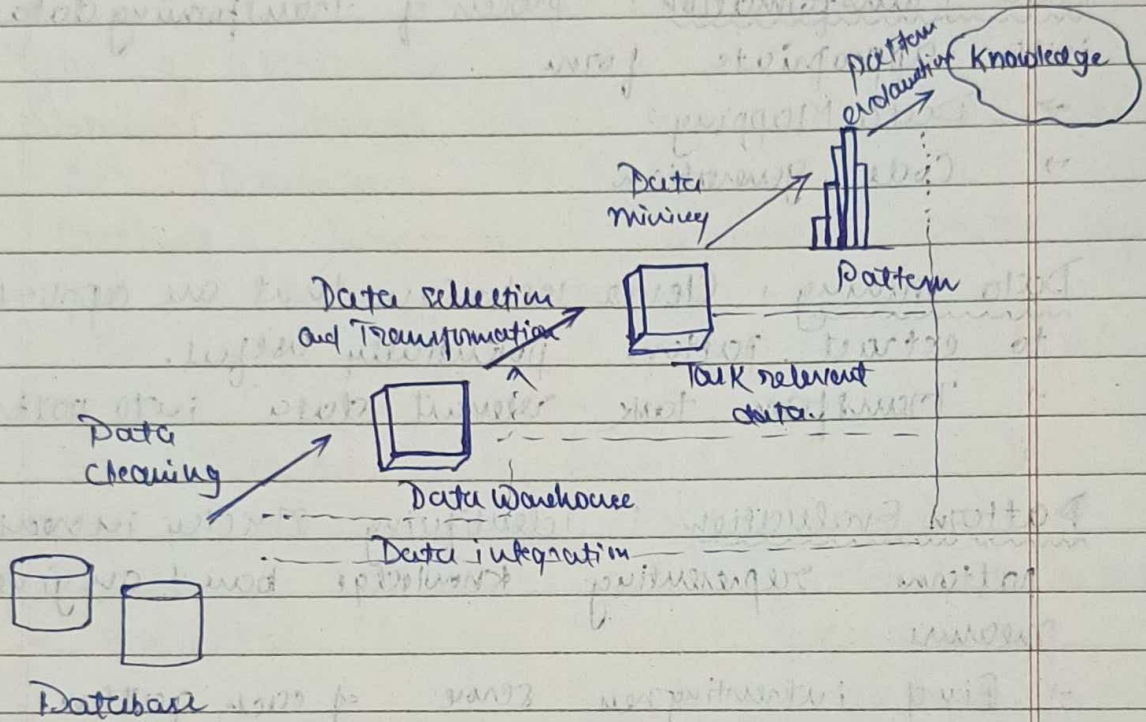


uses → \* Automatic summarization of data  
 \* Extracting essence of information  
 \* Discovering patterns in raw data

Data Mining  
 Extraction of interesting (non-trivial, implicit, previously unknown) pattern or knowledge from large amount of data.

Alternative Name: Knowledge discovery, knowledge extraction, data pattern analysis.

KDD process (Knowledge Discovery)



\* Data cleaning: Removal of noisy and irrelevant data from collection.

→ Missing values

→ Noisy data

↳ Noise refers to modification of original value

\* Data integration: Heterogeneous data are combined from multiple sources.

→ uses Data Migration tools.

→ uses Data Synchronization tools.



\* Data Selection: Data relevant to the analysis is decided and retrieved from data collection.

→ Data Selection using Neural Network

→ Decision Trees

→ Naïve Bayes

→ Clustering, Regression

\* Data Transformation: Process of transforming data into appropriate form.

→ Data Mapping:

→ Code - generation

\* Data mining: clever technique that are applied to extract patterns potentially useful.

→ Transform task relevant data into patterns

\* Pattern Evaluation: Identifying strictly increasing patterns representing knowledge based on given measure:

→ Find interestingness score of each pattern

→ Uses summarization and visualization to make data understandable by the user.

\* Knowledge representation: Techniques which utilizes visualization tools.

→ Reports

→ Tables

\* It is an iterative process where evaluation measures can be enhanced, mining can be refined.



Outlier: A data object that does not comply with the general behaviour of the data.

