CMRIT
CMR INSTITUTE OF TECHNOLOGY, BENGALURU.
ACCREDITED WITH A+ GRADE BY NAAC
CELEBRATING 25 YEARS

| Sub: | Data Mining and Data warehouse | | | Sub Code: | 18CS641/17 CS651 | Branch: | ISE | | |
|------|------|------|------|------|------|------|------|------|------|
| Date: | 23/06/2021 | Duration: | 90 min's | Max Marks: | 50 | Sem/Sec: | VI  A,B&C | | OBE |
| **Answer any FIVE FULL Questions** | | | | | | | | MARKS | CO | RBT |

| | | MARKS | CO | RBT |
|------|------|------|------|------|
| 1a) | Compare the data mining tasks clustering and classification with example. | [4] | CO1 | L2 |

**At least 4 points- 4 marks**

| Classification | Clustering |
|------|------|
| Find a model for class attribute as a function of the values of other attributes. | Given a set of data points, each having a set of attributes, and a similarity measure among them, find clusters such that Data points in one cluster are more similar to one another and Data points in separate clusters are less similar to one another. |
| previously unseen records should be assigned a class as accurately as possible. | Similarity measures like Euclidean Distance, Manhattan are used |
| Mode is built based on training data set and tested with testing data set. | Data set does not have class labels |
| Decision Tree, Naive Bayes, KNN are examples Classifying bank customer, Predict fraudulent cases in credit card transactions. | K-Means, DBSCAN, Agglomerative Clustering Market Segmentation, Grouping the documents, and News |

| | | MARKS | CO | RBT |
|------|------|------|------|------|
| 1b) | Explain the various data quality problems and how to handle them. | [6] | CO1 | L2 |

Data mining focuses on
(1) the detection and correction of data quality problems (called data cleaning.)
(2) the use of algorithms that can tolerate poor data quality.

☐    **Examples of data quality problems:[ 1 Mark]**
•    Noise and outliers
•    missing values
•    duplicate data

**Noise:** Noise refers to modification of original values **[ 1 Mark]**
Examples: distortion of a person's voice when talking on a poor phone and "snow" on television screen
**Outliers :**Outliers are data objects with characteristics that are considerably different than most of the other data objects in the data set**[ 1 Mark]**

**Missing Values: [ 2 Marks]**

   ☐    **Reasons for missing values**
- Information is not collected (e.g., people decline to give their age and weight)
- Attributes may not be applicable to all cases (e.g., annual income is not applicable to children)

   ☐    **Handling missing values**
- Eliminate Data Objects
- Estimate Missing Values
- Ignore the Missing Value During Analysis
- Replace with all possible values (weighted by their probabilities)

**Duplicate Data: [ 1 Mark]**

   ☐    Data set may include data objects that are duplicates, or almost duplicates of one another

   ☐    Major issue when merging data from heterogeneous sources

   ☐    Examples:

Same person with multiple email addresses

   ☐    **Data cleaning**

Process of dealing with duplicate data issues

---

**2a)** Why data preprocessing is significant step data mining? Explain Discretization and Binarization techniques with example.  [2+4]  CO1  L2

Data preprocessing makes the data suitable for data mining

It reduces noise, duplicates, processing time and memory requirements and removes unwanted data. **[ 2 Marks]**

# Binarization: [ 2 Marks]

- **Binarization** maps a continuous or categorical attribute into one or more binary variables. Typically used for association analysis
- **Often convert a continuous attribute to a categorical attribute and then convert a categorical attribute to a set of binary attributes**
- Association analysis needs asymmetric binary attributes

**Table 2.5.** Conversion of a categorical attribute to three binary attributes.

| Categorical Value | Integer Value | $x_1$ | $x_2$ | $x_3$ |
|---|---|---|---|---|
| awful | 0 | 0 | 0 | 0 |
| poor | 1 | 0 | 0 | 1 |
| OK | 2 | 0 | 1 | 0 |
| good | 3 | 0 | 1 | 1 |
| great | 4 | 1 | 0 | 0 |

