

Project - 5 (DATASET: Online Retail) The transactions made by a UK-based, 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 [1]: 1 import pandas as pd
        2 from matplotlib import pyplot as plt
        3 %matplotlib inline
```

```
In [4]: 1 df=pd.read_csv(r"C:\Users\DELL E5490\Documents\OnlineRetail.csv")
        2 df
```

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

541909 rows × 8 columns

```
In [5]: 1 df.head()
```

Out[5]:

	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 [6]:

```
1 df.tail()
```

Out[6]:

	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	France
541905	581587	22899	CHILDREN'S APRON DOLLY GIRL	6	09-12-2011 12:50	2.10	12680.0	France
541906	581587	23254	CHILDRENS CUTLERY DOLLY GIRL	4	09-12-2011 12:50	4.15	12680.0	France
541907	581587	23255	CHILDRENS CUTLERY CIRCUS PARADE	4	09-12-2011 12:50	4.15	12680.0	France
541908	581587	22138	BAKING SET 9 PIECE RETROSPOT	3	09-12-2011 12:50	4.95	12680.0	France

In [7]:

```
1 df['InvoiceNo'].value_counts()
2
```

Out[7]:

```
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 [8]: 1 df['CustomerID'].value_counts()  
        2
```

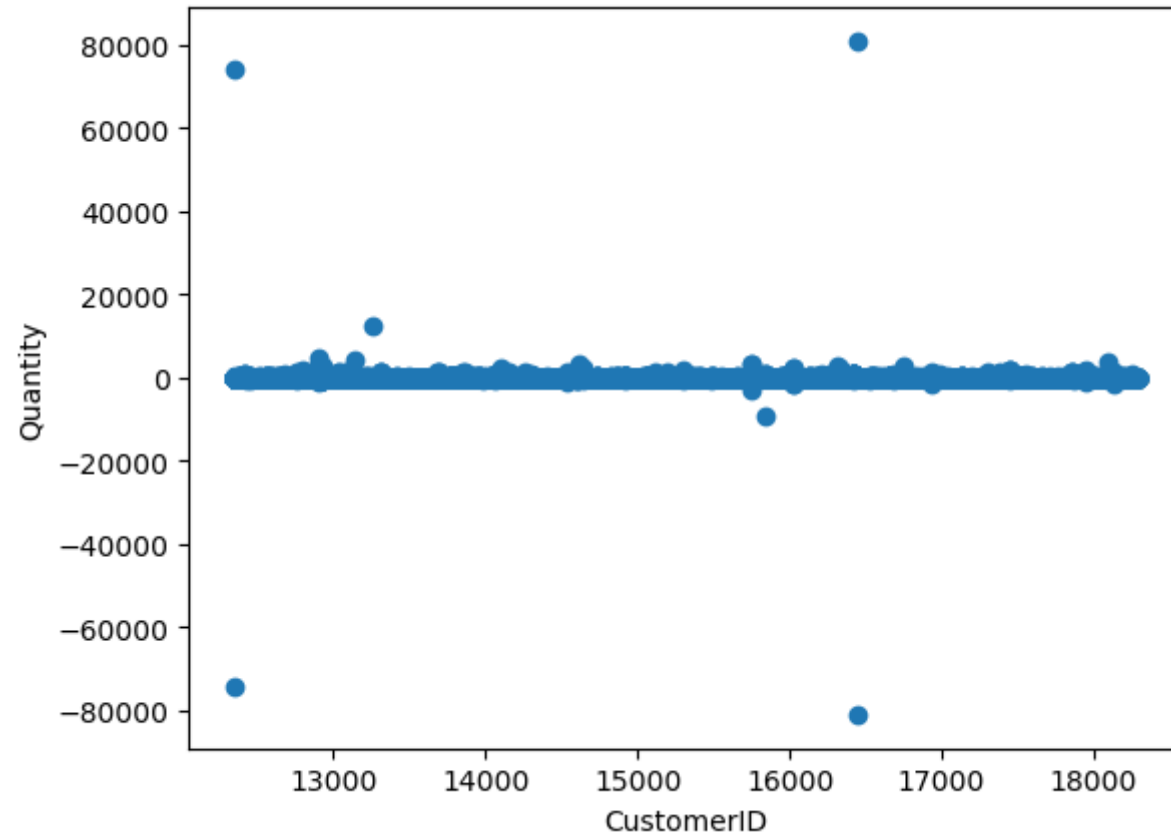
```
Out[8]: 17841.0    7983  
        14911.0    5903  
        14096.0    5128  
        12748.0    4642  
        14606.0    2782  
        ...  
        15070.0     1  
        15753.0     1  
        17065.0     1  
        16881.0     1  
        16995.0     1  
        Name: CustomerID, Length: 4372, dtype: int64
```

