

Statistics

Day - I

Job
data analyst, business analyst, Data scientist

Statistics's Definition: statistics is the science of collecting, organizing and analyzing the data

Data: facts or pieces of information

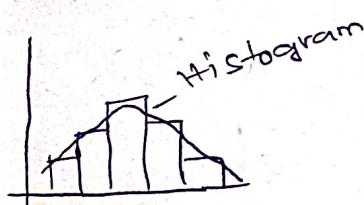
Eg: Ages of students in a class room

$$\{24, 25, 26, 28, 30, 42\} \Rightarrow$$

Statistics

Descriptive stats

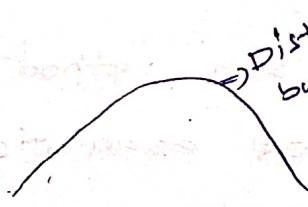
- ① It consists of organizing and summarizing the data.



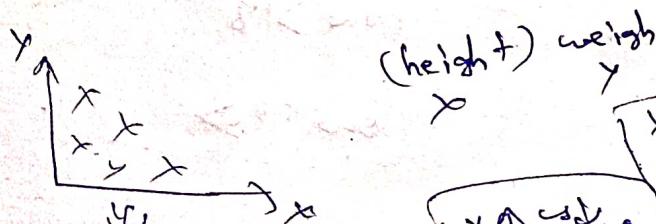
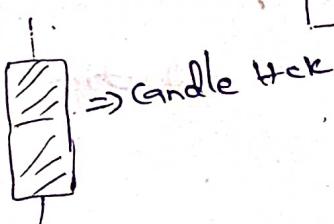
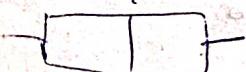
Bar chart



Distribution



Boxplot



$x \uparrow y \uparrow$
 $x \downarrow y \downarrow$

Inferential stats

- * It consists of collecting sample data and making conclusion about population data using some experiments

Hypothesis testing

University \rightarrow 500 people

Class A \rightarrow 60 people

!!

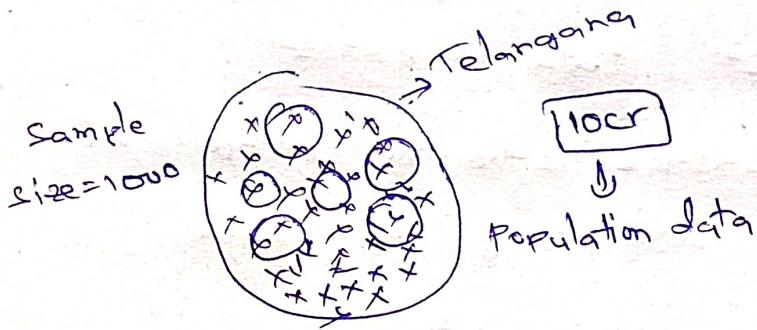
Sample data \Rightarrow Age \Rightarrow Avg age of entire university

!!

Hypothesis testing

c.i \Rightarrow confidence interval

Sample data Vs population data



Hypothesis testing

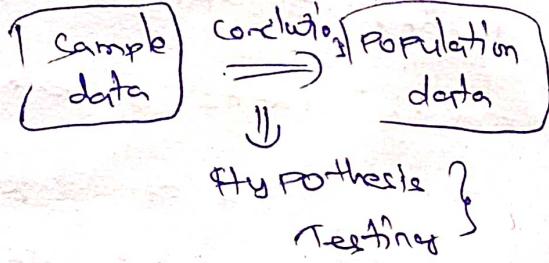
P-value

① Z-test

② T-test

③ Chi-square test

④ F-test



choose a sample data:

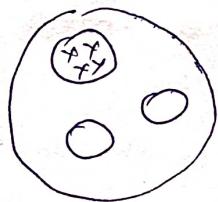
Sampling Techniques

Population (N)

sample (n)

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

$$n=1000$$



- ② Stratified Sampling: is a method of sampling that involves the division of population into smaller subgroups known as strata.

