

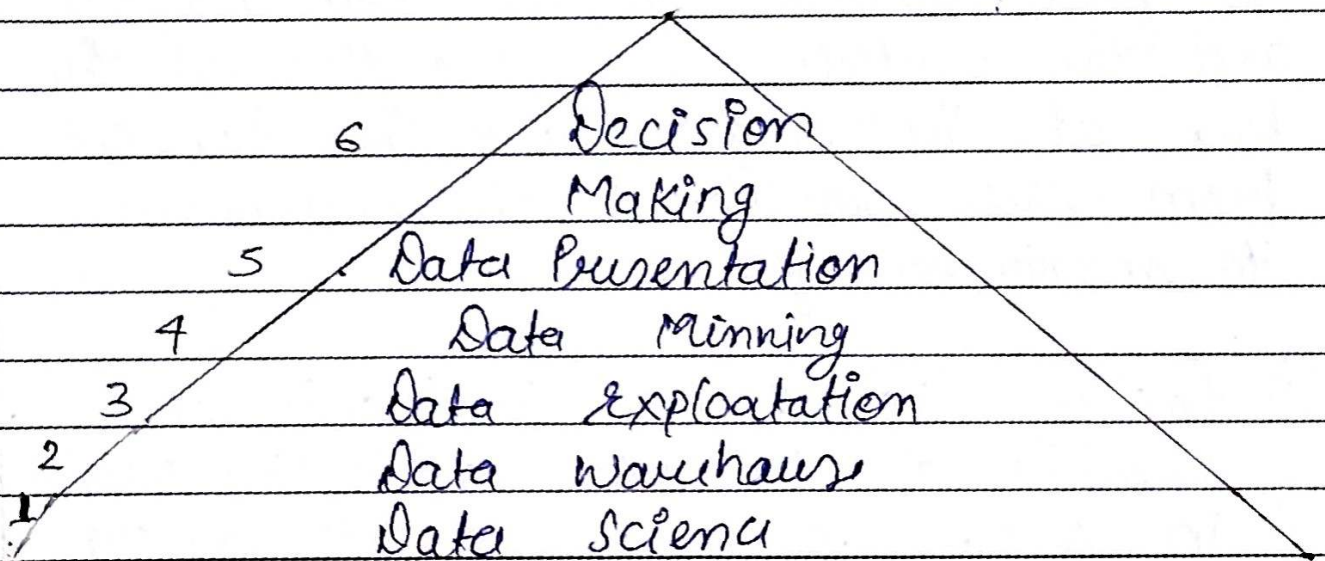
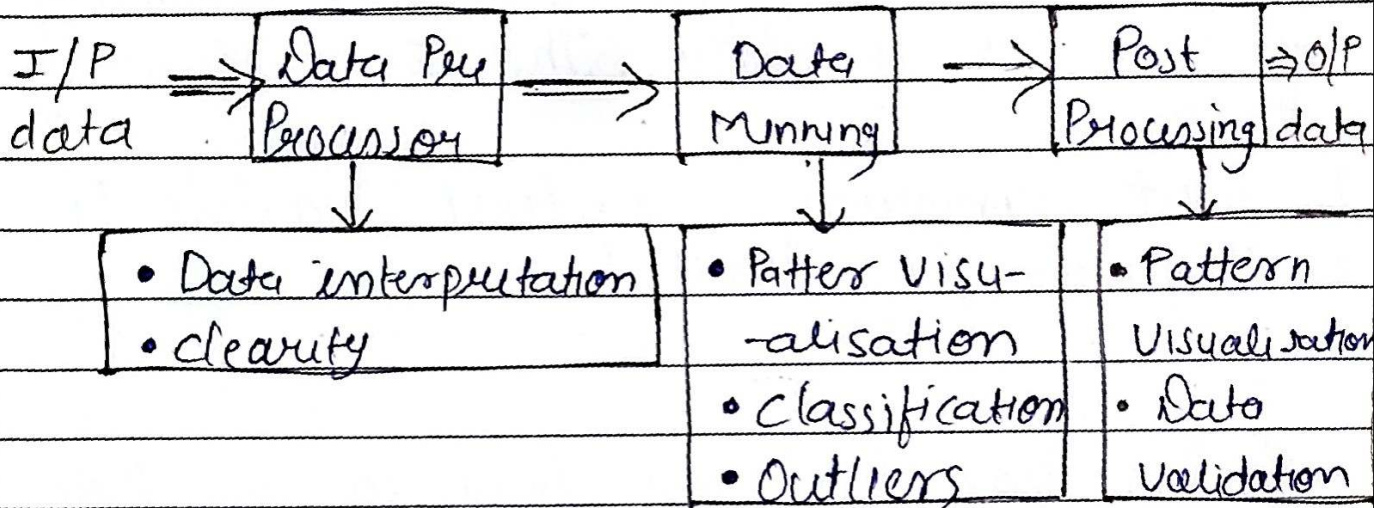
# # Data Science :-

Abstracting information from a source is called data mining.

