

## Lecture - 06

- CHI - SQUARE Test.
- Covariance.
- Pearson Correlation, Coefficient.
- Spearman Rank Correlation.
- Practical Implementation [z-test, t-test, Chi-Square]
- F-Test [ANOVA]

### CHI - SQUARE TEST:

- Chi-Square test Claims about population proportions.
- It is a non parametric test that is performed on Categorical (nominal or ordinal) data.

Problem: [Based on Chi-Square test]

1. In the 2000 Indian Census, the age of the individual in a small town were found to be the following.

Less than 18	Age (18 - 35)	Above 35
20%	30%	50%

In 2010, age of  $n = 500$  individuals were Sampled.

Below are the results.

$< 18$	$18-35$	$> 35$
100	121	288 <del>268</del> 250

Would you conclude the population distribution of ages has changed in the last 10 years? [using  $\alpha = 0.05$ ]

Given,

In 2000,

$< 18$	$18-35$	$> 35$
20%	30%	50%

In 2010,

	$< 18$	$18-35$	$> 35$
2010	121	288	91
2000 with respect to Value	$0.2 \times 500$ $= 100$	$0.3 \times 500$ $= 150$	$0.5 \times 500$ $= 250$

with  
Sample mean

	$<18$	$18-35$	$>35$	
$2010$	(2)	288	91	Observation
$2000$	100	150	250	Expected

(i) Null Hypothesis ( $H_0$ ),

The Observation data meets the distribution of 2000 Census.

(ii) Alternative Hypothesis ( $H_1$ ),

The observation data does not meets the distribution of 2000 Census.

(iii) Significance value  $\alpha=0.05$  (95% of CI)  
[Given].

(iv) Degree of freedom =  $n-1$

not take  $n=500$ , because it is applicable for Categorical values. So, take the category of  $<18, 18-35, >35$ .

3 different Samples.

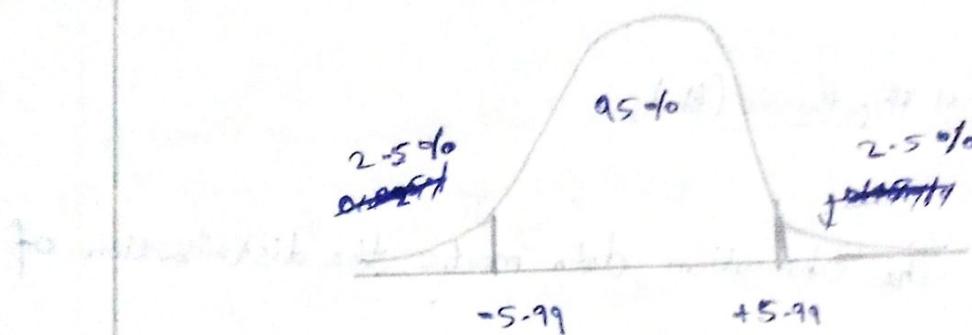
$$= n-1 = 3-1 = 2$$

$$\boxed{df=2}$$

## (V) Decision Boundary,

If the population distribution is high or low, so,

Use two-tail test.  $[df=2, \alpha=0.05]$



$$df=2, \alpha=0.05$$

$\downarrow$  Chi-square table  
Value  
 $5.99$

## (vi) Calculate Test Statistics:

$$\chi^2 = \sum \frac{(f_o - f_e)^2}{f_e}$$

$f_o \rightarrow$  observation value

$f_e \rightarrow$  expected value

$\chi^2 \rightarrow$  Representation of

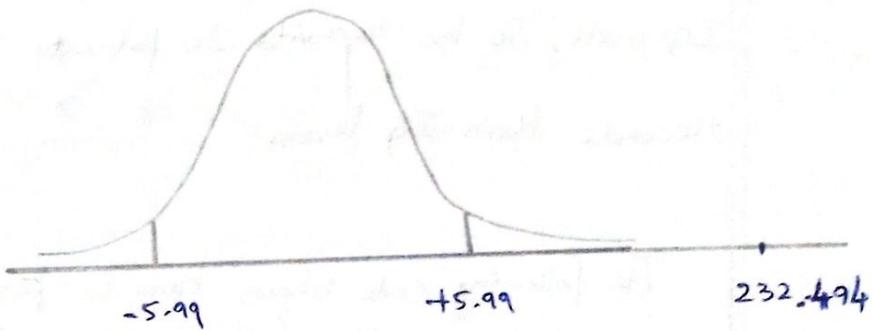
$$= \frac{(121-100)^2}{100} + \frac{(288-150)^2}{150} + \frac{(91-250)^2}{250}$$

$$\chi^2 = 232.494$$

(iii) State Decision:

$$\chi^2 = 232.494 > 5.99 \quad \left\{ \begin{array}{l} \text{Reject the null} \\ \text{hypothesis} \end{array} \right\}$$

So, we accept alternate hypothesis.



