

Statistics

Agenda

- ① Sampling Techniques ✓
- ② Covariance And Correlation
- ③ Probability distribution function
- ④ Probability density function
- ⑤ Probability Mass function
- ⑥ Cumulative Density Function
- ⑦ Types of Distribution

① Sampling Techniques

① Random Sampling : Simple Random Sampling gives each member of the population an equal chance of being chosen for the sample.

Eg: Vaccination Test $\rightarrow 100 \rightarrow$ Randomly Select Couple

Accidental Test $\rightarrow 1000 \rightarrow A \rightarrow 2$ vehicle

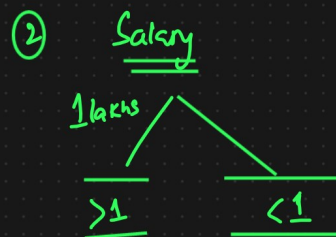
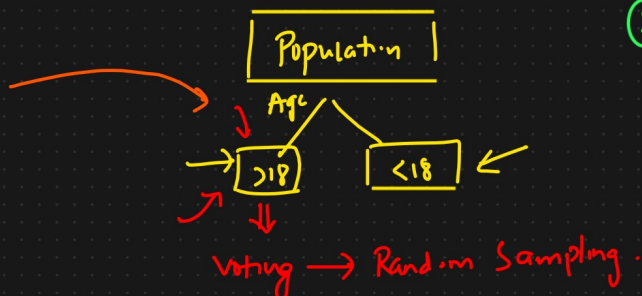
Exit Poll

Average IQ of the school \rightarrow Select 10 people \Rightarrow

② Stratified Sampling [Stratified \rightarrow layers]

Stratified Sampling Involves dividing the population into sub population that may differ in important ways

Eg: Exit Poll





③ Systematic Sampling : Systematic Sampling is a statistical method involving the selection of elements from an ordered sampling frame.

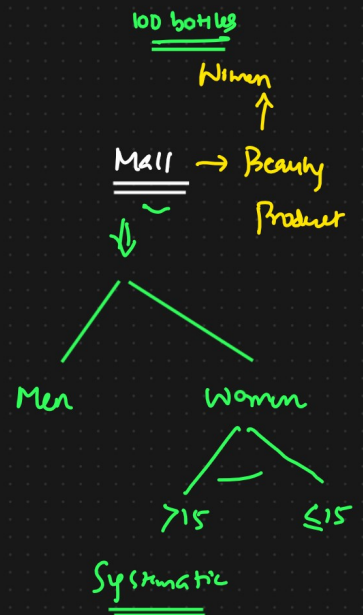
Eg: Airport

A Credit Card Company

Mail → Fill up a form for a card.



n^{th} individual → 7^{th}

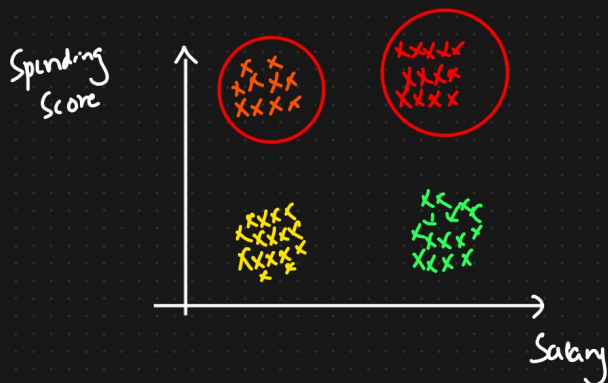


④ Convenience Sampling Assignment

⑤ Purposive Sampling → Judgemental Sampling is a method where researchers decide which members of the target population will be sampled.

⑥ Cluster Sampling

Mail iPhone 28 → 100 → 20%



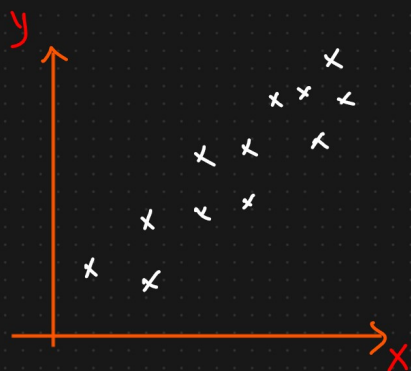
② Covariance And Correlation

X y
2 3
4 5
6 7
8 9

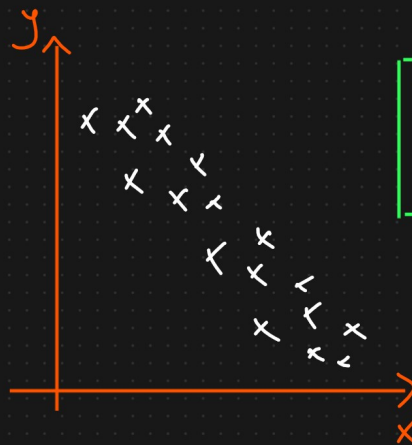
X↑ y↑
X↓ y↓

Quantify
[Relationship between X And y].

↓
Covariance And Correlation.



