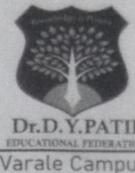


Experiment 1

Page No.	
Date	



Title : DATA WRANGLING I

PROBLEM STATEMENT:-

1. Locate open source data from web.
2. Provide clear description of the data
3. Load the Dataset into pandas data frame.
4. Data preprocessing.
5. Data formatting and data normalization.
6. Turn categorical.

Objectives :-

1. To learn and concepts of python libraries.
2. To learn and understand data science.
3. To understand and practice.

Prerequisite :-

1. Basic python programming.
2. Concept of Data Preprocessing, Data Formatting.

Theory:-

1. Introduction to Big data

Big data means really a big data, it is a collection of large datasets that cannot be processed using traditional computing technique. Big data involves data produced by different devices and applications.



2. Introduction to Dataset.

A dataset is a collection of records, similar to relational database table. Records are similar to table rows, but the columns can contain not only string

3. Python Libraries for Data Science.

a. Numpy :-

One of the most fundamental packages in python. Numpy is general purpose array processing package. It provides high-performance multidimensional array objects and tools to work with arrays.

What can do Numpy?

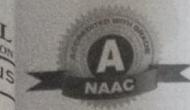
1. Basic array operations: add, multiply
2. Advanced array operations: stack and
3. Work with Date Time
4. Basic slicing and Advanced indexing.

b. Pandas :-

- Pandas is an open-source package that provides high-performance, easy to use data structures and data analysis tools.

What can you do with pandas?

1. Indexing, manipulating, renaming.
2. Update, Add or Delete columns.
3. Impute missing files.



Q. scikit learn :-

- Introduced to the world as Google Summer Code project. Scikit learn is robust machine learning library for python. It features ML features like SVM's, random forests.

What can you do scikit learn?

1. Classification : spam detection.
2. Clustering : Drug response, stock price.
3. Dimensionality reduction.
4. Dimension Model selection.

Description of Dataset :-

The iris dataset was used in R.A.Fisher's classic 1936 paper.

Measurements in Taxonomic problems, and can also be found on UCI machine learning Repository. It includes three iris species with 50 samples each and as well as some properties about each flower.

The sample - 150

The columns in this dataset are:

1. Id
2. SepalLengthCm
3. SepalWidthCm
4. PetalLengthCm
5. Species.

5. Pandas Data frame for Load Dataset

1. The dataset is downloaded from UCI
2. Now read csv file as a DataFrame in python from path where you saved the same. The Iris Dataset is stored in .csv stands for comma-separated values. It is easier to load .csv files in pandas DataFrame and perform various analytical operations on it.

Load Iris.csv into a pandas DataFrame

Syntax:

```
iris = pd.read_csv('csv_file', header=None)
```

- 3) Read the dataset from UCI machine learning repository link and specify columns name to use.

Function	Description
dataset.head(n)	Return first n rows.
dataset.tail(n)	Return last n rows.
dataset.index	The index of data
dataset.columns	The column labels of the dataset

5. dataset.shape

Return a tuple representing the dimensions of dataset.

6. dataset.dtypes

Return the dtypes of dataset

7. dataset[columnname]

Read data column wise.

8. dataset.iloc[5]

Purely integer-location based indexing for selection.

9. dataset[0:3]

Selecting via [] which slices the rows.

10. dataset.iloc[:,n]

A subset of first n columns of original data

11. dataset.iloc[:,m:n]

A subset of the first m rows & first n columns.

- Missing values in DataFrame as a whole.

Syntax: DataFrame.isnull()

- Missing values across each column.

Syntax: DataFrame.isnull.any()

- Count of missing values across each column using isna() + isnull().

Syntax: DataFrame.isnull().sum().sum()

F. Pandas functions for Data Formatting and Normalization:-

The transforming data stage is about converting the data set into a format that can be analyzed or modeled effectively.

Dataframe Function Description Output

- (1) df.dtypes To check the data types
- (2) df['Petal length (cm)'] = df['Petal length (cm)'].astype('int') To change the data type

B1 Data normalization - Mapping all the values onto a uniform scale is involved in data normalization. It is also known as min-max scaling.

Algorithm:

- Step 1: Import pandas and sklearn library
- Step 2: Load Iris dataset in dataframe object df.

```
iris = load_iris()
```

```
df = pd.DataFrame(iris.data, columns = iris.feature_names)
```

Step 3: Print dataset

```
df.head()
```

Step 4: Create minimum & maximum Processor project

Step 5: X-scaled = min-max scaler.

```
fit_transform(X)
```

Step 6: view the database

df normalized.

8. Panda function for handling categorical variables:

- Categorical variables have values that describe 'quality' or 'characteristics' of a data unit.
- Categorical features refer to string type data and can be easily understood by human beings. Therefore, categorical data must be translated into numerical data that can be understood by machine.

Label Encoding : Label Encoding refers to converting the labels into a numeric form so as to convert them into machine readable form. It is an important processing step for structured dataset in supervised learning.

eg:-	Height	Height
	Tall	0
	medium	1
	short	2

Label Encoding on iris dataset:

For iris dataset target column is species.

• Preprocessing: LabelEncoder

• fit_transform(T):



Algorithm:

Step 1: Import pandas and sklearn library for preprocessing.

from sklearn import preprocessing

Step 2: Load the Iris dataset in dataframes object of.

Step 3: observe the unique value for the species column.

df['species'].unique()

Step 4: define LabelEncoder object known how to understand labels

labelEncoder = preprocessing.LabelEncoder

Step 5: Encode labels in column

df['species'] = labelEncoder.fit_transform(df['species'])

Step 6: observe the unique values for the species column.

df['species'].unique()

Conclusion :-

In this way we have explored the functions of Python library for Data Preprocessing, Data wrangling tool and How to handle missing values on Iris dataset.



Experiment 2

Page No.	
Date	



Title : DATA WRANGLING II

Problem statement :-

Create an 'Academic Performance' dataset of students and perform the following operations using python.

1. Scan all variables for missing values and inconsistencies.
2. Scan all numeric variables for outliers.
3. Apply data transformations on least one of the variables.

Objectives / main purpose :-

Students should be able to perform data wrangling operations using python on any open source dataset.

PREREQUISITE :-

1. Basic python programming.
2. Concept of Data preprocessing, Data formatting, Data normalization.

THEORY :-

1. Creation of dataset using MS Excel:-
The dataset created as "CSV" format.
• The name of dataset is studentPerformance.
• The features of dataset are: Math score, Reading score, Writing score, Placement score.

- Number of Instances : 30
- The response variable is Placement offer count.
- Range of values.
- math score [60-80], club JOIN_DATE [01-12]
- Syntax :- RANDBETWEEN (bottom, up)

2. Identification and Handling of NULL :-
 Missing data can occur when no info is provided for one or more items or for while unit. Missing data is very big problem in real life scenarios. Missing data can also refer to as NA.

In pandas missing data is represented

- None
- NaN

Pandas treat None and NaN as essentially interchangeable for indicating missing or NULL values. To facilitate this conversion

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

- Checking for missing values using isnull() and notnull()
- Checking for missing values using isnull()

Algorithm:

Step 1: Import pandas and numpy in order to check missing values in Pandas database.

import pandas as pd

import numpy as np

Step 2: Load the dataset in dataframe object df

df = pd.read_csv('student_performance.csv')

Step 3: Display the dataframe

Step 4: Use isnull() function to check null values in dataset.

df.isnull().

Step 5: To create a series true for NaN values for specific column.

Checking for missing values (using isnull())

1. Algorithm

Step 1: Import pandas and numpy in order to check missing values in Pandas database.

import pandas as pd

import numpy as np

Step 2: Load the dataset in dataframe

df = pd.read_csv('student_performance.csv')

Step 3: Display dataframe.

df

3. Identification and Handling of outliers:
 One of the most imp steps as part of data preprocessing is detecting and treating the outliers as they can negatively affect the statistical analysis and the training process of machine learning algorithm.

1. What are outliers?

We all have heard the idiom "odd one out" which means something unusual in comparison to the others in group.

