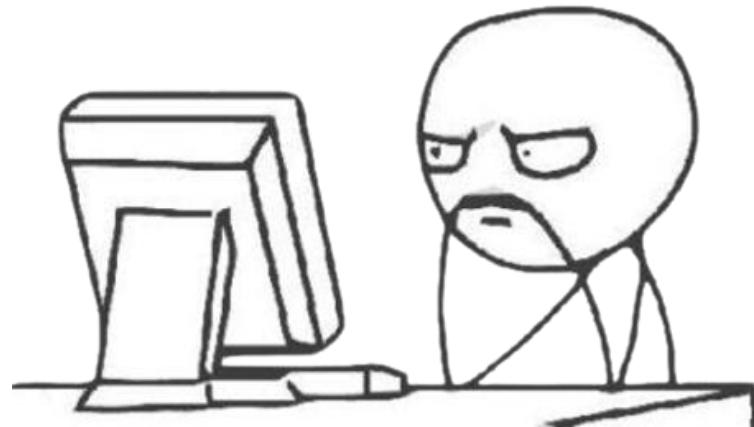# Hypothesis Testing
## Chi-Square Test

# Recap Test (Solve it and share your answers in the chat section only to me. Will discuss the solution at 07:08 AM)

Quiz: A group of 5 patients were treated with medicine A and another group of 7 patients with medicine B. Researchers claim that B is better than A in treating the disease. Test the claim at 5% significance level. (Note: In the below table, higher number is better)

A = [42, 39, 38, 60, 41]
B = [38, 42, 56, 64, 68, 69, 62]

( Two independent Sample T-Test )

$H_0$: $\mu_a \geq \mu_b$
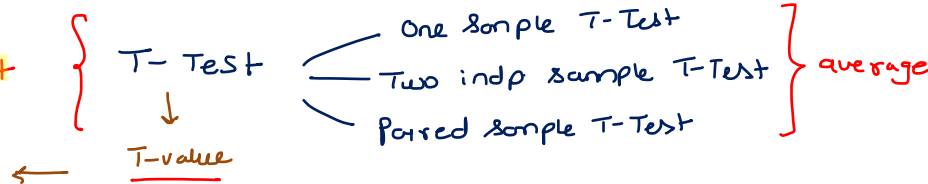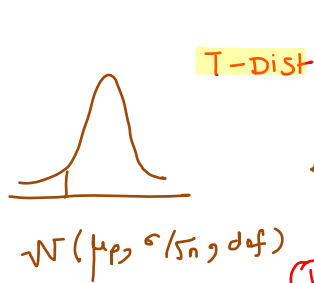
$H_a$: $\mu_a < \mu_b$

① Sample Size (n < 30) ⟶ T-Test

②

a) one sample T-Test : we are given with one group (sample) and a global/population mean. compare Sample mean with population mean

b) Two independent Sample T-Test : Two independent groups. compare avg between the two groups

c) paired Sample T-Test

⟶ one group but the data has been collected at two different time

| | B.T | A.T |
|-----|-----|-----|
| P-1 | | |
| P-2 | | |
| P-3 | | |
| ⋮ | | |

Z-Dist

Z - Test ⟨
  one Sample Z - Test
  Two ind. Sample z - Test  ⎫ average
  one Sample Z.prop.Test
  Two indp Sample Z.prop.Test ⎫ proportion

↓

← Z-value

$N(\mu_P, \sigma/\sqrt{n})$

---

T-Dist

T - Test ⟨
  One Sample T-Test
  Two indp sample T-Test  ⎫ average
  Paired Sample T-Test

↓

T-value

←

$N(\mu_P, \sigma/\sqrt{n}, dof)$

T-Dist is dependent
on degree of
freedom

① Income between males and females

② Test if pill improves the IQ Level

③ Height of a group of people is less than 66 inches

(all target variables are
continuous variable)

T - Test $\Big\}$ continuous
Z - Test $\Big\}$ variable

chi - Square : Hypothesis testing on categorical data

① chi-square test of goodness of fit
② chi-square test of independence

chi-Square value

chi-Square dist

(dof)

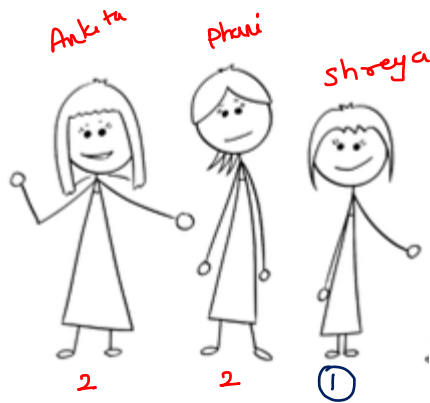# Agenda:

