

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

In [2]: df=pd.read_csv('amazon.csv')

In [3]: df.head()

Out[3]:
   index      Order ID      Date      Status      Fulfillment      Sales Channel      ship-service-level      Category      Size      Courier Status      ..      currency      Amount      ship-city      ship-state      ship-postal-code      ship-country      B2B      fulfilled-by      New      PendingS
0      0      405-90784-673146      04-30-22      Cancelled      Merchant      Amazon.in      Standard      T-shirt      S      On the Way      ..      INR      647.82      MUMBAI      MAHARASHTRA      400050.0      IN      False      Easy Ship      NaN      NaN
1      1      171-9198151-101146      04-30-22      Shipped      Delivered to Buyer      Merchant      Amazon.in      Standard      Shirt      XL      Shipped      ..      INR      406.00      BENGALURU      KARNATAKA      560095.0      IN      False      Easy Ship      NaN      NaN
2      2      404-0608978-727146      04-30-22      Shipped      Amazon      Amazon.in      Expedited      Shirt      XL      Shipped      ..      INR      329.00      NAWALUMBAI      MAHARASHTRA      410210.0      IN      True      NaN      NaN      NaN
3      3      403-9516377-813361      04-30-22      Cancelled      Merchant      Amazon.in      Standard      Blazer      L      On the Way      ..      INR      753.33      PUDUCHERRY      PUDUCHERRY      605000.0      IN      False      Easy Ship      NaN      NaN
4      4      407-1068790-740320      04-30-22      Shipped      Amazon      Amazon.in      Expedited      Trousers      3XL      Shipped      ..      INR      574.00      CHENNAI      TAMIL NADU      600073.0      IN      False      NaN      NaN      NaN

5 rows * 21 columns

In [4]: df.tail()

Out[4]:
   index      Order ID      Date      Status      Fulfillment      Sales Channel      ship-service-level      Category      Size      Courier Status      ..      currency      Amount      ship-city      ship-state      ship-postal-code      ship-country      B2B      fulfilled-by      New      PendingS
128971      128971      406-001380-707107      05-31-22      Shipped      Amazon      Amazon.in      Expedited      T-shirt      M      Shipped      ..      INR      517.0      HYDERABAD      TELANGANA      500013.0      IN      False      NaN      NaN      NaN
128972      128972      407-954769-315288      05-31-22      Shipped      Amazon      Amazon.in      Expedited      Blazer      XXL      Shipped      ..      INR      690.0      HYDERABAD      TELANGANA      500040.0      IN      False      NaN      NaN      NaN
128973      128973      402-238160-650956      05-31-22      Shipped      Amazon      Amazon.in      Expedited      T-shirt      XS      Shipped      ..      INR      159.0      HUBLI      GOA      580050.0      IN      False      NaN      NaN      NaN
128974      128974      409-743056-073512      05-31-22      Shipped      Amazon      Amazon.in      Expedited      T-shirt      S      Shipped      ..      INR      696.0      RAIPUR      CHHATTISGARH      492014.0      IN      False      NaN      NaN      NaN

5 rows * 21 columns

In [5]: df.shape

Out[5]:
(128976, 21)

In [6]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 128976 entries, 0 to 128975
Data columns (total 21 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   index      128976 non-null    int64
 1   Order ID   128976 non-null    object
 2   Date       128976 non-null    object
 3   Status     128976 non-null    object
 4   Fulfillment 128976 non-null    object
 5   Sales Channel 128976 non-null    object
 6   ship-service-level 128976 non-null    object
 7   Category    128976 non-null    object
 8   Size        128976 non-null    object
 9   Courier Status 128976 non-null    object
10   Qty         128976 non-null    int64
11   currency    121176 non-null    object
12   Amount      121176 non-null    float64
13   ship-city   128941 non-null    object
14   ship-state  128941 non-null    object
15   ship-postal-code 128941 non-null    float64
16   ship-country 128941 non-null    float64
17   B2B         128976 non-null    bool
18   fulfilled-by 39263 non-null    float64
19   New         6 non-null         float64
20   PendingS    6 non-null         float64
dtypes: bool(1), float64(1), int64(2), object(14)
memory usage: 37.4+ MB

In [7]: #drop blanks columns
df.drop(['New', 'PendingS'],axis=1,inplace=True)

In [8]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 128976 entries, 0 to 128975
Data columns (total 19 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   index      128976 non-null    int64
 1   Order ID   128976 non-null    object
 2   Date       128976 non-null    object
 3   Status     128976 non-null    object
 4   Fulfillment 128976 non-null    object
 5   Sales Channel 128976 non-null    object
 6   ship-service-level 128976 non-null    object
 7   Category    128976 non-null    object
 8   Size        128976 non-null    object
 9   Courier Status 128976 non-null    object
10   Qty         128976 non-null    int64
11   currency    121176 non-null    object
12   Amount      121176 non-null    float64
13   ship-city   128941 non-null    object
14   ship-state  128941 non-null    object
15   ship-postal-code 128941 non-null    float64
16   ship-country 128941 non-null    float64
17   B2B         128976 non-null    bool
18   fulfilled-by 39263 non-null    object
dtypes: bool(1), float64(2), int64(2), object(14)
memory usage: 37.4+ MB