Obviously, The population distribution of ages was changed and increased in the last 10 years.

Z-test	T-Test	CHI-SQUARE TEST
formula, Point estimate $\pm$ margin of error	Formula, Point estimate $\pm$ margin of error	formula, Test stats $\chi^2 = \sum \frac{(f_o - f_e)^2}{f_e}$
$\bar{x} \pm z_{\alpha/2} \left( \frac{\sigma}{\sqrt{n}} \right) \rightarrow CI$ Test statistics, $Z = \frac{\bar{x} - \mu}{\left[ \frac{\sigma}{\sqrt{n}} \right]}$	$\bar{x} \pm t_{\alpha/2} \left( \frac{s}{\sqrt{n}} \right) \rightarrow CI$ Test statistics, $t = \frac{\bar{x} - \mu}{\left[ \frac{s}{\sqrt{n}} \right]}$	upperbound = $[x^2(1-\alpha), \infty)$ lowerbound = $(0, x^2(\alpha)]$ $CI = [x^2(1-\alpha/2), x^2(\alpha/2)]$

## PRACTICAL: [Python] Z-test

Suppose the IQ in a certain population is normally distributed with a mean  $\mu = 100$  and standard deviation of  $\sigma = 15$ .

A researcher wants to know if a new drug affects IQ levels, so he recruits 20 patients to try it and records their IQ levels.

The following code shows how to perform a One Sample Z-test in Python to determine if the new drug causes a significant difference in IQ levels.

### Program:

```
from statsmodels.stats.weightstats import ztest as ztest  
  
# Enter IQ level for 20 patients.  
  
data = [88, 92, 94, 94, 96, 97, 97, 97, 99, 99,  
        105, 109, 109, 110, 112, 112, 113, 114, 115]  
  
ztest(data, value=100) # Value is Population Mean
```

Output:

(1.5976240527147705, 0.1101266701438426)



Z-test Value



P-value

Conditions:

P-value < Significance value



Reject the null hypothesis

P-value > Significance value



Accept the null hypothesis

Based on P-value, The Computer rejects or accepts the hypothesis.

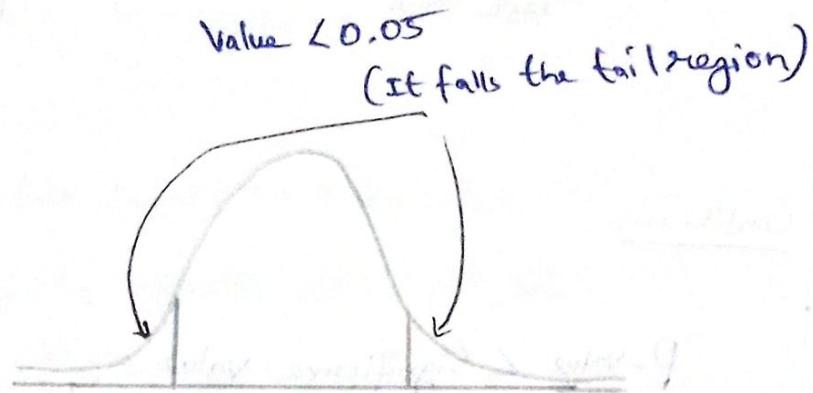
The above Program P-value > Significance Value.

0.11 > 0.05 (default significance value)

So, Accept the null hypothesis.

So, the

So, the new dog cannot have a significant difference in IQ levels.



Co-Variance:

Let us consider, Two features are

X

Y

Weight

Height

50

160

60

170

70

180

75

181

Note,

$X \uparrow Y \uparrow$   
 $X \downarrow Y \downarrow$

Let us consider another one,

X	Y
No of hours	play
Study	
2	6
3	4
4	3
5	2

Note,

$$\begin{array}{l} x \uparrow y \downarrow \\ x \downarrow y \uparrow \end{array}$$

Covariance is the measure of two relationship between two features in the form of Quantity.

formula,

$$\text{Covariance, } \text{Cov}(x, y) = \sum_{i=1}^n \frac{(x_i - \bar{x})(y_i - \bar{y})}{N}$$