## ⇒ Compilation Process :-

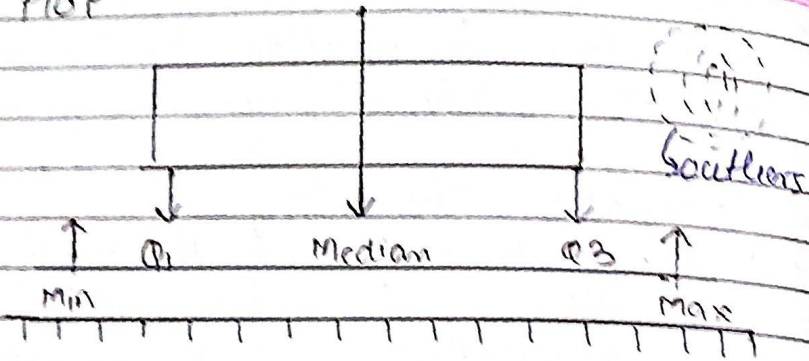
- |                      |                    |
|----------------------|--------------------|
| 1. Lexical analyser  | 2. Syntax analyser |
| 3. Semantic          | 4. Intermediate    |
| 5. Code optimisation | 6. Code generation |

21/08/23



6 D'S

### # Box Plot



A box plot display the S no. Summary of a set of data. The S no. Summary is :-

1. first minimum
2. first Quartile  $Q_1$
3. Minimum first Quartial
- 4.

In a box plot we draw a box from the first quartile to the 3rd quartile and the vertical line goes through the box at the median and the box go from each quartile to the minimum to maximum for ex :-

(Q1) finding the S no. Summaries. A sample of 10 boxes with the weight in gram i.e 25g, 28g, 29g, 30g, 34g, 35g, 37g, 38g, 29g, 35g make a box plot of the data

Step 1 :- Order the data from smallest to largest

25, 28, 29, 29, 30, 34, 35, 35, 37, 38

Step 2 :- find the median

$$\frac{30 + 34}{2} = 32$$

Step 3 :- finding 1st Quartile :-

The first quartile is the median of the data points to the left of the median i.e  $Q_1$

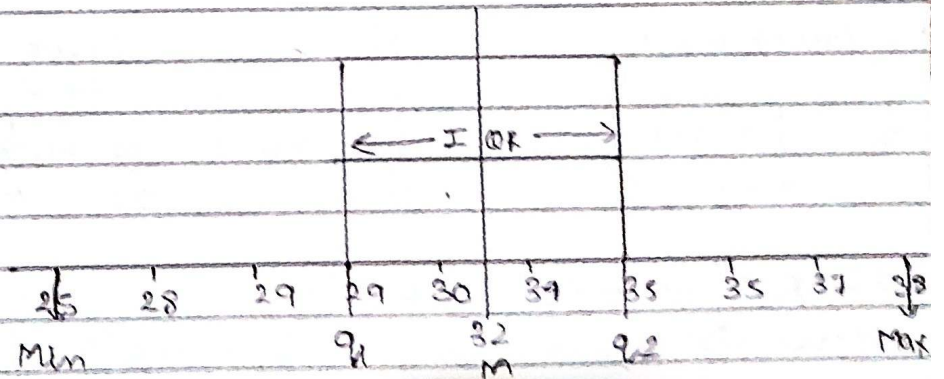
$$\Rightarrow Q_1 = 29$$

Step 4 :- find 3rd Quartile :-

The third quartile is the median of the data points that is right of the median  $\Rightarrow Q_3 = 35$

Step 5 :- find Min and Max

$$\Rightarrow \text{Min} = 25, \text{Max} = 38$$



⇒ Inter Quantile Range (IQR) =  $Q_3 - Q_1$   
 ⇒ IQR =  $35 - 29 = 6$

Q2) Find out the 5 tuple summary with (IQR) let the data range be 199, 201, 236, 271, 278, 283, 291, 301, 303, 341  
 ∴  $n = 11$

