

Statistics = I

Use Case

Product Based Company

IADFC Bank

(A)
Atm

(C)
Atm

(B)
Atm

Q. Should we need to open new Atm (C) or not?

② find the average size of the shark through the world?

Amazon Big Billion Day Sale ~~Init~~ Init?
which month should you select?

Life Cycle of Data Science Project

① Requirement gathering
↓

① Product manager / Proj man
② Business Analyst

② Data Analysts team

- ① Data Analyst
- ② Data scientist
- ③ Big Data engineers
- ④ Cloud Engineers

Where we get data

- ① Internal DataBase
- ③ web scraping

② 3rd Party API's

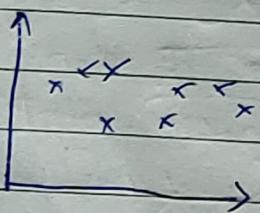
Sent Data to

Big Data Eng → ① MySQL data base
② NoSQL data base

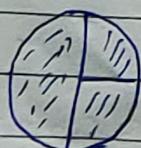
Now start life cycle of DS Project

- 1) EDA
- 2) FE
- 3) Feature Selection
- 4) Model Training
- 5) Hyper Parameter Tuning
- 6) Deploy

Analysis of Data



Descriptive Stats
⇒ Summarising Data



Descriptive Stats

Date _____
Page _____

Age = [12, 13, 14, 18, 20, 25] = Avg of Age

↓
Descriptive stats

↓

Measure of Central Tend

statistics \Rightarrow Statistics is the science of collecting, organising and analysing the data.

Data = Facts or pieces of information.

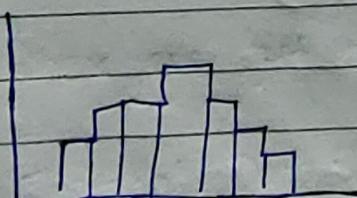
Eg: Ages of students in classroom
[24, 25, 32, 29, 28] \Rightarrow mean, median, standard deviation

Types of statistics

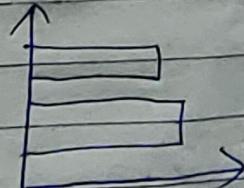
↓
Descriptive stats

↓
Inferential statistics

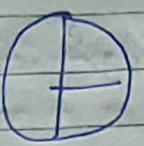
- ① It consists of organising and summarizing the data



Hist



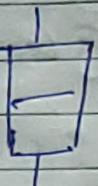
Bar chart



Pie



Distribution



Box plot

Box plot +



Scatter plot

② Inferential Stat

- ① It consists of Counting sample data and making conclusion about population data using some experiments

Hypothesis testing

Example

University

Date _____
Page _____

Class A] → 60 people



Sample data ⇒ [Age] ⇒ Average age of the con-

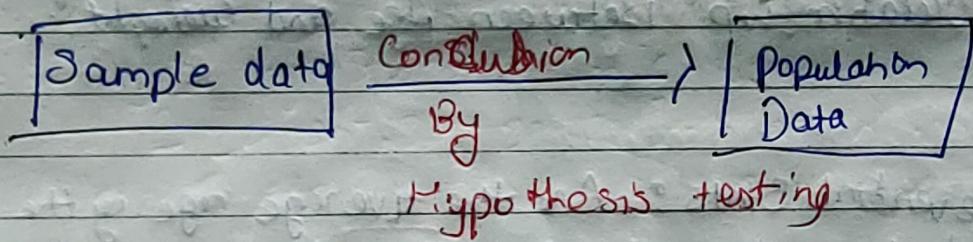
C.I = Confidence Interval

1) Z test

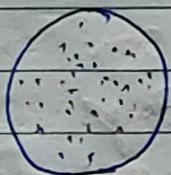
2) t test

3) Chi square test

4) F test



Sample Data vs Population Data



Population Data

Total data is population

Sample data

Some random sample we will take from population.

