# Miniproject-5

## Problem statement: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 ObjectiveUsing 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]:
```python
import pandas as pd
from matplotlib import pyplot as plt
%matplotlib inline
```

In [2]: 
```
df=pd.read_csv(r"C:\Users\LENOVO\Downloads\Online Retail.csv")
df
```

Out[2]:

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

541909 rows × 8 columns

In [3]: `df.describe`

Out[3]: 
```
<bound method NDFrame.describe of          InvoiceNo StockCode
Description  Quantity
0          536365    85123A    WHITE HANGING HEART T-LIGHT HOLDER         6  \
1          536365     71053                   WHITE METAL LANTERN         6
2          536365    84406B        CREAM CUPID HEARTS COAT HANGER         8
3          536365    84029G  KNITTED UNION FLAG HOT WATER BOTTLE         6
4          536365    84029E        RED WOOLLY HOTTIE WHITE HEART.         6
...           ...       ...                                  ...       ...
541904     581587     22613          PACK OF 20 SPACEBOY NAPKINS        12
541905     581587     22899         CHILDREN'S APRON DOLLY GIRL         6
541906     581587     23254        CHILDRENS CUTLERY DOLLY GIRL         4
541907     581587     23255      CHILDRENS CUTLERY CIRCUS PARADE         4
541908     581587     22138        BAKING SET 9 PIECE RETROSPOT         3

               InvoiceDate  UnitPrice  CustomerID         Country
0          01-12-2010 08:26       2.55     17850.0  United Kingdom
1          01-12-2010 08:26       3.39     17850.0  United Kingdom
2          01-12-2010 08:26       2.75     17850.0  United Kingdom
3          01-12-2010 08:26       3.39     17850.0  United Kingdom
4          01-12-2010 08:26       3.39     17850.0  United Kingdom
...                     ...        ...         ...             ...
541904     09-12-2011 12:50       0.85     12680.0          France
541905     09-12-2011 12:50       2.10     12680.0          France
541906     09-12-2011 12:50       4.15     12680.0          France
541907     09-12-2011 12:50       4.15     12680.0          France
541908     09-12-2011 12:50       4.95     12680.0          France

[541909 rows x 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 | Coun |
|---|---|---|---|---|---|---|---|---|
| 541904 | 581587 | 22613 | PACK OF 20 SPACEBOY NAPKINS | 12 | 09-12-2011 12:50 | 0.85 | 12680.0 | Fran |
| 541905 | 581587 | 22899 | CHILDREN'S APRON DOLLY GIRL | 6 | 09-12-2011 12:50 | 2.10 | 12680.0 | Fran |
| 541906 | 581587 | 23254 | CHILDRENS CUTLERY DOLLY GIRL | 4 | 09-12-2011 12:50 | 4.15 | 12680.0 | Fran |
| 541907 | 581587 | 23255 | CHILDRENS CUTLERY CIRCUS PARADE | 4 | 09-12-2011 12:50 | 4.15 | 12680.0 | Fran |
| 541908 | 581587 | 22138 | BAKING SET 9 PIECE RETROSPOT | 3 | 09-12-2011 12:50 | 4.95 | 12680.0 | Fran |

In [6]: `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 [7]: 
```python
df.fillna(method="ffill",inplace=True)
df
```

Out[7]:

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

541909 rows × 8 columns

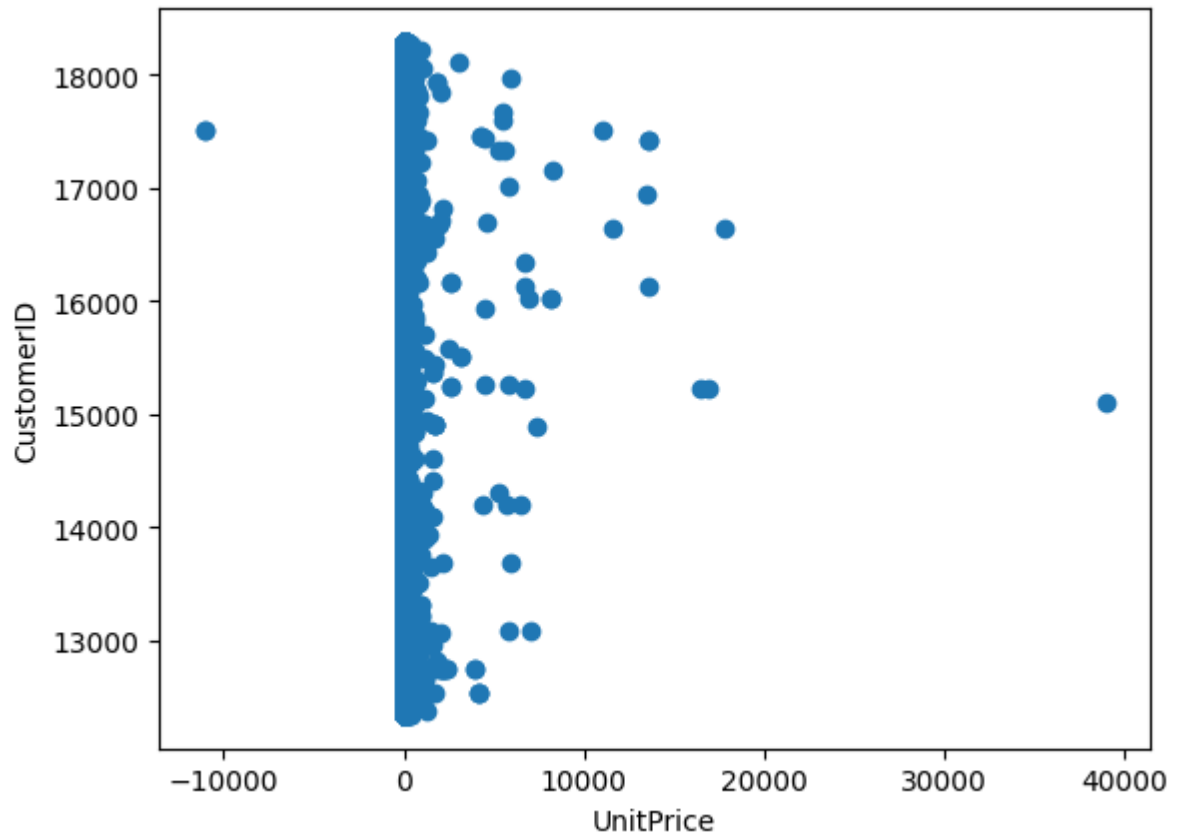In [8]: `df.isna().sum()`

Out[8]: 
```
InvoiceNo      0
StockCode      0
Description    0
Quantity       0
InvoiceDate    0
UnitPrice      0
CustomerID     0
Country        0
dtype: int64
```

In [9]: `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  541909 non-null  object
 3   Quantity     541909 non-null  int64
 4   InvoiceDate  541909 non-null  object
 5   UnitPrice    541909 non-null  float64
 6   CustomerID   541909 non-null  float64
 7   Country      541909 non-null  object
dtypes: float64(2), int64(1), object(5)
memory usage: 33.1+ MB
```

