

P-P1 60, 70, 80, 75, 65, 70, 80, 70, 65, 65

(Ans) sorted Order :— 60, 65, 65, 65, 70, 70, 70, 75, 80, 80

Here,

$$n = 10$$

$$n = 10$$

$$\frac{(n)}{2} \text{ th value} + \frac{(n+1)}{2} \text{ th value}$$

$$\text{mean } (\bar{x}) = \frac{\sum x_i}{n}$$

$$= \frac{700}{10}$$

$$= 70$$

(Ans)

$$= \frac{5 \text{ th value} + 6 \text{ th value}}{2}$$

$$= \frac{70 + 70}{2}$$

$$= 70$$

(Ans)

P-P2

60, 70, 80, 75, 65, 70, 100, 70, 65, 65

$$\text{mean } (\bar{x}) = \frac{720}{10}$$

$$= 72$$

(Ans)

For Truncated mean :

sorted the dataset : 60, 65, 65, 65, 70, 70, 70, 75, 80, 100

For  $P = 1$  so, 60 and 100 will be discarded.

$$\text{mean } (\bar{x}) = \frac{\sum x_{i+1}}{(n-2P)}$$

$$= \frac{560}{10-2}$$

$$= \frac{560}{8}$$

$$= 70$$

(Ans)

Weighted mean :  $x_i = 40, 45, 80, 75$  and 10 have

$$w_i = 1, 2, 3, 4 \text{ and } 5$$

$$(\bar{x}_w) = \frac{\sum_{i=1}^n w_i x_i}{\sum w_i}$$

$$= \frac{(40 \times 1) + (45 \times 2) + (80 \times 3) + (75 \times 4) + (10 \times 5)}{1+2+3+4+5}$$

$$= \frac{720}{15} = 48$$

(Ans)

mean Absolute deviation :

$$MAD = \frac{\sum |x_i - \bar{x}|}{n}$$

dataset 1 : 30, 40, 50, 60, 70

$$\bar{x} = 50$$

$$MAD : \frac{|(30-50)| + |(40-50)| + |(50-50)| + |(60-50)| + |(70-50)|}{5}$$

$$= \frac{20 + 10 + 0 + 10 + 20}{5}$$

$$= \frac{60}{5}$$

$$= 12$$

### Median Absolute deviation:-

$$MAD = \text{median}(|x_1 - m|, |x_2 - m|, \dots, |x_n - m|)$$

dataset :- (2) 0, 25, 50, 75, 100

n = 5 (odd number)

$$\text{(Ans)} \Rightarrow MAD: \text{median}(|10-50|, |25-50|, |50-50|, |75-50|, |100-50|)$$

$$= \text{median}(40, 25, 0, 25, 50)$$

$$= \left(\frac{5+1}{2}\right) \left(\frac{6}{2}\right) 3\text{rd value}$$

$$= \text{median}(0, 25, 25, 50, 50)$$

~~Median  
of 25~~

$$n = 5 (\text{odd})$$

$$= \left(\frac{5+1}{2}\right) \text{th value}$$

$$= \left(\frac{5+1}{2}\right)^{\text{th}}$$

$$= 3\text{rd value}$$

$$= 25 \quad (\text{Ans})$$

### Variance :- (S<sup>2</sup>)

$$S^2 = \frac{1}{n} \sum (x_i - \bar{x})^2$$

Example:- The Height : 600 mm, 470 mm, 170 mm, 430 mm, 500 mm

$$\text{mean } (\bar{x}) : \frac{1970}{5} = 394$$

$$\text{variance } (S^2) : \frac{(600-394)^2 + (470-394)^2 + (170-394)^2 + (430-394)^2 + (500-394)^2}{5}$$

$$= 21704 \quad (\text{Ans})$$

$$\text{standard deviation } (S) = \sqrt{21704}$$

$$= 147.32 \quad (\text{Ans})$$

### sample variance (S<sup>2</sup>) :-

dataset 1 :- 30, 40, 50, 60, 70

$$n = 5$$

$$\bar{x} = \frac{210}{5} = 42$$

$$S^2 = \frac{(30-42)^2 + (40-42)^2 + (50-42)^2 + (60-42)^2 + (70-42)^2}{5-1}$$

$$= \frac{1000}{4}$$

$$= 250 \quad (\text{Ans})$$

standard deviation;  $s = \sqrt{S^2}$

$$= \sqrt{250}$$

$$= 15.82 \quad (\text{Ans})$$

dataset 2 :- 0, 25, 50, 75, 100

$$\bar{x} = \frac{210}{5} = 42$$

$$= 42$$

$$S^2 = \frac{(0-42)^2 + (25-42)^2 + (50-42)^2 + (75-42)^2 + (100-42)^2}{5-1}$$

$$= \frac{1250}{4}$$

$$= 312.5$$

standard deviation :-

$$s = \sqrt{312.5}$$

$$= 17.68 \quad (\text{Ans})$$

### Practice Example :-

dataset :— 8, 9, 10, 10, 10, 10, 11, 11, 11, 12, 13

$$\text{mean } (\bar{x}) : \frac{10.5}{10}$$

$$= 10.5 \quad (\text{Ans})$$

$n = 10$  (Even number)

$$\text{median} = \frac{\left(\frac{n}{2}\right)^{\text{th}} \text{ value} + \left(\frac{n}{2}+1\right)^{\text{th}} \text{ value}}{2}$$

$$= \frac{10+11}{2}$$

$$= 10.5 \quad (\text{Ans})$$

mode : (10, 11) Bi-modal dataset  $(\text{Ans})$

$$\text{Range} : (\text{max} - \text{min})$$

$$= (13 - 8)$$

$$= 5 \quad (\text{Ans})$$

$$\sigma^2 = \frac{(8-10.5)^2 + (9-10.5)^2 + (10-10.5)^2 + (10-10.5)^2 + (10-10.5)^2 + (11-10.5)^2 + (11-10.5)^2 + (12-10.5)^2 + (12-10.5)^2 + (13-10.5)^2}{10}$$

$$\approx 0.75 \quad \therefore \sigma = 1.085$$

$$\sigma = \sqrt{1.085}$$

$$= 1.086 \quad (\text{Ans})$$

~~MID of  
sample C2~~

sample variance

$$S^2 = \frac{(8-10.5)^2 + (9-10.5)^2 + (10-10.5)^2 + (10-10.5)^2 + (10-10.5)^2 + (11-10.5)^2 + (11-10.5)^2 + (12-10.5)^2 + (12-10.5)^2 + (13-10.5)^2}{(10-1)}$$

$$= 2.055$$

$$S = \sqrt{2.055}$$

$$= 1.43 \quad (\text{Ans})$$

Mean Absolute Deviation :— MAD :  $\frac{|8-10.5| + |9-10.5| + |10-10.5| + |10-10.5| + |10-10.5|}{10}$

$$= \frac{5+3+6}{10} = \frac{14}{10} = 1.4 \quad (\text{Ans})$$

Median Absolute Deviation :—

MAD : Median  $(|8-10.5|, |9-10.5|, |10-10.5|, |10-10.5|, |10-10.5|, |11-10.5|)$

$+ |11-10.5|, |11-10.5|, |12-10.5|, |12-10.5|, |13-10.5|$

$n=10$   
 $\frac{\text{sum of value}}{2} = \frac{10+11}{2} = 10.5$   
 $0.5 + 0.5$   
~~0.5 + 0.5~~

$$= \text{median } \{ 2.5, 1.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 1.5, 2.5 \}$$

$$= \text{median } (0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 1.5, 1.5, 2.5, 2.5)$$

