

CST466	DATA MINING	CATEGORY	L	T	P	CREDIT	YEAR OF INTRODUCTION
		PEC	2	1	0	3	2019

Preamble: This course helps the learner to understand the concepts of data mining and data warehousing. It covers the key processes of data mining, data preprocessing techniques, fundamentals and advanced concepts of classification, clustering, association rule mining, web mining and text mining. It enables the learners to develop new data mining algorithms and apply the existing algorithms in real-world scenarios.

Prerequisite: NIL

Course Outcomes: After the completion of the course the student will be able to

CO#	CO
CO1	Employ the key process of data mining and data warehousing concepts in application domains. (Cognitive Knowledge Level: Understand)
CO2	Make use of appropriate preprocessing techniques to convert raw data into suitable format for practical data mining tasks (Cognitive Knowledge Level: Apply)
CO3	Illustrate the use of classification and clustering algorithms in various application domains (Cognitive Knowledge Level: Apply)
CO4	Comprehend the use of association rule mining techniques. (Cognitive Knowledge Level: Apply)
CO5	Explain advanced data mining concepts and their applications in emerging domains (Cognitive Knowledge Level: Understand)

Mapping of course outcomes with program outcomes

	PO 1	PO2	PO3	PO4	PO5	PO6	PO7	PO8	PO9	PO10	PO11	PO12
CO1												
CO2												
CO3												

CO4												
CO5												

Abstract POs defined by National Board of Accreditation			
PO#	Broad PO	PO#	Broad PO
PO1	Engineering Knowledge	PO7	Environment and Sustainability
PO2	Problem Analysis	PO8	Ethics
PO3	Design/Development of solutions	PO9	Individual and team work
PO4	Conduct investigations of complex problems	PO10	Communication
PO5	Modern tool usage	PO11	Project Management and Finance
PO6	The Engineer and Society	PO12	Lifelong learning

Assessment Pattern

Bloom's Category	Continuous Assessment Tests		End Semester Examination Marks (%)
	Test 1 (%)	Test 2 (%)	
Remember	20	20	20
Understand	30	30	30
Apply	50	50	50
Analyze			
Evaluate			
Create			

Mark Distribution

Total Marks	CIE Marks	ESE Marks	ESE Duration
150	50	100	3

Continuous Internal Evaluation Pattern:

Attendance	10 marks
Continuous Assessment Test(Average of Internal Test1&2)	25 marks
Continuous Assessment Assignment	15 marks

Internal Examination Pattern

Each of the two internal examinations has to be conducted out of 50 marks. First series test shall be preferably conducted after completing the first half of the syllabus and the second series test shall be preferably conducted after completing the remaining part of the syllabus. There will be two parts: Part A and Part B. Part A contains 5 questions (preferably, 2 questions each from the completed modules and 1 question from the partly completed module), having 3 marks for each question adding up to 15 marks for part A. Students should answer all questions from Part A. Part B contains 7 questions (preferably, 3 questions each from the completed modules and 1 question from the partly completed module), each with 7 marks. Out of the seven questions, a student should answer any five.

End Semester Examination Pattern:

There will be two parts; Part A and Part B. Part A contains 10 questions with 2 questions from each module, having 3 marks for each question. Students should answer all questions. Part B contains 2 full questions from each module of which student should answer any one. Each question can have a maximum 2 subdivisions and carries 14 marks.

Syllabus**Module – 1 (Introduction to Data Mining and Data Warehousing)**

Data warehouse-Differences between Operational Database Systems and Data Warehouses, Multidimensional data model- Warehouse schema, OLAP Operations, Data Warehouse Architecture, Data Warehousing to Data Mining, Data Mining Concepts and Applications, Knowledge Discovery in Database Vs Data mining, Architecture of typical data mining system, Data Mining Functionalities, Data Mining Issues.

Module - 2 (Data Preprocessing)

Data Preprocessing-Need of data preprocessing, Data Cleaning- Missing values, Noisy data, Data Integration and Transformation, Data Reduction-Data cube aggregation, Attribute subset selection, Dimensionality reduction, Numerosity reduction, Discretization and concept hierarchy generation.

Module - 3 (Advanced classification and Cluster analysis)

Classification- Introduction, Decision tree construction principle, Splitting indices -Information Gain, Gini index Decision tree construction algorithms-ID3, Decision tree construction with presorting-SLIQ, Classification Accuracy-Precision, Recall.