In [10]:
```python
plt.scatter(df["UnitPrice"],df["CustomerID"])
plt.xlabel("UnitPrice")
plt.ylabel("CustomerID")
```

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



In [12]:
```python
from sklearn.cluster import KMeans
km=KMeans()
km
```

Out[12]:  KMeans()

**In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.**
**On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.**

In [13]:
```python
y_predicted=km.fit_predict(df[["UnitPrice","CustomerID"]])
y_predicted
```

```
C:\Users\LENOVO\AppData\Local\Programs\Python\Python311\Lib\site-packages\skl
earn\cluster\_kmeans.py:870: FutureWarning: The default value of `n_init` wil
l change from 10 to 'auto' in 1.4. Set the value of `n_init` explicitly to su
ppress the warning
    warnings.warn(
```
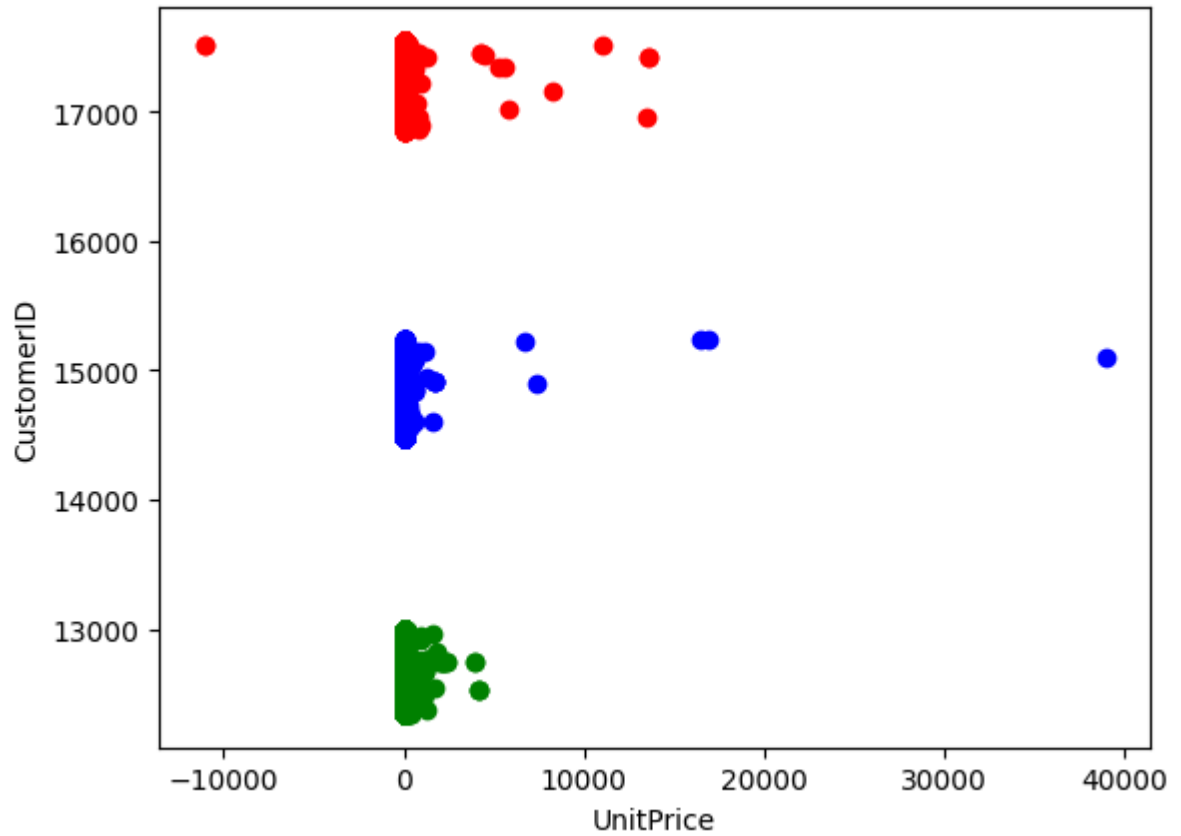
Out[13]: array([5, 5, 5, ..., 1, 1, 1])

In [14]:
```
df["cluster"]=y_predicted
df.head()
```

Out[14]:

| | InvoiceNo | StockCode | Description | Quantity | InvoiceDate | UnitPrice | CustomerID | Country | clu |
|---|---|---|---|---|---|---|---|---|---|
| **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 [15]:
```python
df1=df[df.cluster==0]
df2=df[df.cluster==1]
df3=df[df.cluster==2]
plt.scatter(df1["UnitPrice"],df1["CustomerID"],color="red")
plt.scatter(df2["UnitPrice"],df2["CustomerID"],color="green")
plt.scatter(df3["UnitPrice"],df3["CustomerID"],color="blue")
plt.xlabel("UnitPrice")
plt.ylabel("CustomerID")
```

Out[15]: Text(0, 0.5, 'CustomerID')

In [16]:
```python
from sklearn.preprocessing import MinMaxScaler
scaler=MinMaxScaler()
scaler.fit(df[["CustomerID"]])
df["CustomerID"]=scaler.transform(df[["CustomerID"]])
df.head()
```

Out[16]:

| | InvoiceNo | StockCode | Description | Quantity | InvoiceDate | UnitPrice | CustomerID | Country | cl |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 536365 | 85123A | WHITE HANGING HEART T-LIGHT HOLDER | 6 | 01-12-2010 08:26 | 2.55 | 0.926443 | United Kingdom | |
| 1 | 536365 | 71053 | WHITE METAL LANTERN | 6 | 01-12-2010 08:26 | 3.39 | 0.926443 | United Kingdom | |
| 2 | 536365 | 84406B | CREAM CUPID HEARTS COAT HANGER | 8 | 01-12-2010 08:26 | 2.75 | 0.926443 | United Kingdom | |
| 3 | 536365 | 84029G | KNITTED UNION FLAG HOT WATER BOTTLE | 6 | 01-12-2010 08:26 | 3.39 | 0.926443 | United Kingdom | |
| 4 | 536365 | 84029E | RED WOOLLY HOTTIE WHITE HEART. | 6 | 01-12-2010 08:26 | 3.39 | 0.926443 | United Kingdom | |