Hossain 2021-1-60-071

Finding percentile:— 8, 9, 10, 10, 10, 11, 11, 11, 12, 13, 15, 16, 16, 17, 20

Total data point  $n = 15$

$$25^{\text{th}} \text{ percentile} = 25\% \times 15 \quad | \quad 50^{\text{th}} \text{ percentile} = 50\% \times 15$$

$$\Downarrow \quad = 3.75 + \text{value}$$

$$\textcircled{Q}_1 \quad = 10 + (10 - 10) \times 0.75$$

$$= 10 \quad \underline{\text{Ans}}$$

$$\Downarrow \quad = 7.5 + \text{value}$$

$$\textcircled{Q}_2 \quad = 11 + (11 - 11) \times 0.5$$

$$= 11 \quad \underline{\text{Ans}}$$

$$75^{\text{th}} \text{ percentile} = 75\% \times 15$$

$$\Downarrow \quad = 11.25 + \text{value}$$

$$\textcircled{Q}_3 \quad = 15 + (16 - 15) \times 0.25$$

$$= 15.25 \quad \underline{\text{Ans}}$$

$$\text{IQR} = (\textcircled{Q}_3 - \textcircled{Q}_1)$$

$$= (15.25 - 10)$$

$$= 5.25 \quad \underline{\text{Ans}}$$

$$95^{\text{th}} \text{ percentile} = 95\% \times 15$$

$$= 14.25 + \text{value}$$

$$= 17 + (20 - 17) \times 0.25$$

$$= 17.75 \quad \underline{\text{Ans}}$$

~~Box Plot~~

64, 64, 68, 69, 70, 71, 72, 72, 75, 75, 80, 81, 83, 85

$$\text{Min} = 64 \quad \textcircled{Q}_1 = 25^{\text{th}} \text{ percentile} = 25\% \times 14$$

$$\text{Max} = 85 \quad = 3.5 + \text{value}$$

$$= 68 + (69 - 68) \times 0.5$$

$$= 68.5$$

$$\textcircled{Q}_3 = 75^{\text{th}} \text{ percentile} = 75\% \times 14$$

$$= 10.5 + \text{value}$$

$$= 75 + (80 - 75) \times 0.5$$

$$= 77.5$$

↑ outliers

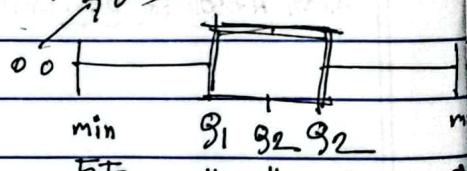
$$\text{IQR} = (77.5 - 68.5)$$

$$= 9$$

$$\text{Upper Extreme: } \textcircled{Q}_3 + \text{IQR} \times 1.5$$

$$= (77.5 + 9 \times 1.5)$$

$$= 91$$



$$\text{Lower Extreme: } \textcircled{Q}_1 - \text{IQR} \times 1.5$$

$$= 55$$

Any values greater than the upper extreme and some smaller than the lower extreme would be called outliers.

Date : 24.1.24

Max: 24.7  
Min: 0 IQR:  $Q_3 - Q_1 \Rightarrow (19.225 - 7.925) \Rightarrow 11.3$

$$\begin{aligned}
 Q_1 &= 25\% \times 11 \\
 &= 2.75 + \text{value} \\
 &= 5.6 + (8.7 - 5.6) \times 0.75 \\
 &= 7.925 \\
 Q_3 &= 75\% \times 11 \\
 &= 8.25 + \text{value} \\
 &= 19.2 + (19.3 - 19.2) \times 0.25 \\
 &= 19.225 \\
 Q_2 &= 50\% \times 11 \\
 &= 5.5 + \text{value} \\
 &= 14.1 + (15 - 14.1) \times 0.5 \\
 &= 14.55
 \end{aligned}$$

Upper Extreme:  $Q_3 + IQR \times 1.5$

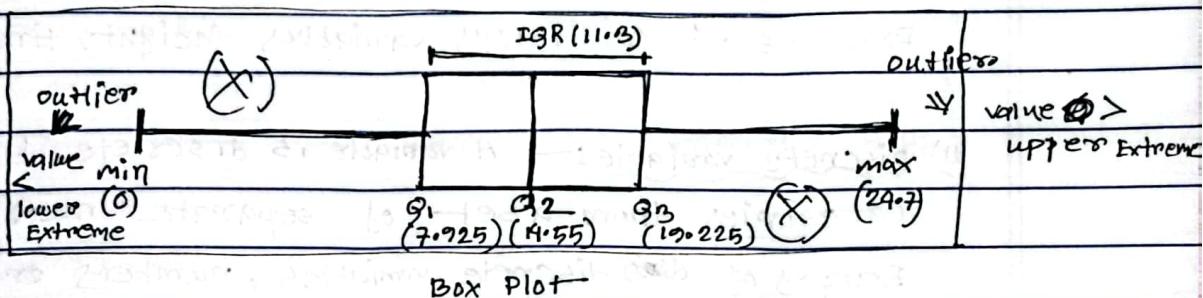
$$= 19.225 + 11.3 \times 1.5$$

$$= 36.175$$

Lower Extreme:  $Q_1 - IQR \times 1.5$

$$= 7.925 - 11.3 \times 1.5$$

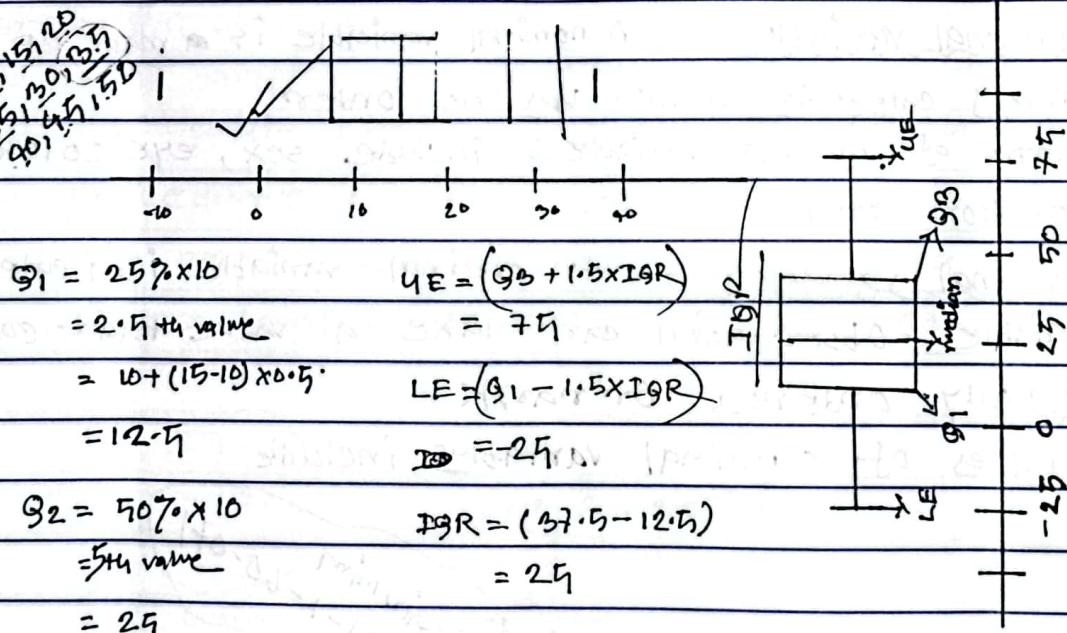
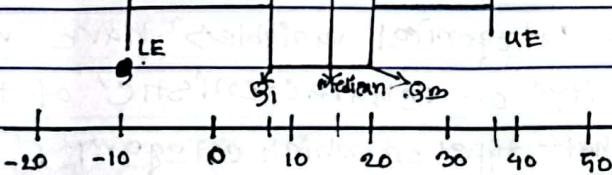
$$= -9.025$$



IQR

Hussain Mim

2021-1-60-071



$$\begin{aligned}
 Q_3 &= 75\% \times 10 \\
 &= 7.5 + \text{value} \\
 &= 37.5 + (40 - 35) \times 0.5 \\
 &= 37.5
 \end{aligned}$$

variables: A variable is any characteristic, number or quantity that can be measured. A variable may also be called a data item.

Age, sex, eye colour, country of birth etc are examples of variables.

Numerical variables:— Numerical variables have values that describe a measurable quantity as a number, like 'how many' or 'How much'.

i) continuous variable:— A continuous variable is a numeric variable. Observations can take any value between a set of certain set of real numbers.

Examples of continuous variables height, time, age etc.

ii) discrete variable:— A variable is discrete if it possible categories from a set of separate numbers.

Examples of discrete variables, number of children in a family, Number of registered cars, Number of boy in our school university.

categorical variables:— Categorical variables have values that describe a 'quality' or 'characteristic' of the data unit like 'what type' or 'which category'.

i) Nominal variables:— A nominal variable is made up of various categories which has no order.

Example of nominal variables include sex, eye colour, religion etc.

ii) ordinal variable:— An ordinal variable is a categorical variable. Observation can take a value that can be logically ordered or rank.

