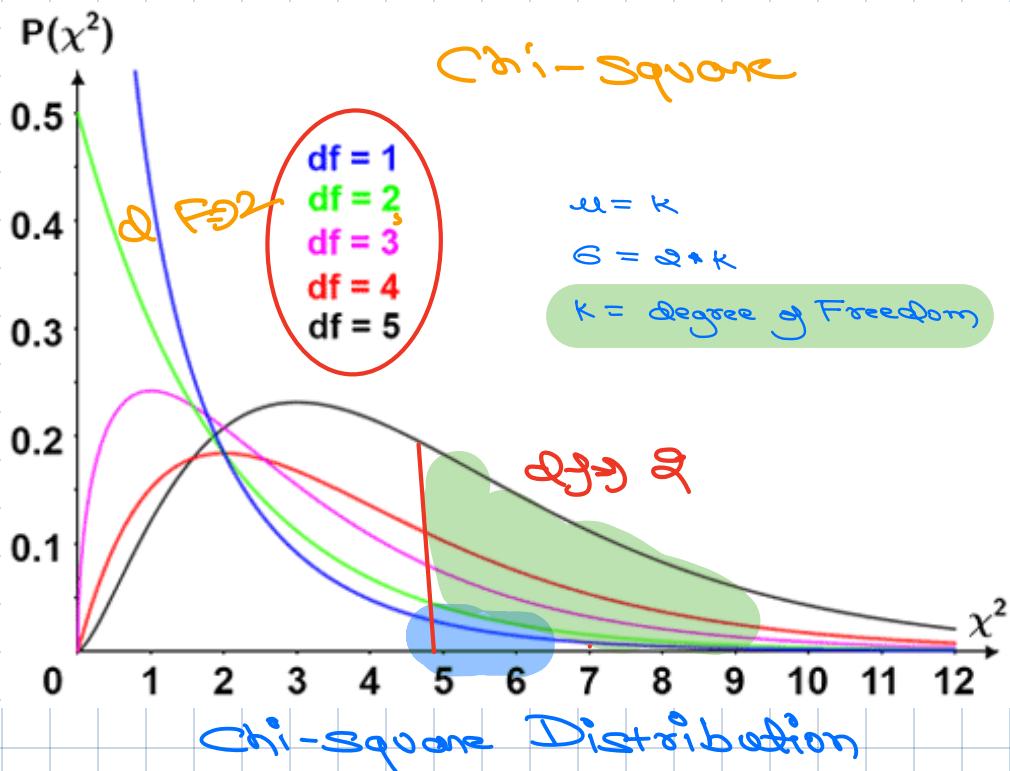


## \* Post Read



1. **Definition:** The chi-square distribution is defined as the distribution of a sum of the squares of  $k$  independent standard normal random variables (where  $k$  is the degrees of freedom).

2. **Degrees of Freedom:** The shape of the chi-square distribution depends on the degrees of freedom (df). As the degrees of freedom increase, the distribution becomes more symmetric and resembles a normal distribution.

### 3. Properties:

- It is positively skewed, especially for lower degrees of freedom.
- Its mean is equal to the degrees of freedom ( $k$ ).
- Its variance is  $2k$ .
- The distribution is used only with positive values, as it represents squared values.

### 4. Applications:

- **Goodness-of-fit tests:** To determine if an observed distribution matches an expected one.
- **Test for independence:** In contingency tables, to see if two categorical variables are related.
- **Variance estimation:** To test hypotheses about the variance of a normally distributed population.

