

# Scikit\_Assignment 01\_Solution

July 3, 2020

## 1 Assignment 01: Evaluate the Ad Budget Dataset of XYZ Firm

*The comments/sections provided are your cues to perform the assignment. You don't need to limit yourself to the number of rows/cells provided. You can add additional rows in each section to add more lines of code.*

*If at any point in time you need help on solving this assignment, view our demo video to understand the different steps of the code.*

**Happy coding!**

---

### 1: Import the dataset

```
[7]: #Import the required libraries
import pandas as pd
```

```
[2]: #Import the advertising dataset
df_adv_dataset = pd.read_csv("Advertising Budget and Sales.csv", index_col=0)
```

### 2: Analyze the dataset

```
[3]: #View the initial few records of the dataset
df_adv_dataset.head()
```

```
[3]:
```

	TV Ad Budget (\$)	Radio Ad Budget (\$)	Newspaper Ad Budget (\$)	Sales (\$)
1	230.1	37.8	69.2	22.1
2	44.5	39.3	45.1	10.4
3	17.2	45.9	69.3	9.3
4	151.5	41.3	58.5	18.5
5	180.8	10.8	58.4	12.9

```
[9]: #Check the total number of elements in the dataset
df_adv_dataset.size
df_adv_dataset.shape
```

```
[9]: 800
```

### 3: Find the features or media channels used by the firm

```
[10]: #Check the number of observations (rows) and attributes (columns) in the dataset
df_adv_dataset.columns
```

```
[10]: Index(['TV Ad Budget ($)', 'Radio Ad Budget ($)', 'Newspaper Ad Budget ($)',
        'Sales ($)'],
        dtype='object')
```

```
[13]: #View the names of each of the attributes
df_adv_dataset[["Newspaper Ad Budget ($)", "Radio Ad Budget ($)", "TV Ad Budget_
↪ ($)"]]
```

```
[13]:      Newspaper Ad Budget ($)  Radio Ad Budget ($)  TV Ad Budget ($)
1              69.2              37.8             230.1
2              45.1              39.3              44.5
3              69.3              45.9              17.2
4              58.5              41.3             151.5
5              58.4              10.8             180.8
..              ...              ...              ...
196             13.8               3.7              38.2
197              8.1               4.9              94.2
198              6.4               9.3             177.0
199             66.2             42.0             283.6
200              8.7               8.6             232.1
```

[200 rows x 3 columns]

### 4: Create objects to train and test the model; find the sales figures for each channel

```
[27]: #Create a feature object from the columns
X_feature = df_adv_dataset[["Newspaper Ad Budget ($)", "Radio Ad Budget ($)", "TV_
↪ Ad Budget ($)"]]
```

```
[28]: #View the feature object
X_feature.head()
```

```
[28]:      Newspaper Ad Budget ($)  Radio Ad Budget ($)  TV Ad Budget ($)
1              69.2              37.8             230.1
2              45.1              39.3              44.5
3              69.3              45.9              17.2
4              58.5              41.3             151.5
5              58.4              10.8             180.8
```

```
[16]: #Create a target object (Hint: use the sales column as it is the response of_
↪ the dataset)
Y_target = df_adv_dataset[['Sales ($)']]
```

```
[17]: #View the target object
      Y_target.head()
```

```
[17]:    Sales ($)
      1      22.1
      2      10.4
      3       9.3
      4      18.5
      5      12.9
```

```
[29]: #Verify if all the observations have been captured in the feature object
      X_feature.shape
```

```
[29]: (200, 3)
```

```
[20]: #Verify if all the observations have been captured in the target object
      Y_target.shape
```

```
[20]: (200, 1)
```

## 5: Split the original dataset into training and testing datasets for the model

```
[30]: #Split the dataset (by default, 75% is the training data and 25% is the testing
      ↪data)
      from sklearn.model_selection import train_test_split
      x_train, x_test, y_train, y_test = train_test_split(X_feature, Y_target,
      ↪random_state = 1)
```

```
[31]: #Verify if the training and testing datasets are split correctly (Hint: use the
      ↪shape() method)
      print(x_train.shape)
      print(x_test.shape)
      print(y_train.shape)
      print(y_test.shape)
```

```
(150, 3)
(50, 3)
(150, 1)
(50, 1)
```

## 6: Create a model to predict the sales outcome

```
[38]: #Create a linear regression model
      from sklearn.linear_model import LinearRegression
      linreg = LinearRegression()
      linreg.fit(x_train,y_train)
```

```
[38]: LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None, normalize=False)
```

```
[39]: #Print the intercept and coefficients  
print(linreg.intercept_)  
print(linreg.coef_)
```

```
[2.87696662]  
[[0.00345046 0.17915812 0.04656457]]
```

```
[40]: #Predict the outcome for the testing dataset  
y_predict = linreg.predict(x_test)  
y_predict
```

```
[40]: array([[21.70910292],  
            [16.41055243],  
            [ 7.60955058],  
            [17.80769552],  
            [18.6146359 ],  
            [23.83573998],  
            [16.32488681],  
            [13.43225536],  
            [ 9.17173403],  
            [17.333853  ],  
            [14.44479482],  
            [ 9.83511973],  
            [17.18797614],  
            [16.73086831],  
            [15.05529391],  
            [15.61434433],  
            [12.42541574],  
            [17.17716376],  
            [11.08827566],  
            [18.00537501],  
            [ 9.28438889],  
            [12.98458458],  
            [ 8.79950614],  
            [10.42382499],  
            [11.3846456 ],  
            [14.98082512],  
            [ 9.78853268],  
            [19.39643187],  
            [18.18099936],  
            [17.12807566],  
            [21.54670213],  
            [14.69809481],  
            [16.24641438],  
            [12.32114579],
```

```
[19.92422501],  
[15.32498602],  
[13.88726522],  
[10.03162255],  
[20.93105915],  
[ 7.44936831],  
[ 3.64695761],  
[ 7.22020178],  
[ 5.9962782 ],  
[18.43381853],  
[ 8.39408045],  
[14.08371047],  
[15.02195699],  
[20.35836418],  
[20.57036347],  
[19.60636679]])
```

## 7: Calculate the Mean Square Error (MSE)

```
[47]: #Import required libraries for calculating MSE (mean square error)  
      from sklearn import metrics  
      import numpy as np
```

```
[49]: #Calculate the MSE  
      print(np.sqrt(metrics.mean_squared_error(y_test,y_predict)))
```

1.404651423032897

```
[ ]:
```