**PROBLEM STATEMENT :Online Retail) The transactionsmade by a UK-based, registered, non-store onlineretailer between December 1, 2010, and December 9,2011, are all included in the transnational data setknown as online retail. The company primarily offersone-of-a-kind gifts for every occasion. The companyhas a large number of wholesalers as clients.CompanyObjectiveUsing the global online retail dataset, we willdesign a clustering model and select the ideal groupof clients for the business to target** ¶

In [1]:
```python
1  import pandas as pd
2  from matplotlib import pyplot as plt
3  %matplotlib inline
```

```
In [8]:  1  df=pd.read_csv(r"C:\Users\Dell\Downloads\OnlineRetail.csv")
         2  df
```

Out[8]:

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

541909 rows × 8 columns

In [3]:
```
1 df.head()
```

Out[3]:

| | 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 [4]:
```
1 df.tail()
```

Out[4]:

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

```
In [5]:    1  df['InvoiceNo'].value_counts()
```

```
Out[5]:  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
```
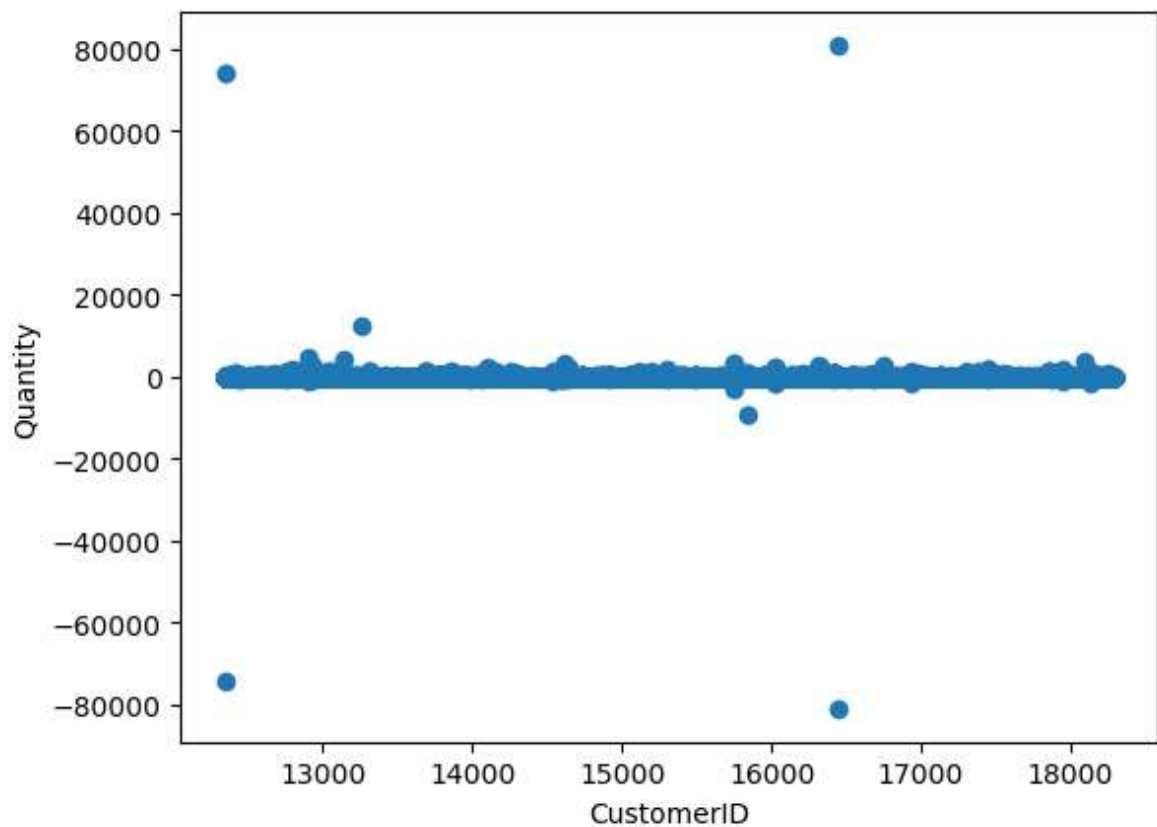
```
In [6]:    1  df['CustomerID'].value_counts()
```

```
Out[6]:  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 [7]:    1  df['Quantity'].value_counts()
```

```
Out[7]:   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]:   1  plt.scatter(df["CustomerID"],df["Quantity"])
          2  plt.xlabel("CustomerID")
          3  plt.ylabel("Quantity")
```

