#Problem statement: The transactions made by a UK-based,we will design a clustering model and select the ideal groupof clients for the business to target

In [1]:

```python
#importing libraries
import pandas as pd
from matplotlib import pyplot as plt
%matplotlib inline
```

In [3]:

```python
df=pd.read_csv(r"C:\Users\jas_m\Downloads\212U5A3305\OnlineRetail.csv")
df
```

Out[3]:

| | InvoiceNo | StockCode | Description | Quantity | InvoiceDate | UnitPrice | CustomerID | |
|---|---|---|---|---|---|---|---|---|
| **0** | 536365 | 85123A | WHITE HANGING HEART T-LIGHT HOLDER | 6 | 01-12-2010 08:26 | 2.55 | 17850.0 | |
| **1** | 536365 | 71053 | WHITE METAL LANTERN | 6 | 01-12-2010 08:26 | 3.39 | 17850.0 | |
| **2** | 536365 | 84406B | CREAM CUPID HEARTS COAT HANGER | 8 | 01-12-2010 08:26 | 2.75 | 17850.0 | |
| **3** | 536365 | 84029G | KNITTED UNION FLAG HOT WATER BOTTLE | 6 | 01-12-2010 08:26 | 3.39 | 17850.0 | |
| **4** | 536365 | 84029E | RED WOOLLY HOTTIE WHITE HEART. | 6 | 01-12-2010 08:26 | 3.39 | 17850.0 | |
| **...** | ... | ... | ... | ... | ... | ... | ... | |
| **541904** | 581587 | 22613 | PACK OF 20 SPACEBOY NAPKINS | 12 | 09-12-2011 12:50 | 0.85 | 12680.0 | |
| **541905** | 581587 | 22899 | CHILDREN'S APRON DOLLY GIRL | 6 | 09-12-2011 12:50 | 2.10 | 12680.0 | |
| **541906** | 581587 | 23254 | CHILDRENS CUTLERY DOLLY GIRL | 4 | 09-12-2011 12:50 | 4.15 | 12680.0 | |
| **541907** | 581587 | 23255 | CHILDRENS CUTLERY CIRCUS PARADE | 4 | 09-12-2011 12:50 | 4.15 | 12680.0 | |
| **541908** | 581587 | 22138 | BAKING SET 9 PIECE RETROSPOT | 3 | 09-12-2011 12:50 | 4.95 | 12680.0 | |

541909 rows × 8 columns

In [4]:

```
df.head()
```

Out[4]:

|   | InvoiceNo | StockCode | Description | Quantity | InvoiceDate | UnitPrice | CustomerID | Country |
|---|-----------|-----------|-------------|----------|-------------|-----------|------------|---------|
| 0 | 536365 | 85123A | WHITE HANGING HEART T-LIGHT HOLDER | 6 | 01-12-2010 08:26 | 2.55 | 17850.0 | United Kingdom |
| 1 | 536365 | 71053 | WHITE METAL LANTERN | 6 | 01-12-2010 08:26 | 3.39 | 17850.0 | United Kingdom |
| 2 | 536365 | 84406B | CREAM CUPID HEARTS COAT HANGER | 8 | 01-12-2010 08:26 | 2.75 | 17850.0 | United Kingdom |
| 3 | 536365 | 84029G | KNITTED UNION FLAG HOT WATER BOTTLE | 6 | 01-12-2010 08:26 | 3.39 | 17850.0 | United Kingdom |
| 4 | 536365 | 84029E | RED WOOLLY HOTTIE WHITE HEART. | 6 | 01-12-2010 08:26 | 3.39 | 17850.0 | United Kingdom |