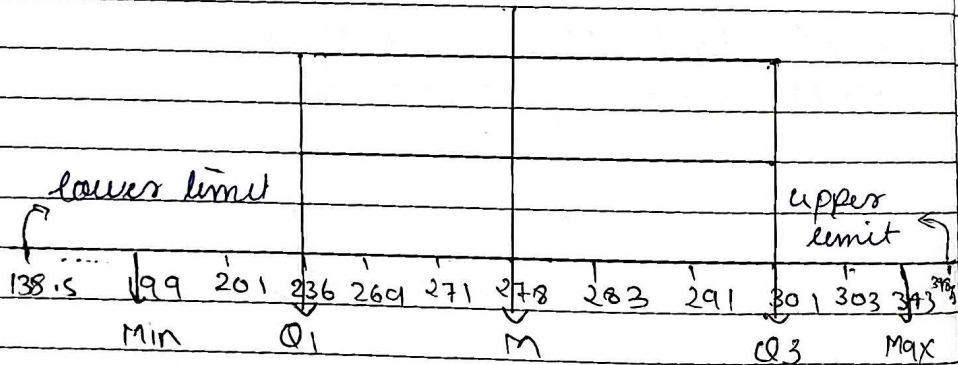
STEP 1 :- 199, 201, 236, 271, 278, 283, 291, 301, 303, 341

STEP 2 :-  $M = 278$

STEP 3 :-  $Q_1 = 236$

STEP 4 :-  $Q_3 = 301$

STEPS :- Min 199, Max 341



⇒ IQR =  $301 - 236 = 65$

Min ⇒ lower limit =  $Q_1 - 1.5(IQR)$

⇒ lower limit =  $236 - 97.5$   
 ⇒ lower limit = 138.5

⇒ Upper limit =  $Q_3 + 1.5(IQR) = 301 + 97.5$   
 ⇒ upper limit = 398.5

# List Square Method

Q1) A producer believes that the sale of his product for that years is related to an economic index and the data for that period is shown in the following tables

Years	Economic Index (x)	Sales (Y) 10000 unit
1	104	2.0
2	100	2.3
3	111	2.1
4	129	2.6
5	126	2.3
6	115	2.4
7	152	2.5
8	161	2.8
9	143	2.6
10	170	3.0
$n = 10$	$\Sigma x = 1311$	$\Sigma Y_n = 24.6$

Q2) Determine the equation of the list square line that describe the relation

b/w the sale and the economic indicator.

ii) Determine the strength of the relation b/w the two variables by computing the value of coefficient of correlation.

iii) If the value of the economic index for 11th year is 175 find out the sales in the 11th year.

Sol<sup>n</sup>  $Y = a + mx \Rightarrow Y_a = na + m \sum x$  — (1)

$\sum x Y_a = a \sum x + m \sum x^2$  square method — (2)

$\Rightarrow 24.6 = 10 \cdot a + m \cdot 1311$

Years	Economic Index (x)	Sales (Y <sub>a</sub> ) 10000 Units	x · Y <sub>a</sub>	x <sup>2</sup>	Y <sub>a</sub> <sup>2</sup>
1	104	2.0	208	10816	4.0
2	100	2.3	230	10000	5.29
3	111	2.1	233.1	12321	4.41
4	129	2.6	335.4	16641	6.76
5	126	2.3	289.8	15876	5.29
6	115	2.4	276	13225	5.76
7	152	2.5	380	23104	6.25
8	161	2.8	450.8	25921	7.84
9	143	2.6	371.8	20449	6.76
10	170	3.0	510	28900	9.0
	$\sum x = 1311$	$\sum Y_a = 24.6$	$\sum x Y_a = 32849$	$\sum x^2 = 177253$	$\sum Y_a^2 = 61.36$

$$r = \frac{n \sum x y - (\sum x)(\sum y)}{\sqrt{n \sum x^2 - (\sum x)^2} \sqrt{n \sum y^2 - (\sum y)^2}}$$

→ coefficient of relation

$\Rightarrow 24.6 = 10 \cdot a + 1311 m$  — (3)

$\Rightarrow 3284.9 = 1311a + 177253 m$  — (4)

Multiply with 1311 on eq<sup>n</sup> (3) and 10 on eq<sup>n</sup> (4)

$\Rightarrow 32250.6 = 13110 a + 1718721 m$

$\Rightarrow 32849.0 = 13110 a + 1772530 m$

$+ 5984 = + 53809 m$

$\Rightarrow m = 0.011$

$\Rightarrow 24.6 = 10a + 14.421$

$\Rightarrow a \approx 1.01$  ans//

1)  $Y_a = 1.002 + 0.011 x$

2)  $r = \frac{n \sum x y - (\sum x)(\sum y)}{\sqrt{n \sum x^2 - (\sum x)^2} \sqrt{n \sum y^2 - (\sum y)^2}}$

$= \frac{10 \times 3284.9 - (1311)(24.6)}{\sqrt{10 \times 177253 - (1311)^2} \sqrt{10 \times 61.36 - (24.6)^2}}$

$= \frac{32849 - 32250.6}{\sqrt{53809} \times \sqrt{8.44}} = \frac{598.4}{673.84}$

