Paper / Subject Code: BEX32 / Elective-I (1) Data Mining & W. Proper / Such Per villes uning Aprilori Algorithm for the following set of a large and and Confidence 80% A.D.E. sper / Subject Code: BERS2 / Election-5 (1) Data Calculate the best splits using Shannon's Ent Chart the resulting docume trees using both place between Bayerian Belief Nictworks and A.C.D.E 200 100 B.C.E about sodes on the following secondly date 400 B.D.E 400 C.D Supervised methods Scott-supery sed methydd 600 A.B.C. 700 Unsupervised methods A.D.E the help of a near labelled diagram of 800 A.B,E 900 data warehouses a pre-requisite for 1000 at Write a short note on Data Mining Query Language. dain with unitable example the folia d) How is data warehouse different from a database? How are they similar? a) Star Schema b) Socialiske Schema inferentiate between OLTP and O Q3 at Briefly outline the major steps of decision tree classification. aplain multidimensional data t M box is Data Cube? Given an b) Given two objects represented by the tuples (22, 1, 42, 10) and (20, 0, 36, 8); a) Compute the Euclidean Distance between the two objects. b) Compute the Manhattan Distance between the two objects. Compute the Minkowski Distance between the two objects, using q=3. What is metadata with rewarehousing? Explain two types of data in cluster analysis. Explain the following Shared n Q.4a) Differentiate between Clustering and Classification. Give a brief application for each. Write a short not b) Consider the following data: a) MOLAP b) ROLAR TID Home Owner Marital Annual Income Class: Loan Status Defaulter S .8a) List and bri 125 N M 100 N S b) Write sho 70 N M N 120 N 6 D N 95 M 60 N D N 220 S 85 M 75 S 90 2 F08749487D39F86420B23A89D6FE11C7

a) 1

(d

Usin

IT 7-(E-I) 1(RC)

### THE REPORT OF THE PARTY OF THE PARTY.

### MODULE-II

401

- 3. a) Why is tree pruning useful in decision tree induction? How does tree pruning work?
  - b) Given two objects A1(3, 10) and A2(2, 5)
    - i) Compute the Euclidean distance between the two objects.
    - ii) Compute the Manhattan distance between the two objects.
    - iii) Compute the Minkowski distance between the two objects (use p = 3).
  - c) Design an efficient method that performs effective naive Bayesian classification over an infinite data stream.
  - d) What are the characteristics for agglomerative and divisive hierarchial clustering?
- 4. a) Explain the K-means algorithm for clustering. Differentiate between K-means and K-medoids.
  - b) Differentiate between Bayesian belief networks and Naive Bayesian classifier.
  - c) Write short notes on :
    - i) Density based methods.
    - ii) Hierarchial methods.

### MODULE-III

- 5. a) Differentiate between OLTP and Data Warehousing.
  - b) What is anamoly detection? List and explain anomaly detection methods.
  - Mention the need, functions and applications of data warehouses in the field of data mining.
- 6. a) With the help of a neat labelled diagram, explain the data warehouse architecture.
  - b) Write note on the following:
    - i) Star Schema
    - ii) Snowflake Scheme
    - iii) Fact Constellation Schema
    - iv) Efficient computation of data cube.

# MERCHANISM MERCHANISM

- a) What is web usag
  - b) Explain in brief th i) Intraquery part

    - ii) Interquery par
  - c) Compare ROLA
- a) With the help of parallel process
  - b) Write short note
    - i) Web content
    - ii) OLAP tools a
  - c) List the four typ

(5+5)

| 770  | ne  | e 1 | Su | bj | ect | C |
|------|-----|-----|----|----|-----|---|
| P 23 | 112 | 200 |    |    |     |   |

With the help of

Write short notes b) Deviation 1) Statistica

0.5

0.6

Q.7

Q.8

98

10

10.

ii)

Write short note Star Sch (i Snowfla

ii) Fact co iii)

Data cu iv)

Write short no Web 1)

ROL ii) MOL iii)

Meta iv)

a) What is the

b) Discuss private e) Explain the

Subject Code: RES 12 / Elective I (1) Data Mining & Warehousing. Studen Martin Named Rest Ken Module II Consers a decision role for the data give. Buys computer Credit rating No Fair No Excellent Hugh Yes Fair Yes High N Fair Heb N Yes Fair Medium No Excellent Low Y Yes Excellent Low Y No 5.40 Fair Low Yes Mediam Fair Yes Low Fair <= 30 Y Medium Yes Excellent >40 Y Medam Yes <= 30 Excellent N Medium 31---40 Yes Fair Y 31-40 High No Excellent

b) Write a short note on DBSCAN.

Ose entropy and information gain.

>40

Medium

a) Consider the objects given below. Assume  $C_1 = x_2$  and  $C_2 = x_8$  perform Kmedoid using Manhattan distance. Note down the Error criterion for the next aeratori choose a random object  $x_7$  as medoid. Note down the inference.  $x_1(2.6) x_2(3.4) x_3(3.8) x_4(4.7) x_5(6.2) x_6(6.4) x_7(7.3) x_8(7.4) x_9(8.5)$ 

b) Jack Mary and Jin are subjected to tests for certain illness they are suffering from The results are as follow.

| Name  | Gender     | Fever | Count | IT.    | 1      |           |        |
|-------|------------|-------|-------|--------|--------|-----------|--------|
| Jack  |            |       | Cough | Test-1 | Test-2 | Test-3    | Test-4 |
| 146.K | Male       | Yes   | N     | P      | N      | N         | 27     |
| Mary  | Female     | Yes   | N     | D      | 1      | - Indiana | N      |
| E1    | Male       | Yes   | D     | 1      | N      | P         | N      |
| 1     | Invite and | _     | P     | N      | N      | N         | N      |

ify whether Jack and Mary suffer from same illness 11)

Justify whether Jack and Jin suffer from same illness Justify whether Jin and Mary suffer from same illness. (iii)

# Paper / Subject Code: BE832 / Elective-I (1) Data Mining & Warehousing.

| 9 |     | Module III   |       |
|---|-----|--|-------|
|   |     | a) With the help of a neat diagram explain data warehouse architecture.  | 10    |
|   |     | With the help of a neat diagram explain data water.  | 10    |
|   | 100 | a) Write short notes on b) Deviation Based Technique   |       |
|   | 0.5 | 1) Based Technique   |       |
|   |     | ii) Statistical Paris  | 15×4= |
|   |     | Write short notes on   | 20    |
|   |     | W Star Schelle   |       |
|   | 0.0 | Complake schema  |       |
|   |     | Fact consideration settlement  |       |
|   |     | iv) Data cube  |       |
|   |     | Module IV  |       |
|   |     |  | 20    |
|   |     | a) Write short notes on  |       |
|   | - 7 | as a straight manny  |       |
|   | Q.7 | 1/ mov AD  |       |
|   |     | M/ - COY AD  |       |
|   |     | Metadata Interchange Intaktive   | 06    |
|   |     | 117  | 10    |
|   |     | What is the significance of user believe with regard to data Mining  | 04    |
|   | Q.8 | a) What is the significance of user behavior mining?     b) Discuss privacy protection technique with regard to data Mining. |       |
|   |     | b) Discuss privacy protection c) Explain the need for OLAP.  |       |
|   |     |  |       |

# 

# B.E. (I.T.) (Semester – VII) (RC) (2007-08) Examination, Nov./Dec. 2017 DATA MINING AND WAREHOUSING (Elective – I)

puration: 3 Hours

al

Total Marks: 100

Instructions: 1) Answer any five questions with at least one from each Module.

2) All questions carry equal marks.

### MODULE - I

a) In real-world data, tuples with missing values for some attributes are a common occurrence. Describe various methods for handling this problem.

5

b) Find Association rules using Apriori algorithm for the following set of transactions. Consider Support = 3 and confidence =80%.

8

| TID  | Item List  |
|------|------------|
| 100  | A, D, E    |
| 200  | A, B, C, E |
| 300  | A, B, D, E |
| 400  | A, C, D, E |
| 500  | B, C, E    |
| 600  | B, D, E    |
| 700  | C, D       |
| 800  | A, B, C    |
| 900  | A, D, E    |
| 1000 | A, B, E    |

c) With the help of neat block diagram, explain the data mining process.

7

 a) Suppose that the data for the analysis include attribute the frequency of the stop words in documents. The values are given in increasing order

13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 30, 33, 33, 35, 35, 35, 35, 36, 40, 45, 46, 52, 70

Apply following methods.

6

- i) Use smoothing by bin with a depth of 3.
- ii) Use min max normalization to transform the value 40 into range from 0.0 to 1.0.
- iii) Use z score normalization to normalize 70.

P.T.O.

- b) List and explain different data mining stores on which data mining can be IT 7 - (E-I) 1 (RC) 2007-08
  - c) Explain pattern interestingness measure as one of the data mining primitives.

THE REAL COLUMN TO THE REAL PROPERTY OF THE PERSON NAMED IN COLUMN TO PERSON NAMED IN COLUMN TO

- d) Write short note on Data Mining Query language.

### MODULE - II

- 3 a) Write short note on following clustering techniques :
  - i) Density based methods
  - ii) Grid based methods.
  - b) Draw decision tree for following data set. Explain steps.

| Draw decisio | n tree to temp | Car Type | Income  | Class |
|--------------|----------------|----------|---------|-------|
| CustID       | Gender         | Cui . yr | Average | CO    |
| 1 2          | М              | Family   |         | C1    |
|              | F.             | Sports   | High    |       |
|              | M              | Luxury   | High    | C1    |
| 3            |                | Family   | Low     | CO    |
| 4            | M              |          | Average | CO    |
| 5            | M              | Sports   |         |       |
| 6            | F              | Luxury   | High    | C1    |
|              | F              | Luxury   | High    | C1    |
| 7            |                | Family   | Low     | C1    |
| 8            | М              | ганну    |         |       |
| 9            | F              | Luxury   | High    | CO    |
| 10           | М              | Sports   | High    | C     |

- c) Explain K nearest neighbor classification with suitable example.
- 4. a) Consider following dataset 1, 2, 6, 7, 8, 10, 15, 17, 20. From 3 cluster using k-medoid algorithm and considering 6, 7, 8 as 3 medoids of 3 clusters respectively.
  - b) What are the requirements of clustering in data mining?
  - c) List and explain different applications of classification and prediction.
  - d) Compare the advantages and disadvantages of eager classification versus lazy classification.

- a) Describe major applica
- b) Write short note on : i) Deviation based t
  - ii) Statistical Based
  - c) Explain with suitable
    - i) Star Schema
      - ii) Snowflake sch
- 6. a) Differentiate between
  - b) Explain different ( c) Explain Multi tier
  - 7. a) Write short not
    - i) Shared me ii) Shared no
    - b) What is signi
    - c) Give some i
    - 8. a) Write short
      - i) HOLAF
      - ii) MOLA
      - b) Explain d
      - c) Explain in

6

5

### ALBERT HER

### B.E. (IT) (Semester - VII) Examination, December 2009 DATA MINING AND WAREHOUSES (E-I)

ation: 3 Hours

Total Marks: 100

Instruction: Attempt any five questions such that at least one question from each Module is selected.

### MODULE-1

a) What is the main difference between the data warehouse and operational database system? With an example explain the three steps in the process of knowledge 10 discovery in databases.

b) Explain the following types of databases:

- i) Relational database
- ii) Transactional database
- iii) Object oriented database
- iv) Object relational database
- v) Text database and multimedia database.
- a) Do we prefer to use the OLAP operations in the multimedia data model? If yes then what are those operations and if no then why are they not preferred?