In [5]:

```
df.tail()
```

Out[5]:

|   | InvoiceNo | StockCode | Description | Quantity | InvoiceDate | UnitPrice | CustomerID | C |
|---|-----------|-----------|-------------|----------|-------------|-----------|------------|---|
| 541904 | 581587 | 22613 | PACK OF 20 SPACEBOY NAPKINS | 12 | 09-12-2011 12:50 | 0.85 | 12680.0 | |
| 541905 | 581587 | 22899 | CHILDREN'S APRON DOLLY GIRL | 6 | 09-12-2011 12:50 | 2.10 | 12680.0 | |
| 541906 | 581587 | 23254 | CHILDRENS CUTLERY DOLLY GIRL | 4 | 09-12-2011 12:50 | 4.15 | 12680.0 | |
| 541907 | 581587 | 23255 | CHILDRENS CUTLERY CIRCUS PARADE | 4 | 09-12-2011 12:50 | 4.15 | 12680.0 | |
| 541908 | 581587 | 22138 | BAKING SET 9 PIECE RETROSPOT | 3 | 09-12-2011 12:50 | 4.95 | 12680.0 | |

In [6]:

```python
df['InvoiceNo'].value_counts()
```

Out[6]:

```
573585      1114
581219       749
581492       731
580729       721
558475       705
            ...
554023         1
554022         1
554021         1
554020         1
C558901        1
Name: InvoiceNo, Length: 25900, dtype: int64
```

In [7]:

```python
df['CustomerID'].value_counts()
```

Out[7]:

```
17841.0    7983
14911.0    5903
14096.0    5128
12748.0    4642
14606.0    2782
           ...
15070.0       1
15753.0       1
17065.0       1
16881.0       1
16995.0       1
Name: CustomerID, Length: 4372, dtype: int64
```

In [8]:

```python
df['Quantity'].value_counts()
```

Out[8]:
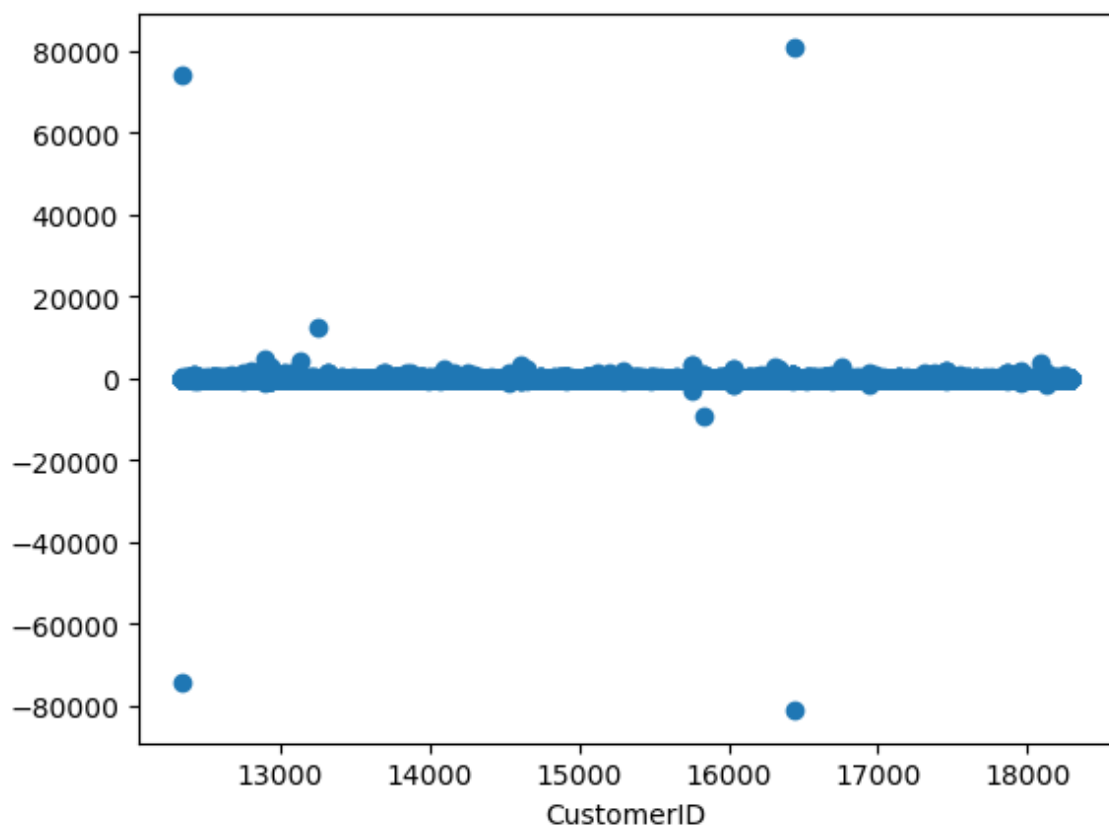
```
 1      148227
 2       81829
 12      61063
 6       40868
 4       38484
          ...
-472         1
-161         1
-1206        1
-272         1
-80995       1
Name: Quantity, Length: 722, dtype: int64
```