Application of Data mining: Healthcare, Fraud detection, Education, Education, Lie Detection, Financial Banking, Manufacturing Engineering.

Data: Collection of data objects and attributes  
↓  
property of characteristic of an object.

- \* Types of attributes: unique
  - Nominal: Things which are fixed or universal.  
Ex: ID number, Zip Code.
  - Ordinal: Ranking or Interval made by oneself accordingly.  
Ex: height category of small, medium, tall?
  - Interval: Fixed universal intervals.  
Ex: Calendar dates, temp. in Celsius or F.
  - Ratio: length, time, temp. in Kelvin.
  - Binary: Attribute with only two states (0 and 1)
    - \* Symmetric Binary: Both outcomes are equally important.
    - \* Asymmetric Binary: Outcomes not equally important.

→ Distinctness:  $= \neq$   
Order:  $< >$   
Addition:  $+$   $-$   
Multiplication:  $*$   $/$

→ Nominal: Distinctness

Ordinal: Distinctness, Order

Interval: " " addition

Ratio: " " Multiplication



- \* Continuous attributes: comes from an infinite set (i.e. real numbers, you can make them as large or small). (floating point variables)
- \* Discrete attributes: comes from a finite set or countably infinite set (i.e. integers). (binary attributes)

### Types of Missing Data

- \* MCAR: (Missing Completely at Random)  
The missingness of the variable will be unpredictable based on analysis of other data from dataset.
- \* MAR: (Missing At Random)  
Missingness is related to other variables.
- \* MNAR: (Missing Not at Random)  
To find cause of missingness of data.

### Data Preprocessing

- \* Aggregation
- \* Sampling
- \* Dimensionality Reduction
- \* Feature Subset Selection

### Euclidean Distance

$$(x_1, y_1), (x_2, y_2) \text{ dist} = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

### Manhattan Distance

$$= |x_2 - x_1| + |y_2 - y_1|$$

$$\text{Minkowski Distance} \rightarrow \text{dist} = \left( \sum_{k=1}^n |p_k - q_k|^r \right)^{1/r}$$

$r = 1$  (Manhattan Dist.)

$r = 2$  (Euclidean Dist.)

$r = \infty$  ("Supremum") This is the max diff. between any component of vectors.



\*  $SMC = M_{11} + M_{00} / (M_{01} + M_{10} + M_{00} + M_{11})$

\* Jaccard (J) =  $M_{11} / (M_{01} + M_{10} + M_{11})$

\* Cosine Similarity:

$$\cos(d_1, d_2) = (d_1 \cdot d_2) / \|d_1\| \|d_2\|$$

Ex:  $d_1 = 1 \ 2 \ 5 \ 0 \ 2$

$d_2 = 3 \ 1 \ 2 \ 7 \ 1$

$$d_1 \cdot d_2 = 1 \cdot 3 + 2 \cdot 1 + \dots + 2 \cdot 1$$

$$\|d_1\| = 1 \cdot 1 + 2 \cdot 2 + \dots + 2 \cdot 2$$

$$\|d_2\| = 3 \cdot 3 + 1 \cdot 1 + \dots + 1 \cdot 1$$

Mean:

Sum of all elements / Total no. of elements

\* Trimmed Mean:

Ex: Trimmed 20% mean for 60, 81, 83, 91, 99

$$\Rightarrow \frac{81 + 83 + 91}{3}$$

Median: Middle most element

Mode: Most repeated: Unimodal

Bimodal

Multimodal

Grouped Data:

Estimated

\* Mean: (Interval midpoint) =  $m$  and frequency,  $f_1, f_2, \dots$

$$\therefore \text{mean} = \frac{m_1 \cdot f_1 + m_2 \cdot f_2 + \dots + m_n \cdot f_n}{f_1 + f_2 + \dots + f_n}$$

\* Estimated median

$$= L + \frac{(n/k) - B}{G} \times w$$

L: lower class boundary of the group containing median

n: Total no. of values (sum of frequencies)

B: cumulative frequency of groups before median group

G: Frequency of median group.



$$\text{Estimated Mode} = 1 + \frac{(f_m - f_{m-1})}{(f_m - f_{m-1}) + (f_m - f_{m+1})} \times W$$

\* Symmetric vs. Skewed Data:

