

Lecture 5: Chi_Square Test

#Agenda

- ① Types of Test
- ② Degrees of Freedom
- ③ Chi-squared Goodness of fit Test
- ④ Chi-squared Test for independence

Types of Test

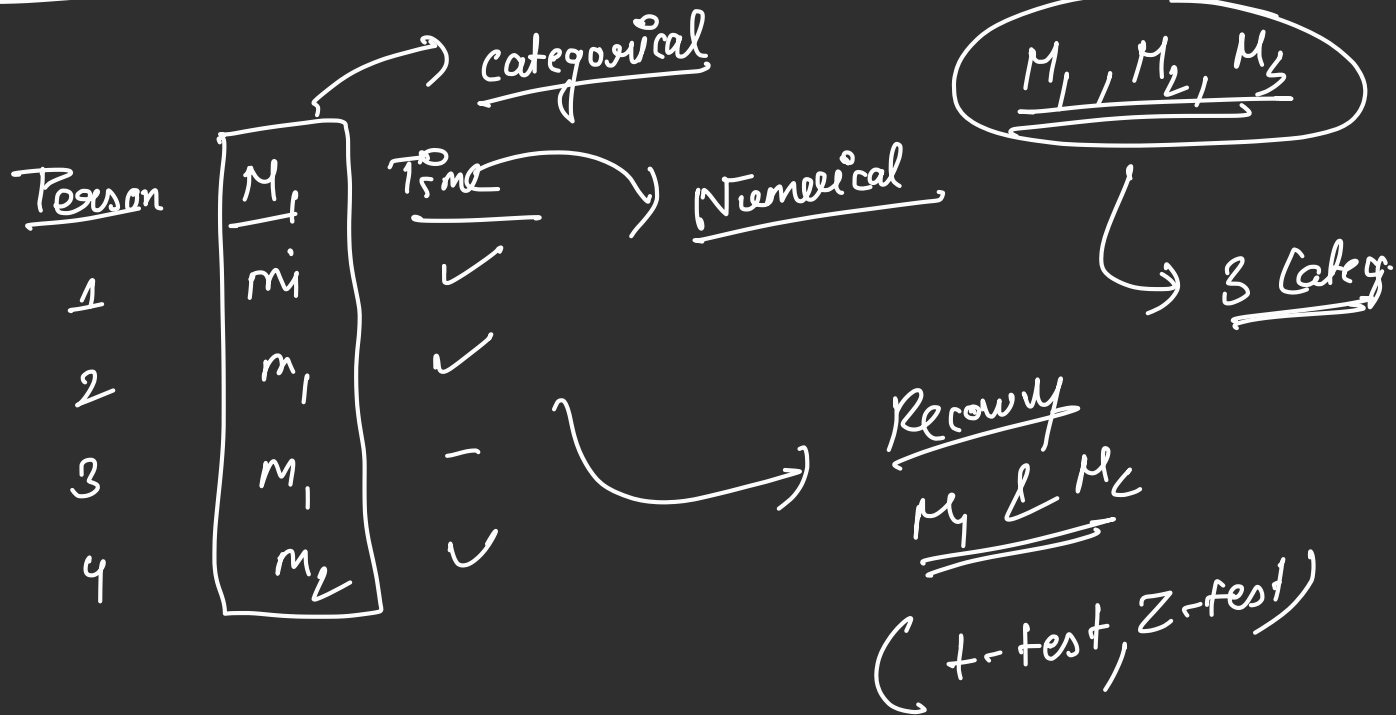
Sofar

Numeric Vs Categorical \rightarrow 2 (T-test)
 \rightarrow >2 (ANOVA)

Categorical Vs Categorical \rightarrow Chi square test

Numeric Vs Numeric \rightarrow Correlation

Numerical Vs Categorical



Degrees of freedom (DOF)