In [9]:

```python
plt.scatter(df["CustomerID"],df["Quantity"])
plt.xlabel("CustomerID")
```

Out[9]:

Text(0.5, 0, 'CustomerID')



In [10]:

```python
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 541909 entries, 0 to 541908
Data columns (total 8 columns):
 #   Column       Non-Null Count   Dtype
---  ------       --------------   -----
 0   InvoiceNo    541909 non-null  object
 1   StockCode    541909 non-null  object
 2   Description  540455 non-null  object
 3   Quantity     541909 non-null  int64
 4   InvoiceDate  541909 non-null  object
 5   UnitPrice    541909 non-null  float64
 6   CustomerID   406829 non-null  float64
 7   Country      541909 non-null  object
dtypes: float64(2), int64(1), object(5)
memory usage: 33.1+ MB
```

In [11]:

```python
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 541909 entries, 0 to 541908
Data columns (total 8 columns):
 #   Column       Non-Null Count   Dtype
---  ------       --------------   -----
 0   InvoiceNo    541909 non-null  object
 1   StockCode    541909 non-null  object
 2   Description  540455 non-null  object
 3   Quantity     541909 non-null  int64
 4   InvoiceDate  541909 non-null  object
 5   UnitPrice    541909 non-null  float64
 6   CustomerID   406829 non-null  float64
 7   Country      541909 non-null  object
dtypes: float64(2), int64(1), object(5)
memory usage: 33.1+ MB
```

In [12]:

```python
df.fillna(method='ffill',inplace=True)
df.isnull().sum()
```

Out[12]:

```
InvoiceNo      0
StockCode      0
Description    0
Quantity       0
InvoiceDate    0
UnitPrice      0
CustomerID     0
Country        0
dtype: int64
```

In [13]:

```python
from sklearn.cluster import KMeans
km=KMeans()
km
```

Out[13]:

```
▾ KMeans
KMeans()
```

In [14]:

```python
y_predicted=km.fit_predict(df[["CustomerID","Quantity"]])
y_predicted
```

C:\Users\jas_m\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:870:
FutureWarning: The default value of `n_init` will change from 10 to 'auto'
in 1.4. Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(