6

b) Explain how data mining is considered as a step in the process of knowledge discovery.

8

c) Discuss the methods for filling in the missing values for the attributes in data cleaning.

### MODULE - II

a) Do we prefer to use the neural network in data mining? What are the limitations and consequences of choosing neural network in decision support systems? 6

| *                       | Paper / Subject Code: TE631 / Data Mining  | TE631                         |
|-------------------------|--|-------------------------------|
| ng class<br>et Class (I | c) What is Spatial Data Mining?  | (8 Marks) (4 Marks) (8 Marks) |
| Computer                | a) What is a Data Warehouse?     b) List and explain differences between Operational Database Systems and Data Warehouses.     | (8 Marks)                     |
|                         | c) Explain in detail various steps of ETL process.  PART C   |                               |
|                         | Answer anyone questions from the following-  | (8 Marks)                     |
|                         | a) List and explain any 3 Major Clustering Approaches  | (6 Marks)                     |
|                         | b) Explain How K-Means Clustering algorithm work.     c) Explain how DBSCAN Density –Based Spatial Clustering algorithm works. | (6 Marks)                     |
|                         |  | (8 Marks)                     |
| Marks) Q.8              | a) Explain different types of OLAP Servers.     b) With a diagram explain Star Schema for Data Warehouse design.               | (6 Marks)                     |
| er.<br>(6 Mar           | b) With a diagram explain State Settlement     Explain OLAP operations of Roll-up and Drill-down.                              | (6 Marks)                     |
|                         |  |                               |

(6 Marks

(6 Marks)

(8 Marks)

(6 Marks)

(6 Marks)

c) What is Spatial

a) What is a Data

b) List and expla Warehouses.

c) Explain in d

Answer anyone q

a) List and e

b) Explain

c) Explain

a) Explai

b) With

c) Expl

(4

Marks)

(6 Mar

(8 Marks)

a) Construct first level of decision tree for classification utilizing the following class Q.3 Construct first level of decision free for the Construct first level of decision free for the Class labelled data, using information Gain as Attribute Selection Method. Target Class Buys\_Computer ="Yes" Or "No"

| No  | Age                   | Income | Student | Credit_ Rating | Class;<br>Buys Comp |
|-----|-----------------------|--------|---------|----------------|---------------------|
| 1   | Youth                 | High   | N       | Fair           | No                  |
| 2   | Youth                 | High   | N       | Excellent      | No                  |
| 3   | Middle Aged           | High   | N       | Fair           | Yes                 |
|     | Senior                | Medium | N       | Fair           | Yes                 |
|     | Senior                | Low    | Y       | Fair           | Yes                 |
|     | Senior                | Low    | Y       | Excellent      | No                  |
|     | Middle Aged           | Low    | Y       | Excellent      | Yes                 |
| _   | Youth                 | Medium | N       | Fair           | No                  |
| = - | Youth                 | Low    | Y       | Fair           | Yes                 |
|     | Senior                | Medium | Y       | Fair           | Yes                 |
|     | Youth                 | Medium | Y       | Excellent      | Yes                 |
|     | Middle Aged           | Medium | N       | Excellent      | Yes                 |
| 1   | Middle Aged<br>senior | High   | Y       | Fair           | Yes                 |
|     | CITIO                 | Medium | N       | Excellent      | No                  |

b) Why is "Naïve Bayes Classifier" called Naïve?

c) From the following data predict value of price for a distance value of 50 using liner

| Distance | 40   | 42   | 100  |      |      |
|----------|------|------|------|------|------|
| (kms)    |      | 1    | 45   | 48   | 52   |
| Price    | 8500 | 8250 | 0000 |      |      |
| (INR)    |      | 0230 | 8000 | 7750 | 6500 |

### PART B

Answer any two questions from the following.

No

- a) What is Outlier Analysis/ Anomaly Detection? Explain different variations of Q.4 (6 Marks
  - b) Explain Graphical/ Visual approaches of anomaly detection. List their limitations.
  - c) Compare Distance Based and Density based approaches of outlier detection. (6 Marks)
- a) Explain why Graph Mining is important 0.5
  - b) Explain briefly how social networks behave. (6 Marks)

(6 Marks)

Write sho

i) (iii

iii)

a) Explain b) With th

ware h c) Are the not?

a) How b) Expla

Disc

Writ

iv)

i) ii) iii)

6-

7.

8.

Module 12 mig decision tree for the data given below. Generate the rules [12 mig ted decision tree.

decision tree.

|                   | -effect #  | d dre-     | 441   | Tiay Con              |
|-------------------|--|------------|---|-----------------------|
| a) Co             | in the constructed   | Hity       | Windy   | No                    |
| 800               | - matut  | e Humidity | False   | No                    |
| - Tank            | Temperatur   | High       | True  | Yes                   |
| Outlook           | Hot  | High       | False   | Yes                   |
| Rainy             | Hot  | High       | False   | Yes                   |
| Overcast          | Hot<br>Mild  | High       | False   | No                    |
|                   |  | Normal     | True  |                       |
| Suntry            | Cool   | Normal     | True  | Yes                   |
| Somy              | Cool   | Normal     | Total Control of the | No                    |
| Sunny<br>Overcast | Cool   | High       | False   | Yes                   |
|                   | Mild   | Normal     | False   |                       |
| Rainy             | Cool   |            | False   | Yes                   |
| Rainy             | Mild   | Normal     | True  | Yes                   |
| Sunny             | Mild   | Normal     | True  | Yes                   |
| Rainy             | Mild   | High       |   | Yes                   |
| Overcast          | The same of the sa | Normal     | False   |                       |
| Overcast          | Hot  | High       | True  | No                    |
| Control           | Mild   | ringo      |   | TO THE REAL PROPERTY. |

Use Entropy and Information gain.

b) Perform clustering [hierarchical] for the matrix given below and draw [8 mks] the dendrogram.

| 1  | 2    | 3    | 4     | 5   |
|----|------|------|-------|-----|
| 0  | 1000 | 7 50 | 1000  | 8 5 |
| 9  | 0    |      | 1-018 | 39  |
| 3  | 7    | 0    | 30    |     |
| 6  | 5    | 9    | 0     |     |
| 11 | 10   | 2    | 8     | 0   |

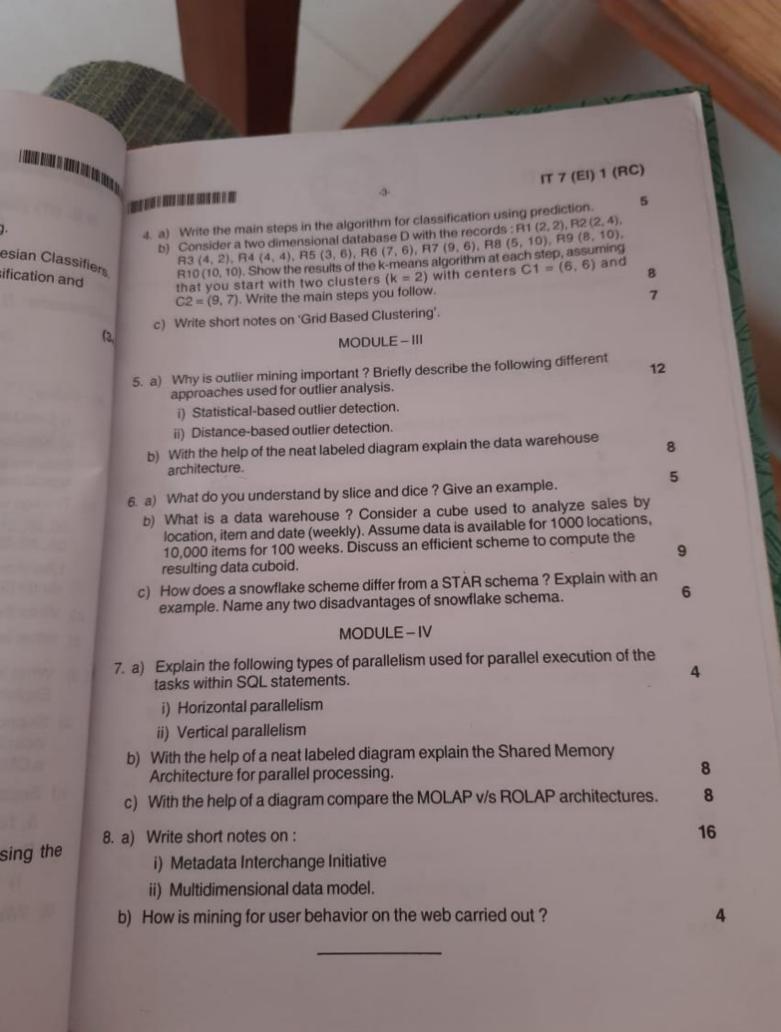
a) Use k-means algorithm to cluster the data into 3-cluster. (3,5,11,12,4,21,32,12,28)

[8 mks]

b) Explain k NN with the help of an example. c) Explain Nominal and Ordinal data.

[8 mks]

[4 mks]



BEB32

### Module III

nks

8.

| a)             | Write short notes on:  Graphical based anomaly detection techniques  Statistical based anomaly detection technique  Distance based anomaly detection technique  Model based anomaly detection technique  | 5 × 4 = 20mks                             |
|----------------|--|---|
| a)<br>b)<br>c) | Explain basic functions of a data warehouse.  With the help of a neat diagram explain overall architecture of a data ware house.  Are the data ware houses a pre-requisite for data mining? Why/Why not?  Module IV                              | [6 mks]<br>[8mks]<br>[6 mks]              |
| a)<br>b)<br>c) | How can mining be used to identify the users behavior on the web? Explain web content mining. Discuss the benefits of data mining for the financial data analysis.  Write short notes on:  Advantages of DB Mines.  Visual and Audio Data Mining | [4 mks]<br>[8 mks]<br>[8 mks]<br>[20 mks] |
|                | iii) Query and Reporting Tools<br>iv) Data Partitioning scheme   |   |

Paper / Subject Code: BE832 / Elective-I (1) Data Mining & Warehousing.

BE832

al No. of Printed Pages 03

### B.E. (Information Technology) Semester- VII (Revised Course 2007-08) **EXAMINATION MAY/JUNE 2019**

Elective-I (1) Data Mining & Warehousing.

Duration : 3 Hours

[Max. Marks :100]

Instructions:

06

06

OR.

0.2

- a) Assume data whenever necessary.
- b) Draw neat labeled diagram using pencil and rules
- c) Answer any five questions by selecting at least one from each module

Module - I

a) You are given a transaction data as shown below from a fast food restaurant. For simplicity we assign the meal items short names [M1-M5].

List of Items Meal- Item List of item Meal- Item {M2, M3} Order 6 {M1, M2, M5} Order 1 (M1, M3) Order 7 (M2, M4) [M1, M2, M3, M5] Order 2 Order 8 (M2, M3) (M1, M2, M3) Order 3 Order 9 {M1, M2, M4} Order 4 {M1, M3} Order 5

For all the min sup 2/9 and min conf=7/9. Apply Apriori and identify all kfrequent itemsets. Find all the strong association rules.

b) Construct the FP-tree for the database above. Consider min sup=2.

10

10

06

10

a) Define maximal and closed frequent itemset identify the above from the database

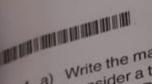
given as Items Transaction ID (A, C, T, W) TI {C, D, W} T2 {A, C, T, W} T3 {A, C, D, W} T4  $\{A, C, D, T, W\}$ T5 {C, D, T} T6

b) Suppose that the data for analysis includes the attribute age. The age values for the data tuples are 13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 25, 25, 30, 33, 33, 35,

Draw the box plot and indicate the (5-No. summary) five number summary. Also indicate the possible outliers







4. 8) (d

BERLENNING.