① Setup I: Salary data

$$P_1 \rightarrow 35L$$

$$P_2 \rightarrow ? \text{ y}$$

$$P_3 \rightarrow ? (x)$$

$$\frac{x + 30 + 35}{3} = 35$$

$$\underline{\underline{x = 40}}$$

$$\text{Avg} \rightarrow 35L$$

$$\frac{35 + x + y}{3} = 35$$

② Schp II: Salary data

$$P_1 \rightarrow 35L$$

$$P_2 \rightarrow 36L$$

$$P_3 \rightarrow ? (x)$$

$$P_4 \rightarrow 30L$$

$$\text{Avg} \rightarrow 37L$$

Q) What is the salary of the 3rd person?

$$\underline{\underline{x = 47}}$$

$$\frac{35 + 36 + x + 30}{4} = 37$$

$$\underline{\underline{x = 47}} \checkmark$$

Generalize

$$\underline{\underline{n-1}}$$

3 person salary should be present \rightarrow

$$\{ \text{DOF} = \underline{\underline{n-1}} \}$$

③ Setup III: Height & Weight

$$n_1, n_2 \leq n$$

	H	W
	73	85
	68	73
	74	
		82
	62	70
Avg	71	81.2

$$\begin{aligned} H &\rightarrow n_1 - 1 \\ W &\rightarrow n_2 - 1 \end{aligned} \quad \left. \vphantom{\begin{aligned} H &\rightarrow n_1 - 1 \\ W &\rightarrow n_2 - 1 \end{aligned}} \right\} \text{ for } \text{ho} \quad \left(\frac{2(n-1)}{\quad} \right) \times$$

$$\underbrace{(n_1 - 1)}_{\text{DOF}_1} + \underbrace{(n_2 - 1)}_{\text{DOF}_2}$$

$$\underline{\text{DOF} = (n_1^\checkmark + n_2^\checkmark - 2)}$$

$$M_1: [36, 26, 32, \dots]$$

$$M_2: [25, 32, 36, \dots]$$

summary

For 1 sample: $\text{DOF} \Rightarrow (n-1)$

For 2 sample: $\text{DOF} \Rightarrow (n_1 + n_2 - 2)$

④ Setup IV: Sachin (Century Vs Victory)
Win

		F	T	
Cent	F	160		314
	T			46
		176	184	360

Row sum

col sum

grand Total

(DOF \rightarrow 1)

$$(1 \times 1) = \underline{\underline{1}}$$

⑤ Setup V: Politicians

4 politicians: A, B, C, D
3 cities: X, Y, Z

$$(\# \text{row} - 1) * (\# \text{col} - 1)$$

\Rightarrow Dof \Rightarrow In this case

$$(3-1) * (4-1) \\ 2 * 3 = \underline{\underline{6}}$$

$$\underline{\underline{\text{Dof} = 6}}$$

	A	B	C	D	
X	90	60	104		349
Y	30	50	51		151
Z					150
	150	150	200	150	650

Two arrays
What will be the DOF?

Array 1 } $(n_1 - 1)$
Array 2 } $(n_2 - 1)$

$$\Rightarrow \underbrace{(n_1 - 1) + (n_2 - 1)}_{\Rightarrow} \quad \checkmark$$

$$M_1 = [62, \dots] \rightarrow \text{Nerve}$$

$$M_2 = [\dots] \rightarrow \text{Nerve}$$

Medice ReT

$(M_1) - 62$
 $M_1 - 63$
 $M_1 -$

} mean()

M_2 $() \rightarrow \text{mean()}$

Coin Toss Example

Objective: Check if a given coin is fair

No of tosses = 50 Times

	H	T	
Expected	25	25	50
Actual	28	22	50

↙

↗

↘

Under H_0 (Coin is fair)

Guarantee !! \rightarrow 25
 \rightarrow 25

	H	T	
Expected	25		50
Actual		25	50
	χ	χ	

✓
DOF = 1

DOF = 2

No of times = 50 Times

✓
 $DOF = (\text{#row} - 1) \times (\text{#col} - 1)$

$\Rightarrow (2-1) \times (2-1)$

$\Rightarrow (1) \checkmark$

No of times = 50 Times

	H	T	
Expected	25	25	50
Actual	28	22	50

Under $H_0 \Rightarrow (28 - 25) = +3$

$(22 - 25) = -3$

$+3 - 3 = 0 \checkmark \checkmark$

Case I: $(28-25) + (22-25)$

$\Rightarrow +3 + (-3) = 0$ { the value becomes zero }

Case II: $(28-25)^2 + (22-25)^2$

$\Rightarrow 9 + 9 = \underline{\underline{18}}$ (we solved the above problem, but it doesn't work when there are many trials/sample)

