# Problem statement

# Data collection

# Importing libraries

In [1]:
```python
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

# Importing dataset

In [2]:
```python
data=pd.read_csv(r"C:\Users\user\Downloads\3_Fitness-1 - 3_Fitness-1.csv")
data
```

Out[2]:

|   | 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 |

# head

In [3]:
```python
# to display first 8 dataset values
da=data.head(8)
da
```

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 |

| | 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 |

# info

In [4]:
```python
# to identify missing values
data.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
```

# describe

In [5]:
```python
# to display summary of the dataset
data.describe()
```

Out[5]:

| | 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 |

# columns

```
In [6]:   # to display headings of the dataset
          data.columns
```

```
Out[6]: Index(['Row Labels', 'Sum of Jan', 'Sum of Feb', 'Sum of Mar',
               'Sum of Total Sales'],
              dtype='object')
```

```
In [7]:   a=data.dropna(axis=1)
          a
```

Out[7]:

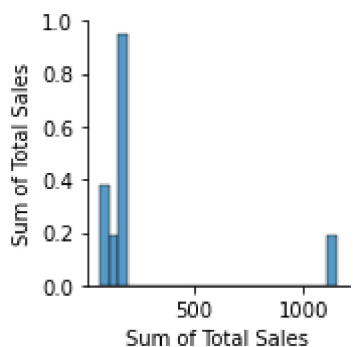|   | Row Labels | Sum of Jan | Sum of Feb | Sum of Mar | Sum of Total Sales |
|---|------------|-----------|-----------|-----------|--------------------|
| 0 | A | 5.62% | 7.73% | 6.16% | 75 |
| 1 | B | 4.21% | 17.27% | 19.21% | 160 |
| 2 | C | 9.83% | 11.60% | 5.17% | 101 |
| 3 | D | 2.81% | 21.91% | 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 [8]:   a.columns
```

```
Out[8]: Index(['Row Labels', 'Sum of Jan', 'Sum of Feb', 'Sum of Mar',
               'Sum of Total Sales'],
              dtype='object')
```

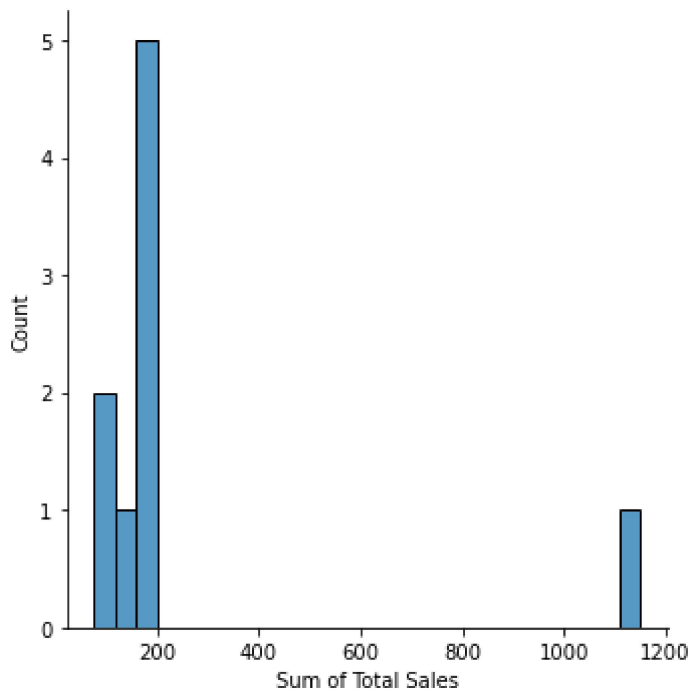# EDA and Visualization

```
In [9]:   sns.pairplot(a)
```

```
Out[9]: <seaborn.axisgrid.PairGrid at 0x1ce812bcf40>
```



# distribution plot

In [11]:
```python
sns.displot(a["Sum of Total Sales"])
```

Out[11]: <seaborn.axisgrid.FacetGrid at 0x1ce811d9a90>



# correlation

In [12]:
```python
dat=data[['Row Labels', 'Sum of Jan', 'Sum of Feb', 'Sum of Mar',
          'Sum of Total Sales']]
sns.heatmap(dat.corr())
```

Out[12]: <AxesSubplot:>



# To train the model-Model Building

```
In [17]:   x=a[['Sum of Total Sales']]
           y=a['Sum of Total Sales']
```

```
In [18]:   # to split my dataset into training and test data
           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 [19]:   from sklearn.linear_model import LinearRegression
           lr= LinearRegression()
           lr.fit(x_train,y_train)
```

```
Out[19]:   LinearRegression()
```
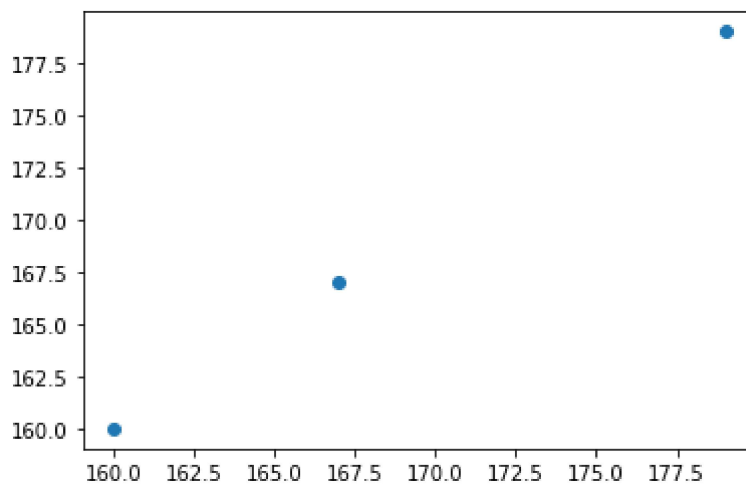
```
In [20]:   print(lr.intercept_)
```

```
5.684341886080802e-14
```

```
In [21]:   coeff=pd.DataFrame(lr.coef_,x.columns,columns=['Co-efficient'])
           coeff
```

Out[21]:

|  | Co-efficient |
|---|---|
| **Sum of Total Sales** | 1.0 |

```
In [22]:   prediction=lr.predict(x_test)
           plt.scatter(y_test,prediction)
```

Out[22]:   <matplotlib.collections.PathCollection at 0x1ce838434f0>



```
In [23]:   print(lr.score(x_test,y_test))
```

```
1.0
```