

Problem Statement

Title: Predicting Sales Based on Advertising Spend

Problem Description:

In today's competitive market, companies are constantly looking for ways to optimize their advertising budgets to maximize sales. Understanding the relationship between advertising expenditure across various media channels and sales performance is crucial for making informed marketing decisions. This dataset contains data on advertising expenditures across three channels: TV, Radio, and Newspaper, along with the resulting sales figures.

Objective:

The objective of this analysis is to build a predictive model that can accurately forecast sales based on the amount of money spent on TV, Radio, and Newspaper advertisements. This model will help in understanding the impact of each advertising channel on sales and in making strategic decisions regarding budget allocation to maximize sales.

```
In [4]: import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
```

```
In [5]: code=pd.read_csv('advertising.csv')
code
```

```
Out[5]:
```

	TV	Radio	Newspaper	Sales
0	230.1	37.8	69.2	22.1
1	44.5	39.3	45.1	10.4
2	17.2	45.9	69.3	12.0
3	151.5	41.3	58.5	16.5
4	190.8	10.8	56.4	17.9
...
195	38.2	3.7	13.8	7.6
196	34.2	4.9	8.1	14.0
197	177.0	9.3	6.4	14.8
198	283.6	42.0	66.2	25.5
199	232.1	8.6	8.7	18.4

200 rows x 4 columns

Dataset Description:

```
In [ ]: The dataset consists of the following columns:
1. TV: Amount of money spent on TV advertising (in thousands of dollars).
2. Radio: Amount of money spent on Radio advertising (in thousands of dollars).
3. Newspaper: Amount of money spent on Newspaper advertising (in thousands of dollars).
4. Sales: Sales of the product. (in thousands of units).
```

Data Understanding

```
In [21]: code.head()
```

```
Out[21]:
```

	TV	Radio	Newspaper	Sales
0	230.1	37.8	69.2	22.1
1	44.5	39.3	45.1	10.4
2	17.2	45.9	69.3	12.0
3	151.5	41.3	58.5	16.5
4	190.8	10.8	56.4	17.9

```
In [22]: code.describe()
```

```
Out[22]:
```

	TV	Radio	Newspaper	Sales
count	200.000000	200.000000	200.000000	200.000000
mean	147.842500	23.264000	30.554000	15.135000
std	85.854236	14.546909	21.779621	5.283592
min	0.700000	0.000000	0.300000	1.600000
25%	74.375000	9.970000	12.700000	11.000000
50%	149.750000	22.900000	25.750000	16.000000
75%	218.825000	36.525000	45.100000	19.000000
max	296.400000	49.600000	114.000000	27.000000

Initial check up

```
In [23]: code.info()
```

```
Out[23]:
```

	TV	Radio	Newspaper	Sales
count	200	200	200	200
mean	147.842500	23.264000	30.554000	15.135000
std	85.854236	14.546909	21.779621	5.283592
min	0.700000	0.000000	0.300000	1.600000
25%	74.375000	9.970000	12.700000	11.000000
50%	149.750000	22.900000	25.750000	16.000000
75%	218.825000	36.525000	45.100000	19.000000
max	296.400000	49.600000	114.000000	27.000000

```
In [ ]: #Checking Questions
#1. tv wise sales
#2. radio wise sales
#3. newspaper wise sales
#4. maximum sales on tv
#5. minimum sales on radio
#6. highest sales on newspaper
```

```
In [29]: #1. tv wise sales
code['TV'].value_counts(normalize=True)*100
```

```
Out[29]:
```

TV	Percentage
199.8	1.0
189.8	1.0
17.2	1.0
177.0	1.0
222.4	1.0
...	...
139.3	0.5
218.8	0.5
199.1	0.5
26.8	0.5
232.1	0.5

Name: proportion, Length: 198, dtype: float64

```
In [28]: #2. radio wise sales
code['Radio'].value_counts(normalize=True)*100
```

```
Out[28]:
```

Radio	Percentage
4.1	1.5
6.7	1.5
13.9	1.0
14.3	1.0
16.9	1.0
...	...
42.8	0.5
14.5	0.5
38.0	0.5
33.8	0.5
4.6	0.5

Name: proportion, Length: 167, dtype: float64

```
In [28]: #3. newspaper wise sales
code['Newspaper'].value_counts(normalize=True)*100
```

```
Out[29]:
```

Newspaper	Percentage
9.3	1.5
25.8	1.5
8.7	1.5
34.6	1.0
8.5	1.0
...	...
27.2	0.5
19.2	0.5
31.3	0.5
66.2	0.5

Name: proportion, Length: 172, dtype: float64

Data Visualization

```
In [7]: sns.histplot(code,x='TV')
```

```
Out[7]:
```

```
In [8]: sns.kdeplot(code,x='Radio')
```

```
Out[7]:
```

```
In [9]: sns.histplot(code,x='Newspaper')
```

```
Out[8]:
```

```
In [10]: #4. maximum sales on tv
sns.lmplot(code,x='TV',y='Sales')
```

```
Out[10]:
```

```
In [10]: #5. minimum sales on radio
sns.scatterplot(code,x='Radio',y='Sales')
```

```
Out[10]:
```

```
In [14]: #6. highest sales on newspaper
sns.scatterplot(code,x='Newspaper',y='Sales')
```

```
Out[14]:
```

```
In [76]: code.shape
```

```
Out[76]: (200, 4)
```

```
In [12]: x=code.iloc[:,0:1]
y=code.iloc[:,1:]
```

train_test_split

```
In [13]: from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.2,random_state=2)
```

```
Out[13]:
```

TV
137
163
111
123
109

```
In [14]: from sklearn.linear_model import LinearRegression
lr=LinearRegression()
```

```
In [17]: lr.fit(x_train,y_train)
```

```
Out[17]:
```

LinearRegression()

```
In [16]: x_train.head()
```

```
Out[16]:
```

TV
137
163
111
123
109

```
In [34]: lr.coef_
```

```
Out[34]: array([0.05835051])
```

```
In [39]: lr.intercept_
```

```
Out[39]: 6.57452993436827
```

```
In [36]: lr.predict([18.54])
```

```
Out[36]: array([7.0713608])
```

```
C:\Users\jyothanconda\lib\site-packages\sklearn\base.py:439: UserWarning: X does not have valid feature names, but LinearRegression was fitted with feature names
warnings.warn(msg)
```

```
In [38]: plt.scatter(code['TV'],code['Sales'])
plt.plot(x_train,lr.predict(x_train),color='red')
plt.xlabel('TV')
plt.ylabel('Sales')
```

```
Out[38]:
```

```
In [19]: plt.scatter(code['Radio'],code['Sales'])
plt.plot(x_train,lr.predict(x_train),color='green')
plt.xlabel('Radio')
plt.ylabel('Sales')
```

```
Out[19]:
```

```
In [20]: plt.scatter(code['Newspaper'],code['Sales'])
plt.plot(x_train,lr.predict(x_train),color='black')
plt.xlabel('Newspaper')
plt.ylabel('Sales')
```

```
Out[20]:
```

Model Evaluation

```
In [21]: from sklearn.metrics import mean_absolute_error,mean_squared_error
y_pred=lr.predict(x_test)
```

```
Out[21]:
```

TV
17.2
17.2
17.2
17.2
17.2

```
In [22]: mean_absolute_error(y_test,y_pred)
```

```
Out[22]: 2.824868390546763
```

```
In [23]: mean_squared_error(y_test,y_pred)
```

```
Out[23]: 6.748219198871345
```

```
In [24]: MSE=mean_squared_error(y_test,y_pred)
RMSE=RMSE**0.5
```

```
Out[24]:
```

actual_y_test	predicted(y_pred)
112	17.1
29	10.5
182	8.7
199	18.4
193	19.6

```
In [25]: plt.scatter(x_test,y_test,color='red',label='actual')
plt.scatter(x_test,y_pred,color='green',label='predicted')
plt.xlabel('TV')
plt.ylabel('Sales')
plt.legend()
plt.show()
```

```
Out[25]:
```

```
In [26]: plt.scatter(x_test,y_test,color='black',label='actual')
plt.scatter(x_test,y_pred,color='blue',label='predicted')
plt.xlabel('Radio')
plt.ylabel('Sales')
plt.legend()
plt.show()
```

```
Out[26]:
```

```
In [27]: plt.scatter(x_test,y_test,color='yellow',label='actual')
plt.scatter(x_test,y_pred,color='violet',label='predicted')
plt.xlabel('Newspaper')
plt.ylabel('Sales')
plt.legend()
plt.show()
```

```
Out[27]:
```

```
In [ ]: #Insights
#we have to import the advertising dataset based on the problem description
#the head() it takes the first five rows of the dataset and describe() it calculates the count,mean,std,min,max of each column of the dataset
#88 above tv are the highest count of the plot and the least count of the tv is 15
#maximum sales on tv is 25 above the sales on tv should be increased
#the scatterplot defines the relationship between two numerical variables and it display the plot
#we have to import linear regression and fit the x_train and y_train
#coefficient,intercept,predict are performed on the dataset by using linear regression
#model evaluation is implemented to predict the mean_absolute_error,mean_squared_error
#actual data and predicted data are performed in MSE and RMSE
```