

Total no. of Pages-1

(14)

Roll no. ....

FIFTH SEMESTER B. Tech. (Information Technology)  
END-SEMESTER EXAMINATION, December- 2021

Course Code- ITITE05  
Course Title- Data warehouse and data mining

Time- 3.00 Hours

Max. Marks- 40

Note: - All questions are compulsory. Attempt only two out of three parts from every question.  
Missing data/ information if any, maybe suitably assumed & mentioned in the answer.

No	Question	Marks	CO
1a	What are the various principles of Data Warehouse?	4	CO1
1b	What are the advantages of using Data Marts?	4	CO1
1c	How is data cleaning different from data transformation?	4	CO3
2a	Explain multi-dimensional modelling with a suitable example of your choice.	4	CO3
2b	What is a data cube? How can it be used in dimensional modelling?	4	CO3
2c	Using "Sales" as a factor, differentiate between Star and Snowflake schema.	4	CO3
3a	Explain some document oriented NoSQL databases with their applications.	4	CO3
3b	What is an OLAP server? Explain the working of ROLAP server with its architecture.	4	CO3
3c	What are the advantages of using a Hybrid OLAP server?	4	CO3
4a	Use the following information to solve for Apriori algorithm with Support=50% and Confidence=60%:	4	CO5
Transaction		List of items	
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	
4b	Using the given information, apply KNN algorithm to classify data point( $X_1=3, X_2=7$ ).	4	CO5

$X_1 = \text{Acid Durability (seconds)}$   
 $X_2 = \text{Strength (kg/square meter)}$   
 $Y = \text{Classification}$

7	7	Bad
7	4	Bad
3	4	Good
1	4	Good

4c	Explain the process of KDD using a suitable example and its diagram.	4	CO2
5a	Differentiate between: (i) KNN and K-means clustering (ii) Genetic Algorithm and Neural Networks	2+2=4	CO4
5b	Write short notes on the following: (i) Supervised Learning (ii) Decision Tree	2+2=4	CO2/ CO4
5c	Cluster the following eight points (with $(x, y)$ representing locations) into three clusters: $A_1(2, 10), A_2(2, 5), A_3(8, 4), A_4(5, 8), A_5(7, 5), A_6(6, 4), A_7(1, 2), A_8(4, 9)$ Initial cluster centers are: $A_1(2, 10), A_4(5, 8)$ and $A_7(1, 2)$ .	4	CO5

## END SEMESTER EXAMINATION December 2021

151

02/20

Course Code: COCSC16

Course Title: Data Mining

Time: 3 Hours

Max. Marks : 40

Note: - Attempt all the five questions. Missing data/ information if any, maybe suitably assumed & mentioned in the answer.

