

## Project - 5 (DATASET: Online Retail) The transactions made by a UKbased, registered, non-store online retailer between December 1, 2010,

and December 9, 2011, are all included in the transnational data set known as online retail. The company primarily offers one-of-a-kind gifts for every occasion. The company has a large number of wholesalers as clients. Company Objective Using the global online retail dataset, we will design a clustering model and select the ideal group of clients for the business to target.

In [21]:

```
import pandas as pd
from matplotlib import pyplot as plt
%matplotlib inline
```

In [22]:

```
df=pd.read_csv(r"C:\Users\Teju\Downloads\OnlineRetail1.csv")
df
```

Out[22]:

	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 [23]:

```
df.head()
```

Out[23]:

	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	Unitec Kingdom
1	536365	71053	WHITE METAL LANTERN	6	01-12-2010 08:26	3.39	17850.0	Unitec Kingdom
2	536365	84406B	CREAM CUPID HEARTS COAT HANGER	8	01-12-2010 08:26	2.75	17850.0	Unitec Kingdom
3	536365	84029G	KNITTED UNION FLAG HOT WATER BOTTLE	6	01-12-2010 08:26	3.39	17850.0	Unitec Kingdom
4	536365	84029E	RED WOOLLY HOTTIE WHITE HEART.	6	01-12-2010 08:26	3.39	17850.0	Unitec Kingdom



In [24]:

```
df.tail()
```

Out[24]:

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



In [26]:

```
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 [27]:

```
df.isnull().sum()
```

Out[27]:

```
InvoiceNo      0
StockCode      0
Description    1454
Quantity       0
InvoiceDate    0
UnitPrice      0
CustomerID    135080
Country        0
dtype: int64
```

In [28]:

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

In [29]:

```
df.isnull().sum()
```

Out[29]:

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

In [30]:

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

Out[30]:

```
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 [31]:

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

Out[31]:

```
17841.0      8644
14911.0      7648
12748.0      6134
14096.0      5412
14606.0      3952
...
15753.0        1
14424.0        1
15562.0        1
13302.0        1
17331.0        1
Name: CustomerID, Length: 4372, dtype: int64
```

In [32]:

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

Out[32]:

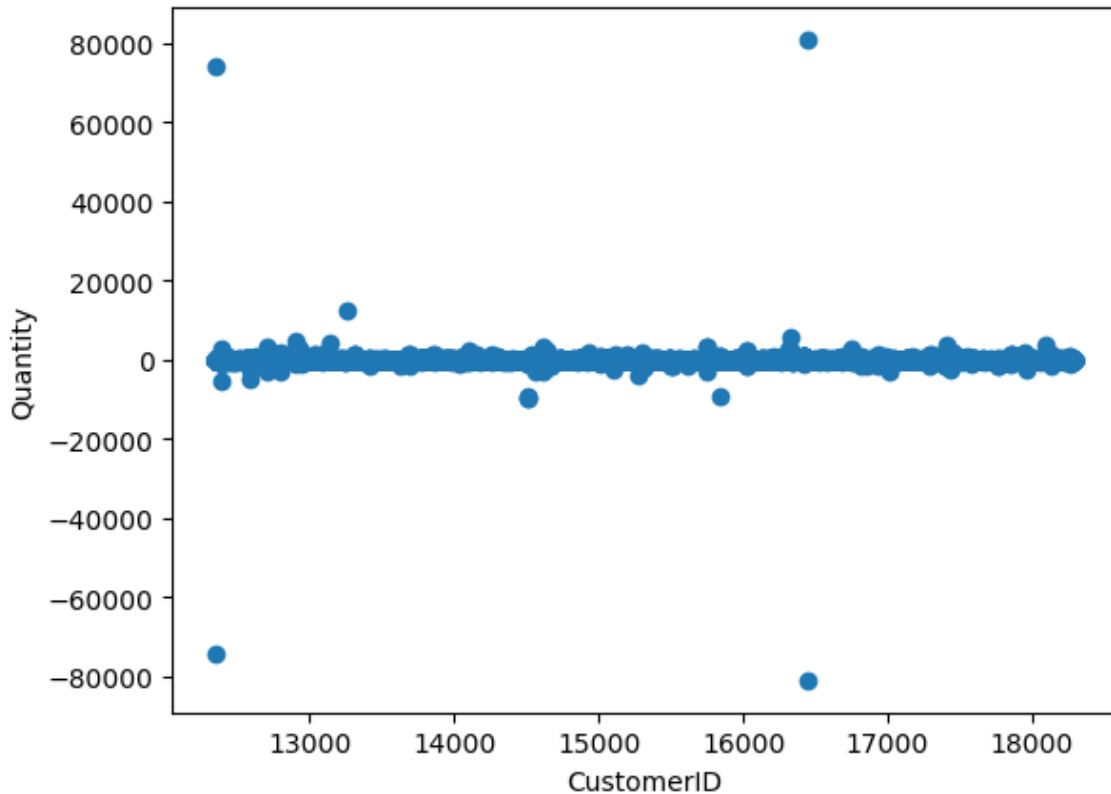
```
1      148227
2      81829
12     61063
6      40868
4      38484
...
-472        1
-161        1
-1206       1
-272        1
-80995      1
Name: Quantity, Length: 722, dtype: int64
```

