# Chi - Squared Test

# Framework

① $H_0$ Vs $H_a$ ⟶ Burden of proof (solid evidence)

    ⟶ default assumption ( in the absence of data)

② Test stat ( from observation)

③ R Vs L Vs Two - Tailed

④ P - value ⟶ Prob of seeing data assuming $H_0$ was true

⑤ Compare p-value with $\alpha$ ⟶ significance level

## Chance Vs Signifance

$\left( \begin{array}{l} \text{avg IQ of} \\ \text{10 students} \rightarrow 101.5 \end{array} \right)$
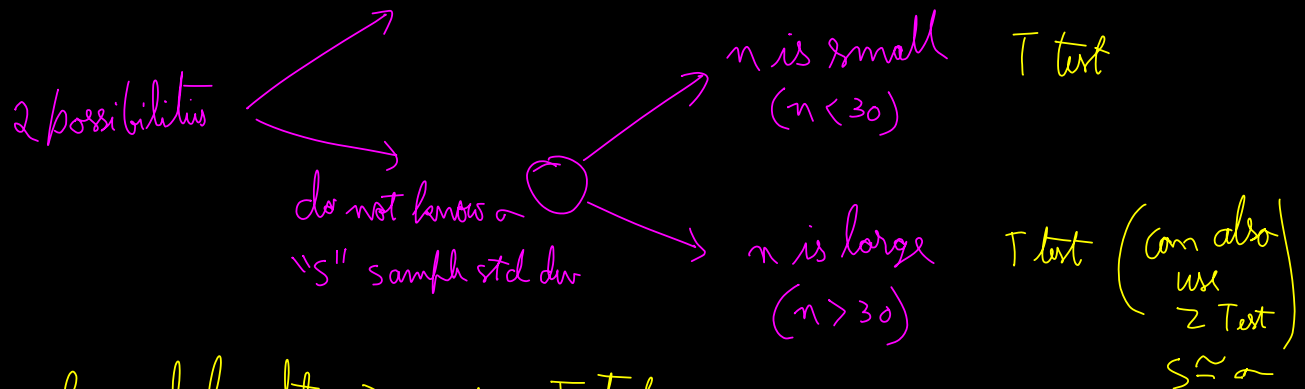
Recovery after taking drug 1 Vs drug 2 } ttest_ind

        Two sets of samples

## About population Std dev "$\sigma$" & no. of samples "$n$"

$$Z\text{-stat} = \frac{X - \mu}{\sigma/\sqrt{n}}$$

$$T\text{-stat} \quad \frac{X - \mu}{s/\sqrt{n}}$$

we know $\sigma$    Z test

2 possibilities

do not know $\sigma$
"$s$" sample std dev

$n$ is small    T test
$(n < 30)$

$n$ is large    T test $\begin{pmatrix} \text{can also} \\ \text{use} \\ \text{Z Test} \end{pmatrix}$
$(n > 30)$          $s \approx \sigma$

By default $\rightarrow$ use T Test

# Degree of freedom

Salary ①    35L , 36L, ? $\xrightarrow{\text{avg}}$ 35L
                 0
                 34 L

From these 3 nos, we only need 2 (if we know avg)

②    35L, 36L, ?, 38L $\xrightarrow{\text{avg}}$ 37L
                 0
                 ↓
                 39L

From 3 out of 4, we can find missing no.
if avg is known

**Q**   If "$n$" nos are given with their mean,
how many will I need to know?

      $n-1$      " degrees of freedom "

H & W

| Inch | Kg |
|------|------|
| 73 | 85 |
| 68 | 73 |
| 74 | 96 |
| 71 | 82 |
| X | X |

DF = 8

Avg    71    81.2

$n_1$   no of  hight values

$n_2$   no of  wright values

& their avg

$n_1 - 1 + n_2 - 1$

$\boxed{n_1 + n_2 - 2}$ $\longrightarrow$ degrees of freedom

# Sachin Century Vs Victory

|  | Win | | |
|--------|-------|------|-----|
|  | False | True | |
| False | (160) | 154 | 314 |
| True | 16 | 30 | 46 |
|  | 176 | 184 | 360 |

Century

df = 1

# Regional support for Politician

4 politicians → A, B, C, D

3 Cities → X, Y, Z

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

df = 6

2 x 3

$(r-1)(c-1)$

2 x 3 = 6

# Coin Toss : fair or Biased     50 tosses

$H_0$: fair coin

$H_a$: Biased coin

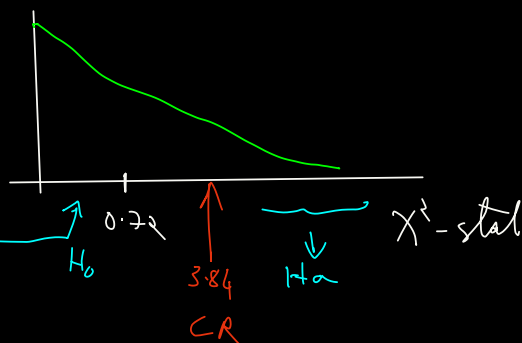|          | Heads | Tails |
|----------|-------|-------|
| Expected | 25    | 25    |
| Observed | (28)  | 22    |

$\chi^2$ statistic : $\dfrac{(28-25)^2}{25} + \dfrac{(22-25)^2}{25} = 0.72$

$\rightarrow df = 1$

Under $H_0$, will this
stat be low or high

Right tailed
because $H_a$ is on right side



$H_0$   0.72   3.84   $H_a$   $\chi^2$-stat
              CR

Online Vs Offline   →   does this preference depend on Gender

Survey observation

|        | M   | W   |     |      |
|--------|-----|-----|-----|------|
| Offline | 527 | 72  | 599 | 66%  |
| Online  | 206 | 102 | 308 | 34%  |
|        | 733 | 174 | 907 |      |

Expected under $H_0$

|        | M   | W   |
|--------|-----|-----|
| offline | 484 | 115 |
| Online  | 249 | 59  |

$H_0$: Offline/online is independent of Gender
$H_a$: depends on gender

$$\frac{(527-484)^2}{484} + \frac{(72-115)^2}{115} + (\quad) + (\quad) \qquad \approx 59$$

Suppose $H_0$ is true, 66% prefer offline,

   Among 733 men, how many are expected to prefer offline?
      → 66% of 733 → 484

   Among 174 women, how many are expected to prefer offline?
      → 66% of 174 → 115

$$\chi^2 : \quad \frac{(527-484)^2}{484} +$$

$$\sum_{i=1}^{4} \frac{(O_i - e_i)^2}{e_i}$$

# Assumptions of Chi Squared

1. Variables are Categorical

2. Observations are independent

3. Each cell is mutually exclusive

4. Expected value in each cell $\left( \begin{array}{c} \text{atleast } 80\% \\ \text{of the cells} \end{array} \right)$ should be greater than 5