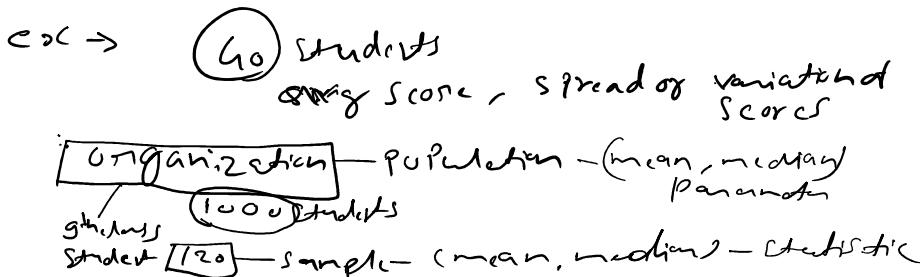


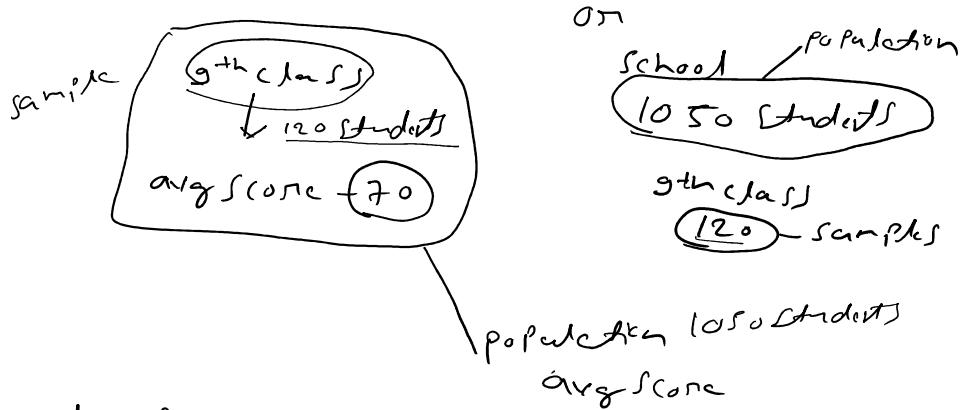
Descriptive Statistics

Collect
+
Organize
+
Summarize (analyze/describe)

Methods - mean, mode, median, standard deviation, variation, range, percentiles.



Inferential Statistics - Inferences or Prediction



12th class - 60 students

Sample size = 120
no. of data points in Sample

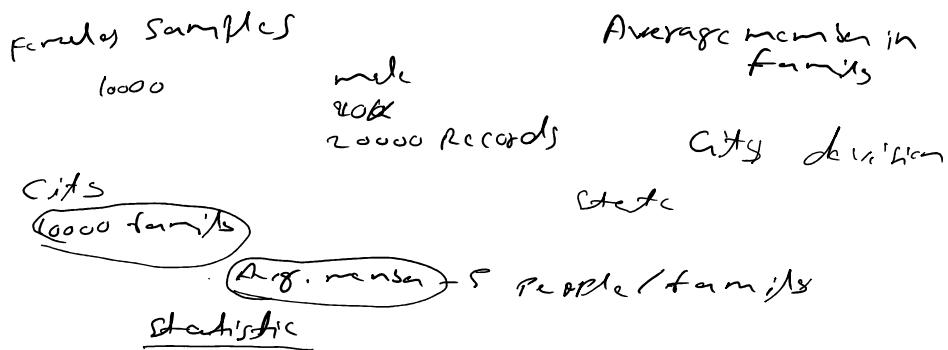
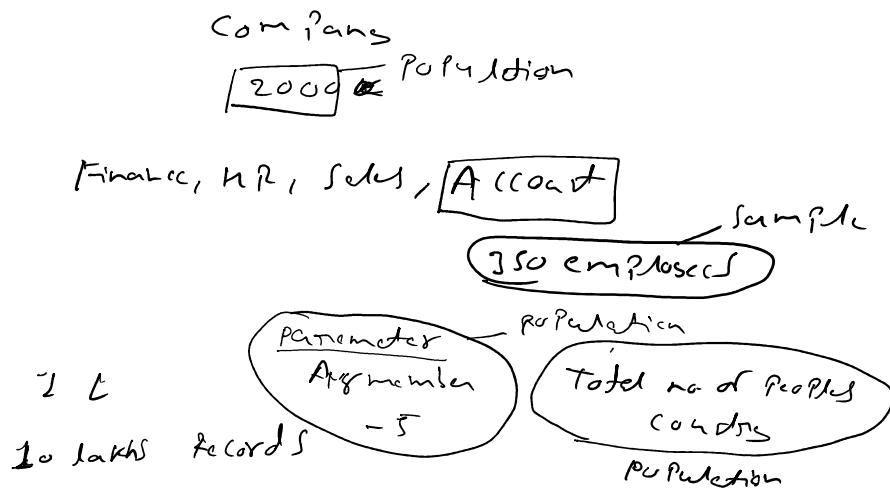
methods: regression analysis, confidence interval,

