

STATISTICS

It is the science of collecting, organizing & analyzing data.

- Data : Facts or pieces of information.

Eg: Height of student, no. of sales in term of revenue

→ Types of Statistics-

① Descriptive statistics

- It consists of organizing, summarizing & visualizing data.

1) Measure of Central Tendency
[Mean, Median, Mode]

2) Measure of Dispersion
[Variance, Z-score, Standard deviation]

3) Different types of Distribution of data

[Histogram, pdf, pmf,]
Gaussian Dist., Binomial
Bernoulli, Exponential,
Poisson, Log Normal
Dist

② Inferential Statistics

- It consists of using data you have measured to form conclusion.

- 1) Z-test
- 2) T-test
- 3) Chi-square test
- 4) Anova test

(hypothesis testing, P-values, Significance value)

8 Feature engineering

EDA: Exploratory Data analysis

DESCRIPTIVE STATISTICS

Date _____
Page _____

① Population & Sample Data

Population (N) : Group or superset of data that you are interested in studying.

Sample (n) : It is a subset of population data.

→ Types of data

Data.

Quantitative

Qualitative

Discrete

Continuous

Nominal

Ordinal



+ve whole number

any value

No Ranks

Ranks

Eg: no. of children
in family

Eg: Average
rainfall

Eg: Gender
(M,F)

Eg: Customer
feedback

- number of people
working in a
company

- length
of river

- Blood grp
- colors
- location

* Scale of Measurement

- 1) Nominal scale data
- 2) Ordinal scale data
- 3) Interval scale data
- 4) Ratio scale data.

1) Nominal Scale Data

- Qualitative / Categorical data
- Order or rank does not matter
- Eg: Gender, color, labels

Eg: Color

red \rightarrow 5 \rightarrow 50 %.

Yellow \rightarrow 3 \rightarrow 30 %.

orange \rightarrow 2 \rightarrow 20 %.

10

2) Ordinal scale data

- rank is important
- order matters
- Difference cannot be measured

3) Interval scale data

- order matters
- difference can be measured
- ratio cannot be measured
- No 0(zero) starting point

Eg: Temperature variable

30F \rightarrow 60:30 = 2:1

60F \rightarrow XX

90F

120 O'F XX

4) Ratio scale data.

- order matter
- difference are measurable (ratio)
- contain "0" starting point

Eg: Students marks in class

3:1

0, 90, 60, 30, 75, 45, 50
3:2

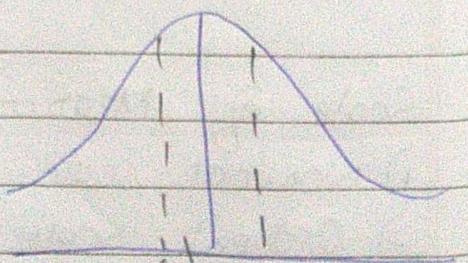
* Measure of central Tendency

It is very imp in EDA [Exploratory Data Analysis]

1) Mean or average

2) Median

3) Mode



To measure this
Central region we
need central tendency

① Mean : Avg of all the elements present in distribution

Population (N)

Sample (n)

random variable : $X = \{1, 1, 2, 2, 3, 3, 4, 5, 5, 6\}$

$$\text{Population mean } (\mu) = \frac{\sum_{i=1}^N X_i}{N} = \frac{[1+1+2+2+3+3+4+5+5+6]}{10} = \frac{32}{10} = 3.2$$

② Median : Whenever there is outlier we use median.

$$X = \{4, 5, 2, 3, 2, 2\}$$

Steps :

1) Sort the random variable X

$$X = \{1, 2, 2, 3, 4, 5\}$$

2) No. of elements $\begin{cases} \text{if } n=\text{even} \Rightarrow \text{middle 2 element} \Rightarrow \text{Avg} \\ \text{if } n=\text{odd} \Rightarrow \text{central element} \end{cases}$

$$= \frac{2+3}{2} = 2.5$$

③ Mode : frequency of maximum occurring element

$$X = \{2, 1, 1, 1, 4, 5, 7, 8, 9, 9, 10\}$$

$$\text{Mode} = 1$$

- In most of scenarios mode is used in categorical value

* Measure of Dispersion

1) Variance → The More the spread the More variance

Population variance

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$

Sample variance

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

x_i → Data points

μ → population mean

N → population size

x_i → Data points

\bar{x} → sample mean

n → sample size

• The sample variance is divided by $(n-1)$ so that we can create an unbiased estimator of population variance.

2) Standard Deviation

Population Std

$$\sigma = \sqrt{\text{Variance}}$$

Sample Std

$$s = \sqrt{\text{sample variance}}$$

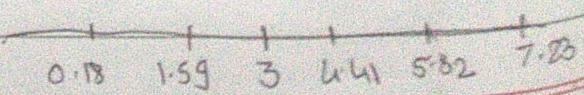
• How many std x_i is away from mean

Eg: $\{1, 2, 3, 4, 5\}$, $\mu = 3$

$$\text{Population variance} = \frac{[(1-3)^2 + (2-3)^2 + (3-3)^2 + (4-3)^2 + (5-3)^2]}{5}$$

$$= \frac{4+1+0+1+4}{5} = \frac{10}{5} = 2$$

$$\sigma = \sqrt{2} = 1.41$$



* Random variables

It is a process of mapping the output of a random process or experiment to a number

Eg: Tossing a coin

$$X = \begin{cases} 0 & \text{: head} \\ 1 & \text{: Tail} \end{cases}$$

Rolling a dice

$$Y = \{2, 1, 4, 2, 3, 5\}$$

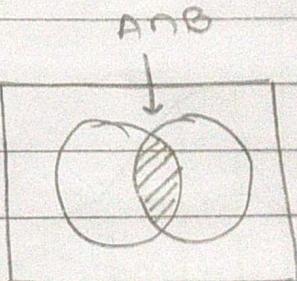
* Sets

$$A = \{1, 2, 3, 4, 5, 6, 7, 8\}$$

$$B = \{3, 4, 5, 6, 7\}$$

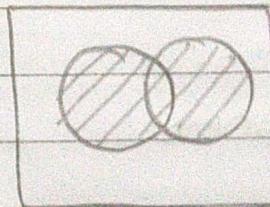
1) Intersection

$$A \cap B = \{3, 4, 5, 6, 7\}$$



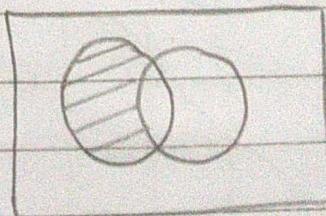
2) Union

$$A \cup B = \{1, 2, 3, 4, 5, 6, 7, 8\}$$



3) Difference

$$A - B = \{1, 2, 8\}$$



4) Subset

$A \rightarrow B$: False

$B \rightarrow A$: True

5) Superset

$A \rightarrow B$: True $\rightarrow A$ has all the elements w.r.t B

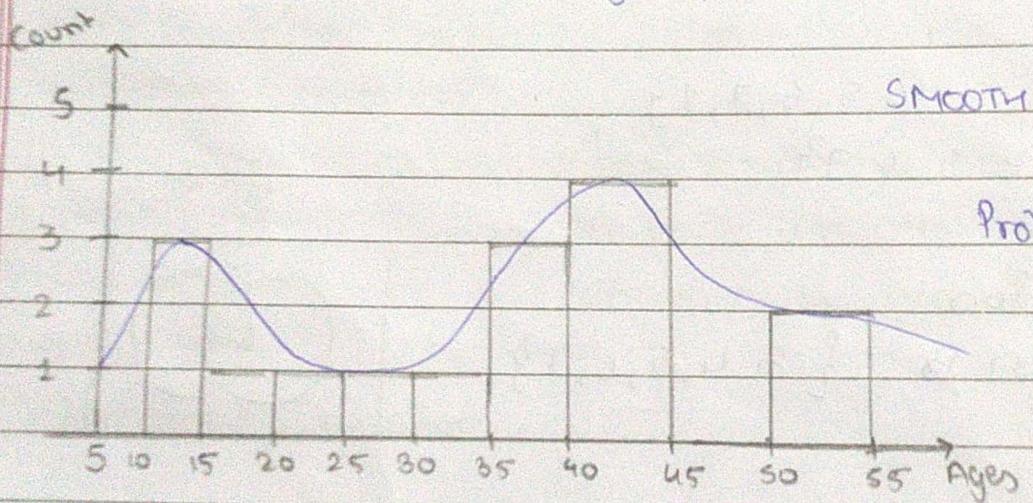
$B \rightarrow A$: False

* Histogram And Skewness

[Frequency] \Rightarrow Distribution.

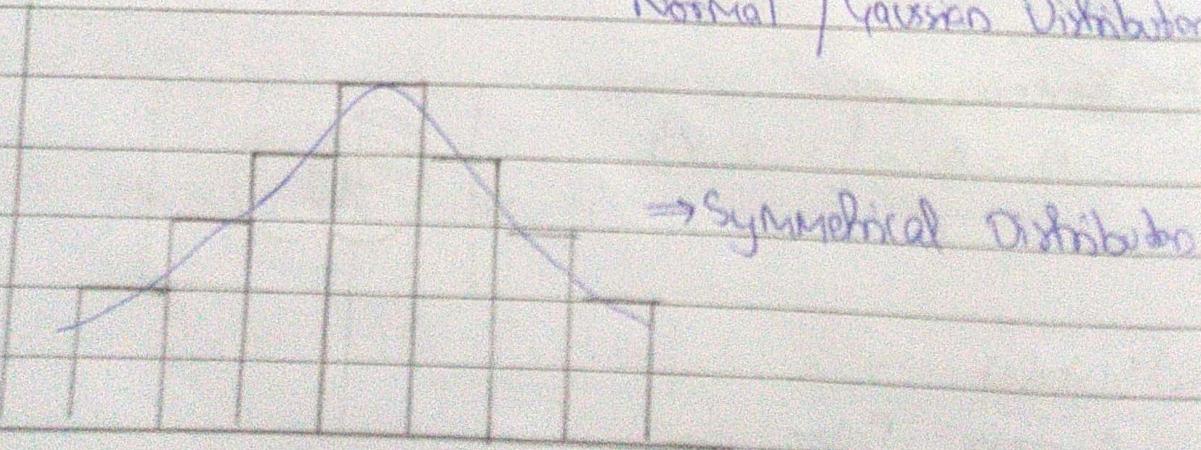
Ages = {10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 50, 51}

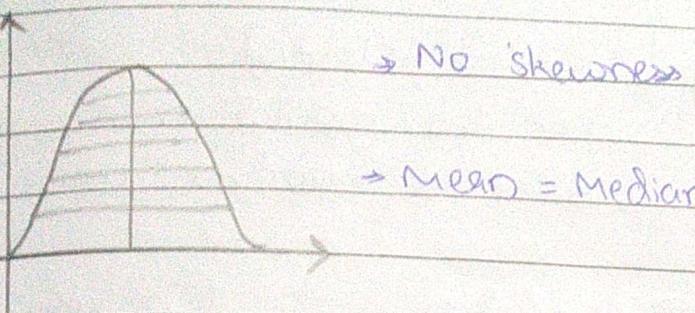
Bins : How many groups we are going to make & what will be the group size.



* Skewness

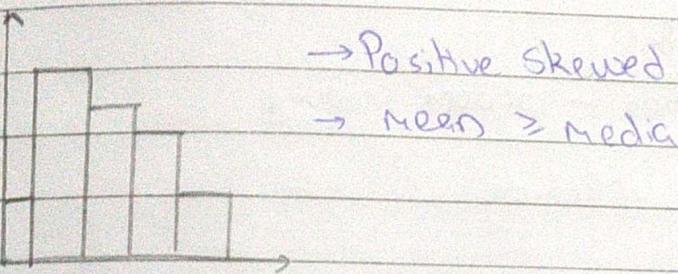
Normal | Gaussian Distribution





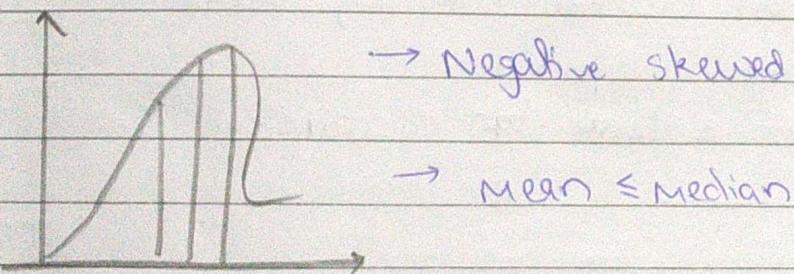
→ Mean = Median = Mode

2) Right skewed



→ Mean \geq Median \geq Mode

3) Left skewed



→ Mean \leq Median \leq Mode

* Sampling Techniques

i) Random sampling : Simple random sampling gives each member of the population an equal chance of being chosen for the sample.

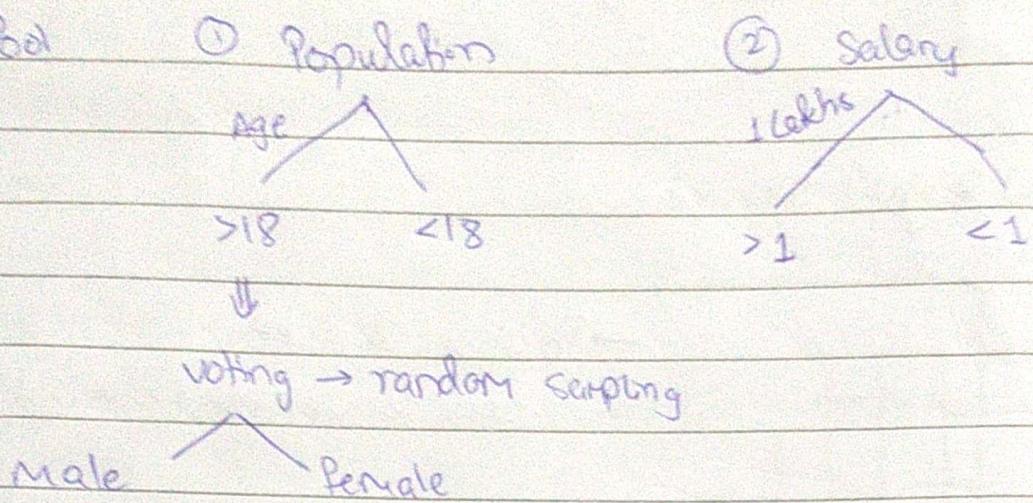
Eg: Vaccination test \rightarrow 100 \rightarrow randomly select couple
Accidental Test \rightarrow 1000 \rightarrow A \rightarrow 2 vehicle

Avg IQ of school \rightarrow select 10 people

2) Stratified Sampling [Stratified \rightarrow layers]

It involves dividing the population into sub-population that may differ in important ways

Eg: Exit Poll



3) Systematic Sampling

It is a statistical method involving the selection of elements from an ordered sampling frame.

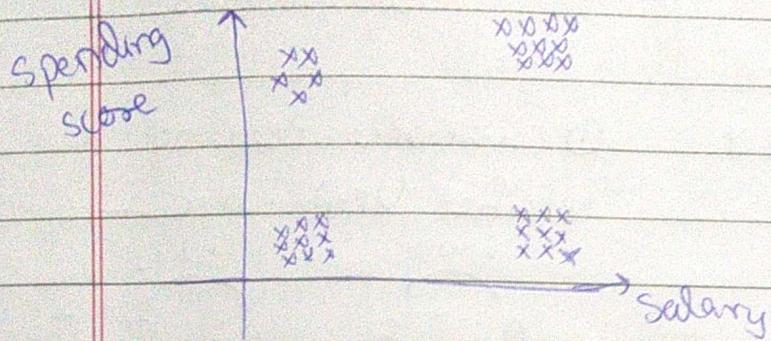
Eg: Mall \rightarrow fill up a form for a cause.

4) Convenience Sampling

5) Purposive Sampling

Judgemental Sampling is a method where researchers decides which members of the target population will be sampled.

6) Cluster Sampling



* Covariance & Correlation

X	Y
2	3
4	5
6	7

Relation between X & Y

$x \uparrow$ $y \uparrow$



$x \downarrow$ $y \downarrow$

Covariance & Correlation

Covariance

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

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

x_i → Data point of x

$$= \frac{\sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})}{n-1}$$

\bar{x} → Sample mean of x

y_i → Data points of y

\bar{y} → Sample mean of y

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

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

$x \uparrow$	$y \uparrow$
$x \downarrow$	$y \downarrow$

+ve covariance

$x \uparrow$	$y \downarrow$
$x \downarrow$	$y \uparrow$

-ve covariance

Advantages

Disadvantages

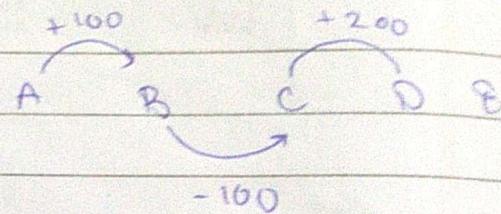
① Relationship between x & y

+ve or -ve values

$A \leftrightarrow B$ $\begin{cases} +\infty \\ -\infty \end{cases}$

① Covariance does not have a

specific limit value



② Pearson Correlation Coefficient

In order to prevent the disadvantage we use it.

$$r_{x,y} = \frac{\text{Cov}(x,y)}{\sigma_x \cdot \sigma_y} \rightarrow -1 \text{ to } +1$$

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

x

y

8

6

$$\bar{x} = 10, \bar{y} = 3.6$$

9

5

$$\sigma = \sqrt{2.5}$$

10

4

$$\sigma_x = 1.53, \sigma_y = 2.073$$

12

2

$$\text{Cov}(x,y) = -3$$

11

3

$$\text{variance} = \sum_{i=2}^n (x_i - \bar{x})^2$$

$$\sigma = \sqrt{\text{variance}}$$

$$r(x, y) = \frac{-3}{1.58 \times 2.073} = \frac{-3}{3.27} = -0.917$$

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

negative correlation $\rightarrow 91.7\%$

1 unit $x \uparrow y \downarrow 0.917$

- Pearson Correlation Coefficient does not work well with non-linear data

③ Spearman Rank Correlation [-1 to 1]

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

Ascending order

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

Spearman Rank Correlation



Non linear relationship