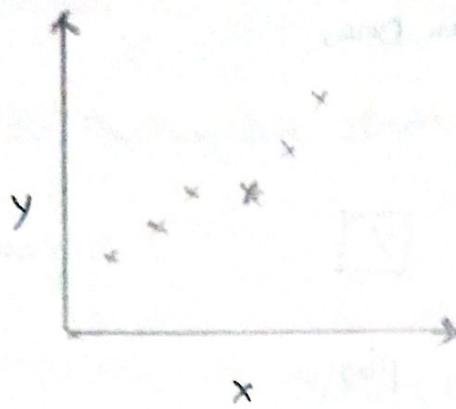
Working with Sample  $(n-1)$

instead of  $N$  in the formula

Apply this formula, we get

$$= +ve (or) -ve$$

Positive                      Negative  
Correlation                  Correlation



positive

$$x \uparrow y \uparrow$$

Correlation

$$x \downarrow y \downarrow$$



Negative

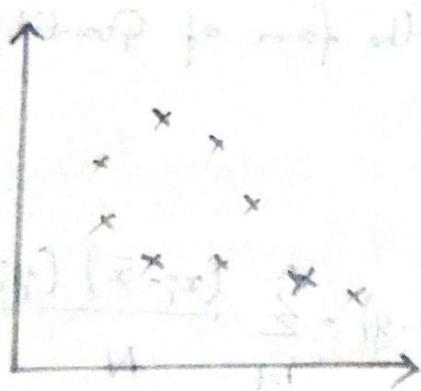
correlation

$$x \uparrow y \downarrow$$

$$x \downarrow y \uparrow$$

$\Downarrow$  (or) say

Negative Covariance



Covariance = 0

[There is no relationship

between the features]

### Disadvantage of Covariance:

1. Return the positive or negative value but the

Values are quite large like +100, +1000, -3000 like

this. So we cannot measure within the particular

range.

2. The magnitude there is no such limit (i.e. between this range to this range). So, it is hard to measure.

Advantage:

We can find the direction is positive or negative.

But the problem is not measurable.

To solve this, we use Pearson Correlation Coefficient.

Pearson Correlation Coefficient:

This is used for measure the relationship within particular range from (-1 to 1).

→ The more towards +1 then say More  
Positively Correlation.

→ The more towards -1 then say More  
Negatively Correlation.

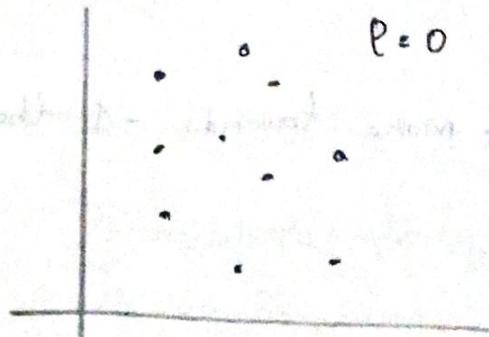
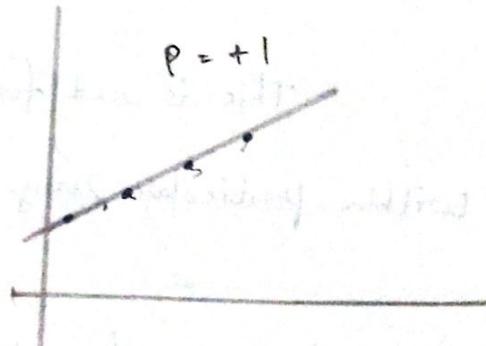
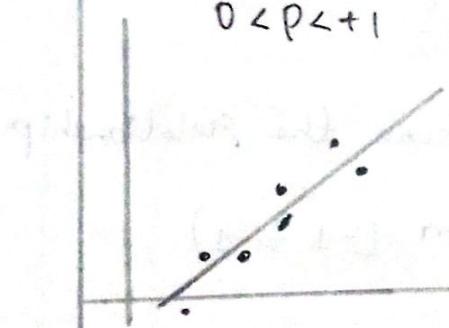
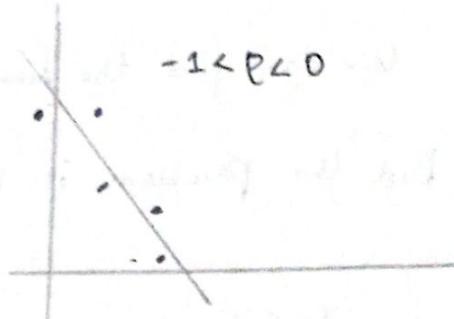
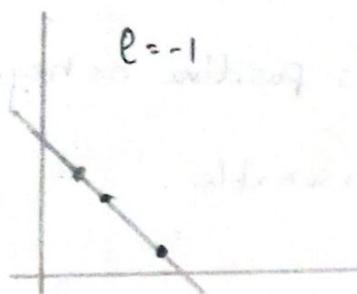
Formula,

$\sigma_x \rightarrow$  sd of feature x

$$\rho(x,y) = \frac{\text{Cov}(x,y)}{\sigma_x \sigma_y}$$

$\sigma_y \rightarrow$  sd of feature y

value = {-1 to 1} (Ans)