⇒ Symmetric: Mean, Median, Mode lies on same line

⇒ ⊕ve Skewed Data: Mode < Median

⇒ ⊖ve Skewed Data: Mode > Median

\* Standard Deviation

$$S_d = \sqrt{\frac{\sum (d_i)^2}{N-1}} \quad \text{where } d_i = X_i - \bar{X}$$

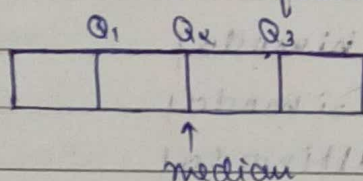
↑
mean

←
Bessel's correction

\* Variance:

$$\sigma^2 = (89)^2$$

\* Quantiles: It divides an ordered data set into four equal parts.



Even Sample Size

$$\text{median} = 71 = Q_2$$

62 63 64 64 70 72 76 77 81 81

$$Q_1 = 64$$

$$Q_3 = 77$$

Odd Sample Size

$$\text{median} = 72$$

63 64 64 70 72 76 77 81 81

$$\frac{64+64}{2} = 64 = Q_1$$

$$\frac{77+81}{2} = 79 = Q_3$$

Interquartile Range (IQR)

$$IQR = (Q_3 - Q_1)$$

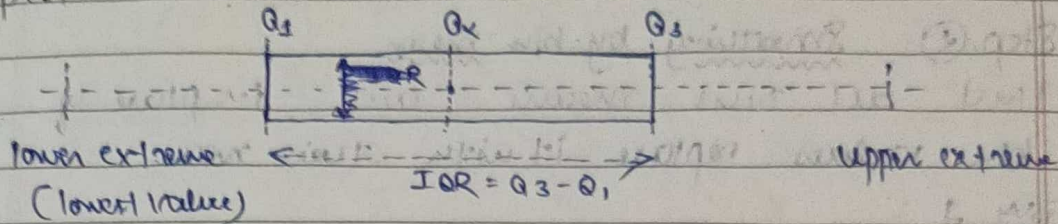
Date: / / 22

Page No.

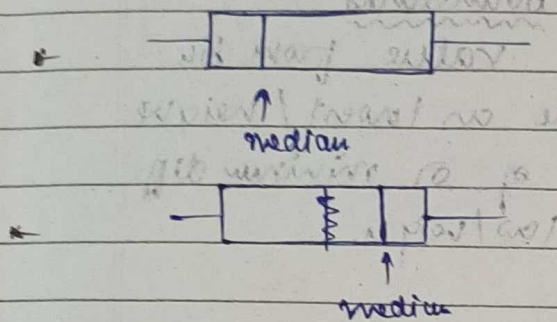


Outliers: Below  $Q_1 - 1.5(IQR)$   
 Above  $Q_3 + 1.5(IQR)$  } This boundary is called Turkey fences

Boxplot :



Boxplot Analysis :



Right / Positive Skewed  
 $(Q_3 - Q_2) > (Q_2 - Q_1)$

Left / Negative Skewed  
 $(Q_3 - Q_2) < (Q_2 - Q_1)$

Normalization

① Decimal Normalization

$$V_i' = \frac{V_i}{10^J} \quad \text{where } J = \text{no. of digits in greatest element of data set}$$

Note: Arrange the given data into ascending order.

② Max-Min Normalization

$$V_i' = \frac{V_i - \min(D)}{\max(D) - \min(D)}$$

D: data set

Main Data Handling

①

① Binning Method: (Applied to sorted data)



→ Step ① Partition into (equi-depth) bins

Bin 1: // Equal beams

Bin 2

→ Step ② Smoothing by bin mean

Find all bins means and in place of bin values replace it with their mean

Bin 1:

Bin 2:

→ Step ③ Smoothing by bin boundaries

Take lowest and highest value from bin then write the bin value as lowest/highest according to nearest element of a minimum diff bet<sup>n</sup> its element. Any low/high.

Ex: 4, 8, 9, 15, 21, 21, 24, 25

Step ① Bin 1: 4, 8, 9, 15

mean = 9

Bin 2: 21, 21, 24, 25

mean = 21.75 = 23

Step ② Bin 1: 9, 9, 9, 9

Bin 2: 23, 23, 23, 23

Step ③ Bin 1: 4, 4, 4, 15

(low = 4, high = 15)

Bin 2: 21, 21, 24, 25

24 is nearer to 25

II

Clustering: It is a method of converting a group of abstract objects into classes of similar objects.

To find centroid:

(A, B)

(x<sub>1</sub>, y<sub>1</sub>)

(x<sub>2</sub>, y<sub>2</sub>)

Centroid = (mean(x<sub>1</sub>, x<sub>2</sub>), mean(y<sub>1</sub>, y<sub>2</sub>))



Ex:

A	2	3
B	6	1
C	1	2
D	3	0

$$A-B \text{ (centroid)} = \left( \frac{2+6}{2}, \frac{3+1}{2} \right) = (4, 2)$$

$$C-D \text{ (centroid)} = (2, 1)$$

Then Square Euclidean Distance

	A	B	C	D
A-B	5	5	9	5
C-D	4	16	2	2

as A is very near to CD then ACD could be a cluster.

$$\text{Centroid of ACD} = \left( \frac{2+1+3}{3}, \frac{3+2+0}{3} \right) = (2, 1.6)$$

	A	B	C	D
B	20	0	26	10
ACD	1.7	16.44	2.1	3.7

B is very far from ACD so it will not be in cluster with ACD.

(II)

Regression: It is a method to determine the relationship between a dependent variable and one or more independent variables.

(1) Linear Regression (best line to fit two variables)

$$y = mx + b \quad \text{but here } m \text{ we will take } a$$

$$= y = ax + b$$

$$a = \frac{n(\sum xy) - \sum x \sum y}{n(\sum x^2) - (\sum x)^2} \quad b = \frac{\sum y \sum x^2 - \sum x \sum xy}{n(\sum x^2) - (\sum x)^2}$$



# Handling Redundant Data In Data Integration

## > Correlation Analysis:

$$r_{A,B} = \frac{\sum (A - \bar{A})(B - \bar{B})}{(N \times \sigma_A \times \sigma_B)} = \frac{\sum (AB) - N\bar{A}\bar{B}}{N\sigma_A\sigma_B}$$

N: no. of instance

$\bar{A}$  and  $\bar{B}$ : are mean

$\sigma_A$  and  $\sigma_B$ : standard deviation

$\sum (AB)$ : Sum of AB cross product.

if A ↑ then B ↑

\*  $r_{A,B} > 0$ : A and B are positively correlated

\*  $r_{A,B} < 0$ : A and B are negatively correlated

\*  $r_{A,B} = 0$ : independent (No correlation)

## > $\chi^2$ (Chi-Square Test)

If A ↑ then B ↓  
vice versa

\* Find Expected value

\* Find  $\chi^2 = \sum \frac{(\text{Observed} - \text{Expected})^2}{\text{Expected}}$

\* Then compare  $\chi^2$  value with critical value of Chi-Square

If  $\chi^2 < \text{critical value}$  → Null Hypothesis (we don't agree)

If  $\chi^2 > \text{critical value}$  → Alternate Hypothesis

Note: Significance level is given in question

Degree of freedom (df) =  $(R-1) * (C-1)$

R: No. of Rows, C: no. of columns

Example:

	Play chess	Not play chess	Sum
LSF	250	200	450
NLSF	50	1000	1050
Sum	300	1200	1500



Expected frequency

$$E(\text{play chess, LSF}) = \frac{\text{count}(\text{play chess}) \times \text{count}(\text{LSF})}{N} = \frac{300 \times 480}{1800} = 90$$

$$E(\text{play chess, NLSF}) = \frac{300 \times 1080}{1800} = 210$$

Expected frequency Table

	Play Chess	Not play chess
LSF	90	360
NLSF	210	840

$$\text{then } \chi^2 = \frac{(250-90)^2}{90} + \frac{(50-210)^2}{210} + \frac{(360-200)^2}{360} + \frac{(840-1000)^2}{840}$$

$$\chi^2 = 507.93$$

Given significance level = 0.001

$$\text{degree (df)} = (2-1) \times (2-1) = 1$$

So from Chi-square Table

Critical value is = 10.828

∴ As computed value (507.93) is more than 10.828, hence, ~~perfect~~ there are correlated.