

### DISCIPLINE SPECIFIC ELECTIVE COURSE: Data Mining -

#### Credit distribution, Eligibility and Pre-requisites of the Course

Course title & Code	Credits	Credit distribution of the course			Eligibility criteria	Pre-requisite of the course (if any)
		Lecture	Tutorial	Practical/ Practice		
<b>Data Mining</b>	<b>4</b>	<b>3</b>	<b>0</b>	<b>1</b>	Pass in Class XII	Programming using Python

#### Learning 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:

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

#### SYLLABUS OF DSE

##### Unit 1 (7 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 (8 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, Distance-based method: K-means algorithm, measures for cluster validation, determine optimal number of clusters. Density-Based Method: DBSCAN Algorithm, Comparison of these two methods

**Unit 4 (8 hours)**

**Association Rule mining:** Transaction data-set, frequent itemset, support measure, rule generation, confidence of association rule, apriori principle, apriori algorithm

**Unit 5 (11 hours)**

**Classification:** Naive bayes classifier, nearest neighbour classifier, decision tree, overfitting, confusion matrix, evaluation metrics and model evaluation

**Text Book:**

1. Tan P.N., Steinbach M, Karpatne A. and Kumar V. Introduction to Data Mining, Second edition, Sixth Impression, Pearson, 2023.

**Additional References:**

1. Han J., Kamber M. and Pei J. *Data Mining: Concepts and Techniques*, 3<sup>rd</sup> edition, 2011, Morgan Kaufmann Publishers.
2. Zaki M. J. and Meira J. Jr. *Data Mining and Machine Learning: Fundamental Concepts and Algorithms*, 2<sup>nd</sup> edition, Cambridge University Press, 2020.
3. Aggarwal C. C. *Data Mining: The Textbook*, Springer, 2015

**Datasets may be downloaded from :**

1. <https://archive.ics.uci.edu/datasets>
2. <https://www.kaggle.com/datasets?fileType=csv>
3. <https://data.gov.in/>
4. <https://ieee-dataport.org/datasets>

**Suggested Practical Exercises**

1. Apply data cleaning techniques on any dataset (e.g., Paper Reviews dataset in UCI repository). Techniques may include handling missing values, outliers and 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 item sets 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 %
4. Use Naive bayes, K-nearest, and Decision tree classification algorithms to build classifiers on any two datasets. Pre-process the datasets using techniques specified in Q2. Compare the Accuracy, Precision, Recall and F1 measure reported for each dataset using the abovementioned classifiers under the following situations:
  - i. Using Holdout method (Random sampling):
    - a) Training set = 80% Test set = 20%
    - b) Training set = 66.6% (2/3rd of total), Test set = 33.3%

- ii. Using Cross-Validation:
  - a) 10-fold
  - b) 5-fold
- 5. Apply simple K-means algorithm for clustering any dataset. Compare the performance of clusters by varying the algorithm parameters. For a given set of parameters, plot a line graph depicting MSE obtained after each iteration.
- 6. Perform density-based clustering algorithm on a downloaded dataset and evaluate the cluster quality by changing the algorithm's parameters

**Project:** *Students should be promoted to take up one project on using dataset downloaded from any of the websites given above and the 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.*