

2018

①

(a) Data KDD - Knowledge discovery in database, It is a process of taking information out of an unknown database & comprises of following steps:-

- 1) Data selection :- We select the data from multiple sources & combine them if needed
- 2) Data preprocessing : It involves cleaning of data like handling missing values, duplicate data etc
- 3) Data transformation : We transform the data according to our needs
- 4) Data mining : Finding hidden patterns from the databases which cannot be seen
- 5) Evaluating patterns.

Data mining is just a step involved in KDD

②

- b) (i) eye color - Nominal
- (ii) grades - Ordinal
- (iii) dates - interval
- (iv) age - ratio

(c) $0 \leq \text{Gini index} \leq 1$

$$1 - \left[\left(\frac{1}{6} \right)^2 + \left(\frac{5}{6} \right)^2 \right]$$

$$= 1 - \left[\frac{1}{36} + \frac{25}{36} \right]$$

$$= \frac{10}{36} = \boxed{\frac{5}{18}}$$

(d)₁₀ ** Graph data structure(e) $p_1(0, 2)$ $p_2(2, 0)$ $p_3(3, 1)$ $p_4(5, 1)$

$$p_1 p_2 = \sqrt{(2)^2 + (2)^2} = 2\sqrt{2}$$

$$p_3 p_4 = \sqrt{(2)^2 + (0)^2} = 2$$

(f) ** (Categorical attribute)(g)₂₀ Interval Ratio

1) Operations like

 $+/-$ can be performed

2) Values can be

specified within a

given range

eg - temp.

1) Operations like

 \times / \div can also be

performed

2) Doesn't make sense
when value is zero

eg - length, height,

age . . .

(h) Eager learner : - Here, a model is prepared first on the basis of the training set and then this model is tested on the test set. eg - decision tree, rule based classifier

5) Lazy learner : Rule based classifiers are lazy learners because they don't build a model rather they learn the whole training set & if records that matches training set comes, they are able to classify them but if any record outside training set comes, they could not classify them as no info is available
10) eg - KNN

(i) Apriori principle : If an itemset is frequent, then all its subsets must also be frequent. On contrary, if an itemset is infrequent, all its supersets are also infrequent.

25) As we know that support count of an itemset never exceeds that of its subset, hence the given statement is true according to Apriori principle.

(j)

$$\{ 18, 21, 22, 25 \}$$

0-1

$$\text{new_min} = 0$$

$$\min = 18$$

$$\text{new_max} = 1$$

$$\max = 25$$

$$n' = \frac{n - \min}{\max - \min} \quad (\text{new_max} - \text{new_min}) + \text{new_min}$$

$$(i) \quad n' = \frac{18 - 18}{25 - 18} \quad () \quad = \underline{\underline{0}}$$

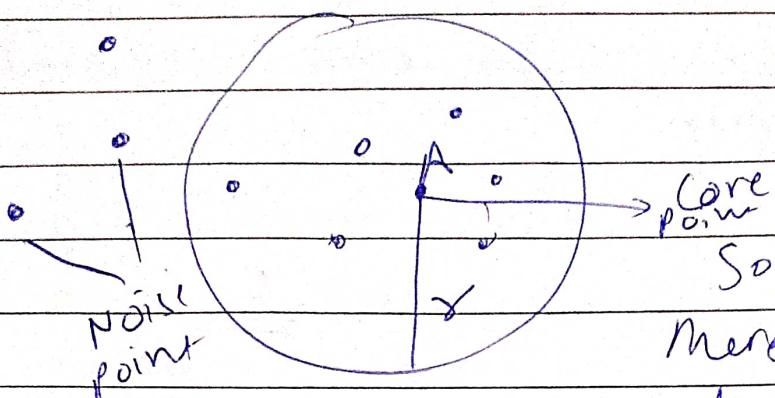
$$ii) \quad n' = \frac{21 - 18}{25 - 18} \quad (1) \quad = \frac{3}{7}$$

$$iii) \quad n' = \frac{22 - 18}{25 - 18} \quad = \frac{4}{7}$$

$$iv) \quad n' = \frac{25 - 18}{25 - 18} \quad = \underline{\underline{1}}$$

(k) (i) Core point: When we take any point as the centre & draw a circle of radius r , & if the no. of points inside the circle including that point exceeds the minimum value, we can term that point as core point.

