

Statistics

Defn:

Statistics is the science of collecting, organizing and analyzing data.

Data: "Facts or pieces of information".

Ex:- Heights of students
IQ of students
Daily Activities
weights or Age of people

Types of Statistics

① Descriptive stats

Defn: It consists of organizing and summarizing data.

- ① Measure of central tendency
[mean, median, mode]
- ② Measure of Dispersion
[variance, standard deviation]
- ③ Different types of distribution of data.

Ex: Histogram,

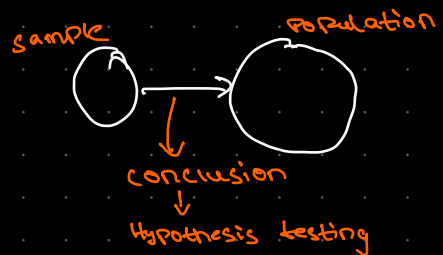
PDF (Probability Distribution Function)

PMF (Probability Mass Function)

CDF (Cumulative Distribution Function)

② Inferential stats.

Defn: It consists of data you have measured to form conclusion.



1) Z-test

2) t-test

3) Chi square test

4) ANOVA

5) F test

conclusion of sample or population

Ex: Lets say there are 20 classes in your college and you have collected the heights of students in the class

Heights are recorded as { 175 cm, 180 cm, 200 cm }

Descriptive:

"What is the average height of students in class" (mean).

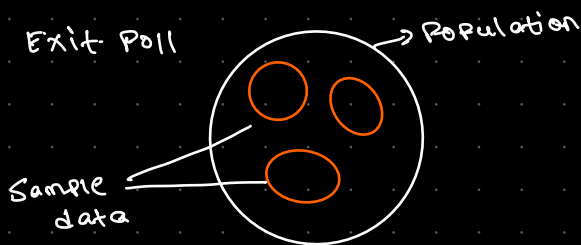
Inferential:

"Are the height of students in the classroom ^{sample} similar to what you expect in the college" (Population)

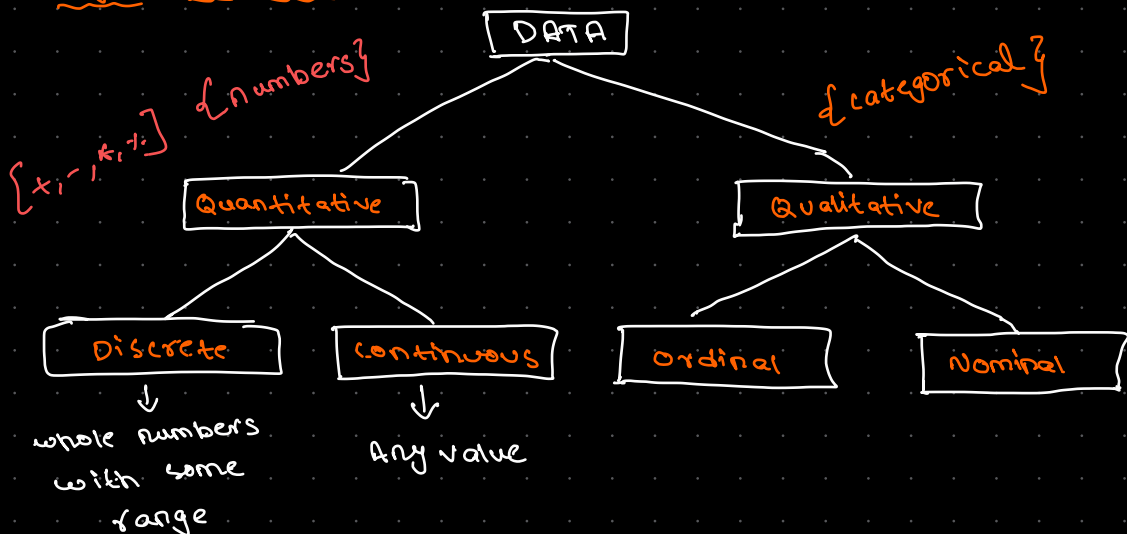
(N) Population and (n) Sample data

Ex:

Exit Poll



Types of Data:



Ex: * Number of bank accounts
 * no. of children in a family
Ex: weight, height, Age, temperature, speed etc...
Ex: Ranks [small, medium, large] [good, better, best]
Ex: gender (M, F) Blood groups (O, B, A...) colour of hair..