*Date _____
Page _____*

~~STUDY~~ 2009

Ex: let's say there are 20 classrooms in a university and you have collected the age of students in one classroom.

Age [21, 20, 18, 34, 17, 22, 25, 26, 23, 22]

Weight [- - - - -]

Descriptive Stats :-) What is the average age of student in the classroom?

2) Relationship Between Age and Weight?

Inferential Stats: Are the average age of the students in the classroom bigger than the average age of the students in the university?

Greater, Equal

11000 Students	50%	50%
	girls	Boys
	9.5%	90%

Population (N)

Sample (n)

Q) Sampling techniques : (n)

1) Simple Random Sampling : Every number of the population (N) has an equal chance of being selected for your sample (n)



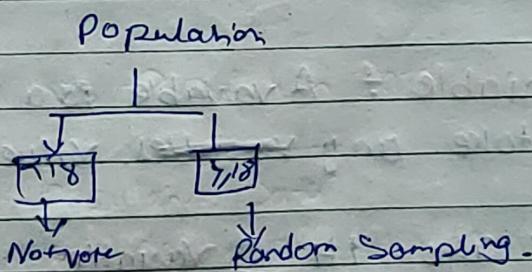
Goit \rightarrow Random Sampling
Doll General Survey

2) Stratified Sampling

Strata \rightarrow layers \rightarrow clusters \rightarrow Groups

Gender [male
Female]

Education [High School
master
post]



3) Systematic sampling \rightarrow [Airport]

[Credit Card] Sale persons

\rightarrow select every nth individual out of population (N)

Convenience Sampling \Rightarrow Only those who are interested in the survey will only participate.

{ Data Science Survey \rightarrow General AI Survey ? }

① Survey Regarding New Technology

convenience Sampling

② RBI Survey \rightarrow [mailed women] \Rightarrow Stratified + Random Sampling

③ Credit Card \Rightarrow Stratified + Random Sampling.

① Variable : A variable is a property that can take any ~~variables~~ values.

Eg : age = 14 variables

age = 25 ages = [24, 25, 26, 27]

age = 68 collection

Two different types of variable

① Quantitative Variables \rightarrow measured numerically [mathematical operation]

Eg: Age, weight, height, rainfall (cm), temp

2) Qualitative Variables → Categorical variables
Because Based on some characteristics they are group together.

Eg: Gender, types of flower, types of movies.

Quantitative Variables

- (1) Discrete variables → whole number
- (2) Continuous variables → continuous number

Discrete Variables

Eg: No of children → whole

Eg: No of Bank Account.

[1, 2, 3, 4]

Continuous Variables

Eg: Height, weight, money in account.

Assessment in Variables?

- 1) what kind of variables is marital status?
→ Categorical Variable
- 2) what kind of variables is Granga given length?
→ Continuous Variable
- 3) what kind of variables is movie duration?
→ Continuous Variable
- 4) what kind of variable is Rinode?
→ Discrete Variable
- 5) what kind of variable is IQ?
→ Discrete Variable. discrete variable

Agenda

- ① Histogram
- ② Measure of Central Tendency
- ③ Measure of Dispersion
- ④ Percentiles And Quartiles
- ⑤ 5 Number Summary (Box plot)

① Histogram

Ages : [10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40] ~~40~~

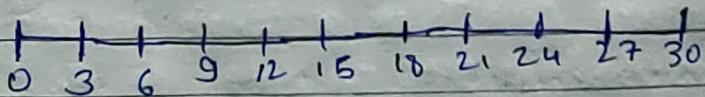
$$\text{Min} = 10$$

$$\text{Max} = 40$$

- ① Sort the Numbers
- ② Bins \rightarrow no. of groups
- ③ Bins size \rightarrow Size of Bins formula = $\frac{\text{Max} - \text{Min}}{\text{Bins}}$