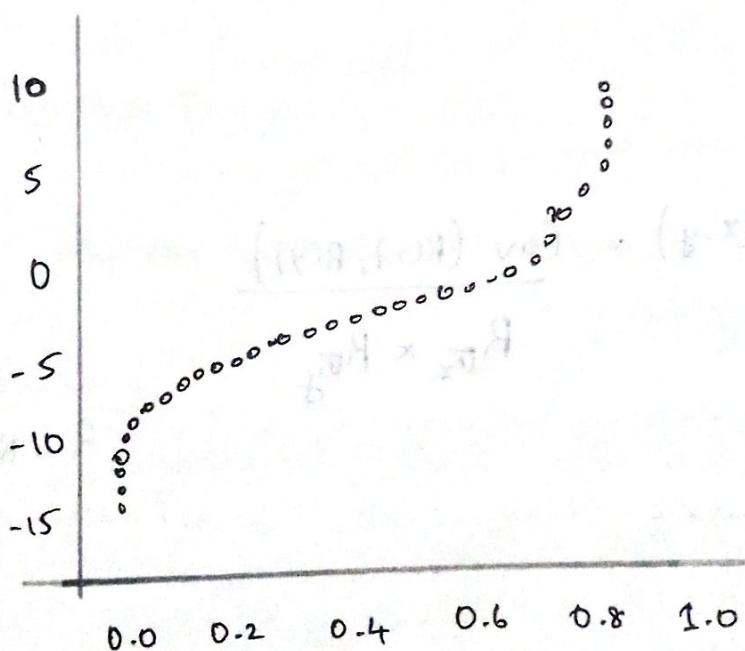


Does not follow straight line we say range from -1  
 $\rho < 0$  (or)  $0 < \rho < +1$  [Range]

### Disadvantage:

→ The Pearson Correlation Coefficient is only works well for Linear properties. Does not fit for Non-Linear properties. So, we jump into Spearman rank Correlation.

### Spearman's rank Correlation Coefficient:



In this Scenario,

$$\text{Spearman Correlation} = 1$$

$$\text{Pearson Correlation} = 0.88$$

Spearman Considered best because, it is non-linear Property. So works well and also show accurate Correlation.

Advantage:

- It is somewhat similar to Pearson Correlation,  
The difference is it works well for non-linear  
property as well.
- It also calculate the relation accurately. (i.e. One  
point overlap with some distance also calculate. But,  
Pearson cannot).

formula,

$$\text{Spear}(x,y) = \frac{\text{cov}(R(x), R(y))}{R_{Rx} \times R_{Ry}}$$

R → Rank

How to calculate?

Let Consider,

x	y	R(x)	R(y)
170	75	2	2
160	62	3	3
150	60	4	4
145	55	5	5
180	85	1	1

x → Height

y → weight

Based on the higher value, Rank will be assigned.

The formula only works with  $R(x)$ ,  $R(y)$ ,  $R_{Rx}$  and  $R_{Ry}$ . It neglect the original feature  $X$  and  $Y$  while performing.

The main thing of the Spearman Rank Correlation is that Capture non-linear properties as well.

PRACTICAL EXAMPLE: [T-test]

The Q/A is previous Z-test Q/A.

Program:

```
ages = [10, 20, 35, 50, 28, 40, 55, 18, 16, 55, 30, 25, 43, 18,  
30, 28, 14, 24, 16, 17, 32, 35, 26, 27, 65, 18, 43, 23,  
21, 20, 19, 70]
```

```
import numpy as np
```

```
ages_mean = np.mean(ages)
```

```
SampleSize = 10
```

```
age_sample = np.random.choice(ages, sample_size)
```

```
from Scipy.stats import ttest_1samp
```

```
ttest_1samp (age_Sample, 30) # 30 population mean
```

Output:

```
Ttest1sampResult (statistic = 0.31410216574, pvalue = 0.7606021)
```

Another Example: [T-Test]

# Ages of the college students (Population)

# 1 class Students Mean of all the ages .

```
import numpy as np
```

```
import Pandas as pd
```

```
import Scipy.stats as stats
```

```
import math
```

```
np.random.Seed(6)
```

```
School_ages = Stats.Poisson.rvs(loc=18, mu=35, size=1500)
```

```
ClassA_ages = Stats.Poisson.rvs(loc=18, mu=30, size=60)
```

# To find School age mean

School-age mean = np.mean(School-ages)

ttest\_1samp (class A-ages, Popmean = School-age mean)

# return (Statistic = -9.604795, Pvalue = 1.139027071016e-13)

# check the hypothesis

if P-Value < 0.05:

Print ("Accept H0")

else:

Print ("Reject H0")

Output: Accept H0 [Pvalue less than 0.05 because it is exponent  
of value]

PRACTICAL FOR CORRELATION:

import Seaborn as sns

df = sns.load\_dataset('iris')

df.head()

```
df.corr() # return correlation for every two features.
```

# To visualize Correlation Use Pairplot .

```
Sns. Pairplot (df)
```