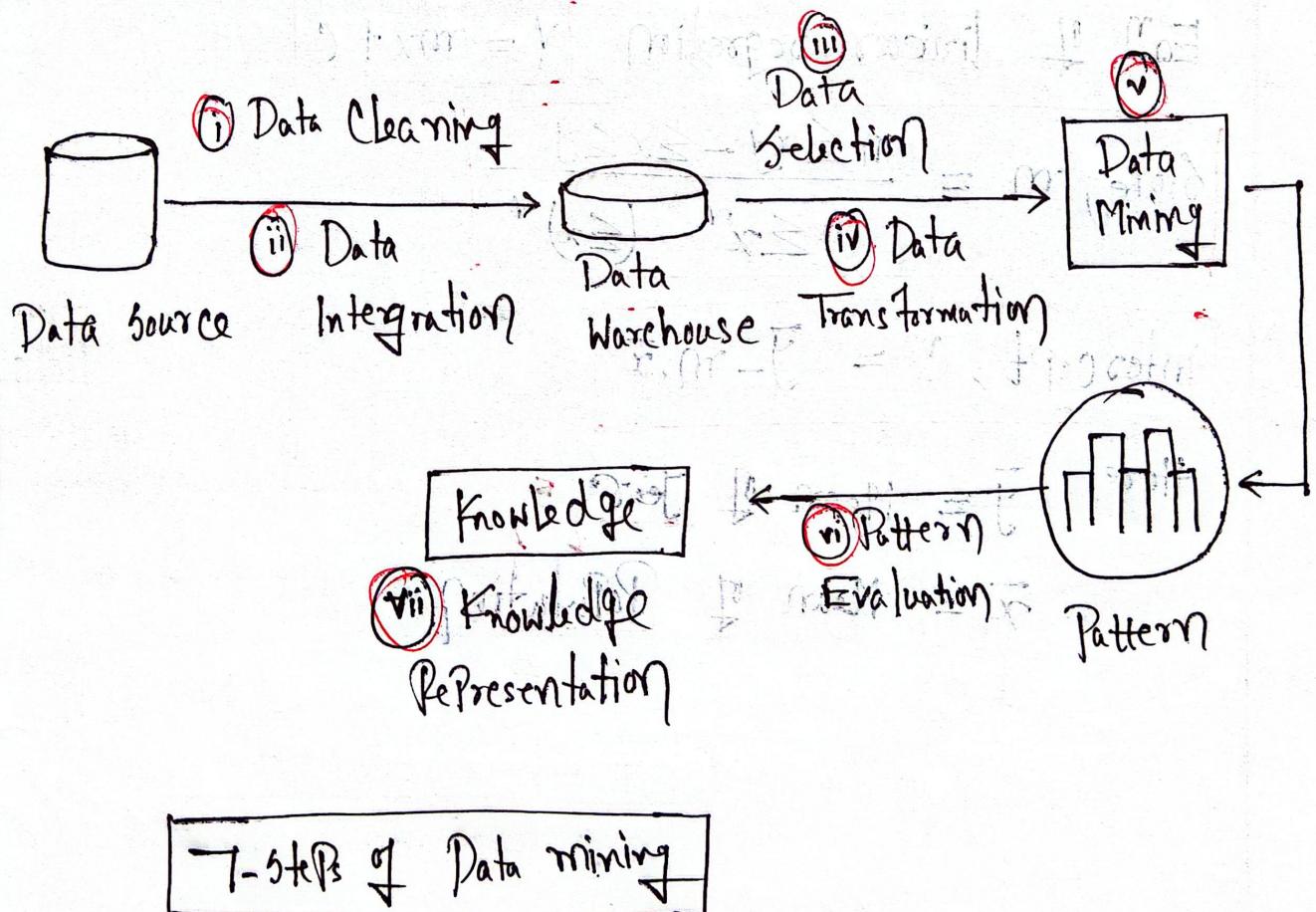


"Machine Learning & Data Mining"

Segment - 01

~~Data mining~~: Data mining is a process of extracting useful information and insights from large data set.

