

Nirma University
Institute of Technology
Computer Engineering Department
Data Mining – 2CSDE71

B. Tech Semester – VI

Academic Year: Even 2023-2024

List of Practical

Laboratory work will be based on the above syllabus with minimum 10 experiments to be incorporated.

Sr. No	List of Experiments	No. of Hours	Mappe d CLO																				
1	Data Domain selection and Identification of Characteristics of selected Dataset of different Formats. Also write a report with following detail a) Selection of data domain b) Define the data domain c) The data source d) Objective e) Define the selection of fields f) Characteristic and behaviors (distribution and inference) of data for each selected field	2	CLO1																				
2	Calculate the dissimilarity value of the following table containing various types of attributes or on the dataset identified in Practical 1. <table><tr><td></td><td>Nominal</td><td>Ordinal</td><td>Numerical</td></tr><tr><td>1</td><td>CodeA</td><td>Excellent</td><td>45</td></tr><tr><td>2</td><td>CodeB</td><td>Fair</td><td>22</td></tr><tr><td>3</td><td>CodeC</td><td>Good</td><td>64</td></tr><tr><td>4</td><td>CodeD</td><td>Excellent</td><td>28</td></tr></table>		Nominal	Ordinal	Numerical	1	CodeA	Excellent	45	2	CodeB	Fair	22	3	CodeC	Good	64	4	CodeD	Excellent	28	2	CLO1
	Nominal	Ordinal	Numerical																				
1	CodeA	Excellent	45																				
2	CodeB	Fair	22																				
3	CodeC	Good	64																				
4	CodeD	Excellent	28																				
3	Identify and implement any three methods to fill the missing values indicated by'?' in the given data set or on the dataset identified in Practical 1. (Analysis of the best method with reasoning) <table><tr><th>Name</th><th>Value</th></tr><tr><td>A</td><td>45</td></tr><tr><td>B</td><td>37</td></tr><tr><td>C</td><td>59</td></tr><tr><td>D</td><td>?</td></tr><tr><td>E</td><td>47</td></tr><tr><td>F</td><td>39</td></tr><tr><td>G</td><td>?</td></tr><tr><td>H</td><td>43</td></tr><tr><td>I</td><td>52</td></tr></table>	Name	Value	A	45	B	37	C	59	D	?	E	47	F	39	G	?	H	43	I	52	2	CLO1
Name	Value																						
A	45																						
B	37																						
C	59																						
D	?																						
E	47																						
F	39																						
G	?																						
H	43																						
I	52																						

	J ?																														
4	<p>Implement the menu driven program for normalization for given dataset from user (by all means) for the following normalization techniques. for example, dataset (age):</p> <p>13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 25, 25, 30, 33, 33, 35, 35, 35, 35, 36, 40, 45, 46, 52, 70.</p> <ul style="list-style-type: none">• Use Min-Max Normalization to transform the value 25 and 52 for age onto the range [0.0, 1.0].• Use z-score normalization to transform the value 35 for <i>age</i>, where the standard deviation of <i>age</i> is 12.94 years.• Use normalization by decimal scaling to transform the value 35 for <i>age</i>.	2	CLO1																												
5	<p>Implement Principal component analysis Dimensionality reduction technique on the following dataset (Do not use any available library of PCA)</p> <table><thead><tr><th>Diastolic BP</th><th>Systolic BP</th><th>Weight</th><th>Height</th></tr></thead><tbody><tr><td>78</td><td>126</td><td>67</td><td>170</td></tr><tr><td>80</td><td>128</td><td>77</td><td>177</td></tr><tr><td>81</td><td>127</td><td>89</td><td>183</td></tr><tr><td>82</td><td>130</td><td>90</td><td>187</td></tr><tr><td>84</td><td>130</td><td>50</td><td>165</td></tr><tr><td>86</td><td>132</td><td>55</td><td>164</td></tr></tbody></table> <p>BP BS</p>	Diastolic BP	Systolic BP	Weight	Height	78	126	67	170	80	128	77	177	81	127	89	183	82	130	90	187	84	130	50	165	86	132	55	164	2	CLO2
Diastolic BP	Systolic BP	Weight	Height																												
78	126	67	170																												
80	128	77	177																												
81	127	89	183																												
82	130	90	187																												
84	130	50	165																												
86	132	55	164																												
6	<p>Identify the frequent patterns and generate strong association rule from the frequent pattern for the following data set (using Apriori and any improved version of Apriori). Keep minimum 50% support and 40% confidence.</p> <table><thead><tr><th>Tid</th><th>Items bought</th></tr></thead><tbody><tr><td>10</td><td>Beer, Nuts, Diaper</td></tr><tr><td>20</td><td>Beer, Coffee, Diaper</td></tr><tr><td>30</td><td>Beer, Diaper, Eggs</td></tr><tr><td>40</td><td>Nuts, Eggs, Milk</td></tr><tr><td>50</td><td>Nuts, Coffee, Diaper, Eggs, Milk</td></tr></tbody></table> <p>* Prepare the analytics report on interesting patterns and correlation analysis of generated output.</p>	Tid	Items bought	10	Beer, Nuts, Diaper	20	Beer, Coffee, Diaper	30	Beer, Diaper, Eggs	40	Nuts, Eggs, Milk	50	Nuts, Coffee, Diaper, Eggs, Milk	2	CLO2																
Tid	Items bought																														
10	Beer, Nuts, Diaper																														
20	Beer, Coffee, Diaper																														
30	Beer, Diaper, Eggs																														
40	Nuts, Eggs, Milk																														
50	Nuts, Coffee, Diaper, Eggs, Milk																														

7	<p>The DBLP dataset (www.informatik.uni-trier.de/ley/db/) consists of over one million entries of research papers published in computer science conferences and journals. Among these entries, there are a good number of authors that have coauthor relationships. (a) Propose a method and implement it to efficiently mine a set of coauthor relationships that are closely correlated (e.g., often coauthoring papers together).</p>	4	CLO2
8	<p>Implement C-Means clustering algorithm cluster the following four points (with (x; y) representing location) into two clusters.</p> <p>A1(1; 2), A2(2; 3), A3(9; 4), A4(10,1)</p> <p>B1(5; 8); B2(7; 5); B3(6; 4);</p> <p>C1(4; 2); C2(4; 9)</p> <p>The distance function is Euclidean distance. Suppose initially we assign A1, B1, and C1 as the centre of each cluster, respectively. Use the k-means algorithm to show the three cluster centres and all the points of clusters after the 2nd round of execution.</p>	4	CLO2
9	<p>Implement Random Forest Algorithm on Car Evaluation Database. It contains examples with the structural information removed, i.e., it directly relates CAR to the six input attributes: buying, maint, doors, persons, lug_boot, safety. Basically, we have to build a classifier to classify a car as 'Unacceptable', 'Acceptable', 'Good' and 'Very Good' based on the attributes.</p> <p>The different attributes values are given as follows:</p> <ol style="list-style-type: none"> 1. buying: vhigh, high, med, low. 2. maint: vhigh, high, med, low. 3. doors: 2, 3, 4, 5, more. 4. persons: 2, 4, more. 5. lug_boot: small, med, big. 6. safety: low, med, high. <p>You can download the dataset from here : https://archive.ics.uci.edu/ml/datasets/Car+Evaluation</p>	4	CLO2 CLO3
10	<p>Explore Weka/Knime tool for KDD process (classification algorithms) in any of the following application.</p> <ol style="list-style-type: none"> 1. Healthcare 2. Government 3. Transportation 4. Education 	4	CLO2 CLO3

	5. E-commerce 6. Entertainment		
--	-----------------------------------	--	--

Suggested Readings^:

1. Jiawei Han and Micheline Kamber, Data mining: Concepts and Techniques, Morgan Kaufmann Publishers.
2. Ian H. Witten and Eibe Frank, Data Mining: Practical Machine Learning Tools and Techniques, Morgan Kaufmann
3. Hand, Mannila, and Smyth., Principles of Data Mining, MIT Press
4. Berry and Linoff, Mastering Data Mining, Wiley
5. Delmater and Hancock, Data Mining Explained, Digital Press

Lesson Plan

Sr. No.	Topics	CLO
1	Motivation and importance, different kinds of data, data mining functionalities,	1
2	classification of data mining systems,	1
3	major issues in data mining	1
4	Data summarization,	1
5	data cleaning,	1
6	data integration and transformation,	1
7	data reduction,	1
8	data discretization and concept hierarchy generation	1
9	Basic concept of Mining Frequent Patterns	2
10	efficient and scalable frequent itemset mining methods	2
11	efficient and scalable frequent itemset mining methods	2
12	Mining various kind of association rules	2
13	Mining various kind of association rules	2
14	from association mining to correlation analysis,	2
15	constraint-based association mining	2
16	Classification vs. prediction, Issues regarding classification and prediction	2
17	Classification by decision tree induction	2
18	Statistical-Based Algorithms, Distance-Based Algorithms	2
19	Decision Tree-Based Algorithms, Neural Network-Based Algorithms, Rule-Based Algorithms,	2
20	Combining Techniques, accuracy and error measures, evaluation of the accuracy of a classifier or predictor.	2
21	Types of data in cluster analysis	2,3
22	overview of major clustering methods,	2,3
23	Probabilistic model based clustering	2,3

24	Clustering high dimensional data,	2,3
25	Clustering Graph and Network data.	2,3
26	Applications of Distributed and parallel Data Mining	2,3
27	Advanced Techniques: Web Mining	2,3
28	Spatial Database Mining	2,3
29	Temporal Mining,	2,3
30	Multimedia Mining.	2,3