In [33]:

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

Out[33]:

Text(0, 0.5, 'Quantity')



In [34]:

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

Out[34]:

```
▼ KMeans
KMeans()
```

In [35]:

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

C:\Users\Teju\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[35]:

```
array([4, 4, 4, ..., 2, 2, 2])
```

In [36]:

```
df["cluster"]=y_predicted
df.head()
```

Out[36]:

	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	Unitec Kingdom
1	536365	71053	WHITE METAL LANTERN	6	01-12-2010 08:26	3.39	17850.0	Unitec Kingdom
2	536365	84406B	CREAM CUPID HEARTS COAT HANGER	8	01-12-2010 08:26	2.75	17850.0	Unitec Kingdom
3	536365	84029G	KNITTED UNION FLAG HOT WATER BOTTLE	6	01-12-2010 08:26	3.39	17850.0	Unitec Kingdom
4	536365	84029E	RED WOOLLY HOTTIE WHITE HEART.	6	01-12-2010 08:26	3.39	17850.0	Unitec Kingdom

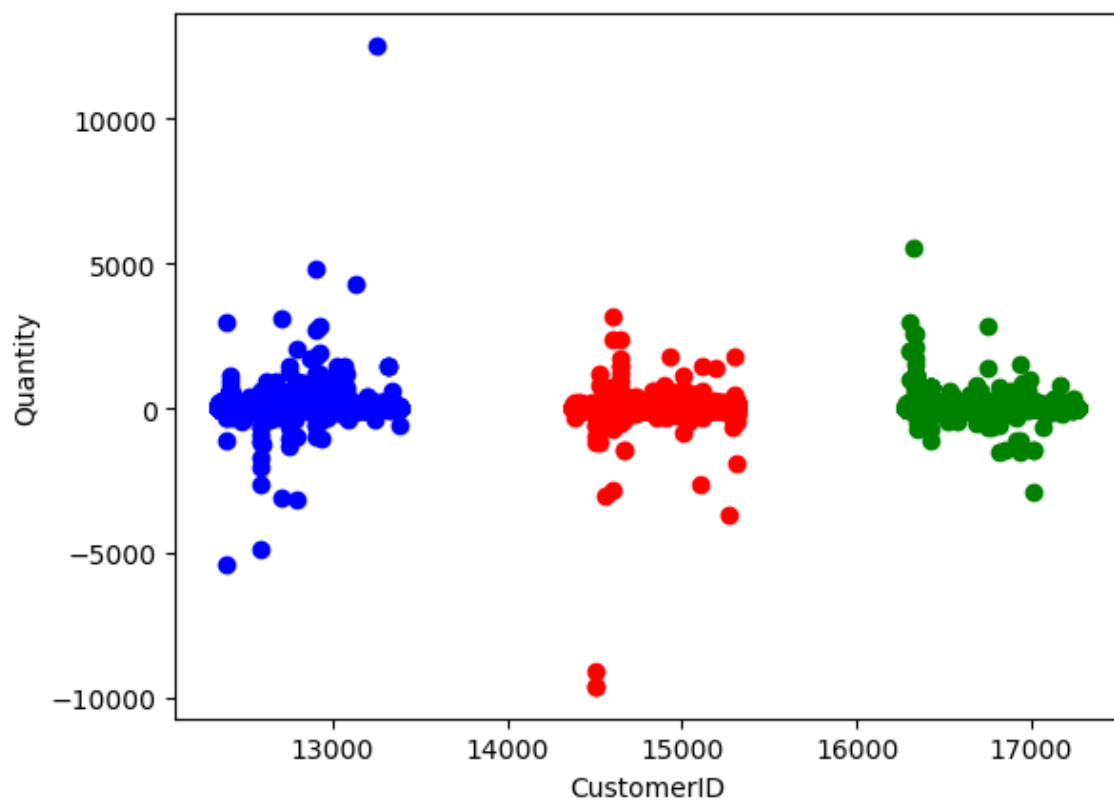


In [37]:

```
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[37]:

Text(0, 0.5, 'Quantity')





In [38]:

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

Out[38]:

	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	Unitec Kingdon
1	536365	71053	WHITE METAL LANTERN	0.500037	01-12-2010 08:26	3.39	17850.0	Unitec Kingdon
2	536365	84406B	CREAM CUPID HEARTS COAT HANGER	0.500049	01-12-2010 08:26	2.75	17850.0	Unitec Kingdon
3	536365	84029G	KNITTED UNION FLAG HOT WATER BOTTLE	0.500037	01-12-2010 08:26	3.39	17850.0	Unitec Kingdon
4	536365	84029E	RED WOOLLY HOTTIE WHITE HEART.	0.500037	01-12-2010 08:26	3.39	17850.0	Unitec Kingdon

In [39]:

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

Out[39]:

	InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID
0	536365	85123A	WHITE HANGING HEART T- LIGHT HOLDER	0.500037	01-12-2010 08:26	2.55	0.926443
1	536365	71053	WHITE METAL LANTERN	0.500037	01-12-2010 08:26	3.39	0.926443
2	536365	84406B	CREAM CUPID HEARTS COAT HANGER	0.500049	01-12-2010 08:26	2.75	0.926443
3	536365	84029G	KNITTED UNION FLAG HOT WATER BOTTLE	0.500037	01-12-2010 08:26	3.39	0.926443
4	536365	84029E	RED WOOLLY HOTTIE WHITE HEART.	0.500037	01-12-2010 08:26	3.39	0.926443
...	...	...	...	...	...	...	...
541904	581587	22613	PACK OF 20 SPACEBOY NAPKINS	0.500074	09-12-2011 12:50	0.85	0.056219
541905	581587	22899	CHILDREN'S APRON DOLLY GIRL	0.500037	09-12-2011 12:50	2.10	0.056219
541906	581587	23254	CHILDRENS CUTLERY DOLLY GIRL	0.500025	09-12-2011 12:50	4.15	0.056219
541907	581587	23255	CHILDRENS CUTLERY CIRCUS PARADE	0.500025	09-12-2011 12:50	4.15	0.056219
541908	581587	22138	BAKING SET 9 PIECE RETROSPOT	0.500019	09-12-2011 12:50	4.95	0.056219

541909 rows × 9 columns



In [40]:

```
km=KMeans()
```

In [41]:

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

C:\Users\Teju\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[41]:

```
array([2, 2, 2, ..., 5, 5, 5])
```

In [42]:

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

Out[42]:

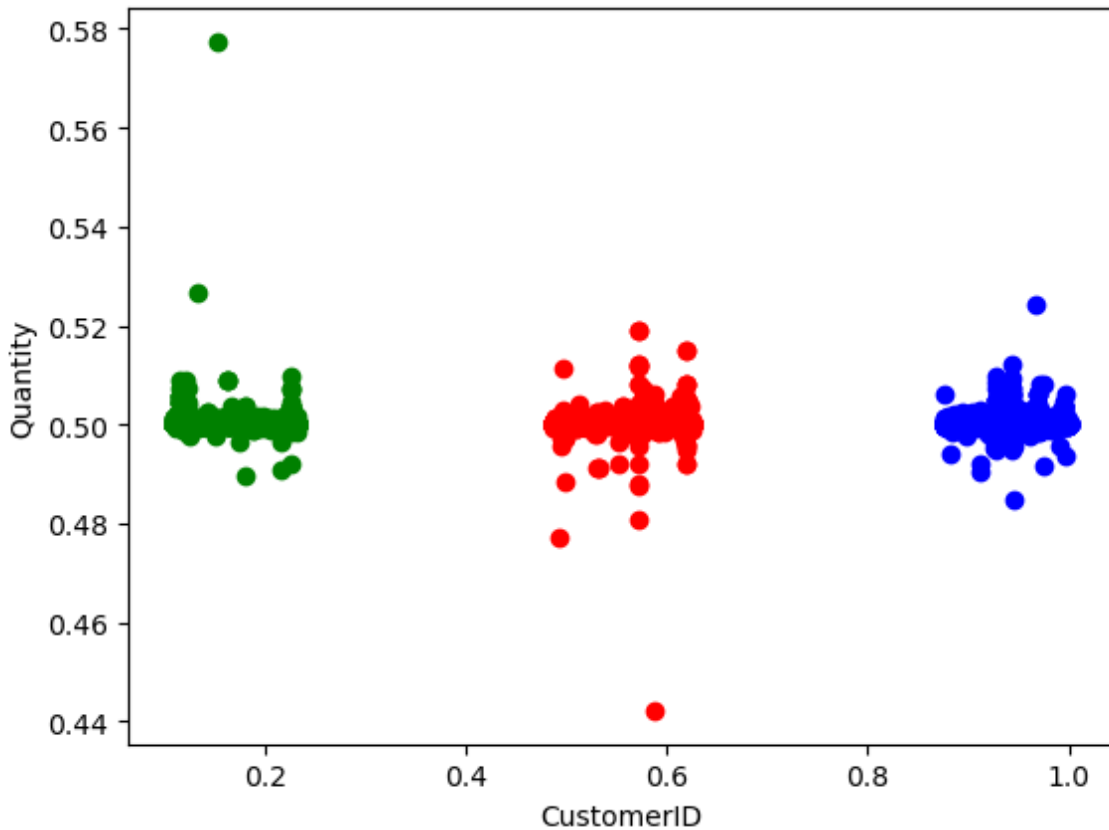
	InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country
0	536365	85123A	WHITE HANGING HEART T- LIGHT HOLDER	0.500037	01-12-2010 08:26	2.55	0.926443	Unitec Kingdon
1	536365	71053	WHITE METAL LANTERN	0.500037	01-12-2010 08:26	3.39	0.926443	Unitec Kingdon
2	536365	84406B	CREAM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

In [43]:

```
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[43]:

Text(0, 0.5, 'Quantity')



In [44]:

```
km.cluster_centers_
```

Out[44]:

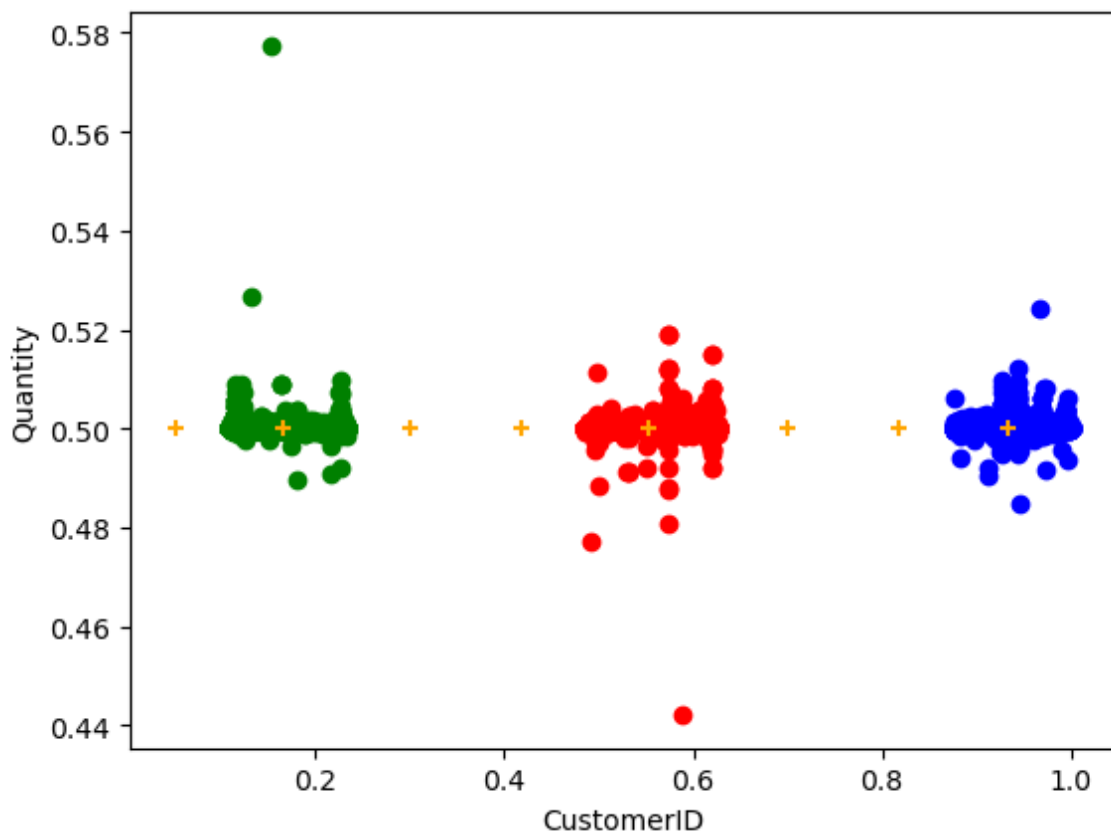
```
array([[0.55364751, 0.50005364],
       [0.16661839, 0.50006044],
       [0.93299701, 0.50005095],
       [0.69989479, 0.50005818],
       [0.29933484, 0.50006096],
       [0.05178062, 0.50006695],
       [0.41821262, 0.50006106],
       [0.81775848, 0.50005991]])
```

