

TBC E-commerce Dataset of Period: June-2016 to Jan-2021

Models used are: EDA - Exploratory Data Analysis Time Series Analysis using Pandas Market Basket Analysis using Apriori Algorithm

Loading Necessary Libraries

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
```

Importing the Datasets

```
In [2]: df=pd.read_csv('Ecommerce Dataset.csv')
df_new=pd.read_csv('Ecommerce Dataset.csv')
```

```
In [3]: df_new.head()
```

SKU	Product	Order number	Date placed	Quantity	Unit price	Total price	Order status	Categories	Customer email	Customer phone
0	2030301005022	The Wolf of Wall Street	2018 16/06/2016 20:43	1	350.0	332.5	Delivered	Books > Novels: Books > Biographies	oscar.nganga@gmail.com	2.547213e+11
1	2030301003312	Decision Points	2018 16/06/2016 20:43	1	1390.0	1320.5	Delivered	Books > Biographies	oscar.nganga@gmail.com	2.547213e+11
2	2030301004296	China must go	2018 16/06/2016 20:43	1	950.0	902.5	Delivered	Books > Novels	oscar.nganga@gmail.com	2.547213e+11
3	2030301001226	Why He's So Last Minute And She's Got It All W...	2027 17/06/2016 01:38	1	350.0	350.0	Cancelled	Books > Novels: Books > Motivation & Self Help	aligulaaa@gmail.com	2.547057e+11
4	2030301002838	What the Dog Saw (Penguin)	2027 17/06/2016 01:38	1	1200.0	1200.0	Cancelled	Books > Novels: Fiction: Books > Motivation &...	aligulaaa@gmail.com	2.547057e+11

```
In [4]: df=df.rename(columns={'Order number': 'Order_number'}, inplace = False)
```

i) Exploration of Data

Shape of the Dataset by Row xColumns

```
In [5]: df.shape
Out[5]: (116171, 3)
```

```
In [6]: df_new.shape
Out[6]: (116171, 12)
```

Value Count on Popularity on the Entire Dataset

Popular Categories by Value

```
In [7]: df_new['Categories'].value_counts()
Out[7]: Books > Children's Books 7345
Stationery > School Stationery 5255
Stationery > Pens, Pencils & Pouches 5117
Art Supplies > Paints & Mediums 4866
Books > Novels > Fiction 4385
...
Books > Novels: Books > Biographies: Books > Sports & Hobbies 1
Digital Books > Afrikaans 1
Text Books > Higher Education > Current Affairs & Politics: Text Books > Higher Education > Social Science 1
Books > Novels > Romance: Books > Novels > Fiction: Books > Africans 1
Text Books > Primary School > Standard 7: Text Books > Primary School > Standard 8: Text Books > Primary School > KPE Revision 1
Name: Categories, Length: 475, dtype: int64
```

Popular Customer by Orders

```
In [8]: df_new['Customer name'].value_counts()
Out[8]: Victoria Kritzell 1180
Julie Wilson 845
HASSAN SUMBA 586
Mr Wilson - Head Teacher 533
Anita 523
...
CHRISTINA BULIDA 1
NOEL OZAROL 1
Olivia Pendergast 1
Antony Munda 1
Timothy Babu 1
Name: Customer name, Length: 19605, dtype: int64
```

#Recognizing these customers may enhance more loyalty.

```
In [9]: #Extract Popular Customers by Orders as a list into Excel for Comparison.
popular_items = df_new['Customer name'].value_counts()
popular_items.to_csv('popular_items.csv', index=True, header=True)
```

Popular Products by Customer Orders

```
In [10]: df_new['Product'].value_counts()
Out[10]: Photocopy Paper A4 AONE 614
A Promised Land (Barack Obama) 338
Hp Ink Cartridge 123 Black 332
Born a crime 312
Pelikan Eraser BR40 296
...
Panya na Chura 1
KS3 Chemistry Study & Question Book - Higher 1
Improve your English Std 7 Workbook 1
Julius Caesar (EAP) 1
Champ Ball pen Red 25s 1
Name: Product, Length: 16523, dtype: int64
```

Recognizing these products may enhance inventory management and knowing products to consider for offers by Marketing