	Question	Marks	CO																		
1	Attempt any 2 parts of the following.																				
	Elaborate various stages of Data Mining Process.	4	CO1																		
	Differentiate classification and Regression for predictive analysis. Explain clustering.	4	CO1																		
	Suppose the fraction of undergraduate students who play football is 15% and the fraction of graduate students who play football is 23%. If one-fifth of the college students are graduate students and the rest are undergraduates, what is the probability that a student who plays football is a graduate student? Also, Suppose 30% of the graduate students live in hostel but only 10% of the undergraduate students live in hostel. If a student plays football and lives in hostel, is he or she more likely to be a graduate or undergraduate student? You can assume independence between students who live in hostel and those who play football.	4	CO1																		
2	Attempt any 2 parts of the following.																				
1	Discuss four techniques to deal with missing data in dataset along with suitable examples.	4	CO1																		
2	Calculate Dissimilarity matrix for given dataset? Object: 1, 2, 3, 4, 5, 6 Values: 40, 50, 42, 21, 30	4	CO2																		
3	If two data objects are given as $x = \{3, 2, 0, 5, 0, 0, 0, 2, 0, 0\}$ and $y = \{1, 0, 0, 0, 0, 0, 0, 1, 0, 2\}$ . Calculate its Cosine Similarity.	4	CO1																		
3	Attempt any 2 parts of the following.																				
a	Consider market basket dataset shown in the following table.	4	CO4																		
	<table border="1"> <thead> <tr> <th>T. ID</th> <th>Items Purchased</th> <th>T. ID</th> <th>Items Purchased</th> <th>T. ID</th> <th>Items Purchased</th> </tr> </thead> <tbody> <tr> <td>1</td> <td>{a,d,e}</td> <td>3</td> <td>{a,d,b}</td> <td>5</td> <td>{b,c}</td> </tr> <tr> <td>2</td> <td>{a,d,b,c}</td> <td>4</td> <td>{a,e}</td> <td>6</td> <td>{a,d,b,e,c}</td> </tr> </tbody> </table>	T. ID	Items Purchased	T. ID	Items Purchased	T. ID	Items Purchased	1	{a,d,e}	3	{a,d,b}	5	{b,c}	2	{a,d,b,c}	4	{a,e}	6	{a,d,b,e,c}		
T. ID	Items Purchased	T. ID	Items Purchased	T. ID	Items Purchased																
1	{a,d,e}	3	{a,d,b}	5	{b,c}																
2	{a,d,b,c}	4	{a,e}	6	{a,d,b,e,c}																
b	Compute frequent pattern generated by FP growth algorithm if Minimum Support is 2.																				
	For following given dataset, generate association rules using Apriori Algorithm. Consider Min Support as 50% and Confidence as 75%.	4	CO4																		
	<table border="1"> <thead> <tr> <th>Transaction ID</th> <th>Items Purchased</th> <th>Transaction ID</th> <th>Items Purchased</th> </tr> </thead> <tbody> <tr> <td>T1</td> <td>{bread, egg, cheese}</td> <td>T3</td> <td>{bread}</td> </tr> <tr> <td>T2</td> <td>{egg, juice}</td> <td>T4</td> <td>{bread, egg}</td> </tr> </tbody> </table>	Transaction ID	Items Purchased	Transaction ID	Items Purchased	T1	{bread, egg, cheese}	T3	{bread}	T2	{egg, juice}	T4	{bread, egg}								
Transaction ID	Items Purchased	Transaction ID	Items Purchased																		
T1	{bread, egg, cheese}	T3	{bread}																		
T2	{egg, juice}	T4	{bread, egg}																		

3c

The following contingency table summarizes supermarket transaction data, where *sandwiches* refers to the transactions containing sandwiches, *-sandwiches* refers to the transactions that do not contain sandwiches, *burgers* refers to the transactions containing burgers, and *-burgers* refers to the transactions that do not contain burgers.

	<i>sandwiches</i>	<i>-sandwiches</i>	<b>Sum row</b>
<i>burgers</i>	2000	500	2500
<i>-burgers</i>	1000	1500	2500
<b>Sum col</b>	3000	2000	5000

- Suppose that the association rule "*sandwiches*  $\rightarrow$  *burgers*" is mined. Given a minimum support threshold of 25% and a minimum confidence threshold of 50%, is this association rule strong?
- Compare the results using lift and  $\chi^2$  correlation measures.

Q4

Attempt any 2 parts of the following.

4a

Consider the following 1-D data.

X	0.5	3.0	4.5	4.6	4.9	5.2	5.3	5.5	7.0	9.5
Label	N	N	P	P	P	N	N	P	N	N

- Classify the data point  $X = 5.0$  according to its 3-, and 5- nearest neighbors.
- What would be the class label if distance-weighted voting approach is used?

4b

Consider the dataset shown in the following table having three attributes A1, A2, and A3. Predict the class label for a test sample ( $A_1 = 0, A_2 = 1, A_3 = 0$ ) using the Naïve Bayes Algorithm.

Instance	A1	A2	A3	Class Label
1.	0	0	0	Class 1
2.	0	0	1	Class 2
3.	0	1	1	Class 2
4.	0	1	1	Class 2
5.	0	0	1	Class 1
6.	1	0	1	Class 1
7.	1	0	1	Class 2
8.	1	0	1	Class 2
9.	1	1	1	Class 1
10.	1	0	1	Class 1

4c

Explain following: GINI Index, Entropy, Average Information Entropy and Information gain.

Q5

Attempt any 2 parts of the following.

5a

Expectation Maximization clustering works on which parameter and how it is done. What is the role of Parsing and soft parsing.

5b

Given points five points C1, C2, C3, C4 and C5 where  $d(C_1, C_2) = 7, d(C_1, C_3) = 4, d(C_1, C_4) = 1, d(C_1, C_5) = 2, d(C_2, C_3) = 8, d(C_2, C_4) = 5, d(C_2, C_5) = 11, d(C_3, C_4) = 9, d(C_3, C_5) = 10$  and  $d(C_4, C_5) = 3$ . Perform Agglomerative clustering (Single Linkage) on these points and find Dendrogram.

5c

Divide the given sample data in three (3) clusters using k-means Algorithm.  
 Height: 183, 171, 167, 176, 180, 177, 180, 180, 182, 183, 185, 185  
 Weight: 70, 56, 60, 72, 84, 76, 71, 68, 69, 77, 72, 74

## MID-SEMESTER EXAMINATION, SEPTEMBER - 2022

Course Code – ITITE05

Course Title – Data Warehouse and Data Mining

Time – 1.5 Hours

Max. Marks – 15

**Note – Attempt all questions. Missing data/information (if any), may be suitably assumed & mentioned in the answer.**

Q. No	Questions	Marks	CO
1a	What is a Data Warehouse? Explain three-tier Data Warehouse architecture with a neat and labeled diagram.	2	CO1
1b	How is a data warehouse different from an operational database? Distinguish between them w.r.t. OLAP and OLTP.	1	CO3
2a	What can you accomplish with data gateway? Explain with diagrammatic illustration the working of data gateway. How is standard mode different from personal mode of data gateway?	2	CO1
2b	What are the steps involved in data cleaning? Briefly explain each one.	1	CO1
3a	What is Database synchronization? Discuss the concepts of Unidirectional and Bidirectional Database synchronization taking suitable example.	2	CO1
3b	Why is market moving from ETL to ELT in the cloud data warehouse? Discuss the pros and cons of ETL vs. ELT.	1	CO1
4a	What do you understand by Dimensional Modeling? Highlight and explain the key steps decided during the design of a Dimensional Model.	2	CO3
4b	State the objectives of JAD session. What is the Business Analyst's role in one?	1	CO3
5a	Suppose that a data warehouse for <i>Big-University</i> consists of the following four dimensions: <i>student</i> , <i>course</i> , <i>semester</i> , <i>instructor</i> and two measures <i>count</i> , <i>avg_grade</i> . When at the lowest conceptual level (e.g. for a given <i>student</i> , <i>course</i> , <i>semester</i> and <i>instructor</i> combination), the <i>avg_grade</i> measure stores the actual course grade of the student. At higher conceptual levels, <i>avg_grade</i> stores the average grade for the given combination.  (i) Draw a snowflake schema diagram for the data warehouse.	1.5	CO3

Total no. of Pages 1

Roll no. ....

## SEMESTER-UG 5th

## MID-SEMESTER EXAMINATION, September, 2022

Course Code- COCSC16, CDCSC16

Course Title- Data Mining

Time- 1.5 Hours

Max. Marks- 15

Note: - Attempt all questions. Missing data/ information (if any) may be suitably assumed & mentioned in the answer.

Q. No.	Question	Marks	CO
1a	What are the differences between the three main types of data warehouse usage: information processing, analytical processing, and data mining? Discuss the motivation behind OLAP mining (OLAM).	2	CO1
1b	Given the attribute name and values of attributes, classify its type (i) Gender-Male, Female (ii) Height-5.8,6.2 (iii) Color-Black, Brown, Grey	1	CO1
2a	A poker-dealing machine is supposed to deal cards at random, as if from an infinite deck. In a test, you counted 1600 cards, and observed the following: Spades 404 , Hearts 420 , Diamonds 400 , Clubs 376 Could it be that the suits are equally likely? Or are these discrepancies too much to be random?	2	CO2
2b	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 2010?	1	CO2
3a	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, 36, 40, 45, 46, 52, 70. (a) What is the mean of the data? What is the median? (b) What is standard deviation? What is the mode of the data?	2	CO3
3b	Suppose that $x$ and $y$ are the first two term-frequency vectors i.e., $x = (5,0,3,0,2,0,0,2,0,0)$ and $y = (3,0,2,0,1,1,0,1,0,1)$ . How similar are $x$ and $y$ ? compute the cosine similarity between the two vectors.	1	CO3
4a	Apply discretization on following data objects using binning: a) [7,5,7,7,7,6,9,8] b) [1,2,4,6,7,8,9,8,4]	2	CO4
4b	Given two sets of $s_1$ and $s_2$ , $S_1 = \{1, 2, 3, 4, 5\}$ , $s_2 = \{4, 5, 6, 7, 8, 9, 10\}$ . find the Jaccard Index and the Jaccard Distance between the two sets.	1	CO4
5a	In a class mark of students in two subjects are enlisted as follows: {55, 11, 36, 72, 41, 91, 96, 63, 29, 44} and {32, 67, 89, 23, 17, 48, 59, 77, 61, 7}. Find the covariance matrix for this data.	2	CO5
5b	Explain the techniques to deal with noise in the dataset.	1	

Total No. of Pages: 2

Roll No. \_\_\_\_\_

V/VII SEM - B. TECH END SEMESTER EXAMINATION, NOV-DEC-2022

Course Code: TTITE05

Course Title: Data Warehouse and Data Mining

Time: 3 Hours

Max. Marks: 40

Note: - Attempt all the five questions. Missing data/ information (if any), may be suitably assumed & mentioned in the answer.

Q. No.	Questions	Marks	CO
Q1	Attempt any two parts of the following: 1a Data warehousing is the only viable means to resolve the information crisis and to provide strategic information. List four reasons to support this assertion and explain them. 1b What is a Data Warehouse? Highlight and discuss the four key characteristics of a Data Warehouse. 1c Explain the major building blocks/components of a Data Warehouse with a well-labeled diagram.	4 4 4	CO1 CO1 CO1
Q2	Attempt any two parts of the following: 2a Elaborately discuss the concepts of the following schemas in Data Warehouse modelling (taking suitable examples with well-labeled diagram for each one): (i) Snowflake schema. (ii) Galaxy schema 2b Suppose that a data warehouse for Gross sales consists of the four dimensions: date, spectator, location and movie, and the two measures, count and charge, where charge is the fare that a spectator pays when watching a movie on a given date. Spectators may be students, adults or seniors, with each category/status having its own charge rate. (i) Draw a star schema diagram for the data warehouse. (ii) Starting with the base cuboid [date, spectator, location, movie], what specific OLAP operations should be performed in order to list the total charge paid by student spectators at PVR_Vegas in the year 2021?	4 (2 + 2)  (2 + 2)	CO3 CO3
2c	Suppose a company would like to design a data warehouse to facilitate the analysis of moving vehicles in an online analytical processing manner. The company registers huge amounts of auto movement data in the format of (Auto_ID, location, speed, time). Each Auto_ID represents a vehicle associated with information, such as vehicle_category, driver_category, etc., and each location may be associated with a street in a city. Assume that a street map is available for the city. (i) Design such a data warehouse to facilitate effective online analytical processing in multidimensional space. (ii) The movement data may contain noise. Discuss how you would develop a method to automatically discover data records that were likely erroneously registered in the data repository.	4 (2 + 2)	CO3