$$\text{Bin} = 10$$

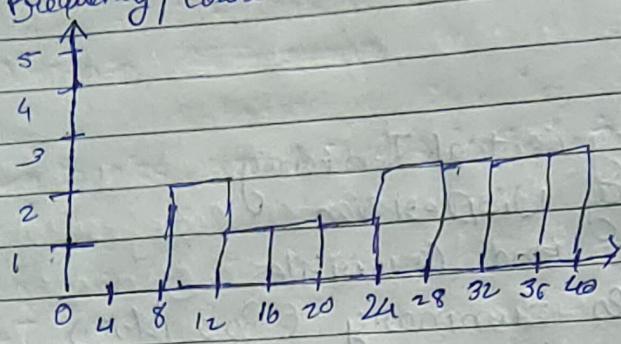
$$\text{Included} \rightarrow \frac{1^{\text{st}} \text{ no} - 10}{10} = \frac{30}{10} = 3$$



while calculating the Binsize do not include 1st number

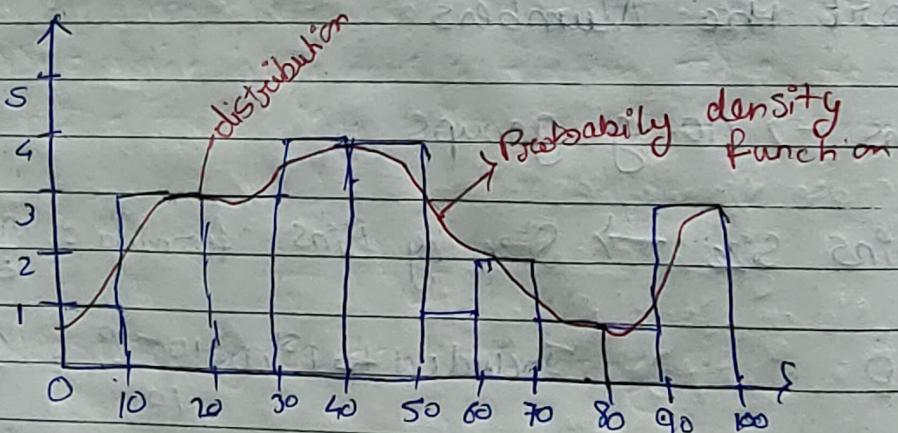
$$\frac{40}{10} = 4$$

Frequency / count

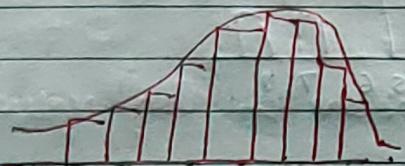
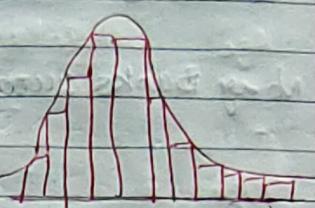
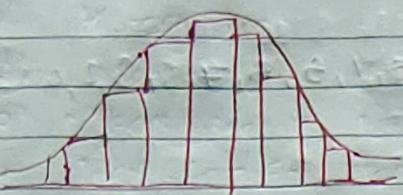


Ages : [10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 50, 51, 65, 67
78, 90, 95, 100]

