

**Aim: Write down to libraries used for the following machine learning applications**

**1. Scikit-learn(sklearn):**

- Scikit-learn is probably the most useful library for machine learning in Python. The sklearn library contains a lot of efficient tools for machine learning and statistical modeling including classification, regression, clustering and dimensionality reduction.
- In python, sklearn is a machine learning package which includes a lot of ML algorithms.

```
[9] import sklearn as sc
    print(sc)

<module 'sklearn' from '/usr/local/lib/python3.7/dist-packages/sklearn/__init__.py'>
```

**2. NumPy:**

- NumPy is a Python library used for working with arrays.
- It also has functions for working in domain of linear algebra, fourier transform, and matrices.
- NumPy was created in 2005 by Travis Oliphant. It is an open source project and you can use it freely.
- NumPy stands for Numerical Python.
- NumPy is a Python library and is written partially in Python, but most of the parts that require fast computation are written in C or C++.
- It is a numeric python module which provides fast maths functions for calculations.
- It is used to read data in numpy arrays and for manipulation purpose.

```
[2] import numpy as np
    print(np.pi)

3.141592653589793
```

**3. Pandas:**

- Pandas is a Python library used for working with data sets.
- It has functions for analyzing, cleaning, exploring, and manipulating data.
- The name "Pandas" has a reference to both "Panel Data", and "Python Data Analysis" and was created by Wes McKinney in 2008.
- Used to read and write different files such as csv.
- Data manipulation can be done easily with dataframes.

```
import pandas as pd
print(pd)

<module 'pandas' from '/usr/local/lib/python3.7/dist-packages/pandas/__init__.py'>
```

**4. Matplotlib:**

- Matplotlib is an amazing visualization library in Python for 2D plots of arrays. Matplotlib is a multi-platform data visualization library built on NumPy arrays and designed to work with the broader SciPy stack. It was introduced by John Hunter in the year 2002.

- One of the greatest benefits of visualization is that it allows us visual access to huge amounts of data in easily digestible visuals. Matplotlib consists of several plots like line, bar, scatter, histogram etc

```
[5] import matplotlib as mp
    print(mp)

<module 'matplotlib' from '/usr/local/lib/python3.7/dist-packages/matplotlib/__init__.py'>
```

## 5. SciPy:

- SciPy is an open-source Python library which is used to solve scientific and mathematical problems. It is built on the NumPy extension and allows the user to manipulate and visualize data with a wide range of high-level commands. As mentioned earlier, SciPy builds on NumPy and therefore if you import SciPy, there is no need to import NumPy.

```
[4] import scipy as sp
    print(sp)

<module 'scipy' from '/usr/local/lib/python3.7/dist-packages/scipy/__init__.py'>
```

| ALGORITHM NAME | NUMPY | PANDAS | SCI-KT-LEARN (SKLEARN) | MATPLOTLIB | XGBOOST |
|----------------|-------|--------|------------------------|------------|---------|
| REGRESSION     |       |        |                        |            |         |
| LINEAR         | Y     | Y      | Y                      | Y          |         |
| LOGISTIC       | Y     | Y      | Y                      | Y          |         |
| KNN            | Y     | Y      | Y                      | Y          |         |
| CLASSIFICATION |       |        |                        |            |         |
| ID3            | Y     | Y      | Y                      | Y          |         |
| C4.5           | Y     | Y      | Y                      | Y          |         |
| NAÏVE BAYES    | Y     | Y      | Y                      | Y          |         |
| CLUSTERING     |       |        |                        |            |         |
| K-MEANS        | Y     | Y      | Y                      | Y          |         |
| K-MEDOID       | Y     | Y      | Y                      | Y          |         |
|                |       |        |                        |            |         |
| SVM            | Y     | Y      | Y                      | Y          |         |
| RANDOM FOREST  | Y     | Y      | Y                      | Y          |         |
| XGBOOST        | Y     | Y      | Y                      | Y          | Y       |

## **Numpy Applications:**

### **1. An alternative for lists and arrays in Python**

Arrays in Numpy are equivalent to lists in python. Like lists in python, the Numpy arrays are homogenous sets of elements. The most important feature of NumPy arrays is they are homogenous in nature.

This differentiates them from python arrays. It maintains uniformity for mathematical operations that would not be possible with heterogeneous elements. Another benefit of using NumPy arrays is there are a large number of functions that are applicable to these arrays.

These functions could not be performed when applied to python arrays due to their heterogeneous nature.

### **2. NumPy maintains minimal memory**

Arrays in NumPy are objects. Python deletes and creates these objects continually, as per the requirements. Hence, the memory allocation is less as compared to Python lists. NumPy has features to avoid memory wastage in the data buffer.

It consists of functions like copies, view, and indexing that helps in saving a lot of memory. Indexing helps to return the view of the original array, that implements reuse of the data. It also specifies the data type of the elements which leads to code optimization.

### **3. Using NumPy for multi-dimensional arrays**

We can also create multi-dimensional arrays in NumPy. These arrays have multiple rows and columns. These arrays have more than one column that makes these multi-dimensional. Multi-dimensional array implements the creation of matrices.

These matrices are easy to work with. With the use of matrices the code also becomes memory efficient. We have a matrix module to perform various operations on these matrices.

### **4. Mathematical operations with NumPy**

Working with NumPy also includes easy to use functions for mathematical computations on the array data set. We have many modules for performing basic and special mathematical functions in NumPy.

There are functions for Linear Algebra, bitwise operations, Fourier transform, arithmetic operations, string operations, etc.

## **Numpy Array Applications**

### **1. Shape Manipulations**

Users can change array dimensions at runtime if the output produces the same number of elements. We apply `np.reshape(...)` function on the array. The reshape function is useful for performing various operations. For eg, we use it when we want to broadcast two dissimilar arrays.

## **2. Array Generation**

We can generate array data set for implementing various functions. We can also generate a predefined set of numbers for the array elements using the `np.arange(...)` function. Reshape function is useful to generate a different set of dimensions.

We can also use the random function to generate an array having random values. Similarly, we can use `linspace` function to generate arrays having similar spacing in elements.

We can create arrays with pre-filled ones or zeroes. The default data type is set to be `float64` but we can edit the data type using `dtype` option.

## **3. Array Dimensions**

Numpy consists of both one and multidimensional arrays. Some functions have restrictions on multidimensional arrays. It is then necessary to transform those arrays into one-dimensional arrays. We can transform multi-dimensional to single dimension using `np.ravel(..)`

## **Numpy Applications with Other Libraries**

### **1. NumPy with Pandas**

Pandas is one of the most important libraries in python for data analysis. Pandas provide high performance, fast analysis, and data cleaning. We use it to manipulate data structures and have data analysis tools.

It consists of a data frame object. It interoperates with NumPy for faster computations. When we use both the libraries together it is a very helpful resource for scientific computations.

### **2. NumPy with Matplotlib**

Matplotlib is a module in NumPy. It is a very helpful tool to work with graphical representations. It consists of a wide range of functions to plot graphs and also manipulate them.

This combination can replace the functionalities of MatLab. It is used to generate the graphs of the results. We enhance it further with the use of graphic toolkits like PyQt and wxPython.

### **3. NumPy with SciPy**

Scipy is an open-source library in Python. It is the most important scientific library in python. It has been built upon the functionalities of NumPy. There are advanced functionalities in SciPy for scientific computations.

We can combine it with NumPy for greater mathematical performance. The combination helps in the implementation of complex scientific operations.

### ALGORITHMS:

#### A) Linear Regression:

It is one of the best statistical models that studies the relationship between a dependent variable (Y) with a given set of independent variables (X). The relationship can be established with the help of fitting a best line.