Case III:

for $n=1000$,

$(625-500)^2 + (375-500)^2 = \checkmark$

(This value is more as the trials keep increasing)

	H	T
Exp	500	500
Actual	650	350

$n = 1000$

✓ $(650 - 500)^2 + (350 - 500)^2 = \uparrow$

Case IV (Final Solution)

	H ✓	T ✓
Expt	25	25
Actual	28	22

$$\Rightarrow \frac{(28-25)^2}{25} + \frac{(22-25)^2}{25} =$$

✓

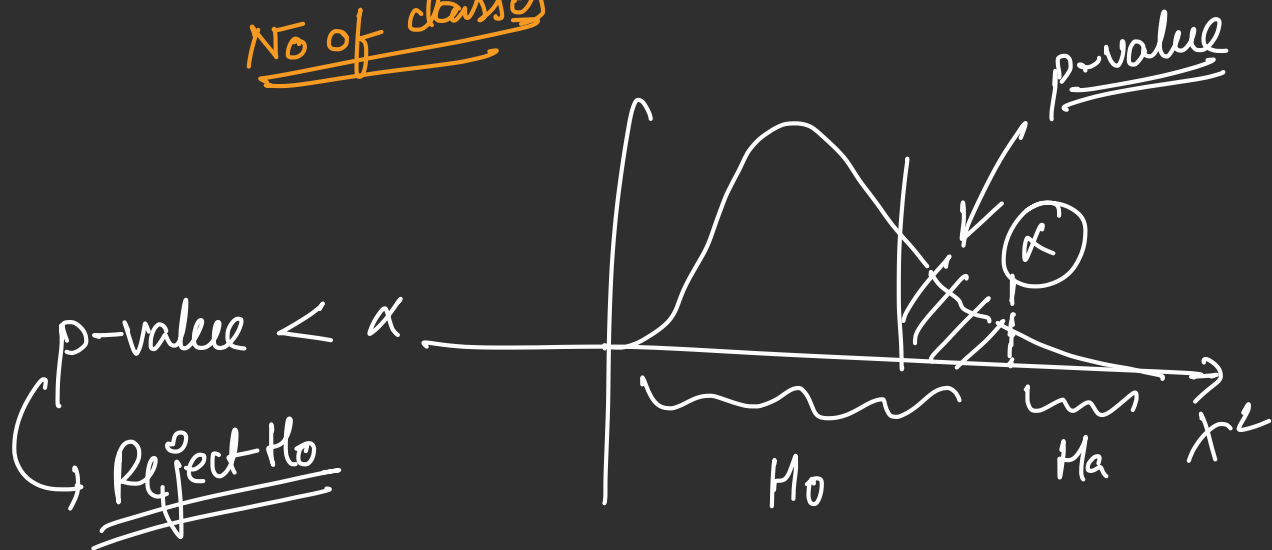
$$\chi^2 = \sum_k \frac{(O - E)^2}{E}$$

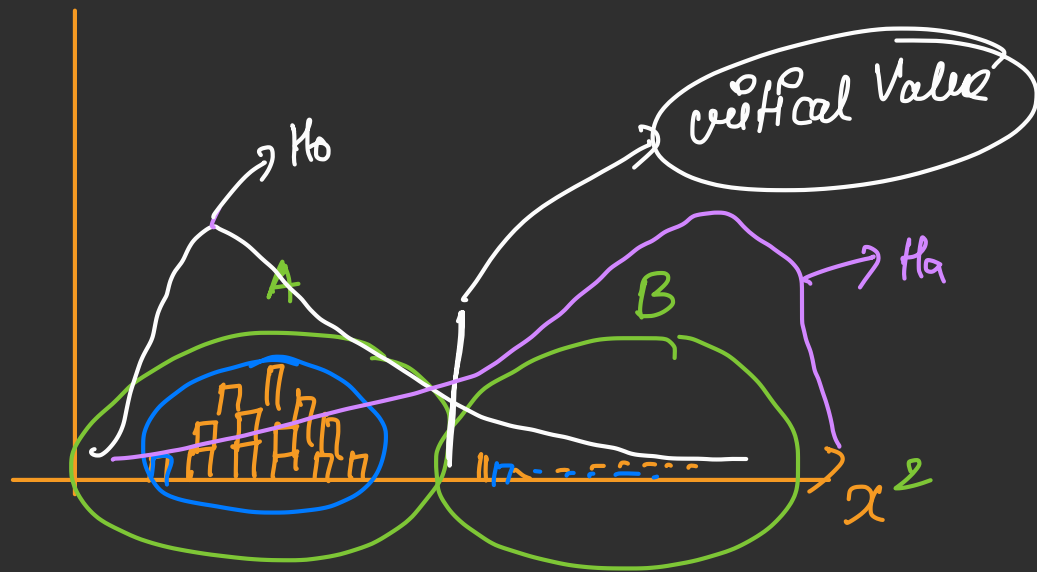
T-statistics

$$\chi^2 = \sum_k \frac{(O - E)^2}{E}$$

↓
No of classes

$\begin{cases} O: \text{observed value} \\ E: \text{Expected} \end{cases}$

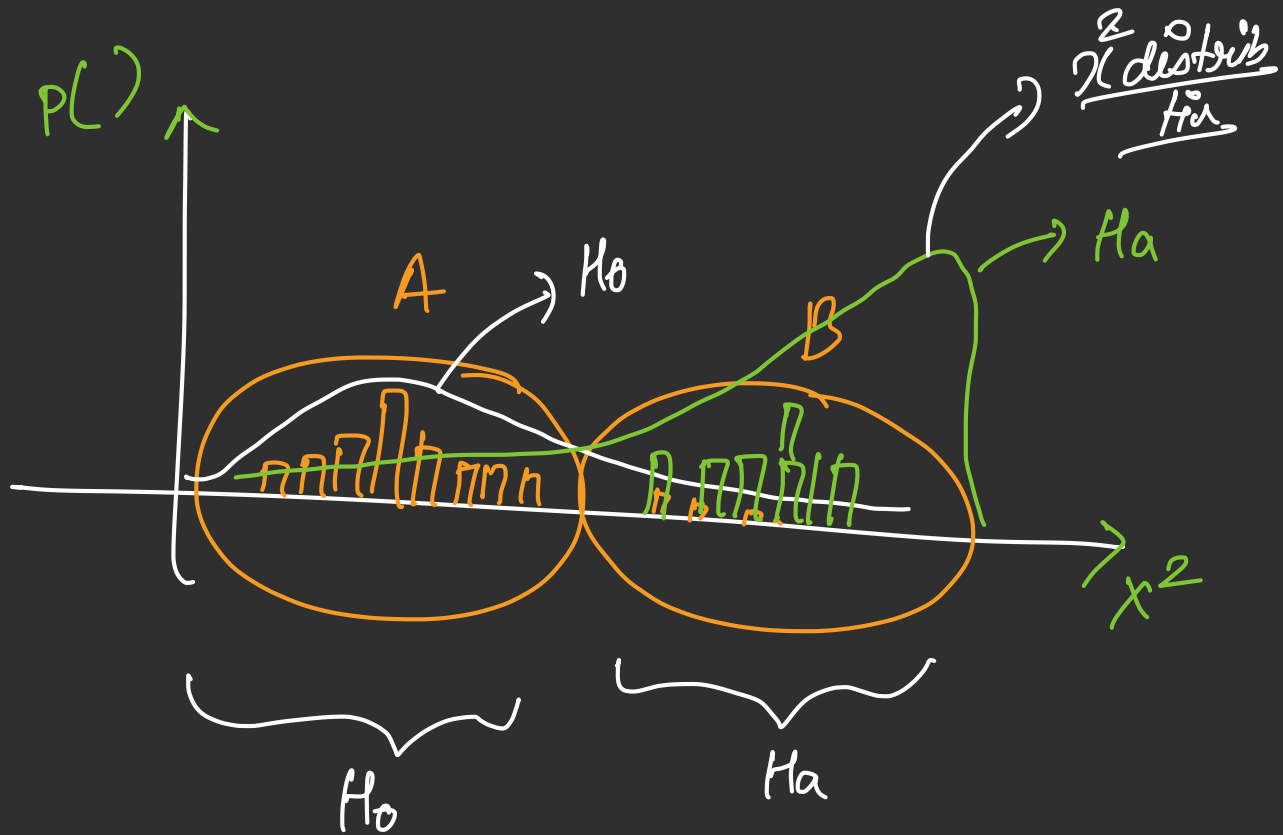




$\chi^2 \downarrow$
0 (smaller)
fair coin

H_0

H_a



Quiz

$$\chi^2 = \sum_k \frac{(O - E)^2}{E}$$

Observed

Expected

Question:

To assess whether a coin is fair or not, we need to compare the expected outcomes with the observed outcomes of tossing the coin.

The expected outcome for a fair coin toss is 50% heads and 50% tails, which corresponds to 25 heads and 25 tails in 50 tosses.

The observed outcomes from 50 coin tosses are 28 heads and 22 tails.