2. Why do they occur?

An outlier may occur due to the variable in data, or due to experimental error, human error.

3. Why do they affect?

In statistics we have three measure of tendencies mean, median, mode.

4. Detecting outliers:-

If our dataset is small we can detect the outlier by at the dataset. But what if we have huge dataset. Below are some techniques:-

- Boxplots
- Scatterplots.
- Z score.
- Inner Quartile Range.

Predicting outliers using Boxplot :-
 It captures the summary of data effectively and efficiently with only a simple box and whisker. Boxplot summarizes sample data using 25th, 50th & 75th percentile.

Algorithm :-

Step 1: import pandas as pd.
 import numpy as np

Step 2: load the dataset

df = pd.read_csv("demo.csv")

Step 3: display the data frame.

Step 4: Select the columns for boxplot.
 Col = ['mathscore', 'reading score', 'writingscore',
 'Placement score']

Step 5: He can now print the outliers
 foreach column with reference

print(np.where(df['mathscore'] > 90))

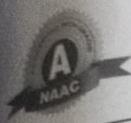
Handling of outliers :-

For removing the outliers, 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.



Page No.	
Date	

Conclusion 8- In this way we have explored the functions of the python library for data identifying and handling outliers. Data transformation techniques are explored with the purpose of creating the new variable and reducing the skewness from data.



Page No.	
Date	

Experiment NO - 3

Title :- Measure of Central Tendency AND Variability.

Problem Statement :-

Perform the following operations on any source dataset.

1. Provide summary statistics (mean, median, minimum, S.D) for a dataset with numeric variables.
2. Write a python program to display some basic statistical details like percentile, mean.

Objective :- To analyze and demonstrate knowledge of statistical data analysis's technique for decision-making.

PREREQUISITE :-

- ① Basic of python programming.
- ② Concept of statistics mean, median, minimum, maximum, S.D.

THEORY :-

The data are summarized in some, but not all ways. We choose descriptive that are either most often reported covered in introductory courses.

Central Tendency :-

- Mean
- Median
- Mode.

- Dispersion :-
- Standard Deviation :-
- Minimum :-
- Maximum :-

Central Tendency :-

$$\text{The mean: } \bar{x} = \frac{\sum x}{N}$$

Here, Σ represents summation.

x represents observation.

N represents the no. of observations.

$$\text{Mean} = \frac{\sum f x}{\sum f}$$

$$\text{Where } \sum f = N$$

The median :-

If the total no. of obs (n) is odd number, then formula given below.

$$\text{Median} = \left(\frac{n+1}{2} \right)^{\text{th}} \text{ obs.}$$

If the total no. of obs (n) is even number, the formula given below.

$$\text{Median} = \left(\frac{n}{2} \right)^{\text{th}} \text{ obs} + \left(\frac{n}{2} + 1 \right)^{\text{th}} \text{ obs}$$

2

- When data presented in the form of frequency distribution.

Find median class, the total Σf .

The median class consist of class in which $(\frac{n}{2})$ is present.

$$\text{Median} = l + \left[\frac{\frac{n}{2} - C}{f} \right] \times h$$

The mode :-

The mode is most frequent occurring obs. Consider case where data is continuous and the value mode be computed using following steps.

a) determine modal class that possessing maximum frequency.

b) calculate mode using formula.

$$\text{Mode} = l + \left[\frac{f_m - f_1}{2f_m - f_1 - f_2} \right] \times h$$

Standard Deviation :-

$$S = \sqrt{\frac{\sum (x - \bar{x})^2}{n-1}}$$

Steps for calculating s.d:-

i. Step 1 : Find Mean.

ii. Step 2 : Find each score's deviation from the mean.



Page No.	
Date	

Step 3: Square each deviation from the mean.

Step 4: Find the sum of squares.

Step 5 : Find the variance

Step 6 : Find square root of variance.

Use of describe function
To display statistical details like mean, S.D etc.

Conclusion:- In this way we have explored the function of python library for calculate Data statistics.

Experiment - 4.

Title - Data analytics - I

Problem statement :- Create a Linear Regression model using Python & R to predict home price prediction using Boston Housing Dataset.

~~Objective :- students should be able to do data analysis using linear regression using python for any open-source dataset~~

Prerequisite :-

1. Basic of python programming
2. Concept of Regression

Contents of Theory :-

1. Linear Regression :- Univariate & Multivariate.
2. Least Square method for linear regression.
3. Measuring performance of LR.
4. Example of LR.
5. Training dataset & Testing dataset.

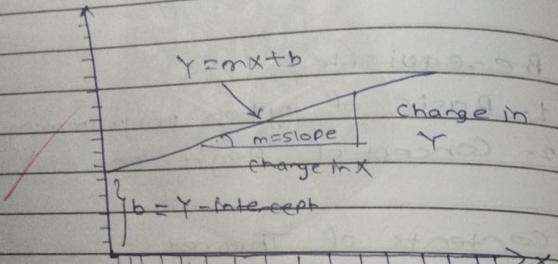
1. Linear Regression :-

It is machine learning algorithm based on supervised learning. It target predicted values on basis of independent variables. It is prefered to find out relationship betn forecasting and variables.

- Linear regression is popular because its cost function is Mean Squared Error which is equal to average squared difference b/w observations actual and predict values.

equation of line like:

$$Y = m \cdot X + b + c$$



Multivariate Regression :-

- It concerns the study of two or more predictor variables.

Usually transformations of the original features into polynomial features.

- A simple linear model $Y = a + bx$; in original feature will be transformed into polynomial features & linear reg applied it and, it something like,

$$Y = a + b_1 X_1 + c_1 X_2$$

- If high degree value used in to the curve becomes over-fitted

- List square Method of Linear Regression
Linear regression involves establishing linear relationship b/w dependent & independent variables. Such relationship is portrayed in the form of eqn.
- A simple linear model is one of which involves only one dependent and one independent variable.
- However, for a simple univariate linear model it can be denoted by the regression eqn.

3. Measuring Performance of Linear Regression - Mean Square Error:-

The mean squared error (MSE) represent the error of estimator or predictive model created based on given set of obs in the sample.. Two or more regression models created using the same data

$$MSE = \frac{1}{n} \sum (Y - \hat{Y})^2$$

An MSE of zero(0) represent the fact that the predictor is a perfect Predictor.

iii) Interpretation of regression line.

Interception 1

For an increase in value of x by 0.6 units there is an increase in value of y in one unit.

Interception 2 :-

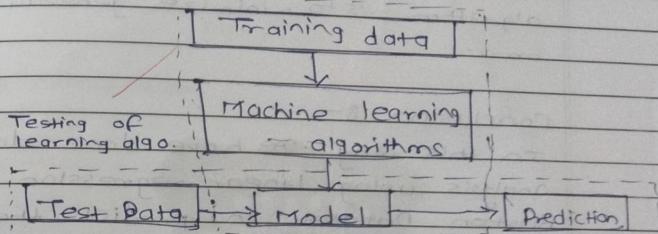
Even if $x=0$ value of independent variable it is expected that value of y is 20 score in XIT standard (y_i) is 0.69 units depending on score in X standard.

Machine learning algorithm has two phases:-
1. Training and 2. Testing.

The input of the training phase is training data, which is passed to any machine learning algo and machine learning model is generated as output of training phase.

The input of testing phase is test data which is passed to the machine learning model and prediction is done.

Training Learning Algorithm



Training & Testing phase in machine learning:-

a) Training phase:-

- Training dataset is provided as input.
- Training dataset is dataset having attributes & class labels.

b) Testing phase:-

- Testing dataset is provided as input to this phase.
- Test dataset for which class label is unknown.
- A test dataset used for which class label is unknown.

c) Generalization :-

- Generalization is the prediction of the future based on the past system.
- It needs to generalize beyond the training data to some future data.



Page No.		
Date		

- The ultimate goal of machine learning algorithm is to minimize generalization error.

Conclusion :-

In this way we have done data analysis using linear regression for Boston Dataset and predict the price of house using features of the Boston Dataset.

continuous Assessment of student :-