Properties: (1) mean = 2 variance = 2
(2) $\nu_1 + \nu_2 = \chi_1^2 + \chi_2^2$

(3) chi-square distribution curve approaches the normal distribution when dof increases

$$\chi^2 = \sum \frac{(\text{observed value} - \text{expected value})^2}{\text{Expected value}}$$

or

$$\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i}$$

(4) A survey on cars had conducted in 2011 and determined that 60% of cars owners have only one car, 28% have two cars and 12% have three or more. Supposing that you have decided to conducting your own survey and have collected the data below, determine whether your data supports the result of the day?

Use a significance level of 0.05 & given that out of 129 car owners, 73 had one car & 38 had two cars?

H_0 : The proportion of car owners with one car $\rightarrow 0.60$, two car $\rightarrow 0.28$, three or more $\rightarrow 0.12$

H_a : The proportions doesn't match the proposed model.

A chi-square goodness of fit test is appropriate because we are examining the distribution of a single categorical variable.

	Observed(O_i)	Expected(E_i)	$O_i - E_i$	$(O_i - E_i)^2$	$/E_i$
One car	73	$0.60 \times 129 \Rightarrow 77.4$	-4.4	19.36	0.250
two car	38	$0.28 \times 129 \Rightarrow 36.1$	1.9	3.61	0.1
three or more car	18	$0.12 \times 129 \Rightarrow 15.5$	2.5	6.25	0.40
total	129				0.7533

$$df \Rightarrow 3 - 1 \Rightarrow 2 \quad \chi^2$$

using table, the critical value of a 0.05 significance level with a $df = 2$ is 5.99

That means 95 times out of 100, a survey that agrees with a sample will have a χ^2 value of 5.99 or less

The chi-square statistic is only 0.7533, so we'll accept Null hypothesis

Code Implementation:

The test is applied when you have two categorical variables from a single population. It's used to determine there's a significant association between variables.

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

```
import seaborn as sns
```

```
pandas as pd
```

```
numpy as np
```

```
df = sns.load_dataset('tips')
```

```
df.head()
```

↳ whether there's any association b/w "sex" & "smoker".

```
alpha = 0.05
```

```
df_table = pd.crosstab(df['sex'], df['smoker'])
```

↳ smoker yes no
sex
male 60 93
female 33 54

```
# Observed values
```

```
obs = df_table.values
```

```
print("Observed values : \n", obs)
```

```
([[60, 93], [33, 54]])
```

val = stats.chisq_contingency(df_table)

(0.0008, 0.92, 1, array([55.01, 93.15, 33.45, 53.54]))

↳ chi-square statistics, pvalues, expected-values

Expected-values = val[3]

```
nos_rows = len(df_table.iloc[0:2, 0])
```

```
nos_cols = len(df_table.loc[0, 0:2])
```

```
df = (nos_rows - 1) * (nos_cols - 1)
```

```
print("Degree of freedom : ", df)
```

```
chi_square = sum([(o - e)**2 / e for o, e  
in zip(obs, expected_values)])
```

```
chi_square_statistic = chi_square[0] +
```

```
chi_square[1].
```

```
(0.0019)
```

critical_value = chi2.ppf(q=1 - alpha, df=df)

print("critical-value is", critical_value)

Output: 3.841

pvalue ↴

p-value = 1 - chi2.cdf(x=chi-square-statistic, df=df)

Output: 0.964

if chi-square-statistic >= critical_value:
print("Reject H0, there's a reln b/w two categorical variables")

else:

print("Retain H0, There's no reln b/w two categorical variables")

if p-value <= alpha:

else: