2053
B.E. (Computer Science and Engineering)
Sixth Semester
Elective – I
CS-605C: Data Mining and Analysis

Time allowed: 3 Hours

Max. Marks: 50

NOTE: *Attempt* <u>*five*</u> *questions in all, including Question No. 1 (Section-A)    which is compulsory and selecting two questions each from Section B-C.*

x-x-x

## Section -A

Q 1(a)  What is Schema Integration?                                                                              (10)

(b)  Differentiate between t-weights and d-weights.

(c)  What are the metrics used to compare the performance of Classification Algorithm?

(d)  What is data discretization?

(e)  What is difference between time series and sequence database?

## Section -B

Q2 (a)  Illustrate the differences between ROLAP, MOLAP and HOLAP. Which architecture is preferred in large   (5)
organization

(b)  What is partial materialization of data warehouse? Why Full materialization is not feasible in large   (5)
dimensional data? Explain by taking a suitable example.

Q3 (a)  Following are the number of customers visited the store in last 30 days,                                 (5)

250, 350, 415, 200; 230, 420, 500, NaN, 450, 50, 340, 260, 90, 470, 530, 60, 380, 440, 560, NaN. 420, 310,
170, 190, 290, 470, 320, 450, 310, 220

Pre-process the data and find first and third quartile of data, show the boxplot and quantile plot, divide them
into equal size 3 bins and smooth these by boundary.

(b)  What is the purpose of Aggregate Fact table? What are the advantages of using these?                       (5)

Q 4 (a)  How data cubes are computed? Describe the different operations that can be performed on data cubes.     (5)

(b)  What are different Data summarization approaches? Explain analytical characterization of data.             (5)

## Section -C

Q5  Describe Apriori Algorithm in detail. How can we reduce the complexity of the algorithm? Explain the   (10)
Support and Confidence metric and its use in Algorithm.

Q6 (a)  Describe the Decision tree algorithm. Explain how Gini Index can be used  to compare the attributes.   (5)

(b)  How lazy learners are different. Explain the k-NN algorithm. How the choice of k affects the outcome of   (5)
algorithm.

Q7 (a)  What are the main advantages of Partition around mediods algorithm over k-means algorithm? How it can be   (5)
further improved?

(b)  What are multimedia databases? Explain the Multimedia mining approaches using an example.               (5)

x-x-x