



Chi-Square Test

- 1) Degree of freedom
 - 2) Chi Square → Goodness of fit
→ Independence Test
-

① 2 Numerical Samples →

Num. vs 2 Cat → Two Sample T Test

More than 2 Cat → Anova

② Cat vs Cat Comparisons →

⇒ IQ Score vs Schools (2 schools)

⇒ Gender vs Rating

⇒ Gender vs Preference Online/Offline

⇒ Is a coin fair?

⇒ Toss Multiple Times →

H	T	→
# of	[50]	[50]

Counts → Category

⇒ Degree of freedom →

lets say → we want to plan a movie week.
→ everyday we want to see a different genre.

1 genre → 7 days

(i) Choices for 1st night → 7

(ii) 2nd night → 6

(iii) 3rd night → 5

(iv) 4th night → 4

(v) 5th night → 3

(vi) 6th night → 2

(vii) 7th night → no choice

Why we need DOF → Chi-Square Test

⇒ 1) Salary of 3 individual

1st → 35 L

2nd → 36 L

3rd → ??

avg salary → 35 L

⇒ 2) 4 people →

35, 36, x, 37 → avg salary → 36

\Rightarrow given you want to know n object, and
you already know one aggregation value
How many objects should you already know to calculate
all of them.

$\Rightarrow \underline{n-1} \leftarrow$ degree of freedom

<u>Q</u>	H	W
	73	85
	68	73
	71	96
	72	82
	62	70
<u>avg</u>	69.2	81.2

$SH, SW \rightarrow 4H, 4W$
and avg
to calculate everything

$$\rightarrow \underline{4+4}$$

\Rightarrow I need 8 values to predict remaining value.

minimum 8 values to predict \rightarrow

minimum no. of variable needed to be known \rightarrow

$$\boxed{n_1 + n_2 - 2} \leftarrow$$

(3) Sachin Century vs Victory

		Victor		
		Yes	No	
				46
C	Yes			46
Yes				314
100		184	176	360

Given margin
aggregation are
known

Only need 1
value

⇒ If both row & column sums are known

→ 1 value → dot = 1

→ (4) Election →

← Politicians

	A	B	C	D	
X					349
Y					151
Z					150
	150	150	200	150	650

$$(\# \text{ rows} - 1) > (\# \text{ col} - 1)$$

$$(3-1)(4-1)$$

$2 \times 3 = 6$ values needs to be known.

⇒ Importance of DOF →

- 1) Chi-Square (Categorical) → Critical values
- 2) Based on DOF → distribution of data changes
- 3) Hypothesis Testing → find dist for H₀
 - ↓
 - Calculate p-value

If you have two arrays with lengths n1 and n2, what is the formula to calculate degrees of freedom for the chi-square test?

$$\Rightarrow \mu = \{ \quad] \leftarrow n_1 \text{ length} \\ \omega = [\quad] \leftarrow n_2 \text{ length}$$

$$n_1^{-1} \quad n_2^{-1} \Rightarrow (n_1 - 1) + (n_2 - 1)$$

$$\Rightarrow [(x_1, y_1), (x_2, y_2), \dots]$$

$$x_1 y_1 = - \quad \cdot$$

$$x_2 y_2 = - \quad \cdot$$

⋮

⇒ Chi-Square Goodness of fit →

Q Is a Coin fair →

⇒ One Categorical variable → 2 diff values.