~~KDD Process~~: (Knowledge Discovery from Data)



~~Principle of KDD~~ & Fuzzy Logic

~~Adv. of KDD~~

- i) Decision support
- ii) Automation
- iii) Discovery of hidden Pattern

~~Dis-Adv. of KDD~~

- i) Data Quality
- ii) Computational Complexity
- iii) Data Security

~~Math Problem:~~

Eq. of linear regression, $y = mx + c$

$$\text{Slope, } m = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2}$$

$$\text{intercept, } c = \bar{y} - m \cdot \bar{x}$$

Here, \bar{y} = Mean of Years

\bar{x} = Mean of Population

~~Principle of KDD~~

SP23(1)

x	y	xy	$(xy)^n$
1980	2.1	4158	3920400
1983	2.4	9759.2	3932289
1986	2.9	5759.4	3944196
1989	3.2	6369.8	3956121
1992	3.8	7569.6	3968664
1995	4.1	8179.5	3980025
1998	4.3	8591.9	3992009
2001	4.7	9404.7	4004001
15924	27.5	59786.6	31697100

(Now, $\bar{x} = \frac{1980 + 1983 + \dots + 2001}{8} = 1990.5$)

$\bar{y} = \frac{2.1 + 2.4 + \dots + 4.7}{8} = 3.44$

$$m = \frac{n \sum xy - \sum x \sum y}{n \sum x^n - (\sum x)^n} = \frac{8 * 59786.6 - 15924 * 27.5}{8 * 31697100 - (15924)^8} = 0.127$$

$$C = 3.44 - 0.127 * 1990.5 = -299.35$$

i) For, Population(2007) = $mx + C$

$$= 0.127 * 2007 + (-299.35)$$

$$= 5.539$$

ii) Population(2010) = $0.127 * 2010 + (-299.35)$

$$= 5.92$$

Types of Data Mining:

- i) Classification
- ii) Association Analysis
- iii) Clustering
- iv) Regression
- v) Anomaly Detection.

Popular ML Algorithms:

- i) Linear Regression
(Supervised Learning)
- ii) Artificial Neural Network
- iii) Decision Tree (classification)

Segment - 02

Similarity & Dissimilarity

- i) Euclidean Distance, $d = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$ (x₁, y₁)
(x₂, y₂)
- ii) Manhattan Distance, $d = |x_1 - x_2| + |y_1 - y_2|$
- iii) Minkowski " " $d = \left\{ (|x_1 - x_2|)^p + (|y_1 - y_2|)^p \right\}^{1/p}$
- iv) supremum " ", $d = \max(|x_1 - x_2|, |y_1 - y_2|)$
- v) Cosine Distance, $d = 1 - \frac{x \cdot y}{\|x\| \|y\|}$ [x, y = vector]

e.g.: For, $x = [2, 3, 5]$ and $y = [1, 9, 6]$

$$d = 1 - \frac{(2 \cdot 1) + (3 \cdot 9) + (5 \cdot 6)}{\sqrt{2^2 + 3^2 + 5^2} \times \sqrt{1^2 + 9^2 + 6^2}} = 0.128$$

~~(*)~~ Similarity between binary vector:

i) Simple Matching Coefficient (SMC) [symmetric]

$$SMC = \frac{\text{number of matching attribute}}{\text{number of attributes}}$$

$$f_{11} + f_{00}$$

$$f_{11} = \frac{f_{11} + f_{00} + f_{01} + f_{10}}{f_{00} + f_{01} + f_{10} + f_{11}}$$

~~(*)~~ ii) Jaccard Co-efficient (J) [Asymmetric binary Attribute]

$$J = \frac{f_{11}}{f_{01} + f_{10} + f_{11}}$$

e.g.: $P = 1 0 0 0 0 0 0 0 0 0$

$$q = 0 0 0 0 0 0 1 0 0 1$$

$$f_{01} = 2$$

$$i) SMC = \frac{f_{11} + f_{00}}{f_{00} + f_{01} + f_{10} + f_{11}} = \frac{0 + 7}{2 + 1 + 7 + 0}$$

$$f_{10} = 1$$

$$f_{00} = 7$$

$$f_{11} = 0$$

$$ii) J = \frac{f_{11}}{f_{01} + f_{10} + f_{11}} = \frac{0}{2 + 1 + 0} = 0$$

~~* Types of Data / Attributes~~

~~Types of Data with example :~~

~~i) Quantitative / Numerical Data~~

(i) Discrete Data → Marks of student, Age.