<b>Q3</b>	Attempt any two parts of the following:																																		
3a	What are the differences between these three main types of data warehouse usage: <i>information processing</i> , <i>analytical processing</i> and <i>data mining</i> ? Discuss the motivation behind <i>OLAM</i> ( <i>OLAP mining</i> ).	4	CO2, CO4																																
3b	Differentiate between <i>ROLAP</i> , <i>MOLAP</i> and <i>HOLAP</i> . Explain the following <i>OLAP operations</i> in the multidimensional data model (taking suitable example with diagram).  (i) <i>Roll-up</i> (ii) <i>Drill-down</i> (iii) <i>Slice for</i> (iv) <i>Dice for</i> (v) <i>Pivot</i>	4  (1½ + 2½)	CO3																																
3c	Give an introduction to NoSQL databases. Compare the applicabilities of CouchDB and MongoDB.	4	CO2, CO3																																
<b>Q4</b>	Attempt any two parts of the following:																																		
4a	Explain how <i>decision tree induction</i> algorithm works. Elucidate the concept of its key factors: <i>Entropy</i> and <i>Information Gain</i> . How are they used to build <i>decision trees</i> ?	4	CO5																																
4b	(i) Explain <i>Classification</i> and <i>Prediction</i> in Data Mining. (ii) Differentiate between eager learners and lazy learners. Give examples for both of them.	4  (2 + 2)	CO2, CO4																																
4c	Consider the following dataset for a <i>classification</i> task. Using <i>k-Nearest Neighbor (KNN)</i> algorithm with <i>Euclidean Distance</i> as the distance metric, predict the class label <i>Fruit Taste</i> for data sample <i>Strawberry</i> with $k=1, 3$ and $5$ . Write down your calculations clearly.	4	CO5																																
	<table border="1"> <thead> <tr> <th>Fruit</th> <th>Sweetness</th> <th>Sourness</th> <th>Fruit Taste</th> </tr> </thead> <tbody> <tr> <td>Lemon</td> <td>1</td> <td>9</td> <td>Sour</td> </tr> <tr> <td>Grapefruit</td> <td>2</td> <td>8</td> <td>Sour</td> </tr> <tr> <td>Orange</td> <td>3</td> <td>7</td> <td>Sour</td> </tr> <tr> <td>Cherry</td> <td>6</td> <td>4</td> <td>Sweet</td> </tr> <tr> <td>Banana</td> <td>9</td> <td>1</td> <td>Sweet</td> </tr> <tr> <td>Grapes</td> <td>8</td> <td>2</td> <td>Sweet</td> </tr> <tr> <td>Strawberry</td> <td>5</td> <td>5</td> <td>?</td> </tr> </tbody> </table>	Fruit	Sweetness	Sourness	Fruit Taste	Lemon	1	9	Sour	Grapefruit	2	8	Sour	Orange	3	7	Sour	Cherry	6	4	Sweet	Banana	9	1	Sweet	Grapes	8	2	Sweet	Strawberry	5	5	?		
Fruit	Sweetness	Sourness	Fruit Taste																																
Lemon	1	9	Sour																																
Grapefruit	2	8	Sour																																
Orange	3	7	Sour																																
Cherry	6	4	Sweet																																
Banana	9	1	Sweet																																
Grapes	8	2	Sweet																																
Strawberry	5	5	?																																
<b>Q5</b>	Attempt any two parts of the following:																																		
5a	What is <i>Cluster Analysis</i> ? Describe the different types of <i>Clustering methods/techniques</i> with examples.	4	CO5																																
5b	For the following transactional dataset, find the <i>frequent itemsets</i> and generate <i>strong association rules</i> using <i>Apriori algorithm</i> . Consider the thresholds as <i>min. support = 40%</i> and <i>min. confidence = 70%</i> . Calculate the <i>confidence</i> of each rule and identify all the <i>strong association rules</i> .	4	CO5																																
	<table border="1"> <thead> <tr> <th>T_ID</th> <th>Items bought</th> </tr> </thead> <tbody> <tr> <td>1</td> <td>{Bread, Butter, Milk}</td> </tr> <tr> <td>2</td> <td>{Bread, Butter}</td> </tr> <tr> <td>3</td> <td>{Beer, Cookies, Diapers}</td> </tr> <tr> <td>4</td> <td>{Milk, Diapers, Bread, Butter}</td> </tr> <tr> <td>5</td> <td>{Beer, Diapers}</td> </tr> </tbody> </table>	T_ID	Items bought	1	{Bread, Butter, Milk}	2	{Bread, Butter}	3	{Beer, Cookies, Diapers}	4	{Milk, Diapers, Bread, Butter}	5	{Beer, Diapers}																						
T_ID	Items bought																																		
1	{Bread, Butter, Milk}																																		
2	{Bread, Butter}																																		
3	{Beer, Cookies, Diapers}																																		
4	{Milk, Diapers, Bread, Butter}																																		
5	{Beer, Diapers}																																		
5c	How <i>FP-tree</i> is better than <i>Apriori Algorithm</i> ? Using <i>FP-Growth</i> algorithm, demonstrate construction of <i>FP-tree</i> for the transactional dataset given above in (b) part of Question 5.	4	CO5																																

