



CAYMET’s SIDDHANT COLLEGE OF ENGINEERING, SUDUMBARE

Department of Computer Engineering

SEM - VI (Year: 2024 - 25)

Subject: DSBDA Laboratory

Year: T. E.

Name:

Roll No.

LIST OF ASSIGNMENTS

Sr. No.	Title of Assignment	Date	Sign
Group A			
1	<p>Data Wrangling, I</p> <p>Perform the following operations using Python on any open source dataset (e.g., data.csv)</p> <ol style="list-style-type: none">1. Import all the required Python Libraries.2. Locate an open source data from the web (e.g., https://www.kaggle.com). Provide a clear description of the data and its source (i.e., URL of the web site).3. Load the Dataset into pandas dataframe.4. Data Preprocessing: check for missing values in the data using pandas isnull(), describe() function to get some initial statistics. Provide variable descriptions. Types of variables etc. Check the dimensions of the data frame.5. Data Formatting and Data Normalization: Summarize the types of variables by checking the data types (i.e., character, numeric, integer, factor, and logical) of the variables in the data set. If variables are not in the correct data type, apply proper type conversions.6. Turn categorical variables into quantitative variables in Python.		
2	<p>Data Wrangling, II</p> <p>Create an “Academic performance” dataset of students and perform the following operations using Python.</p> <ol style="list-style-type: none">1. Scan all variables for missing values and inconsistencies. If there are missing values and/or inconsistencies, use any of the suitable techniques to deal with them.2. Scan all numeric variables for outliers. If there are outliers, use any of the suitable techniques to deal with them.3. Apply data transformations on at least one of the variables. The purpose of this transformation should be one of the following reasons: to change the scale for better understanding of the variable, to convert a non-linear relation into a linear one, or to decrease the skewness and convert the distribution into a normal distribution.		

3	Descriptive Statistics - Measures of Central Tendency and variability Write a Python program to display some basic statistical details like percentile, mean, standard deviation etc. of the species of 'Iris-setosa', 'Iris-versicolor' and 'Iris-versicolor' of iris.csv dataset.		
4	Data Analytics I Create a Linear Regression Model using Python/R to predict home prices using Boston Housing Dataset. The Boston Housing dataset contains information about various houses in Boston through different parameters. There are 506 samples and 14 feature variables in this dataset.		
5	Data Analytics II <ol style="list-style-type: none">1. Implement logistic regression using Python/R to perform classification on Social_Network_Ads.csv dataset.2. Compute Confusion matrix to find TP, FP, TN, FN, Accuracy, Error rate, Precision, Recall on the given dataset.		
6	Data Analytics III <ol style="list-style-type: none">1. Implement Simple Naïve Bayes classification algorithm using Python/R on iris.csv dataset.2. Compute Confusion matrix to find TP, FP, TN, FN, Accuracy, Error rate, Precision, Recall on the given dataset.		
7	Text Analytics <ol style="list-style-type: none">1. Extract Sample document and apply following document preprocessing methods: Tokenization, POS Tagging, stop words removal, Stemming and Lemmatization.2. Create representation of document by calculating Term Frequency and Inverse Document Frequency.		
8	Data Visualization I <ol style="list-style-type: none">1. Use the inbuilt dataset 'titanic'. The dataset contains 891 rows and contains information about the passengers who boarded the unfortunate Titanic ship. Use the Seaborn library to see if we can find any patterns in the data.2. Write a code to check how the price of the ticket (column name: 'fare') for each passenger is distributed by plotting a histogram.		
9	Data Visualization II <ol style="list-style-type: none">1. Use the inbuilt dataset 'titanic' as used in the above problem. Plot a box plot for distribution of age with respect to each gender along with the information about whether they survived or not. (Column names: 'sex' and 'age').2. Write observations on the inference from the above statistics.		
10	Data Visualization III Download the Iris flower dataset or any other dataset into a DataFrame. Scan the dataset and give the inference as: <ol style="list-style-type: none">1. List down the features and their types (e.g., numeric, nominal) available in the dataset.2. Create a histogram for each feature in the dataset to illustrate the feature distributions.3. Create a boxplot for each feature in the dataset.4. Compare distributions and identify outliers.		

Group B			
11	Write a code in JAVA for a simple WordCount application that counts the number of occurrences of each word in a given input set using the Hadoop MapReduce framework on local-standalone set-up.		
12	Design a distributed application using MapReduce which processes a log file of a system.		
13	Write a simple program in SCALA using Apache Spark framework		
Group C (Mini Project)			
1	Use the Tweets.csv dataset and classify tweets into positive and negative tweets.		
2	Develop a movie recommendation model using the scikit-learn library in python.		

Prof. Trupti Rajput
Subject In-charge

N. S. Kulkarni
HOD