* Scales 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

Ex: Gender, colours, labels

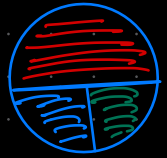
→ order does not matter.

Ex: Favourite colour

Red → 5 (50%)

Blue → 3 (30%)

Green → 2 (20%)



10

② Ordinal scale data:-

- ① categorical data
- ② Ranking and order matters
- ③ Difference cannot be measured

Ex:

small → 1

medium → 2

large → 3

③ Interval scale data:-

- The order matters
- Difference can be measured
- Ratio cannot be measured.
- No "0" starting points.

Temperature

30 F
 60 F
 90 F
 120 F

④ Ratio Scale data

- The order matters
- Differences are measurable including ratios.
- Contains a '0' starting point.

Ex:- Students mark in class

0, 30, 45, 60, 90, 99

Examples:-

- ① marital status (Nominal scale data)
- ② Favourite food based on gender (Nominal)
- ③ IQ measurement (Ratio scale data)
(can also be converted to ordinal)

Descriptive Statistics

① measure of central tendency

i, mean ii, median iii, mode.

① mean:-

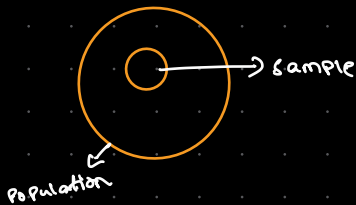
Population (N)

sample (n)

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

$$\text{Population mean } (\mu) = \sum_{i=1}^n \frac{x_i}{N}$$

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



$$\frac{1+1+2+2+3+3+4+5+5+6}{10} = 3.2$$

Population size (N)

sample size (n)

② Median

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

Steps:

① Sort the numbers $\{1, 2, 2, 3, 4, 5\}$

② No. of elements

③ if count == even

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

$$\downarrow$$
$$\frac{2+3}{2}$$

$$= 2.5 \text{ (median)}$$

if count == odd

$$\{1, 2, 2, 3, 4, 5, 6\}$$

\downarrow
median

why median?

* means are affected by outlier *

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

$$\bar{x} = \frac{1+2+3+4+5}{5}$$

$$= 3 \text{ (mean)}$$

$$\text{median} = 3$$

(odd)

$$X = \{1, 2, 3, 4, 5, 100\}$$

$$\bar{x} = \frac{1+2+3+4+5+100}{6}$$

$$= \frac{115}{6} \approx 19 \text{ (mean)}$$

$$\text{median} = \frac{3+4}{2}$$

$$= 3.5$$

Conclusion:

median is used to find central tendency when outlier is present.

③ Mode:

maximum frequency occurring elements

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

$$\text{mode: } 1$$

$$* \{1, 2, 2, 3, 4\}$$

$$\{1, 2\}$$

(depends) *

EDA and Feature Engineering

Age	weight	salary	gender	degree
24	70	40k	m	BE
25	80	70k	f	-
27	95	45k	f	-
24	-	50k	m	PhD
32	-	60k	-	BE
-	60	-	-	ms
-	65	55k	-	BSc
40	22	-	m	BE

If outliers are present,
we can use median or
mean to fill the missing
values (depends on
usecase)

For the categorical
data, we can use
mode to replace
the missing values.