```
In [9]: 1 df['Quantity'].value_counts()
```

```
Out[9]: 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 [10]: 1 plt.scatter(df["CustomerID"],df["Quantity"])
          2 plt.xlabel("CustomerID")
          3 plt.ylabel("Quantity")
          4
```

```
Out[10]: Text(0, 0.5, 'Quantity')
```



In [11]:

```
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 [12]:

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

```
Out[12]: InvoiceNo          0
StockCode          0
Description        1454
Quantity           0
InvoiceDate        0
UnitPrice          0
CustomerID        135080
Country            0
dtype: int64
```

In [13]:

```
1 df.fillna(method='ffill',inplace=True)
2
```

```
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:\Users\DELL E5490\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[15]: array([2, 2, 2, ..., 1, 1, 1])
```

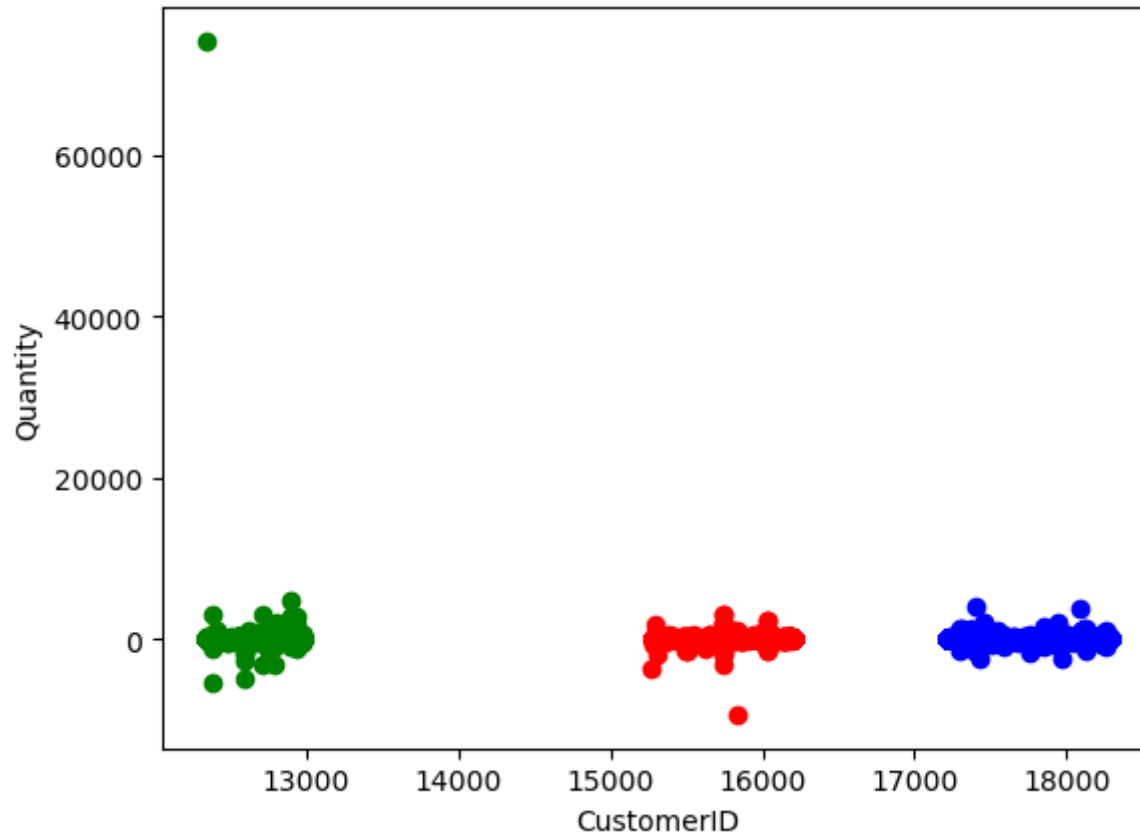
```
In [16]: 1 df["cluster"]=y_predicted
          2 df.head()
```

```
Out[16]:
```