Out[14]:

```
array([0, 0, 0, ..., 2, 2, 2])
```

In [15]:

```python
df["cluster"]=y_predicted
df.head()
```
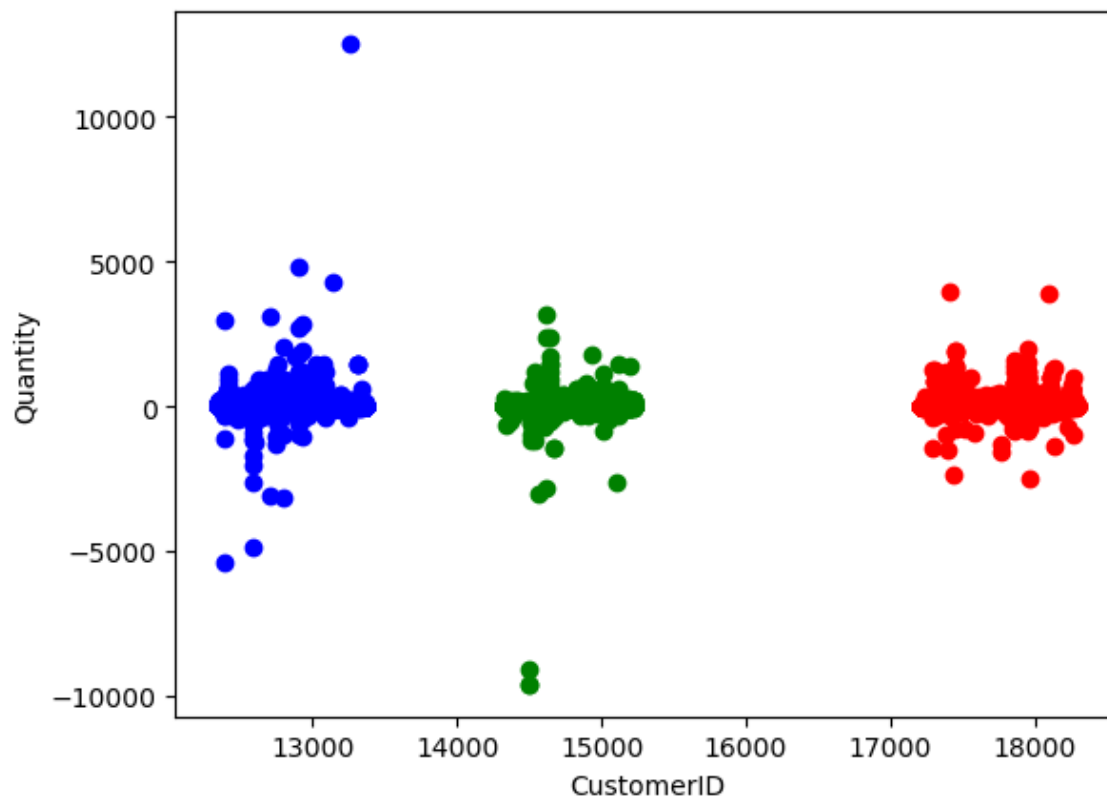
Out[15]:

| | InvoiceNo | StockCode | Description | Quantity | InvoiceDate | UnitPrice | CustomerID | Country |
|---|---|---|---|---|---|---|---|---|
| 0 | 536365 | 85123A | WHITE HANGING HEART T-LIGHT HOLDER | 6 | 01-12-2010 08:26 | 2.55 | 17850.0 | United Kingdom |
| 1 | 536365 | 71053 | WHITE METAL LANTERN | 6 | 01-12-2010 08:26 | 3.39 | 17850.0 | United Kingdom |
| 2 | 536365 | 84406B | CREAM CUPID HEARTS COAT HANGER | 8 | 01-12-2010 08:26 | 2.75 | 17850.0 | United Kingdom |
| 3 | 536365 | 84029G | KNITTED UNION FLAG HOT WATER BOTTLE | 6 | 01-12-2010 08:26 | 3.39 | 17850.0 | United Kingdom |
| 4 | 536365 | 84029E | RED WOOLLY HOTTIE WHITE HEART. | 6 | 01-12-2010 08:26 | 3.39 | 17850.0 | United Kingdom |

In [16]:

```python
df1=df[df.cluster==0]
df2=df[df.cluster==1]
df3=df[df.cluster==2]
plt.scatter(df1["CustomerID"],df1["Quantity"],color="red")
plt.scatter(df2["CustomerID"],df2["Quantity"],color="green")
plt.scatter(df3["CustomerID"],df3["Quantity"],color="blue")
plt.xlabel("CustomerID")
plt.ylabel("Quantity")
```

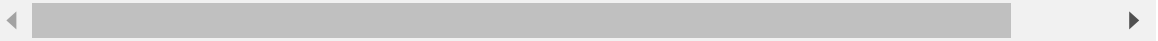Out[16]:

Text(0, 0.5, 'Quantity')

In [17]:

```python
from sklearn.preprocessing import MinMaxScaler
scaler=MinMaxScaler()
scaler.fit(df[["Quantity"]])
df["Quantity"]=scaler.transform(df[["Quantity"]])
df.head()
```

