

**Test: CLA-T2**

**Course Code & Title: 18CSC355T Data Mining and Analytics**

**Year & Sem: III Year / V Sem**

**Date:**

**Duration: 1 Hour**

**Max. Marks: 50**

**Course Articulation Matrix: (to be placed)**

S.No.	Course Outcome	PO1	PO2	PO3	PO4	PO5	PO6	PO7	PO8	PO9	PO10	PO11	PO12
1	CO1	L	H		H	L				L	L		H
2	CO2	M	H		H	L				M	L		H
3	CO3	M	H		H	L				M	L		H
4	CO4	M	H		H	L				M	L		H
5	CO5	H	H		H	L				M	L		H

<b>Part - A</b> <b>(10 x 1 = 10 Marks)</b>						
<b>Instructions: Answer all</b>						
<b>Q. No</b>	<b>Question</b>	<b>Marks</b>	<b>BL</b>	<b>CO</b>	<b>PO</b>	<b>PI Code</b>
<b>1</b>	In some cases, telecommunication companies desire to segment their clients into distinct groups in order to send suitable and related subscription offer. This can be considered as an example of which of the following methods? a. Supervised learning b. <b>Unsupervised learning</b> c. Serration d. Data extraction	<b>1</b>	<b>L1</b>	<b>3</b>	<b>1</b>	<b>1.7.1</b>
<b>2</b>	What is the relation between a candidate and frequent itemsets? (a) A candidate itemset is always a frequent itemset <b>(b) A frequent itemset must be a candidate itemset</b> (c) No relation between these two (d) Strong relation with transactions	<b>1</b>	<b>L2</b>	<b>2</b>	<b>1</b>	<b>1.7.1</b>
<b>3</b>	Which algorithm requires fewer scans of data? <b>(a) FP Growth</b> (b) Naïve Bayes (c) Apriori (d) Decision Tree	<b>1</b>	<b>L2</b>	<b>2</b>	<b>2</b>	<b>1.7.1</b>
<b>4</b>	For the question given below consider the data Transactions : 1. I1, I2, I3, I4, I5, I6 2. I7, I2, I3, I4, I5, I6 3. I1, I8, I4, I5 4. I1, I9, I10, I4, I6 5. I10, I2, I4, I11, I5 With support as 0.6 find all frequent itemsets? <b>(a) &lt;I1&gt;, &lt;I2&gt;, &lt;I4&gt;, &lt;I5&gt;, &lt;I6&gt;, &lt;I1, I4&gt;, &lt;I2, I4&gt;, &lt;I2, I5&gt;, &lt;I4, I5&gt;, &lt;I4, I6&gt;, &lt;I2, I4, I5&gt;</b> <b>(b) &lt;I2&gt;, &lt;I4&gt;, &lt;I5&gt;, &lt;I2, I4&gt;, &lt;I2, I5&gt;, &lt;I4, I5&gt;, &lt;I2, I4, I5&gt;</b> <b>(c) &lt;I11&gt;, &lt;I4&gt;, &lt;I5&gt;, &lt;I6&gt;, &lt;I1, I4&gt;, &lt;I5, I4&gt;, &lt;I11, I5&gt;, &lt;I4,</b>	<b>1</b>	<b>L2</b>	<b>2</b>	<b>2</b>	<b>2.6.3</b>

	I6>, <I2, I4, I5> (d)<I1>, <I4>, <I5>, <I6>															
5	Which of the following is a classification algorithm? a. K-means <b>b. Decision tree</b> c. Apriori d. DBSCAN	1	L1	3	1	1.7.1										
6	Frequency of occurrence of an itemset is called as ____ (a) Support (b) Confidence <b>(c) Support Count</b> (d) Rules	1	L1	2	1	1.7.1										
7	In the example of predicting number of babies based on storks’ population size, number of babies is <b>(a) Outcome</b> (b) Feature (c) Attribute (d) observation	1	L2	3	5	1.7.1										
8	Which of the following is not a type of decision tree node? a. Root node b. Leaf node c. Decision node <b>d. Branch node</b>	1	L2	3	1	1.7.1										
9	How do you calculate Confidence (A -> B)? <b>(a) Support(A ∩ B) / Support (A)</b> (b) Support(A ∩ B) / Support (B) (c) Support(A ∪ B) / Support (A) (d) Support(A ∪ B) / Support (B)	1	L2	2	1	1.7.1										
10	The classification or mapping of a class using a predefined class or group is called: a. Data Sub Structure b. Data Set <b>c. Data Discrimination</b> d. Data Characterisation	1	L2	3	1	1.7.1										
<b>Part – B</b> <b>(5x 5 = 25 Marks)</b>																
1	Consider the data set D. Given the minimum support 2, apply Apriori algorithm on this dataset <table><tr><th>Transaction ID</th><th>Items</th></tr><tr><td>100</td><td>A,C,D</td></tr><tr><td>200</td><td>B,C,E</td></tr><tr><td>300</td><td>A,B,C,E</td></tr><tr><td>400</td><td>B,E</td></tr></table>	Transaction ID	Items	100	A,C,D	200	B,C,E	300	A,B,C,E	400	B,E	5	L3	2	2	2.5.2
Transaction ID	Items															
100	A,C,D															
200	B,C,E															
300	A,B,C,E															
400	B,E															

