PUNE INSTITUTE OF COMPUTER TECHNOLOGY DHANKAWADI, PUNE – 43.

LIST OF LAB EXPERIMENTS

DEPARTMENT: Computer Engineering CLASS: T.E

ACADEMIC YEAR: 2022-23 SEMESTER: II

SUBJECT: Data Science and Big Data Analytics Lab

LAB EXPT. NO	PROBLEM STATEMENT
	GROUP A
	Data Wrangling, I
	Perform the following operations using Python on any open-source dataset (e.g., data.csv) 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' data frame.
	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.
	In addition to the codes and outputs, explain every operation that you do in the above steps and explain everything that you do to import/read/scrape the data set.
2	Data Wrangling II
	Create an "Academic performance" dataset of students and perform the following operations
	 using Python. 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. 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. Reason and document your approach properly.
3	Descriptive Statistics - Measures of Central Tendency and variability
	Perform the following operations on any open-source dataset (e.g., data.csv) 1. Provide summary statistics (mean, median, minimum, maximum, standard deviation)

for a dataset (age, income etc.) with numeric variables grouped by one of the qualitative (categorical) variable. For example, if your categorical variable is age groups and quantitative variable is income, then provide summary statistics of income grouped by the age groups. Create a list that contains a numeric value for each response to the categorical variable.

2. 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.

Provide the codes with outputs and explain everything that you do in this step.

4 Data Visualization I

- 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.

The objective is to predict the value of prices of the house using the given features.

5 Data Visualization II

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')

6 Data Visualization III

Download the Iris flower dataset or any other dataset into a DataFrame. (e.g., https://archive.ics.uci.edu/ml/datasets/Iris).

Scan the dataset and give the inference as:

- 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 box plot for each feature in the dataset. Compare distributions and identify outliers.

7 Text Analytics

- 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 Analytics I

Create a Linear Regression Model using Python/R to predict home prices using Boston Housing Dataset (https://www.kaggle.com/c/boston-housing). 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.

9 Data Analytics II

- 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

10	Data Analytics III
	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.
	Group B- Big Data Analytics – JAVA/SCALA
11	Write a code in JAVA for a simple Word Count application that counts the number of
	occurrences of each word in a given input set using the Hadoop Map-Reduce framework on
	local-standalone set-up.
12	Locate dataset (e.g., sample_weather.txt) for working on weather data which reads the text
	input files and finds average for temperature, dew point and wind speed using the Hadoop
12	Map-Reduce framework on local-standalone set-up.
13	Write a simple program in SCALA using Apache Spark framework
	Group C- Mini Projects/ Case Study – PYTHON/R (Any TWO Mini Project) (Students will select one mini project from 14.15.16)
14	(Students will select one mini project from 14,15,16) Use the following dataset and classify tweets into positive and negative tweets.
14	https://www.kaggle.com/ruchi798/data-science-tweets
15	Develop a movie recommendation model using the scikit-learn library in python.
13	Refer dataset
	https://github.com/rashida048/Some-NLP-Projects/blob/master/movie_dataset.csv
16	Use the following covid vaccine statewise.csv dataset and perform following analytics on the
	given dataset
	https://www.kaggle.com/sudalairajkumar/covid19-in-
	india?select=covid_vaccine_statewise.csv
	a. Describe the dataset
	b. Number of persons state wise vaccinated for first dose in India
	c. Number of persons state wise vaccinated for second dose in India
	d. Number of Males vaccinated
1.7	d. Number of females vaccinated
17	Write a case study to process data driven for Digital Marketing OR Health care systems with
	Hadoop Ecosystem components as shown. (Mandatory)
	HDFS: Hadoop Distributed File SystemYARN: Yet Another Resource Negotiator
	MapReduce: Programming based Data Processing
	 Spark: In-Memory data processing
	PIG, HIVE: Query based processing of data services
	HBase: NoSQL Database (Provides real-time reads and writes)
	Mahout, Spark MLLib: (Provides analytical tools) Machine Learning algorithm
	libraries
	Solar, Lucene: Searching and Indexing
	Question -Answer session with students about all above experiments

Head of Department Dr. G V Kale

Subject Coordinator Mrs. R S Paswan