(ii) Continuous Data → Weight, Height, Temperature.

(iii) Numeric → Digits

~~ii) Qualitative Data~~

(i) Nominal Data → Gender, Color, State of residence.

(ii) Ordinal Data → Ranking, ~~Bool value~~ etc.

(iii) Binary Data

~~Proximity measurement for Binary attribute:~~

Example :

~~(t,i)~~ ~~Object i~~ Col-1 Col-2 Col-3

~~Object j~~ 0 0 1

~~Object k~~ 0 1 0

~~Object l~~ 0 0 0

$$\text{Here } P_{i,j} = \frac{a}{d} \quad \text{For } (i,j) \Rightarrow a = 1 \quad b = 0 \quad c = 1 \quad d = 2$$

$$P_{i,k} = \frac{a}{d} \rightarrow 1 \quad P_{i,l} = \frac{a}{d} \rightarrow 0$$

$$P_{j,k} = \frac{b}{d} \rightarrow 0 \quad P_{j,l} = \frac{b}{d} \rightarrow 0$$

$$P_{k,l} = \frac{c}{d} \rightarrow 1$$

~~Disimilarity~~ for (i, j) :

i) Distance for symmetric binary,

$$d(i, j) = \frac{b+c}{a+b+c+d} = \frac{0+1}{1+0+1+1} = \frac{1}{3} = 0.33$$

ii) Distance for asymmetric binary,

$$d(i, j) = \frac{b+c}{a+b+c} = \frac{0+1}{1+0+1} = \frac{1}{2} = 0.5$$

~~Similarity~~ for (i, j) :

i) Jaccard (Asymmetric Binary)

$$J(i, j) = \frac{a}{a+b+c} = \frac{1}{1+0+1} = 0.5$$

ii) Cosine Co-efficient (symmetric binary)

$$C(i, j) = \frac{a+d}{a+b+c+d} = \frac{2}{2+0+0+0} = 0.67$$

$a = 0$
 $b = 0$
 $c = 0$
 $d = 2$

Thus final value for (i, k) and (j, k)

$$= \frac{1}{P} \left(\sum_{m=1}^P \min_{k=1}^M \max_{i=1}^N \min_{j=1}^N \delta_{ij} \right) = \frac{1}{P} \sum_{m=1}^P \min_{k=1}^M \max_{i=1}^N \min_{j=1}^N \delta_{ij}$$

~~Proximity Measures of Nominal Attributes:~~

$$d(i, j) = \frac{P-m}{P} \quad | P \rightarrow \text{Number of attribute} \\ m \rightarrow \text{Number of matches}$$

$$\begin{aligned} \text{sim}(i, j) &= 1 - d(i, j) \\ &= 1 - \frac{P-m}{P} = \frac{m}{P} \\ &= \frac{m}{P} = \frac{m}{P} = \frac{m}{P} \end{aligned}$$

SP'22 | 1(c)

~~Using Cosine measurement for finding similarity and dissimilarity~~

For $(\text{Doc-1}$ and $\text{Doc-2})$, at this step minimum

$$d_{\text{sim}}(D_1, D_2) = \frac{D_1 \cdot D_2}{\|D_1\| * \|D_2\|} \leftarrow \text{Numerical}$$

$$\begin{aligned} D_1 \cdot D_2 &= 5*3 + 3*2 + 2*1 + 2*1 \\ &= 15 + 6 + 2 + 2 \\ &= 25 \end{aligned}$$

$$\|D_1\| * \|D_2\| = \sqrt{5^2 + 3^2 + 2^2 + 2^2} * \sqrt{3^2 + 2^2 + 1^2 + 1^2 + 1^2}$$

$$= 6.981 * 4.123$$

$$\text{stradef} \rightarrow \approx 26.72$$

$$\frac{\alpha - \beta}{\beta} = (i, i)_b$$

~~$$d_{ij}(D_1, D_2) = \frac{25}{26.72}$$~~

$$(i, i)_{D_1} = 1 - 0.99 = 0.01$$

$$\text{so, } d_{ij}(D_1, D_2) = 1 - 0.99 = 0.01$$

(This similarly calculate for (D_1, D_3) (D_1, D_4) (D_2, D_3) (D_2, D_4))