Examples of ordinal variables include

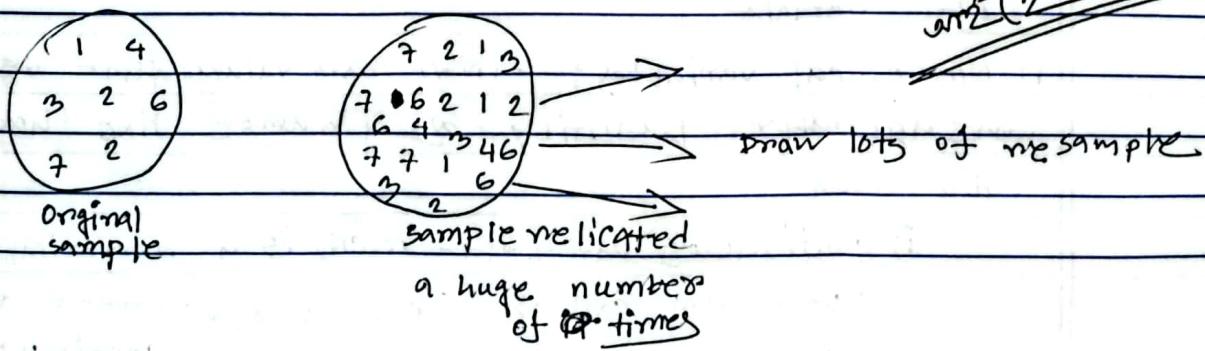
Highest  
021-160-071

There are mainly seven steps involved into it :-

- i) Decide on the objectives :- These objectives may usually require significant data collection and analysis.
- ii) what to measure and how to measure :- measurement generally refers to the assigning of the numbers to indicate different values of variables
- iii) Data collection :- once you know what type of data you need for your statistical study then you can determine whether your data can be gathered.
- iv) Data cleaning :- This is another crucial step in data analysis pipeline is to improved data quality for your existing data.
- v) summarizing :- Exploratory data analysis helps to understand the data better
- vi) Data modeling :-
- vii) optimize and Repeat :- The data analysis is a repeatable process and sometime leads to continuous improvements , both to the business .

The Bootstrap:- one easy and effective way to estimate the sampling distribution of a statistic of a model parameters , is to draw additional samples , with replacement from the sample itself and recalculate the statistic or model each resample

This procedure is called bootstrap and it dose not necessarily involve any assumption about the data or sample statistic ~~being~~ being normally distributed.



The algorithm for a bootstrap sampling of the mean is as follows for a sample of size  $n$ :

1) Draw a sample value, record, replace it

2) Repeat  $n$  times

3) Record the mean of the  $n$  resampled values

4) Repeat steps (1-3)  $R$  times

5) Use the  $R$  results to:

i) calculate their standard deviation

ii) produce a histogram or boxplot

iii) Find a confidence interval.

Confidence Interval: A  $\text{CI}$  is statistics is a range of values that used to estimate an unknown population parameter, such as the mean or proportion. It provides a range of plausible values for the parameter and is calculated from sample data.

For example 95%  $\text{CI}$  for mean weight of a population might be from 60kg to 70kg, indicating that we're 95% confident that the true population mean weight falls within this range.

The width of the interval is influenced by the sample size and the desired level of confidence.

variable: A variable is any characteristic, number, quantity that can be measured. A variable also be called a data item. Age, sex, eye colour, country of birth etc are examples of variable.

i) Numerical variables:  $(\text{NV})$  have values that describe a measurable quantity as a number like 'how many' or 'How much'

ii) continuous variable:  $(\text{CN})$  is a numerical variable.

Observation can take any value between a certain set of real numbers. Examples of continuous variable include weight, time, age.

7.3

 $\frac{2}{4} \times \frac{1}{2}$ 

Date:

Page:

Standard Error — dataset : 85, 90, 88, 92, 87

$$\frac{2}{10} \times \frac{2}{10} \times \frac{4}{10} \times \frac{6}{10}$$

$$\text{Mean } (\bar{x}) = \frac{442}{5}$$

$$= 88.4$$

$$s^2 = \frac{(85 - 88.4)^2 + (90 - 88.4)^2 + (88 - 88.4)^2 + (92 - 88.4)^2 + (87 - 88.4)^2}{5-1}$$

$$= 7.3$$

$$s = \sqrt{7.3}$$

$$= 2.673$$

$$SE = \frac{s}{\sqrt{n}}$$

$$= \frac{2.673}{\sqrt{5}}$$

$$= 1.220 \quad (\text{Ans})$$

ordinal :— observation can take a value that can be logically ordered and

rank. Example ordinal variable include

Health & quality (good, adequate, bad)

Example :— Random Sampling :— dataset 10, 8, 13, 7, 12, 5, 6, 9, 4  $\Rightarrow \frac{2}{4}$

sample size = 4

sample 1 (with replacement) 2, 2, 4, 6 =  $\frac{2+2+4+6}{4}$  with out

sample 2 (with replacement) 9, 9, 10, 8 =  $\frac{9+9+10+8}{4}$

= 9 (Ans)

Box 1  
Find the probability of the picking two black ball from Box 1 when you draw the samples without replacement?

$$P(\text{1st ball is black}) \times P(\text{2nd ball is black})$$

$$\Rightarrow \frac{3}{7} \times \frac{2}{6}$$

$$= \frac{6}{42} \quad (\text{Ans})$$

You draw the samples with replacement

$$\frac{3}{7} \times \frac{3}{7}$$

$$= \frac{9}{49} \quad (\text{Ans})$$

(Ans)

$H_0$ : There is no link between gender and political party

$H_1$ : There is link between gender and political party

Now,

The Expected value for male republican =  $\frac{290 \times 200}{420}$   
 $= 114$

The Expected value for female republican =  $\frac{790 \times 220}{420}$   
 $= 125.72$

The Expected value for male democrat =  $\frac{200 \times 130}{420}$   
 $= 61.91$

The Expected value for female democrat =  $\frac{220 \times 130}{420}$   
 $= 68.01$

The Expected value for male Independent =  $\frac{200 \times 50}{420}$   
 $= 23.81$

The Expected value for female Independent =  $\frac{220 \times 50}{420}$   
 $= 26.19$

\*  $x_e^{\vee} = \frac{(O_i - E_i)}{E_i}$

$x^{\vee}$  is the sum of the all value :  $(1.72 + 1.63 + 1.05 + 8.59 + 1.61 + 1.78)$

we have

$= 8.74$

$\{(r-1) \times (c-1)\}^2$

$\Rightarrow \{(2-1) \times (3-1)\}^2$

$= (1 \times 2)$

$= 2$

Hosking min  
20.21 - 60.02

As we see that for alpha level of 0.05 and two degrees of freedom the critical statistic is 5.991 which is less than our obtain statistic of 8.74  
 obtain statistic > critical statistic

$H_0$  can be rejected.

(Ans)

$  \text{if } O_i - E_i > \alpha \text{ value}$	$  \text{if } O_i - E_i < \alpha$
$O_i > E_i$	Crit. $O_i$
$H_0$ reject	$\therefore H_0$ do not rejected

A confidence interval is a range of values that are used to estimate an unknown population parameter, such as the mean or proportion.

It provides a range of plausible values for the parameters and is calculated from sample data.

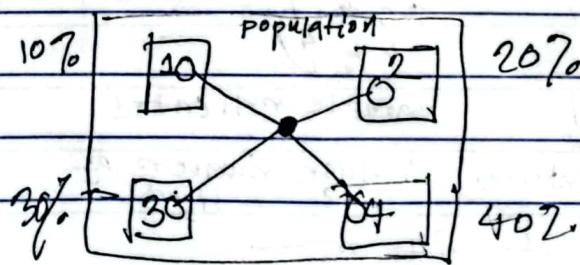
95% ⇒ It means that if we repeat an experiment or survey over and over again, 95% of the time our results will match the results we got from a pop. population.

self-selection bias in ~~statistic~~ occurs when individuals or entities self-select into a sample being studied, leading to a non-random or biased representation of the population.

An Example self-selection bias in political polling. Consider a situation where a poll about a certain candidate on online platform people who actively seek out and participate in such poll might have stronger opinions or be more engaged in the topic compared to general population.

stratified sampling is a sampling technique in statistic where the population is divided into distinct subgroups, called strata, based on certain characteristics that are relevant to the study.

The strata are homogeneous within themselves but different from each other in some key aspects.



