## Chi - Square test

→ Tests for categorical Variable →

$$X^2 = \sum_{ij} \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

→ Take two categorical variable and check whether they share a relationship i.e whether they are independent or dependent

→ It is commonly used for __Testing__ relationships between variables. The null hypothesis of the chi- square test is the no relationship exists on the categorical Variables in the population, They are independent.

null hyp  $H_0$ : They are independent

alt hyp  $H_1$ : They are dependent

## Problem : __Chi- square test of independence__

|  | non smoker | smoker |  |
|---|---|---|---|
| athlete | 14 | 4 | 18 |
| non - athlete | 0 | 10 | 10 |
|  | 14 | 14 |  |

$P(S) = \frac{1}{2}$ ,  $P(S') = 1 - \frac{1}{2} = \frac{1}{2}$

$P(NS) = \frac{1}{2}$ ,  $P(NS') = 1 - \frac{1}{2} = \frac{1}{2}$

$P(A) = \frac{18}{28}$ ,

$P(NA) = \underline{10}$

$H_0$ : S and A are independent

$H_1$ : S and A are dependent

Step 1 : Assume $H_0$ is correct

① $P(S) . P(A) = \frac{1}{2} \times \overset{9}{\cancel{18}}/28$

$= \frac{9}{28} \times \cancel{28} = 9$ people

② $P(S) . P(NA) = \frac{1}{2} \times \frac{10}{28} = \frac{5}{28}$

$= \frac{5}{28} \times 28 = 5$

③ $P(NS) . P(A) = \frac{1}{2} \times \overset{9}{\cancel{18}}/28$

$= \frac{9}{28} \times 28 = 9$ people

④ $P(NS) . P(NA) = \frac{1}{2} \times \overset{5}{\cancel{10}}/28$

$= \frac{5}{28} \times 28 = 5$

Calculated Value

|  | non smoker | smoker |  |
|---|---|---|---|
| athlete | 9 | 9 | 18 |
| non - athlete | 5 | 5 | 10 |
|  | 14 | 14 |  |

Step 2 : Convert this into Chi- squared distribution

$$x^2 = \sum_{ij} \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

$$\Rightarrow \frac{(14-9)^2}{9} + \frac{(0-5)^2}{5} + \frac{(4-9)^2}{9} + \frac{(10-5)^2}{5}$$

$$= 2.78 + 2.78 + 5 + 5 = 15.56 \quad \text{is} \quad \text{Not } 0$$
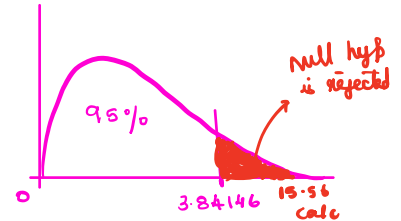
Step 3 :

$$df = \left( row \overset{\text{no}}{\div} - 1 \right) * \left( \overset{\text{no}}{\div} \overset{\text{of}}{colu} - 1 \right) = 1$$

$$\alpha = 0.05$$

$$\therefore \text{ the critical Value} = 3.84146$$

Test statistics $= 15.56 >$ critical value $= 3.84146$


95%
null hyp is rejected
0    3.84146   15.56 calc

∴ Reject the null hypothesis

---

| Degree of freedom | $= ($ No of rows $-1) * ($ No of columns $-1)$ |

EM Algorithm $\Rightarrow$ Expectation Maximization Algorithm

Extremly Important

Normalization

1mb

# Minmax Scaler

⇒ No matter the 'value' of the data after the application of Minmax Scaler the data is converted to $[0, 1]$.

$$Max = a$$
$$min = b$$
$$data = x_i$$

$$x_i \longrightarrow \frac{x_i - b}{a - b}$$

$$a < \frac{x_i - b}{a - b} < 1$$

# Standard Scaler

$$\text{to make } x_1 \cdots x_n \longrightarrow N(\overset{\mu}{0}, \overset{\sigma^2}{1})$$

$$\mu = \bar{x}$$

$$\sigma^2 = \sum \frac{(x_i - \bar{x})^2}{n - 1}$$

normalization ⟹
$$x_1 \rightarrow \frac{x_1 - \mu}{\sigma}$$
$$x_2 \rightarrow \frac{x_2 - \mu}{\sigma}$$
$$\vdots$$
$$x_n \rightarrow \frac{x_n - \mu}{\sigma}$$

# Robust Scaler

Change `$\mu$` mean by median

`$\sigma$` variance by IQR = Intra Quartile Range

i.e $\dfrac{x_i - (median)}{IQR}$

## Algorithms to study

1) Statistical Algorithm

2) ML Algorithms

3) Neural Networks / Deep Learning


1) ## Statistical algorithm

→ old

→ don't need computational efficiency

→ Pedantic ⇒ require many conditions to be satisfied

2) ML Algorithm

→ newer

→ need good - efficient computers
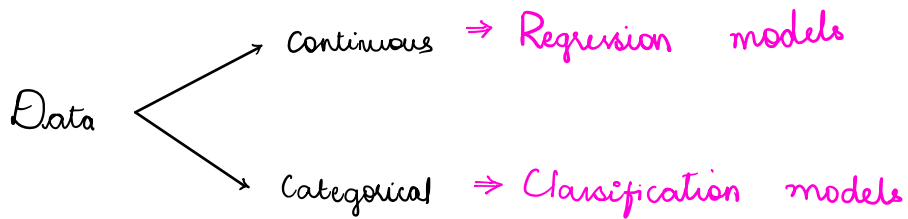
→ Need to satisfy minimal conditions

3) DL Algorithm

→ new

→ Extremely efficient machines

## Statistical Algorithm

⇒ We try to bring out some sort of inference.

⇒ We make predictions.

Data
- → Continuous ⇒ **Regression models**
- → Categorical ⇒ **Classification models**
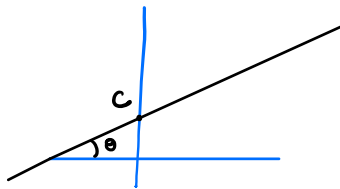
## Regression model

⇒ We will try to predict a continuous variable.

## Linear model

### first look at the Plot

→ Linear model is essentially a straight line

## Regression



$$y = mx + c$$

$$m = \tan\theta$$