	<div><p>min sup = 2</p><table><thead><tr><th>TID</th><th>Items</th></tr></thead><tbody><tr><td>100</td><td>A, C, D</td></tr><tr><td>200</td><td>B, C, E</td></tr><tr><td>300</td><td>A, B, C, E</td></tr><tr><td>400</td><td>B, E</td></tr></tbody></table><p>1<sup>st</sup> scan <math>C_1</math></p><table><thead><tr><th>Itemset</th><th>sup</th></tr></thead><tbody><tr><td>{A}</td><td>2</td></tr><tr><td>{B}</td><td>3</td></tr><tr><td>{C}</td><td>2</td></tr><tr><td>{E}</td><td>3</td></tr></tbody></table><p><math>F_1</math></p><table><thead><tr><th>Itemset</th><th>sup</th></tr></thead><tbody><tr><td>{A}</td><td>2</td></tr><tr><td>{B}</td><td>3</td></tr><tr><td>{C}</td><td>2</td></tr><tr><td>{E}</td><td>3</td></tr></tbody></table><p><math>C_2</math></p><table><thead><tr><th>Itemset</th><th>sup</th></tr></thead><tbody><tr><td>{A,B}</td><td>1</td></tr><tr><td>{A,C}</td><td>2</td></tr><tr><td>{A,E}</td><td>1</td></tr><tr><td>{B,C}</td><td>2</td></tr><tr><td>{B,E}</td><td>3</td></tr><tr><td>{C,E}</td><td>2</td></tr></tbody></table><p><math>F_2</math></p><table><thead><tr><th>Itemset</th><th>sup</th></tr></thead><tbody><tr><td>{A,C}</td><td>2</td></tr><tr><td>{B,C}</td><td>2</td></tr><tr><td>{B,E}</td><td>3</td></tr><tr><td>{C,E}</td><td>2</td></tr></tbody></table><p><math>C_3</math></p><table><thead><tr><th>Itemset</th><th>sup</th></tr></thead><tbody><tr><td>{B,C,E}</td><td>2</td></tr></tbody></table><p><math>F_3</math></p></div>	TID	Items	100	A, C, D	200	B, C, E	300	A, B, C, E	400	B, E	Itemset	sup	{A}	2	{B}	3	{C}	2	{E}	3	Itemset	sup	{A}	2	{B}	3	{C}	2	{E}	3	Itemset	sup	{A,B}	1	{A,C}	2	{A,E}	1	{B,C}	2	{B,E}	3	{C,E}	2	Itemset	sup	{A,C}	2	{B,C}	2	{B,E}	3	{C,E}	2	Itemset	sup	{B,C,E}	2					
TID	Items																																																															
100	A, C, D																																																															
200	B, C, E																																																															
300	A, B, C, E																																																															
400	B, E																																																															
Itemset	sup																																																															
{A}	2																																																															
{B}	3																																																															
{C}	2																																																															
{E}	3																																																															
Itemset	sup																																																															
{A}	2																																																															
{B}	3																																																															
{C}	2																																																															
{E}	3																																																															
Itemset	sup																																																															
{A,B}	1																																																															
{A,C}	2																																																															
{A,E}	1																																																															
{B,C}	2																																																															
{B,E}	3																																																															
{C,E}	2																																																															
Itemset	sup																																																															
{A,C}	2																																																															
{B,C}	2																																																															
{B,E}	3																																																															
{C,E}	2																																																															
Itemset	sup																																																															
{B,C,E}	2																																																															
2	<div><p><b>How does tree pruning approach works?</b></p><p>There are two common approaches to tree pruning:</p><ul style="list-style-type: none"><li>-prepruning</li><li>-Postpruning</li></ul><p>Prepruning approach</p><ul style="list-style-type: none"><li>A tree is “pruned” by halting its construction early e.g., by deciding not to further split or partition the subset of training tuples at a given node</li><li>Upon halting, the node becomes a leaf.</li><li>The leaf may hold the most frequent class among the subset tuples or the probability distribution of those tuples.</li><li>When constructing a tree, measures such as statistical significance, information gain, Gini index, and so on, can be used to assess the goodness of a split.</li><li>If partitioning the tuples at a node would result in a split that falls below a prespecified threshold, then further partitioning of the given subset is halted.</li><li>There are difficulties, however, in choosing an appropriate threshold.</li><li>High thresholds could result in oversimplified trees, whereas low thresholds could result in very little simplification.</li></ul><p>Postpruning Approach</p><ul style="list-style-type: none"><li>The second and more common approach is which removes subtrees from a “fully grown” tree.</li><li>A subtree at a given node is pruned by removing its branches and replacing it with a leaf.</li><li>The leaf is labeled with the most frequent class among the subtree being replaced.</li></ul><div><p>(a)</p></div></div>	5	L3	3	2	1.7.1																																																										
3	<div><p><b>Explain the various metrics for evaluating the classifier performance.</b></p><p><b>Metrics for Evaluating classifier Performance</b></p><ul style="list-style-type: none"><li>It presents measures for assessing how good or how “accurate” your classifier is at predicting the class label of tuples</li></ul></div>	5	L1	2	1	2.6.3																																																										