In [17]:
```python
scaler.fit(df[["UnitPrice"]])
df["UnitPrice"]=scaler.transform(df[["UnitPrice"]])
df.head()
```

Out[17]:

| | InvoiceNo | StockCode | Description | Quantity | InvoiceDate | UnitPrice | CustomerID | Country | cl |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 536365 | 85123A | WHITE HANGING HEART T-LIGHT HOLDER | 6 | 01-12-2010 08:26 | 0.221150 | 0.926443 | United Kingdom | |
| 1 | 536365 | 71053 | WHITE METAL LANTERN | 6 | 01-12-2010 08:26 | 0.221167 | 0.926443 | United Kingdom | |
| 2 | 536365 | 84406B | CREAM CUPID HEARTS COAT HANGER | 8 | 01-12-2010 08:26 | 0.221154 | 0.926443 | United Kingdom | |
| 3 | 536365 | 84029G | KNITTED UNION FLAG HOT WATER BOTTLE | 6 | 01-12-2010 08:26 | 0.221167 | 0.926443 | United Kingdom | |
| 4 | 536365 | 84029E | RED WOOLLY HOTTIE WHITE HEART. | 6 | 01-12-2010 08:26 | 0.221167 | 0.926443 | United Kingdom | |

In [18]:
```python
km=KMeans()
```

In [19]:
```python
y_predicted=km.fit_predict(df[["UnitPrice","CustomerID"]])
y_predicted
```

```
C:\Users\LENOVO\AppData\Local\Programs\Python\Python311\Lib\site-packages\skl
earn\cluster\_kmeans.py:870: FutureWarning: The default value of `n_init` wil
l change from 10 to 'auto' in 1.4. Set the value of `n_init` explicitly to su
ppress the warning
  warnings.warn(
```
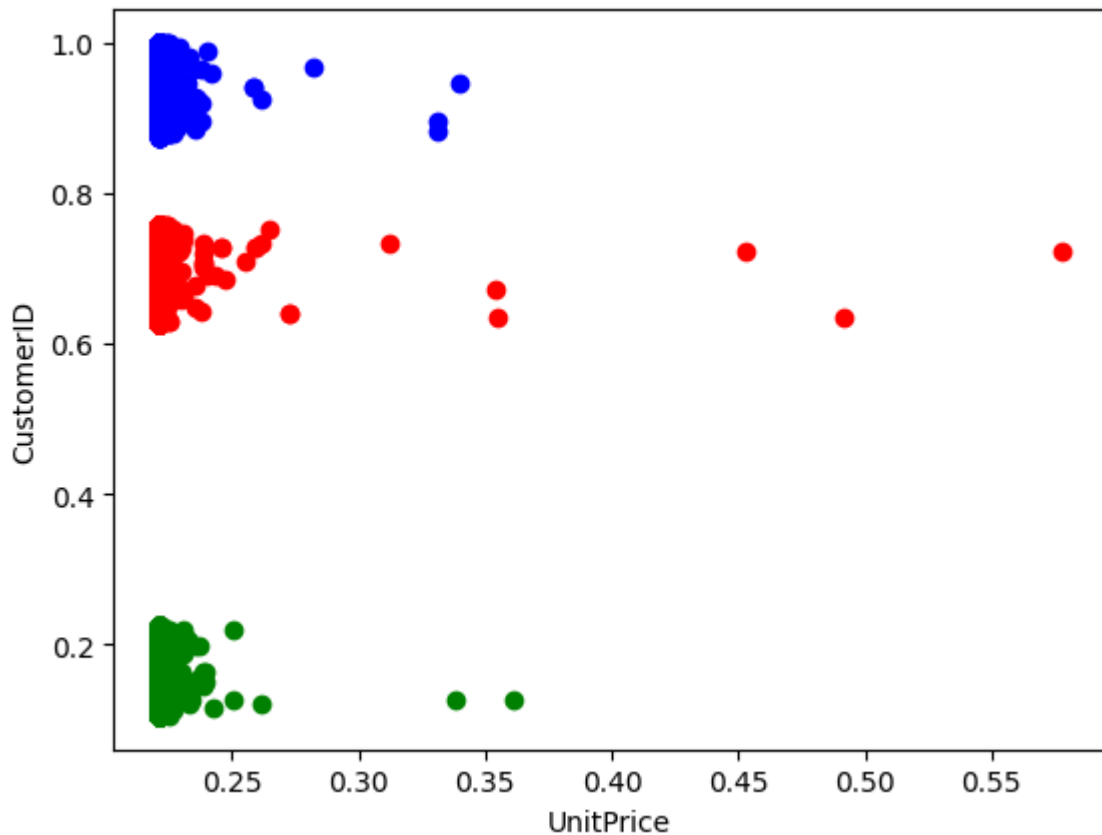
Out[19]: array([2, 2, 2, ..., 5, 5, 5])

