- 7. Display the data stored in a given graph using the Breadth-First Search algorithm.
- 8. Display the data stored in a given graph using the Depth-First Search algorithm.
- Write a program to determine a minimum spanning tree of a graph using the Prim's algorithm.
- 10. Write a program to implement Dijkstra's algorithm to find the shortest paths from a given source node to all other nodes in a graph.
- 11. Write a program to solve the weighted interval scheduling problem.
- 12. Write a program to solve the 0-1 knapsack problem.

For the algorithms at S.No 1, 2 and 3, test run the algorithm on 100 different input sizes varying from 30 to 1000. For each size find the number of comparisons averaged on 10 different input instances; plot a graph for the average number of comparisons against each input size. Compare it with a graph of nlogn.

DSC-A3/DSE: DATA MINING-I

CREDIT DISTRIBUTION, ELIGIBILITY AND PRE-REQUISITES OF THE COURSE

Course title &	Credits	Credit distribution of the course			Eligibility criteria	Pre-requisite of the course
Code		Lecture	Tutorial	Practical/ Practice		(if any)
Data Mining - I	4	3	0	1	Passed 12th class with Mathema tics	Programming using Python

Course Objectives

This course aims to introduce data mining techniques and their application on real-life datasets. The students will learn to pre-process the dataset and make it ready for application of data mining techniques. The course will focus on three main techniques of data mining i.e. Classification, Clustering and Association Rule Mining. Different algorithms for these techniques

will be discussed along with appropriate evaluation metrics to judge the performance of the results delivered.

Learning outcomes

On successful completion of the course, students will be able to:

- Pre-process the data for subsequent data mining tasks
- Apply a suitable classification algorithm to train the classifier and evaluate its performance.
- Apply appropriate clustering algorithm to cluster the data and evaluate clustering quality
- Use association rule mining algorithms and generate frequent item-sets and association rules

Syllabus

Unit 1 (8 hours)

Introduction to Data Mining:

Motivation and Challenges for data mining, Types of data mining tasks, Applications of data mining, Data measurements, Data quality, Supervised vs. unsupervised techniques

Unit 2 (9 hours)

Data Pre-Processing:

Data aggregation, sampling, dimensionality reduction, feature subset selection, feature creation, variable transformation.

Unit 3 (11 hours)

Cluster Analysis:

Basic concepts of clustering, measure of similarity, types of clusters and clustering methods, K-means algorithm, measures for cluster validation, determine optimal number of clusters

Unit 4 (8 hours)

Association Rule Mining:

Transaction data-set, frequent itemset, support measure, rule generation, confidence of association rule, Apriori algorithm, Apriori principle

Unit 5 (9 hours)

Classification:

Naive Bayes classifier, Nearest Neighbour classifier, decision tree, overfitting, confusion matrix, evaluation metrics and model evaluation.

Essential/recommended readings

- 1. Tan P.N., Steinbach M, Karpatne A. and Kumar V. Introduction to Data Mining,2nd edition, Pearson, 2021.
- 2. Han J., Kamber M. and Pei J. Data Mining: Concepts and Techniques, 3 rd edition, 2011, Morgan Kaufmann Publishers.
- 3. Zaki M. J. and Meira J. Jr. Data Mining and Machine Learning: Fundamental Concepts and Algorithms, 2nd edition, Cambridge University Press, 2020.

Additional References

- 1. Aggarwal C. C. Data Mining: The Textbook, Springer, 2015.
- 2. Dunham M. Data Mining: Introductory and Advanced Topics, 1st edition, Pearson Education India, 2006.

Recommended Datasets for:

Classification: Abalone, Artificial Characters, Breast Cancer Wisconsin (Diagnostic)

Clustering: Grammatical Facial Expressions, HTRU2, Perfume data

Association Rule Mining: MovieLens, Titanics

Practicals

- 1. Apply data cleaning techniques on any dataset (e,g, wine dataset). Techniques may include handling missing values, outliers, inconsistent values. A set of validation rules can be prepared based on the dataset and validations can be performed.
- 2. Apply data pre-processing techniques such as standardization/normalization, transformation, aggregation, discretization/binarization, sampling etc. on any dataset

- 3. Run Apriori algorithm to find frequent itemsets and association rules on 2 real datasets and use appropriate evaluation measures to compute correctness of obtained patterns a) Use minimum support as 50% and minimum confidence as 75% b) Use minimum support as 60% and minimum confidence as 60 % I.
- 4. Use Naive bayes, K-nearest, and Decision tree classification algorithms and build classifiers on any two datasets. Divide the data set into training and test set. Compare the accuracy of the different classifiers under the following situations: a) Training set = 75% Test set = 25% b) Training set = 66.6% (2/3rd of total), Test set = 33.3% II. Training set is chosen by i) hold out method ii) Random subsampling iii) Cross-Validation. Compare the accuracy of the classifiers obtained. Data is scaled to standard format.
- 5. Use Simple K-means algorithm for clustering on any dataset. Compare the performance of clusters by changing the parameters involved in the algorithm. Plot MSE computed after each iteration using a line plot for any set of parameters.

Project: Students should be promoted to take up one project on any UCI/kaggle/data.gov.in or a dataset verified by the teacher. Preprocessing steps and at least one data mining technique should be shown on the selected dataset. This will allow the students to have a practical knowledge of how to apply the various skills learnt in the subject for a single problem/project.

Note: Examination scheme and mode shall be as prescribed by the Examination Branch, University of Delhi, from time to time.

DSE-A4/DSE: DATA MINING-II

CREDIT DISTRIBUTION, ELIGIBILITY AND PRE-REQUISITES OF THE COURSE

Course	Credits	Credit di	stribution	of the course	· ·	Pre-requisite of
title & Code		Lecture	Tutorial	Practical/ Practice	criteria	the course (if any)