# Agenda

① Paired T-test

② DOF

③ Chi-square Test

Goodness of Fit

Test for independence

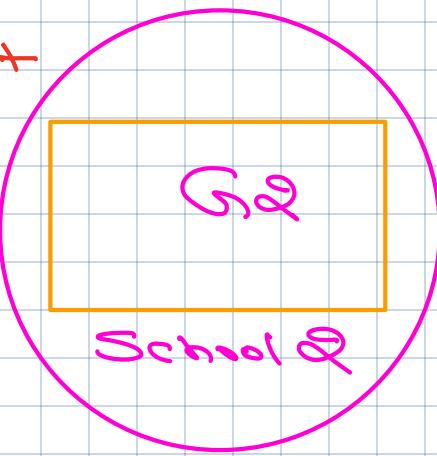
④ Assumption of Chi-square Test

Paired T-test

Independent T-test



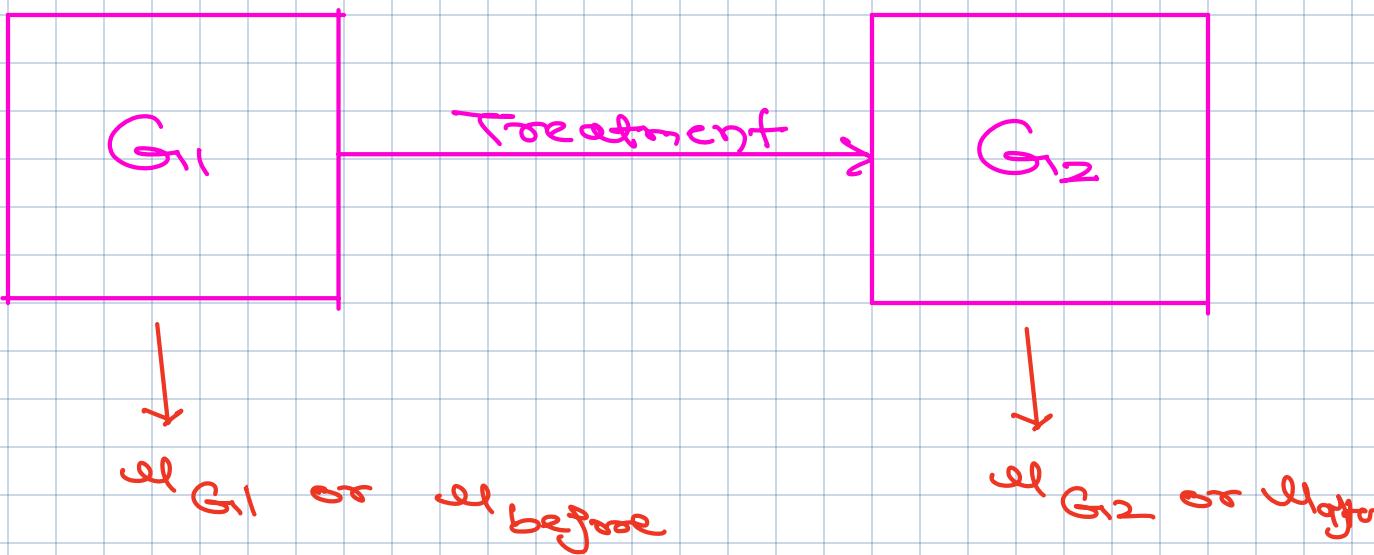
School 1



$$H_0: \mu_{IQ-G1} = \mu_{IQ-G2}$$

$$H_a: \mu_{IQ-G1} \neq \mu_{IQ-G2}$$

## Paired T-test



$$\text{diff} \circ \mu_{\text{after}} - \mu_{\text{before}}$$

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

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

Claim: Problem Solving Session helps Student get better Marks

`df_ps.head()`

	<code>id</code>	<code>test_1</code>	<code>test_2</code>
0	0	40	38
1	1	49	44
2	2	65	69
3	3	59	63
4	4	44	43

before PS      After PS

$$H_0: \mu_{\text{score}_1} = \mu_{\text{score}_2}$$

$$H_a: \mu_{\text{score}_1} \neq \mu_{\text{score}_2}$$

(Two tailed)

$$H_a: \mu_{+1} < \mu_{+2}$$

## Recap

Setup for Relationship of Variables :

- ① Numerical vs Categorical ( $\leq 2$  Categories)  
Z-test, t-test, proportions Z-test
- ② Numerical vs Categorical ( $> 2$  Categories)  
Anova test (In upcoming Sessions)
- ③ Numerical vs Numerical  
Correlation (In upcoming Sessions)
- ④ Categorical vs Categorical  
Chi-Square Test

Chi-Square Test

## Degree of Freedom

- ⑨ Number of independent values that are free to vary when calculating a statistic
- ⑨ It tells us about flexibility or choices we have when trying to make Statistical Decision

Ex-2: If genre of movies for 7 weeks

$$DOF = ?$$

Week 1 : 7 choices

Week 2 : 6 choices

Week 3 : 5 choices

$$DOF \geq 6$$

Week 4 : 4 choices

Week 5 : 3 choices

Week 6 : 2 choices

Week 7 : 1 choice → Only 1 genre available

Ex2 :

## Numerical Data

Empid	Salary
1	36
2	35
3	?

$$\text{mean} = 35$$

$$D.O.F = 2$$

→ 34

Ex2 :

## Categorical Data

Original Data:

Empid	Gender
1	M
2	F
"	g
"	g



Transformed Data:

Empid	M	F
1	1	0
2	0	1
"	0.5	0.5
"	0.5	0.5

Do we need both columns?

3 Categories  $\rightarrow n-1 \Rightarrow 3-1 \Rightarrow 2$

HW:

Empid	Salary
1	36
2	35
3	?
4	30

$$\text{mean} = 37$$

$$D.O.F \ 4-1 \Rightarrow 3$$

3rd Value?