(D_3, D_4) ——————

————— \rightarrow This minimize him

i) Minimum value will be most similar.

ii) Maximum Value \rightarrow Most dissimilar

$$\|D_1\| * \|D_2\|$$

$$1^2 + 2^2 + 3^2 + 2^2 + 1^2 = 15$$

$$1 + 2 + 3 + 2 + 1 =$$

$$10 =$$

~~For Numeric Attribute~~: $(A_i^{(1)}) \rightarrow (A_i^{(n)})$

$d(i,j) = \frac{|i^{th} \text{ Attribute value} - j^{th} \text{ attribute value}|}{\text{Max attribute value} - \text{Min attribute value}}$

~~#~~ For ~~Ordinal~~ Attribute:

Obj	T
1	High
2	Low
3	Medium
4	High

$M_f = \text{Total unique Attribute values}$

~~#~~ Step 1: Give initial value of attributes

$$\text{High} = 1, \text{medium} = 2, \text{Low} = 3$$

~~#~~ Step 2: Normalize the ranking between (0 to 1)

$$z_i = \frac{\text{Initial value of } i^{th} \text{ attribute} - 1}{M_f - 1}$$

$$\text{High} = \frac{1 - 1}{3 - 1} = \frac{0}{2} = 0.0 = (1,0)_{0.5}$$

$$\text{Medium} = \frac{2 - 1}{3 - 1} = \frac{1}{2} = \frac{0.5}{1} = (2,0)_{0.5}$$

$$\text{Low} = \frac{3 - 1}{3 - 1} = 1$$

~~Step 4 : $d(2,1) = (1-0) = 1$~~ ~~indicates diagonal~~

$d(3,1) = (\frac{1}{2} - 0) = \frac{1}{2}$ ~~indicates off-diagonal~~

$d(4,1) = (0 - 0) = 0$ ~~indicates below~~

Indicates engine total = $\frac{1}{2}M$

$$\begin{bmatrix} 0 \\ 1 & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 \\ 0 & 1 & \frac{1}{2} & 0 \end{bmatrix}$$

Indicates number of bits

$$\begin{bmatrix} 1 \\ 1 \\ 2 \\ 2 \\ 2 \end{bmatrix}$$

Indicates P vector Latin min : $\frac{1}{2}M$

~~SP'22 | 2(a)~~

$C = \text{WV}$, $S = \text{number}$, $T = \text{dBfH}$

For Test \rightarrow (Nominal) Indicates off diagonal is 0 bits.

$d(i,j) = \frac{P - \text{Indicates off P vector Latin vector}}{P}$

$$d(2,1) = \frac{1-0}{\frac{1}{2}M} = 1$$

$$d(3,1) = \frac{1-\frac{1}{2}}{\frac{1}{2}} = 0 \Rightarrow \begin{bmatrix} 0 & & \\ \frac{1}{2} & 0 & \\ 0 & 1 & 0 \end{bmatrix} = \text{dBfH}$$

$$d(3,2) = \frac{1-\frac{1}{2}}{\frac{1}{2}} = 0 \Rightarrow \begin{bmatrix} 0 & & \\ \frac{1}{2} & 0 & \\ 0 & \frac{1}{2} & 0 \end{bmatrix} = \text{mild M}$$

$$\therefore P = \frac{1-\theta}{1-\phi} = \text{WV}$$

for Test-3 (Numeric)

$$d(2,1) = \frac{|220 - 490|}{490 - 121} = \frac{220}{319} = 0.69$$

$$d(3,1) = \frac{|121 - 490|}{319} = \frac{369}{319} = 1.16$$

$$d(3,2) = \frac{|121 - 220|}{319} = 0.31$$

for Test-4 (Binary)

$$d(2,1) \Rightarrow a=0, b=1, c=0, d=0$$

$$d(2,1) = \frac{1+0}{0+1+0+0} = 1$$

$$d(3,1) = \frac{0+0}{0+1+0+0} = 0$$

$$d(3,2) = \frac{0+1}{0+1+0+0} = 1$$

Mixed Attribute :

$$\begin{bmatrix} 0 & & \\ (1+.19+1)/3 & 0 & \\ (0+1+0)/3 & (1+.31+1)/3 & 0 \end{bmatrix} = \begin{bmatrix} 0 & & \\ .897 & 0 & \\ .33 & .77 & 0 \end{bmatrix}$$