Example:

Gender	male	Female
--------	------	--------

education

High school
— master
phd

Blood group

Vote — above > 18

below 18

③ Systematic Sampling is a probability sampling method where researchers select member of the population at a regular interval
Ex: by selecting every 15th person on a list of population
* select every nth individual out of population (N)

5th — person → 9th person.

credit card sell

④ Convenience Sampling of only those who are interested in the survey will only participate

Data science surveys → General AI surveys
+ Survey regarding New technology

① Variable: A variable is a property that can take any value.
Ex: age = 14
age = 25
age = 100

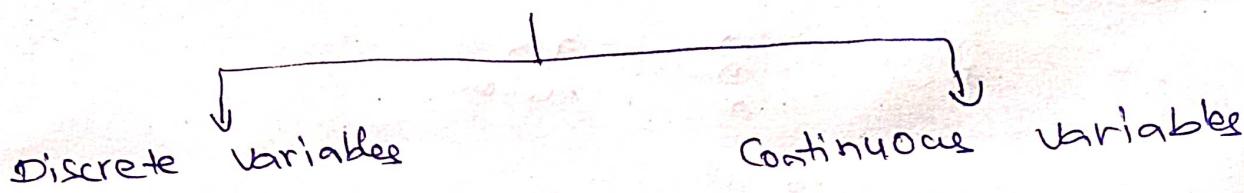
variable
age = [24, 25, 26, 27, 28]

Two different types of variables

① Quantitative Variable
→ measured Numerically & mathematical operations
Ex: age, weight, height, rainfall (cm), temperature

② Qualitative variables or categorical variables
Based on some characteristics they are grouped together.
Eg: Gender, types of flowers, types of movies

Quantitative Variables



Eg: Whole number → fixed

* No. of Bank accounts
2, 1, 2, 2, 4, 5³

* No. of children

Eg: Continuous → decimal value

* Height, weight, age
Rainfall, speed

* married → { married
not married
categorical variable }

Statistics Day-2

Agenda's

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

Histogram is a bar graph-like representation of data that buckets a range of classes into columns along the horizontal axis, the vertical Y-axis represents the number count or percentage of occurrence in the data for each column.

Ex:

Age = 10, 12, 14, 18, 22, 24, 28, 30, 32, 36, 40, 44, 46

- ① Sort the Numbers
- ② Bins \rightarrow No. of groups
- ③ Bins size \rightarrow size of bins

