

24-01-25

Day - 6

Page No.

Date

Today's topics

1. Chi Square
2. Covariance
3. Pearson Correlation coefficient
4. Spearman and 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 proportion  
It is a non parametric test that is performed on categorical (Nominal or Ordinal) data

Q. In the 2000 Indian census, the ages of the individual in a small town were found to be the following:

1. less than 18
2. 18-35
3. > 35

Ans:-

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

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

<18	18-35	>35
121	288	91

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

Answer:-

<18	18-35	735	2000 Population
20%	30%	50%	Expected

<18	18-35	735	$n = 500$
121	288	91	Observed
$500 \times 0.2 = 100$	$500 \times 0.3 = 150$	$500 \times 0.5 = 250$	Expected
100	150	250	

<18	18-35	735	
121	288	91	Observation
100	150	250	Expected

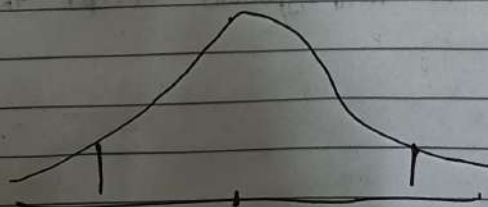
1.  $H_0$  = The data meets the distribution 2000 census

$H_1$  = The data doesn't meet the distribution 2000 census

2.  $\alpha = 0.05$  (5% Confidence Interval)

3. Degree of freedom  $n - 1 = 3 - 1 = 2$

4. Decision Boundary



2 tail test

If  $\chi^2$  is greater than 5.99 reject  $H_0$



5. Calculate test statistics

$$\chi^2 = \sum \frac{(P_o - P_e)^2}{P_e}$$

$P_o$  :- observed  
 $P_e$  :- Expected

$$\chi^2 = \frac{(121 - 100)^2}{100} + \frac{(65 - 40)^2}{150} + \frac{(91 - 20)^2}{200}$$

$$\chi^2 = 4.41 + 126.96 + 101.124$$

$$\chi^2 = 232.494$$

$$\chi^2 = 232.94 > 5.99$$

Reject the Null hypothesis

Conclusion

p-value < Significant value

↓

In ~~both~~ <sup>greater than</sup> Reject the Null hypothesis (0.1)

In less than Accept the Null hypothesis

$$P = 0.1 > 0.05$$

↓

Reject the Null hypothesis

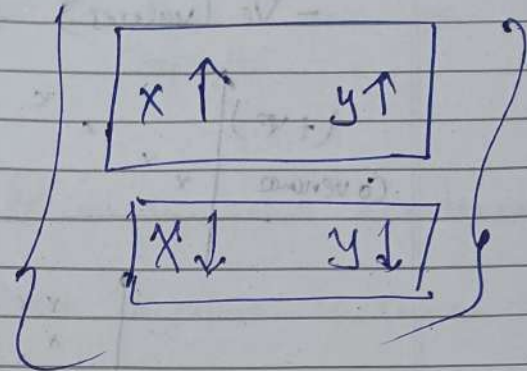
$$P = 0.002 < 0.05$$

↓

Accept the Null hypothesis

## \* Covariance :-

X	y
weight	height
50	160
60	170
70	180
75	181



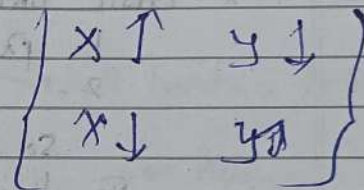
Ex:-

No. of hrs  
Study

No. of hrs  
play

2  
3  
4

6  
4  
4



Quantity relationship between  $x$  &  $y$   
But that time we use

Covariance

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

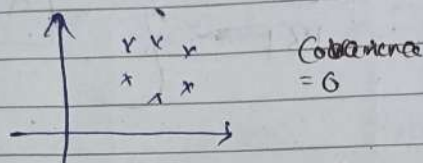
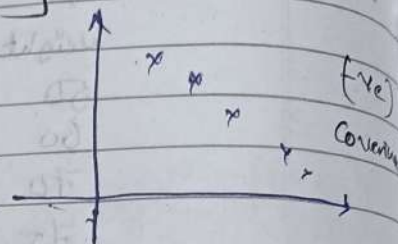
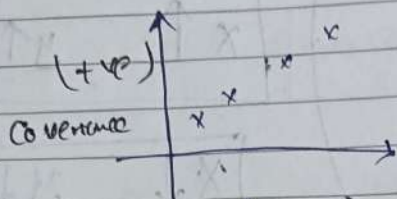
For

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



+ve (values)  $\Rightarrow$   $\begin{matrix} x \uparrow & y \uparrow \\ x \downarrow & y \downarrow \end{matrix}$  } Positive Correlation

-ve (values)  $\Rightarrow$   $\begin{matrix} x \downarrow & y \uparrow \\ x \uparrow & y \downarrow \end{matrix}$  } Negative Correlation



\* Basic disadvantage of Covariance

1. Positive or negative

2. ~~the~~ the value is unlimited (+ve) or (-ve)

So the we are difficult to solve.

In that time we use

\* Pearson Correlation Coefficient

$[-1 \text{ to } 1] \rightarrow$  The more towards  $+1$  more positive correlation

$\rightarrow$  The more towards  $-1$  more negative correlation

Formula :-

$$\rho(x, y) = \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y} \rightarrow [-1, +1]$$

Formula: Pearson

$$\rho_{\text{Pearson}}(x, y) = \frac{\text{Cov}(R(x), R(y))}{R_{\text{max}} \times R_{\text{max}}} = \frac{\text{Cov}(R(x), R(y))}{R_{\text{max}} \times R_{\text{max}}}$$

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

Page No.	
Date	

Q. why do we use Pearson correlation?  
 → It captures "non-linear" properties

Probably  $P \leq 0.05 \rightarrow$  Reject the null hypothesis

5% probably the null hypothesis is correct  
 $\alpha = 0.05$

$P \geq 0.05 \rightarrow$  Accept the null hypothesis