10, 15, 18, 22, (25), 28, 33, 43, 50

 $n=9$  It is odd number.

$$\begin{aligned}\text{median} &= \frac{n+1}{2} \text{ th value} \\ &= \frac{10}{2} \text{ th value} \\ &= 5 \text{ th value} \\ &= 25 \quad (\text{Ans})\end{aligned}$$

10, 15, 18, 22, (25), (28), 30, 33, 50, 55

 $n=10$  It is even numbers

$$\begin{aligned}\text{median} &= \frac{(n)}{2} \text{ th value} + \frac{(n+1)}{2} \text{ th value} \\ &= \frac{5 \text{ th value} + 6 \text{ th value}}{2} \\ &= \frac{25 + 28}{2} \\ &= 26.5 \quad (\text{Ans})\end{aligned}$$

(5), 6, (5), 7, (5), 8, 9, (5)

mode: 5 [Here, 5 is most frequent item]

(5), (6), (5), (6), (5), 8, (6), (5), (6)

mode: 5 and 6 [Here, 5 and 6 both are most frequent item]

(60), (70), (80), (75), (65), (70), (100), (70), (65), (67)

sorted data: 60, 65, 65, 65, 70, 70, 70, 75, 80, 100

$$\begin{aligned}\text{mean} &= \frac{700}{10} \\ &= 70\end{aligned}$$

(60), (70), (80), (75), (65), (70), (100), (70), (65), (67)

sorted data: 10, 60, 65, 65, 65, 70, 70, 70, 75, 80

$$\begin{aligned}\text{mean} &= \frac{730}{10} \\ &= 73\end{aligned}$$

$$\begin{aligned}\text{trimmed mean} &= \frac{560}{10-(2 \times 1)} \\ &= \frac{560}{8} \quad [P=1] \\ &= 70 \quad (\text{Ans})\end{aligned}$$

$$\begin{aligned}\text{trimmed mean} &= \frac{540}{10-(2 \times 1)} \\ &= 67.5 \quad (\text{Ans})\end{aligned}$$

$$x_i = 40, 45, 50, 75, 10$$

$$w_i = 1, 2, 3, 4, 5$$

$$\begin{aligned}\text{weighted mean} &= \frac{\sum x_i w_i}{\sum w_i} \\ &= \frac{(40 \times 1) + (45 \times 2) + (50 \times 3) + (75 \times 4) + (10 \times 5)}{1+2+3+4+5} \\ &= 48 \quad (\text{Ans})\end{aligned}$$

Data → 600, 470, 170, 430, 300

$$\text{mean} = \frac{600 + 470 + 170 + 430 + 300}{5}$$

$$= 394$$

$$\text{variance } (S^2) = \frac{(600-394)^2 + (470-394)^2 + (170-394)^2 + (430-394)^2 + (300-394)^2}{5}$$

$$= 21704$$

$$\text{standard deviation } (S) = \sqrt{21704}$$

$$= 147.32 \quad (\text{Ans})$$

sorted data:— 64, 65, 68, 69, 70, 71, 72, 75, 79, 80, 81, 83, 84

$$25^{\text{th}} \text{ percentile } (Q_1) = 25\% \times 14 \\ = 3.5 \text{ th value} \\ = 68 + (69 - 68) \times 0.5 \\ = 68.5$$

~~MID 0  
JUN 2 15~~

$$50^{\text{th}} \text{ percentile (median)} = 50\% \times 14 \\ = 7^{\text{th}} \text{ value} \\ = 72$$

$$75^{\text{th}} \text{ percentile } (Q_3) = 75\% \times 14 \\ = 10.5 \text{ th value} \\ = 75 + (80 - 75) \times 0.5 \\ = 77.5$$

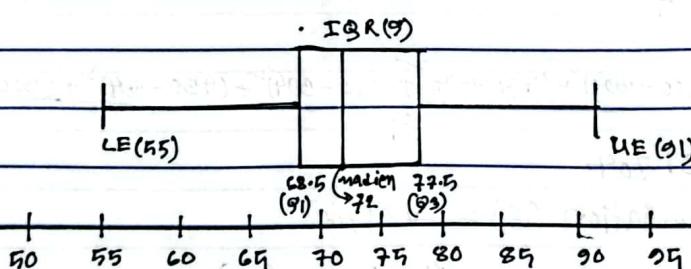
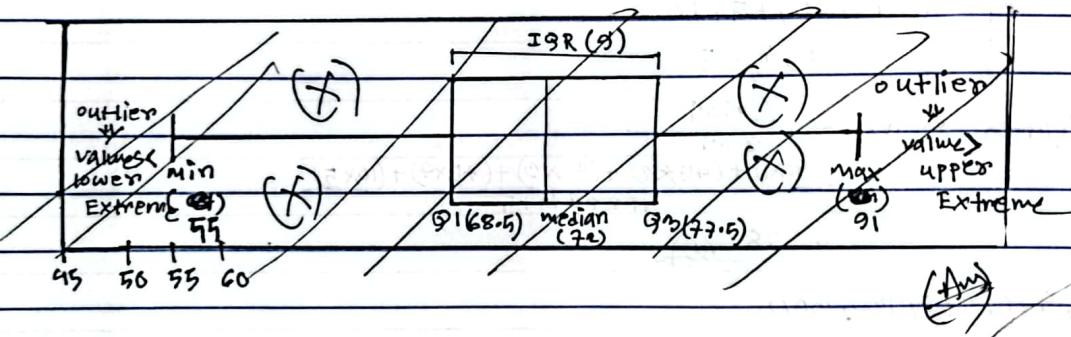
$$\text{IQR} = (77.5 - 68.5) \\ = 9$$

max = 84

min = 64

*(Q3 - Q1)*

upper extreme =  $\{Q_3 + (IQR \times 1.5)\}$  Lower Extreme:  $Q_1 - IQR \times 1.5$   
 $= \{77.5 + (9 \times 1.5)\}$   $= \{68.5 - (9 \times 1.5)\}$   
 $= 86$  55



we can see

no outliers

AZ

Q:

	♂	♀	N <sub>3</sub>
M	40	70	120
F	30	50	80
T	80	120	200

H<sub>0</sub>: There is no significant association between gender and smoking

H<sub>1</sub>: There is a significant association between gender and smoking habits

$$\text{Expected value for male smoker} = \frac{80 \times 120}{200}$$

$$= 48$$

$$\text{Expected value for female smoker} = \frac{80 \times 80}{200}$$

$$= 32$$

$$\text{Expected value for male non-smoker} = \frac{120 \times 120}{200}$$

$$= 72$$

$$\text{Expected value for female non-smoker} = \frac{120 \times 80}{200}$$

$$= 48$$

NOW,

$$x^2_e = \frac{(O_i - E_i)^2}{E_i}$$

$$x^2_e \text{ all values are sum} = (0.084 + 0.125 + 0.025 + 0.046 + 0.080 \cdot 0.084)$$

$$= 0.349$$

$$\text{Degrees of freedom} = (2-1) \times (2-1)$$

$$= 1$$

As we can see for alpha level of 0.05 and 1 degree of freedom.

The critical statistic is 3.841.

And obtain statistic is 0.349

so, critical statistic > obtain statistic

H<sub>0</sub> do not rejected

(Ans)



4.30

$$(2-1) \times (1-1)$$

$$3-1$$

v	i		
c			
b			

	V	C	S	TD
100	120	80	300	
100	100	80	300	

Date :

Page :

H<sub>0</sub>: There is no significant difference in the ice-cream flavor among kids.

H<sub>1</sub>: There is a significant difference in the ice-cream flavor among kids.  
Now,

$$\text{Expected value for vanilla flavor} = \frac{100 \times 300}{300}$$

$$= 100$$

$$\text{Expected value for chocolate flavor} = \frac{120 \times 300}{300}$$

$$= 120$$

$$\text{Expected value for strawberry flavor} = \frac{80 \times 300}{300}$$

$$= 80$$

Now,

$$x_c^2 = \frac{(o_i - E_i)^2}{E_i}$$

$$x_c^2 \text{ all values are sum } (0 + 0 + 0) \\ = 0$$

$$\text{Degrees of freedom} = (3-1)$$

$$= 2$$

As, we can see that alpha level of 0.05 and 2 degrees of freedom.

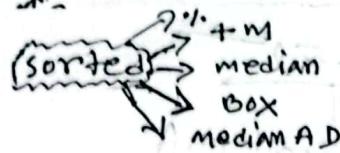
The critici critical statistic is 5.991

The obtain statistic is 0

so, critical statistic > obtain statistic

H<sub>0</sub> do not rejected.