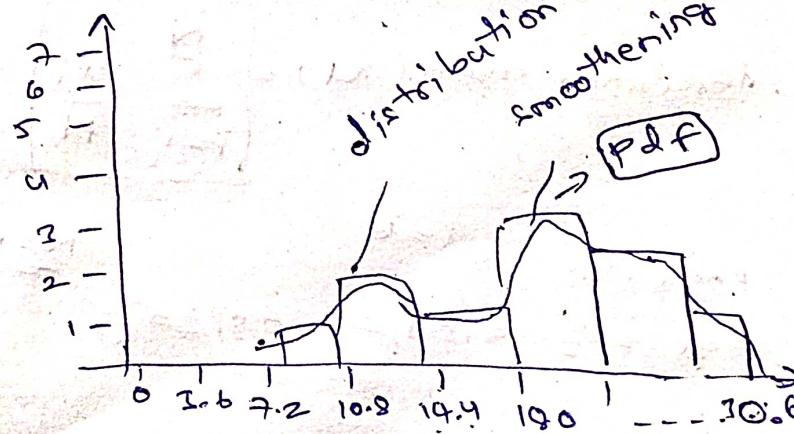
$$\text{bins} = 10$$

$$\text{bin size} = \frac{\text{max-min}}{\text{bins}}$$

$$= \frac{46-10}{10} \Rightarrow \frac{36}{10}$$

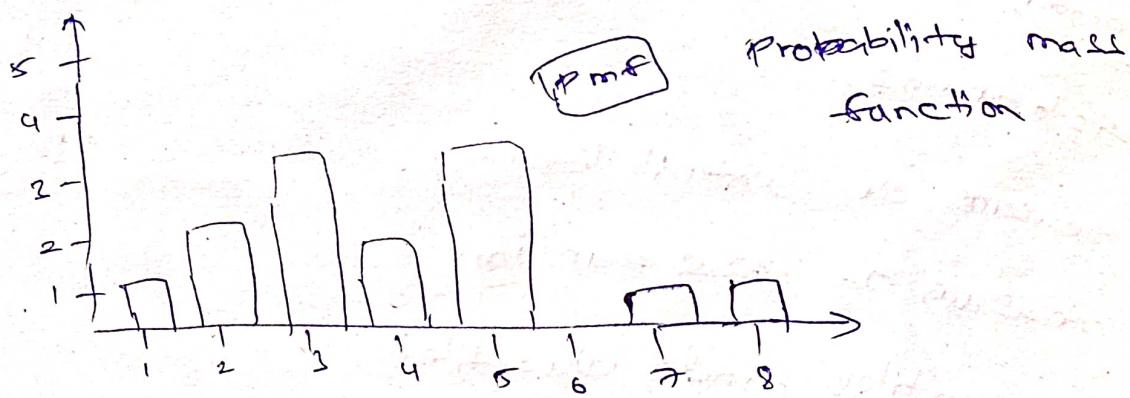
$$\text{bin size} \Rightarrow 3.6$$

(continuous values)



④ Discrete values

No. of Bank accounts = {2, 3, 5, 1, 4, 5, 3, 7, 8, 3, 2, 4, 5}



probability mass function

PDF \Rightarrow Probability density function \Rightarrow Continuous

PMF \Rightarrow Probability mass function \Rightarrow discrete

* Measure of Central Tendency :-

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

① Mean ② Median ③ Mode

① mean :-

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

$$\text{Average/mean} = \frac{1+2+3+4+5}{5} = \underline{\underline{3}}$$

Population (N)

Population mean (\bar{N}) =

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

Sample (n)

Sample mean (\bar{x}) =

$$\boxed{\sum_{i=1}^n \frac{x_i}{n}}$$

Population $N=6$

$$\text{Age} = \{24, 23, 21, 28, 27\}$$

$$\underline{\underline{N=n}}$$

$$n=4$$

$$\{24, 23, 21, 28\}$$

Population

$$\text{mean}(N) = \frac{24+23+21+28+27}{6}$$

$$M = 17.5$$

Sample mean (\bar{x})

$$= \frac{24+23+21+27}{4} \quad \bar{x} = 13.5$$

$$\text{Sample avg} = \{24, 23, 21, 27\}$$

$$\underline{\underline{M = 17.5}}$$

$$\boxed{M \geq \bar{x}} \\ \text{or} \\ \bar{x} \geq M$$

$$\frac{24+23+28+27}{4} \Rightarrow \frac{2^{\cancel{1}}}{\cancel{4}} \cancel{\times}$$

$$\bar{x} = 23$$

`np.nan` \Rightarrow Null value

`NAN` \Rightarrow Not a Number

$$\text{Mean} = \frac{\text{Sum of Observation}}{\text{Total No. of Observations}}$$

<u>Age</u>	<u>Salary</u>	
24	95	
28	50	less Information
29	<u>NAN</u>	
30	60	mean
31	75	
36	80	
<u>NAN</u>	<u>NAN</u>	
39	<u>1200</u>	outlier

* median:

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

$$\{1, 2, 3, 4, 5, \cancel{12}, \cancel{100}\}$$

$$\bar{x} = 19.16$$

$$\bar{x} = 3$$

steps to find out median:

① sort the numbers

② find the central number

→ 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 elements?