(ii) Noise point : All the points which lie completely outside the circle we draw of radius r are termed as noise points



Given :

$$r, \text{ min } = 7$$

So, we can see that there are 7 points inside the circle. Here A is a core point.

(d) Mutually exclusive :- If all the records in the database are triggered by almost one rule, then they are mutually exclusive rules. So, for mutual exclusiveness, two rules should not be triggered by one record.

If rules are not mutually exclusive, then we cannot determine to which class we should classify the record.

Such problem can be resolved by using two methods

i) Ordered method : Here, we give priority to every rule on basis of its accuracy & coverage and classify the record to the rule with higher priority.

2) Unordered method: Here, we see the frequency of every class in dataset and classify the record to the class which is most frequent.

SECTION-B

$$\textcircled{2} \quad \text{a) (i)} \quad \begin{aligned} \{e\} &= 8 \\ \{b, d\} &= 2 \\ \{b, d, e\} &= 2 \\ \{a, b, d, e\} &= 1 \end{aligned}$$

$$\text{(ii)} \quad \{b, d\} \rightarrow \{e\}$$

$$\begin{aligned} \text{Confidence} &= \frac{\sigma(b \cup d \cup e)}{\sigma(b \cup d)} \\ &= \frac{2}{2} = 100\% \end{aligned}$$

$$\{e\} \rightarrow \{b, d\}$$

$$\begin{aligned} \text{Confidence} &= \frac{\sigma(e \cup b \cup d)}{\sigma(e)} \\ &= \frac{2}{8} = 25\% \end{aligned}$$

(iii) No, confidence is not a symmetric measure as confidence for both the parts above is diff

(b) ~~* * metric~~

(3)

(a) Aggregation :- Aggregation is a process where we combine multiple attributes into one for better analysis.

For eg. When we are having elections in country, to calculate voter turn out by state, we sum up all the individual votes obtained in ~~all~~ all the zones of state and there is one ~~total~~ ~~no. of~~ number for every state (Total no. of - votes)

Uses :-

- 1) Reduces the dimension of the dataset
- 2) Avoid analysis on unnecessary info. which is not needed
- 3) We can visualize the data better as we have reduced no. of attributes

(b) ~~noise~~ noise

(i) It is just the unwanted data which occurs in a dataset

Outliers

(i) The data points with characteristics diff from most of the datapoints in dataset

(i) No, noise is never interesting or desirable as it is just the unwanted data which needs to be removed.

- (ii) Yes, outliers can be desirable for example in fraud detection where there is a datapoint way different than most of the datapoints
- 5
(iii) ~~Ans~~ No, noise objects can appear as normal data as well
- 10
(iv) NO, outliers are not always noise because they can be useful sometimes & removing them will be inefficient

(M) (a) Confusion matrix

		Predicted class	
		Class = 1	Class = 0
Actual class	Class = 1	3 TP	3 FN
	Class = 0	1 FP	3 TN

20
Accuracy = $\frac{TP + TN}{TP + TN + FP + FN} = \frac{3+3}{3+3+3+1} = \frac{6}{10} = 0.6$

Error = $\frac{FP + FN}{TP + TN + FP + FN} = \frac{3+1}{10} = 0.4$

b) k-fold cross validation

- (a) When our dataset is divided into k equal partitions, one of them is used as test set & other $k-1$ as training set.
- Whole data can be used as test & training set.

Hold out method: The dataset is divided into two partitions - test set & training set.

- The proportion in which we divide can be 75-25, 50-50 etc. But, it can ignore some of the dataset for testing.

