## Syllabus

### Savitribai Phule Pune University
### Fourth Year of Computer Engineering (2015 Course)
### Elective I

## 410244(D) : Data Mining and Warehousing

| Teaching Scheme :<br>TH : 03 Hours/Week | Credit<br>03 | Examination Scheme :<br>In-Sem (Paper) : 30 Marks<br>End-Sem (Paper) : 70 Marks |
|---|---|---|

**Pre-requisites Courses**

310242-Database Management Systems, 310244 - Information Systems and Engineering Economics

**Companion Course :** 410247- Laboratory Practice II

**Course Objectives**

- To understand the fundamentals of Data Mining.
- To identify the appropriateness and need of mining the data.
- To learn the preprocessing, mining and post processing of the data.
- To understand various methods, techniques and algorithms in data mining.

**Course Outcomes**

On completion of the course the student should be able to :

- Apply basic, intermediate and advanced techniques to mine the data.
- Analyze the output generated by the process of data mining.
- Explore the hidden patterns in the data.
- Optimize the mining process by choosing best data mining technique.

## Course Contents

### Unit I : Introduction                                                                 (08 Hours)

Data Mining, Data Mining Task Primitives, Data : Data, Information and Knowledge; Attribute Types : Nominal, Binary, Ordinal and Numeric attributes, Discrete versus Continuous Attributes; Introduction to Data Preprocessing, Data Cleaning : Missing values, Noisy data; Data integration : Correlation analysis; transformation : Min-max normalization, z-score normalization and decimal scaling; data reduction : Data Cube Aggregation, Attribute Subset Selection, sampling; and Data Discretization : Binning, Histogram Analysis          **(Refer chapter 1)**

### Unit II : Data Warehouse                                                            (08 Hours)

Data Warehouse, Operational Database Systems and Data Warehouses (OLTP Vs OLAP), A Multidimensional Data Model: Data Cubes, Stars, Snowflakes, and Fact Constellations Schemas; OLAP Operations in the Multidimensional Data Model, Concept Hierarchies, Data Warehouse Architecture, The Process of Data Warehouse Design, A three-tier data warehousing architecture, Types of OLAP Servers : ROLAP versus MOLAP versus HOLAP.     **(Refer chapter 2)**

## Unit III : Measuring Data Similarity and Dissimilarity (08 Hours)

Measuring Data Similarity and Dissimilarity, Proximity Measures for Nominal Attributes and Binary Attributes, interval scaled; Dissimilarity of Numeric Data : Minskowski Distance, Euclidean distance and Manhattan distance; Proximity Measures for Categorical, Ordinal Attributes, Ratio scaled variables; Dissimilarity for Attributes of Mixed Types, Cosine Similarity. **(Refer chapter 3)**

## Unit IV : Association Rules Mining (08 Hours)

Market basket Analysis, Frequent item set, Closed item set, Association Rules, a-priori Algorithm, Generating Association Rules from Frequent Item sets, Improving the Efficiency of a-priori, Mining Frequent Item sets without Candidate Generation : FP Growth Algorithm; Mining Various Kinds of Association Rules : Mining multilevel association rules, constraint based association rule mining, Meta rule-Guided Mining of Association Rules. **(Refer chapter 4)**

## Unit V : Classification (08 Hours)

Introduction to : Classification and Regression for Predictive Analysis, Decision Tree Induction, Rule-Based Classification : using IF-THEN Rules for Classification, Rule Induction Using a Sequential Covering Algorithm. Bayesian Belief Networks, Training Bayesian Belief Networks, Classification Using Frequent Patterns, Associative Classification, Lazy Learners-k-Nearest-Neighbor Classifiers, Case-Based Reasoning. **(Refer chapter 5)**

## Unit VI : Multiclass Classification (08 Hours)

Multiclass Classification, Semi-Supervised Classification, Reinforcement learning, Systematic Learning, Wholistic learning and multi-perspective learning. Metrics for Evaluating Classifier Performance : Accuracy, Error Rate, precision, Recall, Sensitivity, Specificity; Evaluating the Accuracy of a Classifier : Holdout Method, Random Sub sampling and Cross-Validation. **(Refer chapter 6)**

❑❑❑

## UNIT II

UNIT V

**Chapter 5 :     Classification          5-1 to 5-34**

**Syllabus :**

Introduction to : Classification and Regression for Predictive Analysis, Decision Tree Induction, Rule-Based Classification : using IF-THEN Rules for Classification, Rule Induction Using a Sequential Covering Algorithm. Bayesian Belief Networks, Training Bayesian Belief Networks, Classification Using Frequent Patterns, Associative Classification, Lazy Learners-k-Nearest-Neighbor Classifiers, Case-Based Reasoning.

### UNIT VI

| Chapter 6 :    Multiclass Classification | 6-1 to 6-8 |
| --- | --- |

**Syllabus :**

Multiclass Classification, Semi-Supervised Classification, Reinforcement learning, Systematic, Learning, Wholistic learning and multi-perspective learning. Metrics for Evaluating Classifier Performance : Accuracy, Error Rate, precision, Recall, Sensitivity, Specificity; Evaluating the Accuracy of a Classifier : Holdout Method, Random Sub sampling and Cross-Validation.