- Consider a t R3 (4, 2), R R10 (10, 10). that you so C2 = (9, 7).
- c) Write short
- Why is out approache i) Statist
  - ii) Distar
  - b) With the h architectu
  - What do 6. a)
    - b) What is location, 10,000 i resulting
    - c) How do example
  - 7. a) Explain tasks v
    - i) Ho
    - ii) Ve
    - b) With th Archit
    - c) With the
    - 8. a) Write
      - i) M
      - ii) N
      - b) How

- 3. a) Clarify the difference between classification and clustering.
  - b) Explain the significance of 'Naive' used by the Naïve Bayesian Classiflers Briefly describe the difference between Naive Bayes Classification and Bayesian Belief Networks.
  - c) Consider the following training data set in Table 1:

IT 7 (EI) 1 (RC)

| Attribute 1 | Attribute 2 | Attribute 3 | Class   |
|-------------|-------------|-------------|---------|
| A           | 70          | True        | Class 1 |
| A           | 90          | True        | Class 2 |
| A           | 85          | False       | Class 2 |
| A           | 95          | False       | Class 2 |
| A           | 70          | False       | Class 1 |
| В           | 90          | True        | Class 1 |
| В           | 78          | False       | Class 1 |
| В           | 65          | True        | Class 1 |
| В           | 75          | False       | Class 1 |
| С           | 80          | True        | Class 2 |
| С           | 70          | True        | Class 2 |
| С           | 80          | False       | Class 1 |
| С           | 80          | False       | Class 1 |
| C           | 96          | False       | Class 1 |

- i) Calculate the gain on Attribute 1 as Gain (x<sub>1</sub>)
- ii) Explain the main steps in construction of the decision tree using the

THE REAL PROPERTY OF THE REAL PROPERTY.

20

8

12

8 12

4

12

IT 7 - (E-I) 1(RC)

# B.E. (IT) (Semester - VII) (RC) Examination, Nov./Dec. 2016 DATA MINING AND WAREHOUSING

Duration: 3 Hours

Total Marks: 100

Instructions: 1) Attempt any five questions, by selecting atleast one question from each Module.

2) Assume necessary data if required.

### MODULE-I

1. a) Explain how the data transformation and integration steps are carried out in 10

b) Explain concept hierarchy generation for numerical data.

20, 20, 21, 22, 22, 25, 25, 25, 25, 13, 15,16, 16, 19, 30, 33, 33, 35, 35, 35, 35, c) The age values for the data tuples are

Suggest a method of data smoothing that can be used on the above data. Show the working.

2. a) How is data warehouse different from a database? How are they similar?

b) List and elaborate the major issues in data mining.

c) Consider the following data:

| CONSIDER THE |                 |  |  |  |  |  |
|--------------|-----------------|--|--|--|--|--|
| T_ID         | Items Purchased |  |  |  |  |  |
| 101          | a, b, e         |  |  |  |  |  |
| 102          | b, d            |  |  |  |  |  |
| 103          | b, c            |  |  |  |  |  |
| 104          | a, b, d         |  |  |  |  |  |
| 105          | a, c            |  |  |  |  |  |
| 106          | b, c            |  |  |  |  |  |
| 107          | a, c            |  |  |  |  |  |
| 108          | a, b, c, e      |  |  |  |  |  |
| 109          | a, b, e         |  |  |  |  |  |
|              |                 |  |  |  |  |  |

Using apriori algorithm, find candidate itemsets and frequent itemsets. (Assume minimum support count = 2).

5

4

6

### 

HI HI WALL

ts of PDM

vcle.

P

# B.E. (IT) Semester - VII (RC) Examination, Nov./Dec. 2015 DATA MINING AND WARE HOUSING

Total Marks: 100

Duration: 3 Hours

Instruction: Answer any five questions by selecting at least one question from each Module.

### MODULE-1

- a) What is Data Mining? Explain the different steps involved in Knowledge Discovery from Data. b) Given the vectors  $\mathbf{x} = (12, -12, 10, 20, 10, -30)$  and  $\mathbf{y} = (-10, 10, -10, 10, 10, -10)$ ,
  - 6

8

- calculate the proximity between them using the following measures: i) Cosine
  - ii) L<sub>2</sub> norm
  - iii) Tanimoto coefficient.
- c) Given the following data:

4, 18, 15, 21, 22, 24, 24, 25, 25, 27, 28, 34.

Try "smoothing by bin medians" and "smoothing by bin boundaries", on above 6 data with bin size of 4.

2. a) Explain the following types of databases:

- i) Multimedia database
- ii) Relational database
- iii) Object oriented database
- iv) Transactional database.
- b) Differential between Data Mart and Data Warehouse.

IT 7 - (E-I) 1(RC) -3-I MANUAL DISTRICT AND RANGE AND REAL PROPERTY AND AN AREA. MODULE-IV 5 ing a) What is web usage mining ? Explain. 8 b) Explain in brief the following: 6 3 i) Intraquery parallelism. 9 c) Compare ROLAP, MOLAP and hybrid servers with its architecture. 8. a) With the help of a neat diagram explain the shared disk architecture for 8 10 parallel processing. 6 b) Write short notes on : i) Web content Mining. 5 ii) OLAP tools and internet. 2 IS c) List the four types of data partitioning techniques. 6 (5+5)5 10 5 10 10