To determine whether the coin is fair, we perform a chi-square test to check if the observed results significantly deviate from the expected results. If the deviation is statistically significant, it may indicate that the coin is not fair.

Step 1: Assumption

H_0 : Coin is fair
 H_a : Coin is biased

Step 2: Distribution?

χ^2 distribution

✓ Step 3: p-value

✓ Step 4: Comparing with α

50 Times

	H	T
Expect	25	25
Actual	28	22

$$\chi^2 = ?? \Rightarrow \left(\frac{(28-25)^2}{25} + \frac{(22-25)^2}{25} \right)$$

$$\Rightarrow \underline{\underline{0.72}}$$

Question:

To assess whether a coin is fair or not, we need to compare the expected outcomes with the observed outcomes of tossing the coin.

The expected outcome for a fair coin toss is 50% heads and 50% tails, which corresponds to 25 heads and 25 tails in 50 tosses.

The observed outcomes from 50 coin tosses are 45 heads and 5 tails.

To determine whether the coin is fair, we perform a chi-square test to check if the observed results significantly deviate from the expected results. If the deviation is statistically significant, it may indicate that the coin is not fair.

Observed → (cars, hit, pub)
expected → (✓, ✓, ✓)

①
②

Goodness of fit → (Observed value
fits
expected distribution