	<ul style="list-style-type: none"><li>Consider the case of where the class tuples are more or less evenly distributed, as well as the case where classes are unbalanced<ul style="list-style-type: none"><li>e.g., where an important class of interest is rare such as in medical tests</li></ul></li><li>They include accuracy (also known as recognition rate), sensitivity (or recall), specificity, precision, <math>F1</math>, and <math>F</math>.</li><li>Note that although accuracy is a specific measure, the word “accuracy” is also used as a general term to refer to a classifier’s predictive abilities.</li><li>Using training data to derive a classifier and then estimate the accuracy of the resulting learned model can result in misleading overoptimistic estimates due to overspecialization of the learning algorithm to the data.</li></ul> <table><tr><th>Measure</th><th>Formula</th></tr><tr><td>accuracy, recognition rate</td><td><math>\frac{TP+TN}{P+N}</math></td></tr><tr><td>error rate, misclassification rate</td><td><math>\frac{FP+FN}{P+N}</math></td></tr><tr><td>sensitivity, true positive rate, recall</td><td><math>\frac{TP}{P}</math></td></tr><tr><td>specificity, true negative rate</td><td><math>\frac{TN}{N}</math></td></tr><tr><td>precision</td><td><math>\frac{TP}{TP+FP}</math></td></tr><tr><td><math>F</math>, <math>F_1</math>, <math>F</math>-score, harmonic mean of precision and recall</td><td><math>\frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}</math></td></tr><tr><td><math>F_\beta</math>, where <math>\beta</math> is a non-negative real number</td><td><math>\frac{(1 + \beta^2) \times \text{precision} \times \text{recall}}{\beta^2 \times \text{precision} + \text{recall}}</math></td></tr></table> <p><b>Meaning of the various measures</b></p> <ul style="list-style-type: none"><li>True positives (TP): These refer to the positive tuples that were correctly labeled by the classifier. Let TP be the number of true positives.</li><li>True negatives (TN): These are the negative tuples that were correctly labeled by the classifier. Let TN be the number of true negatives.</li><li>False positives (FP): These are the negative tuples that were incorrectly labeled as positive (e.g., tuples of class <i>buys_computer = no</i> for which the classifier predicted <i>buys_computer = yes</i>). Let FP be the number of false positives.</li><li>False negatives (FN): These are the positive tuples that were mislabeled as negative (e.g., tuples of class <i>buys_computer = yes</i> for which the classifier predicted <i>buys_computer = no</i>). Let FN be the number of false negatives.</li></ul>	Measure	Formula	accuracy, recognition rate	$\frac{TP+TN}{P+N}$	error rate, misclassification rate	$\frac{FP+FN}{P+N}$	sensitivity, true positive rate, recall	$\frac{TP}{P}$	specificity, true negative rate	$\frac{TN}{N}$	precision	$\frac{TP}{TP+FP}$	$F$ , $F_1$ , $F$ -score, harmonic mean of precision and recall	$\frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$	$F_\beta$ , where $\beta$ is a non-negative real number	$\frac{(1 + \beta^2) \times \text{precision} \times \text{recall}}{\beta^2 \times \text{precision} + \text{recall}}$					
Measure	Formula																					
accuracy, recognition rate	$\frac{TP+TN}{P+N}$																					
error rate, misclassification rate	$\frac{FP+FN}{P+N}$																					
sensitivity, true positive rate, recall	$\frac{TP}{P}$																					
specificity, true negative rate	$\frac{TN}{N}$																					
precision	$\frac{TP}{TP+FP}$																					
$F$ , $F_1$ , $F$ -score, harmonic mean of precision and recall	$\frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$																					
$F_\beta$ , where $\beta$ is a non-negative real number	$\frac{(1 + \beta^2) \times \text{precision} \times \text{recall}}{\beta^2 \times \text{precision} + \text{recall}}$																					
4	<p><b>Suppose the data contain the frequent itemset <math>I = \{I1, I2, I5\}</math>. What are the association rules that can be generated from <math>I</math>? The nonempty subsets of <math>I</math> are <math>\{I1, I2\}</math>, <math>\{I1, I5\}</math>, <math>\{I2, I5\}</math>, <math>\{I1\}</math>, <math>\{I2\}</math> and <math>\{I5\}</math>. Show the resulting association rules and its confidence.</b></p> <p>Consider a database, D, consisting of 9 transactions. Suppose min. support count required is 2 (i.e. <math>\text{min\_sup} = 2/9 = 22\%</math>). Let the minimum confidence required is 70%. We have to first find out the frequent itemset using Apriori algorithm.</p> <p>Then, Association rules will be generated using min. support &amp; min. confidence.</p> <p>Step 1: Generating 1-Itemset Frequent Pattern</p> <table><tr><th>Itemset</th><th>Count</th></tr><tr><td><math>\{I1\}</math></td><td>6</td></tr><tr><td><math>\{I2\}</math></td><td>7</td></tr></table>	Itemset	Count	$\{I1\}$	6	$\{I2\}$	7	5	L2	2	2	2.5.2										
Itemset	Count																					
$\{I1\}$	6																					
$\{I2\}$	7																					