IT 7 - (E-I) 1(RC)

c) Find the frequent item sets with minimum support count of 3 for the following

| transac            | tional | da  | ta:  |
|--------------------|--------|-----|------|
| [LSLISE.           | -      | -   |      |
| · promotion of the | -      | 4 1 | Mill |

| T ID Bread Milk |       | Detergents | Shampoo | Eggs | SOIL Drink |   |
|-----------------|-------|------------|---------|------|------------|---|
| T_ID            | Bread | MIIK       | 1       | 0    | 0          | 0 |
| 1               | 1     | 1          |         | -    | 1          | 0 |
| 2               | 1     | 0          | 1       |      |            | 0 |
| 3               | 0     | 1          | 1       | 1    | 0          | 1 |
| 4               | 1     | 1          | 1       | 1    | 0          | 0 |
| 5               | 1     | 1          | 1       | 0    | 0          | 1 |

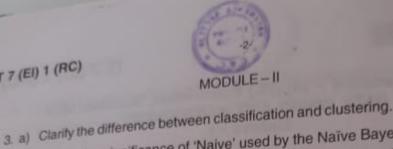
### MODULE-2

- 3. a) Differentiate between Clustering and Classification. Give a brief application for each.

| TI  | ID AGE  | Income | Student | Credit Rating | Class : Buys Compu |
|-----|---------|--------|---------|---------------|--------------------|
|     | 1 Youth | High   | No      | Fair          | No                 |
| 2   | 2 Youth | High   | No      | Excel         | N                  |
| 3   | Mid_Ag  | e High | No      | Fair          | Yes                |
| 4   | Senior  | Medium | No      | Fair          | Yes                |
| 5   | Senior  | Low    | Yes     | Fair          | No                 |
| 6   | Senior  | Low    | Yes     | Excel         | No                 |
| 7   | Mid_age | Low    | Yes     | Excel         | Yes                |
| 8   | Youth   | Medium | No      | Fair          | No                 |
| 9   | Youth   | Low    | Yes     | Fair          | Yes                |
| 0   | Senior  | Medium | Yes     | Fair          | Yes                |
| 1   | Youth   | Medium | Yes     | Excel         | Yes                |
| 2   | Mid_Age | Medium | No      | Excel         | Yes                |
| 3/1 | Mid_Age | High   | Yes     | Fair          |                    |
| 15  | Senior  |        | No      | Excel         | Yes                |

# THE REAL PROPERTY AND ADDRESS OF THE PARTY.

- i) Classify the t Credit Rating
- c) Define Accuracy performance.
- a) Using the data give Income = High,  $L_1$ -norm for 3Nf
  - b) Give the k-med out the drawba
  - c) Which clustering the DBSCAN f
  - a) Is outlier dete the challenge
    - b) Find the outlie assuming tha 17.9, 18.3, 18
    - c) Give the DB outliers?
  - 6. a) Using an exa data.
    - b) Differentiate
    - c) Explain the
  - 7. a) With the h parallel Da
    - b) Explain th
      - i) HOLA
      - ii) ROLA
      - iii) MOLA
    - c) Differenti
    - 8. a) Write she
      - b) Write no



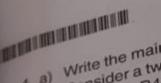
IT 7 (EI) 1 (RC)

Bayesian Belief Networks.



b) Explain the significance of 'Naive' used by the Naïve Bayesian Classifiers

Briefly describe the difference between Naive Bayes Classification and



BEEL MILLEN

Consider a tw 4. 8)

R3 (4, 2), R4 R10 (10, 10). S that you sta C2 = (9, 7).

- c) Write short r
- 5. a) Why is outli approaches i) Statistic
  - ii) Distance
  - b) With the he architectur
  - What do y
    - b) What is a location, it 10,000 ite resulting
    - c) How does example
  - 7. a) Explain tasks wi
    - i) Hori
    - ii) Ver
    - b) With the Archite
    - c) With the
    - 8. a) Write s
      - i) Me
      - ii) Mu
      - b) How is

Table 1 Class Attribute 3 Attribute 2 Attribute 1 Class 1 True 70 A Class 2 True 90 A Class 2 False 85 A Class 2 False 95 A Class 1 False 70 A Class 1 True B 90 Class 1 B 78 False B Class 1 True 65 B 75 False Class 1 C 80 True Class 2 C 70 True Class 2 C 80 False Class 1 C 80 False Class 1 C 96 False Class 1

c) Consider the following training data set in Table 1:

- i) Calculate the gain on Attribute 1 as Gain (x1)
- ii) Explain the main steps in construction of the decision tree using the information in i).

