



DEPARTMENT OF COMPUTER ENGINEERING  
[NBA Accredited]

**EXPERIMENT 5**

<b>Title :</b>	Perform exploratory Data Analysis and preprocessing for the data set.
<b>Problem :</b>	Perform following Data pre-processing activities using WEKA:  <ol style="list-style-type: none"><li>1. Normalization</li><li>2. Discretization</li><li>3. Duplicate Removal</li><li>4. Standardization</li><li>5. Handling of Missing Data.</li></ol>
<b>Theory:</b>	<p>Data pre-processing is a data mining technique which is used to transform the raw data in a useful and efficient format.</p> <p><b>Data Cleaning:</b> The data can have many irrelevant and missing parts. To handle this part, data cleaning is done. It involves handling of missing data, noisy data etc.</p> <p><b>Missing Data:</b> This situation arises when some data is missing in the data. It can be handled in various ways. Some of them are: Ignore the tuples: This approach is suitable only when the dataset we have is quite large and multiple values are missing within a tuple. Fill the Missing values: There are various ways to do this task. You can choose to fill the missing values manually, by attribute mean or the most probable value.</p> <p><b>Data Transformation:</b> This step is taken in order to transform the data in appropriate forms suitable for mining process. This involves following ways:</p> <p><b>Normalization:</b> It is done in order to scale the data values in a specified range (-1.0 to 1.0 or 0.0 to 1.0)</p> <p><b>Attribute Selection:</b> In this strategy, new attributes are constructed from the given set of attributes to help the mining process.</p> <p><b>Discretization:</b> This is done to replace the raw values of numeric attribute by interval levels or conceptual levels.</p> <p><b>Concept Hierarchy Generation:</b> Here attributes are converted from lower level to higher level in hierarchy. For Example-The attribute "city" can be converted to "country".</p>



DEPARTMENT OF COMPUTER ENGINEERING  
[NBA Accredited]

<b>Performance</b>	<ol style="list-style-type: none"><li>1. Load the <i>diabetes.arff</i> file into WEKA</li><li>2. Study the dataset and answer following Number of instances = ? Number of Attributes = ? Types of attributes other than class = ? Type of class Attribute = ?</li><li>3. Study the Class labels. Do you think if there is any class imbalance. What are the class labels? Which class is dominant? Paste the screen shot – depicting the number of classes</li><li>4. Study the spread of at-least three numeric attributes and provide statistical measures.</li><li>5. Discretize Age attribute – choose appropriate number of bins for discretization</li><li>6. Normalize plas attribute using min-max normalization – use normalization range as (0 ,1)</li><li>7. Remove Duplicate instances if any</li><li>8. Impute Missing values using mean value imputation.</li></ol>
<b>Deliverables:</b>	A PDF document with screen shots and appropriate description depicting every activity mentioned in the performance section.
<b>Conclusion:</b>	Students will write the conclusion in their own words summarizing the understandings from the practical performed.