~~* Reasons for missing data:~~

(10marks) Q. - 1

- i Data entry errors
 - ii Technical issues
 - iii Incomplete data collection
 - iv Privacy concern
 - v Outliers : Data Points that fall outside the expected range may be excluded as outliers, even though they are valid.
- ~~Q = 3 + 0 = 3, r = d, o = s~~ $\leftarrow (1.8) b$

~~3. Missing Value Imputation techniques:~~

(1.8) b

- i Mean / median : Replace missing value with average.
 - ii Mode : Use most frequent value in gap.
 - iii Zero imputation : Replace missing value with '0' / Constant
- ~~: Student & berim~~

$$\begin{bmatrix} 0 & 0 \\ 0 & \text{F.R. E.C.} \end{bmatrix} = \begin{bmatrix} 0 & \frac{c(1+0+1)}{3} \\ 0 & \frac{c(1+1+0)}{3} \end{bmatrix}$$

Segment - (3)

#

Apriori Algorithm:

Example: Min Support - 50% $\frac{4}{8} = 50\%$

Confidence \rightarrow 70% $\frac{4}{5} = 80\%$

Tid	Items
100	1 3 4
200	2 3 5
300	1 2 3 5
400	2 5

Step 1:

Step 1:

Items	Support
1	$2/4 = 50\%$
2	$3/4 = 75\%$
3	$3/4 = 75\%$
4	$1/4 = 25\%$
5	$3/4 = 75\%$

As $4 \rightarrow 25\%$ is below support 50%.

So, 4 is eliminated.

Step 2:

(1, 2, 3, 5)

Items	Support	Confidence	Interest
(1, 2)	$1/4 = 25\%$ X	0.75	0.45
(1, 3)	$2/4 = 50\%$	0.75	0.50
(1, 5)	$1/4 = 25\%$ X	0.75	0.25
(2, 3)	$2/4 = 50\%$	0.75	0.25
(2, 5)	$3/4 = 75\%$	0.75	0.50
(3, 5)	$2/4 = 50\%$	0.75	0.25

Step 3:

Items	Support	Confidence	Interest	Order
X (1, 2, 3)	$1/4 = 25\%$	0.75	0.45	1
X (1, 2, 5)	$1/4 = 25\%$	0.75	0.50	2
(2, 3, 5)	$1/4 = 25\%$	0.75	0.25	3
X (1, 3, 5)	$1/4 = 25\%$	0.75	0.25	4

1st frequent subset is 1, 2, 3 → 3rd
most frequently 1, 3, 5

Step 4:

Rules/AV	Support	Confidence
$(2 \wedge 3) \rightarrow 5$	0.112 = 11.2%	100% / $(0.112 - 0.01)$
$(3 \wedge 5) \rightarrow 2$	0.112 = 11.2%	100% / $(0.112 - 0.01)$
$(2 \wedge 5) \rightarrow 3$	0.112 = 11.2%	66% / $(0.112 - 0.01)$
initial 3 $\rightarrow (2 \wedge 5)$	0.112 = 11.2%	66% / $x - b + c$
$2 \rightarrow (3 \wedge 5)$	0.112 = 11.2%	66% / x
$5 \rightarrow (2 \wedge 3)$	0.112 = 11.2%	66% / x

$$\text{Confidence} = \frac{s(A \cup B)}{s(A)}$$

~~fibonacci~~

$$\text{For, } (2 \wedge 3) \rightarrow 5 \quad \text{Confidence} = \frac{s((2 \wedge 3) \cup 5)}{s(2 \wedge 3)}$$

