

Welcome Back

Data  $\Rightarrow$  Base of analytics

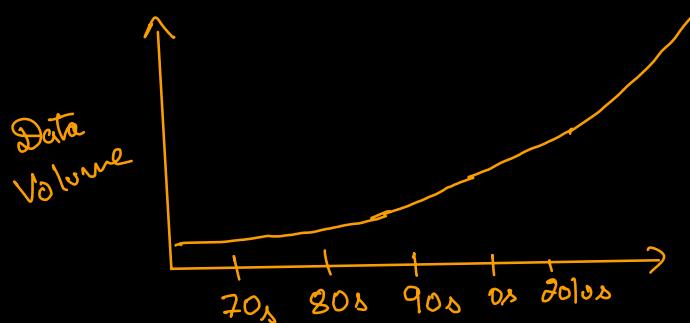
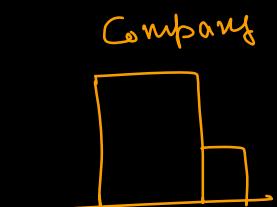
collection of information

any observation

data  $\xrightarrow{\text{science}}$  maths.

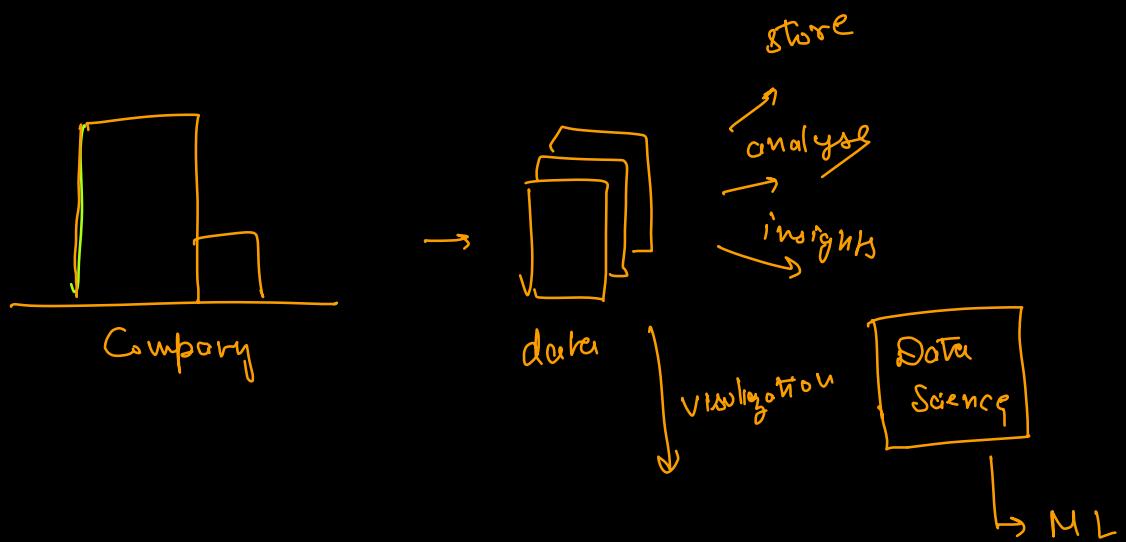
1. Statistics
2. Probability
3. Calculus.

$\hookrightarrow$  it is a concept/technique used to analyse & derive insights from your data.



If data is getting generated,  
we want to use it.

"Data is the new oil"



Tool we can use to analyze/ work on

data  $\Rightarrow$  Python

It makes money & good for business

1. Statistics

↳ Backbone of data science

2. Predictive Analytics

3. Advance Predictive Analysis.

## Statistics

Statistics is the science of conducting studies to collect organize summarize analyze draw conclusion out of your data.

Statistics deal with collective information data, interpreting this data & drawing conclusion from data

Used in many disciplines : Healthcare  
business  
E-commerce  
Education  
Insurance  
Marketing .

## Two areas of Statistics

1. Descriptive

2. Inferential

⇒ Descriptive

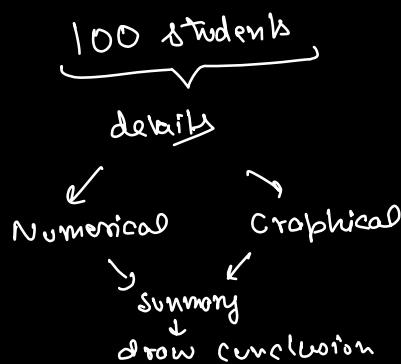
1. help to understand raw data and derive the data acc. to our business problem.

