

# ONLINE RETAIL

## K MEANS CLUSTER

**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 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]: import pandas as pd
        from matplotlib import pyplot as plt
        %matplotlib inline
```

## Data Collection

```
In [3]: df=pd.read_csv(r"C:\Users\91756\Documents\python\online.csv")
df
```

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
...	...	...	...	...	...	...	...	...
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

# Dta Collection

```
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]:
```

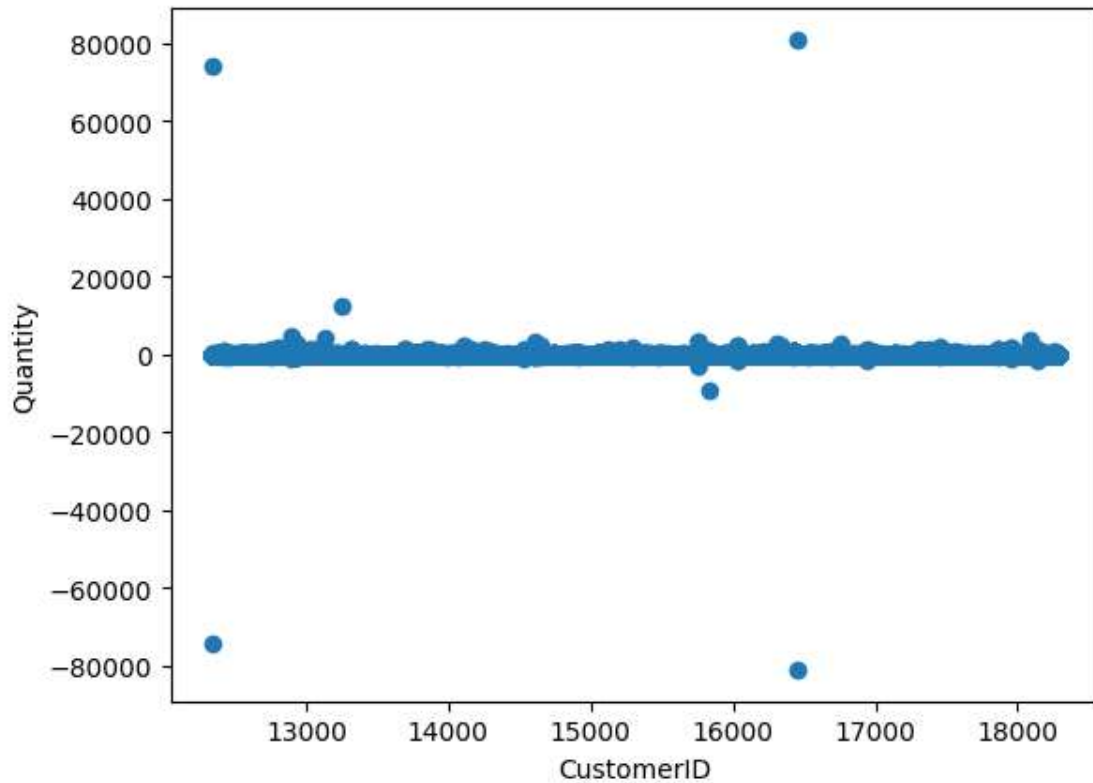
	InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country
541904	581587	22613	PACK OF 20 SPACEBOY NAPKINS	12	09-12-2011 12:50	0.85	12680.0	France
541905	581587	22899	CHILDREN'S APRON DOLLY GIRL	6	09-12-2011 12:50	2.10	12680.0	France
541906	581587	23254	CHILDRENS CUTLERY DOLLY 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 [6]: df['Description'].value_counts()
```

```
Out[6]: Description
WHITE HANGING HEART T-LIGHT HOLDER    2369
REGENCY CAKESTAND 3 TIER              2200
JUMBO BAG RED RETROSPOT               2159
PARTY BUNTING                       1727
LUNCH BAG RED RETROSPOT              1638
...
Missing                               1
historic computer difference?....se   1
DUSTY PINK CHRISTMAS TREE 30CM       1
WRAP BLUE RUSSIAN FOLKART            1
PINK BERTIE MOBILE PHONE CHARM       1
Name: count, Length: 4223, dtype: int64
```

```
In [7]: plt.scatter(df["CustomerID"],df["Quantity"])
plt.xlabel("CustomerID")
plt.ylabel("Quantity")
```

```
Out[7]: Text(0, 0.5, 'Quantity')
```



```
In [8]: 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 [9]: df.isnull().sum()
```

```
Out[9]: InvoiceNo      0
        StockCode     0
        Description  1454
        Quantity     0
        InvoiceDate    0
        UnitPrice     0
        CustomerID   135080
        Country       0
        dtype: int64
```

```
In [10]: df.fillna(method='ffill',inplace=True)
```

```
In [11]: df.isnull().sum()
```

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

```
In [12]: from sklearn.cluster import KMeans
```

```
In [13]: km=KMeans()
        km
```

```
Out[13]: 

▼ KMeans


        KMeans()
```

```
In [14]: y_predicted=km.fit_predict(df[["CustomerID","Quantity"]])
        y_predicted
```

```
C:\Users\91756\AppData\Local\Programs\Python\Python310\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([2, 2, 2, ..., 1, 1, 1])
```

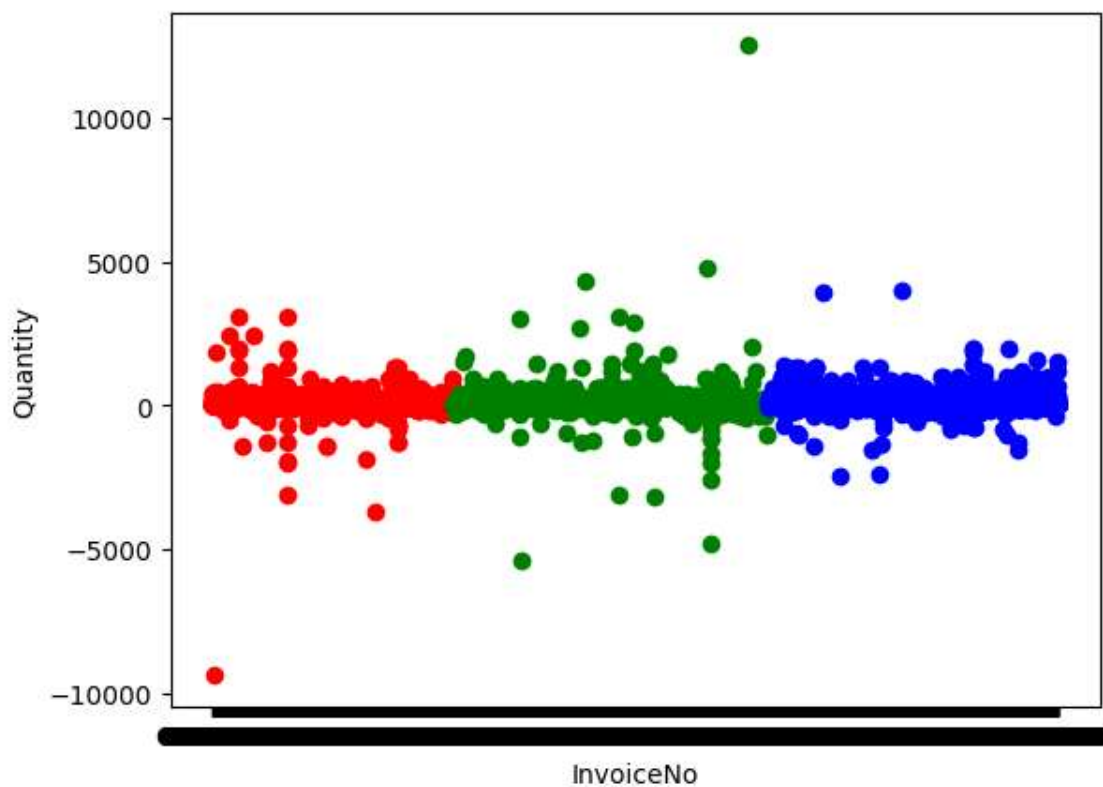
```
In [15]: df["cluster"]=y_predicted
df.head()
```

Out[15]:

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

Out[16]: Text(0, 0.5, 'Quantity')



```
In [17]: from sklearn.preprocessing import MinMaxScaler
```

```
In [18]: scaler=MinMaxScaler()
scaler.fit(df[["Quantity"]])
df["Quantity"]=scaler.transform(df[["Quantity"]])
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 [23]: km=KMeans()
```

```
In [24]: y_predicted=km.fit_predict(df[["CustomerID","Quantity"]])
y_predicted
```

C:\Users\91756\AppData\Local\Programs\Python\Python310\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[24]: array([5, 5, 5, ..., 7, 7, 7])



```
In [25]: df["New Cluster"]=y_predicted
df.head()
```

Out[25]:

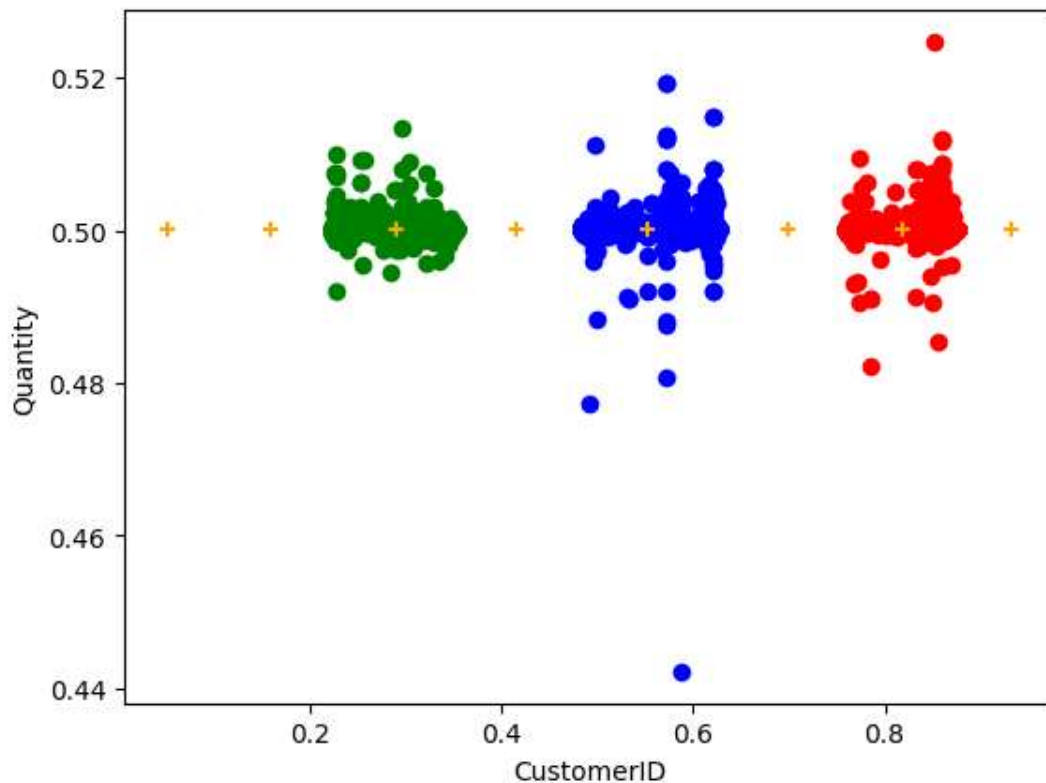
	InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country	cluster	Nº Clust
0	536365	85123A	WHITE HANGING HEART T- LIGHT HOLDER	0.500037	01-12-2010 08:26	2.55	0.926443	United Kingdom	2	
1	536365	71053	WHITE METAL LANTERN	0.500037	01-12-2010 08:26	3.39	0.926443	United Kingdom	2	
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	

```

In [27]: 1=df[df["New Cluster"]==0]
          2=df[df["New Cluster"]==1]
          3=df[df["New Cluster"]==2]
          t.scatter(df1["CustomerID"],df1["Quantity"],color="red")
          t.scatter(df2["CustomerID"],df2["Quantity"],color="green")
          t.scatter(df3["CustomerID"],df3["Quantity"],color="blue")
          t.scatter(km.cluster_centers_[0],km.cluster_centers_[1],color="orange",marker="+")
          t.xlabel("CustomerID")
          t.ylabel("Quantity")

```

Out[27]: Text(0, 0.5, 'Quantity')



```

In [28]: km.cluster_centers_

```

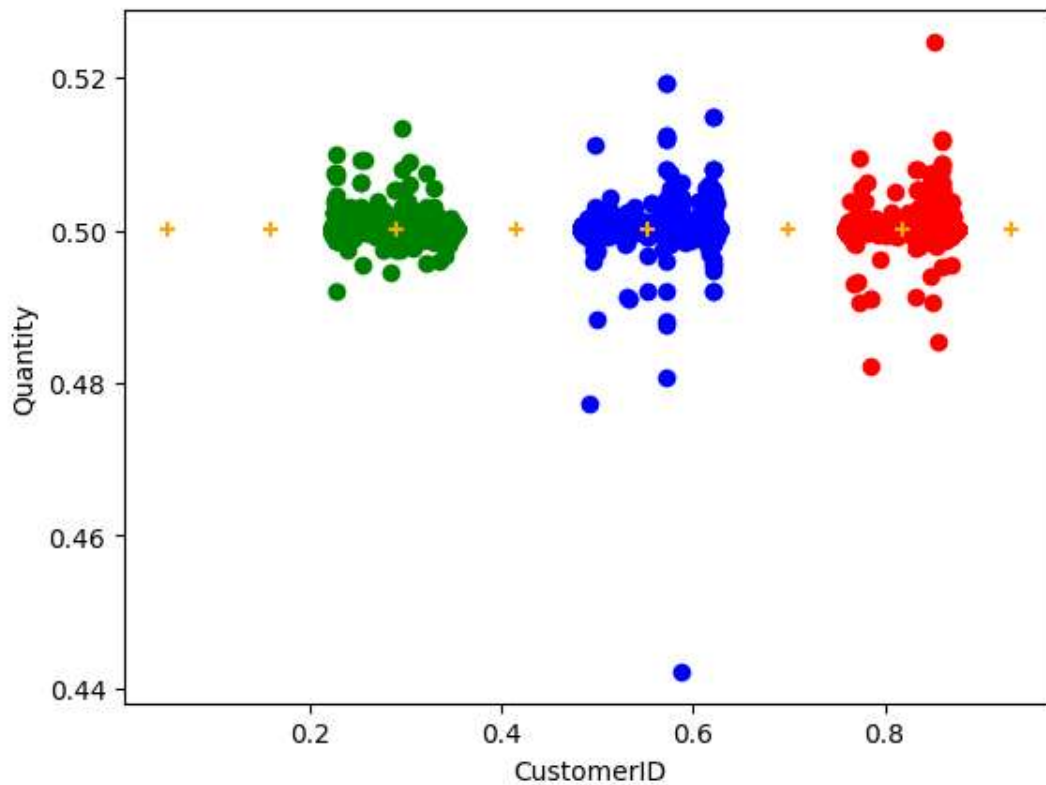
```

Out[28]: array([[0.81756576, 0.50005988],
                [0.29004173, 0.50006579],
                [0.55266146, 0.50005407],
                [0.15890109, 0.50005704],
                [0.69955075, 0.50005827],
                [0.9328779 , 0.50005088],
                [0.41480609, 0.5000595 ],
                [0.05052986, 0.50006666]])

```

```
In [30]: 1=df[df["New Cluster"]==0]
2=df[df["New Cluster"]==1]
3=df[df["New Cluster"]==2]
t.scatter(df1["CustomerID"],df1["Quantity"],color="red")
t.scatter(df2["CustomerID"],df2["Quantity"],color="green")
t.scatter(df3["CustomerID"],df3["Quantity"],color="blue")
t.scatter(km.cluster_centers_[0],km.cluster_centers_[1],color="orange",marker="+")
t.xlabel("CustomerID")
t.ylabel("Quantity")
```

Out[30]: Text(0, 0.5, 'Quantity')



```
In [31]: k_rng=range(1,10)
se=[]
```

```
In [32]: for k in k_rng:
          km=KMeans(n_clusters=k)
          km.fit(df[["CustomerID", "Quantity"]])
          se.append(km.inertia_)
print(se)
plt.plot(k_rng, se)
```

C:\Users\91756\AppData\Local\Programs\Python\Python310\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\91756\AppData\Local\Programs\Python\Python310\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\91756\AppData\Local\Programs\Python\Python310\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\91756\AppData\Local\Programs\Python\Python310\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\91756\AppData\Local\Programs\Python\Python310\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\91756\AppData\Local\Programs\Python\Python310\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\91756\AppData\Local\Programs\Python\Python310\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\91756\AppData\Local\Programs\Python\Python310\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

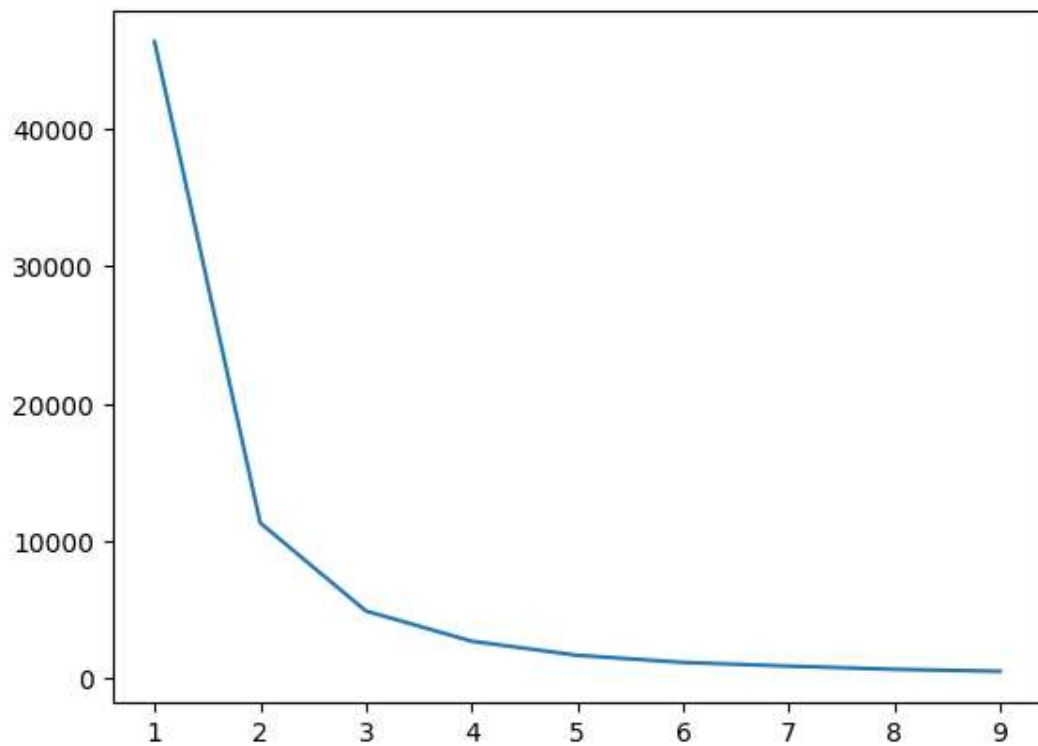
warnings.warn(

C:\Users\91756\AppData\Local\Programs\Python\Python310\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, 4920.125532402079, 2723.5191051894626, 1695.048779139392, 1178.4458741022115, 907.706544409005, 677.2512288808753, 529.6575553383113]

Out[32]: [ <matplotlib.lines.Line2D at 0x2942cbeed40>]



## CONCLUSION ¶

From the above dataset, Online Retail of the data used to take K-Mean cluster method to find the correct form of DataFrame

In [ ]: