

6th DAY LIVE SESSION

- ① CHI SQUARE ✓
- ② Covariance ✓
- ③ Pearson Correlation Coefficient ✓
- ④ Spearman Rank Correlation ✓
- ⑤ Practical Implementation
 - Z-test, t-test, chi square test
- ⑥ F Test (ANOVA)

CHI SQUARE TEST

- ① Chi Square Test claims about population proportions

(non-parametric test usually occurs wrt population proportion, when we are given some proportions of data, we can't specifically use parametric test)

If it is a non parametric test that is performed
on Categorical (nominal or ordinal) data.

why chi square test is used

- Q) In the 2000 Indian Census, the age of the individual in a small town were found to be the following:

Less than 18	18-35	>35
20%	30%	50%

In 2010, age of n=500 individuals were sampled. Below are the results

<18	18-35	>35
121	288	91

Using $\alpha=0.05$, would you conclude the population distribution of ages has changed in the last 10 years?

Ans)

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

{ Population } $\frac{2000}{\text{_____}}$

Expected

<18	$18 - 35$	>35
121	288	91

$n=500$

Observed

<18	$18 - 35$	>35
100	500×0.3	500×0.5
$500 \times 20/100$	$= 150$	$= 250$
121	288	91

Expected

(We are multiplying as we want to find expected here wrt above 2000 population)

<18	$18 - 35$	>35
121	288	91
100	150	250

Observation

Expected

{ Chi Square table }

① H_0 = The data meets the distribution 2000 census (null hypothesis)

H_1 = The data does not meet " " " " (alternate hypothesis)

② $\alpha = 0.05$ (95% C.I)

$df = 2, \alpha = 0.05$

(degree of freedom)

③ Degree of freedom = $n - 1 = 3 - 1 = 2$ // n = no. of categories, here 3 ($<18, 18-35, >35$)

④ Decision Boundary

Now we will look in chi square table for given $df = 2$ and $\alpha = 0.05$, we get 5.99



It is a 2 tail test because data may be less than our distribution or it may be more

If χ^2 is greater than 5.99 reject H_0

(Now we will find chi square and we will reject null hypothesis, if its value is greater than 5.99)

5) Calculate Test Statistics

$$\chi^2 = \sum \frac{(f_0 - f_e)^2}{f_e} \quad // f_0: \text{observed}, f_e: \text{expected}$$

$$= \frac{(121 - 100)^2}{100} + \frac{(288 - 150)^2}{150} + \frac{(91 - 250)^2}{250} \\ \approx 132.94$$

$$\chi^2 = 132.94 > 5.99 \quad \left\{ \begin{array}{l} \text{Reject the Null} \\ \text{Hypothesis} \end{array} \right.$$

If we get 0.11 so

$$0.11 > 0.05$$



Accept the Null

Reject the hypothesis

$$\boxed{0.11 \times 0.05}$$

$$\alpha = 0.05 \quad \boxed{0.01 \quad 0.10} \quad \left\{ \begin{array}{l} \text{Domain} \\ \text{the null hypothesis} \end{array} \right.$$

$$0.002 < 0.05$$

↓ Reject

It is decided by domain expert



I Was Correct Here

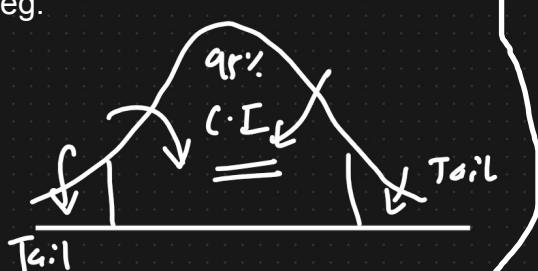
If $\{ \begin{array}{l} \text{P-value} < \text{Significance value} \\ \downarrow \end{array} \}$

Reject the Null Hypothesis.

else

{Accept the Null Hypothesis}

eg:



$$P = 0.11 > 0.05$$

↓ Accept

$$P = \boxed{0.002} < 0.05$$

↓ Reject the Null Hypothesis

② Covariance

X

y

Eg 1:

Weight

Height

50

160

60

170

70

180

75

181

X↑ Y↑

X↓ Y↓

here when x is increasing y is also increasing and when x is decreasing then y is decreasing

Eg 2:

No of hours
Study

play

2

6

3

4

4

3

X↑ Y↓

X↓ Y↑

here when x is increasing y is decreasing and when x is decreasing then y is increasing

Quantity relationship between X & Y

How we can quantify the relationship using numbers: We use covariance

Covariance Formula:

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

+ve

↓

positive correlation

If

X↑ Y↑
X↓ Y↓

↓

X↑ & Y

If

X↓ Y↑
X↑ Y↓

negative correlation

=

+ve or -ve

or

or

Eg 1:

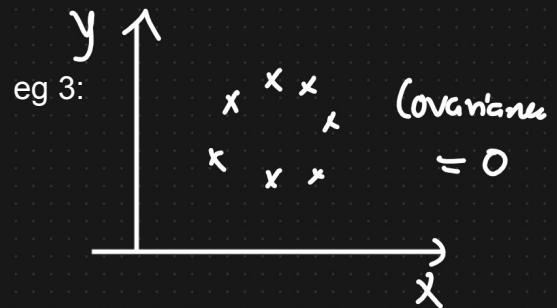
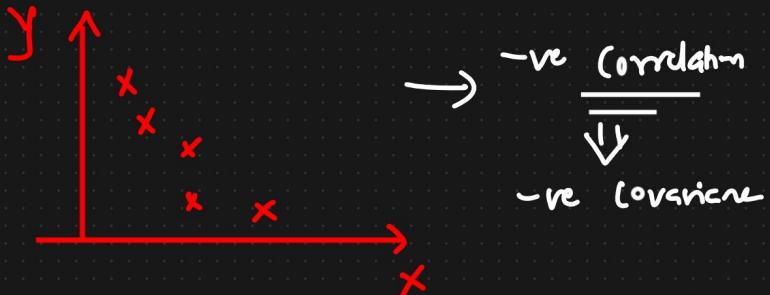


→ +ve Correlation



Suppose we have a dataset which looks like this, this will have +ve correlation as here x is increasing then y is also increasing and when x decr then y decr. Even if we apply above formula we will get +ve value i.e +ve correlation

eg 2:



Disadvantage of Covariance

① Positive OR Negative ✓

It can be +ve or -ve but it doesn't tell difference between 100 and 1000, i.e. magnitude can't be measured

$$\begin{array}{r} +100 \\ -200 \\ \hline -2000 \end{array}$$

+1000
= f Direction

$$-2000$$

Covariance can only measure the directional relationship between two assets. It cannot show the strength of the relationship between assets. The correlation coefficient is a better measure of that strength.

② Pearson Correlation Coefficient

$$(-1 \text{ to } 1)$$

The more towards +1 more positively correlated

The more towards -1 more negatively correlated

(it restricts the value between -1 to 1)

Formula:

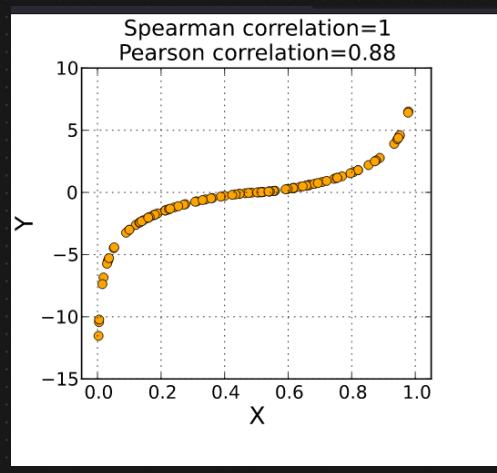
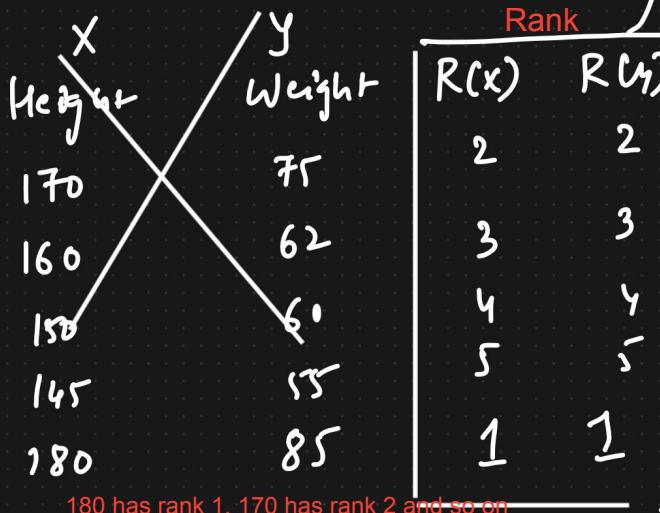
$$f(x,y) = \frac{\text{Cov}(x,y)}{\sigma_x \sigma_y} = \left\{ -1 \text{ to } 1 \right\}$$

Pearson coefficient is good for linear graphs, i.e. when all the points lie on same line, it will give 1 or -1 as correlation

Spearman Correlation:

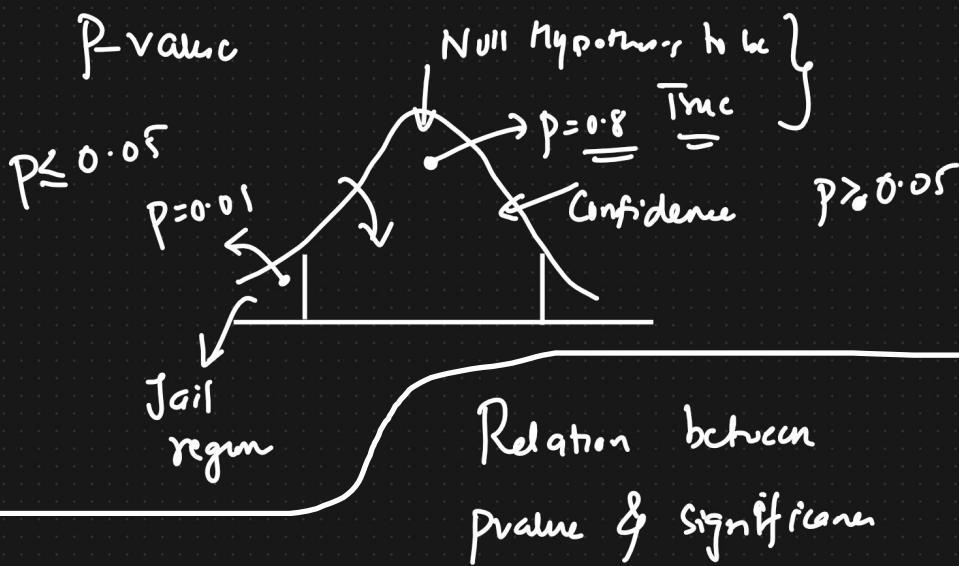
$$\text{Spear}(x,y) = \frac{\text{Cov}(R(x), R(y))}{\sqrt{R(x)} \times \sqrt{R(y)}}$$

But when graph is non-linear and all points are increasing, pearson coefficient will give close to 1 but not 1, here spearman works well and can give +1 correlation:



Non linear properties are captured by spearman rank correlation, that's why we use it

If $P_{\text{value}} \leq 0.05$ → Reject the Null Hypothesis
 probability \downarrow
 $\alpha = 0.05$
 $P \geq 0.05 \rightarrow$ Accept the Null Hypothesis



P_{value} \leftarrow Significant $\boxed{C.I}$
 \hookrightarrow Reject the Null Hypothesis
 P_{value} \leftarrow Accept the Null Hypothesis

(To understand better check next pdf)

① P Value And Significance value