Out[17]:

| | InvoiceNo | StockCode | Description | Quantity | InvoiceDate | UnitPrice | CustomerID | Country |
|---|---|---|---|---|---|---|---|---|
| 0 | 536365 | 85123A | WHITE HANGING HEART T-LIGHT HOLDER | 0.500037 | 01-12-2010 08:26 | 2.55 | 17850.0 | United Kingdom |
| 1 | 536365 | 71053 | WHITE METAL LANTERN | 0.500037 | 01-12-2010 08:26 | 3.39 | 17850.0 | United Kingdom |
| 2 | 536365 | 84406B | CREAM CUPID HEARTS COAT HANGER | 0.500049 | 01-12-2010 08:26 | 2.75 | 17850.0 | United Kingdom |
| 3 | 536365 | 84029G | KNITTED UNION FLAG HOT WATER BOTTLE | 0.500037 | 01-12-2010 08:26 | 3.39 | 17850.0 | United Kingdom |
| 4 | 536365 | 84029E | RED WOOLLY HOTTIE WHITE HEART. | 0.500037 | 01-12-2010 08:26 | 3.39 | 17850.0 | United Kingdom |

In [18]:

```python
scaler.fit(df[["CustomerID"]])
df["CustomerID"]=scaler.transform(df[["CustomerID"]])
df.head()
```

Out[18]:

| | InvoiceNo | StockCode | Description | Quantity | InvoiceDate | UnitPrice | CustomerID | Country |
|---|---|---|---|---|---|---|---|---|
| **0** | 536365 | 85123A | WHITE HANGING HEART T-LIGHT HOLDER | 0.500037 | 01-12-2010 08:26 | 2.55 | 0.926443 | United Kingdom |
| **1** | 536365 | 71053 | WHITE METAL LANTERN | 0.500037 | 01-12-2010 08:26 | 3.39 | 0.926443 | United Kingdom |
| **2** | 536365 | 84406B | CREAM CUPID HEARTS COAT HANGER | 0.500049 | 01-12-2010 08:26 | 2.75 | 0.926443 | United Kingdom |
| **3** | 536365 | 84029G | KNITTED UNION FLAG HOT WATER BOTTLE | 0.500037 | 01-12-2010 08:26 | 3.39 | 0.926443 | United Kingdom |
| **4** | 536365 | 84029E | RED WOOLLY HOTTIE WHITE HEART. | 0.500037 | 01-12-2010 08:26 | 3.39 | 0.926443 | United Kingdom |

In [19]:

```python
km=KMeans()
```

In [20]:

```python
y_predicted=km.fit_predict(df[["CustomerID","Quantity"]])
y_predicted
```

```
C:\Users\jas_m\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:870:
FutureWarning: The default value of `n_init` will change from 10 to 'auto'
in 1.4. Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
```

Out[20]:

```
array([6, 6, 6, ..., 5, 5, 5])
```

In [21]:

```python
df["New Cluster"]=y_predicted
df.head()
```

Out[21]:

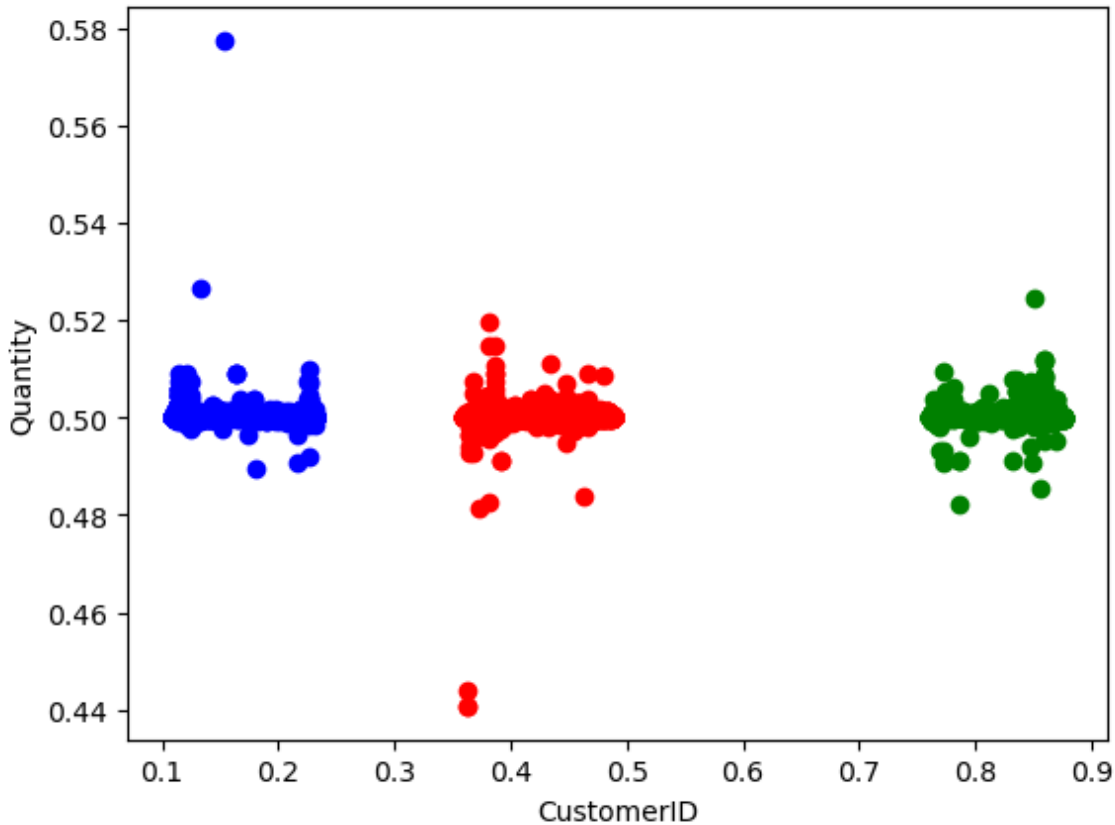| | InvoiceNo | StockCode | Description | Quantity | InvoiceDate | UnitPrice | CustomerID | Country |
|---|---|---|---|---|---|---|---|---|
| **0** | 536365 | 85123A | WHITE HANGING HEART T-LIGHT HOLDER | 0.500037 | 01-12-2010 08:26 | 2.55 | 0.926443 | United Kingdom |
| **1** | 536365 | 71053 | WHITE METAL LANTERN | 0.500037 | 01-12-2010 08:26 | 3.39 | 0.926443 | United Kingdom |
| **2** | 536365 | 84406B | CREAM CUPID HEARTS COAT HANGER | 0.500049 | 01-12-2010 08:26 | 2.75 | 0.926443 | United Kingdom |
| **3** | 536365 | 84029G | KNITTED UNION FLAG HOT WATER BOTTLE | 0.500037 | 01-12-2010 08:26 | 3.39 | 0.926443 | United Kingdom |
| **4** | 536365 | 84029E | RED WOOLLY HOTTIE WHITE HEART. | 0.500037 | 01-12-2010 08:26 | 3.39 | 0.926443 | United Kingdom |

In [22]:

```python
df1=df[df["New Cluster"]==0]
df2=df[df["New Cluster"]==1]
df3=df[df["New Cluster"]==2]
plt.scatter(df1["CustomerID"],df1["Quantity"],color="red")
plt.scatter(df2["CustomerID"],df2["Quantity"],color="green")
plt.scatter(df3["CustomerID"],df3["Quantity"],color="blue")
plt.xlabel("CustomerID")
plt.ylabel("Quantity")
```

Out[22]:

Text(0, 0.5, 'Quantity')