Introduction to clustering-Clustering Paradigms, Partitioning Algorithm- PAM, Hierarchical Clustering-DBSCAN, Categorical Clustering-ROCK

Module 4: (Association Rule Analysis)

Association Rules-Introduction, Methods to discover Association rules, Apriori(Level-wise algorithm), Partition Algorithm, Pincer Search Algorithm, Dynamic Itemset Counting Algorithm, FP-tree Growth Algorithm.

Module 5 (Advanced Data Mining Techniques)

Web Mining - Web Content Mining, Web Structure Mining- Page Rank, Clever, Web Usage Mining- Preprocessing, Data structures, Pattern Discovery, Pattern Analysis. Text Mining-Text Data Analysis and information Retrieval, Basic measures for Text retrieval, Text Retrieval methods, Text Indexing Techniques, Query Processing Techniques.

Text Books

1. Dunham M H, "Data Mining: Introductory and Advanced Topics", Pearson Education, New Delhi, 2003.
2. Arun K Pujari, "Data Mining Techniques", Universities Press Private Limited, 2008.
3. Jaiwei Han and Micheline Kamber, "Data Mining Concepts and Techniques", Elsevier, 2006

Reference Books

1. M Sudeep Elayidom, "Data Mining and Warehousing", 1st Edition, 2015, Cengage Learning India Pvt. Ltd.
2. Mehmed Kantardzic, "Data Mining Concepts, Methods and Algorithms", John Wiley and Sons, USA, 2003.
3. Pang-Ning Tan and Michael Steinbach, "Introduction to Data Mining", Addison Wesley, 2006.

Course Level Assessment Questions

Course Outcome 1 (CO1):

- Explain the OLAP operations in a multidimensional model.
 - Compare the techniques used in ROLAP, MOLAP and HOLAP
- Explain the various data mining issues with respect to mining methodology, user interaction and diversity of data types.
- Suppose that a data warehouse consists of the three dimensions time, doctor, and patient, and the two measures count and charge, where charge is the fee that a doctor charges a patient for a visit.
 - Draw star and snowflake schema diagrams for the data warehouse.
 - Starting with the base cuboid [day; doctor; patient], what specific OLAP operations should be performed in order to list the total fee collected by each doctor in 2004?

Course Outcome 2 (CO2):

- Use the methods below to normalize the following group of data: 100, 200, 300, 400, 550, 600, 680, 850, 1000
 - min-max normalization by setting min = 0 and max = 1
 - z-score normalization
 - Normalization by decimal scaling

Comment on which method you would prefer to use for the given data, giving reasons as to why.
- Identify a suitable dataset from any available resources and apply different preprocessing steps that you have learned. Observe and analyze the output obtained. (Assignment)

Course Outcome 3 (CO3):

- Illustrate the working of ID3 algorithm with the following example

MOTOR	WHEELS	DOORS	SIZE	TYPE	CLASS
NO	2	0	small	cycle	bicycle
NO	3	0	small	cycle	tricycle
YES	2	0	small	cycle	motorcycle
YES	4	2	small	automobile	Sports car
YES	4	3	medium	automobile	minivan
YES	4	4	medium	automobile	sedan
YES	4	4	large	automobile	sumo

- Illustrate the working of K medoid algorithm for the given dataset. $A_1=(3,9)$, $A_2=(2,5)$, $A_3=(8,4)$, $A_4=(5,8)$, $A_5=(7,5)$, $A_6=(6,4)$, $A_7=(1,2)$, $A_8=(4,9)$.

- Take a suitable dataset from available resources and apply all the classification and clustering algorithms that you have studied on original and preprocessed datasets. Analyze the performance variation in terms of different quality metrics. Give a detailed report based on the analysis. (Assignment)

Course Outcome 4 (CO4):

- A database has five transactions. Let min sup = 60% and min con f = 80%.

<i>TID</i>	<i>items_bought</i>
T100	{M, O, N, K, E, Y}
T200	{D, O, N, K, E, Y}
T300	{M, A, K, E}
T400	{M, U, C, K, Y}
T500	{C, O, O, K, I, E}

- Find all frequent item sets using Apriori and FP-growth, respectively. Compare the efficiency of the two mining processes.
 - List all of the strong association rules (with support s and confidence c) matching the following metarule, where X is a variable representing customers, and $item_i$ denotes variables representing items (e.g., “A”, “B”, etc.)

$$\forall x \in transaction, buys(X, item_1) \wedge buys(X, item_2) \Rightarrow buys(X, item_3) \quad [s, c]$$
- Identify and list some scenarios in which association rule mining can be used, and then use at least two appropriate association rule mining techniques in one of the two scenarios. (Assignment)

Course Outcome 5 (CO5):