VIII

13/20

160

No of pages: 1 MID SEMESTER EXAMINATION FEB 2020 Roll No. ....

Course code: ITD08 Course: DATA WAREHOUSE AND DATA MINING Sem: 8

Branch: COE & IT& IC Max. Marks : 15

Note: All questions carry 5 marks. Time: 1.30 hour

---

Q.1 Define the data warehousing. Explain the need and purpose of developing data warehouse. Discuss 3-tier architecture of data warehouse?

OR

Why ETL is consider an essential process in data warehousing. Explain the requirement of ETL process. How ETL process is optimized in term of processing speed?

Q.2 What is Meta data? Discuss the various types of the Meta data and their utilities in the data warehouse?

OR

What are the objectives of the dimensional modeling? List and explain various dimensional modeling schemas in data warehouses. Also Create the star schema for loan department of a banking organization?

Q.3 What do you mean by the OLAP? Give the classification of the OLAP. Also list the OLAP operations with supporting example?

OR

Given the marks of four students as 8, 10, 15 and 20. How to Calculate the Z-Score normalization and Min-max normalization of the given data?

-----End-----

Total No. of Page: 2

Roll No. \_\_\_\_\_

BTech Semester V

CO3

## END SEMESTER EXAMINATION DECEMBER 2022

Course Code: COCSC16, CDCSC16

Course Title: Data Mining

Time: 3 Hours

Max. Marks: 40

Note: - Attempt all the five questions. Missing data/ information if any, maybe suitably assumed &amp; mentioned in the answer

Q. No.	Question	Marks	CO																
Q1	Attempt any 2 parts of the following.																		
1a	Explain various stages of data mining with diagram. Marks of 10 students in Data Mining subject are given as follows: 78, 51, 86, 91, 33, 27, 25, 46, 55, 59. Calculate normalized values of these using Min-Max normalization.	2+2	CO1+ CO3																
1b	Given the confusion matrix below, (a) Calculate the accuracy. Is this a good result? Why? (b) Explain 3 alternative performance measures to support your opinion.	2+2	CO2+ CO4																
	<table border="1"> <tr> <td></td><td>True fraud</td><td>True non-fraud</td></tr> <tr> <td>Predicted fraud</td><td>1</td><td>4</td></tr> <tr> <td>Predicted non-fraud</td><td>9</td><td>986</td></tr> </table>		True fraud	True non-fraud	Predicted fraud	1	4	Predicted non-fraud	9	986									
	True fraud	True non-fraud																	
Predicted fraud	1	4																	
Predicted non-fraud	9	986																	
1c	Regression. Given the graph below, predict y by x, if $MAE < 0.75$ then it's a good result. Use one sentence for each regression method to explain whether it would have good results or not. (a) kNN for regression, $k = 3$ (b) binary tree regression (c) linear regression (d) local regression, $k = 3$ (e) local regression, $k = 20$	4	CO4																
Q2	Attempt any 2 parts of the following.																		
2a	What do you mean by a factless fact table in the context of data warehousing? Differentiate between fact table and dimension table.	4	CO5																
2b	Price of a few items in a store is given as follows. Calculate Dissimilarity matrix for this.	4	CO2+ CO4																
	<table border="1"> <tr> <td>Object No</td><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td><td>6</td><td>7</td></tr> <tr> <td>Price</td><td>28</td><td>45</td><td>64</td><td>93</td><td>59</td><td>67</td><td>38</td></tr> </table>	Object No	1	2	3	4	5	6	7	Price	28	45	64	93	59	67	38		
Object No	1	2	3	4	5	6	7												
Price	28	45	64	93	59	67	38												
c	Predict that an Average build person having high Weight, having good experience and having fair performance will win in wrestling or not.	4	CO3																