**sklearn.linear\_model.LinearRegression** is the module used to implement linear regression.

#### B) Logistic Regression:

Logistic regression, despite its name, is a classification algorithm rather than regression algorithm. Based on a given set of independent variables, it is used to estimate discrete value (0 or 1, yes/no, true/false). It is also called logit or MaxEnt Classifier.

Basically, it measures the relationship between the categorical dependent variable and one or more independent variables by estimating the probability of occurrence of an event using its logistics function.

**sklearn.linear\_model.LogisticRegression** is the module used to implement logistic regression.

#### C) K-NN:

k-NN (k-Nearest Neighbor), one of the simplest machine learning algorithms, is non-parametric and lazy in nature. Non-parametric means that there is no assumption for the underlying data distribution i.e. the model structure is determined from the dataset. Lazy or instance-based learning means that for the purpose of model generation, it does not require any training data points and whole training data is used in the testing phase.

**sklearn.neighbors.NearestNeighbors** is the module used to implement unsupervised nearest neighbor learning. It uses specific nearest neighbor algorithms named BallTree, KDTree or Brute Force. In other words, it acts as a uniform interface to these three algorithms.

#### D) Decision Tree:

Decision trees (DTs) are the most powerful non-parametric supervised learning method. They can be used for the classification and regression tasks. The main goal of DTs is to create a model predicting target variable value by learning simple decision rules deduced from the data features. Decision trees have two main entities; one is root node, where the data splits, and other is decision nodes or leaves, where we get final output.

Decision Tree Algorithms

Different Decision Tree algorithms are explained below –

### **ID3 algorithm stands for Iterative Dichotomiser 3**

It was developed by Ross Quinlan in 1986. It is also called Iterative Dichotomiser 3. The main goal of this algorithm is to find those categorical features, for every node, that will yield the largest information gain for categorical targets.

It lets the tree to be grown to their maximum size and then to improve the tree's ability on unseen data, applies a pruning step. The output of this algorithm would be a multiway tree.

## C4.5

It is the successor to ID3 and dynamically defines a discrete attribute that partition the continuous attribute value into a discrete set of intervals. That's the reason it removed the restriction of categorical features. It converts the ID3 trained tree into sets of 'IF-THEN' rules.

In order to determine the sequence in which these rules should applied, the accuracy of each rule will be evaluated first.

## C5.0

It works similar as C4.5 but it uses less memory and build smaller rulesets. It is more accurate than C4.5.

## CART

It is called Classification and Regression Trees algorithm. It basically generates binary splits by using the features and threshold yielding the largest information gain at each node (called the Gini index).

Homogeneity depends upon Gini index, higher the value of Gini index, higher would be the homogeneity. It is like C4.5 algorithm, but, the difference is that it does not compute rule sets and does not support numerical target variables (regression) as well.

## Classification with decision trees

In this case, the decision variables are categorical.

**Sklearn Module** – The Scikit-learn library provides the module name **DecisionTreeClassifier** for performing multiclass classification on dataset.

1

***criterion*** – string, optional default= "gini"

It represents the function to measure the quality of a split. Supported criteria are "gini" and "entropy". The default is gini which is for Gini impurity while entropy is for the information gain.

### E) KMeans:

This algorithm computes the centroids and iterates until it finds optimal centroid. It requires the number of clusters to be specified that's why it assumes that they are already known. The main logic of this algorithm is to cluster the data separating samples in n number of groups of equal variances by minimizing the criteria known as the inertia. The number of clusters identified by algorithm is represented by 'K'. Scikit-learn have **sklearn.cluster.KMeans** module to perform K-Means clustering. While computing cluster centers and value of inertia, the parameter

named `sample_weight` allows `sklearn.cluster.KMeans` module to assign more weight to some samples.

**F) SVM:**

Support Vector Machine or SVM is one of the most popular Supervised Learning algorithms, which is used for Classification as well as Regression problems. However, primarily, it is used for Classification problems in Machine Learning. The goal of the SVM algorithm is to create the best line or decision boundary that can segregate n-dimensional space into classes so that we can easily put the new data point in the correct category in the future. This best decision boundary is called a hyperplane. SVM chooses the extreme points/vectors that help in creating the hyperplane. These extreme cases are called as support vectors, and hence algorithm is termed as Support Vector Machine

**G) Random Forest:**

Random forest is a supervised learning algorithm which is used for both classification as well as regression. But however, it is mainly used for classification problems. As we know that a forest is made up of trees and more trees means more robust forest. Similarly, random forest algorithm creates decision trees on data samples and then gets the prediction from each of them and finally selects the best solution by means of voting. It is an ensemble method which is better than a single decision tree because it reduces the over-fitting by averaging the result.

**H) XGBoost:**

XGBoost is an optimized distributed gradient boosting library designed to be highly efficient, flexible and portable. It implements machine learning algorithms under the Gradient Boosting framework. XGBoost provides a parallel tree boosting (also known as GBDT, GBM) that solve many data science problems in a fast and accurate way.

**CONCLUSION:**

From this practical, I have successfully learned about libraries which are used for the machine learning applications.



A) Import NumPy as np

✓ [1] import numpy as np  
0s

B) Create an array of 10 zeros

✓ [2] np.zeros(10)  
0s  
array([0., 0., 0., 0., 0., 0., 0., 0., 0., 0.])

C) Create an array of 10 ones

✓ [3] np.ones(10)  
0s  
array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.])

D) Create an array of 10 fives

✓ [4] np.ones(10)\*5  
0s  
array([5., 5., 5., 5., 5., 5., 5., 5., 5., 5.])

E) Create an array of the integers from 10 to 50

✓ [5] np.arange(10,51)  
0s  
array([10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26,  
27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43,  
44, 45, 46, 47, 48, 49, 50])

F) Create an array of all the even integers from 10 to 50

✓ [6] np.arange(10,51,2)  
0s  
array([10, 12, 14, 16, 18, 20, 22, 24, 26, 28, 30, 32, 34, 36, 38, 40, 42,  
44, 46, 48, 50])

**G) Create a 3x3 matrix with values ranging from 0 to 8**

```
0s ✓ np.arange(0, 9).reshape(3,3)
↳ array([[0, 1, 2],
        [3, 4, 5],
        [6, 7, 8]])
```

**H) Create a 3x3 identity matrix**

```
0s ✓ [8] np.eye(3)
      array([[1., 0., 0.],
            [0., 1., 0.],
            [0., 0., 1.]])
```

**I) Use NumPy to generate a random number between 0 and 1**

```
0s ✓ [9] np.random.rand(1)
      array([0.63995357])
```

**J) Use NumPy to generate an array of 25 random numbers sampled from a standard normal distribution**

```
0s ✓ [10] np.random.randn(5,5)
      array([[ 0.00917627,  2.54296586,  0.05459417, -0.16396732,  0.48492437],
            [ 2.59090973, -0.4354634 ,  0.18801321,  1.21906203, -0.09548617],
            [ 0.10950666,  0.73144329,  1.55376488, -1.65816763, -0.30470597],
            [-0.66117415, -0.14516945, -0.00943456, -0.45216566,  1.15621409],
            [-0.13117192,  3.63209751,  0.67736721, -0.06525428,  0.5618076 ]])
```