X↑ y↑
X↓ y↓



X↑ y↓
X↓ y↑

Covariance

$$\text{Cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

$x_i \rightarrow$ Data point of x

$\bar{x} \rightarrow$ Sample mean of x

$y_i \rightarrow$ Datapoints of y

$\bar{y} \rightarrow$ Sample mean of y .

$$\begin{aligned} \text{Var}(x) &= \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1} \\ &= \frac{\sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})}{n-1} \end{aligned}$$

$$\text{Var}(x) = \text{Cov}(x, x) \Rightarrow \underline{\underline{\text{Spread}}}$$

$\text{Cov}(x, y)$

X↑ y↑
X↓ y↓

+ve Covariance

X↑ y↓
X↓ y↑

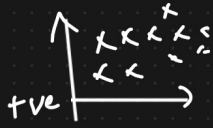
-ve Covariance

| X | y |
|---------------|---------------|
| → 2 | 3 |
| → 4 | 5 |
| 6 | 7 |
| <u>6</u> | <u>7</u> |
| $\bar{x} = 4$ | $\bar{y} = 5$ |

$$\text{Cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}_i)}{n-1}$$

$$= \frac{[(2-4)(3-5) + (4-4)(5-5) + (6-4)(7-5)]}{2}$$

$$= \frac{4 + 0 + 4}{2} = 4 //$$



| X | y |
|-----------------|-----------------|
| 8 | 6 |
| 7 | 5 |
| 6 | 4 |
| 3 | 1 |
| 2 | 0 |
| <u>25</u> | <u>16</u> |
| $\bar{x} = 5.2$ | $\bar{y} = 3.2$ |

+ve
Covariance

6.7

| X | y |
|----------------|-----------------|
| 8 | 6 |
| 9 | 5 |
| 10 | 4 |
| 12 | 2 |
| 11 | 1 |
| <u>50</u> | <u>18</u> |
| $\bar{x} = 10$ | $\bar{y} = 3.6$ |

-ve Covariance

-3



Advantages

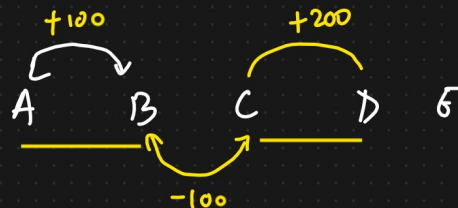
① Relationship between X and y

ve or -ve value

$A \leftrightarrow B$ $\left. \begin{array}{l} + \text{ infinity} \\ - \text{ infinity} \end{array} \right\}$

Disadvantages

① Covariance does not have a specific limit value



② Pearson Correlation Coefficient $[-1 \text{ to } 1]$

$$\rho_{x,y} = \frac{\text{Cov}(x,y)}{\sigma_x \cdot \sigma_y} \Rightarrow -1 \text{ to } 1$$

① The more the value towards $+1$, the more +ve correlated x, y is.

② The more the value towards -1 , the more -ve x, y is

| | | | |
|----------------|-----------------|-------------------------|--|
| x | y | | |
| 8 | 6 | | |
| 9 | 5 | $\text{Cov}(x, y) = -3$ | |
| 10 | 4 | | |
| 12 | 2 | | |
| 11 | 1 | | |
| $\bar{x} = 10$ | $\bar{y} = 3.6$ | | |

| | | | |
|--|--|--|---|
| | | | $\text{Variance} = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$ |
| | | | $\sigma = \sqrt{\text{Variance}}$ |
| | | $\rho_{x,y} = \frac{-3}{1.58 \times 2.073}$ | $\left[\frac{4 + 1 + 0 + 4 + 1}{4} \right] = \frac{10}{4} = 2.5$ |
| | | $= \frac{-3}{3.27}$ | $\sigma = \sqrt{2.5}$ |
| | | $= -0.917 \Rightarrow \text{Negative Correlated} \Rightarrow 91.7\%$ | |

$$x \uparrow \quad y \downarrow \Rightarrow 91.7\%$$

$$1 \text{ unit } x \uparrow \quad y \downarrow \underline{\underline{0.917}}$$

③ Spearman Rank Correlation $[-1 \text{ to } 1]$.

$$r_s = \frac{\text{Cov}(R(x), R(y))}{\sigma(R(x)) * \sigma(R(y))}$$

| X | Y | R(x) | R(y) |
|---|---|------|------|
| 1 | 2 | 2 | 2 |
| 3 | 4 | 3 | 3 |
| 5 | 6 | 4 | 4 |
| 7 | 8 | 5 | 6 |
| 0 | 7 | 1 | 5 |
| 8 | 1 | 6 | 1 |

Ascending Order

{ Non linear Relationship }

Spearman Rank Correlation

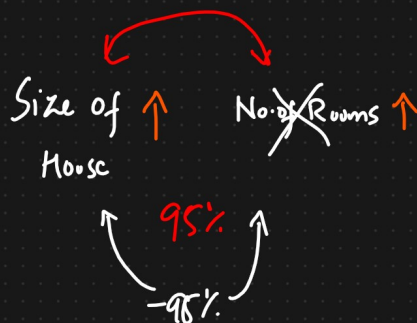


Non linear Relationship

DATA SCIENCE

Feature Selection

Housing Price DATASET



+ve or
-ve Correlation
Location

-ve
Correlation
Hunted

NO
~~No. of people
staying~~

O/p
Price ↑