2. Numerical & Graphical to look for pattern in dataset and summarise the data, to present information in a convenient way

Numerical ⇒ mean, median, s.d., variance

Graphical ⇒ bar charts, histogram, pie charts

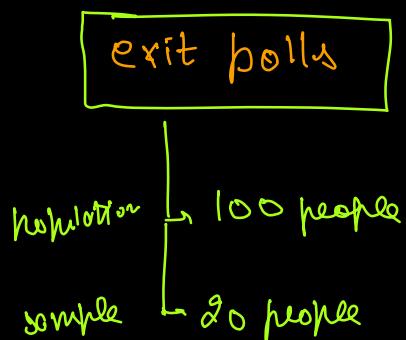
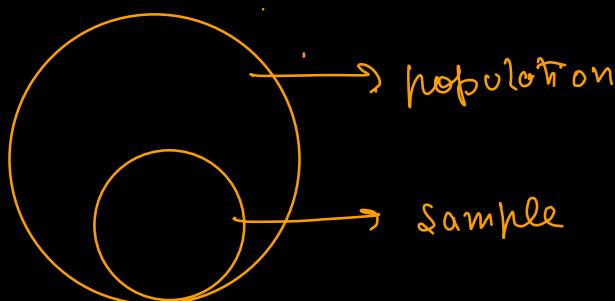
GNI PER CAPITA ANALYSIS USING WORLD BANK 2011 DATA	
Mean	11616.7
Standard Error	1361.2
Median	4200
Mode	430
Standard Deviation	17591.1
	3094480
Sample Variance	35.5
Kurtosis	4.99
Skewness	2.27
Range	88680
Minimum	190
Maximum	88870
Sum	1939990
Count	167



## Inferential Statistics

## Inference

⇒ Using sample data to make estimates, decision, prediction or other generalization about a larger set of data.



& draw conclusion

POLL OF EXIT POLLS					
	SEATS	MAJORITY	PROJECTED RESULTS		
GUJARAT	182	92	BJP 131	Cong + NCP 41	AAP 7
HIMACHAL PRADESH	68	35	BJP 36	Cong 29	AAP 0
DELHI MCD	250	126	AAP 155	BJP 84	Cong 7

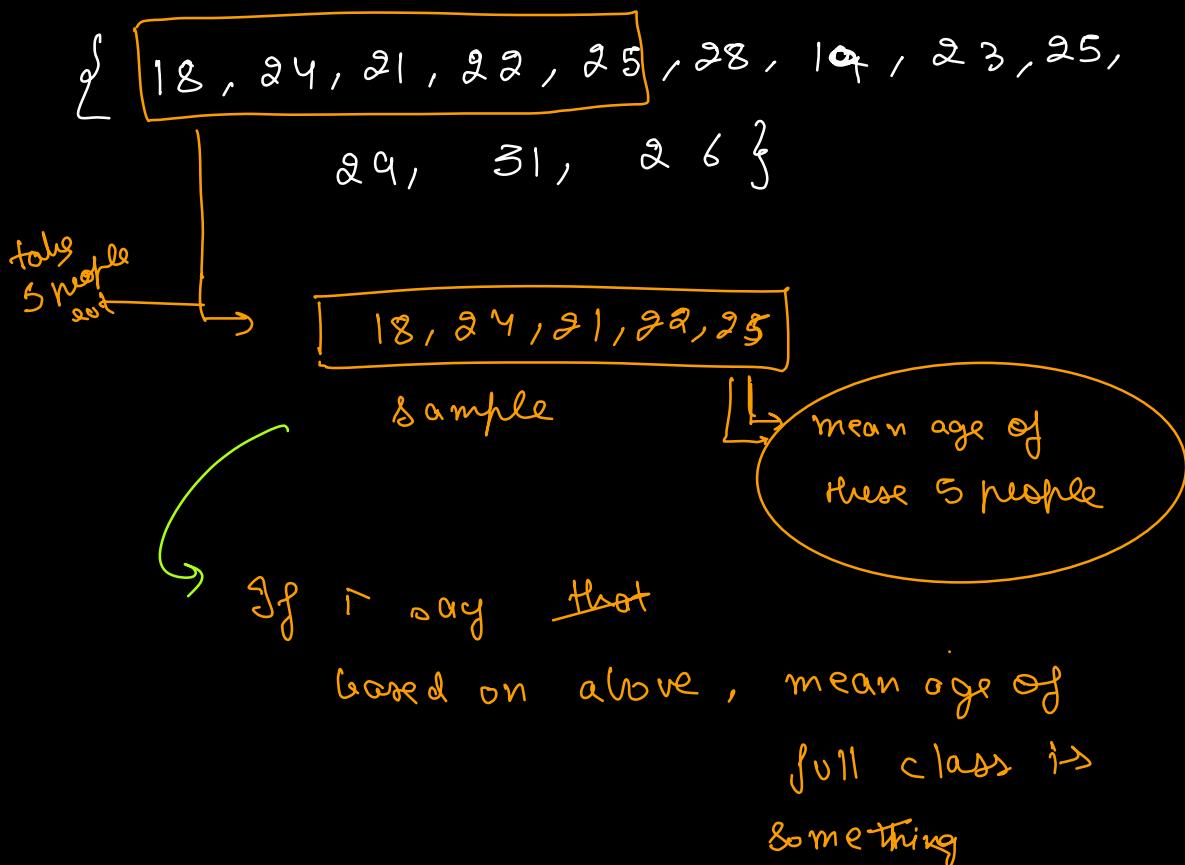
Source: Regional and National Channels  
 Health Warning: Aggregate of exit polls, these polls are not always accurate. Delhi MCD wards 'matched' to ensure a degree of comparability between 2017 & 2022 elections

#PollOfExitPolls 

- With Inferential statistics we are trying to reach conclusion beyond present data

sample's data → population's  
conclusion

### Class mean age example



### exit polls

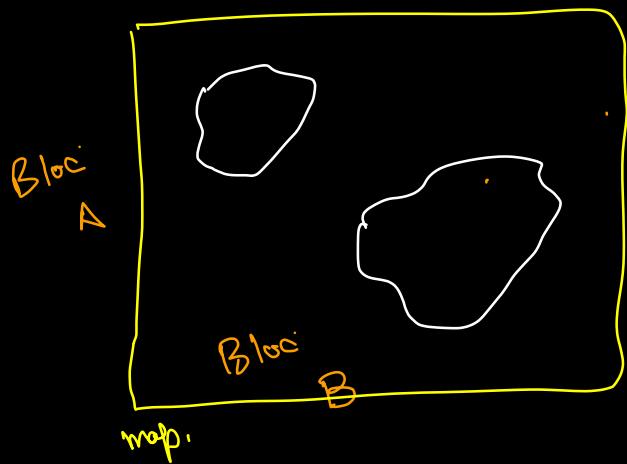
#### 1. Descriptive

2000 people participated

20 was age mean  $\Rightarrow$  mean of india

80% had smartphone

## Real life phone example



Area B

Apple 35%

android 60%

other 5%



because they are  
the same society  
blocks, with same  
no. of people with  
similar char.

⇒ The mobile distribution  
will be same

Grain example ⇒



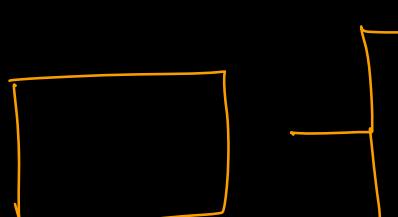
infer for  
whole population

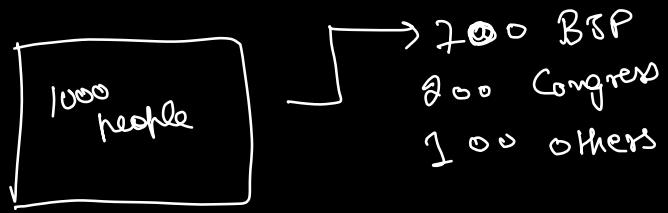
Blood Test  $\Rightarrow$    
 $\frac{\text{of disease.}}{\text{drops}}$        $\Delta\Delta$   
 $\Delta\Delta$   
 $\Delta\Delta$   
 $\Delta\Delta$   
 $\Delta\Delta$   
 $\Delta\Delta$   
 infer about  
 full body

Complete data  $\Rightarrow$  descriptive statistics

sample data  $\Rightarrow$  inferring about the population  
 inferential statistics

don't see data but the end result / use case.

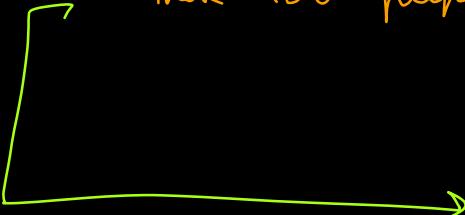
  
 Only this data  
 descriptive  
 Based on finding statistical  
 quantity, you are inferring  
 for larger population  
 Inferential



700/1000 people in this 1000 people

voted for BJP

there 700 people are rich

 70% people will vote  
for BJP in India

All voters of BJP will be  
rich

## Statistics Jargon

Population : A collection or set of individuals or objects or events whose prop. are to be analyzed.

Infinite or finite.

Orders on Amazon:

—————

4:00 —

4:30 —

Infinite pop<sup>r</sup> when new data keeps on getting generated.

iris

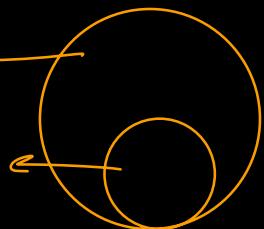
+

Orders on Amazon

titanic

population

sample



Sample — A subset of population.

It should be representative of population.

Variable :- A characteristic about each individual element of a population or sample.

Column.

Variable ↗

	Age	Name	Class	Marks	Pass
s1					
s2					

features, predictions, independent variables.

Experiment : A planned activity whose results yield a set of data.

Parameter: A numerical value summarizing all the data of an entire population

Statistic :- Numerical value summarizing the sample data

# Variable $\Rightarrow$ Types of variable

## Quantitative

measured numericals

no. of rows

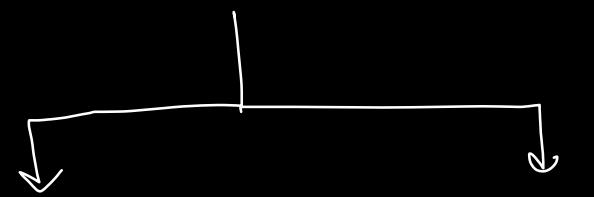
marks

no. of students in class

Age of an individual

Height of an individual

Population of country.



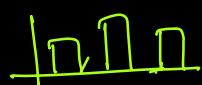
### Discrete

They can take specific, distinct values

Eg.

No. of students

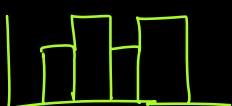
No. of classes



### Continuous

They can take any values  
Eg.  
Height of students

Temperature



## Qualitative [Categorical]

Gender  $\rightarrow M$   
 $\downarrow f$   
 $\downarrow$  Other  
century

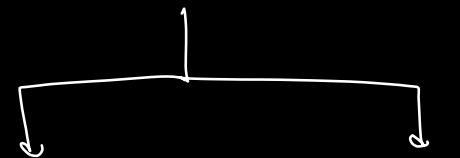
pass or not

eye/hair color

Martial Status

level of education.

Diversity preference



### Nominal

No intrinsic ordering,  
purely categorical

$\Rightarrow$  Gender  $\rightarrow M$

$\Rightarrow$  hair-color

$\downarrow \downarrow$   
B R Brown

### Ordinal

Ordering  
(w/o categories)

$\Rightarrow$  education  
quality factor

$\Rightarrow$  economy  
values

Output of a dice is Categorical

P

Categorical

1 2 3 4 5 6

### Models of iphone

age  
30 #

Categorical

20 25 30 35 40

### Age in diapers

age for vaccination

18, 20, 21  
— 80

Can Quantitative be converted to Qualitative?

2 & 0 - 50 50 - 100 100+  
Small not smart smart smar ters

Money <10L 10L - 50L 50L +  
Economic Poor Middle Rich



1000 players.  
Rankey cricketer  
⇒ 0 15 105  
↓ ↓ ↓

## Central Tendency

A measure of central tendency is  
descriptive statistic that describes  
the average or typical value for  
set of observation

- mean
- ↔ median
- mode
- std
- variance