**Table 2.6.** Conversion of a categorical attribute to five asymmetric binary attributes.

| Categorical Value | Integer Value | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ |
|---|---|---|---|---|---|---|
| awful | 0 | 1 | 0 | 0 | 0 | 0 |
| poor | 1 | 0 | 1 | 0 | 0 | 0 |
| OK | 2 | 0 | 0 | 1 | 0 | 0 |
| good | 3 | 0 | 0 | 0 | 1 | 0 |
| great | 4 | 0 | 0 | 0 | 0 | 1 |

# Discretization[ 2 Marks]
- Discretization is the process of converting a continuous attribute into an ordinal attribute
- Discretization is typically applied to attributes that are used in **classification**
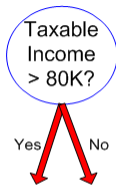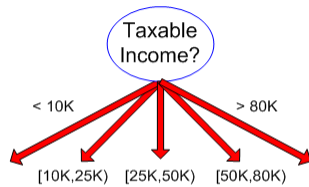
|  | or association analysis. | | | |
|  | • Transformation of a continuous attribute to a categorical attribute involves two subtasks: | | | |
|  | • **deciding how many categories to have and determining how to map the values of the continuous attribute to these categories.** | | | |
|  | • | | | |
|  | (i) Binary split      (ii) Multi-way split | | | |
| 2b) | Compute the cosine similarity of the following two objects that represent two document vectors:<br>x = (2, -7,0,2,0, -3) , y =( -1, 1,- 1,0,0, -1)<br>**x.y - [ 1 mark]**<br>**\|x\| - [ 1 mark]**<br>**\|y \| - [ 1 mark]**<br>**Cosine Similarity = [1 mark]**<br><br>Compute cosine Similarity:<br>$x = (2, -7, 0, 2, 0, -3)$<br>$y = (-1, 1, -1, 0, 0, -1)$<br><br>$x \cdot y = (2 \times -1) + (-7 \times 1) + (0 \times -1) + (2 \times 0) + (0 \times 0) + (-3 \times -1)$<br>$= -6.$<br>$|x| = \sqrt{(2 \times 2) + (-7 \times -7) + (0 \times 0) + (2 \times 2) + (0 \times 0) + (-3 \times -1)}$<br>$= \sqrt{4 + 49 + 4 + 9} = \sqrt{66} = 8.12$<br>$= 8.12$<br>$|y| = \sqrt{((-1) \times -1) + (1 \times 1) + (-1 \times -1) + 0 + 0 + (-1 \times -1)}$<br>$= \sqrt{1 + 1 + 1 + 1} = \sqrt{4} = 2.$<br>$|y| = 2.$<br>$Cos(x, y) = -6 / 8.12 \times 2 = \underline{-0.36}$ Ans. | [4] | CO1 | L3 |
| 3 | Write and Explain Apriori Algorithm for frequent itemset generation with example.<br>**Algorithm with Explanation -[ 5 marks]**<br>• The algorithm initially makes a single pass over the data set to determine the support of each item. Upon completion of this step, the set of all frequent 1-itemsets, $F_1$, will be known (steps 1 and 2).<br><br>• Next, the algorithm will iteratively generate new candidate $k$-itemsets using the frequent $(k-1)$-itemsets found in the previous iteration (step 5). Candidate generation is implemented using a function called apriori-gen, which is described in Section 6.2.3. | [5+5] | CO3 | L2 |

- To count the support of the candidates, the algorithm needs to make an additional pass over the data set (steps 6–10). The subset function is used to determine all the candidate itemsets in $C_k$ that are contained in each transaction $t$. The implementation of this function is described in Section 6.2.4.

- After counting their supports, the algorithm eliminates all candidate itemsets whose support counts are less than *minsup* (step 12).

- The algorithm terminates when there are no new frequent itemsets generated, i.e., $F_k = \emptyset$ (step 13).

---

**Algorithm 6.1** Frequent itemset generation of the *Apriori* algorithm.

1: $k = 1$.
2: $F_k = \{\, i \mid i \in I \wedge \sigma(\{i\}) \geq N \times minsup \,\}$.     {Find all frequent 1-itemsets}
3: **repeat**
4:     $k = k + 1$.
5:     $C_k = \text{apriori-gen}(F_{k-1})$.     {Generate candidate itemsets}
6:     **for each** transaction $t \in T$ **do**
7:         $C_t = \text{subset}(C_k, t)$.     {Identify all candidates that belong to $t$}
8:         **for each** candidate itemset $c \in C_t$ **do**
9:             $\sigma(c) = \sigma(c) + 1$.     {Increment support count}
10:         **end for**
11:     **end for**
12:     $F_k = \{\, c \mid c \in C_k \wedge \sigma(c) \geq N \times minsup \,\}$.     {Extract the frequent $k$-itemsets}
13: **until** $F_k = \emptyset$
14: Result $= \bigcup F_k$.

