

Class : MCA-V SEM.

Mark : 60

Subject : Data Mining and Data Warehousing(MCA-503)

Time : 3 Hours

Q. 1 is compulsory. Answer any four from rest.

Q.1	<p>(a) What is data mining? Is it a simple transformation of technology developed from databases, statistics and machine learning?</p> <p>(b) Name the databases on which data mining techniques can be applied.</p> <p>(c) What is a time series database. Give an example.</p> <p>(d) What is the meaning of Market-Basket Analysis? Give an example.</p> <p>(e) Differentiate between classification and prediction</p> <p>(f) Differentiate between data ware house and data mart.</p> <p>(g) What is DMQL? How it is different from SQL.</p> <p>(h) Define no coupling, loose coupling, semitight coupling and tight coupling</p> <p>(i) How dimensionality reduction is different from data reduction?</p> <p>(j) What is holistic measure? Give examples.</p>	2 x 10																																								
Q. 2	<p>(a) Explain how the evolution of database technology led to Data Mining.</p> <p>(b) Describe the steps involved in data mining when viewed as a process of knowledge discovery.</p>	5+5																																								
Q. 3	<p>Briefly compare the following concepts</p> <p>(a) Snowflake schema, fact constellation, star query model</p> <p>(b) Data cleaning, data transformation, refresh</p> <p>(c) Enterprise warehouse, data mart</p>	3+3+4																																								
Q. 4	<p>Suppose that the values for a given set of data are grouped into intervals. The intervals and corresponding frequencies are as follows</p> <table data-bbox="381 1312 609 1554"> <tr> <td>Age</td> <td>Frequency</td> </tr> <tr> <td>1-5</td> <td>200</td> </tr> <tr> <td>5-15</td> <td>450</td> </tr> <tr> <td>15-20</td> <td>300</td> </tr> <tr> <td>20-50</td> <td>1500</td> </tr> <tr> <td>50-80</td> <td>700</td> </tr> <tr> <td>80-110</td> <td>44</td> </tr> </table> <p>(a) Compute an approximate median value for the data.</p> <p>(b) Plot a histogram using the above data.</p>	Age	Frequency	1-5	200	5-15	450	15-20	300	20-50	1500	50-80	700	80-110	44	5+5																										
Age	Frequency																																									
1-5	200																																									
5-15	450																																									
15-20	300																																									
20-50	1500																																									
50-80	700																																									
80-110	44																																									
Q. 5	<p>Suppose a hospital tested the age and body fat data for 18 randomly selected adults with the following result</p> <table data-bbox="289 1701 1258 1848"> <tr> <td>Age</td> <td>23</td> <td>23</td> <td>27</td> <td>27</td> <td>39</td> <td>41</td> <td>47</td> <td>49</td> <td>50</td> </tr> <tr> <td>%fat</td> <td>9.5</td> <td>26.5</td> <td>7.8</td> <td>17.8</td> <td>31.4</td> <td>25.9</td> <td>27.4</td> <td>27.2</td> <td>31.2</td> </tr> <tr> <td>Age</td> <td>52</td> <td>54</td> <td>54</td> <td>56</td> <td>57</td> <td>58</td> <td>58</td> <td>60</td> <td>61</td> </tr> <tr> <td>%fat</td> <td>34.6</td> <td>42.5</td> <td>28.8</td> <td>33.4</td> <td>30.2</td> <td>34.1</td> <td>32.9</td> <td>41.2</td> <td>35.7</td> </tr> </table> <p>(a) Calculate the mean, median and standard deviation of age and %fat</p> <p>(b) Draw the boxplots for age and %fat</p>	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	4+3+3
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																																	

	(c) Draw a scatter plot based on these two variables.																													
Q.6	<p>Suppose a group of 12 sales price records has been sorted as follows: 5, 10, 11, 13, 15, 35, 50, 55, 72, 92, 204, 215</p> <p>Partition them into three bins by each of the following methods :</p> <p>(a)equal-frequency partitioning (b)equal-width partitioning (c) Clustering</p>	4+3+3																												
Q.7	<p>(a)What is attribute subset selection? Why it is important? What are the different methods available for doing attribute subset selection?</p> <p>(b)What is curse of dimensionality? Discuss the different methods available for dimensionality reduction?</p>	5+5																												
Q.8	<p>(a)Write an algorithm for k-nearest neighbor classification.</p> <p>(b)From the following table shows the midterm and final exam grades obtained for students in a database course</p> <table style="margin-left: 40px;"> <thead> <tr> <th style="text-align: left;">X</th> <th style="text-align: left;">y</th> </tr> <tr> <th style="text-align: left;">Midterm Exam</th> <th style="text-align: left;">Final Exam</th> </tr> </thead> <tbody> <tr><td>72</td><td>84</td></tr> <tr><td>50</td><td>63</td></tr> <tr><td>81</td><td>77</td></tr> <tr><td>74</td><td>78</td></tr> <tr><td>94</td><td>90</td></tr> <tr><td>86</td><td>75</td></tr> <tr><td>59</td><td>49</td></tr> <tr><td>83</td><td>79</td></tr> <tr><td>65</td><td>77</td></tr> <tr><td>33</td><td>52</td></tr> <tr><td>88</td><td>74</td></tr> <tr><td>81</td><td>90</td></tr> </tbody> </table> <p>(i)Use the method of least squares to find an equation for the prediction of a student's final exam grade based on the student's midterm grade in the course. (ii) Predict the final exam grade of a student who received an 86 on the midterm exam.</p>	X	y	Midterm Exam	Final Exam	72	84	50	63	81	77	74	78	94	90	86	75	59	49	83	79	65	77	33	52	88	74	81	90	5+5
X	y																													
Midterm Exam	Final Exam																													
72	84																													
50	63																													
81	77																													
74	78																													
94	90																													
86	75																													
59	49																													
83	79																													
65	77																													
33	52																													
88	74																													
81	90																													