$$\frac{1+2+3+4+5+60}{6} \Rightarrow 19.16$$

Sorted

$$\{0, 1, 2, 3, 4, \boxed{5, 6}, 7, 8, 100, 120\} \quad \text{mean} = 25.6$$

$$\text{median} = \frac{5+6}{2} = 5.5 \quad \text{median} = 5$$

if n is odd = median = $(\frac{n+1}{2})^{\text{th}}$ term

$$\text{if } n \text{ is even: } \text{median} = \frac{(\frac{n}{2})^{\text{th}} \text{ term} + (\frac{n}{2}+1)^{\text{th}} \text{ term}}{2}$$

③ Mode → {most frequent occurring elements}

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

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

123

Date's Types of flower → (categorical variable)

lily

sunflower

Rose

NAN ← Rose

Rose

sunflower

Rose

NAN ← Rose

Measure of Dispersion

① Variance (σ^2)

② Standard deviation (σ)

$$x = \{1, 2, 3, 4, 5\} \quad M=3$$

Variance

Population Variance (σ^2)

$$\sigma^2 = \frac{\sum_{i=0}^N (x_i - M)^2}{N}$$

Sample Variance (s^2)

$$s^2 = \frac{\sum_{i=0}^n (x_i - \bar{x})^2}{n-1}$$

First One

$$\{1, 2, 3, 4, 5, 6, 7, 8, 9, 10\}$$

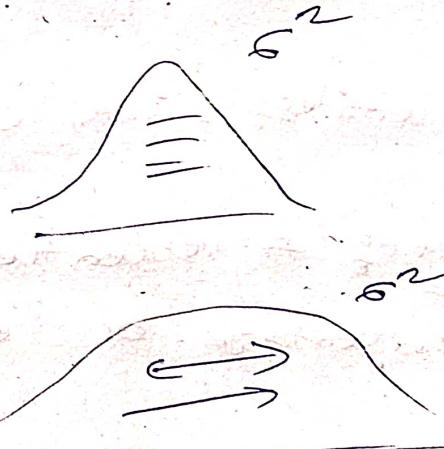
Variance

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

$$M=3$$

$$s^2 = \frac{(1-3)^2 + (2-3)^2 + (3-3)^2 + (4-3)^2 + (5-3)^2}{5}$$

$$\Rightarrow \frac{4+1+0+1+4}{5} = \frac{10}{5} = 2 = s^2$$



Second One

$$\{1, 2, 3, 4, 50, 60, 70, 100\}$$

Variance

$$\{1, 2, 3, 4, 5, 6, 80\}$$

$$M = 14.4$$

$$s^2 = \frac{(1-14.4)^2 + (2-14.4)^2 + \dots + (80-14.4)^2}{7}$$

$$s^2 = 219.10$$

Variance ↑↑ \Rightarrow spread ↑↑

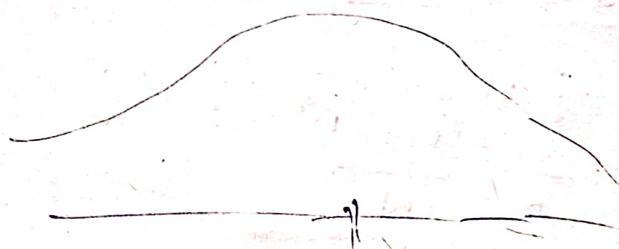
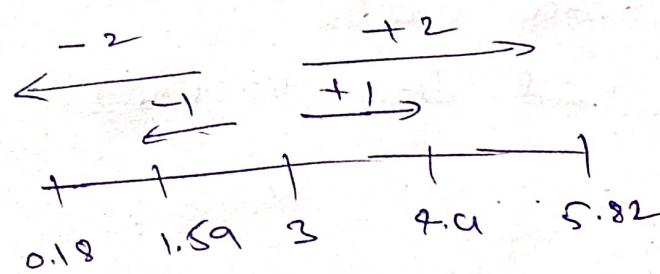
* Standard deviation ($\sqrt{6^2}$) $\Rightarrow 6$

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

$N = 5$

$$\sigma^2 = 2$$

$$6 \times \sqrt{2} = 1.41$$



* percentile and Quantiles%

Percentage = $\{1, 2, 3, 4, 5, 6, 7, 8\}$

$$\text{Percentage of even numbers} = \frac{\text{No. of even numbers}}{\text{Total no. of Numbers}} = \frac{4}{8} = 0.5 \Rightarrow 50\%$$

Percentiles% Gate, Cat, IITs.