|     |       | IT7-(E-1)  |   | 13    |
|-----|-------|--|---|-------|
|     |       | Classify the tuple X = {age = youth, Income = Medium, Student = Yes, youth, Income = Medium, Student = Yes, and Income = Medium, Student = Yes, youth, Income = Medium, Student = Yes, and Income = Medium, Student = Yes, youth, Income = Yes, youth, Income = Medium, Student = Yes, youth, Income =   |   | E     |
| 17  | No.   | 3 Student  |   | 10    |
| ľ   | 1     | lacome Medium  |   | III Z |
|     | 1111  | Youth, Incom Classific Cla   |   |       |
| 1   |       | tuple X = (age Naive Bay netrics for every type tuple X = (age = Youth,  |   |       |
|     |       | i) Classify the Fair). Cla   |   |       |
|     |       | Define Accuracy, Precional President Classify the data using L2  |   | W     |
|     | c)    | i) Classify the tuple X = {age   Youth, Income   Medium, South, Income   Youth, Income   Medium, South, Income   Youth, Income   |   |       |
|     | a)    | i) Classify the tuple X = {age   Naive Bayesian or evaluating   Credit Rating   Fair}; using Naive Bayesian or evaluating   Youth, Credit Rating   Fair}; using Naive Bayesian or evaluating   Youth, Credit Rating   Fair}; using   X = {Age   Youth, Credit Rating   Youth, Cre   |   |       |
| 4   |       | i) Classify the tuple X Credit Rating = Fair); using Nat.  Credit Rating = Fair); using Nat.  Define Accuracy, Precision and Recall; metrics  Define Accuracy, Precision and Recall; metrics  Define Accuracy, Precision and Recall; metrics  Precision and Recall; metrics  Service Age = Your Age = You   |   |       |
|     | to)   | Give the k-mediods k-means algorithm DB30  |   |       |
|     | 0,    | out the draw expering technique is the clustering techniqu   |   |       |
|     | c)    | Using the data given in Que Yes, Great Using the data given in Que Yes, Great Using the High, Student = Yes, Great Income = = Ye   | 6 |       |
|     |       | which clustering techniques which clusters? the DBSCAN find clusters?  MODULE - 3  MODULE - 3  MODULE - 3  Is outlier detection import in the study of Data Mining. What are the some of the challenges of outlier detection? the challenges of outlier detection? the challenges of outlier detection?  The challenges of outlier detection?  The challenges of outlier detection of the challenges of outlier detection?   |   |       |
|     |       | retaction import in the study of the mean-variance metrics, the mean-variance metrics, 18.0,   | c |       |
| 5   | (a)   | Is outlier detection import in the study of Data Mining. White the challenges of outlier detection? The challenges of outlier detection algorithm. How does it compute assuming that it is normal distributed. The challenges of outlier detection algorithm. How does it compute the challenges of outlier detection algorithm. How does it compute the challenges of outlier detection algorithm. How does it compute the challenges of outlier detection algorithm.  | 6 |       |
|     |       | Find the outlier among the listed to the TMIN = (10.7)   | 8 |       |
|     | b)    | assuming that it is normal assuming that it is not a single assuming the single assuming the single assuming that it is not a single assuming the single assuming the single assuming the single assuming the single assuming that it is not a s   | 0 |       |
|     |       | Find the outlier among the liditaributed. $^{1}$ MIN assuming that it is normal distributed. $^{$ |   |       |
|     | c)    | Give the DB( $r$ , $\pi$ ) – outlier detection Algorithms Give the DB( $r$ , $\pi$ ) – outlier detection Algorithms on multi-dimensional outliers?  Using an example explain atleast four OLAP operations on multi-dimensional data.   | 6 |       |
|     |       | an axample explain atleast four OLAr open  | 6 |       |
| 6.  | a)    | data.  | 8 |       |
|     | hl    | data.  Differentiate between OLTP and Data Warehousing.  Differentiate between on building a Data Warehouse.   |   |       |
|     | C)    | Differentiate between OLTP and Data Warehouse.  Explain the different steps in building a Data Warehouse.  |   |       |
|     |       | MODULE-4   |   |       |
|     |       | Shared Disk Architecture for   |   |       |
| 7   | 2)    | With the help of a neat diagram explain the Shared Disk Architecture for   | 8 |       |
| 1 - |       | parallel Databaso pros   | • | ô     |
|     | b)    | Explain the following:   |   |       |
|     | 10.90 | i) HOLAP   |   |       |
|     |       | ii) ROLAP  |   |       |
|     |       | iii) MOLAP.  |   |       |
|     | c)    | Differentiate between multidimensional OLAP and multi-relational OLAF  |   | 6     |
| 8.  | a)    | Write short note on techniques for Web Usage Mining.   |   | 10    |
|     | b)    | Write notes on any two applications of Data Mining techniques.   |   | 10    |
|     |       |  |   |       |
|     |       |  |   |       |

### ONE AND I WHEN HE HE HERE HE I THE

# IT 7 - (E-I) 1(RC)

## B.E. (IT) Semester - VII (RC) Examination, Nov./Dec. 2015 DATA MINING AND WARE HOUSING

Total Marks: 100

Ouration: 3 Hours

Instruction: Answer any five questions by selecting at least one question from each Module.

### MODULE-1

1. a) What is Data Mining? Explain the different steps involved in Knowledge Discovery from Data.

6

- b) Given the vectors x = (12, -12, 10, 20, 10, -30) and y = (-10, 10, -10, 10, 10, -10). calculate the proximity between them using the following measures:
  - i) Cosine
  - ii) L2 norm
  - iii) Tanimoto coefficient.
- c) Given the following data:

4, 18, 15, 21, 22, 24, 24, 25, 25, 27, 28, 34.

Try "smoothing by bin medians" and "smoothing by bin boundaries", on above 6 data with bin size of 4.

2. a) Explain the following types of databases:

- i) Multimedia database
- ii) Relational database
- iii) Object oriented database
- iv) Transactional database.
- b) Differential between Data Mart and Data Warehouse.

IT 7 (EI) 1 (RC) 4. a) Write the main steps in the algorithm for classification using prediction.

b) Consider a two dimensional database D with the records: D1 (2.2). B2 (2.3). b) Consider a two dimensional database D with the records: R1 (2, 2), R2 (2, 4), R3 (4, 2), R4 (4, 4), R5 (2, 6), R6 (7, 6), R7 (0, 6), R9 (6, 10), R9 (7, 6), R9 (7, 6), R9 (8, 10), R9 (9, 10), R9 (9, 10), R9 (10), R9 (10) Consider a two dimensional database D with the records : H1 (2, 2), R2 (2, 4), R3 (4, 2), R4 (4, 4), R5 (3, 6), R6 (7, 6), R7 (9, 6), R8 (5, 10), R9 (8, 10), R10 (10, 10). Show the results of the k-means algorithm at each step. R3 (4, 2), R4 (4, 4), R5 (3, 6), R6 (7, 6), R7 (9, 6), R8 (5, 10), R9 (0, 10), R10 (10, 10). Show the results of the k-means algorithm at each step, assuming that you start with two clusters (k, 2) with contact (1, 6, 6) and DEED STREET OF REAL PROPERTY. that you start with two clusters (k = 2) with centers C1 = (6, 6) and C2 = (9, 7). Write the main stars you follow 8 7 tering. C2 = (9, 7). Write the main steps you follow. Bayesian Classifiers c) Write short notes on 'Grid Based Clustering'. Classification and 5. a) Why is outlier mining important? Briefly describe the following different 12 (3, approaches used for outlier analysis. i) Statistical-based outlier detection. b) With the help of the neat labeled diagram explain the data warehouse 8 5 6. a) What do you understand by slice and dice? Give an example. b) What is a data warehouse ? Consider a cube used to analyze sales by location, item and date (weekly). Assume data is available for 1000 locations, 10,000 items for 100 weeks. Discuss an efficient scheme to compute the 9 c) How does a snowflake scheme differ from a STAR schema? Explain with an 6 example. Name any two disadvantages of snowflake schema. MODULE-IV 7. a) Explain the following types of parallelism used for parallel execution of the 4 tasks within SQL statements. Horizontal parallelism ii) Vertical parallelism b) With the help of a neat labeled diagram explain the Shared Memory Architecture for parallel processing. c) With the help of a diagram compare the MOLAP v/s ROLAP architectures. 8. a) Write short notes on : tree using the i) Metadata Interchange Initiative ii) Multidimensional data model. b) How is mining for user behavior on the web carried out?