Variable column	Hypothesis				Target column
	BHK	Plot size	location	furniture	
1.	2	500	price	Kitchen	10 L
2.	1	450	Bang.	Doors	8 L
circular	3	600	musri	-	20 L

Price - Target column, label

BHK, Plot size, loc, typ. - independent column

1 BHK, 1000 m², Delhi, Doors - Price
Kitchen



- parameters
- 1) numerical characteristics of population
 - 2) population
 - 3) fixed value for entire population
- statistics
- numerical characteristics (Avg, mean, median...) of a sample
- statistics

School (100 students)

Avg Height - 160 cm - Parameter

12th class student 120
Avg Height of 120 students - 155 cm - statistics
 $160 - 155 = 5 \text{ cm}$

12th class
Avg Height of 180 students - 170 cm - statistics
 $160 - 170 = -10 \text{ cm}$

Variable and varies from one sample to another

2) fixed value for population

↳ fixed value for population | Variable and vary from one sample to another

Structured	Semi-structured	Unstructured
tabular form	<u>name</u> : Ram	Images, videos,
4 Rows, 5 columns	<u>age</u> : 19	Audio, text
	<u>id</u> : E01	
	subject: ["math", "bts", "sc", "english"]	
	Add: { country: India, state: "MH", city: "Pune" }	
		JSON, XML
		excel, CSV

(Quantitative Data)

Numerical Data: made of numbers

Age	Weight	Height
22	60 kg	100 cm
15	49	160 cm
20		
29	66	
69	50	
20, 10	100, 50 kg	

Categorical Data: made of words

(Qualitative Data)

Ex: colour	Name	Car brands
Brown	Ram	maruti
Green	Shyam	BMW
Red		Skoda

Numerical Data

↳ continuous Data: you can take any value within a range

19.5, 10 marks
12 days, 10 hrs, 10 min

Age	Weight	Height	Weight	Temp.
0 - 300	0 - 250 cm			
0.5				

	height	new & now fixed value
0 - 300	0 - 250 cm	weight
0.5	0.5	temp.
1.5	1	
30.5	162. cm	20.4 kg
19.5.2	165.5 cm	40.8 kg
	165 cm	35.5 °C

- 2) Discrete data:-
- 1) can take only distincts separate values,
 - 2) counted in whole no.
 - 3) finite no. of possible values

ex.

no. of student	Shoe size
120	2
130	3
90	4
	5
	6
	7
	8
	9
	9.5
	9.5 X

Categorical Data:

- 1) Nominal Data: no hierarchy or order in this data

Ex. :-

colours	fruits	name	gender
Blue	Apple	Ram	Female
Green	Banana	Shyam	male
Red	Pineapple	manoh	

- 2) Ordinal Data: when we get hierarchy or orders or ranking in this data.