- **Degree Of Freedom**
- **Chi-Square Motivations Using Toy Example**
- **Chi-Square Implementation In Business Case**

# Degree Of Freedom

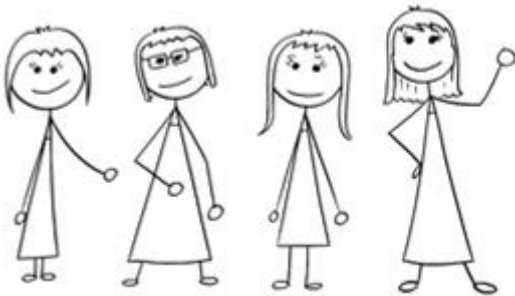Ankita    Phani    Shreya

Sum = 5

2    2    (1)

$$[\ A,\quad P,\quad S\ ]$$

↓

these values
are free to deviate

$$\boxed{dof = (n-1)}$$

when data is a linear series then
degree of freedom $=(n-1)$

Total = 5

Ankita = 2
phani = 2
shreya =

# Degree Of Freedom (Sample Standard-Deviation)

$$\underline{Dof}$$

① linear Series (1 - Dimensional) = $(n-1)$

② multi-dimensional (n - Dimensional) = $\sum\limits_{i=1}^{d} (n-1)i$

③ contengency table / cross-tab = $(r-1) \times (c-1)$

$r = \#$ of categories in rows
$c = \#$ of categories in columns

$\rightarrow 1 - chi2 \cdot cdf(0.72, 1)$

$\chi^2 = 0.72$
$dof = 1$

$0.72$

$chi2 \cdot cdf(0.72, 1)$

$$\boxed{P-value = 1 - chi2 \cdot cdf(\chi^2, dof)}$$

# Degree Of Freedom

S

| Height | Weight |
|--------|--------|
| 73 | 85 |
| 68 | 73 |
| 74 | 96 |
| 71 | 82 |
| X | Y |
| AVG: 71 | 81.2 |

S-1        S-2

8

$S-1 = (n-1)$

$S-2 = (n-1)$

$$S = (n-1) + (n-1) = \sum_{i=1}^{d} (n-1)$$

$(n-1) + (n-1) + (n-1) + \cdots + (n-1)_d$

| Age | income | height | weight |
|-----|--------|--------|--------|
| | | | |
| | | | |
| | | | |

6

$\underbrace{(n-1)}_{age} + \underbrace{(n-1)}_{income} + \underbrace{(n-1)}_{height} + \underbrace{(n-1)}_{weight}$

$(6-1) + (6-1) + (6-1) + (6-1)$

$\underbrace{5 + 5}_{10} + \underbrace{5 + 5}_{10}$

20

# Degree Of Freedom

| | | India Win | | |
|---|---|---|---|---|
| | | False | True | Sum |
| Sachin Century | False | 160 | 154 | 314 |
| | True | 16 | 30 | 46 |
| | Sum | 176 | 184 | 360 |

| | | India Win | | ↓ |
|---|---|---|---|---|
| | | False | True | Sum |
| Sachin Century | False | 160 | 154 | 314 |
| | True | 16 | 30 | 46 |
| → | Sum | 176 | 184 | 360 |

no. of categories in rows → $r = 2$

$c = 2$

↑ no. of categories in columns

Dof = 1

Contengency table / cross-tab

$Dof = (r-1) \times (c-1)$

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

$= \boxed{1}$

# Degree Of Freedom

| Cities | | Political Party | | | | |
| --- | --- | --- | --- | --- | --- | --- |
| | | A | B | C | D | Sum |
| | X | 90 | 60 | 104 | 95 | 349 |
| | Y | 30 | 50 | 51 | 20 | 151 |
| | Z | 30 | 40 | 45 | 35 | 150 |
| | Sum | 150 | 150 | 200 | 150 | 650 |

| Cities | | Political Party | | | | |
| --- | --- | --- | --- | --- | --- | --- |
| | | A′ | B′ | C′ | D′ | Sum |
| | X | 90 | 60 | 104 | – | 349 |
| | Y | 30 | 50 | 51 | – | 151 |
| | Z | – | – | – | – | 150 |
| | Sum | 150 | 150 | 200 | 150 | 650 |

$r = 3$

$c = 4$

$dof = (r-1) \times (c-1)$

$= 2 \times 3$

$= 6$

# Chi-Square Test

1. Goodness Of Fit: If the data is uniformly distributed or not?
2. Test Of Independence: If two categorical variables are independent or not?

