## School of Information Technology and Engineering

### Winter Semester 2022-2023 - Fresher

### Continuous Assessment Test – I

**Programme Name & Branch: MCA**

**Course Name & code:**  Data Mining and Business Intelligence (ITA5007)

**Class Number (s): 0528, 0296, 0530**

**Slot:**  C2+TC2

**Faculty Name (s): Harshita PateL, Dr. Ephzibah E.P., Jagadeesan S.**

**Exam Duration: 90 Min.**                                   **Maximum Marks: 50**

| Q.No. | Question | Max Marks |
|---|---|---|
| 1. | There is a strong linkage between statistical data analysis and data mining. Some people think of data mining as an automated and scalable method for statistical data analysis. Do you agree or disagree with this perception? Present one statistical analysis method that can be automated and/or scaled up nicely by integration with the present data mining methodology. | 10 |
| 2. | Briefly outline how to compute the dissimilarity between objects described by the following:<br>(a) Nominal attributes<br>(b) Asymmetric binary attributes<br>(c) Numeric attributes<br>(d) Term-frequency vectors | 10 |
| 3. | Use these methods to normalize the following group of data:<br>200, 300, 400, 600,1000<br>(a) min-max normalization by setting min = 0 and max = 1<br>(b) z-score normalization<br>(c) normalization by decimal scaling | 10 |
| 4. | Suppose that a hospital tested the age and body fat data for 18 randomly selected adults with the following results: (table below)<br>(a) Calculate the mean, median, and standard deviation of age and %fat.<br>(b) Draw the boxplots for age and %fat.<br>(c) Draw a scatter plot and a q-q plot based on these two variables. | 10 |
| 5. | Consider the following data (in increasing order) for the attribute age: 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.<br>(a) Use smoothing by bin means to smooth these data, using a bin depth of 3. Illustrate your steps. Comment on the effect of this technique for the given data.<br>(b) How might you determine outliers in the data?<br>(c) What other methods are there for data smoothing? | 10 |

Table for Question 4:

| age | 23 | 23 | 27 | 27 | 39 | 41 | 47 | 49 | 50 |
|---|---|---|---|---|---|---|---|---|---|
| %fat | 9.5 | 26.5 | 7.8 | 17.8 | 31.4 | 25.9 | 27.4 | 27.2 | 31.2 |

| age | 52 | 54 | 54 | 56 | 57 | 58 | 58 | 60 | 61 |
|---|---|---|---|---|---|---|---|---|---|
| %fat | 34.6 | 42.5 | 28.8 | 33.4 | 30.2 | 34.1 | 32.9 | 41.2 | 35.7 |