```
In [11]: #Extract Popular Products as a list into Excel for Comparison.
popular_items = df_new['Product'].value_counts()
popular_items.to_csv('popular_items.csv', index=True, header=True)
```

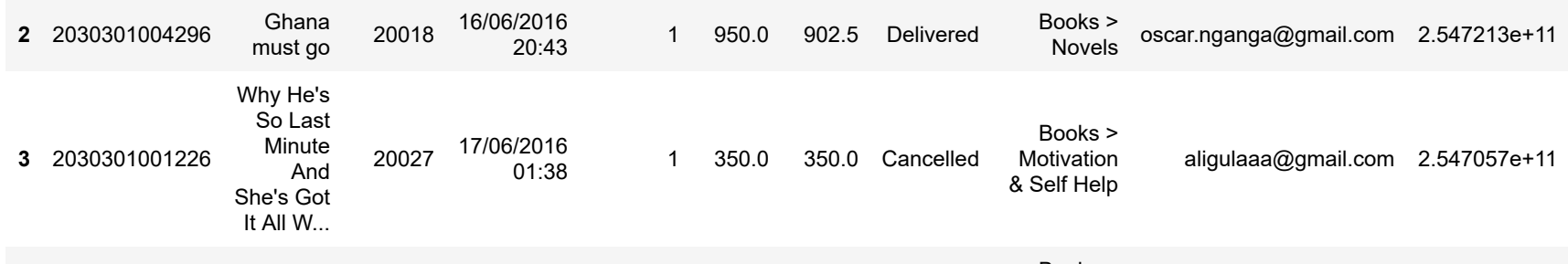
```
In [12]: df_new.columns
Out[12]: Index(['SKU', 'Product', 'Order number', 'Date placed', 'Quantity', 'Unit price', 'Total price', 'Order status', 'Categories', 'Customer email', 'Customer phone', 'Customer name'], dtype='object')
```

Plotting the Order Status by Value

```
In [13]: df_new['Order status'].value_counts()
Out[13]: Delivered 92935
Cancelled 9505
Payment failed 2571
Awaiting payment 4575
Dispatched by EMS 2043
Consolidating 265
Partially paid 182
Dispatched by rider 121
Confirmed 56
Pending 54
Dispatched by pick up location 40
Delivered at customer location 332
Dispatched 22
Name: Order status, dtype: int64
```

Visualizing the Order Status using Bar-Graph

```
In [14]: df_order = df_new['Order status'].value_counts()
df_order.plot(kind='bar', figsize=(18,12))
Out[14]: <matplotlib.axes._subplots.AxesSubplot at 0x1eb38098610>
```



```
In [15]: df1=df.groupby(['Order number','SKU'])[['SKU']].count().T
df1
Out[15]:
Order number 2030301003312 2030301004296 2030301005022 20191410101975 20303010011213 2030301001688 2030301001226 2030301002
SKU 1 1 1 1 1 1 1 1
```

1 rows x 116171 columns

```
In [16]: #Grouping of Order Number by SKU
df1=df.groupby(['Order number','SKU']).count()[['Date placed']]
df1
Out[17]:
```

Date placed	
Order_number	SKU
20018	2030301003312
2030301004296	1
2030301005022	1
20019	20191410101975
2030301001213	1
...	...
378022	2010101001101
2010101001782	1
2010117000646	1
378026	9781118805800
378034	2070701000090
2070701000090	1

116171 rows x 1 columns

```
In [18]: df2=df.groupby(['Order number','SKU']).count()
df2
Out[18]:
```

Date placed	
Order_number	SKU
20018	2030301003312
2030301004296	1
2030301005022	1
20019	20191410101975
2030301001213	1
...	...
378022	2010101001101
2010101001782	1
2010117000646	1
378026	9781118805800
378034	2070701000090
2070701000090	1

116171 rows x 1 columns

#Time Series Analysis

```
In [20]: df_new.head()
Out[20]:
```

SKU	Product	Order number	Date placed	Quantity	Unit price	Total price	Order status	Categories	Customer email	Customer phone
0	2030301005022	The Wolf of Wall Street	2018 16/06/2016 20:43	1	350.0	332.5	Delivered	Books > Novels: Books > Biographies	oscar.nganga@gmail.com	2.547213e+11
1	2030301003312	Decision Points	2018 16/06/2016 20:43	1	1390.0	1320.5	Delivered	Books > Biographies	oscar.nganga@gmail.com	2.547213e+11
2	2030301004296	China must go	2018 16/06/2016 20:43	1	950.0	902.5	Delivered	Books > Novels	oscar.nganga@gmail.com	2.547213e+11
3	2030301001226	Why He's So Last Minute And She's Got It All W...	2027 17/06/2016 01:38	1	350.0	350.0	Cancelled	Books > Novels: Books > Motivation & Self Help	aligulaaa@gmail.com	2.547057e+11
4	2030301002838	What the Dog Saw (Penguin)	2027 17/06/2016 01:38	1	1200.0	1200.0	Cancelled	Books > Novels: Fiction: Books > Motivation &...	aligulaaa@gmail.com	2.547057e+11

#Renaming the Columns of the Dataset

```
EcommerceData.rename(columns={'SKU':'SKU', 'Product':'Product', 'Order number':'OrderNumber', 'Date placed':'OrderDate', 'Quantity':'Quantity', 'Unit price':'UnitPrice', 'Total price':'TotalPrice', 'Order status':'OrderStatus', 'Categories':'Categories', 'Customer email':'CustomerEmail', 'Customer phone':'CustomerPhone', 'Customer name':'CustomerName'}, inplace=True)
```

```
In [24]: EcommerceData.head()
Out[24]:
```

SKU	Product	OrderNumber	OrderDate	Quantity	UnitPrice	TotalPrice	OrderStatus	Categories	CustomerEmail	CustomerPhone
0	2030301005022	The Wolf of Wall Street	2018 16/06/2016 20:43	1	350.0	332.5	Delivered	Books > Novels: Books > Biographies	oscar.nganga@gmail.com	2.547213e+11
1	2030301003312	Decision Points	2018 16/06/2016 20:43	1	1390.0	1320.5	Delivered	Books > Biographies	oscar.nganga@gmail.com	2.547213e+11
2	2030301004296	China must go	2018 16/06/2016 20:43	1	950.0	902.5	Delivered	Books > Novels	oscar.nganga@gmail.com	2.547213e+11
3	2030301001226	Why He's So Last Minute And She's Got It All W...	2027 17/06/2016 01:38	1	350.0	350.0	Cancelled	Books > Novels: Books > Motivation & Self Help	aligulaaa@gmail.com	2.547057e+11
4	2030301002838	What the Dog Saw (Penguin)	2027 17/06/2016 01:38	1	1200.0	1200.0	Cancelled	Books > Novels: Fiction: Books > Motivation &...	aligulaaa@gmail.com	2.547057e+11

#Convert OrderDate to Datetime

```
EcommerceData['OrderDate'] = pd.to_datetime(EcommerceData['OrderDate'])
```

EcommerceData.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 116171 entries, 0 to 116170
Data columns (total 12 columns):
#    Column      Non-Null Count  Dtype
---  -
0    SKU          116171 non-null object
1    Product      115373 non-null object
2    OrderNumber  116171 non-null int64
3    OrderDate   116171 non-null datetime64[ns]
4    Quantity     116171 non-null object
5    UnitPrice    116171 non-null float64
6    TotalPrice   116171 non-null float64
7    OrderStatus  116171 non-null object
8    Categories   115108 non-null object
9    CustomerEmail 116048 non-null object
10   CustomerPhone 114455 non-null float64
11   CustomerName 114786 non-null object
dtypes: datetime64[ns](1), float64(3), int64(2), object(6)
memory usage: 10.4+ MB
```

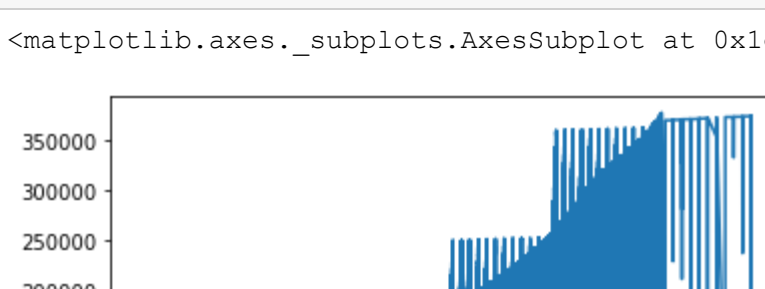
```
In [27]: EcommerceData = EcommerceData.set_index('OrderDate')
```

```
In [28]: EcommerceData.head()
```

20:43:00									Novels	
2016-06-17 01:38:00	2030301001226	Why He's So Last Minute And She's Got It All W...	20027	1	350.0	350.0	Cancelled	Books > Motivation & Self Help	aligulaaa@gmail.com	
								Books >		

```
In [29]: EcommerceData['OrderNumber'].plot(kind='line')
```

```
Out[29]: <matplotlib.axes._subplots.AxesSubplot at 0x1eb38029100>
```

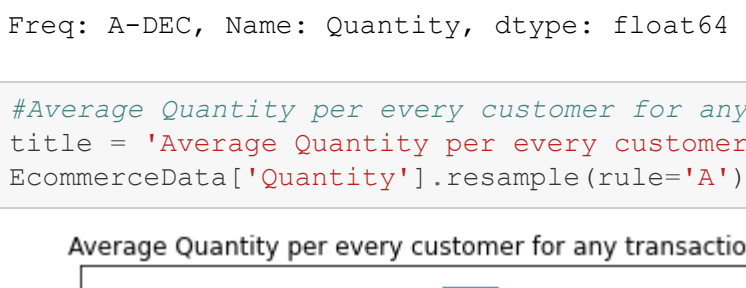


```
In [30]: EcommerceData['Quantity'].resample(rule='A').mean()
```

```
Out[30]:
OrderDate
2016-12-31    1.638517
2017-12-31    1.725082
2018-12-31    1.602872
2019-12-31    1.846176
2020-12-31    1.586022
2021-12-31    1.579149
Freq: A-DEC, Name: Quantity, dtype: float64
```

Average Quantity per every customer for any transaction

```
f = 'Average Quantity per every customer for any transaction'
EcommerceData['Quantity'].resample(rule='A').mean().plot.bar(title=f, color='f77fb4');
```



```
In [32]: EcommerceData['Quantity'].plot(figsize = (15,12))
```

```
Out[32]: <matplotlib.axes._subplots.AxesSubplot at 0x1eb3c143460>
```



```
In [33]:
```

```
In [33]:
```

```
In [33]:
```

```
In [33]:
```

```
In [33]: import matplotlib.pyplot as plt
keys = [SKU for SKU, df in df.groupby('SKU')]
plt.bar(keys,df.groupby('SKU').count()['Order number'])
```



```
In [34]: import matplotlib.pyplot as plt
keys = [SKU for SKU, df in df.groupby('SKU')]
pd.DataFrame(keys,df.groupby('SKU').count()['Order number']).T
```

Order number	1	1	2	1	2	6	1	7
0	GGGNMB7UJQM1	0WPZXVTFSPQN	1K9KCYMLCGLM	1UXIZW2NG5B2	2	2000000000	2010101000004	2010101000012

1 rows x 16861 columns

```
In [35]: #SKU_group = pd.DataFrame(keys,df.groupby('SKU').count()['Order number']).T
#SKU_group.columns = ['Order number']
```

```
In [36]: '''
load apriori and association package from mlxtend.
Used different dataset because mlxtend need data in below format.
'''
```

```
transaction = {
    'itemname': ['apple', 'banana', 'grapes'],
    'transaction': [1, 2, 3],
    'item': [0, 1, 2],
    'quantity': [1, 0, 1]
}
```

```
we could have used above data as well but need to perform operation to bring in this format instead of that used separate data only.
```

```
from mlxtend.frequent_patterns import apriori
from mlxtend.association.rules import association_rules
df1 = pd.read_csv('TBC_Ecommerce.csv', encoding='ISO-8859-1')
df1.head()
```

SKU	Product	Order number	Date placed	Quantity	Unit price	Total price	Order status	Categories	Customer email	Customer phone
0	2030301005022	The Wolf of Wall Street	2018 16/06/2016 20:43	1	350.0	332.5	Delivered	Books > Novels: Books > Biographies	oscar.nganga@gmail.com	2.547213e+11
1	2030301003312	Decision Points	2018 16/06/2016 20:43	1	1390.0	1320.5	Delivered	Books > Biographies	oscar.nganga@gmail.com	2.547213e+11
2	2030301004296	China must go	2018 16/06/2016 20:43	1	950.0	902.5	Delivered	Books > Novels	oscar.nganga@gmail.com	2.547213e+11
3	2030301001226	Why He's So Last Minute And She's Got It All W...	2027 17/06/2016 01:38	1	350.0	350.0	Cancelled	Books > Novels: Books > Motivation & Self Help	aligulaaa@gmail.com	2.547057e+11
4	2030301002838	What the Dog Saw (Penguin)	2027 17/06/2016 01:38	1	1200.0	1200.0	Cancelled	Books > Novels: Fiction: Books > Motivation &...	aligulaaa@gmail.com	2.547057e+11

```
In [37]: # data has many country choose any one for check..
df1.SKU.value_counts().head(15)
```

0	2030301005022	The Wolf of Wall Street	2018 16/06/2016 20:43	1	350.0	332.5	Delivered	Books > Novels: Books > Biographies	oscar.nganga@gmail.com	2.547213e+11
1	2030301003312	Decision Points	2018 16/06/2016 20:43	1	1390.0	1320.5	Delivered	Books > Biographies	oscar.nganga@gmail.com	2.547213e+11
2	2030301004296	China must go	2018 16/06/2016 20:43	1	950.0	902.5	Delivered	Books > Novels	oscar.nganga@gmail.com	2.547213e+11
3	2030301001226	Why He's So Last Minute And She's Got It All W...	2027 17/06/2016 01:38	1	350.0	350.0	Cancelled	Books > Novels: Books > Motivation & Self Help	aligulaaa@gmail.com	2.547057e+11
4	2030301002838	What the Dog Saw (Penguin)	2027 17/06/2016 01:38	1	1200.0	1200.0	Cancelled	Books > Novels: Fiction: Books > Motivation &...	aligulaaa@gmail.com	2.547057e+11

```
In [38]: #convert data in format which it require converting using pivot table and Quantity sum as values. fill 0 if any nan values
basket = pd.pivot_table(data=df1,index='Order number',columns='Product',values='Quantity',aggfunc='sum',fill_value=0)
```

```
In [39]: #Convert data in format which it require converting using pivot table and Quantity sum as values. fill 0 if any nan values
basket = pd.pivot_table(data=df1,index='Order number',columns='Product',values='Quantity',aggfunc='sum',fill_value=0)
```

```
In [40]: basket.head()
```

Order number	"A" Finder Biology Practical Answer booklet	"A" Finder Chemistry Practical Answer booklet	"A" Finder Physics Practical Answer booklet	"A" Finder Physics Practical Answer booklet	"Bantex Clipboard plastic 8859-01 Blue"	"Bantex Clipboard plastic 8859-01 Red"	"Bantex PVC folder 3420-04 Little Red Hen"	"Koko riko namii amemiamashe?"	ladybird Tales - The Little Red Hen	little Oxfords
20018	0	0	0	0	0	0	0	0	...	0
20019	0	0	0	0	0	0	0	0	...	0
20027	0	0	0	0	0	0	0	0	...	0
20029	0	0	0	0	0	0	0	0	...	0
20030	0	0	0	0	0	0	0	0	...	0

5 rows x 16523 columns

```
In [41]: #this to check correctness after binning it to 1 at below code..
basket['What the Dog Saw (Penguin)'].head()
```

Order number	20018	20019	20027	20029	20030
What the Dog Saw (Penguin)	0	0	0	0	0

```
In [42]: # we dont need quantity sum we need either has taken or not so if user has taken that item mark as 1 else he has not taken 0.
def convert_into_binary(x):
    if x > 0:
        return 1
    else:
        return 0
```



```
In [58]: # rules_mltend.rename(columns={'antecedents':'lhs','consequents':'rhs'})

# as based business use case we can sort based on confidence and lift.
rules_mltend[ (rules_mltend['lift'] >= 4) & (rules_mltend['confidence'] >= 0.5) ]
```

Out [58]:

	antecedents	consequents	antecedent support	consequent support	support	confidence	lift	leverage	conviction
0	(Hp Ink Cartridge 650 Black)	(Hp Ink Cartridge 650 Colour)	0.007410	0.005663	0.004910	0.662602	117.005579	0.004868	2.947071
1	(Hp Ink Cartridge 650 Colour)	(Hp Ink Cartridge 650 Black)	0.005663	0.007410	0.004910	0.867021	117.005579	0.004868	7.464276
2	(Hp Ink Cartridge 652 Black)	(Hp Ink Cartridge 652 Colour)	0.005633	0.004850	0.004127	0.732620	151.065400	0.004099	3.721862
3	(Hp Ink Cartridge 652 Colour)	(Hp Ink Cartridge 652 Black)	0.004850	0.005633	0.004127	0.850932	151.065400	0.004099	6.670546
4	(Hp Ink Cartridge 123 Colour)	(Hp Ink Cartridge 123 Black)	0.005814	0.009699	0.004880	0.839378	86.539375	0.004823	6.165420
5	(Hp Ink Cartridge 123 Black)	(Hp Ink Cartridge 123 Colour)	0.009699	0.005814	0.004880	0.503106	86.539375	0.004823	2.000860

```
In [59]: rules_mltend.shape
```

```
Out [59]: (6, 9)
```

```
In [ ] :
```