IT 7- (E-I) 1(RC)

c) Find the frequent item sets with minimum support count of 3 for the following

| transa | ctional u | T.   | Detergents | Shampoo | Eggs | Soft Drink |
|--------|-----------|------|------------|---------|------|------------|
| T ID   | Bread     | Milk | Detergents | Shampoo | -99" | - Think    |
| 1      | 1         | 1    | 1          | 0       | 0    | 0          |
| 2      | 1         | 0    | 1          | 1       | 1    | 0          |
| 3      | 0         | 1    | 1          | 1       | 0    | 1          |
| 4      | 1         | 1    | 1          | 1       | 0    | 0          |
| 5      | 1         | 1    | 1          | 0       | 0    | 1          |

### MODULE-2

- 3. a) Differentiate between Clustering and Classification. Give a brief application for each.

| ľ   | TID AGE |         | Incom   | e Studen | Credit Rating | Class : Buys Com |  |
|-----|---------|---------|---------|----------|---------------|------------------|--|
|     | 1       | Youth   | High    | No       | Fair          | No -             |  |
|     | 2       | Youth   | High    | No       | Excel         | N                |  |
|     | 3       | Mid_Ag  | ge High | No       | Fair          | Yes              |  |
| L   | 4       | Senior  | Mediun  | n No     | Fair          | Yes              |  |
| L   | 5       | Senior  | Low     | Yes      | Fair          | No               |  |
|     | 6       | Senior  | Low     | Yes      | Excel         | No               |  |
| 100 | 7       | Mid_age | Low     | Yes      | Excel         | Yes              |  |
| 8   | 3       | Youth   | Medium  | No       | Fair          | No               |  |
| 9   | 1       | Youth   | Low     | Yes      | Fair          | Yes              |  |
| 10  | )   5   | Senior  | Medium  | Yes      | Fair          |                  |  |
| 11  | Y       | outh    | Medium  | Yes      |               | Yes              |  |
| 12  | M       | id_Age  | Medium  | No       | Excel         |                  |  |
| 3   |         | d_Age   |         | INO      | Excel         | Yes              |  |
|     | -       | 1000    | High    | Yes      | air           |                  |  |
| 1   | Senior  |         | Medium  | No       | Excel         |                  |  |
|     |         |         |         |          | -1001         | No               |  |

# MATERIAL PROPERTY.

THE RESIDENCE OF THE PARTY.

- i) Classify the tuple X = Credit Rating = Fair);
- c) Define Accuracy, Precisi performance.
- a) Using the data given in Qu Income = High, Student L1-norm for 3NN classifi
  - b) Give the k-mediods Algo out the drawback of k-m
  - c) Which clustering technic the DBSCAN find cluste
- 5. a) Is outlier detection imp the challenges of outlie
  - b) Find the outlier among to assuming that it is nor 17.9, 18.3, 18.4, 18.5,
  - c) Give the DB(r, π) οι outliers?
- 6. a) Using an example exp data.
  - b) Differentiate between
  - c) Explain the different :
- 7. a) With the help of a ne parallel Database pr
  - b) Explain the following
    - i) HOLAP
    - ii) ROLAP
    - iii) MOLAP.
  - c) Differentiate between
- 8. a) Write short note on
  - b) Write notes on any

ris fail?

he ERP

mentation

tem

in



# IT 7 (EI) 1 (RC)

B.E. (IT) (Semester - VII) (RC) Examination, Nov./Dec. 2014

DATA MINING AND WAREHOUSES

Elective - I

ration: 3 Hours Instructions: 1) Attempt any five questions by selecting atleast one from each Module. 2) Assume suitable data if necessary. MODULE-I 1. a) Define: i) Data Mining ii) Knowledge Discovery State the difference between the two. Draw a complete labeled diagram of a 8 typical data mining system. 6 b) The age values for the data tuples are: 20, 20, 21, 22, 22, 25, 25, 25, 25, 13, 15, 16, 16, 19, 30, 33, 33, 35, 35, 35, 35, 36, 40, 45, 46, 52, 70. Use min-max normalization to transform the values 20, 30, 40 and 70 in the range [0.0, 1.0]. c) Write the four factors which test the interestingness of the patterns. 2 d) What is Iceberg Query? Give its general syntax. What is association Rule Mining? What are the two steps in the process? 2 a) 6 Explain. b) Suppose that we have sales data given by Address and the Address fields include House Number, Street Name, City, State, Pincode and Country. Write a DMQL statement for expressing the concept hierarchy. 4

c) Suppose a group of 12 sales price records has been sorted as follows:

ii) Equi-width partitioning.

Partition them into three bins by each of the following methods.

5, 10, 11, 13, 15, 35, 50, 55, 72, 92, 204, 215.

d) What is numerosity reduction using regression?

i) Equidepth partitioning

P.T.O.

2