(Ans)



# (60, 70, 80, 75, 65, 70, 80, 70, 65, 65)

Sorted order: 60, 65, 65, 65, 70, 70, 70, 75, 80, 80

- mean = 70 (Ans)

- median =  $\frac{(\frac{n}{2})\text{th value} + (\frac{n+1}{2})\text{th value}}{2}$

$$= \frac{70+70}{2}$$

$$= \frac{140}{2} = 70 \text{ (Ans)}$$

- trimmed mean =  $\frac{560}{10-(2 \times 1)} \quad [\text{For } P=1]$   
= 70 (Ans)

# 60, 70, 80, 75, 65, 70, 10, 70, 65, 65

- trimmed mean =  $\frac{740}{10-(2 \times 1)}$

$$= 67.5 \text{ (Ans)}$$

→ # (5), (6), (5), (6), (5), 8, (6), (5), (6)

- mode: 5 and 6 [Home, 5 and 6 both

same most frequent item]

### Mean absolute deviation —

# 0, 25, 50, 75, 100 mean =  $\frac{250}{5}$

MAD =  $\frac{\sum |x - \bar{x}|}{n}$  = 50

$$= \frac{|10-50| + |25-50| + |50-50| + |75-50| + |100-50|}{5}$$

$$= \frac{50+25+0+25+50}{5}$$

$$= \frac{140}{5}$$

$$= 28 \text{ (Ans)}$$

dataset: 50, 75, 100, 0, 25

sorted: — 0, 25, (50) 75, 100

$$\frac{n+1}{2} \text{ th value}$$

$$= 3 \text{ th value}$$

# 30, 40, 50, 60, 70 mean =  $\frac{250}{5}$

MAD =  $\frac{|30-50| + |40-50| + |50-50| + |60-50| + |70-50|}{5} = \frac{20}{5}$

$$= \frac{20+10+0+10+20}{5}$$

$$= \frac{60}{5}$$

$$= 12 \text{ (Ans)}$$

median (10-50, 125-50, 150-50)

175-50, 1100-50)

median (40, 25, 0, 25, 50)

median (0, 25, (25), 40, 50)

median = 25 (Ans)

# 20, 40, 50, 60, 70, 30, 50, 40, 60

n = 9 odd

median =  $\frac{5+1}{2}$

$$= 3 \text{ th value}$$

$$= 50$$

median (10, 10, 10, 20, 20)

median (0, 10, 10, 20, 20)

median = 10 (Ans)

$$\checkmark \quad \sigma^2 = \frac{(600-394)(x_i - \bar{x})}{n}$$

Data: 600, 470, 170, 430, 300

$$mean: \rightarrow \frac{600 + 470 + 170 + 430 + 300}{5}$$

$$= \frac{1970}{5}$$

$$\checkmark \quad \sigma^2 = \frac{(600-394) + (470-394) + (170-394) + (430-394) + (300-394)}{5}$$

$$= 21704$$

$$\sigma = \sqrt{21704}$$

$$= 147.32 \quad (\text{Ans})$$

# 30, 40, 50, 60, 70

$$mean = \frac{290}{5}$$

$$= 58$$

$$\checkmark \quad s^2 = \frac{(30-58)^2 + (40-58)^2 + (50-58)^2 + (60-58)^2 + (70-58)^2}{5-1}$$

$$= 290$$

$$s = \sqrt{290}$$

$$= 17.82 \quad (\text{Ans})$$

Highest min

# 0, 25, 50, 75, 100

$$mean = 50$$

$$\checkmark \quad s^2 = \frac{(0-50)^2 + (25-50)^2 + (50-50)^2 + (75-50)^2 + (100-50)^2}{5-1}$$

$$= 1562.5$$

$$s = \sqrt{1562.5}$$

$$= 39.5 \quad (\text{Ans})$$

sorted order: 8, 9, 10, 10, 10, 11, 11, 11, 12, 13

$$mean = \frac{105}{10}$$

$$= 10.5 \quad (\text{Ans})$$

$$n = 10 \text{ (even no)}$$

$$= \frac{\sum_{i=1}^{n/2} \text{value}_i + (\frac{n}{2}+1) + \text{value}_{n/2}}{2}$$

$$= \frac{5+\text{value}_5 + 6+\text{value}_6}{2}$$

$$= \frac{10+11}{2}$$

$$= 10.5 \quad (\text{Ans})$$

$$MAD: \frac{10-10.5 + 11-10.5}{10}$$

$$MAD: \frac{18-10.5 + 19-10.5}{10}$$

$$= 1.1$$

$$\text{median } (18-10.5, 19-10.5, 110-10.5)$$

$$(10-10.5) + (11-10.5), 111-10.5$$

$$11-10.5, 112-10.5, 113-10.5$$

$$\text{median } (2.5, 1.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5)$$

mode = 10 and 11 [here, 10 and 11 both are most frequent item] =  $\frac{2.5 + 1.5}{2} = 2$  AM

$$\text{variance } (s^2) = \frac{(8-10.5)^2 + (9-10.5)^2 + (10-10.5)^2 + (11-10.5)^2 + (12-10.5)^2 + (13-10.5)^2}{10-1}$$

$$= 18.5$$

$$= 2.056 \quad (\text{Ans})$$

$$s = \sqrt{2.056} = 1.43$$

$H_0$ : There is no link between gender and political party preference

$H_1$ : There is link between gender and political party preference

$$\text{Expected value for male republican} = \frac{240 \times 200}{420}$$

$$= 114.28$$

$$\text{Expected value for female republican} = \frac{120 \times 220}{420}$$

$$= 125.72$$

$$\text{Expected value for male democrat} = \frac{130 \times 200}{420}$$

$$= 61.91$$

$$\text{Expected value for female democrat} = \frac{130 \times 220}{420}$$

$$= 68.09$$

~~MIDO / 9/2022~~

$$\text{Expected value for male independent} = \frac{50 \times 200}{420}$$

$$= 23.81$$

$$\text{Expected value for female independent} = \frac{50 \times 220}{420}$$

$$= 26.19$$

Now,

$$\chi^2 = \frac{(O_i - E_i)^2}{E_i}$$

$$\therefore \chi^2 \text{ are all sum of values} = (1.79 + 1.63 + 1.07 + 0.97$$

$$+ 1.61 + 1.46)$$

$$\text{Degree degrees of freedom } (2-1) \times (3-1) = 8 - 1 = 7$$

$$= \frac{(1 \times 2)}{2}$$

we can see that alpha level of 0.05 and 2 degrees of freedom

the critical statistic is 5.991  
we obtain statistic is 8.91

~~critical~~ < ob so,  $H_0$  rejected (Ans)

## Chi-square test

standard error— The SE is a single metric that sum up the variability the sampling distribution for a statistic

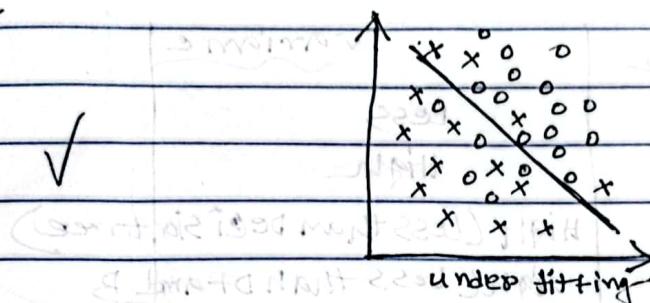
$$SE = \frac{s}{\sqrt{n}} \rightarrow (\text{sample})$$

~~MD~~ ~~SD~~

~~NP~~

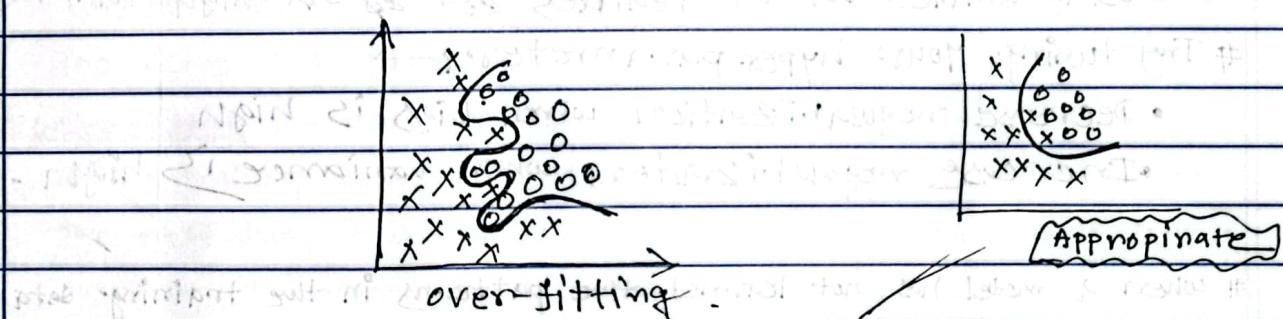
Under-fitting :- when a model has not learned learned the pattern in the training data well and is unable to generalization well on the new data , it is known as under-fitting. An under-fitting model has poor performance will result in underfittable. under-fitting occurs due to high bias and low variance . It is too simple model and too much regularization in the model .

Important



hossein mim  
2021-1-60-071

over-fitting :- over-fitting is a phenomenon that occurs when a machine learning model is constrained constraint to training set and not able to perform well on unseen data. over-fitting model is very complex and too little regularization over-fitting model occurs due to low bias and high variance.



Appropriate

Bias :- we can define bias the error between average model prediction and ground truth.

- A model which High-Bias would not match the data set closely
- A model Low-Bias will closely match the training data set

• If we have a model which is too complex it will overfit the data set . If we have a model which is too simple it will underfit the data set .

variance refers to the changes in the model when using different portions of the training data set. It is the variability in the model prediction - how much the machine learning function can adjust depending on the given data set.