⑤

(a) Partition based clustering

- Makes use of centroid model in which data points are clustered on the basis of their distance from centroid.

e.g. K-means

Hierarchical based clustering

- Nesting clustering (A cluster is clustered into another cluster).
represented by dendograms

** ADV & DISADV

(b) Simple random sampling Randomly select the sample from dataset

(c) (i) Birds $\{ RL \text{ is triggered} \}$

(ii) Mammals $\{ R3 \text{ is triggered} \}$

(6) (a) ABE

A	$\rightarrow B$	$AB \rightarrow E$	$\text{Conf} = 2/3 = 66\%$
B	$\rightarrow A$	$AE \rightarrow B$	$\text{Conf} = 2/2 = 100\%$
A	$\rightarrow E$	$BE \rightarrow A$	$\text{Conf} = 2/4 = 50\%$
E	$\rightarrow A$	$E \rightarrow AB$	$\text{Conf} = 2/4 = 50\%$
B	$\rightarrow E$	$B \rightarrow AE$	$\text{Conf} = 2/5 = 40\%$
E	$\rightarrow B$	$A \rightarrow BE$	$\text{Conf} = 2/4 = 50\%$

(b) Strong association rules

$AB \rightarrow E$

$AE \rightarrow B$

$BE \rightarrow A$

$E \rightarrow AB$

$A \rightarrow BE$

(7) b)

$$2, 4, 10, 12, 3, 20, 30, 11, 25$$

$k=3$

$$\mu_1 = 2 \quad \mu_2 = 4 \quad \mu_3 = 6$$

Points	Distance from μ_1	Distance from μ_2	Distance from μ_3	Cluster
2	0	2	4	C ₁
4	2	0	2	C ₂
10	8	6	4	C ₃
12	10	8	6	C ₃
3	1	1	3	C ₁
20	18	16	14	C ₃
30	28	26	24	C ₃
11	9	7	5	C ₃
25	22	21	19	C ₃

$$C_1 = \{2, 3\}$$

$$C_2 = \{4\}$$

$$C_3 = \{10, 12, 20, 30, 11, 25\}$$

$$\mu'_1 = 2.5$$

$$\mu'_2 = 4$$

$$\mu'_3 = 18$$

Points	$D(\mu_1)$	$D(\mu_2)$	$D(\mu_3)$	Cost
2	0.5	2	16	C_1
4	1.5	0	14	C_2
5 10	7.5	6	8	C_2
12	9.5	8	6	C_3
3	0.5	1	15	C_2 C_1
20	17.5	16	2	C_3
30	27.5	26	12	C_3
10 11	8.5	7	7	C_2
25	22.5	21	7	C_3

$$C_1 = \{ 2, 3 \}$$

$$C_2 = \{ 4, 10, 11, 25 \}$$

$$C_3 = \{ 12, 20, 30 \}$$

$$\mu_1' = 0 \quad \cancel{2.5} \quad 2.5 \approx 2$$

$$\mu_2' = 0 \quad \cancel{8} \quad \approx 8$$

$$\mu_3' = \cancel{22.5} \quad \cancel{25} \quad \approx 22$$

2 8 22

Points	$D(\mu_1)$	$D(\mu_2)$	$D(\mu_3)$	Cluster
2	0	6	20	C ₁
4	2	4	18	C ₁
10	8	2	12	C ₂
12	10	4	10	C ₂
3	1	5	19	C ₁
20	18	12	2	C ₃
30	28	22	8	C ₃
11	9	3	11	C ₂
25	23	17	3	C ₃

$$C_1 = \{2, 4, 3\}$$

$$C_2 = \{10, 12, 11\}$$

$$C_3 = \{20, 30, 25\}$$

$$\mu_1' = 3$$

$$\mu_2' = 11$$

$$\mu_3' = 25$$

3

11

25

Points	$D(\mu_1)$	$D(\mu_2)$	$D(\mu_3)$	cluster
2	1	9	23	C_1
4	1	7	21	C_1
10	7	1	15	C_2
12	9	1	13	C_2
3	0	8	22	C_1
20	17	9	5	C_3
30	27	19	5	C_3
11	8	0	14	C_2
25	22	14	0	C_3

$$C_1 = \{2, 4, 3\}$$

$$C_2 = \{10, 12, 11, 4\}$$

$$C_3 = \{20, 30, 25\}$$

15

$$\mu_1' = 3$$

$$\mu_2' = 11$$

$$\mu_3' = 25$$

Now, as values are stagnate, we stop the process.