
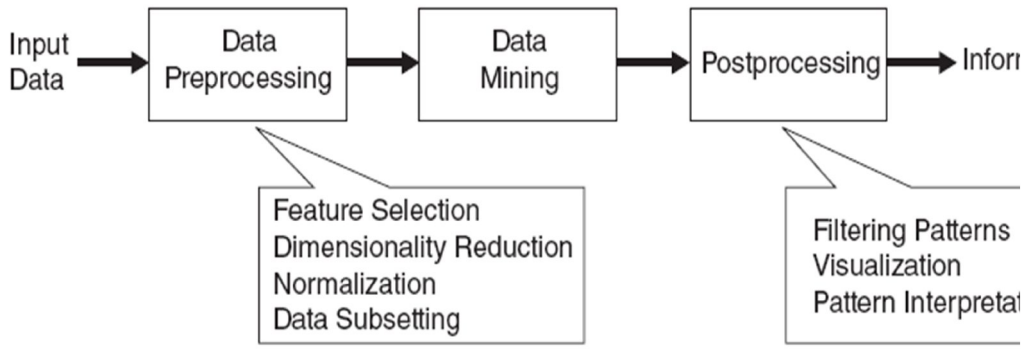


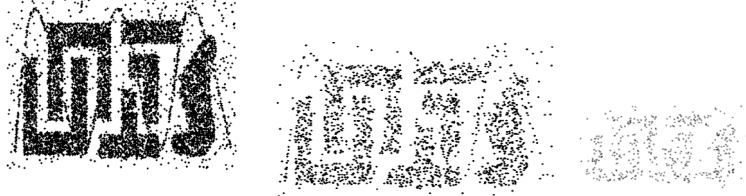
USN					
1	P	E	M	C	A
	<b>PESIT Bangalore South Campus</b> Hosur Road, 1km before Electronic City, Bengaluru -560100 <b>Department of Master of Computer Application</b>				

### INTERNAL ASSESSMENT TEST I

<b>Date</b> : 19-02-2019	<b>Max Marks: 40</b>
<b>Subject &amp; Code:</b> Data Warehousing and Data Mining (17MCA442)	<b>Sem &amp; Sec: 4th Sem MCA</b>
<b>Name of Faculty: Dr. Arti Arya</b>	<b>Time: 11:30AM – 1:00PM</b>

**Note:** Answer FIVE full questions. Select one question from each part.

Part I		
Q 1	<p>i) Briefly explain any two motivating challenges in Data Mining Scalability, Heterogeneous and complex data, High dimensional data, Data ownership and distribution, Non-traditional analysis.</p> <p>Any two can be explained.</p> <p>ii) With a block diagram, explain the process of knowledge discovery and data mining.</p> <div style="text-align: center;">  <pre> graph LR     A[Input Data] --&gt; B[Data Preprocessing]     B --&gt; C[Data Mining]     C --&gt; D[Postprocessing]     D --&gt; E[Information]     B --- B1[Feature Selection Dimensionality Reduction Normalization Data Subsetting]     D --- D1[Filtering Patterns Visualization Pattern Interpretation]           </pre> </div> <p><b>Figure 1.1.</b> The process of knowledge discovery in databases (KDD).</p>	4 4
OR		
Q 2	<p>What are the different methods of handling missing values in a dataset. Differentiate between noise and outlier.</p> <ul style="list-style-type: none"> <li>• Handling missing values             <ul style="list-style-type: none"> <li>– Eliminate Data Objects</li> <li>– Estimate Missing Values</li> <li>– Ignore the Missing Value During Analysis</li> <li>– Replace with all possible values (weighted by their probabilities)</li> </ul> </li> </ul> <p>Noise::Can be wrong measurement. May be distortion of a value. May be addition of spurious objects.</p> <p>Outliers: Outliers are data objects with characteristics that are considerably different than most of the other data objects in the data set. These can be meaningful for user</p>	8