- Consider an e-mail database that stores a large number of electronic mail (e-mail) messages. It can be viewed as a semi structured database consisting mainly of text data. Discuss the following.
 - How can such an e-mail database be structured so as to facilitate multidimensional search, such as by sender, by receiver, by subject, and by time?
 - What can be mined from such an e-mail database?
 - Suppose you have roughly classified a set of your previous e-mail messages as junk, unimportant, normal, or important. Describe how a data mining system may take this as the training set to automatically classify new e-mail messages or unclassified ones.
- Precision and recall are two essential quality measures of an information retrieval system.
 - Explain why it is the usual practice to trade one measure for the other.
 - Explain why the F-score is a good measure for this purpose.

- (c) Illustrate the methods that may effectively improve the F-score in an information retrieval system.
3. Explain HITS algorithm with an example.

Model Question Paper**QP CODE:****Reg No:** _____**Name:** _____**PAGES : 4****APJ ABDUL KALAM TECHNOLOGICAL UNIVERSITY****EIGHTH SEMESTER B.TECH DEGREE EXAMINATION, MONTH & YEAR****Course Code: CST466****Course Name: Data Mining****Max.Marks:100****Duration: 3 Hours****PART A****Answer All Questions. Each Question Carries 3 Marks**

1. Differentiate between OLTP and OLAP.
2. Compare the techniques of ROLAP, MOLAP and HOLAP
3. Explain Concept hierarchy with an example.
4. Explain heuristic methods of attribute subset selection techniques.
5. Consider a two-class classification problem of predicting whether a photograph contains a man or a woman. Suppose we have a test dataset of 10 records with expected outcomes and a set of predictions from our classification algorithm.

	Expected	Predicted
1	man	woman
2	man	man
3	woman	woman
4	man	man
5	woman	man
6	woman	woman
7	woman	woman
8	man	man
9	man	woman
10	woman	woman

Calculate precision, recall of the data.

6. Given two objects represented by the tuples (22,1,42,10) and (20,0, 36,8). Compute the Euclidean and Manhattan distance between the two objects.
7. The pincer search algorithm is a bi-directional search, whereas the level wise algorithm is a unidirectional search. Express your opinion about the statement.
8. Define support, confidence and frequent set in association data mining context.
9. Distinguish between focused crawling and regular crawling.
10. Describe any two-text retrieval indexing techniques. (10x3=30)

Part B

(Answer any one question from each module. Each question carries 14 Marks)

11. (a) Suppose a data warehouse consists of three measures: customer, account and branch and two measures count (number of customers in the branch) and balance. Draw the schema diagram using snowflake schema and star schema. (7)
- (b) Explain three- tier data warehouse architecture with a neat diagram. (7)

OR

- 12 (a) Illustrate different OLAP operations in multidimensional data model (7)
- (b) Describe different issues in data mining (7)
- 13 (a) Suppose that the data for analysis includes the attribute age. The age values for the data tuples are (in increasing order) 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. (8)
- (a) Use min-max normalization to transform the value 35 for age onto

the
range [0-1].

- (b) Use z-score normalization to transform the value 35 for age, where the standard deviation of age is 12.94 years.
- (c) Use normalization by decimal scaling to transform the value 35 for age.
- (d) Use smoothing by bin means to smooth the above data, using a bin depth of 3. Illustrate your steps. Comment on the effect of this technique for the given data.

- (b) With proper illustration, explain how PCA can be used for dimensionality reduction? Explain (6)

OR

- 14 (a) Suppose a group of 12 sales price records has been sorted as follows: 5, 10, 11, 13, 15, 35, 50, 55, 72, 92, 204, 215. Sketch examples of each of the following sampling techniques: SRSWOR, SRSWR, cluster sampling, stratified sampling. Use samples of size 5 and the strata “youth,” “middle-aged,” and “senior.” (8)
- (b) Partition the above data into three bins by each of the following methods: (6)
- (i) equal-frequency (equi-depth) partitioning
 - (ii) equal-width partitioning
- 15 (a) Explain the concept of a cluster as used in ROCK. Illustrate with examples (9)
- (b) Consider the following dataset for a binary classification problem. (5)

A	B	Class Label
T	F	+
T	T	+
T	T	+
T	F	-
T	T	+
F	F	-
F	F	-
F	F	-
T	T	-
T	F	-

Calculate the gain in Gini index when splitting on A and B respectively. Which attribute would the decision tree induction algorithm choose?

OR

- 16 (a) For a sunburn dataset given below, find the first splitting attribute for the decision tree by using the ID3 algorithm. (10)

Name	Hair	Height	Weight	Lotion	Class
Sarah	Blonde	Average	Light	No	Sunburn
Dana	Blonde	Tall	Average	Yes	None
Alex	Brown	Tall	Average	Yes	None
Annie	Blonde	Short	Average	No	Sunburn
Emily	Red	Average	Heavy	No	Sunburn
Pete	Brown	Tall	Heavy	No	None
John	Brown	Average	Heavy	No	None
Katie	Blonde	Short	Light	Yes	None

- (b) Explain the working of SLIQ algorithm. (4)
- 17 (a) Illustrate the working of Pincer Search Algorithm with an example. (7)
- (b) Describe the working of dynamic itemset counting technique? Specify when to move an itemset from dashed structures to solid structures? (7)

OR

- 18 (a) A database has six transactions. Let min_sup be 60% and min_conf be 80%. (9)

TID	items_bought
T1	I1, I2, I3
T2	I2, I3, I4
T3	I4, I5
T4	I1, I2, I4
T5	I1, I2, I3, I5
T6	I1, I2, I3, I4

Find frequent itemsets using FP Growth algorithm and generate strong association rules from a three item dataset.

- (b) Write partitioning algorithm for finding large itemset and compare its efficiency with apriori algorithm (5)

- 19 (a) Describe web content mining techniques. (7)
- (b) Write an algorithm to find maximal frequent forward sequences to mine log traversal patterns. Illustrate the working of this algorithm. (7)

OR

- 20 (a) Explain how web structure mining is different from web usage mining and web content mining? Write a CLEVER algorithm for web structure mining. (7)
- (b) Describe different Text retrieval methods. Explain the relationship between text mining and information retrieval and information extraction. (7)

Teaching Plan

No	Contents	No. of lecture hours (36 Hrs)
Module 1(Introduction to Data Mining and Data Warehousing) (Text3) (6 hours)		
1.1	Data warehouse-Differences between Operational Database Systems and Data Warehouses, Multidimensional data model- Warehouse schema	1
1.2	OLAP Operations	1
1.3	DataWarehouse Architecture, Data Warehousing to Data Mining	1
1.4	Datamining Concepts and Applications, Knowledge Discovery in Database Vs Data mining	1
1.5	Architecture of typical data mining system,Data Mining Functionalities	1
1.6	Data Mining Functionalities, Data Mining Issues	1
Module 2(Data Preprocessing) (6 hours) (Text3)		
2.1	Data Preprocessing: Need of Data Preprocessing, Data Cleaning- Missing values, Noisy data.	1
2.2	Data integration	1
2.3	Data transformation	1
2.4	Data Reduction-Data cube aggregation, Attribute subset selection	1
2.5	Data Reduction-Dimensionality reduction	1

2.6	Numerosity reduction, Discretization and concept hierarchy generation	1
Module 3(Advanced classification and Cluster analysis)(9 hours)(Text2,Text3)		
3.1	Classification- Introduction, Decision tree construction principle, Splitting indices-Information Gain, Gini index	1
3.2	Decision Tree- ID3	1
3.3	Decision Tree- ID3	1
3.4	Decision tree construction with presorting- SLIQ	1
3.5	Accuracy and error measures, evaluation	1
3.6	Introduction to clustering, Clustering Paradigms	1
3.7	Partitioning Algorithm- PAM	1
3.8	Hierarchical Clustering-DBSCAN	1
3.9	Categorical Clustering-ROCK	1
Module 4(Association Rule Analysis) (8 hours) (Text2,Text3,Text1)		
4.1	Association Rules: Introduction, Methods to discover association rules	1
4.2	A priori algorithm (Level-wise algorithm)	1
4.3	A priori algorithm (Level-wise algorithm)	1
4.4	Partition Algorithm	1
4.5	Pincer Search Algorithm	1
4.6	Pincer Search Algorithm	1
4.7	Dynamic Itemset Counting Algorithm	1
4.8	FP-tree Growth Algorithm	1
Module 5(Advanced Data Mining Techniques) (7 hours) (Text1, Text3)		
5.1	Web Mining - Web Content Mining	1
5.2	Web Structure Mining- Page Rank	1
5.3	Web Structure Mining –Clever algorithm	1
5.4	Web Usage Mining- Preprocessing, Data structures	1

5.5	Web Usage Mining -Pattern Discovery, Pattern Analysis	1
5.6	Text Mining-Text Data Analysis and information Retrieval, Basic measures for Text retrieval	1
5.7	Text Retrieval methods, Text Indexing Techniques Query Processing Techniques	1

