# MODULE 4: Linear Regression - Assignment

In this assignment, we will show our understanding on how to apply the Linear Regression model to a selected dataset in Python.

Please, find the dataset on Canvas in the FILES section.

Description of the `advertising` data:

Advertising data sales (in thousands of units) for a particular product advertising budget (in thousands of dollars) for TV.

## STEP 1: Pre-process the dataset

In this step, you need to do the following:

- Import the necessary libraries
- Load the dataset
- Explore the dataset
- Summarize the dataset
- Analyze the dataset

## Import the necessary libraries

```
1 #Import the libraries
2 import pandas as pd
3 import numpy as py
4 import seaborn as sn
5 import matplotlib.pyplot as plt
6
```

## Load the dataset

Comments:

```
1 #Load dataset
2 df = pd.read_csv("/content/advertising.csv")
```

## ⌄ Explore the dataset

**Note:**

Here we can also check outliers, missing values, etc. You can apply all the DATA
CLEANING skills that you learn in the Advanced Python course.

```
1 #Show the head of the data
2 df.head()
3
```

|   | TV | Sales |
|---|------|-------|
| 0 | 230.1 | 22.1 |
| 1 | 44.5 | 10.4 |
| 2 | 17.2 | 12.0 |
| 3 | 151.5 | 16.5 |
| 4 | 180.8 | 17.9 |

Next steps:    [ New interactive sheet ]

```
1 #Show the shape of the data
2 df.shape
3
```

```
(200, 2)
```

```
1 #Get the info of the dataset with "advertising.info()"
2 df.info()
3
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 200 entries, 0 to 199
Data columns (total 2 columns):
 #   Column  Non-Null Count  Dtype
---  ------  --------------  -----
 0   TV      200 non-null    float64
 1   Sales   200 non-null    float64
dtypes: float64(2)
memory usage: 3.3 KB
```

In the next step we will find the variable tyoes which we can also find them using `.info()` as we did in the last step.

```
1 #Show the types of the variables
2 var_types = df.dtypes
3 var_types
4
```

|       | 0       |
|-------|---------|
| **TV** | float64 |
| **Sales** | float64 |

**dtype:** object

```
1 #Show the total null values in the data frame
2 df.isnull().sum()
3
```

|       | 0 |
|-------|---|
| **TV** | 0 |
| **Sales** | 0 |

**dtype:** int64

## Summarize the data

```
1 #Use describe() for summarizing the data
2 #Round it to 2 decimals places
3 round(df.describe(),2)
4
```
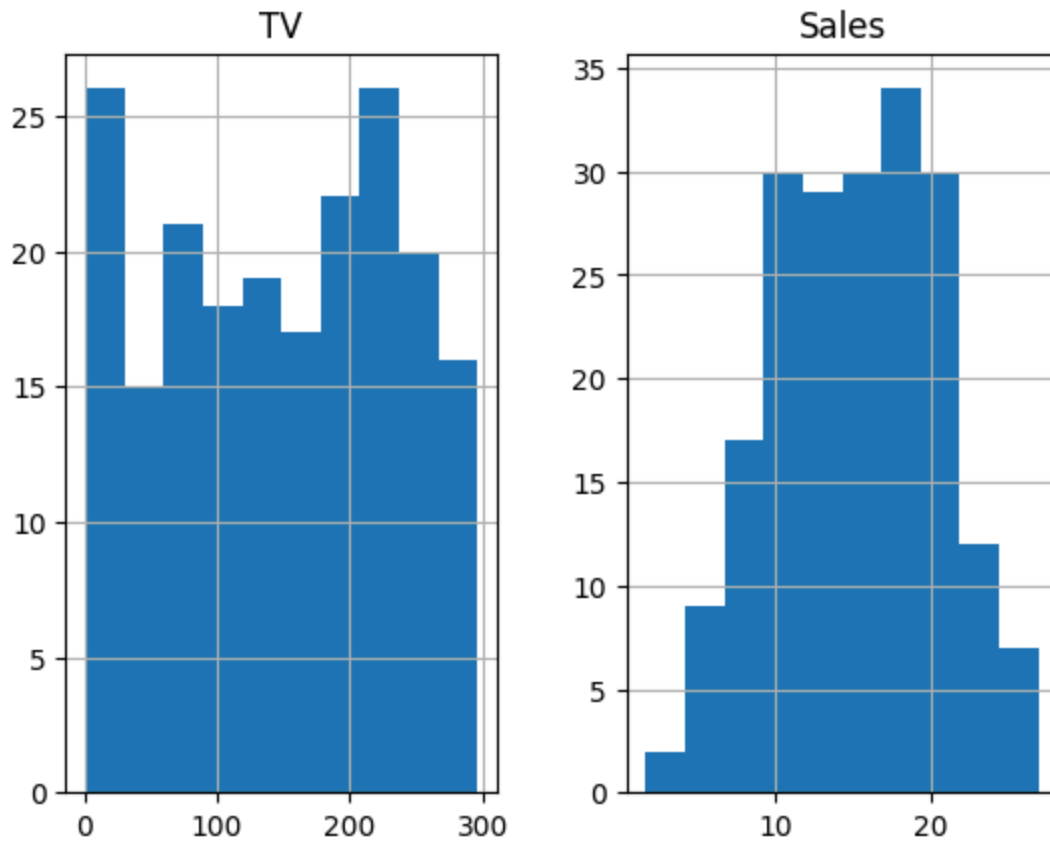
|       | TV | Sales |
|-------|--------|--------|
| count | 200.00 | 200.00 |
| mean | 147.04 | 15.13 |
| std | 85.85 | 5.28 |
| min | 0.70 | 1.60 |
| 25% | 74.38 | 11.00 |
| 50% | 149.75 | 16.00 |
| 75% | 218.82 | 19.05 |
| max | 296.40 | 27.00 |

## Analyze the data

### Visualize the data
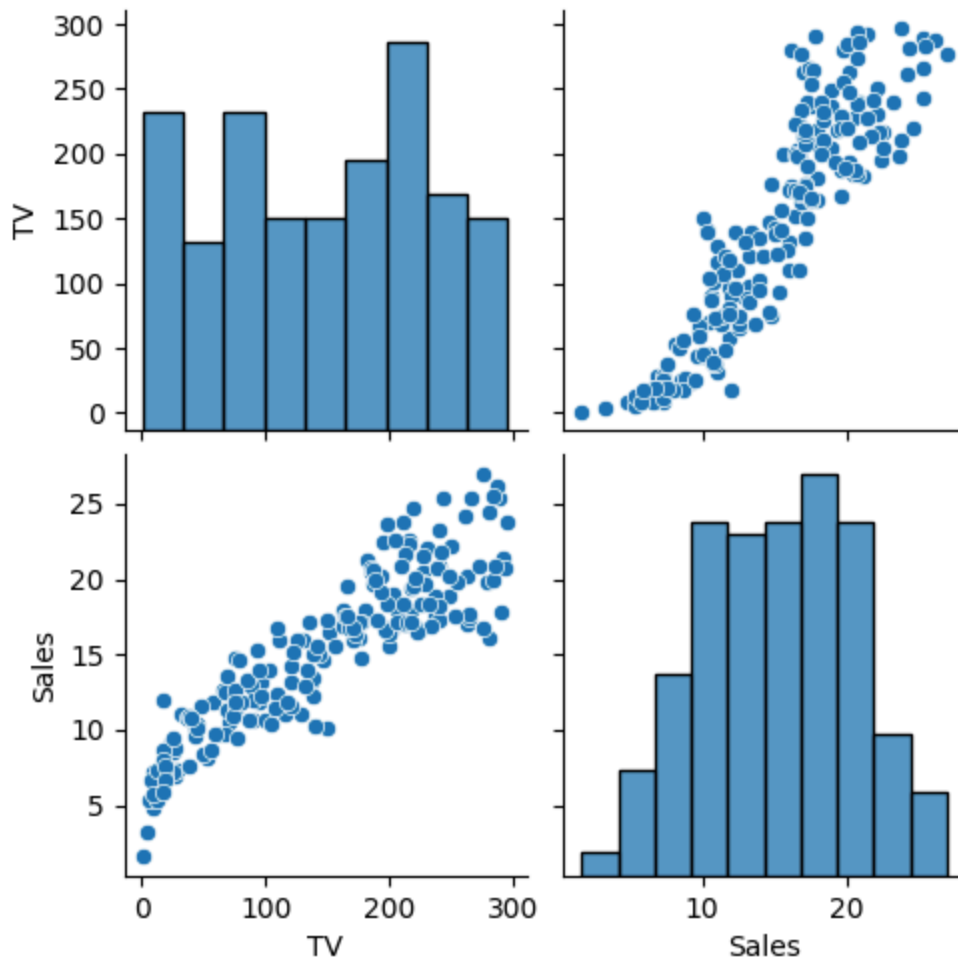
```
1 #Plot Histograms
2 hist_visual = df.hist(bins=10)
3 hist_visual
4
```

```
array([[<Axes: title={'center': 'TV'}>,
        <Axes: title={'center': 'Sales'}>]], dtype=object)
```
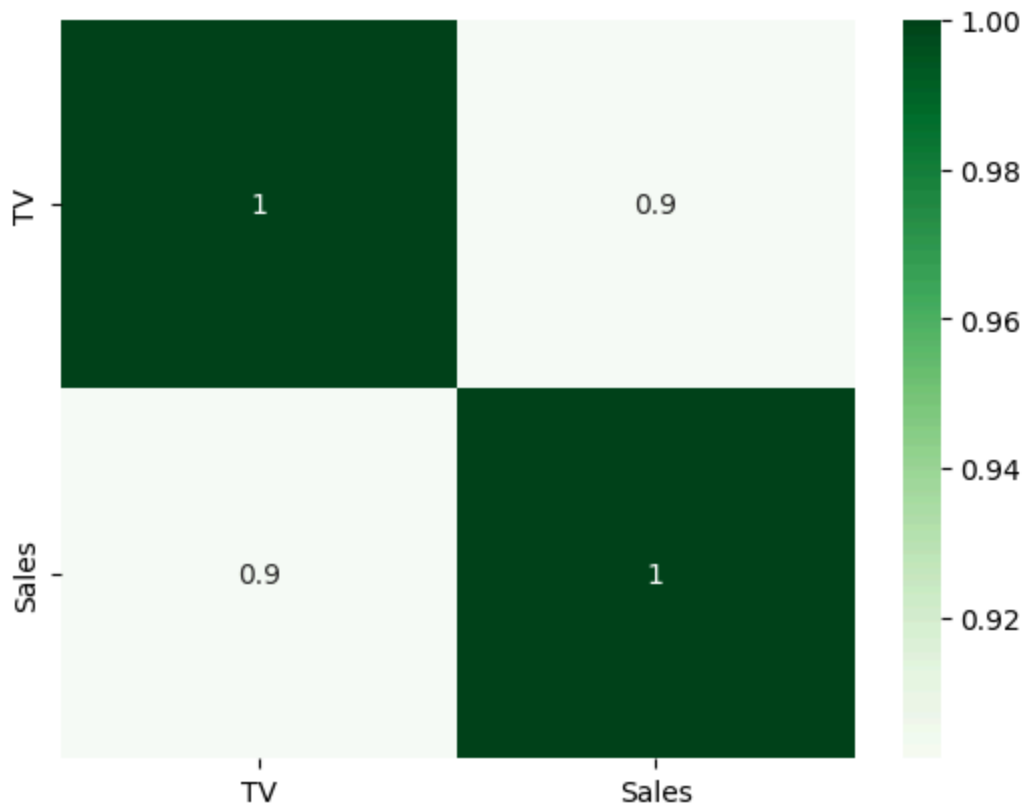


```
1 #Plot a pair plot to see the relationship between the variables
2 sn.pairplot(df)
3
```

```
<seaborn.axisgrid.PairGrid at 0x7bc7b9d62f90>
```



```
1 #Check correlation between TV and Sales
2 corr_matrix = df.corr()
3
4 #Heatmap
5 sn.heatmap(corr_matrix, annot=True, cmap= 'Greens')
6
```