---

# Example with explanation- [ 5 marks]



Figure 6.5. Illustration of frequent itemset generation using the *Apriori* algorithm.

| | | | | |
|---|---|---|---|---|
| 4a) | A database has five transactions. Let min sup=60% and min conf =80%. <br> <br> TID      items bought | [8] | CO3 | L3 |

| | |
|---|---|
| T100 | {M, O, N, K, E, Y} |
| T200 | {D, O, N, K, E, Y} |
| T300 | {M, A, K, E} |
| T400 | {M, U, C, K, Y} |
| T500 | {C, O, O, K, I, E} |

Find all frequent item sets using Apriori algorithm.

**Computation of C1 = 2 marks**
**Computation of F1 = 1 marks**
**Computation of C2 = 1 marks**
**Computation of F2 = 1 marks**
**Computation of C3 = 1 marks**
**Computation of F3 = 2 marks**



| 4b) | What is anti-monotone property with respect to support of an itemset? Explain with example | [2] | CO3 | L2 |
|---|---|---|---|---|

**Property: 1 mark**
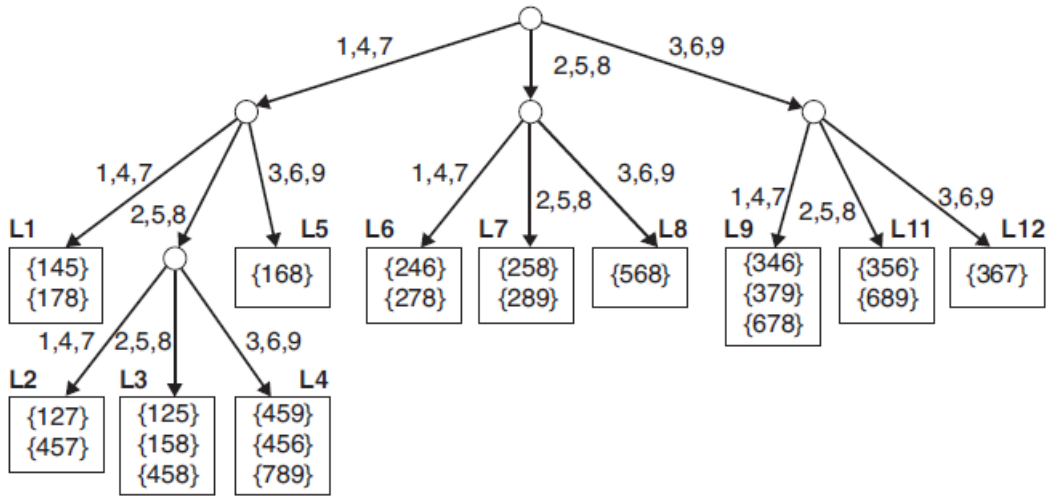**Example: 1 mark**

Apriori principle holds due to the following property of the support measure:

$$\forall X, Y : (X \subseteq Y) \Rightarrow s(X) \geq s(Y)$$

- Support of an itemset never exceeds the support of its subsets
- This is known as the anti-monotone property of support

Example:
Consider {a,b,c} itemset. Assume Support count (a,b,c)=5    support count(a,b)
>=Support count(a,b,c) Number of times a,b,c occuring in Transactions never
exceeds Support count(a,b)

| 5 | Construct the Hash tree structure for the given set of candidate item sets. Given transaction {1, 3, 4, 5, 8}, which of the hash tree leaf nodes will be visited when finding the candidates of the transaction?<br>{145}{178}{127}{457}{125}{158}{458}{459}{456}{789}{168}{246}{278}{258}{289}{568}{346}{379}{678}{356}{689}{367} | [10] | CO3 | L3 |
|---|---|---|---|---|

**Construction of Hash Tree:  5 marks**

```
                              ○
              1,4,7      /    |    \      3,6,9
                       /   2,5,8   \
                    ○         ○         ○
          1,4,7 / |3,6,9  1,4,7/|3,6,9  1,4,7/|3,6,9
            2,5,8                2,5,8         2,5,8
      L1    /   ○    L5   L6 / L7 |    \L8  L9/ L11  \L12
   {145}      / | \   {168}  {246} {258} {568} {346} {356} {367}
   {178}     /  |  \         {278} {289}       {379} {689}
        1,4,7 2,5,8 3,6,9                      {678}
      L2   /  L3 |  \ L4
   {127}  {125}  {459}
   {457}  {158}  {456}
          {458}  {789}
```

Given a transaction that contains items {1, 3, 4, 5, 8}, The leaf nodes visited are L1, L3, L5, L9, and L11 will be visited when finding the candidates of the transaction. **5 marks**

| 6 | Explain the Apriori rule Generation algorithm with example | [10] | CO2 | L2 |
|---|---|---|---|---|

**Algorithm: 5 marks**

**Example: 5 marks**

**Algorithm 6.2** Rule generation of the *Apriori* algorithm.

1: **for** each frequent $k$-itemset $f_k$, $k \geq 2$ **do**
2:   $H_1 = \{i \mid i \in f_k\}$      {1-item consequents of the rule.}
3:   call ap-genrules($f_k, H_1$.)
4: **end for**

**Algorithm 6.3** Procedure ap-genrules($f_k, H_m$).

1: $k = |f_k|$    {size of frequent itemset.}
2: $m = |H_m|$    {size of rule consequent.}
3: **if** $k > m + 1$ **then**
4:   $H_{m+1} =$ apriori-gen($H_m$).
5:   **for** each $h_{m+1} \in H_{m+1}$ **do**
6:     $conf = \sigma(f_k)/\sigma(f_k - h_{m+1})$.
7:     **if** $conf \geq minconf$ **then**
8:       **output** the rule $(f_k - h_{m+1}) \longrightarrow h_{m+1}$.
9:     **else**
10:       **delete** $h_{m+1}$ from $H_{m+1}$.
11:     **end if**
12:   **end for**
13:   call ap-genrules($f_k, H_{m+1}$.)
14: **end if**

## Rule generation:   Min conf = 50%

Assume $k = 3$.

$f_3 = \{B, C, E\}$ frequent item set.

Alg 6.2 produces
$$H_1 = \{B\} \{C\} \{E\}$$

Alg 6.3 ap-genrules

① $k = |f_k| = 3$   ② $m = |H_m| = 1$.
   $k = 3$.

③ $3 > 2$

for $H_{m+1} = H_2 = $ apriori_gen $(H_1)$

$$H_2 = \{B, C\} \{B, E\} \{C, E\}.$$

$\sigma(B,C,E) = 3$
$\sigma(E) = 5$
$\sigma(C) = 4$
$\sigma(B) = 5$.

④ for $\{BC\} \Rightarrow$ conf $= \sigma(BCE) / \sigma(\underset{f_k - h_{m+1}}{BCE - BC}) = \dfrac{3}{\sigma(E)} = \dfrac{3}{5} = 60\% > $ minconf   (50%)

⑤ ∴ $E \to BC$ is valid rule. —①

for $(B,E)$ conf $= \sigma(BCE) / \sigma(BCE - BE) = \dfrac{\sigma(BCE)}{\sigma(C)} = \dfrac{3}{4} = 75\% > $ min conf%   (50%)

$C \to BE$ is valid rule. →②

for $(CE)$ conf $= \sigma(BCE) / \sigma(BCE - CE) = \dfrac{\sigma(BCE)}{\sigma(B)} = \dfrac{3}{5} = 60\% > $ min conf   (50%)

∴ $B \to CE$ is valid rule →③

### Recursive call:
$$\text{ap-genrules} \left( \underset{f_k}{\{B, C, E\}}, \underset{H_{m+1}}{\{\{BC\}, \{B, E\}\{CE\}\}} \right)$$

2nd call

$k = |f_k| = 3$.
$m = |H_m| = 2$ (size of the rule consequent)
$3 > 3$
$k > 2 + 1$
   false
      End if.