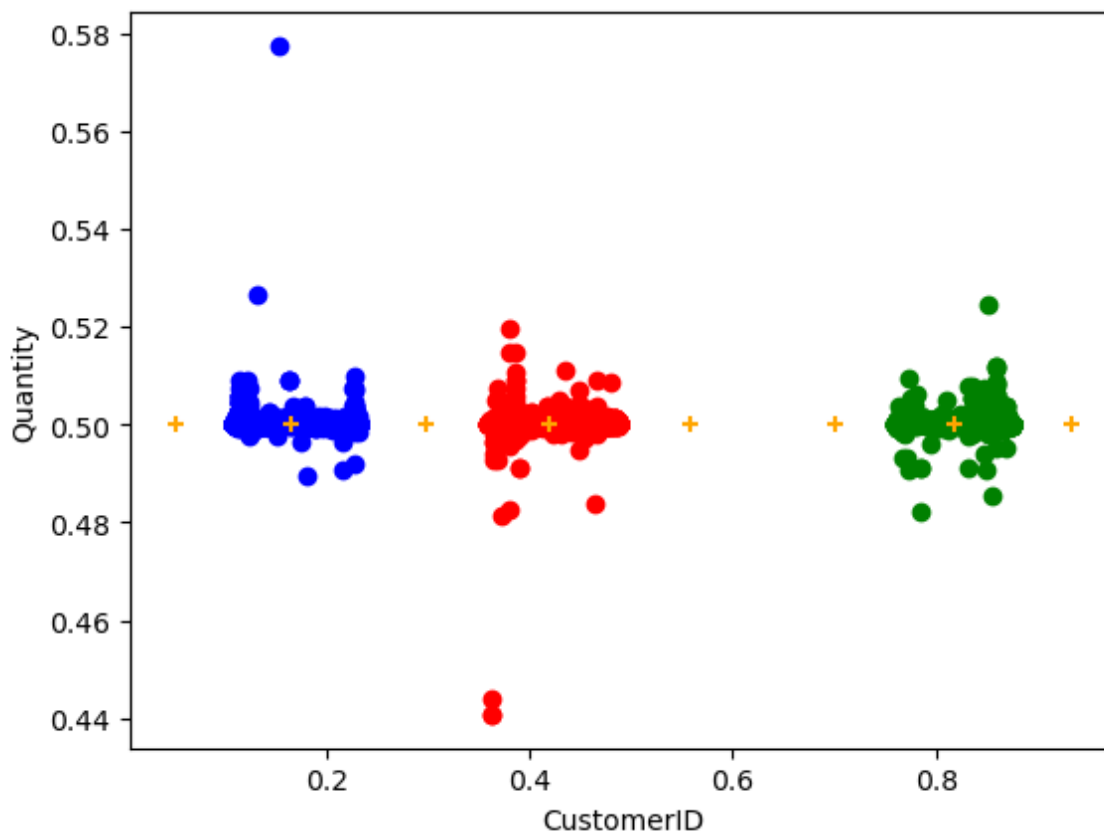In [23]:

```python
km.cluster_centers_
```

Out[23]:

```
array([[0.42034012, 0.50006104],
       [0.81846395, 0.50006031],
       [0.165866  , 0.50006058],
       [0.29887766, 0.50006072],
       [0.70173498, 0.50005811],
       [0.05156814, 0.50006705],
       [0.93301334, 0.50005098],
       [0.55740616, 0.50005319]])
```

In [24]:

```python
df1=df[df["New Cluster"]==0]
df2=df[df["New Cluster"]==1]
df3=df[df["New Cluster"]==2]
plt.scatter(df1["CustomerID"],df1["Quantity"],color="red")
plt.scatter(df2["CustomerID"],df2["Quantity"],color="green")
plt.scatter(df3["CustomerID"],df3["Quantity"],color="blue")
plt.scatter(km.cluster_centers_[:,0],km.cluster_centers_[:,1],color="orange",marker="+")
plt.xlabel("CustomerID")
plt.ylabel("Quantity")
```