⇒ Toss the Coin 50 times → 28 H, 22 T

⇒ ① $H_0 \rightarrow$ Coin is fair

$H_a \rightarrow$ Coin is unfair

② If H_0 is true → 25 H, 25 T

H_0 is true

	H	T	
Exp	25	25	50
Actual	28	22	50

$\text{dof} = \frac{1}{2}$

$\hat{A} =$

→ 500 tries

	250	250	
Exp	250	250	
Act	280	220	

$\hat{B} =$

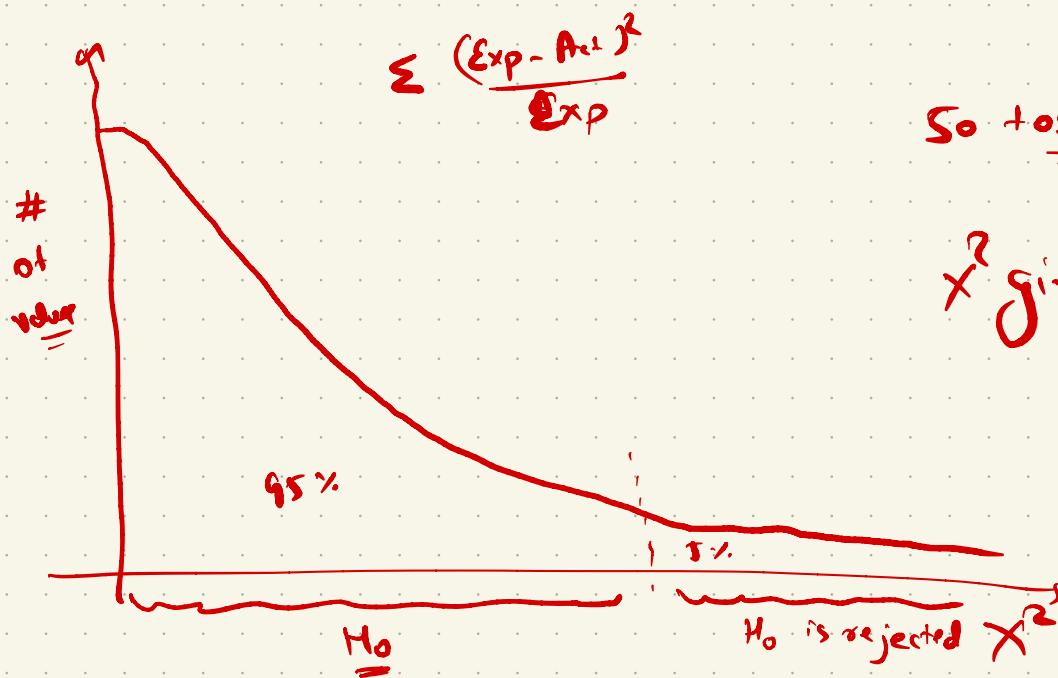
	250	250	
Exp	250	250	
Act	253	247	

$\hat{C} =$

$$\chi^2 = \sum \frac{(Exp - Act)^2}{Act}$$

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

If I plot a graph of χ^2
 If H_0 is true χ^2 \rightarrow small $\leftarrow H_0$ is true
 \rightarrow large $\leftarrow H_0$ is rejected



χ^2 = Chi statistic

↳ dist' of χ^2 \rightarrow Chi Square distribution

Based on $\alpha = 0.05 \rightarrow$ Critical χ^2 value
any test/experiment \rightarrow if they have more
extreme data \rightarrow reject H_0 .

$$\Rightarrow \chi^2 = \sum \frac{[Exp - Ac]^2}{Exp}$$

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

Actual \equiv observed

Exp \equiv Theoretical

$\Rightarrow \chi^2 \propto$ diff²
if diff \uparrow $\chi^2 \uparrow$

Z Score
T Stats
 χ^2 Stats } Base on this P-value
calculate

Compare p-value against α .

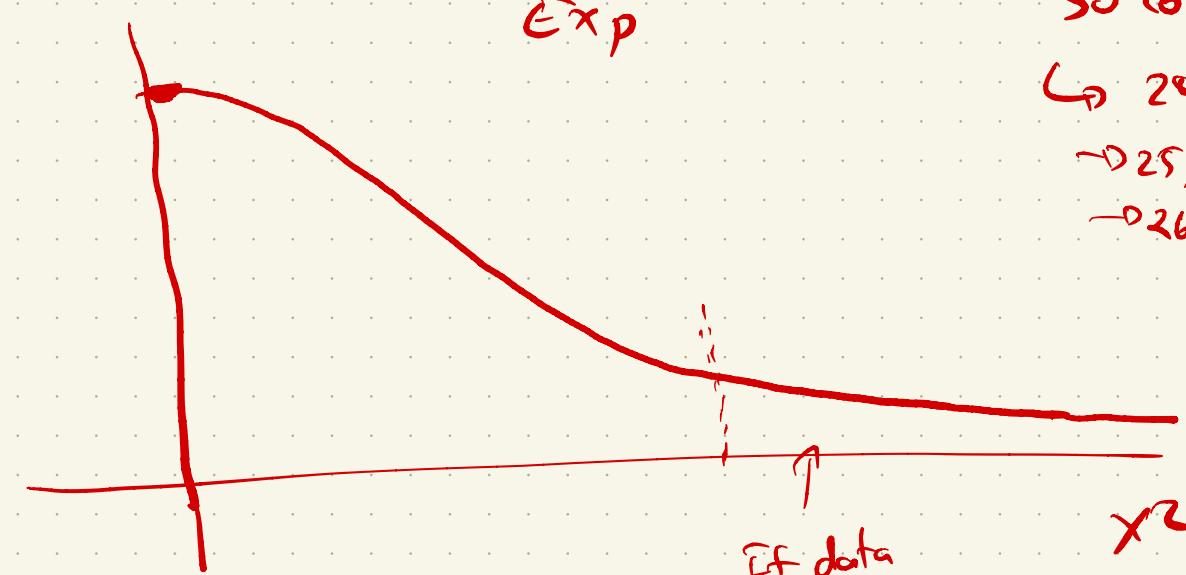
Z-Score \rightarrow How many std dev away from mean the observed data is

