|    | **Unit 1** |
|----|-----------|
| 1  | What is data mining and data warehouse? |
| 2  | Compare data base processing Vs. data mining processing. |
| 3  | Explain applications of data mining in detail. |
| 4  | Explain all data mining models and tasks. |
| 5  | What is KDD? Explain with diagram. |
| 6  | Write a short note on visualization. |
| 7  | Discuss the issues in Data mining. |
| 8  | What is fuzzy logic? Explain in brief with example. |
| 9  | Define following terms:<br>1. Information retrieval<br>2. Precision<br>3. Recall<br>4. Similarity<br>5. Granularities<br>6. Facts<br>7. Roll ups<br>8. Drill down |
| 10 | Define cube and explain with example. |
| 11 | Write a short note on star schema. |
| 12 | Explain characteristics of data warehouse. |
| 13 | Discuss the ways to improve the performance of data warehouse applications. |
| 14 | Write a short note on OLAP operations. |
| 15 | Write a short note on point estimation. |
| 16 | Compute mean, variance and standard deviation for (1,3,4,6,5). |
| 17 | Define<br>1. Mean<br>2. Median<br>3. Mode<br>4. Variance<br>5. Standard deviation (Sample/Population)<br>6. Bias<br>7. MSE |
| 18 | Compute mean, median and mode for (15, 10, 18, 20, 28, 32). |
| 19 | What is Jackknife estimate technique? |
| 20 | Find out Jackknife estimate for variance X={1, 5, 6}<br>Mean X= {5, 6, 6}. |
| 21 | Estimate P that maximizes the likelihood that the given sequence of heads and tails would occur for {H, H, H, T, T} Note: Assume coin with H and T equally likely. |

| 22 | Estimate the missing data and continues until convergence using Expectation Maximization $\{1, 5, 10, 4, *, *\}$. (Guess $\mu^0=3$) |
|---|---|
| 23 | Prove that $X_{11}$ belongs to class $h_2$ using Bayes theorem. |

| ID | Income | Credit | Class | $x_i$ |
|---|---|---|---|---|
| 1 | 4 | Excellent | $h_1$ | $x_4$ |
| 2 | 3 | Good | $h_1$ | $x_7$ |
| 3 | 2 | Excellent | $h_1$ | $x_2$ |
| 4 | 3 | Good | $h_1$ | $x_7$ |
| 5 | 4 | Good | $h_1$ | $x_8$ |
| 6 | 2 | Excellent | $h_1$ | $x_2$ |
| 7 | 3 | Bad | $h_2$ | $x_{11}$ |
| 8 | 2 | Bad | $h_2$ | $x_{10}$ |
| 9 | 3 | Bad | $h_3$ | $x_{11}$ |
| 10 | 1 | Bad | $h_4$ | $x_9$ |

| 24 | Write a short note on Hypothesis testing. |
|---|---|
| 25 | Find Chi square statistics for Observed value = $\{51, 95, 67, 78, 88\}$ Expected value=76 |
| 26 | Write a short note on linear regression. |
| 27 | Write a short note on non-linear regression. |
| 28 | Explain correlation in detail. |
| 29 | Find correlation between Ice cream sales Vs temperature |

| Temperature | Ice Cream Sales (in |
|---|---|
| 14.2 | 215 |
| 16.4 | 325 |
| 11.9 | 185 |
| 15.2 | 332 |
| 18.5 | 406 |
| 22.1 | 522 |

| 30 | Write a short on similarity measures. |
|---|---|
| | **Unit 2** |
| 31 | Explain the need of data pre-processing. |
| 32 | List and explain major tasks in data processing. |
| 33 | Explain terms Quartile and Inter-Quartile range. |
| 34 | What are Box plot and Quantile plot? |
| 35 | What is histogram and scatter plot? |
| 36 | Write a short note on data cleaning tasks. |
| 37 | Explain Binning with example. |
| 38 | Explain Data aggregation, generalization and smoothing. |
| 39 | Write a short note on data transformation. |
| 40 | Write a short note on data normalization. |
| 41 | Define 1. Association rule |

| | | | |
|---|---|---|---|
| | | 2. Support<br>3. Confidence | |
| w) | 42 | Explain apriori algorithm with example. | |
| | 42a | Write a short note on Association rule mining. | |
| | | **Unit 3** | |
| | 43 | What is classification? Discuss the issues. | |
| | 44 | What is prediction? Discuss the issues. | |
| | 45 | Write a short note on decision tree. | |
| | 46 | Write a short note on Bayesian classifier. | |
| | 47 | Write a short note on Rule based classifier. | |
| | 48 | Write a short note on Neural network classifier. | |
| | 49 | Write a short note on Support Vector Machine. | |
| | 50 | Define coverage and accuracy in rule based classifier. | |
| | 51 | Explain triggering and firing of rules. | |
| | 52 | Explain rule based and class based ordering. | |
| | 53 | Discuss "The accuracy on its own is not a reliable estimate of rule | |
| | 54 | Consider a training set that contains 100 positive examples and 400<br>negative examples for each of the following<br>candidate rule. R1 : A $\longrightarrow$ + (covers 4 positive and<br>one negative examples) R2 : B $\longrightarrow$ + (covers 30<br>positive and 10 negative examples) R3 : C $\longrightarrow$ + | |
| | 55 | Consider a training set that contains 100 positive examples and 400<br>negative examples for each of the following<br>candidate rule. R1 : A $\longrightarrow$ + (covers 4 positive and<br>one negative examples) R2 : B $\longrightarrow$ + (covers 30<br>positive and 10 negative examples) R3 : C $\longrightarrow$ + | |
| | 56 | Consider a training set that contains 100 positive examples and 400<br>negative examples for each of the following<br>candidate rule. R1 : A $\longrightarrow$ + (covers 4 positive and<br>one negative examples) R2 : B $\longrightarrow$ + (covers 30<br>positive and 10 negative examples) R3 : C $\longrightarrow$ + | |
| | 57 | Write the difference between classification and clustering. | |
| | 58 | Explain supervised and unsupervised learning. | |
| | 59 | Explain term pruning and overfitting. | |
| | 60 | Find information gain for "income" in following data | |

| 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 |
| 6 | senior | low | yes | excellent | no |
| 7 | middle_aged | low | yes | excellent | yes |
| 8 | youth | medium | no | fair | no |
| 9 | youth | low | yes | fair | yes |
| 10 | senior | medium | yes | fair | yes |
| 11 | youth | medium | yes | excellent | yes |
| 12 | middle_aged | medium | no | excellent | yes |
| 13 | middle_aged | high | yes | fair | yes |
| 14 | senior | medium | no | excellent | no |

| | |
|---|---|
| 61 | Write a short note on Gini index. |
| 62 | Classify the following tuple using Naïve Bayesian classifier. X=(age=youth, income=low , student=yes, credit_rating=fair) using following training data. |

| 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 |
| 6 | senior | low | yes | excellent | no |
| 7 | middle_aged | low | yes | excellent | yes |
| 8 | youth | medium | no | fair | no |
| 9 | youth | low | yes | fair | yes |
| 10 | senior | medium | yes | fair | yes |
| 11 | youth | medium | yes | excellent | yes |
| 12 | middle_aged | medium | no | excellent | yes |
| 13 | middle_aged | high | yes | fair | yes |
| 14 | senior | medium | no | excellent | no |

| 63 | Find out the population for the year 2013 using linear regression. |
|----|--------------------------------------------------------------------|

| 2005 | 2006 | 2007 | 2008 | 2009 | 2013 |
|------|------|------|------|------|------|
| 12 | 19 | 28 | 35 | 45 | ? |

| 64 | Write a short note on confusion matrix. |
|----|-----------------------------------------|

| 65 | Classify $X_1=4$, $X_2=7$ using K-nearest neighbour (assume k=3). |
|----|------------------------------------------------------------------|

| X | Y | Class |
|---|---|-------|
| 7 | 7 | B |
| 7 | 4 | B |
| 3 | 4 | G |
| 1 | 4 | G |