Build	Weight	Experience	Performance	Win?
Average	Medium	Good	Fair	No
Good	High	Moderate	Excellent	Yes
Normal	Low	Moderate	Fair	No
Good	Medium	Good	Fair	No
Good	High	Moderate	Fair	Yes
Average	High	Moderate	Excellent	Yes
Normal	Medium	Good	Fair	No
Normal	High	Moderate	Excellent	Yes
Average	Low	Good	Fair	No
Normal	Medium	Good	Fair	No
Average	High	Moderate	Excellent	No

Q3 Attempt any 2 parts of the following.

- 3a Explain following with one example with respect to itemset mining:
- Frequent itemset
  - Minimum support
  - Histogram
  - Interesting measures

- 3b Trace the results of using the Apriori algorithm on the grocery store example with support threshold  $s=33.34\%$  and confidence threshold  $c=60\%$ . Show the candidate and frequent itemsets for each database scan. Enumerate all the final frequent itemsets. Also indicate the association rules that are generated and highlight the strong ones, sort them by confidence.

Transaction ID	Items
T1	HotDogs, Buns, Ketchup
T2	HotDogs, Buns
T3	HotDogs, Coke, Chips
T4	Chips, Coke
T5	Chips, Ketchup
T6	HotDogs, Coke, Chips

3c Use the following transactional with support threshold=3

Item	Frequency
T1	{E,K, M,N, O,Y}
T2	{D,E,K,N, O,Y}
T3	{A,E,K,M}
T4	{C,K,M,U,Y}
T5	{C,E,I,K,O,O}

- a.) Build a frequent pattern tree (FP-Tree). Show for each transaction how the tree evolves.  
 b.) Use Fp-Growth to discover the frequent itemsets from this FP-tree

Q4 Attempt any 2 parts of the following.

- 4a Given a decision tree, you have the option of (a) converting the decision tree to rules and then pruning the resulting rules, or (b) pruning the decision tree and then converting the pruned tree to rules. What advantage does (a) have over (b)?

- 4b State 2 differences and 2 similarities between bagging and boosting. Suppose you are given 'n' predictions on test data by 'n' different models ( $M_1, M_2, \dots, M_n$ ) respectively. Which method(s) can be used to combine the predictions of these models?

- 4c Dataset collected for whether information over a course of period is given below. Using the ID3 algorithm calculate the outcome.

4 CO1+  
CO5

4 CO2

4 CO2

4 CO2

2+2 CO2+  
CO4

4 CO3

Build	Weight	Experience	Performance	Win?
Average	Medium	Good	Fair	No
Good	High	Moderate	Excellent	Yes
Normal	Low	Moderate	Fair	No
Good	Medium	Good	Fair	No
Good	High	Moderate	Fair	Yes
Average	High	Moderate	Excellent	Yes
Normal	Medium	Good	Fair	No
Normal	High	Moderate	Excellent	Yes
Average	Low	Good	Fair	No
Normal	Medium	Good	Fair	No
Average	High	Moderate	Excellent	No