<Axes: >



```
1 #Using the method "kendall" when using "".corr"
2 df.corr(method= 'kendall')
3
```
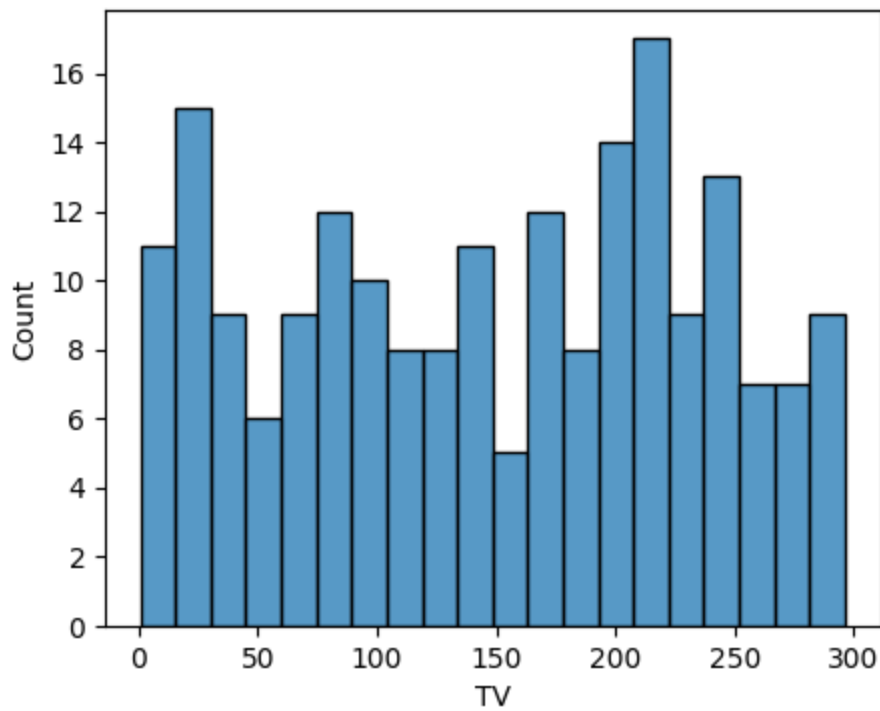
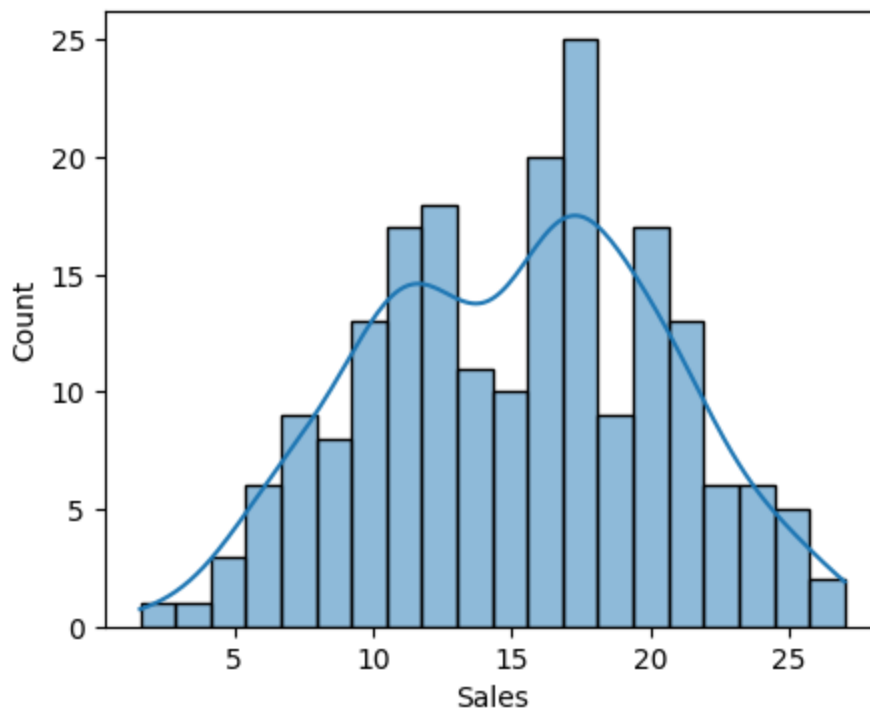|       | TV       | Sales    |
|-------|----------|----------|
| TV    | 1.000000 | 0.727994 |
| Sales | 0.727994 | 1.000000 |

```
1 #Do an outlier analysis for TV
2 plt.figure(figsize=(5, 4))
3 sn.histplot(df['TV'], bins=20)
4 plt.show()
5
```

```
1 #Do an outlier analysis for Sales(the target variable)
2 plt.figure(figsize=(5, 4))
3 sn.histplot(df['Sales'], bins=20, kde=True)
4 plt.show()
5
```
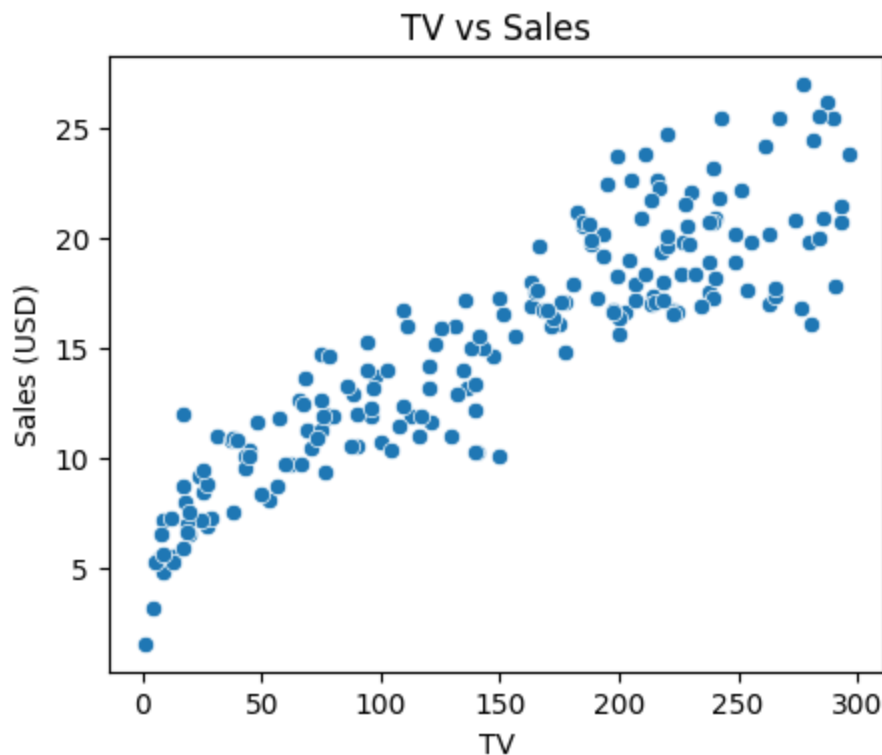


**Question: Do you see any considerable outliers?**

**Answer:** Yes there are a few

**Plotting Salary vs. YearsExperience**

```
1 #Plot a scatter plot to see how Sales is related to TV variable
2 plt.figure(figsize=(5,4))
3 sn.scatterplot(x='TV', y='Sales', data=df)
4 plt.title('TV vs Sales')
5 plt.xlabel('TV')
6 plt.ylabel('Sales (USD)')
7 plt.show()
```



**Question: What can you see here about the correlation?**

More TVs = More Sales

## STEP 2: Apply the Machine learning Model:

Here we will apply the ML model:

- Import the necessary libraries
- Build the model
- Display the results

## ⌄ Import the necessary libraries

```
1 #Import the library for splitting train/test data
2 from sklearn.model_selection import train_test_split
3
4 #Import for using the ML model
5 from sklearn.linear_model import LinearRegression
6 import statsmodels.api as sm #Linear regression with statsmodels
7
```

## ⌄ Buil the model

First we will prepare our dataframes for the x-array and y-array

```
1 #Create a X and Y data frames
2 X = df['TV']              #used to make predictions
3 Y = df['Sales']           #we'll predict this variable
```

Split the dataset into training and test dataset: Training 70% and test 30% sets

```
1 #Split the data into training and test sets
2 X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=
3
```

```
1 #Print the shape of the training and test data for X and Y
2 print("The shape of X_train is: ", X_train.shape)
3 print("The shape of X_test is: ", X_test.shape)
4 print("The shape of Y_train is: ", Y_train.shape)
5 print("The shape of Y_test is: ", Y_test.shape)
```
```
The shape of X_train is:  (140,)
The shape of X_test is:  (60,)
The shape of Y_train is:  (140,)
The shape of Y_test is:  (60,)
```

```
1 #The model requires a 2-D array
2 #Make the changes necessary for this
3 x_train = X_train.values.reshape(-1, 1)
4 x_test = X_test.values.reshape(-1, 1)
5
6 print("x_train shape is: ", x_train.shape)
7 print("x_test shape is: ", x_test.shape)
```
```
x_train shape is:  (140, 1)
x_test shape is:  (60, 1)
```

We build our linear regression model. We already import the necessary library.

```
1 #Set Linear Regression Model
2 model = LinearRegression()
3
```

Train the model with the training data. We use `.fit()` for that.

```
1 #Fit the training data
2 response = model.fit(x_train, Y_train)
3
```

## ⌄ Display Results

```
1 #Get the intercept and coefficient of the model
2 intercept = response.intercept_
3
4 coeff = response.coef_
5
```

```
1 #Print the intercept and coefficient of the model
2 print("The intercept is: ", intercept)
3 print("The coefficient is: ", coeff)
```

```
The intercept is:  6.928475121355413
The coefficient is:  [0.056007]
```

```
1 #Print the standard form using the intercept and coeeficient
2 print("Then, we have: y= %d + %d * x"  %(intercept, coeff))
3
```

```
Then, we have: y= 6 + 0 * x
/tmp/ipython-input-3284058802.py:2: DeprecationWarning: Conversion of an
  print("Then, we have: y= %d + %d * x"  %(intercept, coeff))
```