Out[24]:

```
Text(0, 0.5, 'Quantity')
```



In [25]:

```python
k_rng=range(1,10)
sse=[]
```

In [27]:

```python
for k in k_rng:
    km=KMeans(n_clusters=k)
    km.fit(df[["CustomerID","Quantity"]])
    sse.append(km.inertia_)#km.inertia_ will give you the value of sum of square error
print(sse)
plt.plot(k_rng,sse)
plt.xlabel("K")
plt.ylabel("Sum of Squared Error")
```

```
C:\Users\jas_m\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:870:
FutureWarning: The default value of `n_init` will change from 10 to 'auto'
in 1.4. Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
C:\Users\jas_m\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:870:
FutureWarning: The default value of `n_init` will change from 10 to 'auto'
in 1.4. Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
C:\Users\jas_m\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:870:
FutureWarning: The default value of `n_init` will change from 10 to 'auto'
in 1.4. Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
C:\Users\jas_m\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:870:
FutureWarning: The default value of `n_init` will change from 10 to 'auto'
in 1.4. Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
C:\Users\jas_m\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:870:
FutureWarning: The default value of `n_init` will change from 10 to 'auto'
in 1.4. Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
C:\Users\jas_m\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:870:
FutureWarning: The default value of `n_init` will change from 10 to 'auto'
in 1.4. Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
C:\Users\jas_m\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:870:
FutureWarning: The default value of `n_init` will change from 10 to 'auto'
in 1.4. Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
C:\Users\jas_m\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:870:
FutureWarning: The default value of `n_init` will change from 10 to 'auto'
in 1.4. Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
C:\Users\jas_m\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:870:
FutureWarning: The default value of `n_init` will change from 10 to 'auto'
in 1.4. Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(

[46374.84553398474, 11336.0653054853, 4915.900467978989, 2723.51910518953,
1695.07109745742, 1178.4495409566046, 903.555546770782, 688.4704309720291,
528.5500057548226]
```
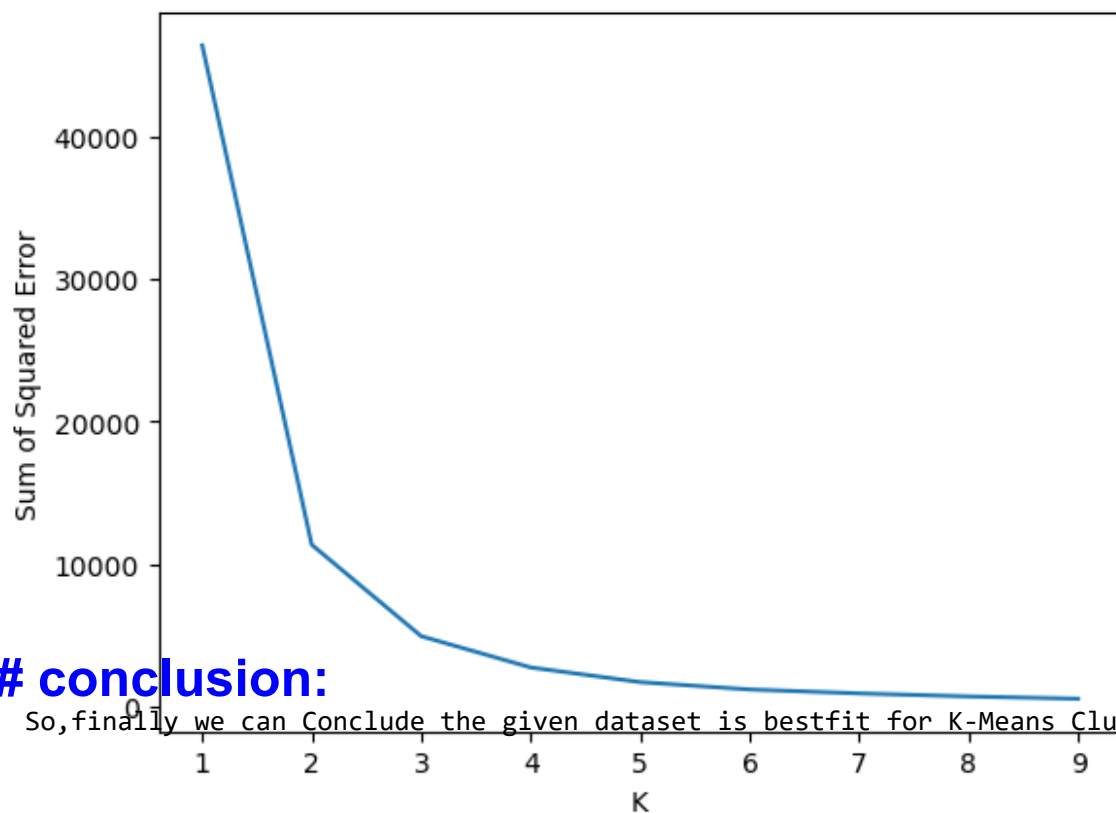
Out[27]:

```
Text(0, 0.5, 'Sum of Squared Error')
```

# conclusion:

So,finally we can Conclude the given dataset is bestfit for K-Means Clustering