① <u>Goodness of fit</u> : we test if the data across various categories is uniformly distributed or not?

Toss of a fair coin
(50)

- Head
- Tail

|  | observed | Expected |
|---|---|---|
| Head | 45 | 25 |
| Tail | 5 | 25 |

If this observed result is Similar to the expected result or not?

we want to test if the observed data is uniformly dist. between the categories.

Test if the attendance across various Stadiums is uniformly dist. or not or is there a particular Stadium with higher or lower attendance than expected.

| | observed attendance | expected attendance |
|---|---|---|
| Delhi | 29k | 25k |
| Chennai | 30k | 25k |
| Bengalore | 45k | 25k |
| Kolkata | 12k | 25k |

↓
categorical

(observed = expected)

# Coin Toss

we fail to reject null

categorical

|  | Head | Tails |
|---|---|---|
| Expected | 25 | 25 |
| Actual | 28 | 22 |

we reject null

|  | Head | Tails |
|---|---|---|
| Expected | 25 | 25 |
| Actual | 45 | 05 |

$H_0$: Expected outcome = observed outcome
$H_a$: Expected outcome $\neq$ observed outcome

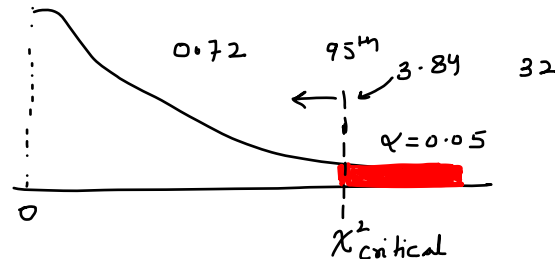$$\chi^2 = \sum_{i=1}^{n} \frac{(O-C)^2}{C}$$

chi-square value / statistic

Difference between actual and expected

Deviation of observed data from expected data

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

$$\chi^2 = 0.72$$

$$\chi^2 = [0, \infty]$$

$$\chi^2 = \left[\frac{(45-25)^2}{25}\right] + \left[\frac{(05-25)^2}{25}\right]$$

$$\chi^2 = 32$$

$$\chi^2_{critical} = (dof, 1-\alpha)$$
$$= (dof = 1, 1-\alpha = 0.95)$$



0.72   95$^{th}$   3.84   32

$\alpha = 0.05$

$\chi^2_{critical}$

**Example**

observed

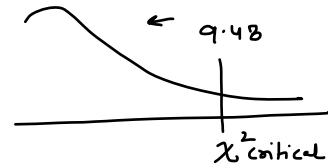|  | Actual | Expected |
|---|---|---|
| Chennai | 50 | 50 |
| Bangalore | 60 | 50 |
| Delhi | 40 | 50 |
| Punjab | 47 | 50 |
| Kolkata | 53 | 50 |
| Total | 250 | 250 |

Attendance at various stadiums hosting IPL is reported as given in the below table. Is the attendance uniform across all stadiums or there was stadium with significantly high attendance? $\alpha = 0.05$

$H_0$ : Actual attendance = Expected attendance

$H_a$ : Actual attendance $\neq$ expected attendance

$$x^2 = \sum_{i=1}^{2} \frac{(o-e)^2}{e}$$

$$= \frac{(50-50)^2}{50} + \frac{(60-50)^2}{50} + \frac{(40-50)^2}{50} + \frac{(47-50)^2}{50} + \frac{(53-50)^2}{50}$$

$$= x^2 = 4.36$$

← 9.48

we fail to reject null

(No such stadium where the attendance was signi higher or lower)

$x^2$ critical

total attendance = 250

no. of stadiums = 5

what would be the uniform attendance = 250/5

assumption → all stadiums would receive the same attendance with no bias or preferences

# chi-square Test of independence  (Two-categorical var.)

objective → To test if two categorical feature are dependent or independent on each other?

**Prod. cat | gender**

| Prod. cat | gender |
|---|---|
| A | M ← |
| A | F ← |
| B | M |
| A | M |
| C | F |
| B | F |
| A | M |

Each sales happening is independent of the other

↘ (m, F)

(A, B, C)

(observed contingency table)

|  | m | f |  |
|---|---|---|---|
| A | 3 | 1 | 4 ← |
| B | 1 | 1 | 2 |
| C | 0 | 1 | 1 |
|  | 4 | 3 | 7 |

(Expected frequencies)

|  | m | f |  |
|---|---|---|---|
| A | = 2 | 2 | 4 |
| B | 1 | 1 | 2 |
| C | 1 | 0 | 1 |
|  | **4** | 3 | 7 |

$\left(\dfrac{\text{Row total}}{\text{Table total}}\right)$

→ column total

$\dfrac{(\text{Row total}) \times (\text{column total})}{\text{Table total}}$

$= \dfrac{4 \times 4}{7}$

If we randomly select a customer what's the prob that customer will buy product A without looking of their gender =

$P(A) = 4/7 = 0.57$
$P(B) = 2/7 = 0.28$
$P(C) = 1/7 = 0.14$

⎫ irrespective of gender

**④**
0.57 × 4 → A ≈ 2.28 ≈ 2
0.28 × 4 → B = 1.12 ≈ 1
0.14 × 4 → C = 0.56 ≈ 1

**③**
0.57 × 3 → A = 1.74 ≈ 2
0.28 × 3 → B = 0.84 ≈ 1
0.14 × 3 → C = 0.42 ≈ 0

**Gender Impact on Offline/Online Purchase**
**Is the "Type of Purchase" influenced by "Gender"**

$$\chi^2 = \frac{(527-484)^2}{484} + \frac{(72-115)^2}{115} + \frac{(206-242)^2}{242}$$

$$+ \frac{(102-57)^2}{57}$$

$$\chi^2 =$$

observed

$dof = (r-1)\times(c-1)$
$= (2-1)\times(2-1)$
$= \textcircled{1}$

expected

|  |  | Gender | | |
|---|---|---|---|---|
|  |  | Male | Female | Sum |
| Type Of Purchase | Offline | 527 | 72 | 599 |
|  | Online | 206 | 102 | 308 |
|  | Sum | 733 | 174 | 907 |

|  |  | Gender | | |
|---|---|---|---|---|
|  |  | Male | Female | Sum |
| Type Of Purchase | Offline | 484 | 115 | 599 |
|  | Online | 242 | 57 | 308 |
|  | Sum | 733 | 174 | 907 |

a customer can fall in one cell with no common elements

$P(offline) = 599/907 = \underline{0.66}$

$P(online) = 308/907 = \underline{0.33}$

expected values

733 →
offline = $0.66 \times 733 = 483.78 \approx 484$
online = $0.33 \times 733 = 241.89 \approx 242$

174 →
offline = $0.66 \times 174 = 114.84 \approx 115$
online = $0.33 \times 174 = 57.42 \approx 57$

# Assumptions of Ch-2 Test

① variables are categorical

② observation are independent

③ each cell is mutually exclusive

④ Expected value in each cell $> 5$

chi-2 distribution

dependent on dof

① 1-Dimension

$dof = (n-1)$

② n-Dimension

$dof = \sum_{i=1}^{q} (n-1)$

③ Cross-tab $dof = (r-1) \times (c-1)$

## chi-Square test

① chi-Square goodness of fit → if the data across various categories is uniformly dist. or not?

② chi-square test of independence → if two categorical variables are dependent or independent

### functions

① chi2 → Critical chi-square = chi2.ppf $(dof, 1-\alpha)$

p-value = $(1 - chi2 \cdot cdf(X, dof))$

② chisquare → $(O, e)$

③ chi2-contingency =

$$X^2 = \sum_{i=1}^{n} \frac{(o-e)^2}{e}$$

↓ chi-square statistics

→ deviation of the observed data from the expected data

Quiz:

$$\begin{array}{c|c} a-1 & a-2 \\ \hline & \\ | & | \\ n_1 & n_2 \end{array}$$

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

   $dof = (n_1 - 1) + (n_2 - 1)$

1. In a chi-square test, what does the chi-statistic represent? → Diff. between observed and expected

1. 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? goodness of fit

   | | |
   |---|---|
   | car | 240 |
   | bic | 166 |
   | | 100 |

   ③

1. 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). ③ 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?

   $dof = (r-1) p (c-1) = (2) p (2) = ④$

1. 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? Test of independence

1. When testing the independence of two categorical variables, what are the assumptions of the chi-square test?