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

Out[9]:
index      0
Order ID    0
Date        0
Status      0
Fulfillment 0
Sales Channel 0
ship-service-level 0
Category     0
Size         0
Courier Status 0
Qty          0
currency     0
Amount       0
ship-city    0
ship-state   0
ship-postal-code 0
ship-country 0
B2B          0
fulfilled-by 89713
dtype: int64

In [10]: sns.heatmap(df.isnull())
plt.show()

In [11]: df.drop(['fulfilled-by'],axis=1,inplace=True)
df.dropna(inplace=True)

In [12]: df.isnull().sum()

Out[12]:
index      0
Order ID    0
Date        0
Status      0
Fulfillment 0
Sales Channel 0
ship-service-level 0
Category     0
Size         0
Courier Status 0
Qty          0
currency     0
Amount       0
ship-city    0
ship-state   0
ship-postal-code 0
ship-country 0
B2B          0
dtype: int64

In [13]: df.shape

Out[13]:
(121143, 18)

In [14]: df['ship-postal-code'].df['ship-postal-code'].astype(int)

In [15]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 121143 entries, 0 to 128975
Data columns (total 18 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   index      121143 non-null    int64
 1   Order ID   121143 non-null    object
 2   Date       121143 non-null    object
 3   Status     121143 non-null    object
 4   Fulfillment 121143 non-null    object
 5   Sales Channel 121143 non-null    object
 6   ship-service-level 121143 non-null    object
 7   Category    121143 non-null    object
 8   Size        121143 non-null    object
 9   Courier Status 121143 non-null    object
10   Qty         121143 non-null    int64
11   currency    121143 non-null    object
12   Amount      121143 non-null    float64
13   ship-city   121143 non-null    object
14   ship-state  121143 non-null    object
15   ship-postal-code 121143 non-null    int32
16   ship-country 121143 non-null    object
17   B2B         121143 non-null    bool
dtypes: bool(1), float64(1), int32(1), int64(2), object(13)
memory usage: 16.3+ MB

In [16]: df['date']>pd.to_datetime(df['date'])

In [17]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 121143 entries, 0 to 128975
Data columns (total 18 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   index      121143 non-null    int64
 1   Order ID   121143 non-null    object
 2   Date       121143 non-null    datetime64[ns]
 3   Status     121143 non-null    object
 4   Fulfillment 121143 non-null    object
 5   Sales Channel 121143 non-null    object
 6   ship-service-level 121143 non-null    object
 7   Category    121143 non-null    object
 8   Size        121143 non-null    object
 9   Courier Status 121143 non-null    object
10   Qty         121143 non-null    int64
11   currency    121143 non-null    object
12   Amount      121143 non-null    float64
13   ship-city   121143 non-null    object
14   ship-state  121143 non-null    object
15   ship-postal-code 121143 non-null    int32
16   ship-country 121143 non-null    bool
dtypes: bool(1), datetime64[ns](1), float64(1), int32(1), int64(2), object(12)
memory usage: 16.3+ MB

In [18]: df.head()

   index      Order ID      Date      Status      Fulfillment      Sales Channel      ship-service-level      Category      Size      Courier Status      Qty      currency      Amount      ship-city      ship-state      ship-postal-code      ship-country      B2B
0      0      405-90784-673146      2022-04-30      Cancelled      Merchant      Amazon.in      Standard      T-shirt      S      On the Way      0      INR      647.82      MUMBAI      MAHARASHTRA      400050.0      IN      False
1      1      171-9198151-101146      2022-04-30      Shipped      Delivered to Buyer      Merchant      Amazon.in      Standard      Shirt      XL      Shipped      1      INR      406.00      BENGALURU      KARNATAKA      560095.0      IN      True
2      2      404-0608978-727146      2022-04-30      Shipped      Amazon      Amazon.in      Expedited      Shirt      XL      Shipped      1      INR      329.00      NAWALUMBAI      MAHARASHTRA      410210.0      IN      True
3      3      403-9516377-813361      2022-04-30      Cancelled      Merchant      Amazon.in      Standard      Blazer      L      On the Way      0      INR      753.33      PUDUCHERRY      PUDUCHERRY      605000.0      IN      False
4      4      407-1068790-740320      2022-04-30      Shipped      Amazon      Amazon.in      Expedited      Trousers      3XL      Shipped      1      INR      574.00      CHENNAI      TAMIL NADU      600073.0      IN      False

In [19]: df.rename(columns={'Qty':'Quantity'},inplace=True)

In [20]: df.describe()

   index      Quantity      Amount      ship-postal-code
count  121143.000000      121143.000000      121143.000000      121143.000000
mean      6486.321956      0.961262      648.578974      400020.740007
std      37320.415404      0.214278      281.196986      191301.588170
min           0.000000      0.000000      0.000000      110001.000000
25%      32284.500000      1.000000      448.000000      382421.000000
50%      64477.000000      1.000000      605.000000      500032.000000
75%      96602.500000      1.000000      788.000000      600020.000000
max      128974.000000      9.000000      5941.000000      999999.000000

In [20]: df.describe(include=object)

   Order ID      Status      Fulfillment      Sales Channel      ship-service-level      Category      Size      Courier Status      currency      ship-city      ship-state      ship-country
count      121143      121143      121143      121143      121143      121143      121143      121143      121143      121143      121143
unique      128941      12      2      1      2      9      11      3      1      8997      68      1
top      171-5077375-201050      Standard      Amazon      Amazon.in      Expedited      T-shirt      M      BENGALURU      MAHARASHTRA      IN
freq      12      7368      83629      121143      82713      47038      20965      19498      12679      21084      121143

In [20]: sns.countplot(x='Size',data=df)

In [20]: ax=plt.xlabel('ax.containers(8)')
plt.show()

Note:From above graph you can see that most people buys M-size

In [42]: df1=df.groupby(['Size'],as_index=False)['Quantity'].sum()
df1

Out[42]:
   Size      Quantity
6   M      20513
9   L      19567
8   XL      19511
10  XXL      16217
7   S      15559
0   3XL      13341
9   XS      9635
4   Free      2061
0XL      887
5XL      512
1   4XL      396

In [44]: sns.barplot(x='Size',y='Quantity',data=df1)

Out[44]:
<AxesSubplot: xlabel='Size', ylabel='Quantity'>

most of quantity buys m-size in the sales

In [23]: sns.countplot(x='Courier Status',data=df,hue='Status')

<AxesSubplot: xlabel='Courier Status', ylabel='count'>

most of couriers are T-shirt and shirt

In [24]: plt.figure(figsize=(5,5))
sns.countplot(df['Size'])
plt.show()

<Figure size 500x500 with 8 Axes>

In [25]: df['Category'].df['Category'].astype(str)

In [26]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 121143 entries, 0 to 128975
Data columns (total 18 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   index      121143 non-null    int64
 1   Order ID   121143 non-null    object
 2   Date       121143 non-null    datetime64[ns]
 3   Status     121143 non-null    object
 4   Fulfillment 121143 non-null    object
 5   Sales Channel 121143 non-null    object
 6   ship-service-level 121143 non-null    object
 7   Category    121143 non-null    object
 8   Size        121143 non-null    object
 9   Courier Status 121143 non-null    object
10   Qty         121143 non-null    int64
11   currency    121143 non-null    object
12   Amount      121143 non-null    float64
13   ship-city   121143 non-null    object
14   ship-state  121143 non-null    object
15   ship-postal-code 121143 non-null    int32
16   ship-country 121143 non-null    bool
dtypes: bool(1), datetime64[ns](1), float64(1), int32(1), int64(2), object(12)
memory usage: 16.3+ MB

In [47]: plt.hist(df['Category'],bins=20,color='olive',edgecolor='black')
plt.xticks(rotation=45)
plt.show()

most of the buyers are T-shirt and shirt

In [28]: p=df['B2B'].value_counts(normalize=True)
p

Out[28]:
False      0.993033
True       0.006967
Name: B2B, dtype: float64

In [30]: plt.pie(p,labelindex.autoscale('%0.2f%%'))
plt.show()

False      99.30%
True       0.70%

In [31]: plt.figure(figsize=(8,5))
sns.countplot(x='Category',y='Size',data=df,color='brown')
plt.show()

In [32]: plt.figure(figsize=(10,5))
sns.countplot(x='ship-state',data=df)
plt.xticks(rotation=90)
plt.show()

In [33]: h=df['ship-state'].value_counts()[0:5]
h

Out[33]:
MAHARASHTRA      23884
KARNATAKA         49398
TAMIL NADU        10813
TELANGANA         36008
UTTAR PRADESH     9956
Name: ship-state, dtype: int64

In [34]: v=list(df['ship-state'].value_counts()[0:5].keys())
v

Out[34]:
['MAHARASHTRA', 'KARNATAKA', 'TAMIL NADU', 'TELANGANA', 'UTTAR PRADESH']

In [35]: plt.figure(figsize=(10,5))
plt.bar(v,h,color='g',edgecolor='black')
plt.show()

most of buyers are maharashtra state

conclusion

The data analysis reveals that the business has a significant customer base in maharashtra state ,mainly serves retailers, fulfills orders through amazon,experiences high demand for T-shirts and see M-size as the prefers most of buyers.

In [ ]:
```