<u>Algorithm</u>	<u>Bias</u>	<u>Variance</u>
• Linear regression	High	Less
• Decision tree	Low	High
• Bagging	Low	High (less than Decision tree)
• Random forest	Low	High less than DT and B

# Get more than data:

- when we are have high variance

# Try different features:

- Adding features helps fix high Bias
- using smaller set of features fix high variance

# Try tuning your hyper parameters:

- Decrease regularization when bias is high

- Increase regularization when variance is high

# when a model has not learned the patterns in the training data well and is unable to generalization well on the new data, it is known as under-fitting.

# when a model has not learned the pattern in the training data well and is unable to generalization well on the new data.

# over-fitting is a phenomenon that occurs when a machine learning model is ~~constraint~~ constraint to training set and not able to perform well on its unseen data.

# over-fitting is a phenomenon that occurs when a ML model is constraint to training set and not able to perform well on unseen data.

confusion matrix

		predicted class	
		Crocodile	Alligator
Actual class	Crocodile	440 (TP)	80 (FN)
	Alligator	60 (FP)	420 (TN)

• Accuracy =  $\frac{TP + TN}{TP + TN + FP + FN}$   
 $= \frac{440 + 420}{440 + 420 + 80 + 60}$   
 $= 0.86$   
 $= 86\%$

• Recall =  $\frac{TP}{TP + FN}$   
 $= \frac{440}{440 + 80}$   
 $= 0.85$   
 $= 85\%$

• Recall of Crocodile =  $\frac{TP}{TP + FN}$   
 $= \frac{440}{440 + 80}$   
 $= 0.85$   
 $= 85\%$

• Recall of Alligator =  $\frac{TN}{TN + FP}$   
 $= \frac{420}{420 + 60}$   
 $= 0.875$   
 $= 87.5\%$

• Precision for Crocodile =  $\frac{TP}{TP + FP}$   
 $= \frac{440}{440 + 60}$   
 $= 0.88$   
 $= 88\%$

• Precision for Alligator =  $\frac{TN}{TN + FN}$   
 $= \frac{420}{420 + 80}$   
 $= \frac{420}{500}$   
 $= 0.84$   
 $= 84\%$

\* The accuracy is at 86% meaning it correctly predicts classes about 78% of the time.

\* This suggests that the classifier performs better in identifying Crocodile instances than Alligator instances, especially in terms of precision and recall.

~~F1-score~~

$$\text{F1-score for Crocodile} = 2 \times \frac{0.85 \times 0.88}{0.85 + 0.88} \\ = 0.864 \\ = 86.4\%$$

~~$$\text{F1-score for Alligator} = 2 \times \frac{0.875 \times 0.84}{0.875 + 0.84} \\ = 0.857 \\ = 85.7\%$$~~

## Inverse Matrix:-

Time consumption  
AII

Date :

Page :

$$A = \begin{vmatrix} 1 & 1 & 1 \\ 1 & -1 & 2 \\ 3 & 5 & -7 \end{vmatrix} \quad \vec{r} = \begin{bmatrix} 8 \\ 6 \\ 14 \end{bmatrix}$$

$$Ax = B$$

$$x = A^{-1}B$$

$$|A| = \begin{vmatrix} 1 & 1 & 1 \\ 1 & -1 & 2 \\ 3 & 5 & -7 \end{vmatrix} \Rightarrow 16$$

For Alligator:-

$$\text{1) Accuracy} = \frac{TN + FP}{TN + FP + FN + TP} = \frac{420 + 80}{420 + 80 + 440 + 60} = 0.5$$

$$= 50\% \quad (\text{AM})$$

$$\text{2) Recall} = \frac{TN}{TN + FN} = \frac{420}{420 + 80} = 0.84$$

$$= 84\% \quad (\text{AM})$$

$$\text{3) Precision} = \frac{TP}{TP + FP} = \frac{420}{420 + 60} = 0.875$$

$$= 87.5\% \quad (\text{AM})$$

$$\text{4) F1-score} = 2 \times \frac{0.84 \times 0.875}{0.84 + 0.875} = 0.857$$

(This suggest that the classifier perform better in identifying "x" instances than "y" instances, especially in terms of precision and recall.)

	predicted	
	Crocodile	Alligator
Actual		
Crocodile (P)	440 (TP)	60 (FN)
Alligator (N)	80 (FP)	420 (TN)

# For Crocodile:-

$$\text{1) Accuracy} = \frac{TP + TN}{TP + TN + FN + FP} = \frac{440 + 420}{440 + 420 + 80 + 60} = 0.86 = 86\% \quad (\text{AM})$$

$$\text{2) Recall} = \frac{TP}{TP + FN} = \frac{440}{440 + 80} = 0.88 = 88\% \quad (\text{AM})$$

$$\text{3) Precision} = \frac{TP}{TP + FP} = \frac{440}{440 + 60} = 0.846 = 84.6\% \quad (\text{AM})$$

$$\text{4) F1-score} = 2 \times \frac{0.846 \times 0.88}{0.846 + 0.88} = 0.852 = 85.2\% \quad (\text{AM})$$

higher CS  
Recall ✓

- In multiplication is associative justify :-

$$x = \begin{bmatrix} A & B \\ C & D \end{bmatrix}, y = \begin{bmatrix} E & F \\ G & H \end{bmatrix}, z = \begin{bmatrix} I & J \\ K & L \end{bmatrix}$$

$$\{(x \times y) \times z\} = \{x \times (y \times z)\}$$

Now,

$$(x \times y) = \begin{bmatrix} AE + BG & AF + BH \\ CE + DG & CF + DH \end{bmatrix}$$



$$\{(x \times y) \times z\} = \begin{bmatrix} AE + BG & AF + BH \\ CE + DG & CF + DH \end{bmatrix} \times \begin{bmatrix} I & J \\ K & L \end{bmatrix}$$

$$= \begin{bmatrix} AEI + BGJ & AEJ + BGK \\ CEI + DGI & CEJ + DGL \end{bmatrix} + \begin{bmatrix} AFK + BHK & AFL + BHL \\ CFK + DHK & CFL + DHL \end{bmatrix}$$

$$(y \times z) = \begin{bmatrix} EI + FK & EJ + FL \\ GI + HK & GJ + HL \end{bmatrix}$$

$$\{x \times (y \times z)\} = \begin{bmatrix} A & B \\ C & D \end{bmatrix} \times \begin{bmatrix} EI + FK & EJ + FL \\ GI + HK & GJ + HL \end{bmatrix} \times \begin{bmatrix} I & J \\ K & L \end{bmatrix}$$

$$= \begin{bmatrix} AEI + AFK + BGJ + BHK & AEJ + AFL + BGK + BHL \\ CEI + CFK + DGI + DHK & CEJ + CFL + DGL + DHL \end{bmatrix}$$

$$\text{so, } \{x \times (y \times z)\} = \{(x \times y) \times z\} \quad (\text{Ans})$$

Matrix multiplication is associative.

~~Firat  
কর্মসূচি~~

~~(Ans)~~

\* In not commutative justify

$$x = \begin{bmatrix} a & b \\ c & d \end{bmatrix}, y = \begin{bmatrix} e & f \\ g & h \\ i & j \end{bmatrix}$$

$$(xy) \neq (yx)$$

$$(xy) = \begin{bmatrix} a & b & p \\ c & d & q \end{bmatrix} \quad x = \begin{bmatrix} e & f \\ g & h \\ i & j \end{bmatrix}$$

$$\begin{bmatrix} ae + bg + pi & af + bh + pj \\ ce + dg + qi & cf + dh + qj \end{bmatrix}$$

$$yx = \begin{bmatrix} e & f \\ g & h \\ i & j \end{bmatrix} \quad x = \begin{bmatrix} a & b & p \\ c & d & q \end{bmatrix}$$

$$\begin{bmatrix} ea + fc & eb + fd & ep + cq \\ ga + hc & gb + hd & gp + hr \\ ia + jc & ib + jd & ip + jq \end{bmatrix}$$

$$so, (xy) \neq yx$$

we can say that the matrix multiplication is not

commutative

~~Final  
Opn Cr~~

- Addition
- Inverse
- Subtraction

$$A = \begin{bmatrix} 1 & 2 & -3 \\ -6 & -8 & 10 \\ -1 & 9 & 4 \end{bmatrix} \quad B = \begin{bmatrix} 0 & -9 & 6 \\ 3 & 4 & -2 \\ 0 & 8 & -1 \end{bmatrix}$$

রূপ এবং স্থিতিশীল আর্ডার একই হচ্ছে ৩x3

$$(A+B) = \begin{bmatrix} 1 & -7 & 3 \\ -3 & -4 & 8 \\ -1 & 17 & 3 \end{bmatrix} \quad (A-B) = \begin{bmatrix} 1 & -11 & -9 \\ -9 & -12 & 12 \\ -1 & 1 & 5 \end{bmatrix}$$

(Ans) (Ans)