Bin: 10

Bin Size $\frac{100}{10} \leftarrow$ no of Bins

Types of histogram

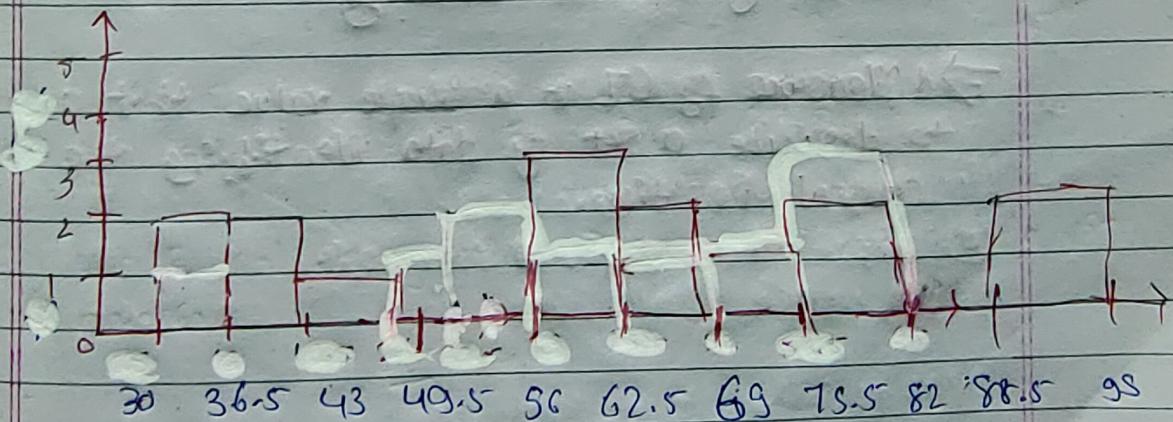


Ex: 2

weight: [30, 35, 38, 42, 46, 58, 59, 62, 63, 68, 75, 77, 80, 90, 95]

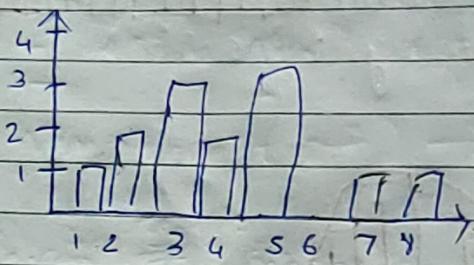
$$\text{Bins} = 10$$

$$\text{Bins Size} = \frac{95 - 30}{10} = \frac{65}{10} = 6.5$$



② Discrete ~~Continuous~~

No. of Banks' account = [2, 3, 4, 5, 1, 5, 3, 7, 8, 3, 2, 4, 5]



Probability mass function we use in Discrete

PDF = Probability density function \rightarrow continuous

pmf = Probability mass function \rightarrow Discrete

② Measure of Central Tendency

\Rightarrow A measure of CT is a single value that attempts to describe a set of data identifying the central position.

Let $x = [1, 2, 3, 4, 5]$

2) Mean : Average / mean $= \frac{1+2+3+4+5}{5} = \frac{15}{5} = 3$

Thus, 3 specify the central position of 5 the set by x
 Population mean (μ)

$$\mu = \frac{\sum_{i=1}^N x_i}{N}$$

$$\text{Sample mean } (\bar{x}) = \left[\frac{\sum_{i=1}^n x_i}{n} \right]$$

Population Age = [24, 23, 21, 28, 27]
 $\therefore N = 6$

$$\text{Population mean } (\mu) = \frac{24+23+21+28+27}{6} \\ = \frac{125}{6} = \mu = 20.83$$

$$\text{Sample mean } (\bar{x}) = \left[\frac{24+23+21+27}{4} \right]$$

$$\bar{x} = 23.5$$

$$\therefore [\mu > \bar{x}] \quad [\bar{x} > \mu]$$

Practical Application [feature Engineering]

Age	Salary	Family Size
-	-	Nan
-	-	-
Nan	-	-
-	Nan	-
-	-	Nan
Nan	-	-

We can't mean on Nan values.

- ② median : we use when there is outliers in Data.

$$Gx [1, 2, 3, 4, 5]$$

$$[1, 2, 3, 4, 5, 100]$$

$$\bar{x} = 3$$

$$\bar{x} = 19.16$$

Steps to find out median

- ① Sort the numbers
- ② Find the central Number.

- (1) If the no. of elements are even we find the average of central elements.
- (2) If the no. of elements are odd we find the central element.