② Measure of Dispersion: ^{→ spread} [spread of the data]

- 1) variance (σ^2)
- 2) standard deviation (σ)

① Variance:-

Population variance

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

$x_i \rightarrow$ Data Points

$\mu \Rightarrow$ Population mean

$N \rightarrow$ Population size

Sample variance

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

$x_i \rightarrow$ Data Points

$\bar{x} \rightarrow$ Sample mean

$n \rightarrow$ sample size

* Assignment:

why we divide sample variance by $n-1$?

(unbiased estimator of population variance)

Example:

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

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

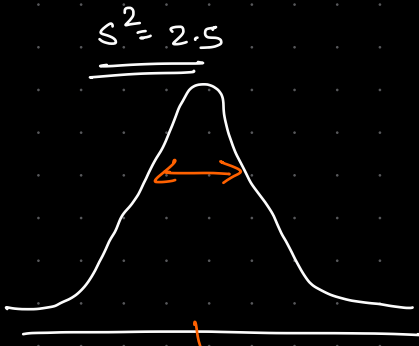
x_i	\bar{x}	$(x_i - \bar{x})^2$
1	3	4
2	3	1
3	3	0
4	3	1
5	3	4
$\bar{x} = 3$		$\sum (x_i - \bar{x})^2 = 10$

$$s^2 = \frac{10}{4} = \frac{5}{2} = 2.5$$

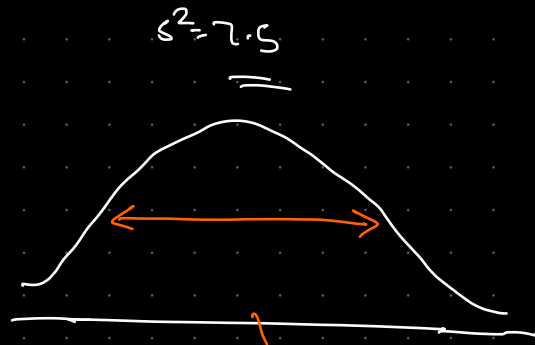
consider x distribution and y distribution

lets say,

$x = \{ \dots \}$



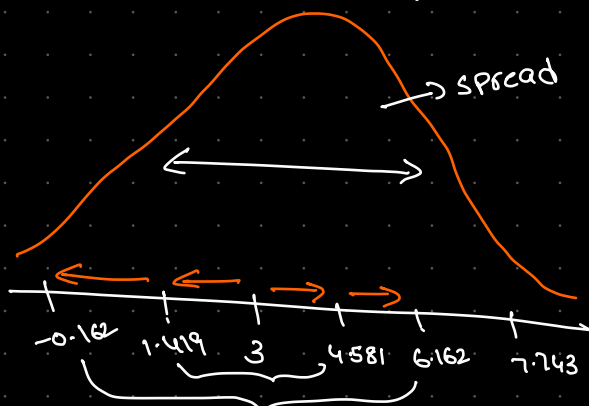
$y = \{ \dots \}$



② Standard deviation

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

Sample std $s = \sqrt{s^2}$



→ (Standard deviation talks about how much far element is away from mean)

→ (Variance is std variation square)

Assignment:

① Variance Error

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

$$\bar{X} = 3$$

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

$$= \sqrt{2.5}$$

$$S = 1.581$$

$$3.00$$

$$1.581$$

$$4.581$$

$$1.581$$

$$6.162$$

$$1.581$$

$$7.743$$

* Random variable:

$$\text{Linear Algebra} \left\{ \begin{array}{l} x + y = 7 \\ 8 = y + x \end{array} \Rightarrow x = 2 \right\} \text{variables}$$

$$\boxed{y = 6}$$

* Random Variable is a Process of mapping the output of a random process or experiment to a number *

Ex: Tossing a coin {head, tail} → Process

$$X = \begin{cases} 0 & \text{if head} \\ 1 & \text{if tail} \end{cases}$$

Rolling a Dice $\{1, 2, 3, 4, 5, 6\}$

$Y = \{\text{sum of rolling of dice 7 times}\}$

This helps in finding

$P_r(Y > 15) = \underline{\hspace{2cm}}$ (probability of x to be greater than 15).

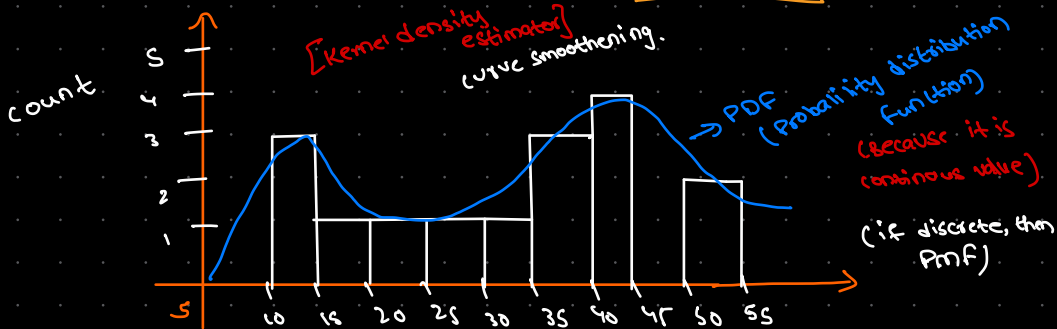
$P_r(Y < 10) : \underline{\hspace{2cm}}$ (probability of x to be less than 10).

* Histograms and skewness \rightarrow [Frequency]

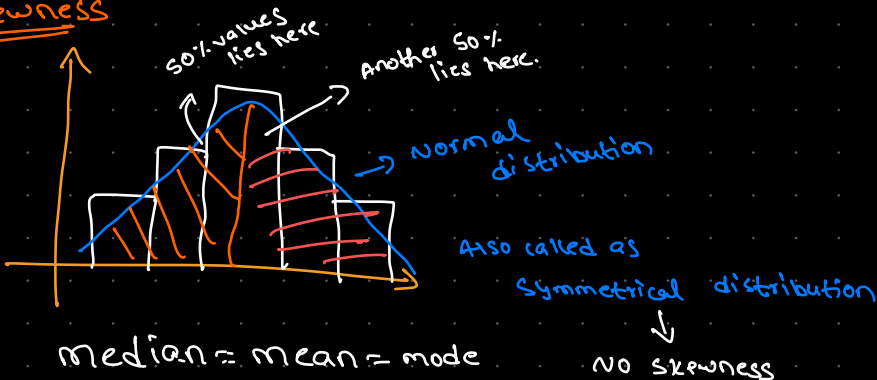
Ages: $\{10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 50, 51\}$

$\frac{50}{10} = 5 \Rightarrow$ (bin size)

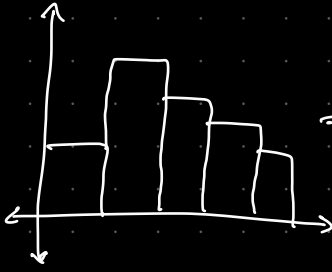
no. of bins = 10



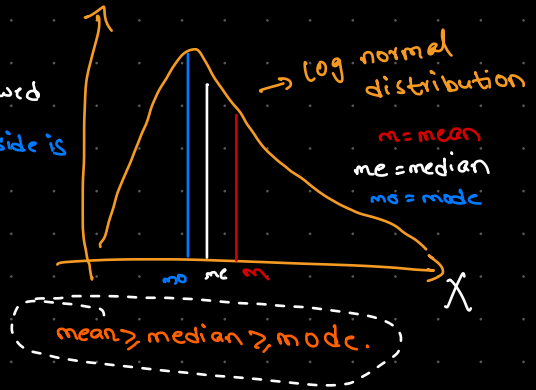
Skewness



* Right Skewed



⇒ Positive Skewed
(Right hand side is elongated)



* Left Skewed