Out[9]:  Text(0, 0.5, 'Quantity')



```
In [10]:   1  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]: 1  df.isnull().sum()
```

```
Out[11]: InvoiceNo            0
         StockCode            0
         Description       1454
         Quantity             0
         InvoiceDate          0
         UnitPrice            0
         CustomerID      135080
         Country              0
         dtype: int64
```

```
In [12]: 1  df.fillna(method='ffill',inplace=True)
```

```
In [13]: 1  df.isnull().sum()
```

```
Out[13]: InvoiceNo       0
         StockCode       0
         Description     0
         Quantity        0
         InvoiceDate     0
         UnitPrice       0
         CustomerID      0
         Country         0
         dtype: int64
```

```
In [14]: 1  from sklearn.cluster import KMeans
         2  km=KMeans()
         3  km
```

```
Out[14]: ▼ KMeans

         KMeans()
```

```
In [15]: 1  y_predicted=km.fit_predict(df[["CustomerID","Quantity"]])
         2  y_predicted
```

C:\ProgramData\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:870: F
utureWarning: 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[15]: array([1, 1, 1, ..., 0, 0, 0])
```

```
In [16]:    1  df["cluster"]=y_predicted
            2  df.head()
```
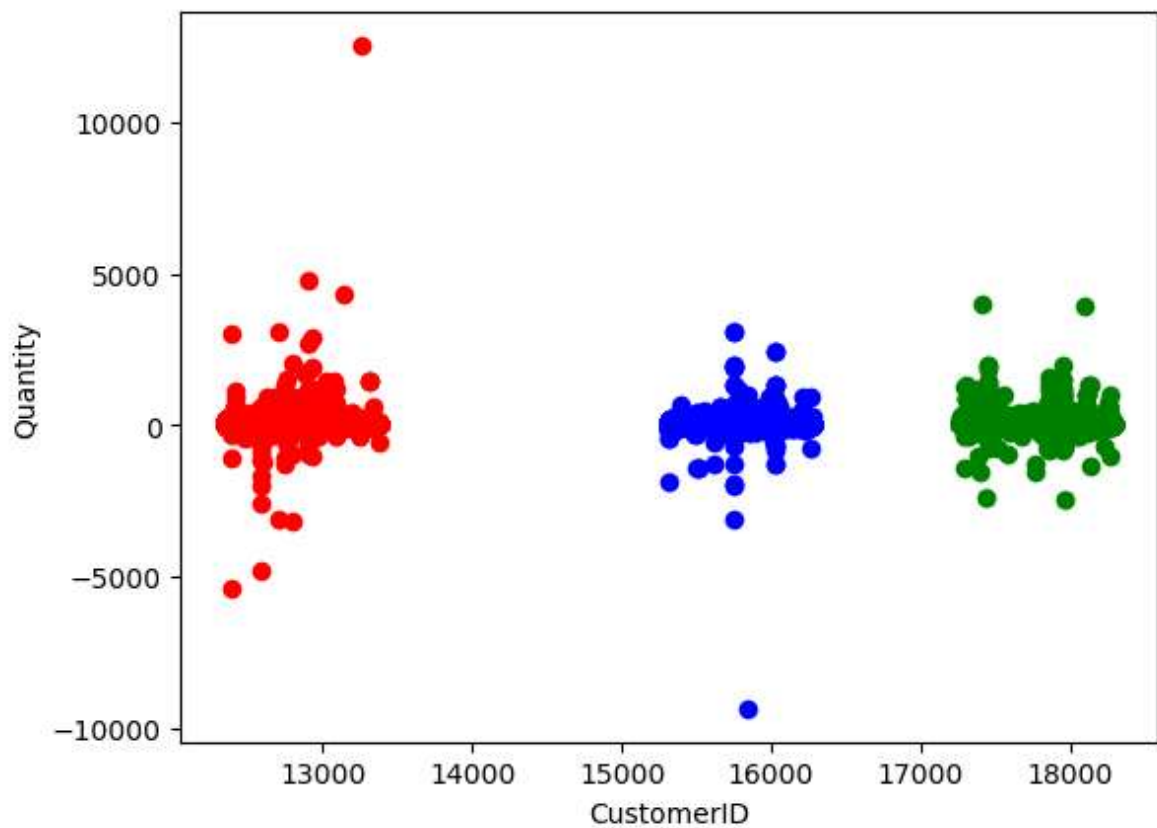
Out[16]:

| | InvoiceNo | StockCode | Description | Quantity | InvoiceDate | UnitPrice | CustomerID | Country | c |
|---|---|---|---|---|---|---|---|---|---|
| **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 | |

```
1  df1=df[df.cluster==0]
2  df2=df[df.cluster==1]
3  df3=df[df.cluster==2]
4  plt.scatter(df1["CustomerID"],df1["Quantity"],color="red")
5  plt.scatter(df2["CustomerID"],df2["Quantity"],color="green")
6  plt.scatter(df3["CustomerID"],df3["Quantity"],color="blue")
7  plt.xlabel("CustomerID")
8  plt.ylabel("Quantity")
```

Out[17]: Text(0, 0.5, 'Quantity')

In [18]:
```python
from sklearn.preprocessing import MinMaxScaler
scaler=MinMaxScaler()
scaler.fit(df[["Quantity"]])
df["Quantity"]=scaler.transform(df[["Quantity"]])
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 | 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 [19]:  1  scaler.fit(df[["CustomerID"]])
          2  df["CustomerID"]=scaler.transform(df[["CustomerID"]])
          3  df.head()
```

Out[19]:

| | InvoiceNo | StockCode | Description | Quantity | InvoiceDate | UnitPrice | CustomerID | Country | c |
|---|---|---|---|---|---|---|---|---|---|
| 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 | |

# K-MEANS CLUSTERING

```
In [20]:  1  km=KMeans()
```

```
In [21]:  1  y_predicted=km.fit_predict(df[["CustomerID","Quantity"]])
          2  y_predicted
```

```
C:\ProgramData\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:870: F
utureWarning: 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[21]:  array([6, 6, 6, ..., 7, 7, 7])

In [22]:
```
1 df["New Cluster"]=y_predicted
2 df.head()
```

Out[22]:

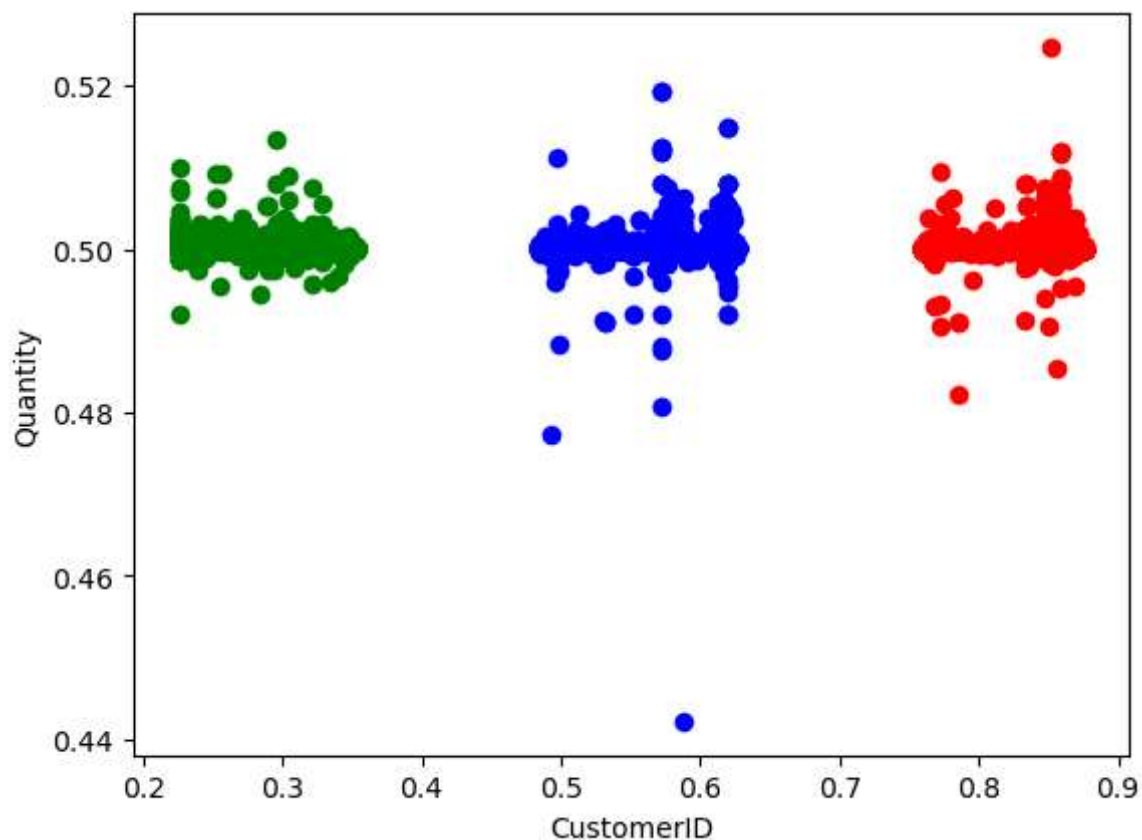| | InvoiceNo | StockCode | Description | Quantity | InvoiceDate | UnitPrice | CustomerID | Country | c |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 536365 | 85123A | WHITE HANGING HEART T-LIGHT HOLDER | 0.500037 | 01-12-2010 08:26 | 2.55 | 0.926443 | United Kingdom | |
| 1 | 536365 | 71053 | WHITE 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 | |

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

Out[23]: Text(0, 0.5, 'Quantity')



In [24]:
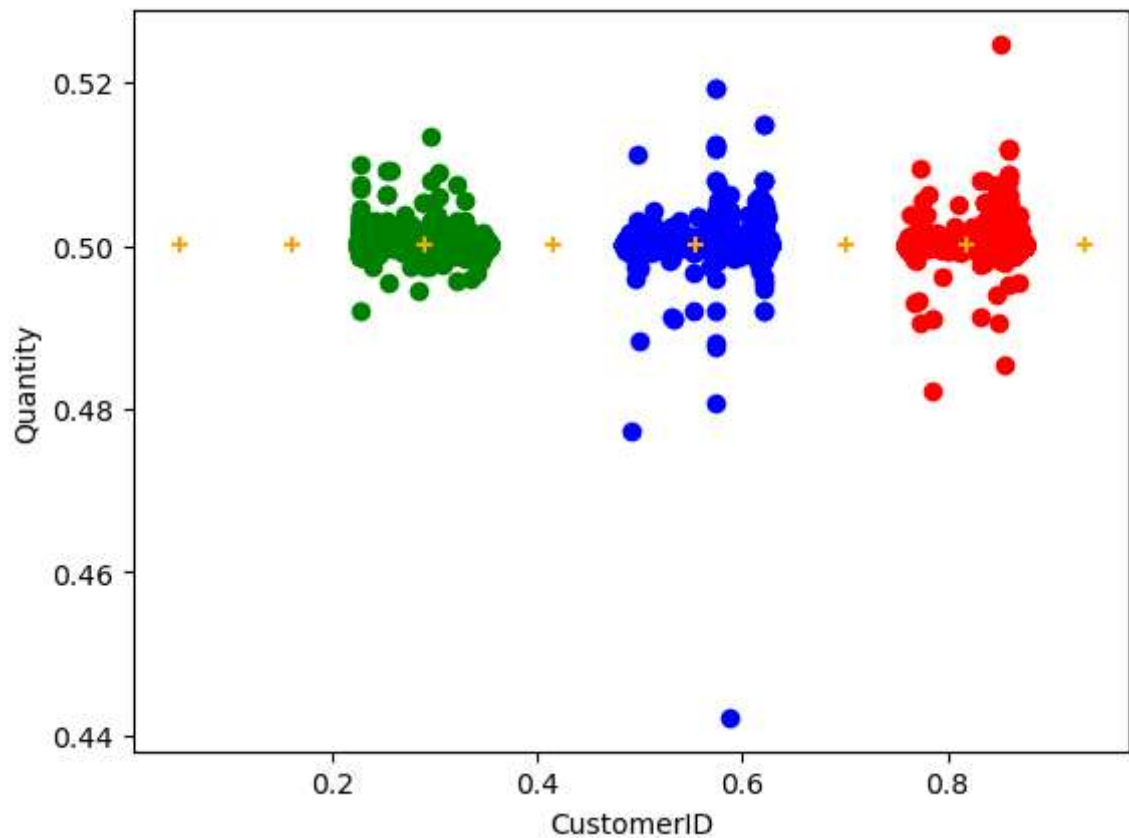
```
1 km.cluster_centers_
```

Out[24]: array([[0.81855044, 0.50006026],
               [0.29101351, 0.50006566],
               [0.5534736 , 0.50005383],
               [0.1603687 , 0.50005698],
               [0.70060666, 0.50005781],
               [0.41539441, 0.50005966],
               [0.93308721, 0.50005101],
               [0.05119252, 0.50006679]])

```
1  df1=df[df["New Cluster"]==0]
2  df2=df[df["New Cluster"]==1]
3  df3=df[df["New Cluster"]==2]
4  plt.scatter(df1["CustomerID"],df1["Quantity"],color="red")
5  plt.scatter(df2["CustomerID"],df2["Quantity"],color="green")
6  plt.scatter(df3["CustomerID"],df3["Quantity"],color="blue")
7  plt.scatter(km.cluster_centers_[:,0],km.cluster_centers_[:,1],color="oran
8  plt.xlabel("CustomerID")
9  plt.ylabel("Quantity")
```

