

(decision making process)

Probability and Statistics

Introduction to Statistics

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

Data: "facts & pieces of information"
eg: Height of students in classroom.
IQ of student in classroom.

Types of Statistics

Descriptive Statistics

It consist of organizing and summarizing data.

- (1) Measure of Central tendency.
- (i) Mean
 - (ii) Median
 - (iii) Mode

- (2) Measure of dispersion
- (i) Variance
 - (ii) Standard deviation

- (3) Diff type of distribution of data.
- eg: Histogram
Probability distribution function.
Probability Mass function.

Inferential Statistics

It consist of using data you have measured to form conclusion.

from Sample data



concluding



- (i) Z-test
- (ii) t-test

(iii) Chi Square test

Hypothesis Testing
 H_0, H_1
P-value
Significance value

g. There are 20 Statistics classes at your University and you have collected the heights of student in the class. Heights are recorded as
{ 175, 180, 140, 140, 125, 160, 135, 190 }

Descriptive Question

"What is the ^{avg} height of the entire classroom?"

Inferential Question

"Are the height of the students in the classroom similar to what you expect in the entire University?"

→ Population data

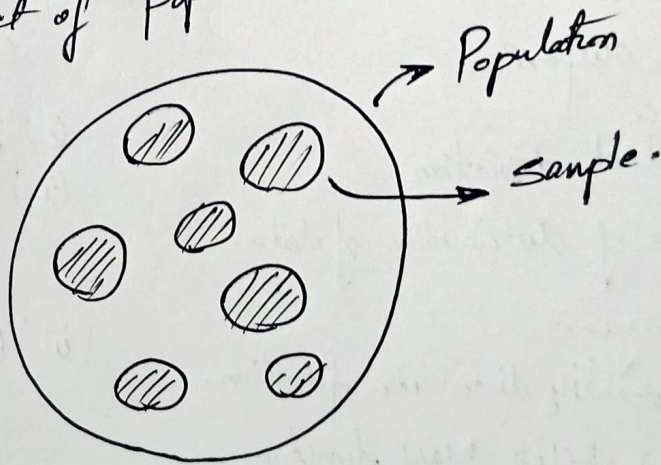
→ Sample data

Population data & Sample data.

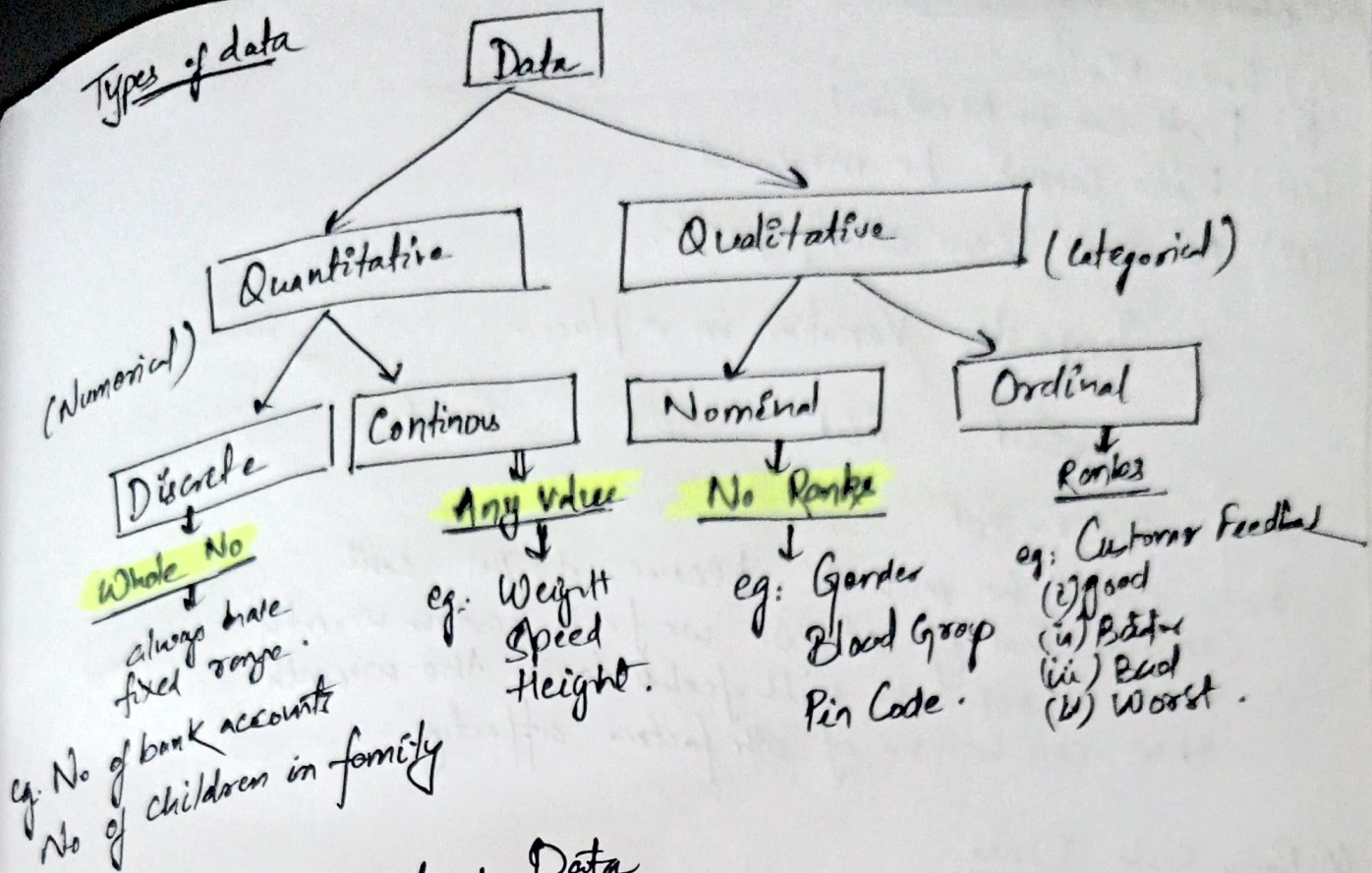
Population : The group you are interested in studying.