	<p> {I3}                6  {I4}                2  {I5}                2 </p> <p> The above table is L1.  In the first iteration of the algorithm, each item is a member of the set of candidate.  The set of frequent 1-itemsets, L1, consists of the candidate 1-itemsets satisfying minimum support.  Step 2: Generating 2-Itemset Frequent Pattern  To discover the set of frequent 2-itemsets, L2, the algorithm uses L1 Join L1 to generate a candidate set of 2-itemsets, C2. Next, the transactions in D are scanned and the support count for each candidate itemset in C2 is accumulated (as shown in the middle table).  The set of frequent 2-itemsets, L2, is then determined, consisting of those candidate 2-itemsets in C2 having minimum support.  <math display="block">L1 \times L1 = \{I1, I2, I3, I4, I5\}.</math> Since <math>L2 = L1 \text{ join } L1</math> then <math>\{I1, I2, I3, I4, I5\} \text{ join } \{I1, I2, I3, I4, I5\}.</math>  It becomes <math>\rightarrow C2 = [ \{I1, I2\}, \{I1, I3\}, \{I1, I4\}, \{I1, I5\}, \{I2, I3\}, \{I2, I4\}, \{I2, I5\}, \{I3, I4\}, \{I3, I5\}, \{I4, I5\} ].</math>  Now we need to check the frequent itemsets with min support count.  Then we get <math>\rightarrow (C2 * C2) L2 = [ \{I1, I2\}, \{I1, I3\}, \{I1, I5\}, \{I2, I3\}, \{I2, I4\}, \{I2, I5\} ].</math>  Similarly, We do it for L3.  Step 3: Generating 3-Itemset Frequent Pattern  <math>L2 = [ \{I1, I2\}, \{I1, I3\}, \{I1, I5\}, \{I2, I3\}, \{I2, I4\}, \{I2, I5\} ].</math>  <math>L3 = L2 \text{ JOIN } L2 \text{ i.e.}</math>  <math>C3 = [ \{I1, I2, I3\}, \{I1, I2, I5\} ].</math>  Now, the Join step is complete and the Prune step will be used to reduce the size of C3. Prune step helps to avoid heavy computation due to large Ck.  Procedure Step 1: Find Items starting with I2 in B  It gives <math>\{I1, I2, I3\}, \{I1, I2, I4\}, \{I1, I2, I5\}.</math>  Step 2: Find Items starting with I3 in B  It gives NIL, Similarly I4, I5.  Step 3: Find out infrequent items sets using min support count and remove them.  Based on the Apriori property that all subsets of a frequent itemset must also be frequent, we can determine that four latter candidates cannot possibly be frequent. How?  For example, lets take <math>\{I1, I2, I3\}.</math> The 2-item subsets of it are <math>\{I1, I2\}, \{I1, I3\} \&amp; \{I2, I3\}.</math> Since all 2-item subsets of <math>\{I1, I2, I3\}</math> are members of L2, We will keep <math>\{I1, I2, I3\}</math> in C3.  Lets take another example of <math>\{I2, I3, I5\}</math> which shows how the pruning is performed. The 2-item subsets are <math>\{I2, I3\}, \{I2, I5\} \&amp; \{I3, I5\}.</math>  BUT, <math>\{I3, I5\}</math> is not a member of L2 and hence it is not frequent violating Apriori Property. Thus We will have to remove <math>\{I2, I3, I5\}</math> from C3.  Therefore, <math>C3 = \{ \{I1, I2, I3\}, \{I1, I2, I5\} \}</math> after checking for all members of the result of Join operation for Pruning.    Now, the transactions in D are scanned in order to determine L3, consisting of those candidates 3-itemsets in C3 having minimum support.  Step 4: Generating 4-Itemset Frequent Pattern  The algorithm uses L3 Join L3 to generate a candidate set of 4-itemsets, C4. Although the join results in <math>\{ \{I1, I2, I3, I5\} \},</math> this itemset is pruned since its subset <math>\{ \{I2, I3, I5\} \}</math> is not frequent.  Thus, <math>C4 = \phi(\text{Null}),</math> and algorithm terminates, having found all </p>					
--	---	--	--	--	--	--

	<p>of the frequent items. This completes our Apriori Algorithm.</p> <p>Step 5: Generating Association Rules From Frequent Itemsets</p> <p>Let the minimum confidence threshold is, say 70%.</p> <p>The resulting association rules are shown below, each listed with its confidence.</p> <p>R1: <math>I1 \wedge I2 \rightarrow I5</math></p> <p>Confidence = <math>sc\{I1, I2, I5\} / sc\{I1, I2\} = 2/4 = 50\%</math>.</p> <p>R1 is Rejected.</p> <p>R2: <math>I1 \wedge I5 \rightarrow I2</math></p> <p>Confidence = <math>sc\{I1, I2, I5\} / sc\{I1, I5\} = 2/2 = 100\%</math>.</p> <p>R2 is Selected.</p> <p>R3: <math>I2 \wedge I5 \rightarrow I1</math></p> <p>Confidence = <math>sc\{I1, I2, I5\} / sc\{I2, I5\} = 2/2 = 100\%</math>.</p> <p>R3 is Selected.</p> <p>R4: <math>I1 \rightarrow I2 \wedge I5</math></p> <p>Confidence = <math>sc\{I1, I2, I5\} / sc\{I1\} = 2/6 = 33\%</math>.</p> <p>R4 is Rejected.</p> <p>R5: <math>I2 \rightarrow I1 \wedge I5</math></p> <p>Confidence = <math>sc\{I1, I2, I5\} / \{I2\} = 2/7 = 29\%</math></p> <p>R5 is Rejected.</p> <p>R6: <math>I5 \rightarrow I1 \wedge I2</math></p> <p>Confidence = <math>sc\{I1, I2, I5\} / \{I5\} = 2/2 = 100\%</math>.</p> <p>R6 is Selected.</p>					
5	<p><b>Explain Naive Baye's Classification.</b></p> <p>The naïve Bayesian classifier, or simple Bayesian classifier, works as follows:</p> <ol style="list-style-type: none"> <li>Let <math>D</math> be a training set of tuples and their associated class labels. As usual, each tuple is represented by an <math>n</math>-dimensional attribute vector, <math>X = (x_1, x_2, \dots, x_n)</math>, depicting <math>n</math> measurements made on the tuple from <math>n</math> attributes, respectively, <math>A_1, A_2, \dots, A_n</math>.</li> <li>Suppose that there are <math>m</math> classes, <math>C_1, C_2, \dots, C_m</math>. Given a tuple, <math>X</math>, the classifier will predict that <math>X</math> belongs to the class having the highest posterior probability, conditioned on <math>X</math>. That is, the naïve Bayesian classifier predicts that tuple <math>X</math> belongs to the class <math>C_i</math> if and only if <math display="block">P(C_i X) &gt; P(C_j X) \quad \text{for } 1 \leq j \leq m, j \neq i.</math> <p>Thus we maximize <math>P(C_i X)</math>. The class <math>C_i</math> for which <math>P(C_i X)</math> is maximized is called the <i>maximum posteriori hypothesis</i>. By Bayes' theorem (Equation (6.10)),</p> <math display="block">P(C_i X) = \frac{P(X C_i)P(C_i)}{P(X)}. \quad (6.11)</math> </li> <li>As <math>P(X)</math> is constant for all classes, only <math>P(X C_i)P(C_i)</math> need be maximized. If the class prior probabilities are not known, then it is commonly assumed that the classes are equally likely, that is, <math>P(C_1) = P(C_2) = \dots = P(C_m)</math>, and we would therefore maximize <math>P(X C_i)</math>. Otherwise, we maximize <math>P(X C_i)P(C_i)</math>. Note that the class prior probabilities may be estimated by <math>P(C_i) =  C_{i,D} / D </math>, where <math> C_{i,D} </math> is the number of training tuples of class <math>C_i</math> in <math>D</math>.</li> <li>In order to reduce computation in evaluating <math>P(X C_i)</math>, the naïve assumption of class conditional independence is made.</li> </ol> $  \begin{aligned}  P(X C_i) &= \prod_{k=1}^n P(x_k C_i) \\  &= P(x_1 C_i) \times P(x_2 C_i) \times \dots \times P(x_n C_i).  \end{aligned}  $ <ol style="list-style-type: none"> <li>If <math>A_k</math> is categorical, then <math>P(x_k C_i)</math> is the number of tuples of class <math>C_i</math> in <math>D</math> having the value <math>x_k</math> for <math>A_k</math>, divided by <math> C_{i,D} </math>, the number of tuples of class <math>C_i</math> in <math>D</math>.</li> <li>If <math>A_k</math> is continuous-valued, then we need to do a bit more work, but the calculation is pretty straightforward. A continuous-valued attribute is typically assumed to have a Gaussian distribution with a mean <math>\mu</math> and standard deviation <math>\sigma</math>, defined by <math display="block">g(x, \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, \quad (6.13)</math> <p>so that</p> <math display="block">P(x_k C_i) = g(x_k, \mu_{C_{i,A_k}}, \sigma_{C_{i,A_k}}). \quad (6.14)</math> </li> </ol>	5	L2	3	2	2.6.3