In [45]:

```
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[45]:

Text(0, 0.5, 'Quantity')



In [46]:

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

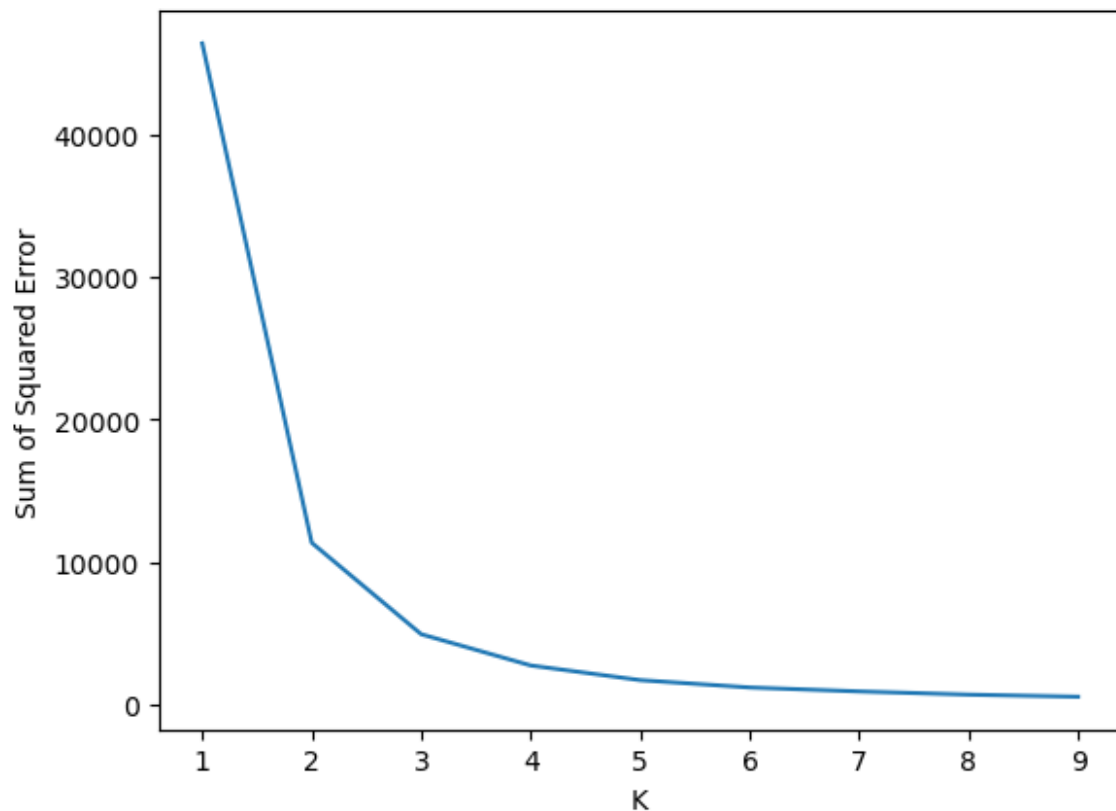
In [47]:

```
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\Teju\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\Teju\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\Teju\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\Teju\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\Teju\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\Teju\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\Teju\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\Teju\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(  
[46374.84553398485, 11336.065305485563, 4916.080763206521, 2723.5191051894  
626, 1695.0544265243825, 1178.5923367697794, 902.5867924503457, 676.549156  
1302019, 529.7115169779566]
```

Out[47]:

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



based on

the given data. In the

above dataset we will take customer id and quantity based on that we make the clusters. When the K-value is

low error rate is more and the K-value is high error rate is very high. So, finally we can conclude the above dataset is best fit for K-Means.

In [ ]: