GURU GHASIDAS VISHWAVIDYALAY, CENTRAL UNIVERSITY, BILASPUR

**BSC VI END SEMESTER EXAM 2022 (Online)**

**INSTRUCTION : On the top of the answer sheet write Date, your Name, Class, Semester, Name of Subject, University Roll Number.** <span style="color:red">**Total pages of the answer sheet strictly limited to 4sheets only.**</span>

<span style="color:red">**One mark will be deducted per minute in case of late submission**</span>

Class : BSc-VI SEM                                                                 Mark : 70

Subject:  Data Mining                                                          Time : 10am-12noon

Answer any fourteen(14). Calculator is allowed.

| Q. No. | | Mark |
|---|---|---|
| 1 | Differentiate between cluster analysis, outlier analysis and evolution analysis. | 5 |
| 2 | What are the major challenges of mining a large amount of data in comparison with mining a small amount of data? | 5 |
| 3 | Suppose that a data warehouse for GGU University consists of the following four dimensions: student, course, semester and instructor and two measures count and avg_grade. When at the lowest conceptual level the avg_grade measure stores the actual course grade of the student. At higher conceptual levels, avg_grade stores the average grade for the given combination. <br> (a)Draw a snow-flake schema for the data warehouse | 5 |
| 4 | For the above drawn schema of Q. 3  write the DMQL | 5 |
| 5 | Explain the indexing techniques used in data warehouse with suitable examples. | 5 |
| 6 | 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 <br> (a)Use min-max normalization to transform the value 35 for age onto the range [0.0 1.0] <br> (b)Use z-score normalization to transform the value 35 for age, where the standard deviation of age is 12.94. | 5 |
| 7 | Find the IDFT value of $X(k) = \{2, -j, 0, j\}$ | 5 |
| 8 | How  dimensions of data can be reduced using transform based techniques? | 5 |
| 9 | A 2x2 contingency table is given as follows : <br> Are gender and preferred_reading correlated? Check by Chi-square test. Given <br> $value\ of\ \chi^2 = 10.828\ for\ \alpha = 0.001\ and\ \deg ree\ of\ freedom = 1$ <table><tr><td></td><td>Male</td><td>Female</td><td>total</td></tr><tr><td>Fiction</td><td>50</td><td>1000</td><td>1050</td></tr><tr><td>Non-fiction</td><td>250</td><td>200</td><td>450</td></tr><tr><td>Total</td><td>300</td><td>1200</td><td>1500</td></tr></table> | 5 |
| 10 | Find the dissimilarity matrix for the following mixed types variables | 5 |

| Object | Test-1 (Categorical) | Test-2 (Ordinal) | Test-3 (Ratio scaled) |
|---|---|---|---|
| 1 | Code-A | Excellent | 445 |
| 2 | Code-B | Fair | 22 |
| 3 | Code –C | Good | 164 |
| 4 | Code-A | Excellent | 1210 |

| 11 | What are the different training methods available for classification or prediction? | 5 |
|---|---|---|
| 12 | Explain the k-mean clustering algorithm for IRIS data set. | 5 |
| 13 | Plot the dendogram using agglomerative algorithm of hierarchical clustering with single linkage for the following points: <br><br> A 3 3.5 <br><br> B 4 4 <br><br> C 3 4 <br><br> D 5 5 <br><br> E 1.5 1.5 <br><br> F 1 1 | 5 |
| 14 | Write the APRIORI Algorithm of association rule mining. | 5 |
| 15 | Explain the artificial neural network based time series prediction. | 5 |
| 16 | How to choose the Root node in decision tree based classifier? | 5 |
| 17 | Using the below given table and Naïve Bayesian Classification algorithm find out the class of the new tuple, <br> X={age=senior, income=medium, student=yes, credit rating =fair} | 5 |
| 18 | Suppose that the data mining task is to cluster the following nine points into three clusters. <br> $A_1(2, 10)$, $A_2( 2, 5)$, $A_3(8,4)$, $B_1(5, 8)$, $B_2(7, 5)$, $B_3(6, 4)$, $C_1(1, 2)$, $C_2( 4, 9)$, $C_3(3,4)$ <br> The distance function is Euclidean distance. <br> Suppose initially we assign $A_2$, $B_2$ and $C_2$ as the center of each cluster respectively. Use the k-means algorithm to show the three cluster centers after the first round of execution | 5 |

For question 17:

| RID | Age | Income | Student | Credit rating | Class : Buys computer |
|---|---|---|---|---|---|
| 1 | Youth | High | No | Fair | No |
| 2 | Youth | High | No | Excellent | No |
| 3 | Middle aged | High | No | Fair | Yes |
| 4 | Senior | Medium | No | Fair | Yes |
| 5 | Senior | Low | Yes | Fair | Yes |