



	A	B	C	D
a)	i	ii	iii	iv
b)	ii	iii	iv	i
c)	iv	iii	ii	i
d)	ii	iv	iii	i

5. Assertion (A): The data mining environment is the part of an organization whose core competency is data mining [K<sub>2</sub>]  
Reason (R): Organization have different sources of business intelligence data  
a) both A and R are individually true and R is the correct explanation of A      b) both A and R are individually true but R is not the correct explanation of A  
c) A is true but R is false      d) A is false but R is true.
6. Which one of the following is not correct statement? [K<sub>2</sub>]  
a) Entropy is one of the dimensionality reduction technique      b) Attribute subset selection reduces the data set size by removing irrelevant or relevant attributes  
c) Smoothing works to remove noise from the data      d) Histograms use binning to approximate data distribution and are a popular form of data reduction
7. The basic features of clustering are \_\_\_\_\_ [K<sub>2</sub>]  
i) The best number of clusters is not known  
ii) There may not be any a priori knowledge concerning the clusters  
iii) Cluster results are dynamic  
a) i and ii only      b) i and iii only  
c) ii and iii only      d) i, ii and iii
8. Assertion (A) : Data warehouse refers to a data repository that is maintained separately from an organization's operational database [K<sub>2</sub>]  
Reason (R) : Data warehouse stores information from present data.  
a) both A and R are individually true and R is the correct explanation of A      b) both A and R are individually true but R is not the correct explanation of A  
c) A is true but R is false      d) A is false but R is true.
9. Document Clustering Analysis is one of the \_\_\_\_\_ approaches. [K<sub>2</sub>]  
a) Text Mining      b) WWW Mining  
c) Spatial data Mining      d) Multimedia data Mining
10. Sequence the operations for cleansing the data [K<sub>2</sub>]  
(i) Locates and identifies individual data elements in the source files and then isolates these data elements in the target files  
(ii) Transform data into its consistent format using both standard and custom business rules  
(iii) Corrects individual data components using sophisticated data algorithms and secondary data sources  
(iv) Searching and matching records within and across the standardized data, based on predefined business rules to eliminate duplications



**PART D (4 x 10 = 40 Marks)**

27. Write Apriori algorithm and using the algorithm find all frequent itemsets for the following database (min\_sup=20%) [K<sub>4</sub>]

List of item IDs									
Transaction ID	A1	A2	A3	A4	A5	A6	A7	A8	A9
T1	1	0	0	0	1	1	0	1	0
T2	0	1	0	1	0	0	0	1	0
T3	0	0	0	1	1	0	1	0	0
T4	0	0	1	0	0	0	0	0	0
T5	0	0	0	0	1	1	1	0	0
T6	0	1	1	1	0	0	0	0	0
T7	0	1	0	0	0	1	1	0	1
T8	0	0	0	0	1	0	0	0	0

28. Write the advantages and disadvantages of k-means clustering as against model-based clustering. For the data set {2, 3, 5, 8, 13, 21, 34, 55, 89, 144, 233, 377}, use the following techniques to form two clusters. [K<sub>4</sub>]
- (i) k-means with initial centroids {1} and {378}
  - (ii) k-means with initial centroids {21} and {34}
29. Discuss and elaborate the current trends in data mining. Apprise the applications of data mining for financial data analysis. [K<sub>3</sub>]
30. Explain the three tier data warehousing architecture with a neat diagram. [K<sub>3</sub>]

\*\*\*\*\*