

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
```

```
df=pd.read_excel('uk_retail_file.xlsx')
```

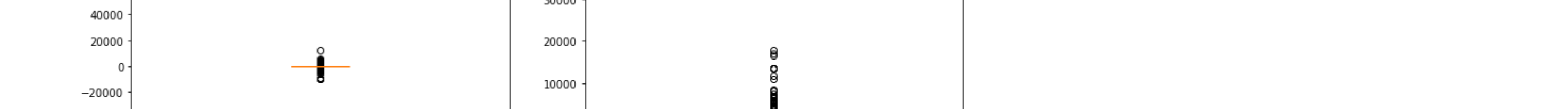
```
df.head()
```

	InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country
0	536365	85120A	WHITE HANGING HEART T-LIGHT HOLDER	6	2010-01-12 08:26:00	2.55	17850.0	United Kingdom
1	536365	71063	WHITE METAL LANTERN	6	2010-01-12 08:26:00	2.10	17908.0	United Kingdom
2	536365	84406B	CREAM CUPID HEARTS COAT HANGER	8	2010-01-12 08:26:00	2.75	17850.0	United Kingdom
3	536365	850205	KNITTED UNION JACK HOOT WATER BOTTLE	6	2010-01-12 08:26:00	2.10	17908.0	United Kingdom
4	536365	129428E	RED WOOLLY HOTTIE WHITE HEART	6	2010-01-12 08:26:00	3.30	17850.0	United Kingdom

1. Perform Basic EDA

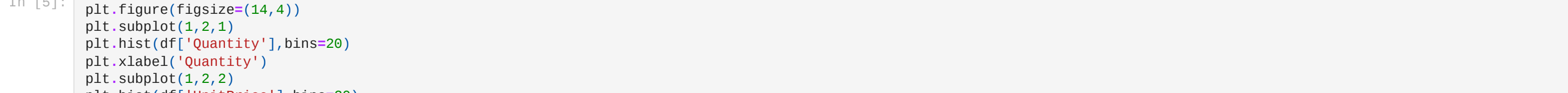
```
a. Boxplot – All Numeric Variables
```

```
plt.figure(figsize=(14,4))
plt.subplot(2,2,1)
plt.hist(df['Quantity'])
plt.xlabel('Quantity')
plt.subplot(2,2,2)
plt.hist(df['UnitPrice'])
plt.xlabel('UnitPrice')
```



b. Histogram – All Numeric Variables

```
plt.figure(figsize=(14,4))
plt.subplot(2,2,1)
plt.hist(df['Quantity'],bins=20)
plt.xlabel('Quantity')
plt.subplot(2,2,2)
plt.hist(df['UnitPrice'],bins=20)
plt.xlabel('UnitPrice')
```



c. Distribution Plot – All Numeric Variables

```
plt.figure(figsize=(14,4))
plt.subplot(2,2,1)
plt.distplot(df['Quantity'])
plt.xlabel('Quantity')
plt.subplot(2,2,2)
plt.distplot(df['UnitPrice'])
plt.xlabel('UnitPrice')
```



d. Aggregation for all numerical Columns

```
df.describe().round(1)
```

	Quantity	UnitPrice	CustomerID
count	541009.0	541009.0	426820.0
mean	9.6	4.6	15287.7
std	218.1	96.8	1713.6
min	-80995.0	-11062.1	12346.0
25%	1.0	1.2	12953.0
50%	3.0	2.1	35152.0
75%	10.0	4.1	16791.0
max	80995.0	38970.0	16281.0

e. Unique Values across all columns

```
df.nunique()
```

InvoiceNo	25988
StockCode	4078
Description	4223
Quantity	722
InvoiceDate	23268
UnitPrice	1638
CustomerID	4372
Country	38
dtype:	int64

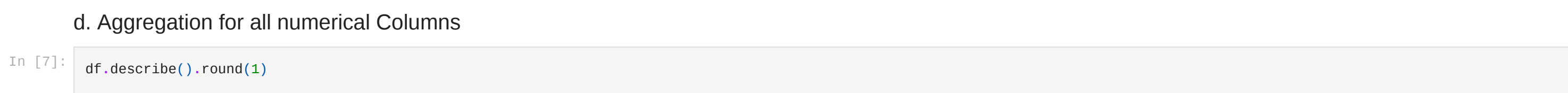
f. Duplicate values across all columns

```
df[df.duplicated()]
```

InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country
517	536409	21866	UNION JACK FLAG LUGGAGE TAG	1	2010-01-12 11:45:00	1.25	17908.0
527	536409	22866	HAND WARMER SCOTTY DOG DESIGN	1	2010-01-12 11:45:00	2.10	17908.0
537	536409	22900	SET 2 TEA TOWELS LOVE LONDON	1	2010-01-12 11:45:00	2.95	17908.0
539	536409	22211	SCOTTY DOG HOT WATER BOTTLE	1	2010-01-12 11:45:00	4.95	17908.0
543	1301421	22327	ROUND SNACK BOXES SET OF 4 DRILLS	1	2010-01-12 11:45:00	2.95	17908.0
...
541675	581538	22068	BLACK PIRATE TREASURE CHEST	1	2011-09-12 11:34:00	0.39	14466.0
541689	581538	23118	BOX OF 6 MINI VINTAGE CRACKERS	1	2011-09-12 11:34:00	2.49	14466.0
541692	581538	22992	REVOLVER WOODEN RULER	1	2011-09-12 11:34:00	1.95	14466.0
541699	581538	22094	WICKER STAR	1	2011-09-12 11:34:00	2.10	14466.0
541701	581538	23343	JUMBO BAG VINTAGE CHRISTMAS	1	2011-09-12 11:34:00	2.08	14466.0
...
526	rows < 8	columns					

g. Correlation – Heatmap - All Numeric Variables

```
sns.heatmap(df.corr())
```



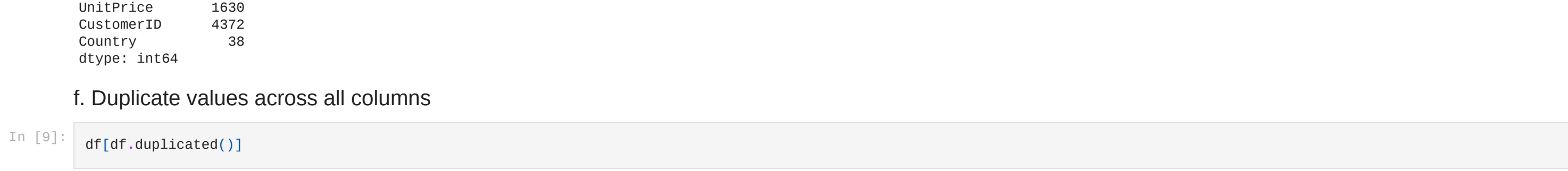
h. Regression Plot - All Numeric Variables

```
sns.regplot(df['Quantity'],df['UnitPrice'])
plt.title('Quantity vs UnitPrice')
```



i. Bar Plot – Every Categorical Variable vs every Numerical Variable

```
plt.figure(figsize=(14,8))
sns.barplot(data=df,y='Country',x='Quantity',orient='h')
plt.show()
```



```
plt.figure(figsize=(14,8))
sns.barplot(data=df,y='Country',x='UnitPrice',orient='h')
plt.show()
```



```
plt.figure(figsize=(16,9))
sns.barplot(data=df,y='Country',x='CustomerID',orient='h')
plt.show()
```



j. Pair plot - All Numeric Variables

Not Working

```
sns.pairplot(df,vars=['Quantity','UnitPrice'])
plt.show()
```

k. Line chart to show the trend of data - All Numeric/Date Variables

```
df['InvoiceDate']=pd.to_datetime(df['InvoiceDate'])
```

```
sns.lineplot(df['InvoiceDate'].dt.month,df['Quantity'])
plt.xlabel('Months')
```



```
sns.lineplot(df['InvoiceDate'].dt.month,df['UnitPrice'])
plt.xlabel('Months')
```



l. Plot the skewness - All Numeric Variables

```
df['Skewed Data'] = pd.DataFrame(df.skew(axis=1,skipna=True))
df['Skewed Data'] = pd.DataFrame(df.skew(axis=1,skipna=True))
```

```
sns.histplot(df['Skewed Data'],bins=20)
```



```
df.drop('Skewed Data',axis=1,inplace=True)
```

2. Check for missing values in all columns and replace them with the appropriate metric

```
(MeanMedianMode)
```

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

InvoiceNo	8
StockCode	8
Description	1454
Quantity	8
InvoiceDate	8
UnitPrice	8
CustomerID	135488
Country	8
dtype:	int64

```
df[df['Quantity']==0]
```

InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country
-----------	-----------	-------------	----------	-------------	-----------	------------	---------

```
df[df['UnitPrice']==0]
```

InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country
-----------	-----------	-------------	----------	-------------	-----------	------------	---------

0	Austria	1
	France	1
	Finland	1
	Bahrain	1
	Greece	1
	Hong Kong	1

Replacing the missing values of the 'Description' column with the Mode,since it is a Categorical Column:-

```
df['Description'].fillna(df['Description'].mode()[0],inplace=True)
```

Replacing the missing values of the 'CustomerID' column with the Mode,since it is a Discrete Column:-

```
df['CustomerID'].fillna(df['CustomerID'].mode()[0],inplace=True)
```

Since the 'UnitPrice' column has some 0.0 values which could not be detected with df.isnull().sum. So we need to replace such values with the Median of the column.

```
df['UnitPrice'].replace(0.0,df['UnitPrice'].median(),inplace=True)
```

3. Remove duplicate rows

```
df.drop(df[df.duplicated()].index,axis=0,inplace=True)
```

4. Remove rows which have negative values in Quantity column

```
df.drop(df[df['Quantity']<0].index,axis=0,inplace=True)
```

5. Add the columns - Month, Day and Hour for the invoice

```
df['Month']=df['InvoiceDate'].dt.month
```

```
df['Day']=df['InvoiceDate'].dt.day
```

```
df['Hour']=df['InvoiceDate'].dt.hour
```

6. How many orders made by the customers?

```
df.groupby('CustomerID')['InvoiceNo'].count().sort_values(ascending=False)
```

CustomerID	148998
17841.0	148998
14811.0	5672
14886.0	5111
12748.0	4413
14886.0	2877
...	...
15948.0	1
15823.0	1
15902.0	1
15753.0	1
12346.0	1
Name:	InvoiceNo, dtype: int64

7. TOP 5 customers with higher number of orders

```
df.groupby('CustomerID')['InvoiceNo'].count().sort_values(ascending=False).head(5)
```

CustomerID	148998
17841.0	5672
14886.0	5111
12748.0	4413
14886.0	2877
Name:	InvoiceNo, dtype: int64

8. How much money spent by the customers?

```
df['Money_spent']=df['Quantity']*df['UnitPrice']
```

```
df.groupby('CustomerID')['Money_spent'].sum()
```

CustomerID	77183.68
12346.0	4328.68
12348.0	1787.24
12349.0	1757.56
12358.0	334.40
...	...
18288.0	180.68
18281.0	80.82
18282.0	178.05
18283.0	2845.53
18287.0	1837.28
Name:	Money_spent, Length: 4339, dtype: float64

9. TOP 5 customers with highest money spent

```
df.groupby('CustomerID')['Money_spent'].sum().sort_values(ascending=False).head(5)
```

CustomerID	1895420.67
17841.0	281404.18
18182.0	259657.38
17458.0	184890.79
16446.0	168472.58
Name:	Money_spent, dtype: float64

10. How many orders per month?

```
df.groupby('Month')['InvoiceNo'].count()
```

Month	38992
1	3998
2	3749
3	3689
4	4146
5	4632
6	4632
7	4632
8	4632
9	4632
10	5528
11	1421
12	7895
Name:	InvoiceNo, dtype: int64

11. How many orders per day?

```
df.groupby('Day')['InvoiceNo'].count()
```

Day	12547
1	1022
2	1209
3	13876
4	15649
5	14589
6	14574
7	14574
8	14574
9	14574
10	22725
11	18603
12	16712
13	17948
14	17489
15	15344
16	15317
17	22381
18	18502
19	18502
20	18502
21	18502
22	15795
23	17395
24	16351
25	15797
26	16351
27	14448
28	16076
29	15603
30	15089
31	10350
Name:	InvoiceNo, dtype: int64

12. How many orders per hour?

```
df.groupby('Hour')['InvoiceNo'].count().sort_values()
```

Hour	1
7	379
29	779
19	3428
18	7696
8	6802
17	27488
9	33738
10	47679
16	130714
11	15592
14	15384
13	70067
15	75811
12	70066
Name:	InvoiceNo, dtype: int64

13. How many orders for each country?

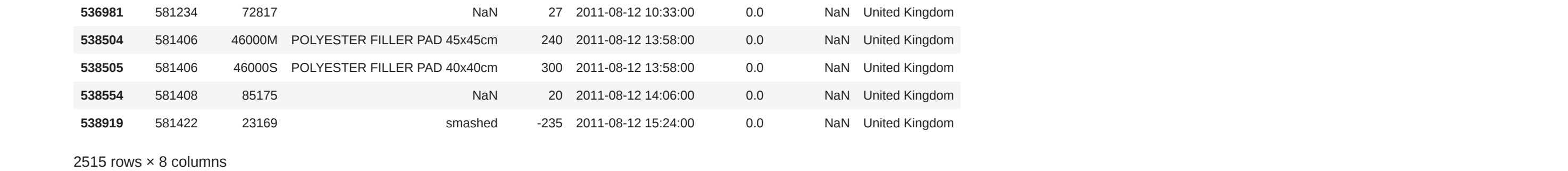
```
df.groupby('Country')['InvoiceNo'].count().sort_values(ascending=False)
```

Country	481143
United Kingdom	481143
Germany	288644.420
France	288644.420
EIRE	228862.560
Spain	209827.450
Netherlands	209827.450
Belgium	209827.450
Switzerland	1959
Portugal	1188
Australia	1188
Norway	1188
Italy	1188
Channel Islands	747
Finland	688
Cyprus	688
Sweden	688
Unspecified	442
Austria	398
Denmark	398
Poland	398
Japan	321
Hong Kong	286
Singapore	222
Iceland	187
USA	179
Canada	151
Greece	145
Mexico	112
United Arab Emirates	68
European Community	68
Lebanon	68
Lithuania	68
European Community	250
Brunei	1145.600
RSA	1884.398
Czech Republic	825.740
Bahrain	754.148
Saudi Arabia	145.200
Name:	InvoiceNo, dtype: int64

14. Orders trend across months

```
plt.plot(df.groupby('Month')['InvoiceNo'].count())
plt.xlabel('No of Orders')
plt.title('Orders trend across months')
plt.show()
```

```
<function matplotlib.pyplot.show(close=None, block=None)>
```



15. How much money spent by each country?

```
df.groupby('Country')['Money_spent'].sum().sort_values(ascending=False)
```

Country	9128175.654
Netherlands	286644.420
EIRE	288644.420
Germany	228862.560
France	209827.450
Australia	1959
Spain	6185.440
Switzerland	1959
Belgium	209827.450
Japan	37416.378
Norway	3016.600
Portugal	33881.050
Finland	2524.000
Singapore	21279.290
Channel Islands	2644.000
Denmark	1895.340
Italy	17483.240