## General Rule of DoF:

- - - - - 3

$\bar{x}$  = number

- Given a set of  $N$  numbers
- If you know average of  $n$ -numbers, what is minimum number of Data-points you need to know to Extract Full Set

$$\text{DOF} \Rightarrow n-1$$

for 1-D array with length  $n$

Ex-3

H	W
73	85
68	73
74	96
71	82
62	70
71	81.2
$n_1$	

$$\text{DOF} = n_1 - 1 + n_2 - 1$$

$$\Rightarrow n_1 + n_2 - 2$$

$$\Rightarrow 2n - 2$$

For arrays with length  $n_1$  and  $n_2$  and means  $\bar{x}_1$  and  $\bar{x}_2$   $\Rightarrow$

Ex-4

## 2D Contingency Table

Sachin century and Victory

Won

	F	T	
L	2	?	314
T	?	?	46
	176	184	360

Contingency Table

$D.o.F = ?$

	F	T	
L	314-184	154	314
T	16	30	46
	176	184	

$D.o.F = 1$

$(n_1-1) \times (n_2-1)$

Q3  
Q4  
Q5  
Q6

$\therefore (n_1-1 + n_2-1)$

Ex-5: Given data of Regional Support for 4 Politicians A, B, C, D from Survey of 3 cities X, Y, Z

	A	B	C	D	Total
X	90	60	104	?	349
Y	30	42	51	?	151
Z	?	?	?	?	150
Total	150	150	300	150	650

# Rows = 3  
# Columns = 4

df = 6

$$(\# \text{rows} - 1) \times (\# \text{columns} - 1)$$

$$(3-1) \times (4-1) \Rightarrow 2 \times 3 = 6$$

Why D.F.?

- ① D.F. is crucial for Chi-Square tests as it is important parameter in Chi-Square distribution i.e. Dist of Test Statistic under Null Hypothesis

② It controls Variability in data.

③ Higher D.o.F leads to higher critical values requiring larger test statistics to reject Null Hypothesis for given ( $\alpha$ ).

### Chi-Square Goodness of Fit Test

Given a categorical Variable, you can test if there is any significant difference b/w observed Values and Expected Values

Example: Coin Toss : 50 trials

	H	T
Observed	28	22
Expected	25	25

$$D.O.F \Rightarrow 1$$

For  $H_0 \rightarrow$  Expected Heads  $\Rightarrow 25$   
Expected Tails  $\Rightarrow 25$

## Observation:

We can observe some deviation in  
Actual vs Expected

$$\frac{(O(H) - E(H))^2}{E(H)}$$

diff b/w Heads

$$\frac{(O(T) - E(T))^2}{E(T)}$$

diff b/w Tails

$$\frac{(Q5 - 25)^2}{25}$$

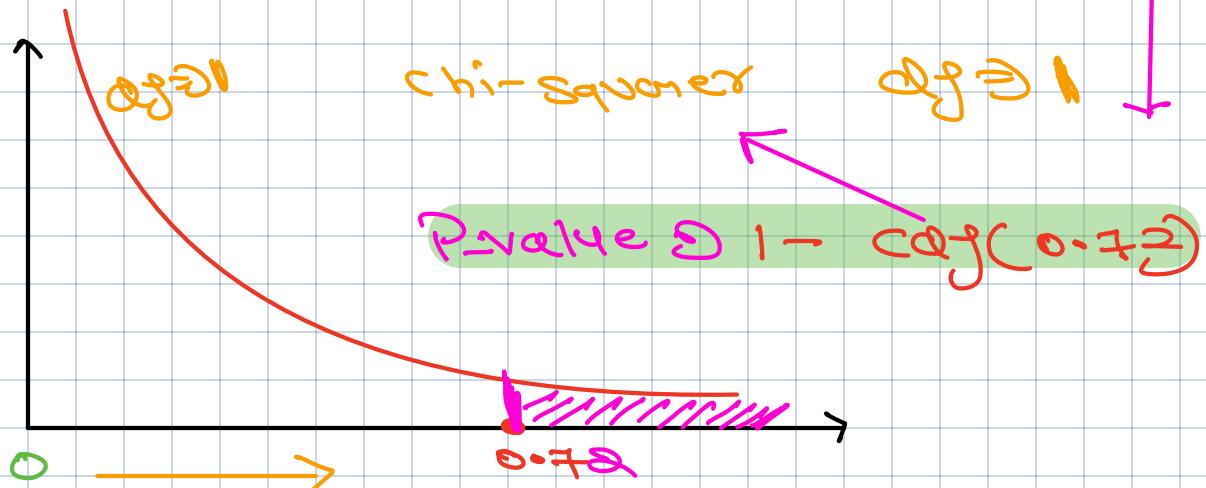
+

$$\frac{(Q2 - 25)^2}{25} \rightarrow \frac{18}{25}$$

Chi-square Statistic  $\Rightarrow 18/25 \approx 0.72$

DOF 9 - 7

What would chi-square statistic look like  
when plotted under  $H_0$  ?



