KJ's Educational Institutes

# Trinity Academy of Engineering, Pune

## Department of Computer Engineering

**S.E. (Computer Engineering)**

**DATA SCIENCE AND BIG DATA ANALYTICS**

**Staff: Mrs. Nilufar Zaman**

**(2019 course)**

| Teaching Scheme: | Examination Scheme: | Credit: |
|---|---|---|
| **Practical: 4Hrs/ Week/Batch** | **Term Work: 25 Marks** | **Practical: 4** |

# DATA SCIENCE AND BIG DATA ANALYTICS

## <u>List of Assignments</u>

1. **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 insult(), 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.
   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.

   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 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.

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

## 5. 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.

## 6. 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,
3. Recall on the given dataset.

## 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 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.

## 9. 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')

Write observations on the inference from the above statistics.

10. 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.


11. Locate dataset (e.g., sample_weather.txt) for working on weather data which reads the textinput files and finds average for temperature, dew point and wind speed.


12. Write a simple program in SCALA using Apache Spark framework.


13. Use the following dataset and classify tweets into positive and negative tweets.
https://www.kaggle.com/ruchi798/data-science-tweets


14 . Develop a movie recommendation model using the scikit-learn library in python.
Refer dataset:
https://github.com/rashida048/Some-NLP-Projects/blob/master/movie_dataset.csv

**Assignment 02:**

**Title of the Assignment: Data Wrangling, II**

Create an "Academic performance" dataset of students and perform the following operations using Python.

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.

Reason and document your approach properly.

-------------------------------------------------------------------------------------------------------
---

**Objective of the Assignment:** Students should be able to perform the data wrangling operation using Python on any open source dataset

**Prerequisite:**

1. Basic of Python Programming

2. Concept of Data Preprocessing, Data Formatting , Data Normalization and Data Cleaning.

-------------------------------------------------------------------------------------------------------
-

**Contents for Theory:**

1. Identification and Handling of Null Values

2. Identification and Handling of Outliers

3. Data Transformation for the purpose of :

     a. To change the scale for better understanding

     b. To decrease the skewness and convert distribution into normal distribution

-------------------------------------------------------------------------------------------------------
--

**Theory:**

**1. Identification and Handling of Null Values**

Missing Data can occur when no information is provided for one or more items or for a whole unit. Missing Data is a very big problem in real-life scenarios. Missing Data can also refer to as NA(Not Available) values in pandas. In DataFrame sometimes many datasets simply arrive with missing data, either because it exists and was not collected or it never existed. For Example, Suppose different users being surveyed may choose not to share their income, some users may choose not to share the address in this way many datasets went missing.

**In Pandas missing data is represented by two value:**

a) None: None is a Python singleton object that is often used for missing data in Python code.

b) NaN : NaN (an acronym for Not a Number), is a special floating-point value recognized by all systems that use the standard IEEE floating-point representation.

c) Pandas treat None and NaN as essentially interchangeable for indicating missing or null values. To facilitate this convention, there are several useful functions for detecting, removing, and replacing null values in Pandas DataFrame :

- isnull()
- notnull()
- dropna()
- fillna()
- replace()

1. **Checking for missing values using isnull() and notnull()**

- **Checking for missing values using isnull()**

  - In order to check null values in Pandas DataFrame, isnull() function is used. This function return dataframe of Boolean values which are True for NaN values.

- **Checking for missing values using notnull()**

  - In order to check null values in Pandas Dataframe, notnull() function is used. This function return dataframe of Boolean values which are False for NaN values.

2. **Filling missing values using dropna(), fillna(), replace()**

- **In order to fill null values in a datasets, fillna(), replace() functions are used.**

  - These functions replace NaN values with some value of their own. All these functions help in filling null values in datasets of a DataFrame.

  - **For replacing null values with NaN**

    missing_values = ["Na", "na"]

    df = pd.read_csv("StudentsPerformanceTest1.csv", na_values =missing_values)

  - **Deleting null values using dropna() method**

    - In order to drop null values from a dataframe, dropna() function is used. This function drops Rows/Columns of datasets with Null values in different ways.

      ➢ Dropping rows with at least 1 null value

> ➢ Dropping rows if all values in that row are missing
> ➢ Dropping columns with at least 1 null value.
> ➢ Dropping Rows with at least 1 null value in CSV file

### 3. Handling of Outliers:

- For removing the outlier, one must follow the same process of removing an entry from the dataset using its exact position in the dataset because in all the above methods of detecting the outliers end result is the list of all those data items that satisfy the outlier definition according to the method used.

- Below are some of the methods of treating the outliers
  - o Trimming/removing the outlier
  - o Quantile based flooring and capping
  - o Mean/Median imputation

### 4. Data Transformation for the purpose of :

- Data transformation is the process of converting raw data into a format or structure that would be more suitable for model building and also data discovery in general. The process of data transformation can also be referred to as extract/transform/load (ETL). The extraction phase involves identifying and pulling data from the various source systems that create data and then moving the data to a single repository. Next, the raw data is cleansed, if needed. It's then transformed into a target format that can be fed into operational systems or into a data warehouse, a date lake or another repository for use in business intelligence and analytics applications. The transformation The data are transformed in ways that are ideal for mining the data. The data transformation involves steps that are.

- **Smoothing:** It is a process that is used to remove noise from the dataset using some algorithms It allows for highlighting important features present in the dataset. It helps in predicting the patterns.

- **Aggregation:** Data collection or aggregation is the method of storing and presenting data in a summary format. The data may be obtained from multiple data sources to integrate these data sources into a data analysis description. This is a crucial step since the accuracy of data analysis insights is highly dependent on the quantity and quality of the data used.

- **Generalization:** It converts low-level data attributes to high-level data attributes using concept hierarchy. For Example Age initially in Numerical form (22, 25) is converted into categorical value (young, old).

- **Normalization:** Data normalization involves converting all data variables into a given range. Some of the techniques that are used for accomplishing normalization are:
  - i. **Min–max normalization:** This transforms the original data linearly.

- **Z-score normalization:** In z-score normalization (or zero-mean normalization) the values of an attribute (A), are normalized based on the mean of A and its standard deviation.

- **Normalization by decimal scaling**: It normalizes the values of an attribute by changing the position of their decimal points.

- **Attribute or feature construction.**

New attributes constructed from the given ones: Where new attributes are created & applied to assist the mining process from the given set of attributes. This simplifies the original data & makes the mining more efficient.

**Conclusion:** In this way we have explored the functions of the python library for Data Identifying and handling the outliers. Data Transformations Techniques are explored with the purpose of creating the new variable and reducing the skewness from datasets.

**Assignment Question:**

**1. Explain the methods to detect the outlier.**
**2. Explain data transformation methods**
**3. Write the algorithm to display the statistics of Null values present in the dataset.**
**4. Write an algorithm to replace the outlier value with the mean of the variable.**