Defn% A percentile is a value below which a certain percentage of observation lie

99 percentage% If means the person has got better marks than 99% of the entire students

Dakent% 2, 2, 3, 4, 5, 5, 5, 6, 7, 8, 8, 8, 8, 9, 9, 10, 11, 11, 12

what is the percentage rank of 10

$$\text{Percentile rank of } x = \frac{\# \text{No. of values below } x}{n} = \frac{16}{20} = 80 \text{ percent} \\ = 0.8$$

* what is the value of that exists at 25 percentile

$$\begin{aligned} \text{Value} &= \frac{\text{Per Centile}}{100} * (n+1) \\ &= \frac{25}{100} * 20 = 5^{\text{th}} \text{ Index} \\ &\quad \text{or } p = \frac{5}{2} \end{aligned}$$

* Number Summary:

① minimum

② first Quartile (25 percentile) (Q_1)

③ median

④ Third Quartile (75 percentile) (Q_3)

⑤ maximum

Boxplot

→ Remove
the outlier
→ r_c

3, 1, 2, 12, 2, 22, 2, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 12, 22

[Lower fence \leftrightarrow Higher fence]

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

$$\text{IQR} = Q_3 - Q_1$$

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

Inter Quartile Range
(IQR)

$$Q_1 = \frac{25}{100} \times 21 = 5.25 \text{ Index} = 3$$

$$Q_3 = \frac{75}{100} \times 21 = 15.75 \text{ Index} = \frac{8+7}{2} = 7.5$$

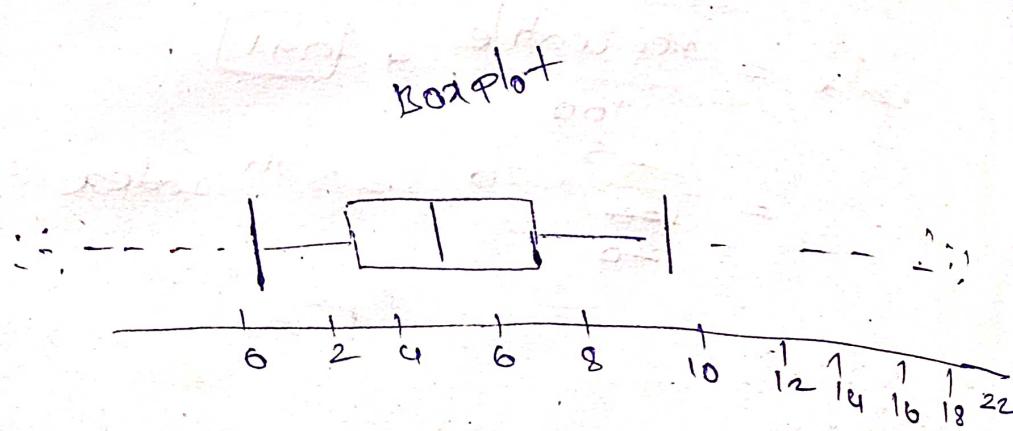
$$\text{Lower fence} = 3 - (1.5)(4.5) = -3.65$$

$$\text{Higher fence} = 7.5 + (1.5)(4.5) = 14.25$$

$\{1, 2, 2, 2, 2, 2, 3, 4, 5, 5, 5, 6, 6, 6, 7, 8, 8, 9, 29\}$

5 Number Summary

- ① minimum = 1
- ② $Q_1 = 3$
- ③ median = 5
- ④ $Q_3 = 7.5$
- ⑤ maximum = 9



To Treat + outliers