Out[25]: Text(0, 0.5, 'Quantity')



In [26]:

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

```
In [27]:   1  for k in k_rng:
           2      km=KMeans(n_clusters=k)
           3      km.fit(df[["CustomerID","Quantity"]])
           4      sse.append(km.inertia_)
           5  #km.inertia_ will give you the value of sum of square error
           6  print(sse)
           7  plt.plot(k_rng,sse)
           8  plt.xlabel("K")
           9  plt.ylabel("Sum of Squared Error")
```

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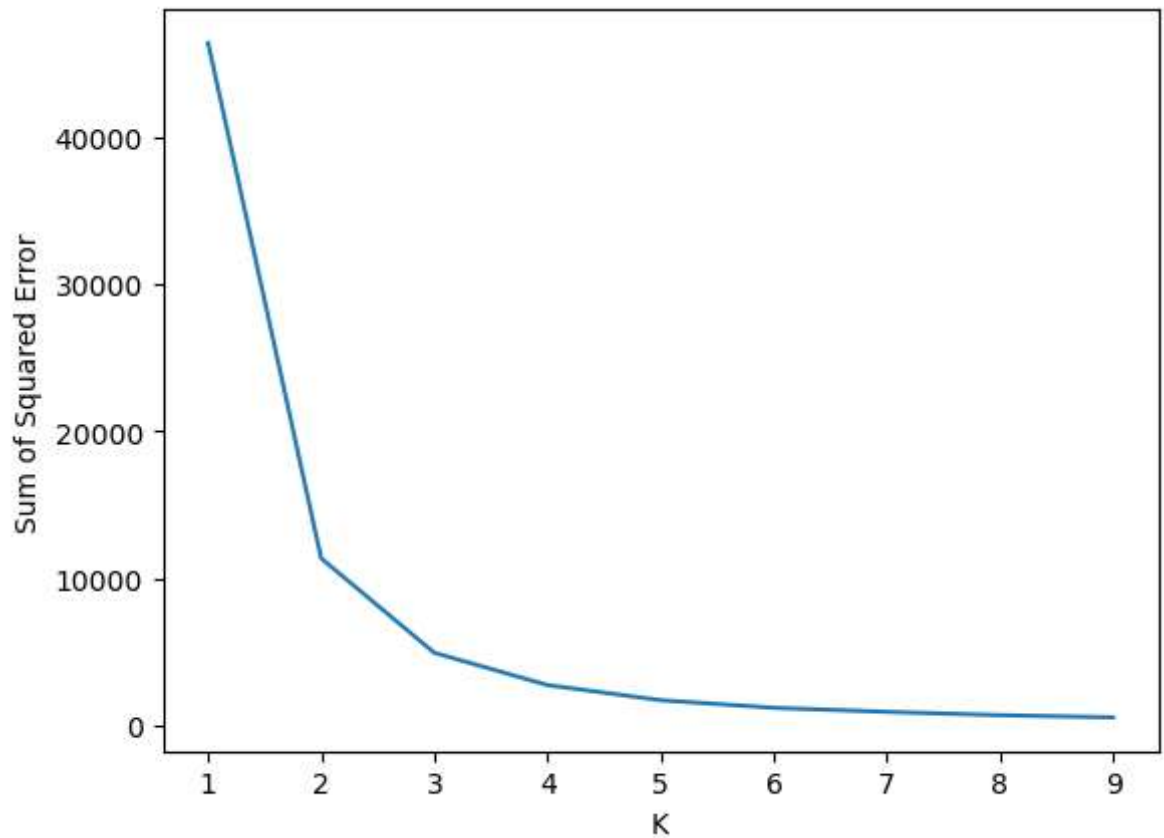1.4. Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
C:\ProgramData\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:870: F
utureWarning: The default value of `n_init` will change from 10 to 'auto' in
1.4. Set the value of `n_init` explicitly to suppress the warning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

[46374.84553398474, 11336.065305485301, 4918.443482888961, 2723.51910518953,
1695.069310119926, 1178.4435998084673, 902.8136802248464, 676.5837674800985,
528.3644172245184]

Out[27]:   Text(0, 0.5, 'Sum of Squared Error')
```

**CONCLUSION :For the above "Online retail"
dataset we use "K-means clustering" to divide
that data in to different clusters**