Ex : Even so we find average of central element
 $[1, 2, 3, 4, \underline{5, 6}, 7, 8, 100, 120]$

$$\text{median} = \frac{s+6}{2} = 5.5$$

(3) MODE : Most frequent occurring elements

$[1, 2, 2, \underline{3, 3, 3}, 4, 5]$ mode = 3 $[1, 2, 3, 2, 2, 3, 3, 4, 5]$ mode = 2, 3

Data Sets

Type of flower

- | | | |
|---------------|---------------|----------|
| (1) Lily | (4) Nun | (7) Rose |
| (2) Sunflower | (5) Rose | (8) Nun |
| (3) Rose | (6) Sunflower | |

Note = MODE most of times we in Categorical Data or Variables

③ Measures of Dispersion

① Variance (σ^2) → talks about spread of Data

② Standard deviation (σ)

① Variance

Population Variance (σ^2)

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

Let's take a example

① [1, 2, 3, 4, 5, 6, 7, 8, 9, 10] ② [1, 2, 3, 4, 50, 60, 70]
If we find variance of both data then variance will be high on 2nd data set

How to calculate [1, 2, 3, 4, 5]

$$\mu = 3$$

$$\sigma^2 = [(1-3)^2 + (2-3)^2 + (3-3)^2 + (4-3)^2 + (5-3)^2] \\ = 4 + 1 + 0 + 1 + 4 = \frac{10}{5} = 2$$

① Sample Variance (S^2)

$$\Rightarrow S^2 = \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{n-1} \right)^2$$

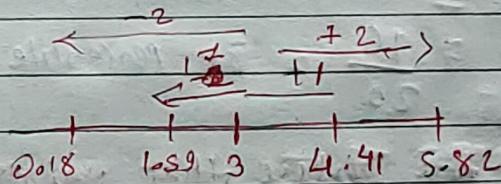
② Standard deviation ($\sqrt{S^2} = \sigma$)

$[1, 2, 3, 4, 5]$

$$\therefore M = 3$$

$$S^2 = 2$$

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



→ Percentiles and Quartiles

Percentile : $[1, 2, 3, 4, 5, 6, 7, 8]$

Percentage of even no. : $\frac{\text{No. of even numbers}}{\text{Total no. of numbers}} \times 100\%$

$$= \frac{4}{8} = 0.5 = 50\%$$

Percentiles: A percentile is a value below which a certain percentage of observation lie.

99 Percentile: It means the person has got better marks than 99% of the entire student.

Ex:

Datasets: [2, 2, 3, 4, 5, 5, 6, 7, 8, 8, 8, 8, 8, 9, 9, 10, 11]

→ What is the percentile rank of 10?

percentile rank of 36 = $\frac{\# \text{ no. of value below}}{n}$

$$= \frac{16}{20} = 80 \text{ percentile.}$$

→ What is the value that exists at 25 Percentile?

~~→ Value = Percentile $\times \frac{n}{100}$~~

$$= \frac{25}{100} \times 20 = 5^{\text{th}} \text{ Index}$$

Index always start with 0th.

~~formula.~~

$$\text{Value} = \frac{\text{Percentile} \times n}{100}$$

→ 5 number Summary

- ① minimum
- ② ~~first~~ first Quartile (25 percentile) Q1
- ③ median (50 percentile) Q2
- ④ third Quartile (75 percentile) Q3 } un use to Remove outliers
- ⑤ maximum

[1, 2, 2, 2, 3, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 6, 7, 8, 8, 9, 27]

We create fence were value ranging between
for finding outliers

[lower fence \leftrightarrow higher fence]

$$\text{Lower fence} = Q1 - 1.5 [IQR], IQR = Q3 - Q1$$

\downarrow
Inter Quartile Range

$$\text{Higher Fence} = Q3 + 1.5 [IQR]$$

$$Q1 = \frac{25}{100} \times n+1 = \frac{25}{100} \times 21 = 5.25 = 3$$

There is no 5.25 index so we will take average of 5th and 6th index

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