~~X1~~

$$(A+B) = \begin{bmatrix} 3 & 1 & -1 \\ 2 & 3 & 4 \\ -4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 1 & 3 & 1 \\ 4 & 2 & 0 \\ 1 & 6 & 9 \end{bmatrix} \Rightarrow \begin{bmatrix} 3 & 4 & 0 \\ 3 & 7 & 6 \\ -3 & 11 & 15 \end{bmatrix}$$

~~X2~~

~~(Ans)~~

$$(A-B) = \begin{bmatrix} 2 & -2 & -2 \\ 1 & -1 & 2 \\ -5 & -1 & -3 \end{bmatrix}$$

(Ans)

$$6A + 3B = \begin{bmatrix} 18 & 6 & -6 \\ 12 & 18 & 24 \\ -24 & 30 & 36 \end{bmatrix} + \begin{bmatrix} 3 & 9 & 3 \\ 3 & 12 & 6 \\ 3 & 18 & 27 \end{bmatrix} = \begin{bmatrix} 21 & 15 & -3 \\ 15 & 30 & 30 \\ -21 & 48 & 33 \end{bmatrix}$$

(Ans)

$$\begin{array}{l} \xrightarrow{1 \rightarrow} \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \quad \begin{bmatrix} 3 & 2 \\ 4 & 4 \\ -1 & 6 \end{bmatrix} \\ \xrightarrow{2 \rightarrow} \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \quad \begin{bmatrix} 3 & 2 \\ 4 & 4 \\ -1 & 6 \end{bmatrix} \\ \xrightarrow{3 \rightarrow} \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \quad \begin{bmatrix} 3 & 2 \\ 4 & 4 \\ -1 & 6 \end{bmatrix} \end{array} \quad * \text{ এটি সর্বোচ্চ সম্ভব } \rightarrow$$

$$\begin{bmatrix} 13+8-3 & 3+8+18 \\ 12+20+(-6) & 8+20+36 \end{bmatrix} = \begin{bmatrix} 18 & 28 \\ 26 & 64 \end{bmatrix}$$

(Ans)

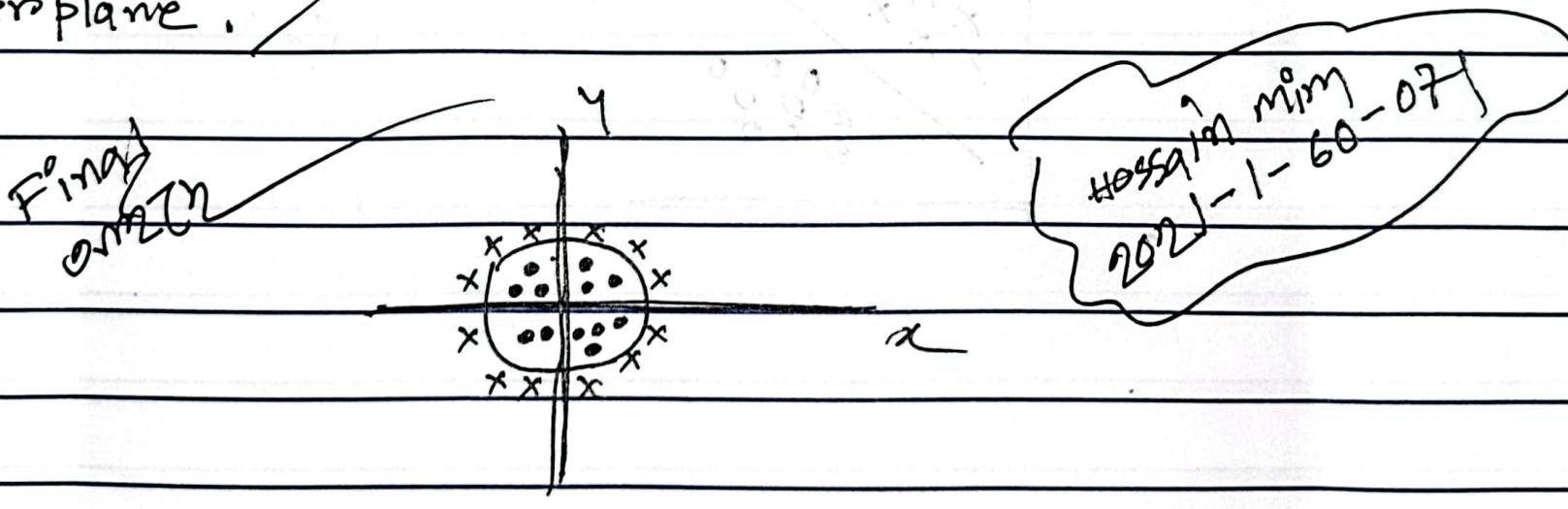
$$A = \begin{bmatrix} 1 & 3 & 8 \\ 2 & 4 & 6 \end{bmatrix}$$

$$A^T = \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 8 & 6 \end{bmatrix}$$

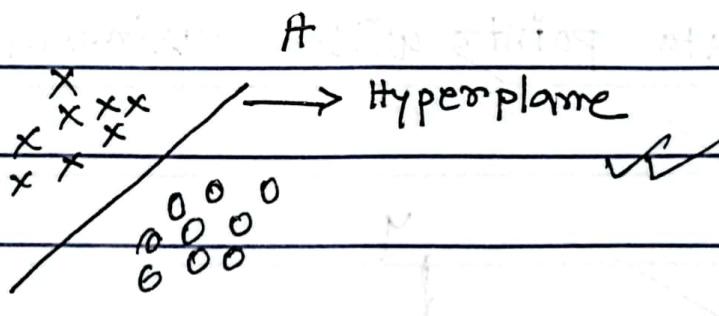
Ans.

The SVM kernel is a function that takes low dimensional input space and transform it into a high dimensional space that means it converts non-separable problem to a separable problem. It is mostly useful in non-linear data separation problem. It simply performs some extremely complex data transformation, then it finds out the process the data based on true labels or output.

It works by finding the best possible boundaries than can separate two classes data points with maximum margin, also it known as hyperplane.



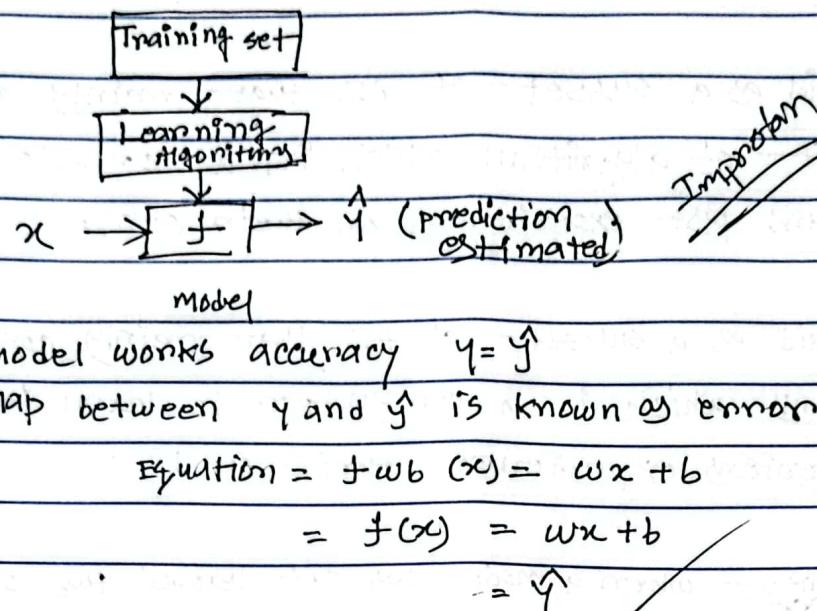
The support vector machine is defined as a machine learning algorithm that uses supervised learning model to solved complex classification, regression and outlier detection problem by performing the optimal data transformation that determine boundaries between the data points based on the labels or output



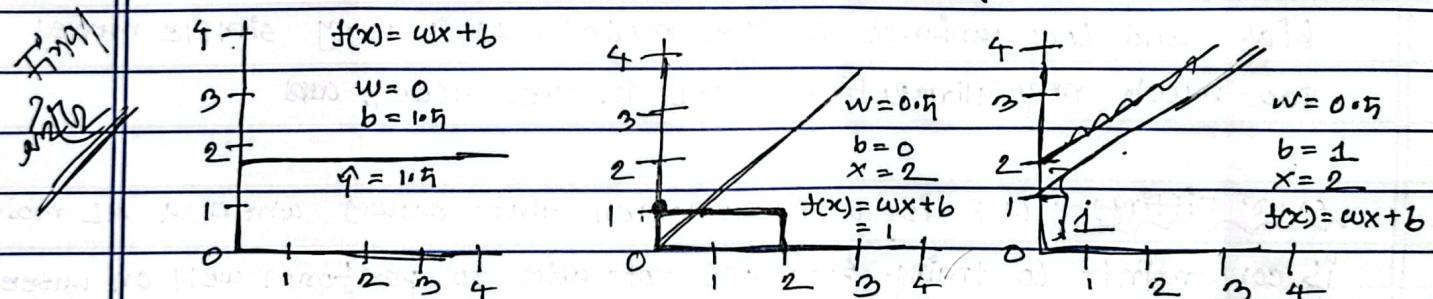
# ML is said as a subset of AI that mainly contrasts with development of algorithm which in a computer to learn from the data and past experiences on their own.

