

Problem Statement:

A real estate agent want to help to predict the house price for regions in USA.He gave us the dataset to work on to use Linear Regression modelCreate a Model that helps him to estimate of what the house would sell for

```
In [2]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [3]: df=pd.read_csv("fitness.csv")
df
```

```
Out[3]:
```

	Row Labels	Sum of Jan	Sum of Feb	Sum of Mar	Sum of Total Sales
0	A	5.62%	7.73%	6.16%	75
1	B	4.21%	17.27%	19.21%	160
2	C	9.83%	11.60%	5.17%	101
3	D	2.81%	21.91%	7.88%	127
4	E	25.28%	10.57%	11.82%	179
5	F	8.15%	16.24%	18.47%	167
6	G	18.54%	8.76%	17.49%	171
7	H	25.56%	5.93%	13.79%	170
8	Grand Total	100.00%	100.00%	100.00%	1150

```
In [4]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 9 entries, 0 to 8
Data columns (total 5 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Row Labels            9 non-null      object
1   Sum of Jan             9 non-null      object
2   Sum of Feb             9 non-null      object
3   Sum of Mar             9 non-null      object
4   Sum of Total Sales     9 non-null      int64
dtypes: int64(1), object(4)
memory usage: 488.0+ bytes
```

```
In [5]: df.head()
```

```
Out[5]:
```

	Row Labels	Sum of Jan	Sum of Feb	Sum of Mar	Sum of Total Sales
0	A	5.62%	7.73%	6.16%	75

	Row Labels	Sum of Jan	Sum of Feb	Sum of Mar	Sum of Total Sales
1	B	4.21%	17.27%	19.21%	160
2	C	9.83%	11.60%	5.17%	101
3	D	2.81%	21.91%	7.88%	127
4	E	25.28%	10.57%	11.82%	179

Data cleaning and Pre-Processing

In [6]:

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 9 entries, 0 to 8
Data columns (total 5 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Row Labels            9 non-null      object
1   Sum of Jan            9 non-null      object
2   Sum of Feb            9 non-null      object
3   Sum of Mar            9 non-null      object
4   Sum of Total Sales    9 non-null      int64
dtypes: int64(1), object(4)
memory usage: 488.0+ bytes
```

In [7]:

```
df.describe()
```

Out[7]:

	Sum of Total Sales
count	9.000000
mean	255.555556
std	337.332963
min	75.000000
25%	127.000000
50%	167.000000
75%	171.000000
max	1150.000000

In [8]:

```
a= df.dropna(axis='columns')
a
```

Out[8]:

	Row Labels	Sum of Jan	Sum of Feb	Sum of Mar	Sum of Total Sales
0	A	5.62%	7.73%	6.16%	75
1	B	4.21%	17.27%	19.21%	160
2	C	9.83%	11.60%	5.17%	101

	Row Labels	Sum of Jan	Sum of Feb	Sum of Mar	Sum of Total Sales
3	D	2.81%	21.91%	7.88%	127
4	E	25.28%	10.57%	11.82%	179
5	F	8.15%	16.24%	18.47%	167
6	G	18.54%	8.76%	17.49%	171
7	H	25.56%	5.93%	13.79%	170
8	Grand Total	100.00%	100.00%	100.00%	1150

In [9]:

```
a.columns
```

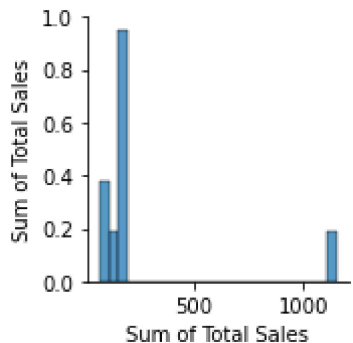
```
Out[9]: Index(['Row Labels', 'Sum of Jan', 'Sum of Feb', 'Sum of Mar',  
             'Sum of Total Sales'],  
            dtype='object')
```

EDA and VISUALIZATION

In [10]:

```
sns.pairplot(a)
```

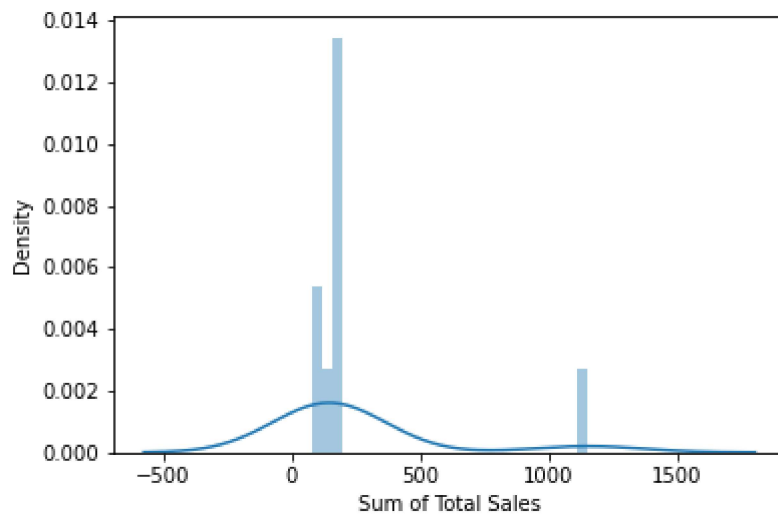
```
Out[10]: <seaborn.axisgrid.PairGrid at 0x23c490137c0>
```



In [12]:

```
sns.distplot(a['Sum of Total Sales'])
```

```
Out[12]: <AxesSubplot:xlabel='Sum of Total Sales', ylabel='Density'>
```



```
In [13]: df1=a[['Row Labels', 'Sum of Jan', 'Sum of Feb', 'Sum of Mar',
               'Sum of Total Sales']]
```

Plot Using Heat Map

```
In [14]: sns.heatmap(df1.corr())
```

```
Out[14]: <AxesSubplot:>
```



To Train The Model-Model Building

we are going to train Linera Regression Model;We need to split out data into two variables x and y where x is independent variable(input) and y is dependent on x(output) we could ignore address column as it required for our model

```
In [15]: x=df1[['Sum of Total Sales','Sum of Total Sales' ]]
         y=df1['Sum of Total Sales']
```

To Split my dataset into training and test data

```
In [16]: from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.3)
```

```
In [17]: from sklearn.linear_model import LinearRegression
lr= LinearRegression()
lr.fit(x_train,y_train)
```

Out[17]: LinearRegression()

```
In [18]: lr.intercept_
```

Out[18]: -1.1368683772161603e-13

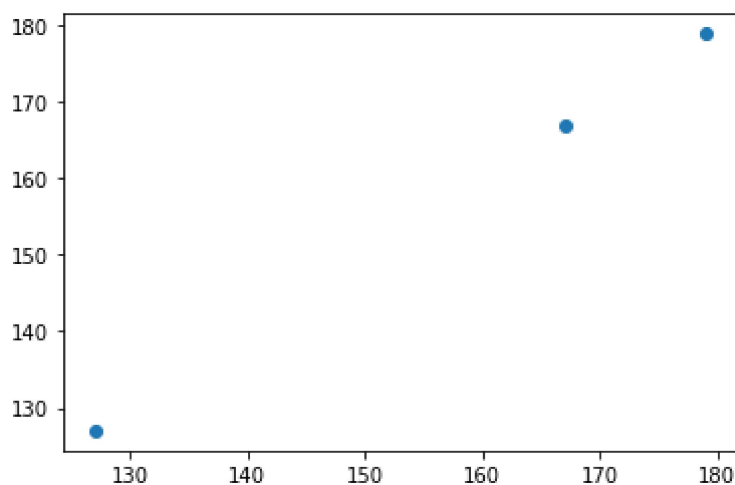
```
In [19]: coeff = pd.DataFrame(lr.coef_,x.columns,columns=['Co-efficient'])
coeff
```

Out[19]:

	Co-efficient
Sum of Total Sales	0.5
Sum of Total Sales	0.5

```
In [20]: prediction = lr.predict(x_test)
plt.scatter(y_test,prediction)
```

Out[20]: <matplotlib.collections.PathCollection at 0x23c4b32bfd0>



```
In [21]: lr.score(x_test,y_test)
```

Out[21]: 1.0