	sometimes. e.g. fraud detection, Intrusion detection.	
<b>Part II</b>		
Q 3	Write the formula for Minkowski's distance and show that $d(A,B) = \sqrt[r]{ A-B  +  B-A }$ satisfies the properties of a metric.  $dist = \left( \sum_{k=1}^n  p_k - q_k ^r \right)^{\frac{1}{r}}$ Refer classwork for the rest of the ans.	8
OR		
Q 4	What are the various types of sampling? Explain with an example. Give an example to show that an inappropriate size of sample may result in loss of information. <ul style="list-style-type: none"> <li>• Sampling is the main technique employed for data selection.</li> <li>• <u>Simple Random Sampling</u></li> <li>• <u>Random Sampling without replacement</u></li> <li>• <u>Random Sampling with replacement</u></li> <li>• <u>Stratified sampling</u></li> <li>• <u>Progressive sampling</u></li> </ul> 8000 points                      2000 points                      500 points  	8
<b>Part III</b>		
Q 5	Compute $L_1$ , $L_2$ and $L_{\infty}$ norm for the following data points: P1(5,0,2,1) and P2(3,4,0,7). Also, Compute similarity between two document vectors $d1=(1,0,0,1,2,5)$ and $d2=(4,7,0,0,1,2)$ using SMC and Jaccard's coefficient.  Ans: Use Euclidean distance for $L_2$ : $\sqrt{(5-3)^2 + (-4)^2 + 2^2 + (-6)^2} = 7.74$ $L_1=14, L_{\infty}=6$	8
OR		
Q 6	What are the various applications of data mining. Write a short note on 'Curse of dimensionality' and discuss any one way of reducing dimensionality.  Applications in the area of Financial Analysis, Fraud detection, Intrusion detection, Predicting stock prices, predicting health concerns, Market basket analysis etc. <ul style="list-style-type: none"> <li>• When dimensionality increases, data becomes increasingly sparse in the space that it occupies. That is for classification, there are not enough data points for creation of a model that confidently assigns a class to all possible objects. Definitions of density and distance between points, which is critical for clustering and outlier detection, become less meaningful.</li> <li>• PCA is a linear algebra technique for cont. attributes that finds new attributes called as Principal components that are linear combination of original attr. Orthogonal to each other. Captures max. variation in the data. SVD- also a DR technique related to PCA</li> </ul>	8

<b>Part IV</b>		
Q 7	<p>Give an architecture for feature subset selection with diagram.</p> <ul style="list-style-type: none"> <li>- 1. A search strategy that generates new subsets of feature.</li> <li>- 2. A measure for evaluating a subset</li> <li>- 3. A stopping criteria</li> <li>- 4. A validation procedure</li> </ul>	<b>8</b>
<b>OR</b>		
Q 8	<p>Differentiate between i) Discretization and Binarization with an example ii) Feature extraction and feature construction.</p> <p><b>Binarization:</b> Both <i>cont.</i> and discrete attr. may be transformed to binary attr.</p> <p><b>Discretization:</b> Transformation of a continuous attribute into categorical attribute....</p> <p><b>Feature Extraction:</b> The creation of new set of features from the original raw data. The feature extraction process results in a much smaller and richer set of attributes. highly domain-specific. The techniques for FE, developed for one field are often not applicable to other fields.</p> <p><b>Construction:</b> Assuming there are n features A1, A2,..... An. After feature construction we may have additional m features An+1, An+2,.... An+m. All new constructed features are defined in terms of original features as such no inherently new information is added through feature construction</p>	<b>4+</b> <b>4</b>
<b>Part V</b>		
Q 9	<p>Compute the correlation coefficient between two variables x and y defined as x=(3,6,0,3,6) and (1,2,2,4,5) and interpret the result.</p> $Corr(x, y) = \frac{cov(x, y)}{sd(x).sd(y)} = \frac{s_{xy}}{s_x s_y}$ $cov(x, y) = s_{xy} = \frac{1}{n-1} \sum_{k=1}^n (x_k - \bar{x})(y_k - \bar{y})$ $sd_x = s_x = \sqrt{\left(\frac{1}{n-1} \sum_{k=1}^n (x_k - \bar{x})^2\right)}$ $sd_y = s_y = \sqrt{\left(\frac{1}{n-1} \sum_{k=1}^n (y_k - \bar{y})^2\right)}$ $\bar{x} = \frac{1}{n} \sum_{k=1}^n x_k$ $\bar{y} = \frac{1}{n} \sum_{k=1}^n y_k$	<b>8</b>

OR		
Q 1 0	<p>What are the different data mining tasks? Explain with brief examples.</p> <ul style="list-style-type: none"><li>• Prediction Methods<ul style="list-style-type: none"><li>◦ Use some variables to predict unknown or future values of other variables.</li><li>◦ The attribute to be predicted- target/dependent variable.</li><li>◦ Independent/explanatory variable- used for making predictions.</li></ul></li><li>• Description Methods<ul style="list-style-type: none"><li>◦ Derive human-interpretable patterns that summarize the underlying relationships in data.</li><li>◦ These are exploratory in nature &amp; requires post processing techniques</li></ul></li></ul>	<b>8</b>