In [20]:
```python
df["New Cluster"]=y_predicted
df.head()
```

Out[20]:

| | InvoiceNo | StockCode | Description | Quantity | InvoiceDate | UnitPrice | CustomerID | Country | clu |
|---|---|---|---|---|---|---|---|---|---|
| **0** | 536365 | 85123A | WHITE HANGING HEART T-LIGHT HOLDER | 6 | 01-12-2010 08:26 | 0.221150 | 0.926443 | United Kingdom | |
| **1** | 536365 | 71053 | WHITE METAL LANTERN | 6 | 01-12-2010 08:26 | 0.221167 | 0.926443 | United Kingdom | |
| **2** | 536365 | 84406B | CREAM CUPID HEARTS COAT HANGER | 8 | 01-12-2010 08:26 | 0.221154 | 0.926443 | United Kingdom | |
| **3** | 536365 | 84029G | KNITTED UNION FLAG HOT WATER BOTTLE | 6 | 01-12-2010 08:26 | 0.221167 | 0.926443 | United Kingdom | |
| **4** | 536365 | 84029E | RED WOOLLY HOTTIE WHITE HEART. | 6 | 01-12-2010 08:26 | 0.221167 | 0.926443 | United Kingdom | |

```
In [21]: df1=df[df["New Cluster"]==0]
         df2=df[df["New Cluster"]==1]
         df3=df[df["New Cluster"]==2]
         plt.scatter(df1["UnitPrice"],df1["CustomerID"],color="r")
         plt.scatter(df2["UnitPrice"],df2["CustomerID"],color="green")
         plt.scatter(df3["UnitPrice"],df3["CustomerID"],color="blue")
         plt.xlabel("UnitPrice")
         plt.ylabel("CustomerID")
```
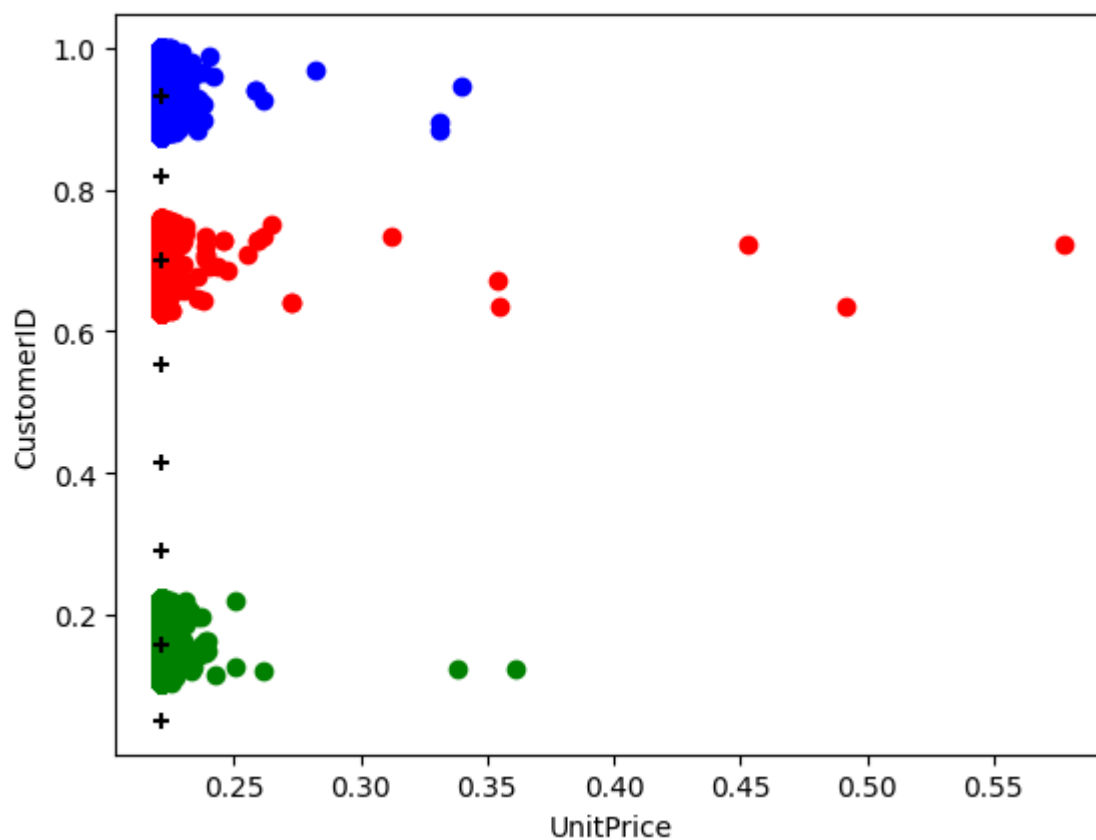
Out[21]: Text(0, 0.5, 'CustomerID')



```
In [22]: km.cluster_centers_
```

