

Course Code:	Course Title	Credit
CSC504	Data Warehousing and Mining	3

<b>Prerequisite: Database Concepts</b>	
<b>Course Objectives:</b>	
1.	To identify the significance of Data Warehousing and Mining.
2.	To analyze data, choose relevant models and algorithms for respective applications.
3.	To study web data mining.
4.	To develop research interest towards advances in data mining.
<b>Course Outcomes:</b> At the end of the course, the student will be able to	
1.	Understand data warehouse fundamentals and design data warehouse with dimensional modelling and apply OLAP operations.
2.	Understand data mining principles and perform Data preprocessing and Visualization.
3.	Identify appropriate data mining algorithms to solve real world problems.
4.	Compare and evaluate different data mining techniques like classification, prediction, clustering and association rule mining
5.	Describe complex information and social networks with respect to web mining.

Module	Content	Hrs
<b>1</b>	<b>Data Warehousing Fundamentals</b>	<b>8</b>
	Introduction to Data Warehouse, Data warehouse architecture, Data warehouse versus Data Marts, E-R Modeling versus Dimensional Modeling, Information Package Diagram, Data Warehouse Schemas; Star Schema, Snowflake Schema, Factless Fact Table, Fact Constellation Schema. Update to the dimension tables. Major steps in ETL process, OLTP versus OLAP, OLAP operations: Slice, Dice, Rollup, Drilldown and Pivot.	
<b>2</b>	<b>Introduction to Data Mining, Data Exploration and Data Pre-processing</b>	<b>8</b>
	Data Mining Task Primitives, Architecture, KDD process, Issues in Data Mining, Applications of Data Mining, Data Exploration: Types of Attributes, Statistical Description of Data, Data Visualization, Data Preprocessing: Descriptive data summarization, Cleaning, Integration & transformation, Data reduction, Data Discretization and Concept hierarchy generation.	
<b>3</b>	<b>Classification</b>	<b>6</b>
	Basic Concepts, Decision Tree Induction, Naïve Bayesian Classification, Accuracy and Error measures, Evaluating the Accuracy of a Classifier: Holdout & Random Subsampling, Cross Validation, Bootstrap.	
<b>4</b>	<b>Clustering</b>	<b>6</b>
	Types of data in Cluster analysis, Partitioning Methods ( <i>k</i> -Means, <i>k</i> -Medoids), Hierarchical Methods (Agglomerative, Divisive).	
<b>5</b>	<b>Mining frequent patterns and associations</b>	<b>6</b>
	Market Basket Analysis, Frequent Item sets, Closed Item sets, and Association Rule, Frequent Pattern Mining, Apriori Algorithm, Association Rule Generation, Improving the Efficiency of Apriori, Mining Frequent Itemsets without candidate generation, Introduction to Mining Multilevel Association Rules and Mining Multidimensional Association Rules.	

<b>6</b>	<b>Web Mining</b>	<b>5</b>
	Introduction, Web Content Mining: Crawlers, Harvest System, Virtual Web View, Personalization, Web Structure Mining: Page Rank, Clever, Web Usage Mining.	

#### **Textbooks:**

1	Paulraj Ponniah, “ <i>Data Warehousing: Fundamentals for IT Professionals</i> ”, Wiley India.
2	Han, Kamber, “ <i>Data Mining Concepts and Techniques</i> ”, Morgan Kaufmann 2 <sup>nd</sup> edition.
3	M.H. Dunham, “ <i>Data Mining Introductory and Advanced Topics</i> ”, Pearson Education.

#### **References:**

1	Reema Theraja, “ <i>Data warehousing</i> ”, Oxford University Press 2009.
2	Pang-Ning Tan, Michael Steinbach and Vipin Kumar, “ <i>Introduction to Data Mining</i> ”, Pearson Publisher 2 <sup>nd</sup> edition.
3	Ian H. Witten, Eibe Frank and Mark A. Hall, “ <i>Data Mining</i> ”, Morgan Kaufmann 3 <sup>rd</sup> edition.

#### **Assessment:**

##### **Internal Assessment:**

Assessment consists of two class tests of 20 marks each. The first-class test is to be conducted when approx. 40% syllabus is completed and second-class test when additional 40% syllabus is completed. Duration of each test shall be one hour.

##### **End Semester Theory Examination:**

1	Question paper will comprise of total six questions.
2	All question carries equal marks
3	Questions will be mixed in nature (for example, If Q.2 part (a) from module 3 then part (b) can be from any module other than module 3)
4	Only Four questions need to be solved.
5	In question paper weightage of each module will be proportional to the number of respective lecture hours as mentioned in the syllabus.

##### **Useful Links**

1	<a href="https://onlinecourses.nptel.ac.in/noc20_cs12/preview">https://onlinecourses.nptel.ac.in/noc20_cs12/preview</a>
2	<a href="https://www.coursera.org/specializations/data-mining">https://www.coursera.org/specializations/data-mining</a>