



Mid Semester Examination

Oct 2023

CET2012B - Data Engineering Concepts

Schedule ID: 21573

| | | | |
|----------------|---------------------------------------|------------|--------------------|
| Faculty/School | Faculty of Engineering and Technology | Term | Semester V |
| Program | TY BTech CSE | Duration | 1 Hours 30 Minutes |
| Specialization | | Max. Marks | 50 |

Instructions to the Candidate:

1. Write the PRN on the top right-hand corner of the question paper.
2. Draw neat diagrams.
3. Assume suitable data, if necessary.
4. Solve any 5 questions.

Section 1 (5 X 10 Marks)

Answer any 5 questions

| | | | | |
|---|---|----------------|-------------|-------------|
| 1 | <p><input checked="" type="radio"/> Explain the difference between structured data and unstructured data. Given a set of unstructured data, propose a technique to convert it into structured data.</p> <p><input checked="" type="radio"/> Apply the concept of data pre-processing to a real-world dataset by listing the specific steps involved in preparing the data for analysis.</p> | 10 marks 10 | CO1 | Applying |
| 2 | <p><input checked="" type="radio"/> Explain the characteristics and applications of time series data.</p> <p><input checked="" type="radio"/> Summarize the role of software engineer, data engineer and a data scientist in IT industry.</p> | 10 marks 10 | CO1, CO2 | Remembering |
| 3 | <p>a. Describe the Data Engineering Lifecycle Process with the help of a diagram. Illustrate the process with the help of a real-time example.</p> <p><input checked="" type="radio"/> Differentiate between discrete data and continuous data. Give any real world example of each type.</p> | 10 marks 5 | CO1, CO2 | Applying |
| 4 | <p><input checked="" type="radio"/> Analyse various reasons of having missing values in a dataset. Explain any one technique to identify and rectify it.</p> <p><input checked="" type="radio"/> Elaborate the various data integration techniques in data warehouse.</p> | 10 marks 5 | CO2 | Analysing |

| | | | | | | | | | | | | | | | | | | | | | | |
|--------------|---|--------------|-----|-----|-----|------|-----|-----------|-----------|----|----|-----|-----|------------|---|-------------------|-----|------------|----|-------------------|-------------|------------|
| 5 | <p>a Evaluate the binning technique on following dataset using equi-depth, means and boundaries methods using bin depth of 4.</p> <p>The values for data sets are (in ascending order) 13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 25, 25, 30, 33.</p> <p>b For the following set of data</p> <ol style="list-style-type: none">1. Draw a scatter plot2. Calculate Pearson product moment correlation coefficient3. Asses the statistical significance of your value and interpret your results <table border="1"><tr><td>Students</td><td>A</td><td>B</td><td>C</td><td>D</td><td>E</td></tr><tr><td>DEC Marks</td><td>1</td><td>3</td><td>6</td><td>10</td><td>12</td></tr><tr><td>AIES Marks</td><td>5</td><td>13</td><td>25</td><td>41</td><td>49</td></tr></table> | Students | A | B | C | D | E | DEC Marks | 1 | 3 | 6 | 10 | 12 | AIES Marks | 5 | 13 | 25 | 41 | 49 | 10 marks 6 | CO1, CO2 | Evaluating |
| Students | A | B | C | D | E | | | | | | | | | | | | | | | | | |
| DEC Marks | 1 | 3 | 6 | 10 | 12 | | | | | | | | | | | | | | | | | |
| AIES Marks | 5 | 13 | 25 | 41 | 49 | | | | | | | | | | | | | | | | | |
| 6 | <p>a The Table shows the distribution of interest paid to bank investors in a particular year. Draw a histogram to illustrate the data.</p> <table border="1"><tr><td>Interest(\$)</td><td>25-</td><td>30-</td><td>40-</td><td>60-</td><td>80-</td><td>110-</td></tr><tr><td>Frequency</td><td>18</td><td>55</td><td>140</td><td>124</td><td>96</td><td>0</td></tr></table> <p>b. Explain following sampling techniques:</p> <ol style="list-style-type: none">1. Random2. Systematic3. Stratified4. Cluster <p>Suppose there are 200 boys and 100 girls in 'Global Talent' school and you need to draw a sample of 60 students. Which sampling techniques could be chosen, explain how this could be done?</p> | Interest(\$) | 25- | 30- | 40- | 60- | 80- | 110- | Frequency | 18 | 55 | 140 | 124 | 96 | 0 | 10 marks 5 | CO2 | Evaluating | | | | |
| Interest(\$) | 25- | 30- | 40- | 60- | 80- | 110- | | | | | | | | | | | | | | | | |
| Frequency | 18 | 55 | 140 | 124 | 96 | 0 | | | | | | | | | | | | | | | | |