or not)

Test for Independence ✓✓

Survey: **Gender** impacts online/offline purchases

↓ ✓
Preferences vs Gender ✓
(Online/offline)

Preference

Observed Values			
	M	W	
Offline	527	72	599
Online	208	102	308
	733	174	907

H_0 : Gender & preference are independent

H_a : Gender & preference are dependent

Preference

Observed Values

	M	W	
Offline	527	72	599
Online	208	102	308
	733	174	907

$$① 907 \times 66\% = 599$$

$$② 907 \times 34\% = 308$$

Expected Values

Preference

	M	W	
Offline	484✓	115✓	<u>599</u> 66%
Online	249✓	59✓	<u>308</u> 34%
	733	174	907

$$\left\{ \chi^2 = \sum_K \frac{(\underline{O} - E)^2}{E} \right\}$$

Quiz

	Online	print	TV
Buy			
Not Buy			

→ Test of Independence
because it says, if there is any relationship b/w
these categories.

Difference

① Coin toss \rightarrow fit the expected distribution

(Goodness of fit) ✓

(when expected distribution is known)

✓ ② preference vs gender \rightarrow Testing for independence

(χ^2 contingency)

Assumptions

- ✓ ① Variables are categorical
- ✓ ② Observations are independent (Political example)
- ✓ ③ Each cell is mutually exclusive (Only 1 choice example)
- ✓ ④ Expected value in each cell is