P-val

Chi- $\chi^2$  stat  $\Theta$   
 $\chi^2$

$$\sum_{i=1}^c \frac{(O(c) - E(c))^2}{E(c)}$$

$$P\text{-value} \Theta 1 - \text{ChiCDF}(\text{chi2\_stat}, \text{dof})$$


Q Can n change in Between?

No and DOF can't change b/c  
Experiment

Q How do we conclude the result i.e. coin  
Biased or Not?

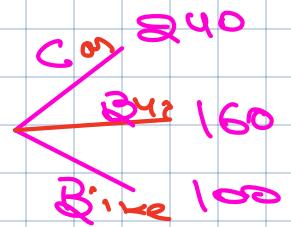
$$\alpha \Theta 0.05$$

If P-value <  $\alpha$ :

reject  $H_0$

\* To Do H-W

Python Example for Coin Toss



Step 1: Setup Hypothesis

$H_0$ :

$H_a$ :



Step 2: Pick distribution and  $\alpha$

Step 3: Perform Test and calculate test statistics

Step 4: Compare p-value and reject / fail to reject Null Hypothesis

## Chi-Square Test of Independence

- Given two variables you want to test if they are related or independent

### Marketing Firm Example

- Imagine you are running a marketing campaign.
- There are two modes of shopping
  - Offline
  - Online

Goal: Increase Online sales by running Targetted Campaign/ads

### Consultancy Claim:

- Focus on Women
- Females have higher chance of shopping online than males

In other words

There is relationship between  
'Buying Pattern and Gender'

```
graph TD; BP[Buying Pattern] --> Offline[Offline]; BP --> Online[Online]; G[Gender] --> M[M]; G --> F[F]
```

Can we test and Evaluate this?

Gender vs Purchase

	M	W	
Offline	524	72	599
Online	206	102	308
	730	174	907

$H_0$  : Gender and preference are independent

$H_a$  : Gender and preference are dependent

① Total % of population that prefers

① Offline ②  $\frac{599}{907} \rightarrow 66\%$

② Online  $\frac{308}{907} \rightarrow 34\%$

If gender has no impact

① Prefer offline

Expect men % 66%  $\rightarrow 66\% \cdot 733 = 484$

Expect Women % 66%  $\rightarrow 66\% \cdot 174 = 115$

② Prefer Online

Expect men % 34%  $\rightarrow 34\% \cdot 733 = 249$

Expect Women % 34%  $\rightarrow 34\% \cdot 174 = 59$

Chi- $\chi^2$  stat  $\rightarrow$

$\chi^2$

$$\sum_{i=1}^c \frac{(O(c) - E(c))^2}{E(c)}$$

P-value  $\rightarrow 1 - \text{chi.sq}(chisq, \text{dof})$

Chi-sq stat

$\chi^2$

②

$$\begin{aligned} & \frac{(527 - 484)^2}{484} \\ & + \\ & \frac{(-72 - 115)^2}{115} \\ & + \\ & \frac{(206 - 249)^2}{249} \\ & + \\ & \frac{(102 - 59)^2}{59} \end{aligned}$$

\*

$$P\text{-value} \Leftrightarrow 1 - \text{Chi-sq. stat}(\chi^2_{df})$$

A marketing manager wants to determine if there is a relationship between the type of advertising (online, print, or TV) and the purchase decision (buy or not buy) of a product.

Var 2

Var 1

The manager collects data from 300 customers and records their advertising exposure and purchase decisions.

What statistical test should the manager use to analyze this data?

# Assumptions

# Chi<sup>2</sup> test

## 1. Random Sample

- The data points for each group in your analysis must have come from a simple random sample.
- This is important because if your groups were not randomly determined then your analysis will be incorrect.
- In statistical terms this is called bias, or a tendency to have incorrect results because of bad data.

## 2. Variables are categorical

## 3. Mutually Exclusive Groups

- The two groups of your categorical variable should be mutually exclusive.
- For example, if your categorical variable is hungry (yes/no), then your groups are mutually exclusive, because one person cannot belong to both groups at once.

## 4. Observations are independent

- Each of your observations (data points) should be independent.
- This means that each value of your variables doesn't "depend" on any of the others.
- For example, this assumption is usually violated when there are multiple data points over time from the same unit of observation (e.g. subject/customer/store).
- Because the data points from the same unit of observation are likely to be related or affect one another.
- A different test must be used if the researcher's data consists of paired samples, such as in studies in which a parent is paired with his or her child.

## 5. The value of the cell expected should be **5 or more** in at least 80% of the cells, and no cell should have an expected of less than one