education level	(class)	Classification
high school	8 th	excellent - 5
college	9 th	good - 4
graduation	10 th	fair - 3
master		poor - 2
		bad - 1

to talk

- 1) business problem - To find out best comfortable weather location all over world

- 2) longitude / latitude / highest temp. / lowest temp. / forest %

population	air humidity
71	81

Population

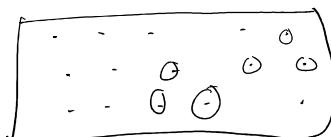
3) Sample size = 10000 cities

4) max. temp., min. temp. - Descriptive statistics

Probability Sampling - It randomly selects a sample from the population, ensuring that each member of the population has an equal chance of being selected in sample.

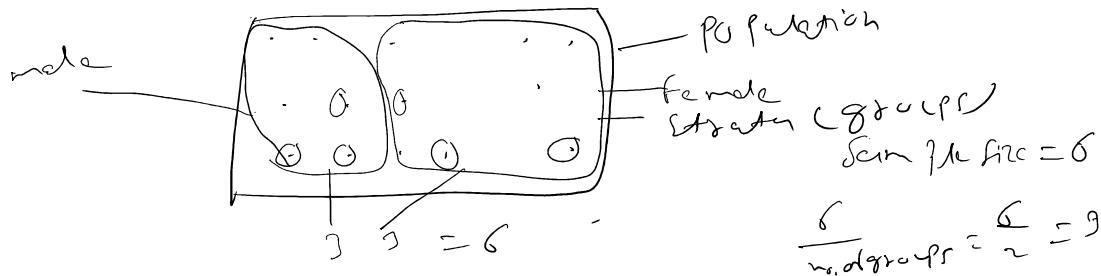
i) Simple Random Sampling - every datapoint has an equal chance of being selected

15 players Cricket team

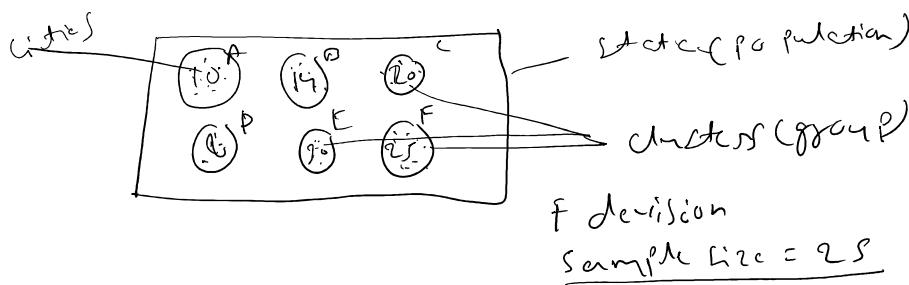


Sample size = 6

ii) Stratified Sampling :- Population is divided into groups (strata) and samples are randomly selected from each stratum.



iii) Clustering Sampling - population is divided into clusters, and entire clusters are randomly selected.



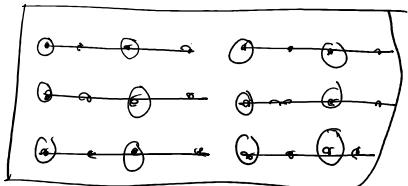
Making a probability decision

Jabaliyan

(12 cities)

- ii) Systematic Sampling - every n^{th} datapoints selected from a list after a random starting point.

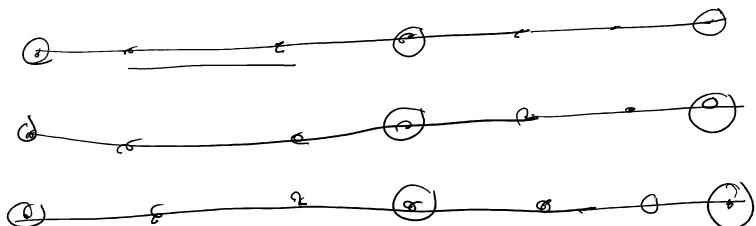
$c \in \mathbb{Z}^+$



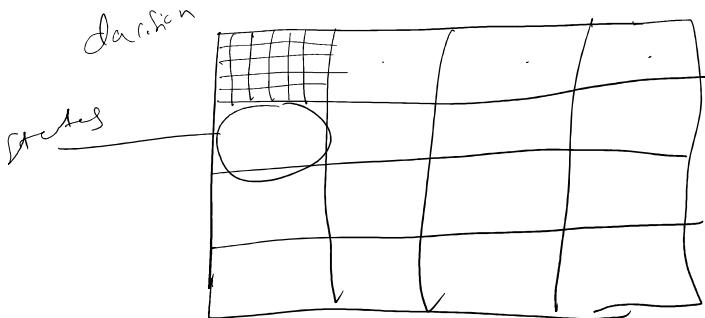
dance team = 12 students

example

restaurant feedback



- 5) multistage sampling -



countries
+
states
+
division
+
cits - people
+
ac

Non-probability sample - not every data points of the populations has an equal chance of being selected.

- 1) Convenience Sampling - select sample on availability or accessibility
- 2) Judgement Sampling - knowledge or judgments
- 3) Purposive Sampling -

Research wants to identify disease in dogs.
 dogs horse
 horse

$$\text{Variance} = \sqrt{9.01 \text{ min}^2}$$

standard deviation $\sigma = 3.$ min

$$\text{mean} = 14.4$$

$$A = [10, 18, 12, 15, 17]$$

(10) (18) Range = $18 - 10 = 8$

$\frac{-4.4 + (-2.4) + 0.6 + 2.6 + 7.6}{5}$

 $= \frac{\sigma}{\sqrt{n}} = 0$

$$\text{Avg Variation} = \text{Variance} = \frac{(x - \bar{x})^2}{n}$$

$$= \frac{(10 - 14.4)^2 + (18 - 14.4)^2 + (12 - 14.4)^2 + \dots}{5}$$

total no. of datapoint

$$\text{Variance} = \sqrt{9.01 \text{ min}^2} = 3.01$$

$$\text{std dev } \sigma = \frac{3 \text{ min}}{\sqrt{5}}$$

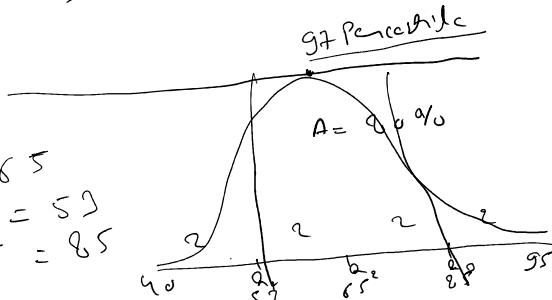
$$\text{Percentile} : = \frac{\text{number of values below } x \text{ (datapoint)}}{\text{total no. of values}} \times 100$$

original $[45, 49, 40, 53, 65, 71, 61, 79, 85, 91]$

sorted $[40, 45, 49, 53, 61, 65, 71, 79, 85, 91]$

$$\text{Percentile of } 71 = \frac{6}{10} \times 100 = 60 \text{ Percentile}$$

1000 dogs



$$Q_2 = 50^{\text{th}} \text{ Percentile} = 65$$

$$Q_1 = 25^{\text{th}} \text{ Percentile} = 53$$

$$Q_3 = 75^{\text{th}} \text{ Percentile} = 85$$

low kr/steal

$$[10K, 20K, 25K, 16K, 30K, 104]$$

... 100K

10K Km/sec

$$[10K, 12K, 15K, 17K, 18K, \text{ (10K)}]$$

Reasons for outlier

- 1) Human mistake
- 2) Machinery/measuring tool mistake
- 3) Experimental mistake