**Unit 4**

| 66 | List all the requirements of clustering Data mining. |
|----|------------------------------------------------------|
| 67 | Write a short note on type of data in clustering analysis. |
| 68 | Compute Euclidean and Manhattan distance for $X_1$ (1, 2) and $X_2$ (3,6). |
| 68 a | Compute Euclidean and Manhattan distance for $X_1$ (1, 2) and $X_2$ (4,6). |
| 69 | Compute<br>1. Similarity between A and B<br>2. Similarity between C and B<br>3. Similarity between A and C and comment on the most similar tuples. |

| Name | Gender | F | C | T1 | T2 | T3 |
|------|--------|---|---|----|----|----|
| A | F | Y | N | P | P | N |
| B | M | Y | N | Y | N | P |
| C | F | Y | P | P | N | N |

| 70 | Write a short note on K-means clustering. | |
|----|-------------------------------------------|---|
| 71 | Write a short note on K-medoids clustering. | |
| 72 | Write a short note on partitioning approach. | |
| 73 | Write a short note on Hierarchical approach. | |
| 74 | Write a short note on DBSCAN. | |
| 75 | List and discuss major clustering approaches. | |
| 76 | Write a short note on ROCK. | |
| 77 | Explain agglomeration and divisive approach. | |
| 78 | Apply hierarchical clustering using single linkage to following data. A (1,1), B(1.5,1.5), C(3,4), D(4,4), E(3,3.5) | |
| 79 | What are outliers? How to find out? Write the applications. | |
| | **Unit5** | |
| 80 | What is graph mining and social network? | |
| 81 | What are multimedia and spatial databases? | |
| 82 | Explain set and listed valued attribute with example. | |
| 83 | Explain set and complex structure valued attribute. | |
| 84 | What is spatial aggregation and approximation? Explain with example. | |
| 85 | Define plan, plan database and plan mining. | |
| 86 | Explain the types of dimensions in spatial data cube. | |
| 87 | Explain measures in spatial data cube. | |
| 88 | Discuss approaches for similarity based retrieval in image database. | |
| 89 | Write a short note on mining association in multimedia data. | |
| 90 | Write a short note on text mining. | |
| 91 | Define 1. Term frequency 2. Term frequency matrix 3. Relative term frequency 4. Inverse document frequency | |
| 92 | Compute TF, IDF and TF-IDF for t2 in d2 for following data. | |

| document/term | $t_1$ | $t_2$ | $t_3$ | $t_4$ | $t_5$ | $t_6$ | $t_7$ |
|---------------|-------|-------|-------|-------|-------|-------|-------|
| $d_1$ | 0 | 4 | 10 | 8 | 0 | 5 | 0 |
| $d_2$ | 5 | 19 | 7 | 16 | 0 | 0 | 32 |
| $d_3$ | 15 | 0 | 0 | 4 | 9 | 0 | 17 |
| $d_4$ | 22 | 3 | 12 | 0 | 5 | 15 | 0 |
| $d_5$ | 0 | 7 | 0 | 9 | 2 | 4 | 12 |

| 92 (a) | Compute TF, IDF and TF-IDF for t3 in d4 for following data. | |
|---|---|---|
| | <table><tr><td>document/term</td><td>$t_1$</td><td>$t_2$</td><td>$t_3$</td><td>$t_4$</td><td>$t_5$</td><td>$t_6$</td><td>$t_7$</td></tr><tr><td>$d_1$</td><td>0</td><td>4</td><td>10</td><td>8</td><td>0</td><td>5</td><td>0</td></tr><tr><td>$d_2$</td><td>5</td><td>19</td><td>7</td><td>16</td><td>0</td><td>0</td><td>32</td></tr><tr><td>$d_3$</td><td>15</td><td>0</td><td>0</td><td>4</td><td>9</td><td>0</td><td>17</td></tr><tr><td>$d_4$</td><td>22</td><td>3</td><td>12</td><td>0</td><td>5</td><td>15</td><td>0</td></tr><tr><td>$d_5$</td><td>0</td><td>7</td><td>0</td><td>9</td><td>2</td><td>4</td><td>12</td></tr></table> | |
| 93 | Discuss "Web poses grate challenges for effective resources and knowledge discovery" | |