$\Rightarrow r = 0.89$  ans//

$$3) Y_a = 1.002 + 0.011(x)$$

$$= 1.002 + 0.011(175)$$

$$\Rightarrow Y_a = 2.927$$

### # Time Series Analysis

(Ques 1) Find the quarterly sales of the fifth year by suitable forecasting method technique for the data gives 4 years also make adjustment for expected seasonal variation

Years	Quarter(x)	Sales 1000 units	$xY_a$	$x^2$	$Y' = 1.80 + 0.185x$
1	1	1.0	1	1	$Y_1 = 1.985$
	2	3.0	6	4	$Y_2 = 2.17$
	3	4.0	12	9	$Y_3 =$
	4	2.0	8	16	
2	5	1.0	5	25	
	6	3.0	18	36	
	7	5.0	35	49	
	8	3.0	24	64	
3	9	2.0	18	81	
	10	4.0	40	100	
	11	6.0	66	121	
	12	2.0	24	144	
4	13	2.0	26	169	
	14	5.0	70	196	

15	70	105	225		
16	4.0	64	256		
$\Sigma X = 136$	$\Sigma Y_a = 54$	522	1496		

$$n = 16$$

$\Rightarrow$  We know that

$$\Sigma Y_a = na + m \Sigma X$$

$\Rightarrow$  putting values

$$54 = 16a + m \times 136$$

$$54 = 16a + 136m \quad \text{--- (1)}$$

$\Rightarrow$  now,

$$\Sigma xY_a = a \Sigma x + m \Sigma x^2$$

$\Rightarrow$  putting values

$$522 = a \times 136 + m \times 1496$$

$$522 = 136a + 1496m \quad \text{--- (2)}$$

$\Rightarrow$  multiplying eq<sup>n</sup> (1) with 11 and then eq<sup>n</sup> (2) - eq<sup>n</sup> (2)

$$594 = 176a + 2496m$$

$$522 = 136a + 1496m$$

$$72 = 40a$$

$$\Rightarrow a = \frac{72}{40} = 1.80$$

$\Rightarrow$  put  $a$  in eq<sup>n</sup> (1)

$$54 = 16 \times 1.80 + 136m$$

$$54 = 28.8 + 136m$$

$$54 - 28.8 = 136m$$

$$25.2 = 136m$$

$$\Rightarrow m = \frac{25.2}{136} \Rightarrow 0.185$$

So  $y' = 1.80 + 0.185x$

now for 5th year, we have to find quarterly sales of 5th year i.e for  $x = 17, 18, 19, 20$  so

$\Rightarrow y^{17} = 1.80 + 0.18(17) = 4.945$

$\Rightarrow y^{18} = 1.80 + 0.18(18) = 5.13$

$\Rightarrow y^{19} = 1.80 + 0.18(19) = 5.315$

$y^{20} = 1.80 + 0.18(20) = 5.5$

(Ques 2) A machine shop produces steel pins the width of the 100 pins was checked after manufacturing and data was recorded as follows

width (in mm)	freq (f)	$fix_f$	$x^2$	$fx^2$	Cell midpoint (x)
9.50 - 9.51	6	57.03	90.39	542.04	9.505
9.52 - 9.53	2	19.05	90.72	181.44	9.525
9.54 - 9.55	20	190.09	91.10	1822	9.545
9.56 - 9.57	32	306.08	91.48	2927	9.565
				.36	
9.58 - 9.59	22	210.87	91.87	2021	9.585
				.14	
9.60 - 9.61	8	76.14	92.25	738	9.605
9.62 - 9.63	6	57.75	92.64	555	9.625
				.84	

9.64 - 9.65	4	38.58	93.02	372.08	9.645
	$\Sigma f_i = 100$	957.1	733.42	9159.9	

$\rightarrow$  find out arithmetic mean, standard deviation and Variance (H.W)

Formula:- 1) Mean :-  $\bar{x} = \frac{\Sigma f_i x_i}{\Sigma f_i}$

2) Standard deviation:-  $\sigma = \sqrt{\frac{\Sigma f_i x^2 - (\bar{x})^2}{n}}$

3) Variance =  $\sigma^2$

28/08/23

Sol<sup>n</sup> ① Mean =  $\bar{x} = \frac{\Sigma f_i x_i}{\Sigma f_i} = \frac{957.1}{100} = 9.571$

$\Rightarrow$  Mean = 9.571

② Standard deviation:-  $\sigma = \sqrt{\frac{\Sigma f_i x^2 - (\bar{x})^2}{n}}$

$\Rightarrow \sigma = \sqrt{\frac{9160.12 - (9.571)^2}{100}}$

$\Rightarrow \sigma = \sqrt{91.601 - 91.604}$

$\Rightarrow \sigma = \sqrt{-0.03} \Rightarrow \sigma = 0.031$

③ Variance =  $\sigma^2 = 0.03$

(Ques 3)

The measured values are  $344^\circ, 338^\circ, 342^\circ, 335^\circ$  and  $336^\circ$  their values are constitute the first sub group

# frequency distribution:-  
 frequency distribution is a tabulation of data obtained from measurement arranged in ascending or descending order acc. to size and find out ~~determine~~ the range of the data, arithmetic mean, mid point, median & mode

$$\Rightarrow \sigma = \sqrt{\frac{(335-339)^2 + (336-339)^2 + (338-339)^2 + (342-339)^2 + (344-339)^2}{5}}$$

$$\Rightarrow \sigma = \sqrt{\frac{16 + 9 + 1 + 9 + 25}{5}} = \sqrt{\frac{60}{5}}$$

$$\Rightarrow \sigma = \sqrt{12} \Rightarrow \sigma = 3.46 \text{ Std deviation}$$

$$\Rightarrow \sigma^2 = 12 \rightarrow \text{Variance}$$

Ques 3) 5 theoretical Control are listed to determine the temp the measure are 344, 338, 342, 335 and 336 their values are constitute the 1st subgroup compute Arithmetic mean, median, range, variance, std. deviation

Ques 4) The mean and the standard deviation of a sample of 100 observation was calculated as 40 and 5.1 respectively while company with the original data it was found that by mistake a figure of 40 was mis-copied as 50 for one observation calculate the correct mean and std. deviation.

Sol<sup>n</sup> arranging  $f$  in ascending order

$$\Rightarrow \del{344}, 335, 336, 338, 342, 344$$

$$\Rightarrow \text{Median} = 338$$

$$\Rightarrow AM(\bar{x}) = \frac{335 + 336 + 338 + 342 + 344}{5} = 339$$

$$\Rightarrow \text{Range} = 344 - 335 = 9$$

$$\Rightarrow \text{std deviation} = \sqrt{\frac{(X_1 - \bar{x})^2 + (X_2 - \bar{x})^2 + \dots + (X_n - \bar{x})^2}{n}}$$

Sol<sup>n</sup>  $\therefore AM(\bar{x}) = \frac{\sum X}{N}$   
 $\therefore \sigma = \sqrt{\frac{\sum X^2}{N} - (\bar{x})^2}$

$$\textcircled{1} \bar{x} = \frac{\sum X}{N} \Rightarrow 40 = \frac{\sum X}{100}$$

$$\Rightarrow \sum X = 4000$$

$$\Rightarrow \sum X = 4000 - 50 + 40 = 3990$$

$$\Rightarrow \sum X = 3990$$

$$\bar{x} = 39.90$$

$$\Rightarrow \text{std deviation } (\sigma) = \sqrt{\frac{\sum x^2 - (\bar{x})^2}{N}}$$

$$\Rightarrow S.D = \sqrt{\frac{\sum x^2 - (40)^2}{100}}$$

$$\Rightarrow (S.D)^2 = \frac{\sum x^2 - 1600}{100}$$

$$\Rightarrow 26.01 = \frac{\sum x^2 - 1600}{100}$$

$$\Rightarrow (26.01 + 1600) \times 100 = \sum x^2$$

$$\Rightarrow \sum x^2 = 162601 \times 100$$

$$\Rightarrow \sum x^2 = 162601 - 50^2 + 40^2$$

$$162601 - 2500 + 1600$$

$$\sum x^2 = 162601 - 900$$

$$\sum x^2 = 161701 \rightarrow \text{correct}$$

$$\text{Correct std. dev.} = \sqrt{\frac{\sum x^2 - (\bar{x})^2}{N}}$$

$$= \sqrt{\frac{161701 - (39.90)^2}{100}}$$

$$= \sqrt{\frac{161701 - 1592.01}{100}}$$

$$\Rightarrow \sqrt{1617.01 - 1592.01}$$

$$\Rightarrow \sqrt{25}$$

$$\Rightarrow 5 // \rightarrow \text{correct}$$

## # Hypothesis Testing

→ state  $H_0$  as well as  $H_1$

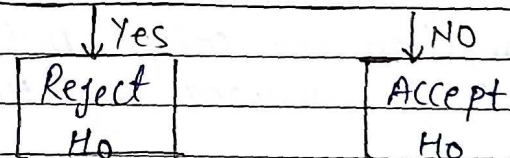
→ specify the value of significant (or the  $\alpha$  value)

→ Decide the correct sampling distribution

→ sample a ~~regular~~ <sup>random</sup> sample & work out an appropriate value for sample data

→ Calculate the probability if  $H_0$  is true

→ Is this probability equal to or similar



run this risk of committing Type I error

run this risk of committing Type II error



# flow diagram of Hypothesis testing  
 To test a hypothesis testing to tell on the basis of the data the researchers are collected whether or not the hypothesis seems to be valid.  
 In hypothesis testing the main question is whether to accept the null hypothesis or not to accept the null hypothesis. Hypothesis testing can also be defined in the following flow diagram for better understanding

where  $H_0$  is a null hypothesis and  $H_a$  denoted as alternative hypothesis  
 → A tentative assumption is made about the parameter or distribution this assumption is called Null hypothesis and is denoted by ( $H_0$ )

→ An alternative hypothesis which the opposite of null hypothesis. In alternative hypothesis maintain the relation b/w two measured variables. The hypothesis testing procedure involves using sampling data to determine whether or not  $H_0$  can be rejected. The ideal procedure in hypothesis testing leads to acceptance of

$H_0$  when  $H_0$  is true, and the rejection of  $H_0$  when  $H_0$  is false. The hypothesis testing is based on the sampling information and possibilities of error must be considered

- A Type I error corresponds to rejecting  $H_0$ , where  $H_0$  is actually true.
- And a Type II error corresponds to accepting  $H_0$ , where  $H_0$  is false.

The probability of making a Type I error is denoted by  $\alpha$ , and a probability of making a Type II error is denoted by  $\beta$

★ Z test :-

$$Z = \frac{\bar{X} - \mu_{H_0}}{\sigma_p / \sqrt{n}}$$

→ Population normal infinite

- Sample size may be large or small but variance of the population is known
- $H_a$  may be one sided or two sided in such cases Z test is used

↳ Population normal finite

- Sample size may be large or small but variance of the population is known

- $H_0$  may be one sided or two sided

in such case the z test is called finite population multiplier

$$Z = \frac{\bar{X} - \mu_{H_0}}{(\sigma_p / \sqrt{n}) \cdot (\sqrt{(N-n)/(N-1)})}$$

↳ Population normal infinite

- Sample size small, variance of the population unknown

- $H_0$  may be one sided or two sided.

in such cases T test is used

$$t = \frac{\bar{X} - \mu_{H_0}}{\sigma_s \sqrt{n}}$$

↳ Population normal finite

- Sample size small, variance of the population unknown

- $H_0$  may be one sided or two sided in such cases t test is used

$$t = \frac{\bar{X} - \mu_{H_0}}{(\sigma_s / \sqrt{n}) \sqrt{(N-n)/(N-1)}}$$

↳ Population may not be normal but

- Same size, large variance of the population may be unknown ~~or~~ known and  $H_0$  may be one sided or two sided

in such situation, Z test is used. (we can use either case 1 or case 2 formula)

for variance :- 
$$\sigma_s = \sqrt{\frac{\sum (X_i - \bar{X})^2}{(n-1)}}$$

(Ques) A sample of 400 male students is found to have a mean height is 67.47 inches can it be reason regarded as a sample from a large population with mean height 67.39 inches and standard deviation is 1.30 inches test at 5% level of significance

Soln  $H_0 = \mu_{H_0} = 67.39''$

$H_a = \mu_{H_a} \neq 67.39''$

$\bar{X} = 67.47''$ ,  $\sigma_p = 1.30$ ,  $n = 400$

Date: / / Page no: \_\_\_\_\_

$$Z = \frac{67.47 - 67.39}{1.51 / \sqrt{400}} = \frac{0.08}{0.065}$$

$$Z = 1.23$$

Considering  $H_0$  is two sided for given question

Date: / / Page no: \_\_\_\_\_

554 test using student whether the breaking strength of any lot may be taken to be 578 kg with 5% level of significance. Verify degree of freedom, arithmetic mean, standard deviation and t test result

Ans

Ans) The specimen of a copper wire drawn from a large set and have the following breaking strength in (weight in kg) 570, 572, 578, 568, 572, 578, 570, 572, 596,

# Data Warehousing :- Data warehousing governing all the data, and control multidimensional space data. The construction of data warehouse involves

- Data cleaning
- Data integration
- Data transformation

Data warehouse provides online analytical processing tool (OLAP) for interactive analysis of multidimensional data which facility effective data generalisation and data mining.

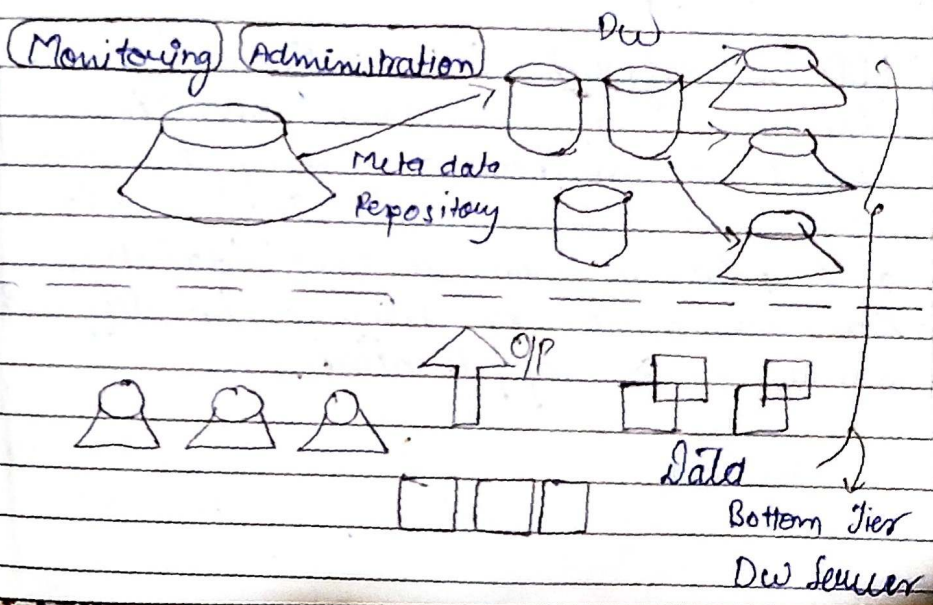
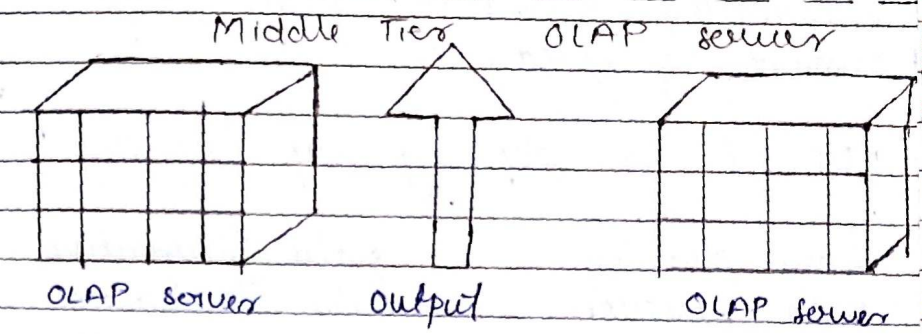
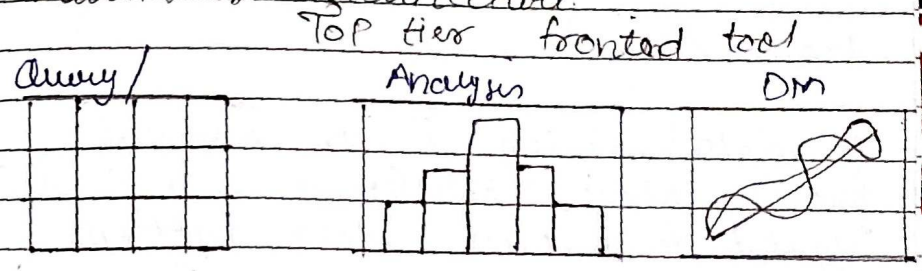
Many other data mining functions such as association, classification, prediction and clustering can be integrated with OLAP operation to enhance interactive mining of knowledge in multiple level of abstraction.

Data warehousing providing architecture and tool for business execution to systematically organising understand and use their data to make strategic decision. It is a subject oriented integrated, and time variant and non volatile collection of data in support and management decision marketing process.

# Compare OLAP & OLTP processing.  
The major distinguish feature of OLAP and OLTP are summarised as follows

features	OLTP	OLAP
user and system oriented	customer oriented and used for transaction & query processing	Market oriented and used for data analysis
Data contained	Manage current data and easily used for decision making	Large amount of data historical data summarisation, aggregation facility
Database design	E-R diagram data model, appl oriented and DP design	star & snow fall model and subject oriented data-base design
Access Pattern	Atomic transaction, concurrency control	Read only operations and complex queries
Users	Customer Oriented	Analyst, manager, it is market oriented

# Multitier Architecture

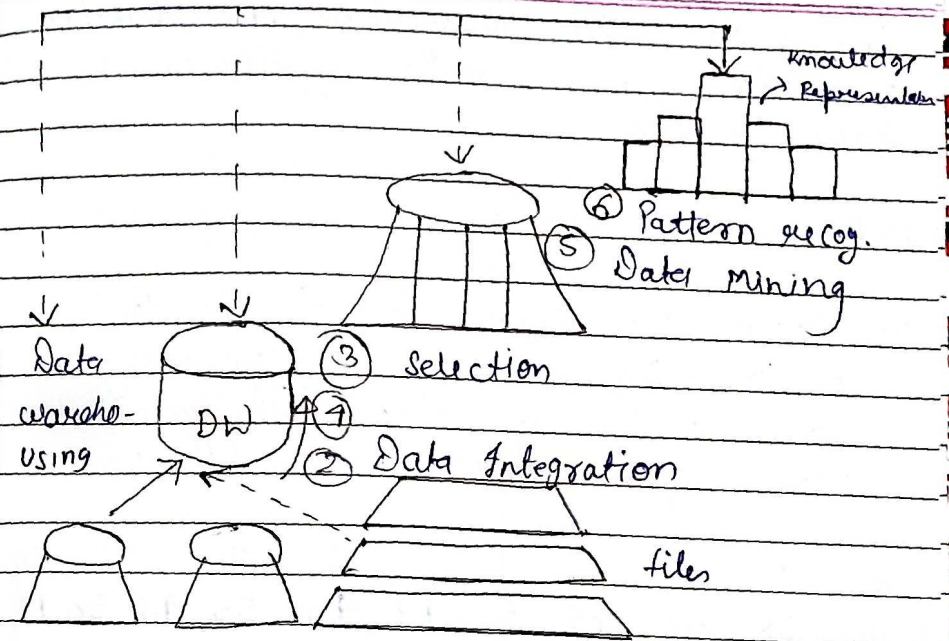


# # KDD Process (Knowledge Discovery from Data)

A data mining is an essential process where intelligence method are applied to extract data patterns. Mining is the process of discovering interesting knowledge from large amount of data.

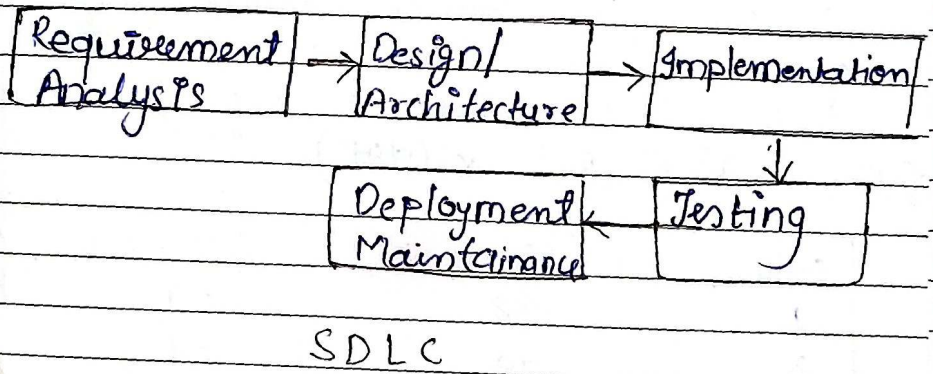
KDD is a step by step processing

1. Data cleaning
2. Data Integration
3. Data selection
4. Data Transformation
5. Data mining
6. Pattern evaluation
7. Knowledge Representation



⇒ 5 phases of software development: 05/09/2023

KDD is referred to as Knowledge Discovery in database and is defined as method of finding, transforming, and refining meaningful data and patterns from a raw database. In order to be utilised in different domains or application.



★ Internal 1st topic :- KDD process, Numerical on time series, Numerical on least square method, Theory on

box plot with Numerical, Introduction of data science with application of data science

#Chi-Square Test:-

(Ques) Apply chi-square test in the context of sampling analysis for comparing a variance to a theoretical variance weight of 10 students as follows:-

S.No	1	2	3	4	5	6	7	8	9	10
weight	38	40	45	53	47	43	55	48	52	49

Can we say that the variance of the distribution of weight for all students from which above sample of 10 students was = 20 kgs, test this at 5% level of significance

Sol<sup>n</sup>  $\chi^2 = \frac{\sigma_s^2}{\sigma_p^2} \times (dof)$

S.No	Weight (x <sub>i</sub> ) in Kg	(x <sub>i</sub> - $\bar{x}$ )
1	38	-9
2	40	-7
3	45	-2
4	53	6
5	47	0
6	43	-4

7	55	8	
8	48	1	
9	52	5	
10	49	2	

$\sum x_i = 470$   
 $\bar{x} = \frac{\sum x_i}{n} = \frac{470}{10} = 47$

$\Rightarrow \sigma_s = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}}$   
 $\Rightarrow \sigma_s = \sqrt{\frac{(-9)^2 + (-7)^2 + (-2)^2 + (6)^2 + 0^2 + (-4)^2 + (8)^2 + (1)^2 + (5)^2 + (2)^2}{9}}$   
 $\Rightarrow \sigma_s = \sqrt{\frac{280}{9}} \Rightarrow \sigma_s = \sqrt{31.11}$

Now,  $\sigma_s^2 = \text{variance} = 31.11$

Now,  $\chi^2 = \frac{31.11 \times 9}{20} = 13.99$

dof = 9  
 at 5%  $\chi^2 = 16.928$

Null hypothesis accepted