$$A = [10, 12, 15, 17, 18, \text{ (10K)}]$$

~~Max~~
57.2
~~Min~~
x

$$\mu = 14.4$$

$$\sigma = 3$$

$$\text{range} = 18 - 10 = 8$$

$$= 8$$

$$\text{after } 57.2$$

$$\mu = 57.2$$

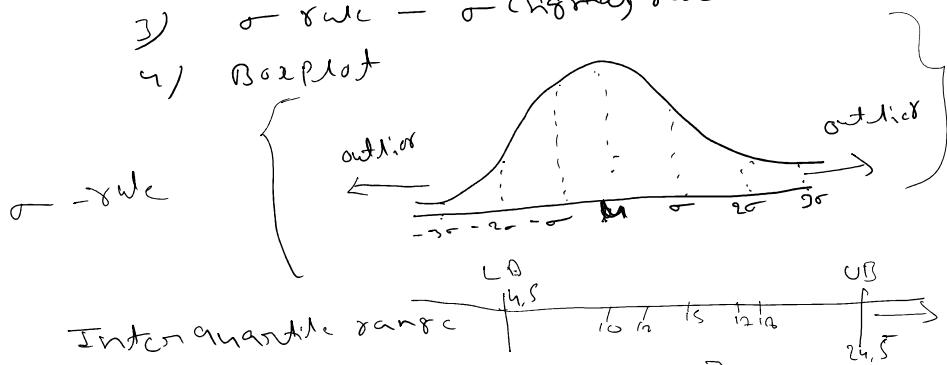
$$\sigma = 9$$

$$\text{range} = 108 - 10$$

$$= 90$$

Methods for identifying outliers - $\sigma = \sqrt{\frac{(x-\mu)^2}{n}}$

- 1) Interquartile range
- 2) σ -rule method
- 3) σ -rule - σ (signal) rule
- 4) Boxplot



$$[10, 12, 15, 17, 18]$$

$$\text{IQR} = Q_3 - Q_1, \quad Q_1 = 25^{\text{th}} \text{ Percentile} = \frac{25}{100} \times 8 = 1.25 \approx 2^{\text{nd}}$$

$$= 12$$

$$\frac{75}{100} \times 8 = 3.5 \approx 4^{\text{th}}$$

$$Q_3 = 75^{\text{th}} \text{ Percentile} = 17$$

threshold value

$$\text{Upper boundary UB} = Q_3 + 1.5(Q_3 - Q_1) = 17 + 1.5(17 - 12) = 17 + 7.5 = 24.5$$

$$= 17 + 1.5(5)$$

$$\text{Lower boundary LB} = Q_1 - 1.5(IQR) = 12 - 1.5(5)$$

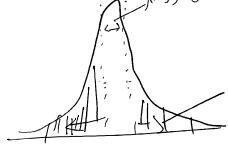
$$= 4.5$$

$$\text{UB} = Q_3 + 1.5(IQR)$$

$$\text{LB} = Q_1 - 1.5(IQR)$$

$K > 3$

less data points



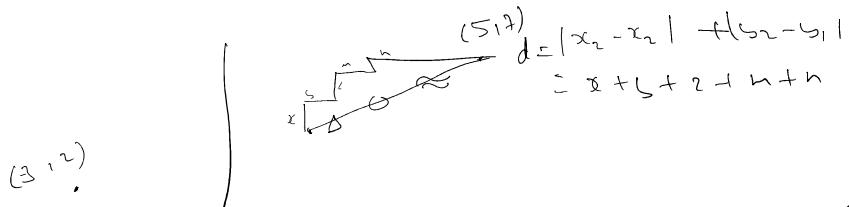
$K = 3$

mean = median
= mode

mean \neq median
mode

$$d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2} = \sqrt{9 + 4}$$

$$d = \sqrt{5^2 + 2^2} = \sqrt{29} \approx 5.4$$



(3, 4)

$$d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2} = \sqrt{(-1 - (-3))^2 + (2 - (-1))^2} = \sqrt{4 + 9} = \sqrt{13}$$

(1, 2)

(-3, -1)