# ML is said as a subset of AI that mainly contrasts with development of algorithms which in a computer to learn from the data and past experiences on their own.

Training set — Training set is the major portion of original dataset which is used for ~~testing~~ training the ML model!



Cost function — Cost function ~~cost~~ quantifies the error between predicted and target values and present that error in the form of a single real number

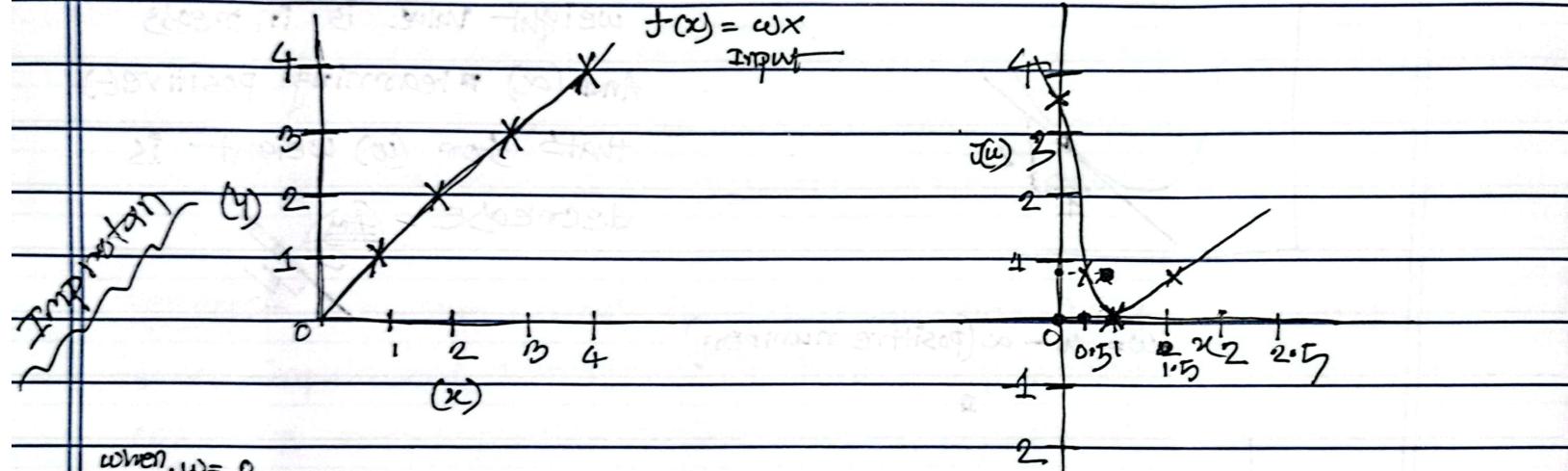


Cost function —

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2$$

$m$  = number of training example.

## Regression problem



when,  $w=0$

$$f(x) = wx$$

$$= (0 \times 1) \Rightarrow 0 \quad = (0 \times 2) \Rightarrow 0 \quad = (0 \times 3) \Rightarrow 0 \quad = (0 \times 4) \Rightarrow 0$$

$$\text{Now, } J(w) = \frac{1}{2 \times 4} \{ (0-1)^2 + (0-2)^2 + (0-3)^2 + (0-4)^2 \}$$

$$= 3.75$$

when  $w=1.5$

$$= 1.5 \times 1 \Rightarrow 1.5 \times 2 \\ = 1.5 \quad \quad \quad = 1.5$$

$$= 1.5 \times 3 \Rightarrow 1.5 \times 4 \\ = 4.5 \quad \quad \quad = 6$$

when  $w=0.5$

$$\Rightarrow (0.5 \times 1) = (0.5 \times 2) \Rightarrow (0.5 \times 3) \Rightarrow (0.5 \times 4) \\ = 0.5 \quad \Rightarrow 1.0 \quad = 1.5 \quad = 2.0$$

$$\text{Now, } J(w) = \frac{1}{2 \times 4} \{ (0.5-1)^2 + (1.0-2)^2 + (1.5-3)^2 + (2-4)^2 \}$$

$$= 0.94$$

$$J(w) = \frac{1}{2 \times 4} \{ (1.5-1)^2 + (3-2)^2 + (4.5-3)^2 + (6-4)^2 \}$$

$$\approx 0.94$$

when,  $w=1$

$$\Rightarrow (1 \times 1) \Rightarrow (1 \times 2) \Rightarrow (1 \times 3) \Rightarrow (1 \times 4) \\ = 1 \quad \Rightarrow 2 \quad = 3 \quad = 4$$

$$\Rightarrow J(w) = \frac{1}{2 \times 4} \{ (1-1)^2 + (2-2)^2 + (3-3)^2 + (4-4)^2 \}$$

$$= 0$$

so, we can see

lowest  $w$  value is

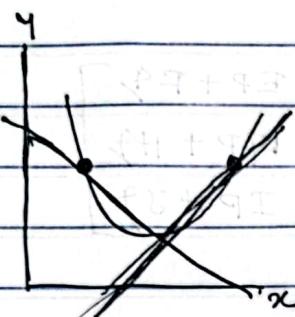
$w=1$   
Best regression line  
 $w=1$

GID Also Problem:-

we know that

$$w = w - \alpha \frac{\partial}{\partial w} J(w, b)$$

$\downarrow$  weight  $\downarrow$  learning rate



Assum,

$$w = 2$$

$$\alpha = -1$$

$$w = \{2 - (-1)\}^2$$

$$= 2 + 1$$

$$= 3$$

Assum,

$$w = 2$$

$$\alpha = 1$$

$$w = (2 - 1)^2$$

$$= 1$$

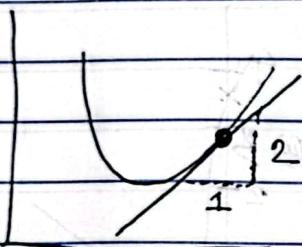
so, we can say that the learning rate

( $\alpha$ ) is negative that is for

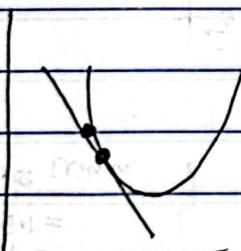
weight value is increasing

And ( $\alpha$ ) is learning positive (+)

that's for ( $w$ ) weight is  
decrease



$$w = w - \alpha (\text{positive number})$$



$$w = w - \alpha (\text{negative number.})$$

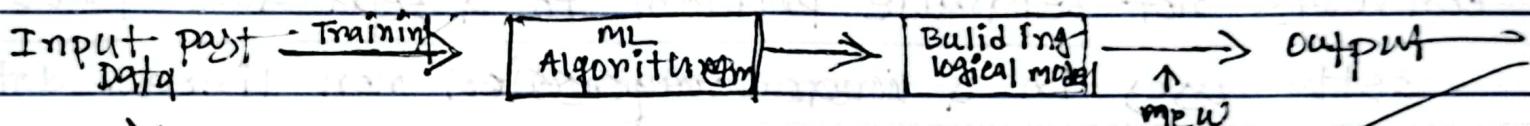
Hessian min  
1-60-opt

- If ( $\alpha$ ) learning rate is too small " the GID is may be slow

- If ( $\alpha$ ) learning rate is too large

The GID may be overshoot, never minimum  
and fail to converge

# machine learning is said as a subset of AI that mainly concern with development which in a computer to learn from the data and past ~~experiences~~ experiences their own.



Supervised learning:- supervised learning is a type of machine learning method in which we provided sample labeled data to the machine learning system in order to train it and on that basis, it predicts the output.

The goal of supervised learning is map input data with the output data. The supervised learning is based on supervision.

supervised learning can be grouped further in categories

- 1) classification
- 2) Regression

The example of SL is spam filtering, image recognition system.

unsupervised learning is a learning method in which a machine learning without any supervision. The training is provided to the machine with the ~~set~~ of data that has not been labeled and algorithm needs to act on that data without any supervision. The goal of usl is to restructure the input data into new features with similar patterns. In ~~an~~ unsupervised learning we do not have a predetermined result.

- i) clustering
- ii) Association.

For example finding out which customers made similar product purchases.