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

```
In [17]: 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]: 1 from sklearn.preprocessing import MinMaxScaler
2 scaler=MinMaxScaler()
3 scaler.fit(df[["Quantity"]])
4 df["Quantity"]=scaler.transform(df[["Quantity"]])
5 df.head()
```

Out[18]:

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

```
In [19]: 1 scaler.fit(df[["CustomerID"]])
2 df["CustomerID"]=scaler.transform(df[["CustomerID"]])
3 df
```

Out[19]:

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

541909 rows × 9 columns

K-MeansClustering

```
In [20]: 1 km=KMeans()  
        2
```

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

C:\Users\DELL E5490\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[21]: array([4, 4, 4, ..., 2, 2, 2])

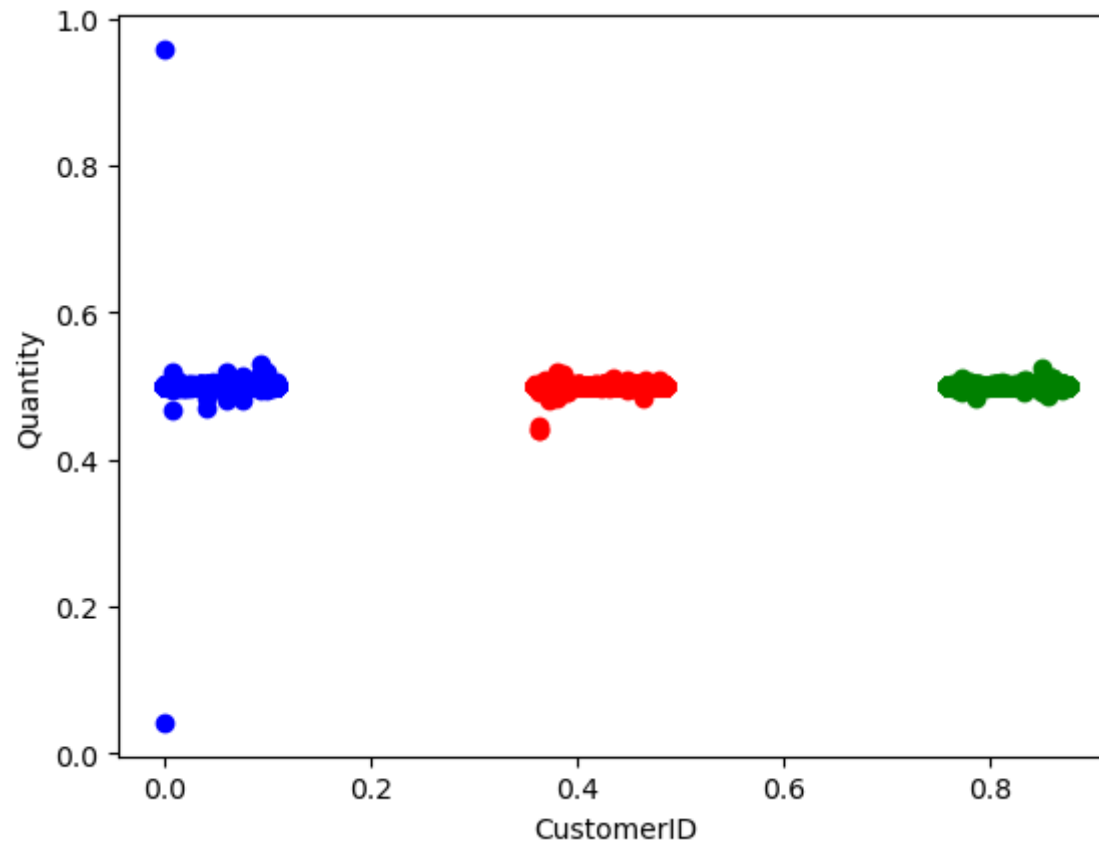
```
In [22]: 1 df["New Cluster"]=y_predicted  
        2 df.head()
```

Out[22]:

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

```
In [23]: 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')

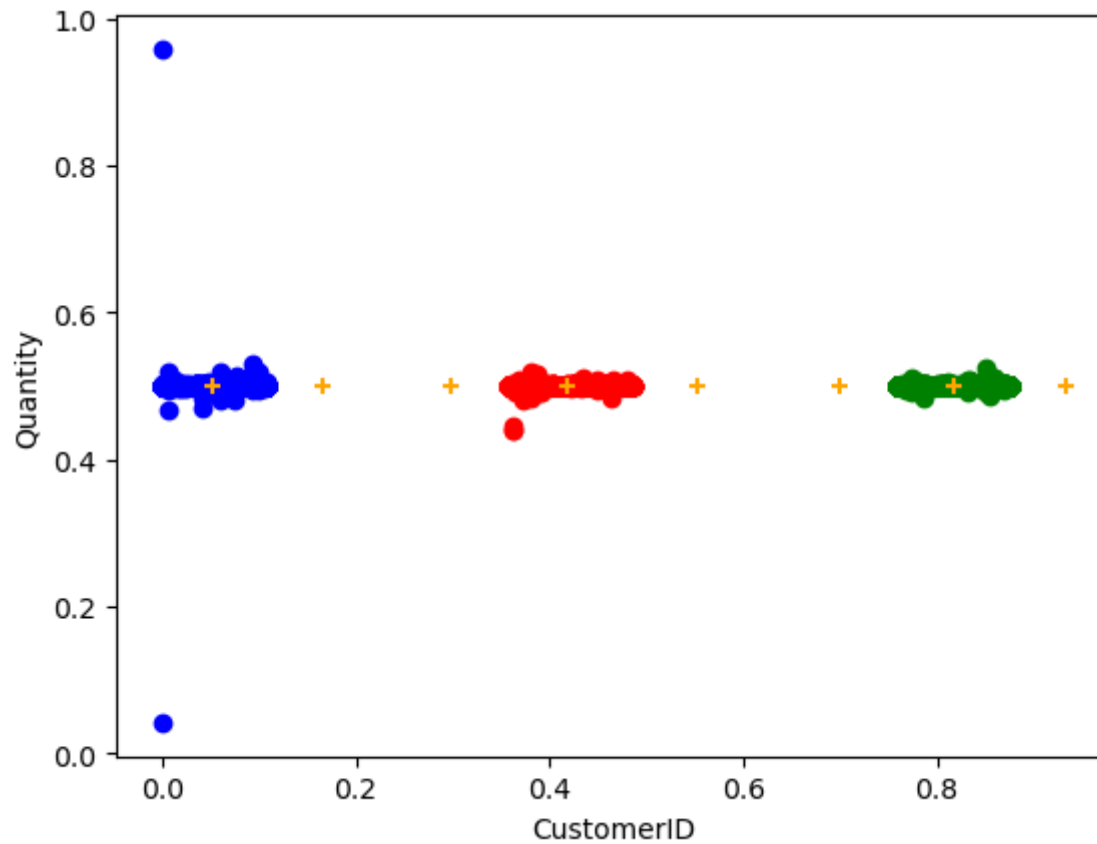


```
In [24]: 1 km.cluster_centers_
```

```
Out[24]: array([[0.41790198, 0.50006105],  
                [0.81756576, 0.50005988],  
                [0.05156814, 0.50006705],  
                [0.55346613, 0.50005369],  
                [0.9328779 , 0.50005088],  
                [0.29827906, 0.50006068],  
                [0.69973228, 0.50005831],  
                [0.16561407, 0.50006062]])
```

```
In [25]: 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="orange",marker="+")
8 plt.xlabel("CustomerID")
9 plt.ylabel("Quantity")
10
```

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:\Users\DELL E5490\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\DELL E5490\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\DELL E5490\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\DELL E5490\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\DELL E5490\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\DELL E5490\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\DELL E5490\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\DELL E5490\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

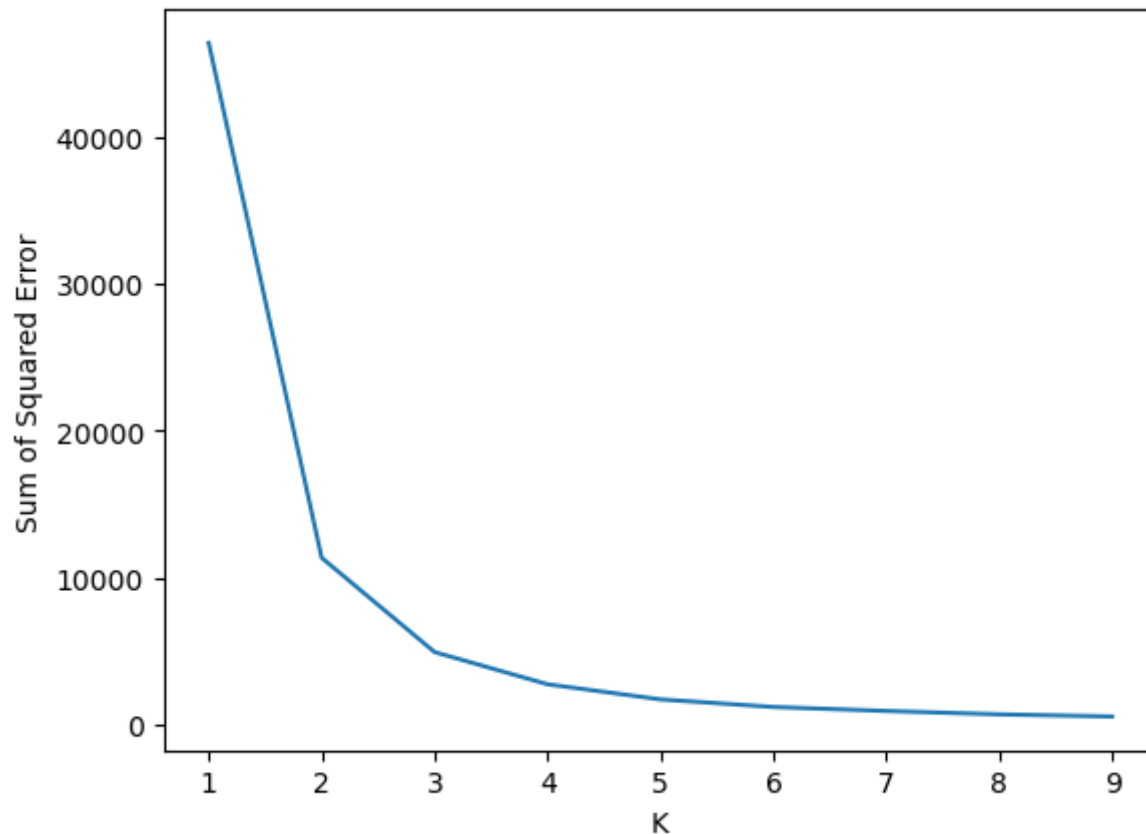
warnings.warn(

C:\Users\DELL E5490\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, 11336.065305485301, 4915.889395801374, 2723.519105189529, 1695.0392229313316, 1178.6175555318514, 907.7098539670895, 677.3921023858019, 530.3184466596015]

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



CONCLUSION

For the given dataset we use K-means Clustering and done the grouping 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



In []:

1