

```
In [2]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
```

```
In [18]: df=pd.read_excel(r"C:\Users\rushg\Downloads\Customer Segmentation\Coustomer Segmentation.xlsx")
data1=df.sample(100000)
data= data1.drop(['Description'], axis=1)
data
```

Out[18]:

	InvoiceNo	StockCode	Quantity	InvoiceDate	UnitPrice	CustomerID	Country
	99518	544774	22138	3	2011-02-23 11:32:00	4.95	17711.0 United Kingdom
	452476	575337	78124	1	2011-11-09 14:11:00	1.25	17867.0 United Kingdom
	81212	543114	21535	6	2011-02-03 13:26:00	2.55	14156.0 EIRE
	191689	553385	10125	1	2011-05-16 15:53:00	0.85	16710.0 United Kingdom
	484170	577523	23660	12	2011-11-20 13:33:00	1.65	12597.0 Spain
	...	...	...	...	...	...	...
	441690	574624	23531	6	2011-11-06 11:20:00	6.95	17769.0 United Kingdom
	181617	552493	85099B	1	2011-05-09 16:21:00	4.13	NaN United Kingdom
	411656	572227	23349	12	2011-10-21 14:01:00	1.25	15051.0 United Kingdom
	202755	554494	72351A	3	2011-05-24 14:00:00	2.10	16327.0 United Kingdom
	26589	538518	22713	12	2010-12-12 16:14:00	0.42	14505.0 United Kingdom

100000 rows × 7 columns

```
In [19]: data.describe()
```

Out[19]:

	Quantity	InvoiceDate	UnitPrice	CustomerID
count	100000.000000	100000	100000.000000	75079.000000
mean	8.586130	2011-07-04 11:30:00.095999744	4.381572	15286.509783
min	-80995.000000	2010-12-01 08:26:00	0.000000	12347.000000
25%	1.000000	2011-03-28 11:56:15	1.250000	13956.000000
50%	3.000000	2011-07-20 12:29:00	2.080000	15150.000000
75%	10.000000	2011-10-19 10:48:00	4.130000	16791.000000
max	4000.000000	2011-12-09 12:50:00	13541.330000	18287.000000
std	261.711147	NaN	62.033860	1710.279544

```
In [20]: data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Index: 100000 entries, 99518 to 26589
Data columns (total 7 columns):
#   Column          Non-Null Count  Dtype
---  -
0   InvoiceNo        100000 non-null object
1   StockCode        100000 non-null object
2   Quantity         100000 non-null int64
3   InvoiceDate       100000 non-null datetime64[ns]
4   UnitPrice        100000 non-null float64
5   CustomerID       75079 non-null  float64
6   Country          100000 non-null object
dtypes: datetime64[ns](1), float64(2), int64(1), object(3)
memory usage: 6.1+ MB
```

```
In [21]: data['CustomerID'] = data.CustomerID.fillna(data.CustomerID.median())
```

```
In [22]: data.isnull().sum()
```

Out[22]:

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

```
In [55]: data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Index: 100000 entries, 99518 to 26589
Data columns (total 7 columns):
#   Column          Non-Null Count  Dtype
---  -
0   InvoiceNo        100000 non-null   object
1   StockCode       100000 non-null   object
2   Quantity        100000 non-null   int64
3   InvoiceDate      100000 non-null   datetime64[ns]
4   UnitPrice       100000 non-null   float64
5   CustomerID      100000 non-null   float64
6   Country         100000 non-null   int64
dtypes: datetime64[ns](1), float64(2), int64(2), object(2)
memory usage: 6.1+ MB
```

In [ ]:

In [35]: `X=data[['Quantity','UnitPrice']]`

In [38]: `X`

Out[38]:

	Quantity	UnitPrice
99518	3	4.95
452476	1	1.25
81212	6	2.55
191689	1	0.85
484170	12	1.65
...	...	...
441690	6	6.95
181617	1	4.13
411656	12	1.25
202755	3	2.10
26589	12	0.42

100000 rows × 2 columns

In [39]: `from sklearn.cluster import KMeans`

In [40]: `k1=KMeans(n_clusters=8)`

In [41]: `k1.fit(X)`

C:\Users\rushg\anaconda3\Lib\site-packages\sklearn\cluster\\_kmeans.py:1412: 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  
super().\_check\_params\_vs\_input(X, default\_n\_init=10)

Out[41]:

▼ KMeans

KMeans()

In [42]: `k1.labels_`

Out[42]: `array([0, 0, 0, ..., 0, 0, 0])`

In [43]: `from sklearn.metrics import silhouette_score`

In [54]: `silhouette_score(X,k1.labels_)`

Out[54]: `0.9485064877787178`

In [44]:

```
wcss=[]
for i in range(1,50):
    k2=KMeans(n_clusters=i)
    k2.fit(X)
    wcss.append(k2.inertia_)
```

C:\Users\rushg\anaconda3\Lib\site-packages\sklearn\cluster\\_kmeans.py:1412: 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  
super().\_check\_params\_vs\_input(X, default\_n\_init=10)  
C:\Users\rushg\anaconda3\Lib\site-packages\sklearn\cluster\\_kmeans.py:1412: 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  
super().\_check\_params\_vs\_input(X, default\_n\_init=10)  
C:\Users\rushg\anaconda3\Lib\site-packages\sklearn\cluster\\_kmeans.py:1412: 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  
super().\_check\_params\_vs\_input(X, default\_n\_init=10)  
C:\Users\rushg\anaconda3\Lib\site-packages\sklearn\cluster\\_kmeans.py:1412: 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

[illegible]



```
Out[47]: KMeans
KMeans(n_clusters=5)
```

```
In [48]: k2.labels_
```

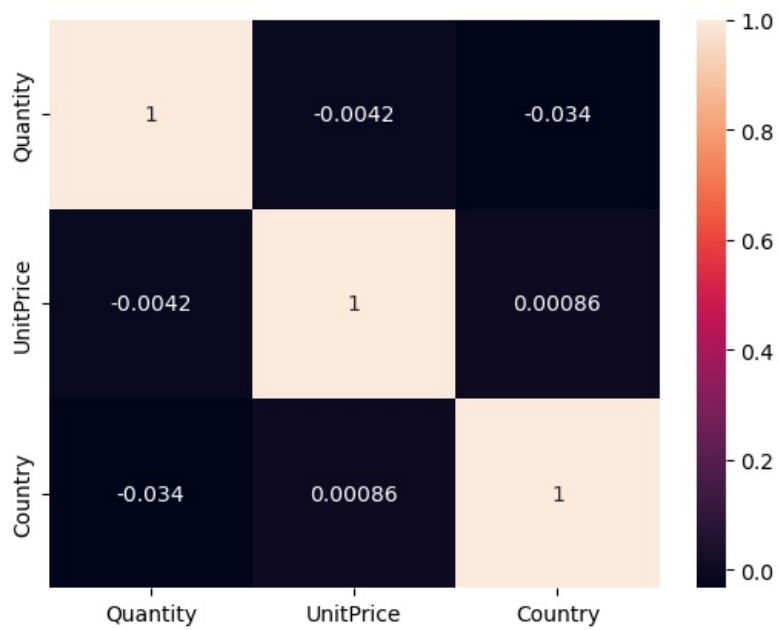
```
Out[48]: array([0, 0, 0, ..., 0, 0, 0])
```

```
In [49]: silhouette_score(X,k2.labels_)
```

```
Out[49]: 0.9957121310742266
```

```
In [49]: #plt.figure(figsize=(15,10))
sns.heatmap(X.corr(),annot=True)
```

```
Out[49]: <Axes: >
```



```
In [ ]:
```

```
In [ ]:
```