~~A B~~

$$= \frac{\frac{2x}{2} / \frac{2}{2}}{\frac{2x}{2} / \frac{10}{2}} = 100\%$$

As the first 2 is above Confidence level 100%

so, $(2 \wedge 3) \rightarrow 5$ & $(3 \wedge 5) \rightarrow 2$ final association rules.

$$T_{2,0} = \frac{0\bar{1} - \bar{0}\bar{0}}{0\bar{0} - \bar{1}\bar{0}} = \bar{x}$$

$$T = \frac{0\bar{1} - \bar{1}\bar{0}}{0\bar{0} - \bar{1}\bar{0}} = \bar{x}$$

~~i) Min-Max Normalization~~

$$\bar{x} = \frac{x - \text{Min}}{\text{Max} - \text{Min}}$$

(Subtract Min from each value)

~~ii) Z-score Normalization:~~

$$z = \frac{x - \text{Mean}}{\text{std}}$$

(Subtract Mean from each value)

x = original value
 Min = Min value of Column
 Max = Max value of Column
 Mean = Mean of Column
 std = Standard deviation of Column

~~Aut'22 | 3.5~~

~~Humidity~~ \rightarrow $(\text{A}) \rightarrow (\text{QVA}) \rightarrow \text{embroid}$

$$\text{std} = \sqrt{\frac{(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + (x_3 - \bar{x})^2 + (x_4 - \bar{x})^2}{4}}$$

$\bar{x} = \frac{53 + 69 + 50 + 51}{4} = 56$

Soln: Min-Max normalization works in 2 steps with 2 formulas

$$\text{Value } \bar{x}_1 = \frac{53 - 50}{69 - 50} = 0.21$$

$\bar{x} \leftarrow (0.21) \quad 0 \leftarrow (0.21), 0.21$

$$\bar{x}_2 = \frac{69 - 50}{69 - 50} = 1$$

$$\bar{x}_3 = \frac{50 - 50}{69 - 50} = 0$$

Dress Code in Islam

$$\bar{x}_3 = \frac{50 - 50}{69 - 50} = 0$$

$$\bar{x}_4 = \frac{51 - 50}{19} = 0.07$$

Humidity

0.21
0.21
1
0
0.07

Z-score:

$$\text{mean} = \frac{53 + 53 + 69 + 50 + 51}{5} = 54.2 = \mu$$

$$\text{std} = \sqrt{\frac{1}{n} ((x_0 - \mu)^2 + (x_1 - \mu)^2 + (x_2 - \mu)^2 + \dots + (x_n - \mu)^2)}$$

$$= 3.99$$

$$Z_0 = \frac{53 - 59.2}{3.94} = -1.6$$

$$Z_1 = -0.3$$

$$Z_2 = \frac{69 - 59.2}{3.94} = 2.48$$

~~model fit measure of what?~~

~~Five Number Summary:~~

- (i) Min
- (ii) Q_1
- (iii) Median
- (iv) Q_3
- (v) Max

~~Exercise:~~

Data set $\Rightarrow 1, 3, 5, 6, 9, 10, 12, 15, 17, 20$

$$\text{Q} = \frac{Q_1 + Q_3}{2} = \frac{9 + 17}{2} = 13$$

$$P_{10} = \frac{Q_1 + P_1}{2} = \frac{9 + 10}{2} = 9.5$$

~~Five Number Summary:~~

$$\text{Min} = 1$$

$$Q_1 = \frac{3+5}{2} = 4$$

$$\text{Median} = \frac{9+10}{2} = 9.5$$

$$Q_3 = \frac{15+17}{2} = 16$$

$$\text{Max} = 20$$

ii) IQR Method :

$$\text{IQR} = Q_3 - Q_1$$
$$= 16 - 4 = 12$$

$$\text{Lower bound} = Q_1 - (1.5 * \text{IQR})$$
$$= 4 - (1.5 * 12)$$
$$= -14$$

$$\text{Upper bound} = Q_3 + (1.5 * \text{IQR})$$
$$= 16 + (1.5 * 12)$$
$$= 34$$

iii) Box-Plot of Data