χ^2 \rightarrow How much diff from the expected data is the observed value.

p-value \rightarrow $1 - \text{norm.cdf}(Z\text{-Score})$

p-value \rightarrow $1 - \text{Chi2.cdf}(\chi^2)$

$$\Rightarrow \chi^2 = \sum \frac{(Exp - Act)^2}{Exp}$$



so cont

$\rightarrow 28, 26 \rightarrow \dots$

$\rightarrow 25, 25 \rightarrow \dots$

$\rightarrow 26, 24 \rightarrow \dots$

If data χ^2

lies in the extremes

right \rightarrow reject H_0

	n	T	
Exp	25	25	norm
Act	28	22	<u>chi²</u>

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

$$\chi^2 = \frac{18}{25} = 0.72$$

Corresponding P -value

$$1 - \text{chi}^2_0 \text{ Cdf}(0.72, df=1)$$

= P -value

A researcher is studying the preferences of people in a city for three different modes of transportation: car, bicycle, and public transit.

The researcher surveyed 500 individuals and found that 240 prefer cars, 160 prefer bicycles, and 100 prefer public transit.

The researcher wants to know if there is a significant difference between the observed preferences and the expected preferences based on historical data.

Which statistical test should the researcher use?

500

\Rightarrow	<u>Obs</u>	240	160	100
<u>Exp</u>		200	200	100

Goodness of fit

→ How well does observed value fit the expected dist'.

Good fit → Small χ^2 →

Chi-Square Test of Independence

⇒ 2 Categorical variable → Check if they are related.

Q: females are more likely to buy online clothes.

⇒ gender affects the preference of channel
M ⤵ ⤵ F → Offline Online

		M	F	Gender
		527	72	599
		206	102	308
prefers		733	174	1907

Observed data

⇒ H_0 → Gender has no relation with channel.

H_a →

⇒ 907 people surveyed

599 → offline

308 → online

733 M
174 F

$$\Rightarrow \frac{599}{907} = 0.66$$

$$\Rightarrow \frac{308}{907} = 0.334$$

\rightarrow first calculate population distⁿ

$$M \rightarrow 733 \rightarrow 16\% = 484$$

$$34\% = 115$$

$$F \rightarrow 174 \rightarrow 66\% = 249$$

$$34\% = 59$$

Exp values

	M	ω	
off	484	249	599
on	115	59	308
	733	174	907

Obs values

	M	ω		
off	527	72	599	608
on	206	102	308	34%
	733	174	907	

\Rightarrow 907 \rightarrow 733 M
 \rightarrow 174 F

$$599 \rightarrow (599) \times \left(\frac{733}{907} \right) = \underline{\underline{484}}$$

$$\chi^2 = \frac{(484 - 527)^2}{484} + \frac{(249 - 272)^2}{272} + \dots$$

1 - Chi-sq. Cdf (χ^2 , df = 1)

= P-value

A market researcher is exploring the connection between age group (under 25, 25-40, over 40) and smartphone brand preference (Brand A, Brand B, Brand C). The researcher collects data from 600 respondents and plans to perform a chi-square independence test.

How many degrees of freedom are associated with this test?

	A	B	C
< 25			
25-40			
> 40			

$$= (\# \text{ rows} - 1)(\# \text{ cols} - 1)$$

$$= 2 \times 2$$

$$= 4$$

⇒ Aerotit →

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.

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?

⇒ GOF → If we want to check how close to exp we are seeing

✓ Independence → Relationship b/w 2 Cat Variable

⇒ Assumption for Chi-Square Test

- 1) Categorical variable
 - 2) Each cell is mutually Exclusive ^{+oo}
 - 3) Observations are independent
 - 4) Each cell ≥ 5 .
-

Collab Link : <https://colab.research.google.com/drive/1unOC7zE0BH5VsYfOmQ-f-8ED-ZB1E2n4?usp=sharing>

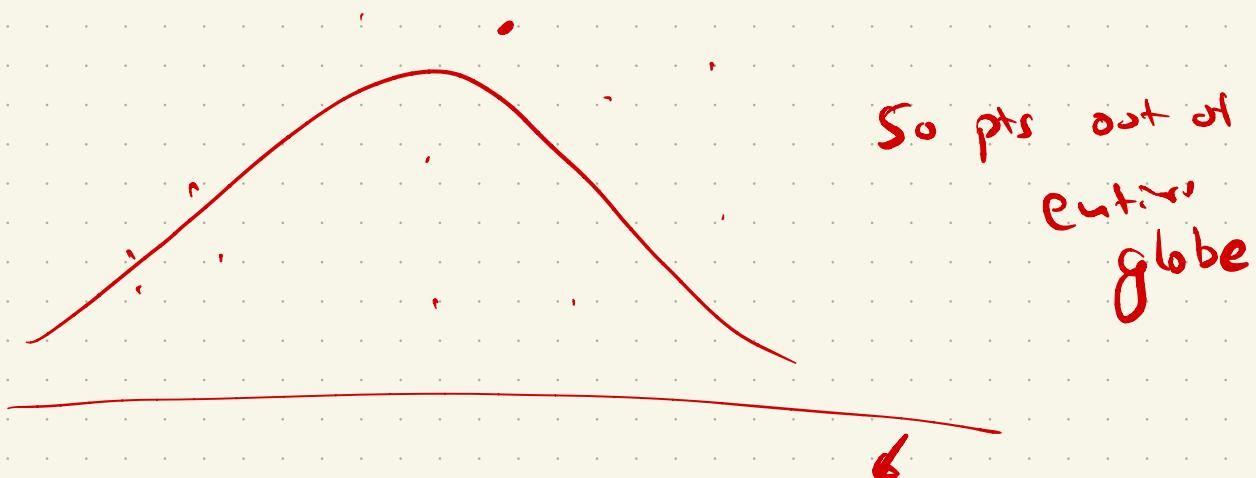
\Rightarrow Rather than checking for actual date.

↓
we check dat of SM \rightarrow
 \hookrightarrow Normal dist'

What is the mean population.

\Rightarrow So people Test \rightarrow IQ

[100, 120, 130, ... 101, ...]
So values.



\Rightarrow 50 people \rightarrow Set 1 \rightarrow avg of s_1 , \leftarrow
50 people \rightarrow Set 2 \rightarrow avg of s_2 , \leftarrow

1000s s_3 —

$M_1, M_2, \dots, M_{1000} \rightarrow$ always follows
normal dist'

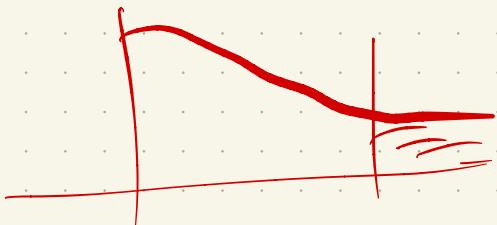
with overall $\mu = M_{\text{pop}}$

$$\sigma = \frac{\sigma_{\text{pop}}}{\sqrt{n}}$$

$$\Rightarrow \sum \frac{(Exp - Obs)^2}{Exp}$$

Exp \rightarrow Counts
Obs \rightarrow

If H_0 true \rightarrow close to 0



\Rightarrow Chi Squared

$$\Rightarrow \chi^2 \Rightarrow 1 - \text{chi2.cdf}(\chi^2) = p\text{-val}$$

$$\Rightarrow \sum \frac{(Exp - Obs)^2}{Exp}$$

\Rightarrow Chisquare \rightarrow Just give array for Exp \rightarrow
for Obs \rightarrow

⇒ Chi-Contingency →

Observed
value

+ H_0 is that data is independent

			R_1
			R_2
C_1	C_2	C_3	T

⇒ Create a Expected table



Collab : <https://colab.research.google.com/drive/1unOC7zE0BH5VsYfOmQ-f-8ED-ZB1E2n4?usp=sharing>