#### 9.2. Data Science – Machine Learning – Linear Regression Example

```
Program
           Loading salary dataset
           demo1.py
Name
           import pandas as pd
           dataset = pd.read_csv('Salary_Data.csv')
           print(dataset.head())
Output
               YearsExperience
                                 Salary
                            1.1
                                 39343.0
                           1.3
                                 46205.0
                                 37731.0
                                 43525.0
                                 39891.0
                            2.2
```

```
Program
               Preparing the data
               demo2.py
Name
               import pandas as pd
               dataset = pd.read_csv('Salary_Data.csv')
               X = dataset.iloc[:, :-1].values
               y = dataset.iloc[:, 1].values
               print(X)
               print(y)
Output
                          46205. 37731. 43525. 39891. 56642. 60150. 54445. 64445.
                  57189. 63218. 55794. 56957. 57081. 61111. 67938. 66029. 83088. 81363. 93940. 91738. 98273. 101302. 113812. 109431. 105582. 116969.
                 112635. 122391. 121872.]
```

```
Splitting the dataset
Program
            demo3.py
Name
            import pandas as pd
            from sklearn.model selection import train test split
            dataset = pd.read_csv('Salary_Data.csv')
            X = dataset.iloc[:, :-1].values
            y = dataset.iloc[:, 1].values
            X_train, X_test, y_train, y_test = train_test_split(X, y, test_size =
            1/3, random_state = 0)
            print("X_train")
            print(X_train)
            print()
            print("X_test")
            print(X_test)
            print()
            print("y_train")
            print(y_train)
            print()
            print("y_test")
            print(y_test)
```

## Output

```
X train
[[ 2.9]
[ 5.1]
[ 3.2]
[ 4.5]
[ 8.2]
[ 6.8]
[ 1.3]
[ 10.5]
[ 3. ]
[ 2.2]
[ 5.9]
[ 6. ]
[ 3.7]
[ 3.2]
[ 9. ]
[ 2. ]
[ 1.1]
[ 7.1]
[ 4. 9]
[ 4. ]]

X test
[[ 1.5]
[ 10.3]
[ 4.1]
[ 3.9]
[ 9.5]
[ 8.7]
[ 9.6]
[ 4. ]
[ 5.3]
[ 7.9]]

y train
[ 56642. 66029. 64445. 61111. 113812. 91738. 46205. 121872. 60150. 39801. 81363. 93940. 57189. 54445. 105582. 43525. 39343. 98273. 67938. 56957.]

y test
[ 37731. 122391. 57081. 63218. 116969. 109431. 112635. 55794. 83088. 101302.]
```

### Program Name

Training the model

demo4.py

import pandas as pd

from sklearn.model\_selection import train\_test\_split
from sklearn.linear\_model import LinearRegression

dataset = pd.read\_csv('Salary\_Data.csv')

X = dataset.iloc[:, :-1].values
y = dataset.iloc[:, 1].values

X\_train, X\_test, y\_train, y\_test = train\_test\_split(X, y, test\_size = 1/3, random state = 0)

print("Training the model")

regressor = LinearRegression()
regressor.fit(X\_train, y\_train)

#### Output

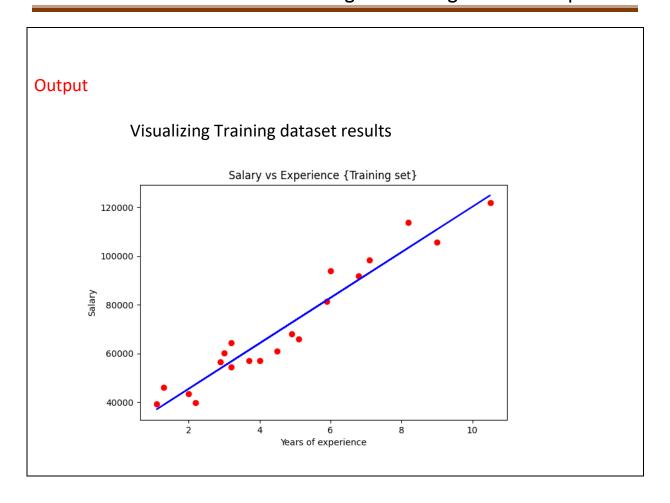
Training the model

```
Program
            Predicting the salaries
Name
            demo5.py
            import pandas as pd
            from sklearn.model selection import train test split
            from sklearn.linear_model import LinearRegression
            dataset = pd.read_csv('Salary_Data.csv')
            X = dataset.iloc[:, :-1].values
            y = dataset.iloc[:, 1].values
            X_train, X_test, y_train, y_test = train_test_split(X, y, test_size =
            1/3, random state = 0)
            regressor = LinearRegression()
            regressor.fit(X_train, y_train)
            print("Predicting the salaries")
            y_pred = regressor.predict(X_test)
            print()
            print(y_pred)
Output
```

# Predicting the salaries

```
[ 40835.10590871 123079.39940819 65134.55626083 63265.36777221 115602.64545369 108125.8914992 116537.23969801 64199.96201652 76349.68719258 100649.1375447 ]
```

```
Program
             Plotting training dataset
Name
             demo6.py
             import pandas as pd
             import matplotlib.pyplot as plt
             from sklearn.model_selection import train_test_split
             from sklearn.linear model import LinearRegression
             dataset = pd.read csv('Salary Data.csv')
             X = dataset.iloc[:, :-1].values
             y = dataset.iloc[:, 1].values
             X_train, X_test, y_train, y_test = train_test_split(X, y, test_size =
             1/3, random_state = 0)
             regressor = LinearRegression()
             regressor.fit(X_train, y_train)
             y_pred = regressor.predict(X_test)
             print("Visualizing Training dataset results ")
             plt.scatter(X_train, y_train, color = 'red')
             plt.plot(X_train, regressor.predict(X_train), color = 'blue')
             plt.title('Salary vs Experience {Training set}')
             plt.xlabel('Years of experience')
             plt.ylabel('Salary')
             plt.show()
```



```
Program
            Plotting test dataset
Name
            demo7.py
            import pandas as pd
            import matplotlib.pyplot as plt
            from sklearn.model_selection import train_test_split
            from sklearn.linear model import LinearRegression
            dataset = pd.read csv('Salary Data.csv')
            X = dataset.iloc[:, :-1].values
            y = dataset.iloc[:, 1].values
            X_train, X_test, y_train, y_test = train_test_split(X, y, test_size =
            1/3, random_state = 0)
            regressor = LinearRegression()
            regressor.fit(X_train, y_train)
            y_pred = regressor.predict(X_test)
            print("Visualizing Training dataset results ")
            plt.scatter(X_test, y_test, color = 'red')
            plt.plot(X_train, regressor.predict(X_train), color = 'blue')
            plt.title('Salary vs Experience {Test set}')
            plt.xlabel('Years of experience')
            plt.ylabel('Salary')
            plt.show()
```