Sample : A subset of Population.

eg: Exit Poll



Types of data



Scale of Measurement of Data

- (1) Nominal Scale Data
 - (2) Ordinal Scale Data
 - (3) Interval Scale Data
 - (4) Ratio Scale Data
- (1) Nominal Scale Data
- (i) Qualitative
 - (ii) Categorical Data
 - (iii) Order/Rank doesn't matter

eg: Survey on favorite color.

(Here order doesn't matter on what color people like.)

- (2) Ordinal Scale Data:
- (i) Ranking is Important
 - (ii) Order Matters
 - (iii) Difference cannot be measured

eg: An feedback form

- (i) Best
- (ii) Good
- (iii) Bad
- (iv) Worst

→ Just on the basis of ranks we cannot calculate differences → like in this case we don't know the reason for this ranking

(3) Interval Scale of Data

(i) Order Matters

(ii) Diff can be Measured

(iii) Ratio Cannot be measured

(iv) No true Zero starting point.

e.g. Temperature Variation in a place.

30°F 60°F 90°F 120°F

Diff: $60 - 30 = 30^\circ\text{F}$

Ratio cannot be measured because it's just cont
Conclude that if at 30°F we feel certain warmth
then at 60°F we will feel double the warmth
there can be no of other factors affecting.

(4) Ratio Scale Data

① The Order matters

② Diff are measurable (include ratio)

③ Common '0' starting point.

e.g. grade of students

0, 90, 60, 30, 75, 40, 50

Measure of Central Tendency

- (i) Mean or Average
- (ii) Median
- (iii) Mode

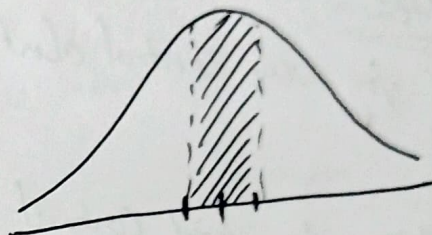
Mean Population (N)

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

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

Sample Mean (n)

$$\text{Sample Mean } (\bar{x}) = \frac{\sum_{i=1}^n X_i}{n}$$



Median

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

Step 1 Sort the random variable X . : $\{1, 2, 3, 4, 5\}$

2 No of elements count : 5

3 if $(\text{Count} \% 2) \neq 0$

find central element

if $\text{count} = 6$ $\{1, 2, \boxed{2, 3}, 4, 5\}$

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

or if $(\text{Count} \% 2) = 0$
find central element
 $2.5 \rightarrow$ 3
Median

Why Median?

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

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

Shift

if now introduce another data.

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

(pulling)

because it
does not
belong to
the
distribution.

$$\bar{x} = \frac{1+2+3+4+5+100}{6} = \frac{115}{6} = 19.1$$

(mean)

Now Median

try the central data in ① case $\rightarrow 3$

$$\textcircled{2} \text{ Case} = \frac{3+4}{2} = 3.5$$

* Median to find Central tendency
when outliers present.

Mode

frequency \rightarrow Max^m frequency.

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

Max^m frequency of an data: 1 to, Mode = 1

Application of Mean, Median, Mode in Feature Engineering

	Weight	Salary	Gender	Degree
Age				BE
24	70	40K	M	-
25	80	70K	F	-
27	95	45K	M	-
24	-	80K	-	PhD
32	-	60K	-	Master
-	60	-	-	Bsc
-	65	55K	M	-
-	72	-	F	-
40				

Numerical Value → We can find mean and fill null values.

→ If outlier → median value.

Categorical Value: We can fill it with mode values.

Measure of Dispersion (Spread of the data)

- (1) Variance
- (2) Standard deviation.

Version 1

Population Variance (σ^2)

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

x_i = Data points
 μ = population mean
 N = population size

Sample Variance

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

↓
 grp

x_i
 \bar{x} = sample mean
 n = sample size

Q Why we divide sample variance by $(n-1)$?

In order to create unbiased estimator } Bessel's correction

Q What is unbiased and biased estimator?

Unbiased: Like expected value is equal to the true value of the parameter

Biased: expected value is not equal to the true value due to bias

eg. If we pick 20 random students from a school and find avg

Unbiased } there is a good chance our estimator will give avg height true to the actual height

Biased } If we pick up students only from the basketball team where students are typically taller our avg value will not be equal to true value.