**K) Create the following matrix:**

```
0s ✓ [11] np.arange(0.01, 1.01, 0.01).reshape(10,10)
      array([[0.01, 0.02, 0.03, 0.04, 0.05, 0.06, 0.07, 0.08, 0.09, 0.1 ],
            [0.11, 0.12, 0.13, 0.14, 0.15, 0.16, 0.17, 0.18, 0.19, 0.2 ],
            [0.21, 0.22, 0.23, 0.24, 0.25, 0.26, 0.27, 0.28, 0.29, 0.3 ],
            [0.31, 0.32, 0.33, 0.34, 0.35, 0.36, 0.37, 0.38, 0.39, 0.4 ],
            [0.41, 0.42, 0.43, 0.44, 0.45, 0.46, 0.47, 0.48, 0.49, 0.5 ],
            [0.51, 0.52, 0.53, 0.54, 0.55, 0.56, 0.57, 0.58, 0.59, 0.6 ],
            [0.61, 0.62, 0.63, 0.64, 0.65, 0.66, 0.67, 0.68, 0.69, 0.7 ],
            [0.71, 0.72, 0.73, 0.74, 0.75, 0.76, 0.77, 0.78, 0.79, 0.8 ],
            [0.81, 0.82, 0.83, 0.84, 0.85, 0.86, 0.87, 0.88, 0.89, 0.9 ],
            [0.91, 0.92, 0.93, 0.94, 0.95, 0.96, 0.97, 0.98, 0.99, 1.  ]])
```

**L) Create an array of 20 linearly spaced points between 0 and 1:**

```
✓ [12] np.linspace(0,1,20)
```

0s

```
array([0.          , 0.05263158, 0.10526316, 0.15789474, 0.21052632,  
       0.26315789, 0.31578947, 0.36842105, 0.42105263, 0.47368421,  
       0.52631579, 0.57894737, 0.63157895, 0.68421053, 0.73684211,  
       0.78947368, 0.84210526, 0.89473684, 0.94736842, 1.          ])
```

**Numpy Indexing and Selection****M) Now you will be given a few matrices, and be asked to replicate the resulting matrix outputs:**

```
✓ [13] mat = np.arange(1,26).reshape(5,5)  
      mat
```

0s

```
array([[ 1,  2,  3,  4,  5],  
       [ 6,  7,  8,  9, 10],  
       [11, 12, 13, 14, 15],  
       [16, 17, 18, 19, 20],  
       [21, 22, 23, 24, 25]])
```

```
✓ [15] mat1 = np.array([12,13,14,15,17,18,19,20,22,23,24,25]).reshape(3,4)  
      mat1
```

0s

```
array([[12, 13, 14, 15],  
       [17, 18, 19, 20],  
       [22, 23, 24, 25]])
```

```
✓ [17] mat1[1][3]
```

0s

```
20
```

```
✓ [19] mat2 = np.arange(2,13,5).reshape(3,1)  
      mat2
```

0s

```
array([[ 2],  
       [ 7],  
       [12]])
```

```
[21] mat3 = np.arange(21,26)
      mat3

array([21, 22, 23, 24, 25])
```

```
✓ [23] mat4 = np.arange(16,26).reshape(2,5)
0s      mat4

array([[16, 17, 18, 19, 20],
       [21, 22, 23, 24, 25]])
```

N) Get the sum of all the values in mat

```
✓ [24] np.sum(mat)
0s

325
```

O) Get the standard deviation of the values in mat

```
✓ [25] np.std(mat)
0s

7.211102550927978
```

P) Get the sum of all the columns in mat

```
✓ [26] np.sum(mat,axis=0)
0s

array([55, 60, 65, 70, 75])
```

### CONCLUSION:

From this practical, I have successfully learned about numpy library in python.

**PANDAS:****A) IMPORTING PANDAS & READING DATASET:**

```
import pandas as pd
```

```
sal = pd.read_csv('Salaries.csv')
```

```
sal.head()
```

|   |   | Id | EmployeeName      | JobTitle                                       | BasePay   | OvertimePay | OtherPay  | Benefits | TotalPay  | TotalPayBenefits | Year | Notes | Agency        | Status |
|---|---|----|-------------------|--|-----------|-------------|-----------|----------|-----------|------------------|------|-------|---------------|--------|
| 0 | 1 |    | NATHANIEL FORD    | GENERAL MANAGER-METROPOLITAN TRANSIT AUTHORITY | 167411.18 | 0.00        | 400184.25 | NaN      | 567595.43 | 567595.43        | 2011 | NaN   | San Francisco | NaN    |
| 1 | 2 |    | GARY JIMENEZ      | CAPTAIN III (POLICE DEPARTMENT)                | 155966.02 | 245131.88   | 137811.38 | NaN      | 538909.28 | 538909.28        | 2011 | NaN   | San Francisco | NaN    |
| 2 | 3 |    | ALBERT PARDINI    | CAPTAIN III (POLICE DEPARTMENT)                | 212739.13 | 106088.18   | 16452.60  | NaN      | 335279.91 | 335279.91        | 2011 | NaN   | San Francisco | NaN    |
| 3 | 4 |    | CHRISTOPHER CHONG | WIRE ROPE CABLE MAINTENANCE MECHANIC           | 77916.00  | 56120.71    | 198306.90 | NaN      | 332343.61 | 332343.61        | 2011 | NaN   | San Francisco | NaN    |
| 4 | 5 |    | PATRICK GARDNER   | DEPUTY CHIEF OF DEPARTMENT,(FIRE DEPARTMENT)   | 134401.60 | 9737.00     | 182234.59 | NaN      | 326373.19 | 326373.19        | 2011 | NaN   | San Francisco | NaN    |

**B) Use the .info() method to find out how many entries there are**

```
sal.info()
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 148654 entries, 0 to 148653
```

```
Data columns (total 13 columns):
```

| #  | Column           | Non-Null Count  | Dtype   |
|----|------------------|-----------------|---------|
| 0  | Id               | 148654 non-null | int64   |
| 1  | EmployeeName     | 148654 non-null | object  |
| 2  | JobTitle         | 148654 non-null | object  |
| 3  | BasePay          | 148045 non-null | float64 |
| 4  | OvertimePay      | 148650 non-null | float64 |
| 5  | OtherPay         | 148650 non-null | float64 |
| 6  | Benefits         | 112491 non-null | float64 |
| 7  | TotalPay         | 148654 non-null | float64 |
| 8  | TotalPayBenefits | 148654 non-null | float64 |
| 9  | Year             | 148654 non-null | int64   |
| 10 | Notes            | 0 non-null      | float64 |
| 11 | Agency           | 148654 non-null | object  |
| 12 | Status           | 0 non-null      | float64 |

```
dtypes: float64(8), int64(2), object(3)
```

```
memory usage: 14.7+ MB
```

C) What is the average BasePay ?

```
sal['BasePay'].mean()
```

```
66325.4488404877
```

D) What is the highest amount of OvertimePay in the dataset ?

```
sal['OvertimePay'].max()
```

```
245131.88
```

E) What is the job title of JOSEPH DRISCOLL ? Note: Use all caps, otherwise you may get an answer that doesn't match up (there is also a lowercase Joseph Driscoll).

```
sal[sal['EmployeeName']=='JOSEPH DRISCOLL']['JobTitle']
```

```
24    CAPTAIN, FIRE SUPPRESSION  
Name: JobTitle, dtype: object
```

F) How much does JOSEPH DRISCOLL make (including benefits)?

```
sal[sal['EmployeeName']=='JOSEPH DRISCOLL']['TotalPayBenefits']
```

```
24    270324.91  
Name: TotalPayBenefits, dtype: float64
```

G) What is the name of highest paid person (including benefits)?

```
ind = sal['TotalPayBenefits'].idxmax()  
sal.loc[ind]['EmployeeName']
```

```
'NATHANIEL FORD'
```

H) What is the name of lowest paid person (including benefits)? Do you notice something strange about how much he or she is paid?

```
ind = sal['TotalPayBenefits'].idxmin()  
sal.iloc[ind]
```

```
Id                148654  
EmployeeName      Joe Lopez  
JobTitle          Counselor, Log Cabin Ranch  
BasePay           0.0  
OvertimePay       0.0  
OtherPay          -618.13  
Benefits          0.0  
TotalPay          -618.13  
TotalPayBenefits  -618.13  
Year              2014  
Notes             NaN  
Agency           San Francisco  
Status            NaN  
Name: 148653, dtype: object
```

i) What was the average (mean) BasePay of all employees per year? (2011-2014) ?

```
sal.groupby('Year').mean()['BasePay']
```

```
Year  
2011    63595.956517  
2012    65436.406857  
2013    69630.030216  
2014    66564.421924  
Name: BasePay, dtype: float64
```

j) How many unique job titles are there?

```
sal['JobTitle'].nunique()
```

```
2159
```

k) What are the top 5 most common jobs?

```
sal['JobTitle'].value_counts().head()
```

```
Transit Operator    7036  
Special Nurse      4389  
Registered Nurse   3736  
Public Svc Aide-Public Works  2518  
Police Officer 3    2421  
Name: JobTitle, dtype: int64
```

- L) How many Job Titles were represented by only one person in 2013? (e.g. Job Titles with only one occurrence in 2013?)

```
(sal[sal['Year']==2013]['JobTitle'].value_counts()==1).sum()
```

202

- M) How many people have the word Chief in their job title?

```
def chief_string(title):  
    if 'chief' in title.lower().split():  
        return True  
    else:  
        return False  
sum(sal['JobTitle'].apply(lambda x:chief_string(x)))
```

477

- N) Is there a correlation between length of the Job Title string and Salary?

```
sal['title_len']=sal['JobTitle'].apply(len)  
sal[['TotalPayBenefits','title_len']].corr()
```

|                  | TotalPayBenefits | title_len |
|------------------|------------------|-----------|
| TotalPayBenefits | 1.000000         | -0.036878 |
| title_len        | -0.036878        | 1.000000  |



#### CONCLUSION:

From this practical, I have successfully learned about pandas library in python.



## **SCIKIT-LEARN:**

### **A) IMPORTING LIBRARIES, DATASET & READING DATASET:**

```
import sklearn
# load the iris dataset as an example
from sklearn.datasets import load_iris
iris = load_iris()

# store the feature matrix (X) and response vector (y)
X = iris.data
y = iris.target

# store the feature and target names
feature_names = iris.feature_names
target_names = iris.target_names

# printing features and target names of our dataset
print("Feature names:", feature_names)
print("Target names:", target_names)

# X and y are numpy arrays
print("\nType of X is:", type(X))

# printing first 5 input rows
print("\nFirst 5 rows of X:\n", X[:5])

Feature names: ['sepal length (cm)', 'sepal width (cm)', 'petal length (cm)', 'petal width (cm)']
Target names: ['setosa' 'versicolor' 'virginica']

Type of X is: <class 'numpy.ndarray'>

First 5 rows of X:
[[5.1 3.5 1.4 0.2]
 [4.9 3. 1.4 0.2]
 [4.7 3.2 1.3 0.2]
 [4.6 3.1 1.5 0.2]
 [5. 3.6 1.4 0.2]]
```

### **B) SPLITTING DATASET INTO TRAINING & TESTING DATASET:**

```
# splitting X and y into training and testing sets
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.4, random_state=1)

# printing the shapes of the new X objects
print('X Train Data: ',X_train.shape)
print('X Test Data: ',X_test.shape)

# printing the shapes of the new y objects
print('Y Train Data: ',y_train.shape)
print('Y Test Data: ',y_test.shape)
```

```
X Train Data: (90, 4)
X Test Data: (60, 4)
Y Train Data: (90,)
Y Test Data: (60,)
```

**C) NORMALIZATION:**

```
# import module
from sklearn.preprocessing import StandardScaler
```

```
# create data
data = [[11, 2], [3, 7], [0, 10], [11, 8]]
```

```
# compute required values
scaler = StandardScaler()
model = scaler.fit(data)
scaled_data = model.transform(data)
```

```
# print scaled data
print(scaled_data)
```

```
[[ 0.97596444 -1.61155897]
 [-0.66776515  0.08481889]
 [-1.28416374  1.10264561]
 [ 0.97596444  0.42409446]]
```

**SCIPY:**

**A) SPARSE MATRIX:**

```
# import necessary modules
from scipy import sparse
# Row-based linked list sparse matrix
A = sparse.lil_matrix((1000, 1000))
print(A)
```

```
A[0,:100] = np.random.rand(100)
A[1,100:200] = A[0,:100]
A.setdiag(np.random.rand(1000))
print(A)
```

|         |                      |
|---------|----------------------|
| (0, 0)  | 0.19036855447545786  |
| (0, 1)  | 0.33776231739432405  |
| (0, 2)  | 0.9749038090665604   |
| (0, 3)  | 0.8942583911186659   |
| (0, 4)  | 0.08515061548006031  |
| (0, 5)  | 0.21431379796449723  |
| (0, 6)  | 0.8676941823129108   |
| (0, 7)  | 0.07258083664498649  |
| (0, 8)  | 0.5557847427639602   |
| (0, 9)  | 0.7459339793334242   |
| (0, 10) | 0.3857750726237723   |
| (0, 11) | 0.8832801032093758   |
| (0, 12) | 0.6813002996002154   |
| (0, 13) | 0.6296641368561785   |
| (0, 14) | 0.16101650416423774  |
| (0, 15) | 0.4873200634402388   |
| (0, 16) | 0.1806188884055946   |
| (0, 17) | 0.28177347650237183  |
| (0, 18) | 0.19815175912132565  |
| (0, 19) | 0.4342894381136917   |
| (0, 20) | 0.6364320992619823   |
| (0, 21) | 0.21080807130051238  |
| (0, 22) | 0.020340881969629243 |
| (0, 23) | 0.6116070621487552   |
| (0, 24) | 0.9183747764032624   |
| (0, 25) | 0.031451450420754035 |
| (0, 26) | 0.3871320802870829   |

**B) CONVERT THIS MATRIX TO COMPRESSED SPARSE ROW FORMAT**

from scipy.sparse import linalg

# Convert this matrix to Compressed Sparse Row format.

A.tocsr()

A = A.tocsr()

b = np.random.rand(1000)

ans = linalg.spsolve(A, b)

# it will print ans array of 1000 size

print(ans)

```

-2.35319469e+02 -1.49842949e+02 9.73797985e-01 8.01662749e-01
2.56669454e-01 3.23166234e-01 2.86738795e-01 1.84657087e+00
4.48405102e-01 9.84638653e-02 4.44847188e-01 2.88717900e-01
4.44918552e-02 4.58318163e-01 1.43384483e+00 1.99350174e+00
1.21888495e+00 1.57496174e+00 6.70136481e-01 3.59897327e+01
1.28297548e+00 1.96483683e+00 3.49994061e+00 5.23922832e-01
2.06604672e+00 1.84576603e+00 1.56911917e-01 1.64625356e+00
1.54586678e+00 3.70630437e-01 9.22407476e-01 1.25307475e+01
1.27201263e+00 8.69771911e-01 5.71008317e-01 2.36020605e+00
2.52519006e+00 2.86415071e-01 4.45305112e-01 1.88676711e+00
9.66215505e-01 5.54360157e-01 5.13487903e+00 1.35910055e+00
1.65365834e+00 6.86115676e+00 1.43536015e+00 1.22939297e+00
3.39476772e-01 1.33685644e+00 3.62356948e+00 2.53105349e+00
1.06194083e+02 4.53304440e+00 3.69673226e-01 1.44402506e+00
9.76192380e-01 7.12327777e-01 1.37028728e+00 9.17760860e-01
6.89348469e-01 2.13957351e+00 9.28235879e+00 9.67851414e-01
2.79828865e-01 4.13236913e-01 1.06858116e+00 1.79529384e+00
6.55611974e+00 8.26922782e-01 3.55375831e+01 5.22998358e-02
8.53147778e-01 6.09557326e-01 1.60171443e+00 1.13057652e+00
1.99948661e+00 3.45236527e-01 5.45060752e+00 9.86027796e-01
5.54949985e-02 1.34151479e+00 6.30784392e-01 1.93189946e+00
2.39506893e+01 1.89743431e-01 3.16347084e-01 8.24358555e-02
4.38889999e-01 2.73188918e+00 7.30921368e-01 1.07383848e+00
4.00896310e-02 1.03398614e+00 7.16866466e-01 3.02780953e-02
2.44974913e+00 1.46451142e+00 2.79801309e-01 2.32484949e+00
3.28744365e-01 1.37912374e-01 1.71207446e-01 1.45246897e-01
9.09435792e-01 1.22430643e+00 2.83132823e+00 7.70920938e-01

```

**C) DETERMINANT:**

```

# import numpy library
import numpy as np
A = np.array([[1,2,3],[4,5,6],[7,8,9]])
# importing linalg function from scipy
from scipy import linalg

# Compute the determinant of a matrix
linalg.det(A)

```

3.0

**D) DOT PRODUCT:**

```

P, L, U = linalg.lu(A)
print(P)
print(L)
print(U)
# print LU decomposition
print(np.dot(L,U))

```

```
[[0. 1. 0.]
 [0. 0. 1.]
 [1. 0. 0.]]
[[1.      0.      0.      ]
 [0.14285714 1.      0.      ]
 [0.57142857 0.5     1.      ]]
[[7.      8.      8.      ]
 [0.      0.85714286 1.85714286]
 [0.      0.      0.5     ]]
[[7. 8. 8.]
 [1. 2. 3.]
 [4. 5. 6.]]
```

#### E) EIGEN VECTOR:

```
eigen_values, eigen_vectors = linalg.eig(A)
print(eigen_values)
print(eigen_vectors)
```

```
[15.55528261+0.j -1.41940876+0.j -0.13587385+0.j]
[[-0.24043423 -0.67468642  0.51853459]
 [-0.54694322 -0.23391616 -0.78895962]
 [-0.80190056  0.70005819  0.32964312]]
```

#### F) INTEGRATION:

```
from scipy import integrate
f = lambda y, x: x*y**2
i = integrate.dblquad(f, 0, 2, lambda x: 0, lambda x: 1)
# print the results
print(i)
```

```
(0.6666666666666667, 7.401486830834377e-15)
```

#### CONCLUSION:

From this practical, I have learned and implemented scikit-learn & scipy libraries in python.

## MATPLOTLIB:

### A) IMPORT MATPLOTLIB:

```
import numpy as np
x = np.arange(0,100)
y = x*2
z = x**2

import matplotlib.pyplot as plt
%matplotlib inline
```

### B) Follow along with these steps:

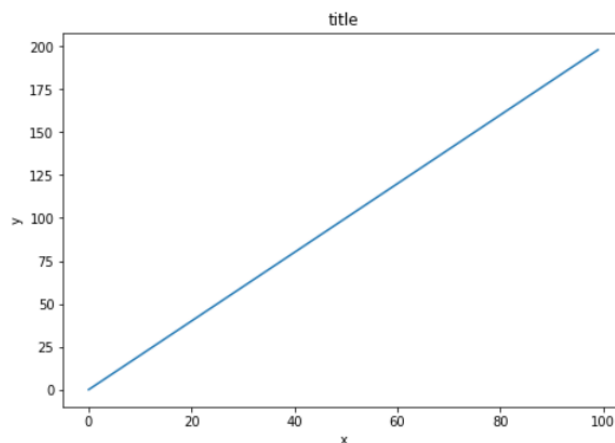
Create a figure object called fig using plt.figure()

Use add\_axes to add an axis to the figure canvas at [0,0,1,1]. Call this new axis ax.

Plot (x,y) on that axes and set the labels and titles to match the plot below:

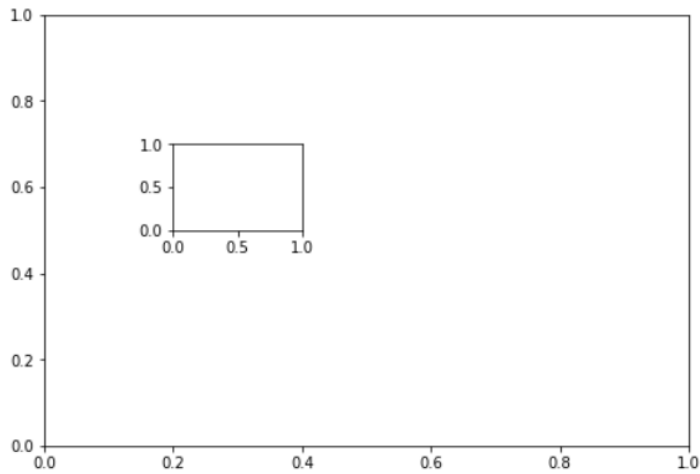
```
fig = plt.figure()
ax= fig.add_axes([0,0,1,1])
ax.plot(x,y)
ax.set_title('title')
ax.set_xlabel('x')
ax.set_ylabel('y')
```

Text(0, 0.5, 'y')



### C) Create a figure object and put two axes on it, ax1 and ax2. Located at [0,0,1,1] and [0.2,0.5,.2,.2] respectively.

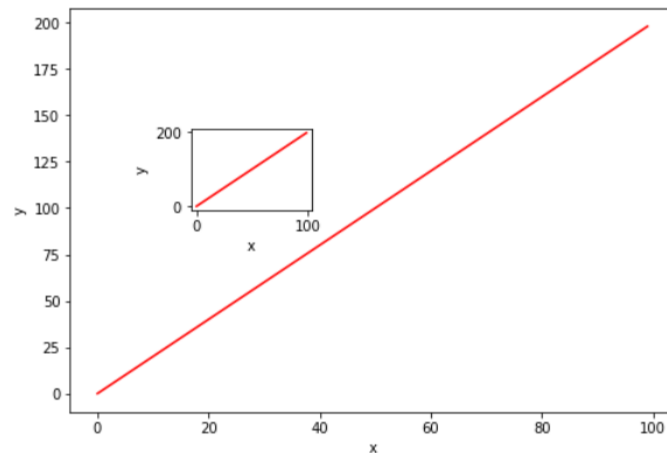
```
fig = plt.figure()
ax1 = fig.add_axes([0,0,1,1])
ax2 = fig.add_axes([0.2,0.5,.2,.2])
```



D) Now plot (x,y) on both axes. And call your figure object to show it.

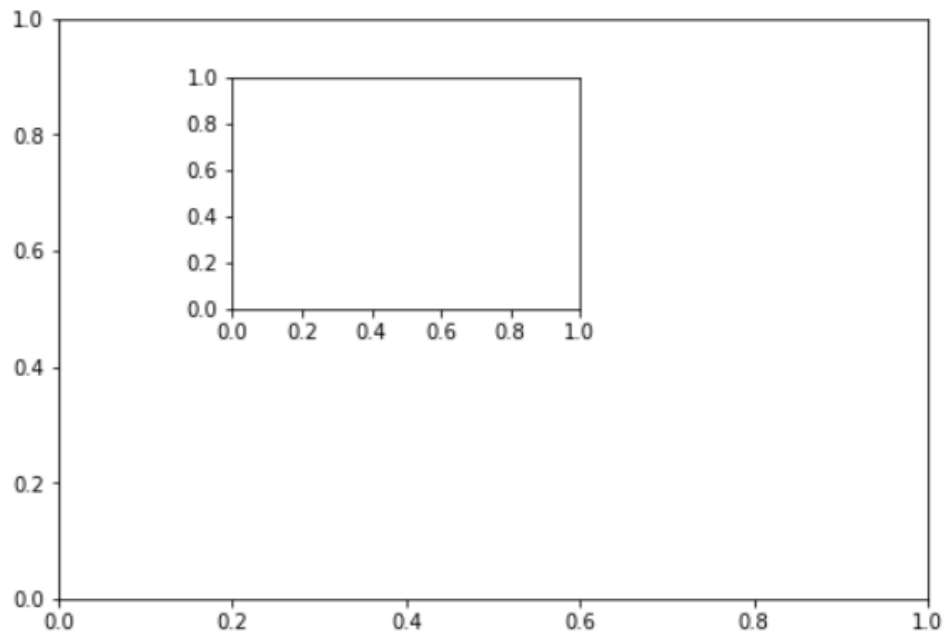
```
fig = plt.figure()
ax1 = fig.add_axes([0,0,1,1])
ax1.plot(x,y,'red')
ax1.set_xlabel('x')
ax1.set_ylabel('y')
ax2 = fig.add_axes([0.2,0.5,.2,.2])
ax2.plot(x,y,'red')
ax2.set_xlabel('x')
ax2.set_ylabel('y')
```

Text(0, 0.5, 'y')



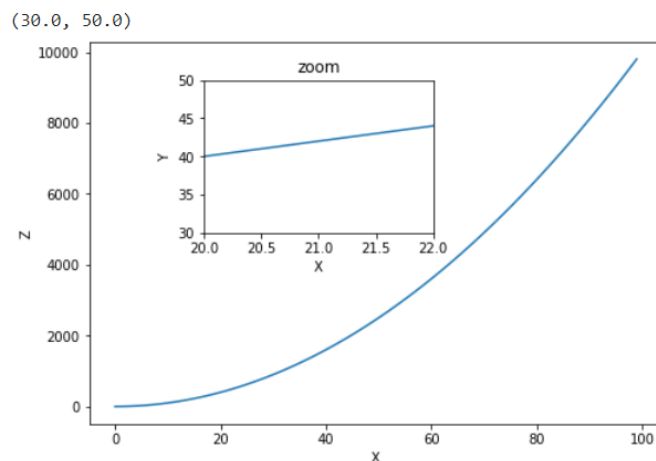
E) Create the plot below by adding two axes to a figure object at [0,0,1,1] and [0.2,0.5,.4,.4]

```
fig = plt.figure()
ax1 = fig.add_axes([0,0,1,1])
ax2 = fig.add_axes([0.2,0.5,.4,.4])
```



F) Now use x,y, and z arrays to recreate the plot below. Notice the xlims and y limits on the inserted plot:

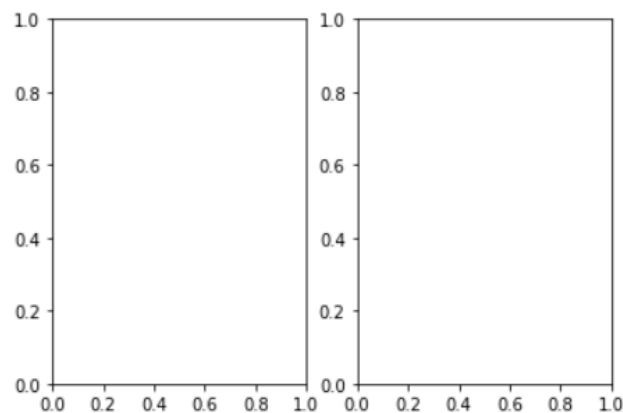
```
fig = plt.figure()
ax1 = fig.add_axes([0,0,1,1])
ax1.plot(x,z)
ax1.set_xlabel('X')
ax1.set_ylabel('Z')
ax2 = fig.add_axes([0.2,0.5,.4,.4])
ax2.plot(x,y)
ax2.set_title('zoom')
ax2.set_xlabel('X')
ax2.set_ylabel('Y')
ax2.set_xlim(20,22)
ax2.set_ylim(30,50)
```





- G) Use `plt.subplots(nrows=1, ncols=2)` to create the plot below.

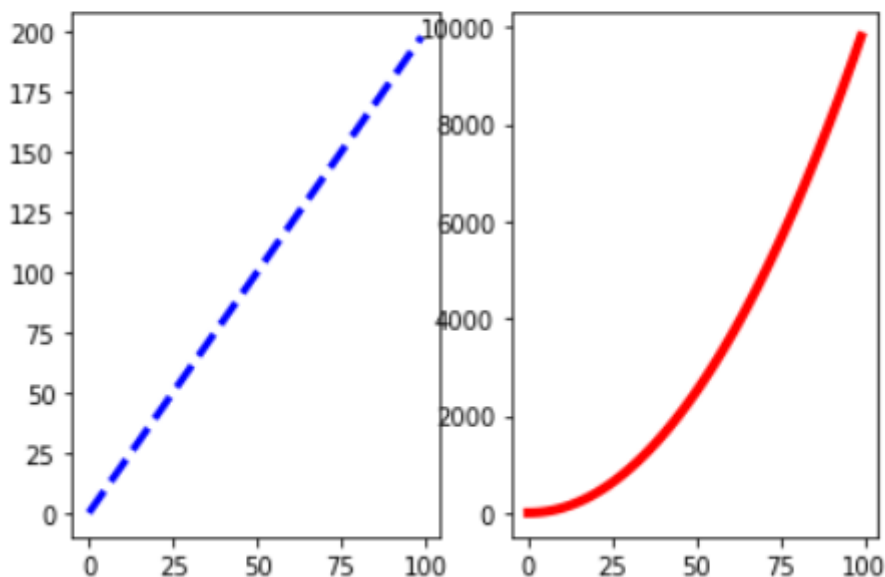
```
fig, axes = plt.subplots(nrows=1, ncols=2)
```



- H) Now plot  $(x,y)$  and  $(x,z)$  on the axes. Play around with the linewidth and style

```
fig, axes = plt.subplots(nrows=1, ncols=2)
axes[0].plot(x,y,color='blue',lw=3,ls='--')
axes[1].plot(x,z,color='red',lw=4)
```

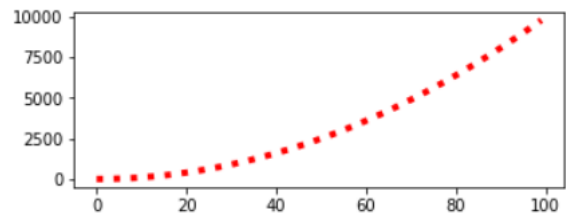
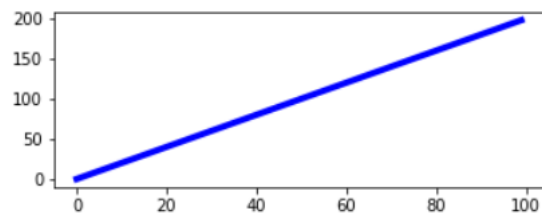
[<matplotlib.lines.Line2D at 0x7f24a91d71d0>]



- I) See if you can resize the plot by adding the `figsize()` argument in `plt.subplots()` are copying and pasting your previous code.

```
fig, axes = plt.subplots(figsize=(12,2),nrows=1, ncols=2)
axes[0].plot(x,y,color='blue',lw=4,ls='-')
axes[1].plot(x,z,color='red',lw=4,ls=':')
```

[<matplotlib.lines.Line2D at 0x7f24a88f6590>]



### CONCLUSION:

From this practical, I have successfully learned about matplotlib library in python.

**SEABORN:****A) IMPORTING LIBRARIES AND DATA:**

```
import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline
```

```
sns.set_style('whitegrid')
```

```
titanic = sns.load_dataset('titanic')
```

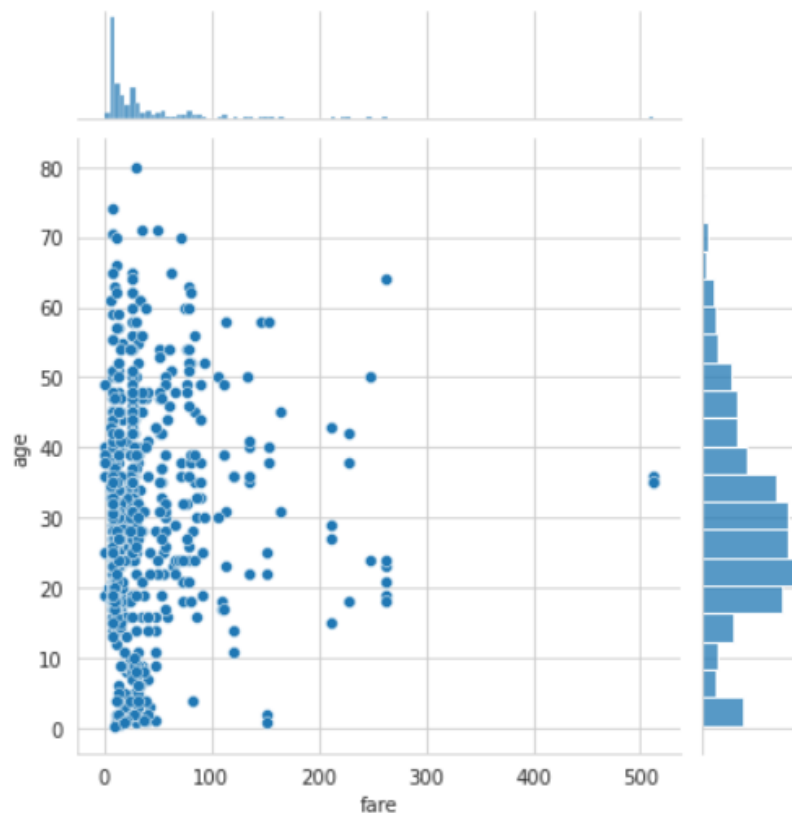
```
titanic.head()
```

|   | survived | pclass | sex    | age  | sibsp | parch | fare    | embarked | class | who   | adult_male | deck | embark_town | alive | alone |
|---|----------|--------|--------|------|-------|-------|---------|----------|-------|-------|------------|------|-------------|-------|-------|
| 0 | 0        | 3      | male   | 22.0 | 1     | 0     | 7.2500  | S        | Third | man   | True       | NaN  | Southampton | no    | False |
| 1 | 1        | 1      | female | 38.0 | 1     | 0     | 71.2833 | C        | First | woman | False      | C    | Cherbourg   | yes   | False |
| 2 | 1        | 3      | female | 26.0 | 0     | 0     | 7.9250  | S        | Third | woman | False      | NaN  | Southampton | yes   | True  |
| 3 | 1        | 1      | female | 35.0 | 1     | 0     | 53.1000 | S        | First | woman | False      | C    | Southampton | yes   | False |
| 4 | 0        | 3      | male   | 35.0 | 0     | 0     | 8.0500  | S        | Third | man   | True       | NaN  | Southampton | no    | True  |

**B) JOINPLOT [FARE V/S AGE]:**

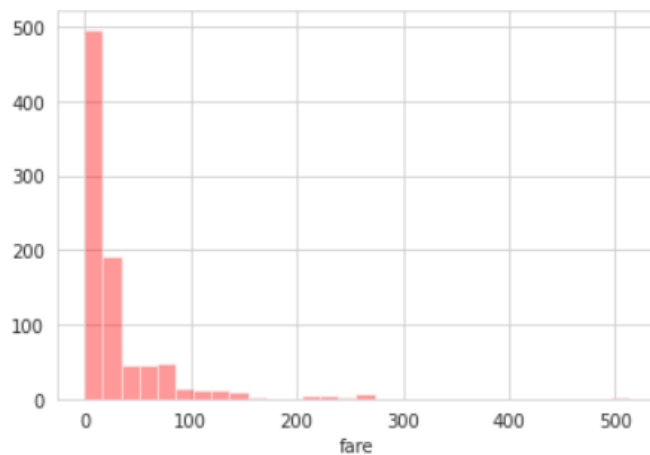
```
sns.jointplot(x='fare',y='age',data=titanic)
```

```
<seaborn.axisgrid.JointGrid at 0x7fd921864090>
```



C) DISPLOT:

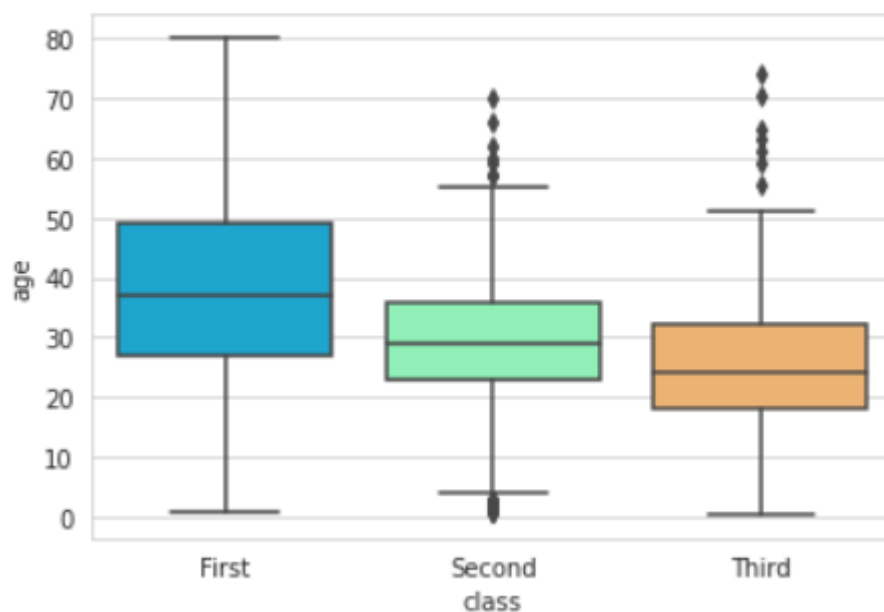
```
sns.distplot(titanic['fare'],bins=30,kde=False,color='red')  
  
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py  
warnings.warn(msg, FutureWarning)  
<matplotlib.axes._subplots.AxesSubplot at 0x7fd91d018550>
```



D) BOXPLOT:

```
sns.boxplot(x='class',y='age',data=titanic,palette='rainbow' )
```

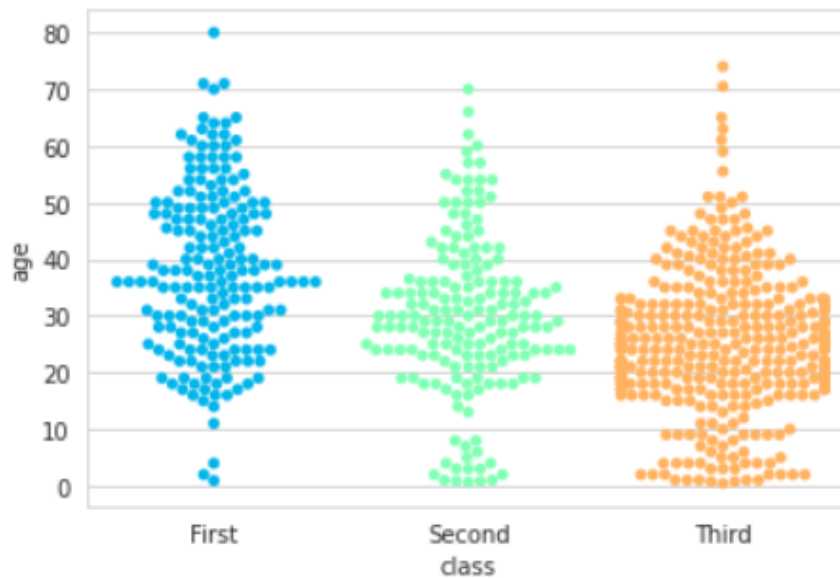
```
<matplotlib.axes._subplots.AxesSubplot at 0x7fd91cf76cd0>
```



E) SWARMPLOT [CLASS V/S AGE]:

```
sns.swarmplot(x='class',y='age',data=titanic,palette='rainbow')
```

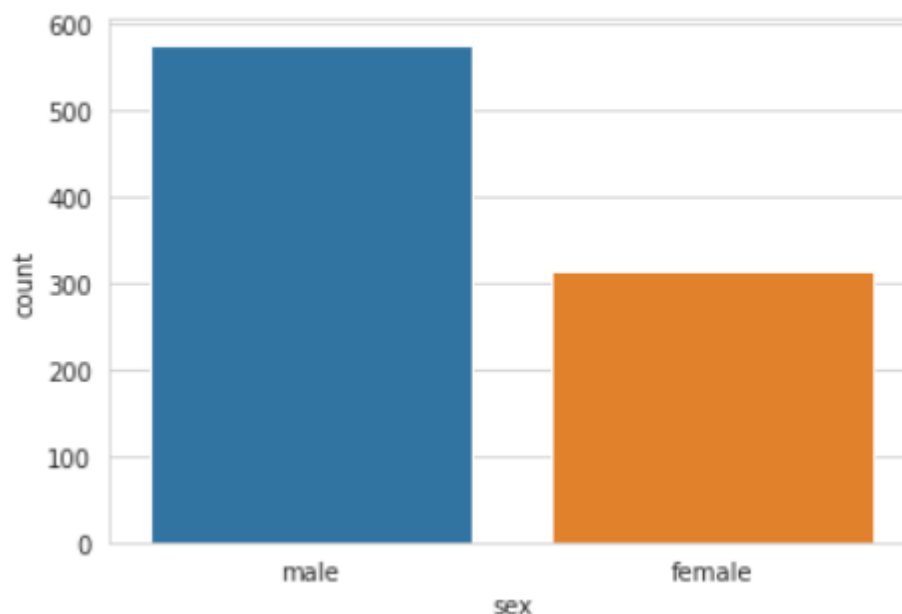
```
/usr/local/lib/python3.7/dist-packages/seaborn/categorical.py:1296:  
warnings.warn(msg, UserWarning)  
<matplotlib.axes._subplots.AxesSubplot at 0x7fd91ce78f90>
```



F) COUNTPLOT:

```
sns.countplot(x='sex',data=titanic)
```

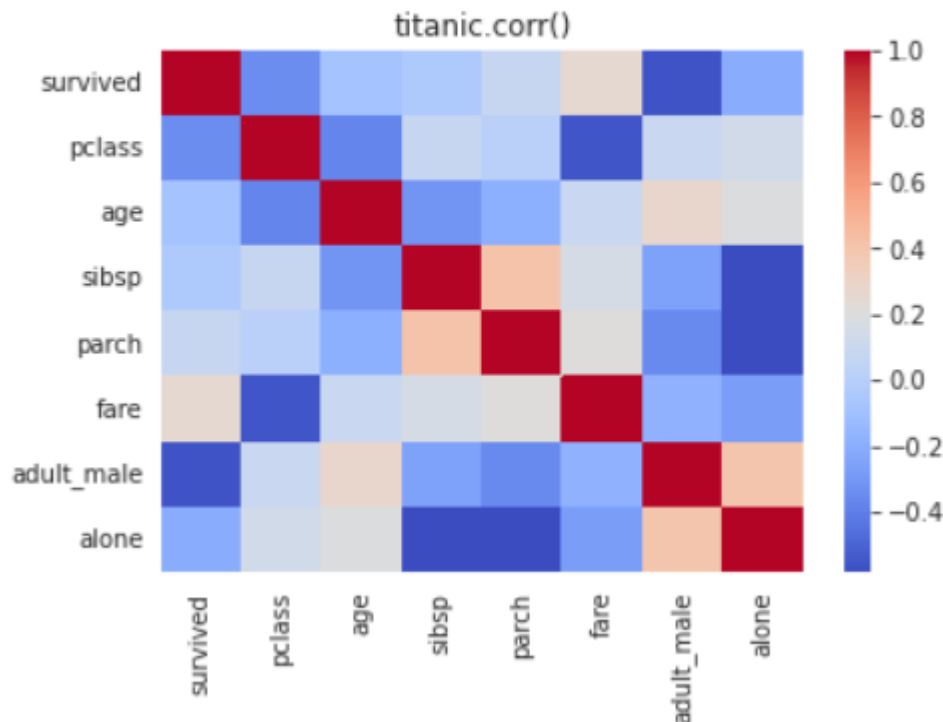
```
<matplotlib.axes._subplots.AxesSubplot at 0x7fd91ce78fd0>
```



G) HEATMAP:

```
sns.heatmap(titanic.corr(),cmap='coolwarm')
plt.title('titanic.corr()')
```

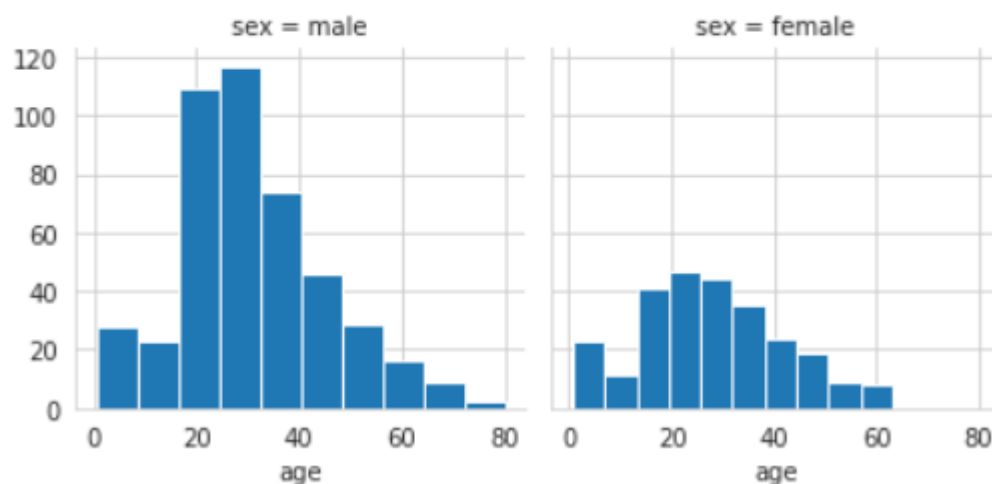
```
Text(0.5, 1.0, 'titanic.corr()')
```



H) FACEGRID:

```
g = sns.FacetGrid(data=titanic,col='sex')
g.map(plt.hist,'age')
```

```
<seaborn.axisgrid.FacetGrid at 0x7fd91cd04450>
```



CONCLUSION:

From this practical, I have successfully learned about seaborn library in python.