Out[22]: array([[0.2211979 , 0.69959256],
                [0.2211822 , 0.15893953],
                [0.2211783 , 0.93299701],
                [0.22118451, 0.41407851],
                [0.22119899, 0.5523978 ],
                [0.22120352, 0.05056666],
                [0.22119949, 0.81774878],
                [0.22119073, 0.28944002]])
```

In [23]:
```python
df1=df[df["New Cluster"]==0]
df2=df[df["New Cluster"]==1]
df3=df[df["New Cluster"]==2]
plt.scatter(df1["UnitPrice"],df1["CustomerID"],color="red")
plt.scatter(df2["UnitPrice"],df2["CustomerID"],color="green")
plt.scatter(df3["UnitPrice"],df3["CustomerID"],color="blue")
plt.scatter(km.cluster_centers_[:,0],km.cluster_centers_[:,1],color="black",ma
plt.xlabel("UnitPrice")
plt.ylabel("CustomerID")
```

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

```
In [24]: k_rng=range(1,10)
         sse=[]
         for k in k_rng:
             km=KMeans(n_clusters=k)
             km.fit(df[["UnitPrice","CustomerID"]])
             sse.append(km.inertia_) #km.inertia_ will give you the value of sum of squ
         print(sse)
         plt.plot(k_rng,sse)
         plt.xlabel("K")
         plt.ylabel("Sum of Squared Error")
```

