

## Chi-Square

- A chi-square ( $\chi^2$ ) is a statistic that checks for patterns or relationships in categorical variables.
- A categorical variable is a non-numeric characteristic like - Gender, language spoken.
- What kind of question can a chi-square answer?  
One way Chi-Square - One variable → Is a six-sided die fair?

1 2 3 4 5 6

Rolls 20 23 21 25 24 32

$\chi^2 = 6.0$ ,  $p < 0.5$ , indicates significant relationship between groups.

$\chi^2 = 6.0$ ,  $p > 0.5$  indicates insignificant relationship between groups.

Q) Why not just look at the number count?

→ We can't be sure if it is reliable or created by chance. If we do again can we get same exact numbers.  
Use Inferential statistics and use chi-square.

Example - Does gender vary across educational majors?

	Engineering	Business	Psychology
Female	2	2	3
Male	3	2	2

If there is no relationship, we would expect gender to be evenly spread across majors.

If there is a relationship, we would expect gender to be unevenly spread across majors.

If there is No relationship.. Tables are known as contingency table.

	Engineering	Business	Psychology
Female	50	50	50
Male	50	50	50

Uneven Relationship Contingency table.

	Engineering	Business	Psychology
Female	40	55	55
Male	60	45	45

$$\chi^2 = \frac{(\text{Expected} - \text{Observed})^2}{\text{Expected}} \Rightarrow \text{Calculate each group}$$

So, is 6 is a real difference? will it hold next time

If new set of data comes?

$$i) \text{Engg - Female} = \frac{(50-40)^2}{50} = \frac{100}{50} = 2$$

$$ii) \text{Engg - Male} = \frac{(50-60)^2}{50} = \frac{100}{50} = 2$$

So we will have 6 group.

$$\chi^2 = 2 + 2 + 0.5 + 0.5 + 0.5 + 0.5$$

$$\chi^2 = 6$$

If  $p = 0.05$ , 5% chance we would get these result with random data.

Accepted range normally is,  $p \leq 0.5$  or else reject

Sample size - Each group should have minimum count of 5.

Types of Chi-Square - i) Test of Independence ii) Goodness of fit.

i) Test of Independence - Test for a relationship between two categorical variables.

Eg - Testing relationship between gender & major.

If say any variable has a influence on other categorical variable.

ii) Goodness of fit - Compares categorical values in your sample to a known or hypothesized value. Eg - Comparing gender & major at your university to the gender & major distribution nationwide.

Limitations & Assumptions - i) It don't need normal distribution ii) At least one value in each cell of contingency table.

iii) It doesn't tell which levels of variable are driving the effect.

iv) All data should be independent. One variable should not affect other variable.