	<p>5. In order to predict the class label of <math>X</math>, <math>P(X C_i)P(C_i)</math> is evaluated for each class <math>C_i</math>. The classifier predicts that the class label of tuple <math>X</math> is the class <math>C_i</math> if and only if</p> $P(X C_i)P(C_i) > P(X C_j)P(C_j) \quad \text{for } 1 \leq j \leq m, j \neq i. \quad (6.15)$ <p>In other words, the predicted class label is the class <math>C_i</math> for which <math>P(X C_i)P(C_i)</math> is the maximum.</p>																																																																																																																																												
<p style="text-align: center;"><b>Part – C</b> <b>(4 x 10 = 10 Marks)</b></p>																																																																																																																																													
1	<p><b>A database has five transactions. Let the minimum support &amp; confidence, min sup=2, min confi=80%</b></p> <table><thead><tr><th>TID</th><th>ITEMS</th></tr></thead><tbody><tr><td>T1</td><td>{1,2,3,4,5,6}</td></tr><tr><td>T2</td><td>{7,2,3,4,5,6}</td></tr><tr><td>T3</td><td>{1,8,4,5}</td></tr><tr><td>T4</td><td>{1,9,0,4,6}</td></tr><tr><td>T5</td><td>{0,2,2,4,5}</td></tr></tbody></table> <p>Find the frequent itemsets and generate the association rules using Apriori algorithm</p> <div><div><table><thead><tr><th>TID</th><th>Items</th></tr></thead><tbody><tr><td>T1</td><td>{1,2,3,4,5,6}</td></tr><tr><td>T2</td><td>{7,2,3,4,5,6}</td></tr><tr><td>T3</td><td>{1,8,4,5}</td></tr><tr><td>T4</td><td>{1,9,0,4,6}</td></tr><tr><td>T5</td><td>{0,2,2,4,5}</td></tr></tbody></table><p>support = <math>60\% = \frac{60}{100} \times 5 = 3</math> confidence = <math>80\%</math></p></div><div><p>1. Finding Frequent Itemset</p><table><thead><tr><th>L1</th><th>Item</th><th>Count</th></tr></thead><tbody><tr><td>1</td><td>3</td><td></td></tr><tr><td>2</td><td>3</td><td></td></tr><tr><td>4</td><td>5</td><td></td></tr><tr><td>5</td><td>4</td><td></td></tr><tr><td>6</td><td>3</td><td></td></tr></tbody></table><p>→</p><table><thead><tr><th>L2</th><th>Item</th><th>Count</th></tr></thead><tbody><tr><td>1,4</td><td>3</td><td></td></tr><tr><td>2,4</td><td>3</td><td></td></tr><tr><td>2,5</td><td>3</td><td></td></tr><tr><td>4,5</td><td>4</td><td></td></tr><tr><td>4,6</td><td>3</td><td></td></tr></tbody></table><p>→</p><table><thead><tr><th>Candidate 1 - Itemset</th><th>Item</th><th>Count</th></tr></thead><tbody><tr><td>1</td><td>3</td><td></td></tr><tr><td>2</td><td>3</td><td></td></tr><tr><td>3</td><td>2</td><td></td></tr><tr><td>4</td><td>5</td><td></td></tr><tr><td>5</td><td>4</td><td></td></tr><tr><td>6</td><td>3</td><td></td></tr><tr><td>7</td><td>1</td><td></td></tr><tr><td>8</td><td>1</td><td></td></tr><tr><td>9</td><td>1</td><td></td></tr><tr><td>0</td><td>2</td><td></td></tr></tbody></table><p>→</p><table><thead><tr><th>Candidate 2 - Itemset</th><th>Itemset</th><th>Count</th></tr></thead><tbody><tr><td>1,2</td><td>1</td><td></td></tr><tr><td>1,4</td><td>3</td><td></td></tr><tr><td>1,5</td><td>2</td><td></td></tr><tr><td>1,6</td><td>2</td><td></td></tr><tr><td>2,4</td><td>3</td><td></td></tr><tr><td>2,5</td><td>3</td><td></td></tr><tr><td>2,6</td><td>2</td><td></td></tr><tr><td>4,5</td><td>4</td><td></td></tr><tr><td>4,6</td><td>3</td><td></td></tr><tr><td>5,6</td><td>2</td><td></td></tr></tbody></table><p>→</p><table><thead><tr><th>Candidate 3 - Itemset</th><th>Item</th><th>Count</th></tr></thead><tbody><tr><td>2,4,5</td><td>3</td><td></td></tr><tr><td>4,5,6</td><td>2</td><td></td></tr></tbody></table><p>Frequent Itemset {2,4,5}</p></div></div> <p>2. Generate Association Rule using Min_conf = 80%.</p> <p>The frequent Itemset - {2,4,5}</p> <p>write possible subset of frequent itemset</p> <p>{2,4}, {4,5}, {2,5}, {2}, {4}, {5}, {}, {2,4,5} → Eliminate the empty subset</p> <p>R1: 2 ∧ 4 → 5 (Accepted)</p> $\text{conf} = \frac{\text{sup}(A \cup B)}{\text{sup}(A)} = \frac{\text{sup}(2,4,5)}{\text{sup}(2,4)} = \frac{3}{3} = 1 (100\%)$ <p>R2: 4 ∧ 5 → 2 (Eliminated)</p> $\frac{\text{sup}(4,5,2)}{\text{sup}(4,5)} = \frac{3}{4} = 0.75 (75\%)$ <p>R3: 2 ∧ 5 → 4 (Accepted)</p> $\frac{\text{sup}(2,5,4)}{\text{sup}(2,5)} = \frac{3}{3} = 1 (100\%)$ <p>R4: 2 → 4 ∧ 5 (Accepted)</p> $\frac{\text{sup}(2,4,5)}{\text{sup}(2)} = \frac{3}{3} = 1 (100\%)$ <p>R5: 4 → 2 ∧ 5 (Eliminated)</p> $\frac{\text{sup}(4,2,5)}{\text{sup}(4)} = \frac{3}{5} = 0.6 (60\%)$ <p>R6: 5 → 2 ∧ 4 (Eliminated)</p> $\frac{\text{sup}(5,2,4)}{\text{sup}(5)} = \frac{3}{4} = 0.75 (75\%)$ <p>Final Rules (That satisfies both supports &amp; confidence)</p> <ol style="list-style-type: none"><li>2 ∧ 4 → 5</li><li>2 ∧ 5 → 4</li><li>2 → 4 ∧ 5</li></ol>	TID	ITEMS	T1	{1,2,3,4,5,6}	T2	{7,2,3,4,5,6}	T3	{1,8,4,5}	T4	{1,9,0,4,6}	T5	{0,2,2,4,5}	TID	Items	T1	{1,2,3,4,5,6}	T2	{7,2,3,4,5,6}	T3	{1,8,4,5}	T4	{1,9,0,4,6}	T5	{0,2,2,4,5}	L1	Item	Count	1	3		2	3		4	5		5	4		6	3		L2	Item	Count	1,4	3		2,4	3		2,5	3		4,5	4		4,6	3		Candidate 1 - Itemset	Item	Count	1	3		2	3		3	2		4	5		5	4		6	3		7	1		8	1		9	1		0	2		Candidate 2 - Itemset	Itemset	Count	1,2	1		1,4	3		1,5	2		1,6	2		2,4	3		2,5	3		2,6	2		4,5	4		4,6	3		5,6	2		Candidate 3 - Itemset	Item	Count	2,4,5	3		4,5,6	2		10	L3	2	1	2.5.2
TID	ITEMS																																																																																																																																												
T1	{1,2,3,4,5,6}																																																																																																																																												
T2	{7,2,3,4,5,6}																																																																																																																																												
T3	{1,8,4,5}																																																																																																																																												
T4	{1,9,0,4,6}																																																																																																																																												
T5	{0,2,2,4,5}																																																																																																																																												
TID	Items																																																																																																																																												
T1	{1,2,3,4,5,6}																																																																																																																																												
T2	{7,2,3,4,5,6}																																																																																																																																												
T3	{1,8,4,5}																																																																																																																																												
T4	{1,9,0,4,6}																																																																																																																																												
T5	{0,2,2,4,5}																																																																																																																																												
L1	Item	Count																																																																																																																																											
1	3																																																																																																																																												
2	3																																																																																																																																												
4	5																																																																																																																																												
5	4																																																																																																																																												
6	3																																																																																																																																												
L2	Item	Count																																																																																																																																											
1,4	3																																																																																																																																												
2,4	3																																																																																																																																												
2,5	3																																																																																																																																												
4,5	4																																																																																																																																												
4,6	3																																																																																																																																												
Candidate 1 - Itemset	Item	Count																																																																																																																																											
1	3																																																																																																																																												
2	3																																																																																																																																												
3	2																																																																																																																																												
4	5																																																																																																																																												
5	4																																																																																																																																												
6	3																																																																																																																																												
7	1																																																																																																																																												
8	1																																																																																																																																												
9	1																																																																																																																																												
0	2																																																																																																																																												
Candidate 2 - Itemset	Itemset	Count																																																																																																																																											
1,2	1																																																																																																																																												
1,4	3																																																																																																																																												
1,5	2																																																																																																																																												
1,6	2																																																																																																																																												
2,4	3																																																																																																																																												
2,5	3																																																																																																																																												
2,6	2																																																																																																																																												
4,5	4																																																																																																																																												
4,6	3																																																																																																																																												
5,6	2																																																																																																																																												
Candidate 3 - Itemset	Item	Count																																																																																																																																											
2,4,5	3																																																																																																																																												
4,5,6	2																																																																																																																																												
2	<p><b>Write the algorithm steps to generate decision tree from the training tuples of the data partition.</b></p>	10	L3	3	2	2.6.2																																																																																																																																							



	<p><b>Algorithm: Generate_decision_tree.</b> Generate a decision tree from the training tuples of data partition <math>D</math>.</p> <p><b>Input:</b></p> <ul style="list-style-type: none"> <li>■ Data partition, <math>D</math>, which is a set of training tuples and their associated class labels;</li> <li>■ <i>attribute_list</i>, the set of candidate attributes;</li> <li>■ <i>Attribute_selection_method</i>, a procedure to determine the splitting criterion that “best” partitions the data tuples into individual classes. This criterion consists of a <i>splitting_attribute</i> and, possibly, either a <i>split point</i> or <i>splitting subset</i>.</li> </ul> <p><b>Output:</b> A decision tree.</p> <p><b>Method:</b></p> <ol style="list-style-type: none"> <li>(1) create a node <math>N</math>;</li> <li>(2) if tuples in <math>D</math> are all of the same class, <math>C</math> then</li> <li>(3)     return <math>N</math> as a leaf node labeled with the class <math>C</math>;</li> <li>(4) if <i>attribute_list</i> is empty then</li> <li>(5)     return <math>N</math> as a leaf node labeled with the majority class in <math>D</math>; // majority voting</li> <li>(6) apply <i>Attribute_selection_method</i>(<math>D</math>, <i>attribute_list</i>) to find the “best” <i>splitting_criterion</i>;</li> <li>(7) label node <math>N</math> with <i>splitting_criterion</i>;</li> <li>(8) if <i>splitting_attribute</i> is discrete-valued and       multiway splits allowed then // not restricted to binary trees</li> <li>(9)     <i>attribute_list</i> ← <i>attribute_list</i> – <i>splitting_attribute</i>; // remove <i>splitting_attribute</i></li> <li>(10) for each outcome <math>j</math> of <i>splitting_criterion</i>       // partition the tuples and grow subtrees for each partition</li> <li>(11)     let <math>D_j</math> be the set of data tuples in <math>D</math> satisfying outcome <math>j</math>; // a partition</li> <li>(12)     if <math>D_j</math> is empty then</li> <li>(13)         attach a leaf labeled with the majority class in <math>D</math> to node <math>N</math>;</li> <li>(14)     else attach the node returned by <i>Generate_decision_tree</i>(<math>D_j</math>, <i>attribute_list</i>) to node <math>N</math>;</li> <li>(15)     endfor</li> <li>(15) return <math>N</math>;</li> </ol> <hr/> <p>Basic algorithm for inducing a decision tree from training tuples.</p>					
3	<p><b>a) What are the advantages of FP-Growth algorithm?</b></p> <p><b>b) Discuss the applications of association analysis.</b></p> <p><b>Ans a) Advantages of FP Growth Algorithm</b></p> <ul style="list-style-type: none"> <li>○ This algorithm needs to scan the database twice when compared to Apriori, which scans the transactions for each iteration.</li> <li>○ The pairing of items is not done in this algorithm, making it faster.</li> <li>○ The database is stored in a compact version in memory.</li> <li>○ It is efficient and scalable for mining both long and short frequent patterns.</li> <li>○ No candidate generation, no candidate test</li> <li>○ Uses compact data structure called FP-Tree</li> <li>○ Eliminates repeated database scan</li> <li>○ Basic operation is counting and FP-tree building</li> </ul> <p><b>b) Applications of association analysis</b></p> <p>The association rule learning is the important technique of machine learning, and it is employed in Market Basket analysis, Web usage mining, continuous production, etc. In market basket analysis, it is an adequate used by several big retailers to find the relations among items.</p> <p>Association rules were originally transformed from point-of-sale data that represent what products are purchased together. Although its roots are in linking point-of-sale transactions, association rules can be used external the retail market to find relationships among types of “baskets.”</p> <p>There are various applications of Association Rule which are as follows –</p> <ul style="list-style-type: none"> <li>• Items purchased on a credit card, such as rental cars and hotel rooms, support insight into the following product that customer are likely to buy.</li> <li>• Optional services purchased by tele-connection users (call waiting, call forwarding, DSL, speed call, etc.) support decide how to bundle these functions to maximize revenue.</li> <li>• Banking services used by retail users (money industry accounts, CDs, investment services, car loans, etc.)</li> </ul>	10	L1	2	1	1.7.1

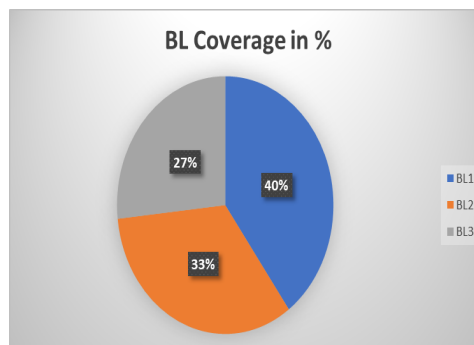
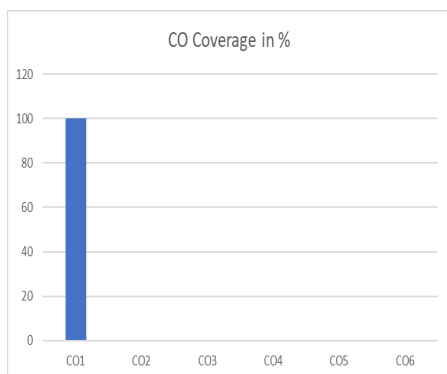


	<p>recognize users likely to needed other services.</p> <ul style="list-style-type: none"> <li>Unusual group of insurance claims can be an expression of fraud and can spark higher investigation.</li> <li>Medical patient histories can supports expressions of likely complications based on definite set of treatments.</li> </ul>					
4	<p><b>Illustrate any two of the attribute selection measure with an example.</b></p> <p><b>Attribute Selection Measures</b></p> <ul style="list-style-type: none"> <li>An attribute selection measure is a heuristic for selecting the splitting criterion that “best” separates a given data partition, <math>D</math>, of class-labeled training tuples into individual classes.</li> <li>If we were to split <math>D</math> into smaller partitions according to the outcomes of the splitting criterion, ideally each partition would be pure (i.e., all the tuples that fall into a given partition would belong to the same class).</li> <li>Attribute selection measures are also known as splitting rules because they determine how the tuples at a given node are to be split.</li> <li>Three popular attribute selection measures <ul style="list-style-type: none"> <li>information gain</li> <li>gain ratio, and</li> <li>gini index.</li> </ul> </li> </ul> <p><b>Information gain</b></p> <ul style="list-style-type: none"> <li>ID3 uses information gain as its attribute selection measure.</li> <li>The attribute with the highest information gain is chosen as the splitting attribute for node <math>N</math>.</li> <li>Where <math>p_i</math> is the probability that an arbitrary tuple in <math>D</math> belongs to class <math>C_i</math> and is estimated by <math> C_i, D / D </math>.</li> <li>A log function to the base 2 is used, because the information is encoded in bits.</li> <li><math>Info(D)</math> is just the average amount of information needed to identify the class label of a tuple in <math>D</math>.</li> <li><math>Info(D)</math> is also known as the entropy of <math>D</math>.</li> <li>The expected information needed to classify a tuple in <math>D</math> is given by <math display="block">Info(D) = - \sum_{i=1}^m p_i \log_2(p_i),</math> </li> <li>How much more information would we still need (after the partitioning) in order to arrive at an exact classification? This amount is measured by <math display="block">Info_A(D) = \sum_{j=1}^v \frac{ D_j }{ D } \times Info(D_j).</math> </li> <li>The term <math> D_j / D </math> acts as the weight of the <math>j</math>th partition.</li> <li><math>Info_A(D)</math> is the expected information required to classify a tuple from <math>D</math> based on the partitioning by <math>A</math>.</li> <li>The smaller the expected information (still) required, the greater the purity of the partitions.</li> <li>Information gain is defined as the difference between the original information requirement (i.e., based on just the proportion of classes) and the new requirement (i.e., obtained after partitioning on <math>A</math>). <math display="block">Gain(A) = Info(D) - Info_A(D).</math> </li> <li>The class label attribute, <i>buys computer</i>, has two distinct values (namely, {yes, no}).</li> <li>There are two distinct classes (that is, <math>m = 2</math>).</li> <li>Let class <math>C_1</math> correspond to <i>yes</i> and class <math>C_2</math> correspond to <i>no</i>.</li> <li>There are nine tuples of class <i>yes</i> and five tuples of class <i>no</i>.</li> </ul> <p>A (root) node <math>N</math> is created for the tuples in <math>D</math>.</p>	10	L2	3	2	2.5.2

	<p> <math display="block">Info(D) = -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right) = 0.940 \text{ bits.}</math> </p> <ul style="list-style-type: none"> <li>• Compute the expected information requirement for each attribute.</li> <li>• Age category youth – 2 yes &amp; 3 no, Middle aged – 4 yes &amp; 0 no, Senior – 3 yes &amp; 2 no.</li> </ul> <p> <math display="block">Info_{age}(D) = \frac{5}{14} \times \left(-\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5}\right) + \frac{4}{14} \times \left(-\frac{4}{4} \log_2 \frac{4}{4} - \frac{0}{4} \log_2 \frac{0}{4}\right) + \frac{5}{14} \times \left(-\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5}\right) = 0.694 \text{ bits.}</math> </p> <p> <math display="block">Gain(age) = Info(D) - Info_{age}(D) = 0.940 - 0.694 = 0.246 \text{ bits.}</math> </p> <ul style="list-style-type: none"> <li>• Compute <math>Gain(income) = 0.029</math> bits, <math>Gain(student) = 0.151</math> bits, and <math>Gain(credit \text{ rating}) = 0.048</math> bits.</li> <li>• Age has the highest information gain among the attributes, it is selected as the splitting attribute.</li> </ul> <p><b>Gain ratio</b></p> <ul style="list-style-type: none"> <li>• It prefers to select attributes having a large number of values.</li> <li>• C4.5, a successor of ID3, uses an extension to information gain known as <i>gain ratio</i>.</li> <li>• It applies a kind of normalization to information gain using a “split information” value defined analogously with <math>Info(D)</math> as</li> </ul> <p> <math display="block">SplitInfo_A(D) = -\sum_{j=1}^v \frac{ D_j }{ D } \times \log_2\left(\frac{ D_j }{ D }\right).</math> </p> <ul style="list-style-type: none"> <li>• It differs from information gain, which measures the information with respect to classification that is acquired based on the same partitioning.</li> <li>• The gain ratio is defined as</li> </ul> <p> <math display="block">GainRatio(A) = \frac{Gain(A)}{SplitInfo(A)}.</math> </p> <p> <math display="block">SplitInfo_A(D) = -\frac{4}{14} \times \log_2\left(\frac{4}{14}\right) - \frac{6}{14} \times \log_2\left(\frac{6}{14}\right) - \frac{4}{14} \times \log_2\left(\frac{4}{14}\right).</math> </p> <p> <math display="block">= 0.926.</math> </p> <p> <math display="block">Gain(income) = 0.029,</math> </p> <p> <math display="block">GainRatio(income) = 0.029/0.926 = 0.031.</math> </p>					
--	--	--	--	--	--	--

**\*Program Indicators are available separately for Computer Science and Engineering in AICTE examination reforms policy.**

#### **Course Outcome (CO) and Bloom's level (BL) Coverage in Questions**



**Approved by the Audit Professor/Course Coordinator**