```
C:\Users\LENOVO\AppData\Local\Programs\Python\Python311\Lib\site-packages\skl
earn\cluster\_kmeans.py:870: FutureWarning: The default value of `n_init` wil
l change from 10 to 'auto' in 1.4. Set the value of `n_init` explicitly to su
ppress the warning
  warnings.warn(
C:\Users\LENOVO\AppData\Local\Programs\Python\Python311\Lib\site-packages\skl
earn\cluster\_kmeans.py:870: FutureWarning: The default value of `n_init` wil
l change from 10 to 'auto' in 1.4. Set the value of `n_init` explicitly to su
ppress the warning
  warnings.warn(
C:\Users\LENOVO\AppData\Local\Programs\Python\Python311\Lib\site-packages\skl
earn\cluster\_kmeans.py:870: FutureWarning: The default value of `n_init` wil
l change from 10 to 'auto' in 1.4. Set the value of `n_init` explicitly to su
ppress the warning
  warnings.warn(
C:\Users\LENOVO\AppData\Local\Programs\Python\Python311\Lib\site-packages\skl
earn\cluster\_kmeans.py:870: FutureWarning: The default value of `n_init` wil
l change from 10 to 'auto' in 1.4. Set the value of `n_init` explicitly to su
ppress the warning
  warnings.warn(
C:\Users\LENOVO\AppData\Local\Programs\Python\Python311\Lib\site-packages\skl
earn\cluster\_kmeans.py:870: FutureWarning: The default value of `n_init` wil
l change from 10 to 'auto' in 1.4. Set the value of `n_init` explicitly to su
ppress the warning
  warnings.warn(
C:\Users\LENOVO\AppData\Local\Programs\Python\Python311\Lib\site-packages\skl
earn\cluster\_kmeans.py:870: FutureWarning: The default value of `n_init` wil
l change from 10 to 'auto' in 1.4. Set the value of `n_init` explicitly to su
ppress the warning
  warnings.warn(
C:\Users\LENOVO\AppData\Local\Programs\Python\Python311\Lib\site-packages\skl
earn\cluster\_kmeans.py:870: FutureWarning: The default value of `n_init` wil
l change from 10 to 'auto' in 1.4. Set the value of `n_init` explicitly to su
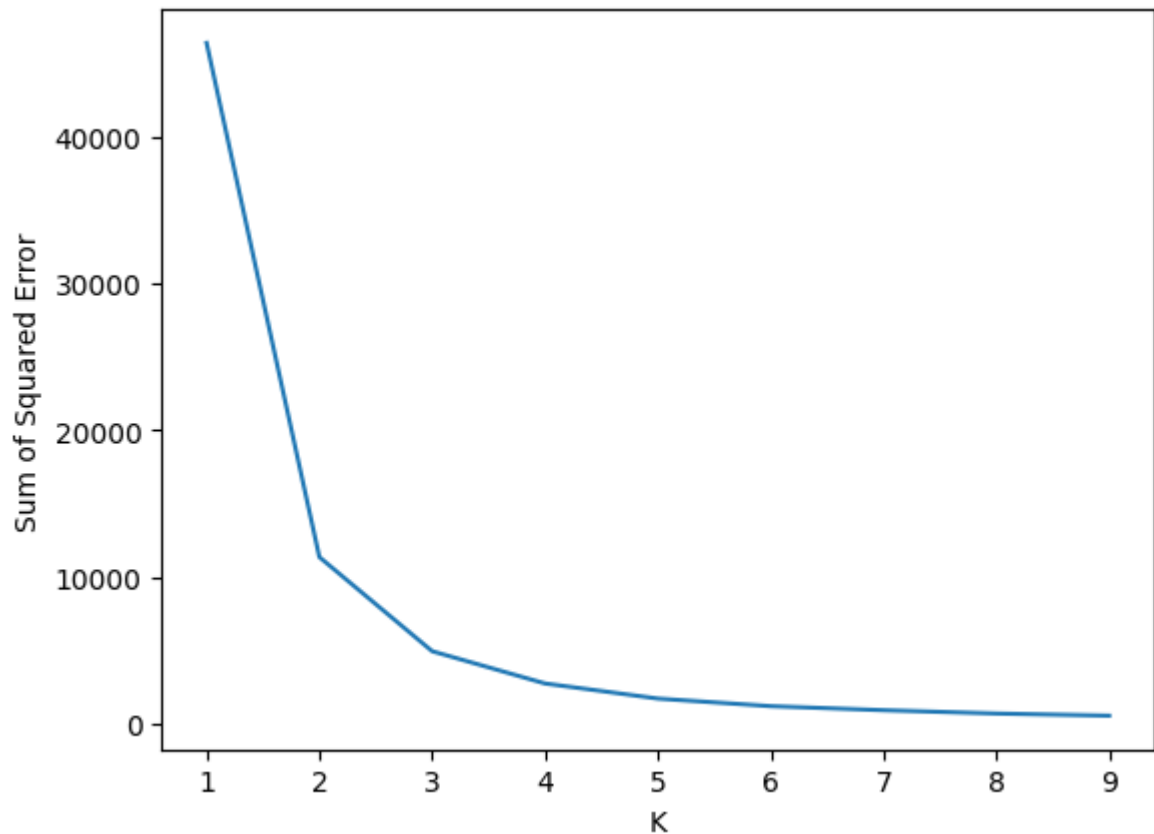ppress the warning
  warnings.warn(
C:\Users\LENOVO\AppData\Local\Programs\Python\Python311\Lib\site-packages\skl
earn\cluster\_kmeans.py:870: FutureWarning: The default value of `n_init` wil
l change from 10 to 'auto' in 1.4. Set the value of `n_init` explicitly to su
ppress the warning
  warnings.warn(

[46375.89020547866, 11337.110496294004, 4922.070144891326, 2724.563781877139
7, 1696.5129440430824, 1179.5364277462713, 903.6402640978915, 677.89745890048
93, 531.9668573302044]
```

Out[24]: Text(0, 0.5, 'Sum of Squared Error')

**Conclusion:For the above dataset we use K-means Clustering and performed the grouping based on the given data,In the above dataset we will take unit price and customerid 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 bestfit for K-Means.**

In [ ]: