

DataSet: Amazon_Sales_Analysis[kaggle].

Discription: Analyze Amazon sales data to check the buyers preferred choice in the sales

Importing Libraries

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

Double-click (or enter) to edit

Load Dataset

```
df=pd.read_csv('/content/Amazon Sale Report.csv')
df.head(3)
```

↗

	index	Order ID	Date	Status	Fulfilment	Sales Channel	ship-service-level	Category	Size	Courier Status	..
0	0	405-8078784-5731545	04-30-22	Cancelled	Merchant	Amazon.in	Standard	T-shirt	S	On the Way	
1	1	171-9198151-1101146	04-30-22	Shipped - Delivered to Buyer	Merchant	Amazon.in	Standard	Shirt	3XL	Shipped	
2	2	404-0687676-7273146	04-30-22	Shipped	Amazon	Amazon.in	Expedited	Shirt	XL	Shipped	

3 rows × 21 columns

```
df.shape
```

↗

```
(128976, 21)
```

```
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   Fulfilment         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  object
17  B2B                128976 non-null  bool
18  fulfilled-by       39263 non-null  object
19  New                0 non-null      float64
20  PendingS           0 non-null      float64
dtypes: bool(1), float64(4), int64(2), object(14)
memory usage: 19.8+ MB
```

row 19&20 are outliers

Cleaning And Analysis

```
#drop unrelated/blank columns
df.drop(['New','PendingS'], axis=1, inplace=True)
```

```
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   Fulfilment            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 object
17   B2B                   128976 non-null bool
18   fulfilled-by          39263 non-null object
dtypes: bool(1), float64(2), int64(2), object(14)
memory usage: 17.8+ MB
```

Checking null value

```
pd.isnull(df)
```

```
Show hidden output
```

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

```
index                0
Order ID              0
Date                  0
Status                0
Fulfilment            0
Sales Channel         0
ship-service-level    0
Category              0
Size                  0
Courier Status        0
Qty                   0
currency              7800
Amount                7800
ship-city             35
ship-state            35
ship-postal-code      35
ship-country          35
B2B                   0
fulfilled-by          89713
dtype: int64
```

null values detect

- 1. currency
- 2. Amount
- 3. fulfilled_by

```
df.dropna(inplace=True)
```

```
df.columns
```

```
Index(['index', 'Order ID', 'Date', 'Status', 'Fulfilment', 'Sales Channel',
      'ship-service-level', 'Category', 'Size', 'Courier Status', 'Qty',
      'currency', 'Amount', 'ship-city', 'ship-state', 'ship-postal-code',
      'ship-country', 'B2B', 'fulfilled-by'],
      dtype='object')
```

```
# change data type
df['ship-postal-code']=df['ship-postal-code'].astype('int')
```

```
#use describe() for specific columns
df[['Qty','Amount']].describe()
```



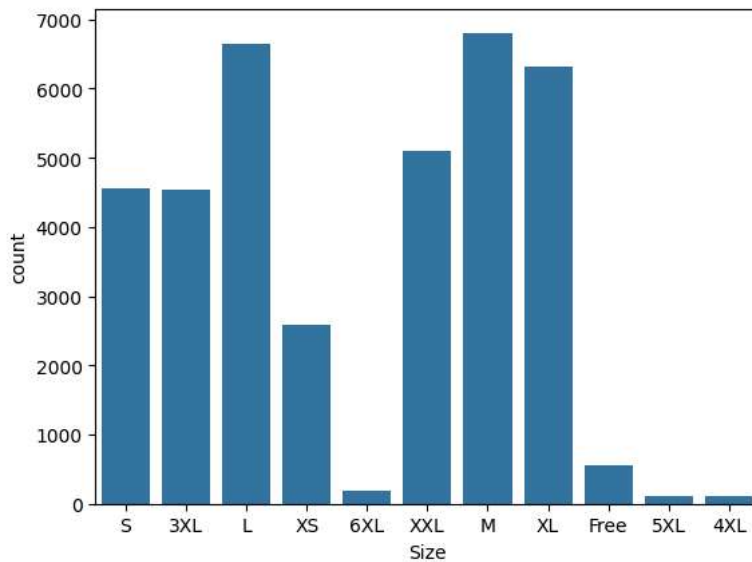
	Qty	Amount
count	37514.000000	37514.000000
mean	0.867383	646.553960
std	0.354160	279.952414
min	0.000000	0.000000
25%	1.000000	458.000000
50%	1.000000	629.000000
75%	1.000000	771.000000
max	5.000000	5495.000000



Exploratory Data Analysis & Visualization

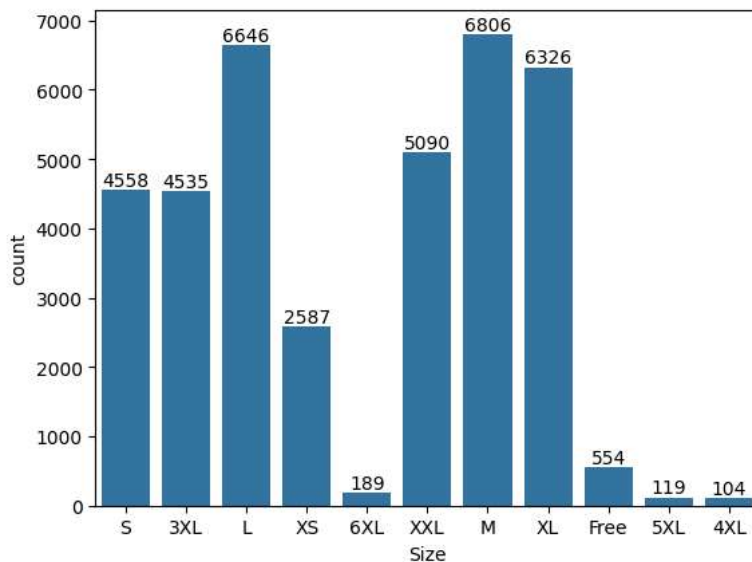
Size

```
ax=sns.countplot(x='Size' ,data=df)
```



```
ax=sns.countplot(x='Size' ,data=df)
```

```
for bars in ax.containers:
    ax.bar_label(bars)
```



Note: From above Graph you can see that most of the people buys M-Size

Group By

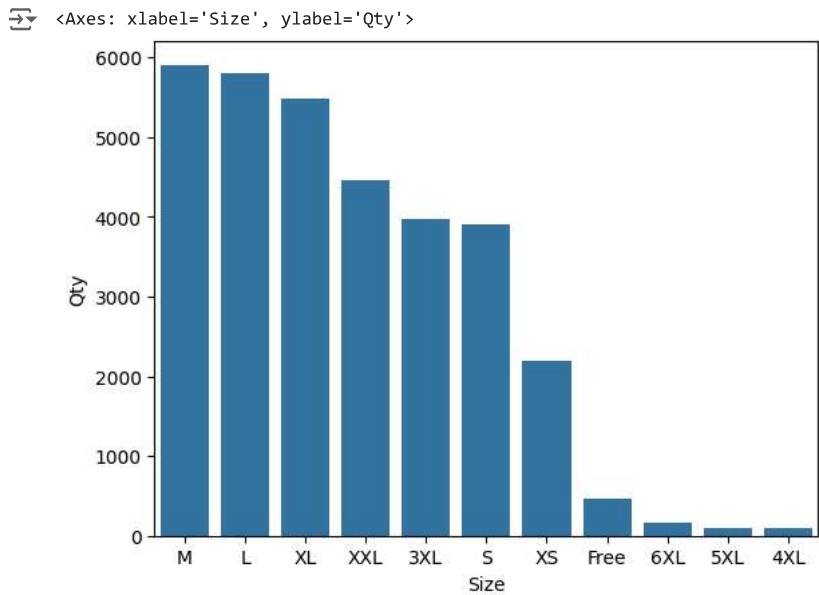
The groupby() function in pandas is used to group data based on one or more columns in a DataFrame

```
df.groupby(['Size'], as_index=False)['Qty'].sum().sort_values(by='Qty',ascending=False)
```

	Size	Qty
6	M	5905
5	L	5795
8	XL	5481
10	XXL	4465
0	3XL	3972
7	S	3896
9	XS	2191
4	Free	467
3	6XL	170
2	5XL	104
1	4XL	93

```
S_Qty=df.groupby(['Size'], as_index=False)['Qty'].sum().sort_values(by='Qty',ascending=False)
```

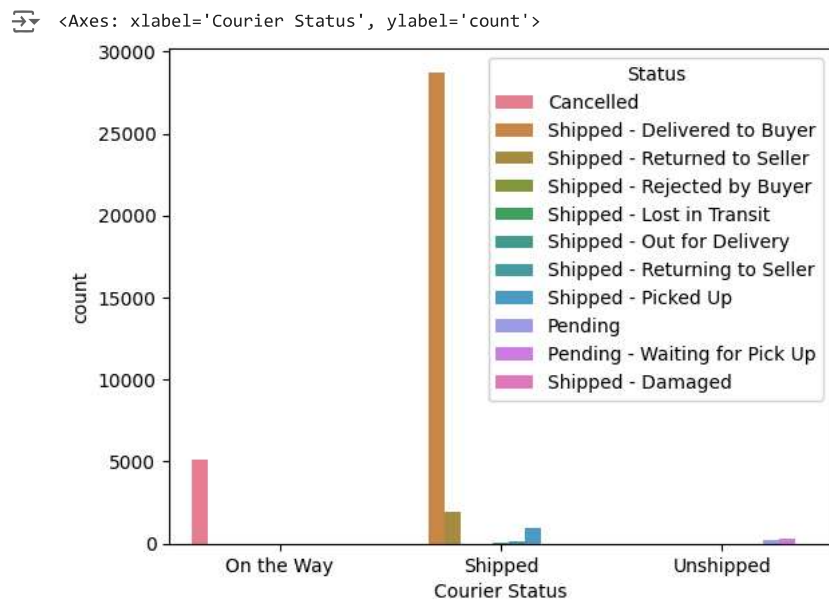
```
sns.barplot(x='Size',y='Qty', data=S_Qty)
```



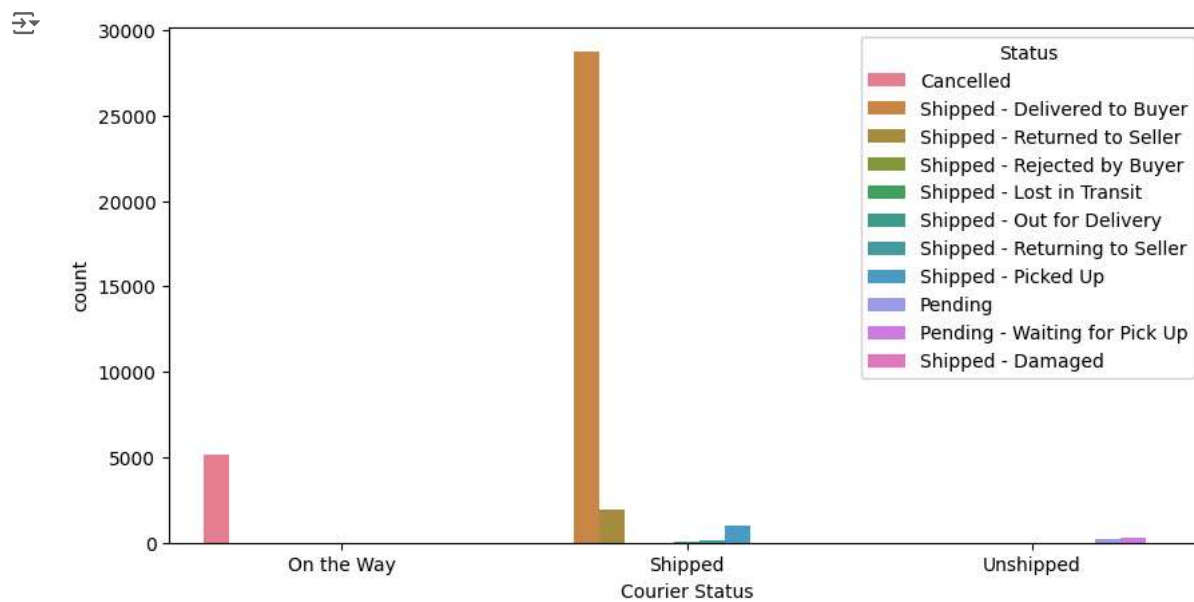
Note: From above Graph you can see that most of the Qty buys M-Size in the sales

Courier Status

```
sns.countplot(data=df, x='Courier Status',hue= 'Status')
```

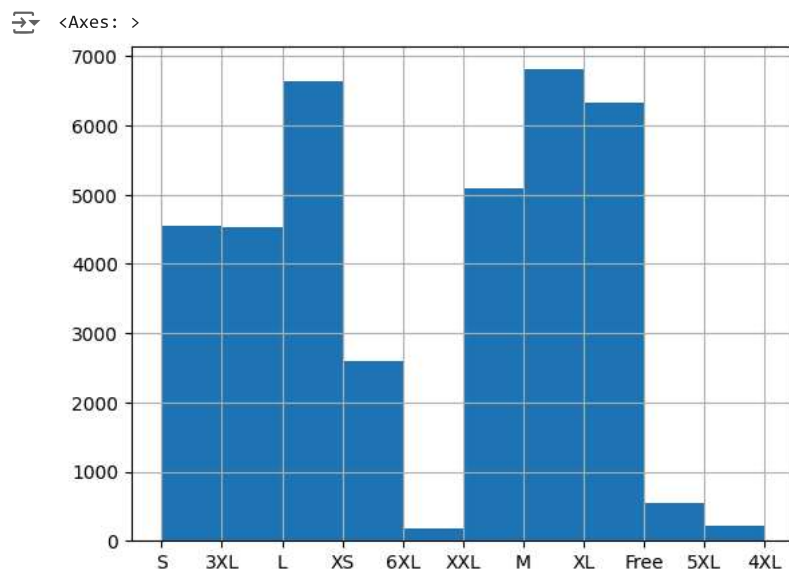


```
plt.figure(figsize=(10,5))  
ax=sns.countplot(data=df, x='Courier Status',hue= 'Status')  
plt.show()
```

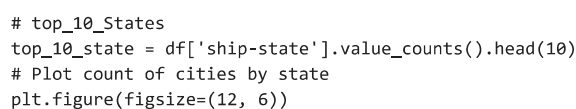
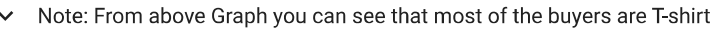


▼ Note: From above Graph the majority of the orders are shipped through the courier.

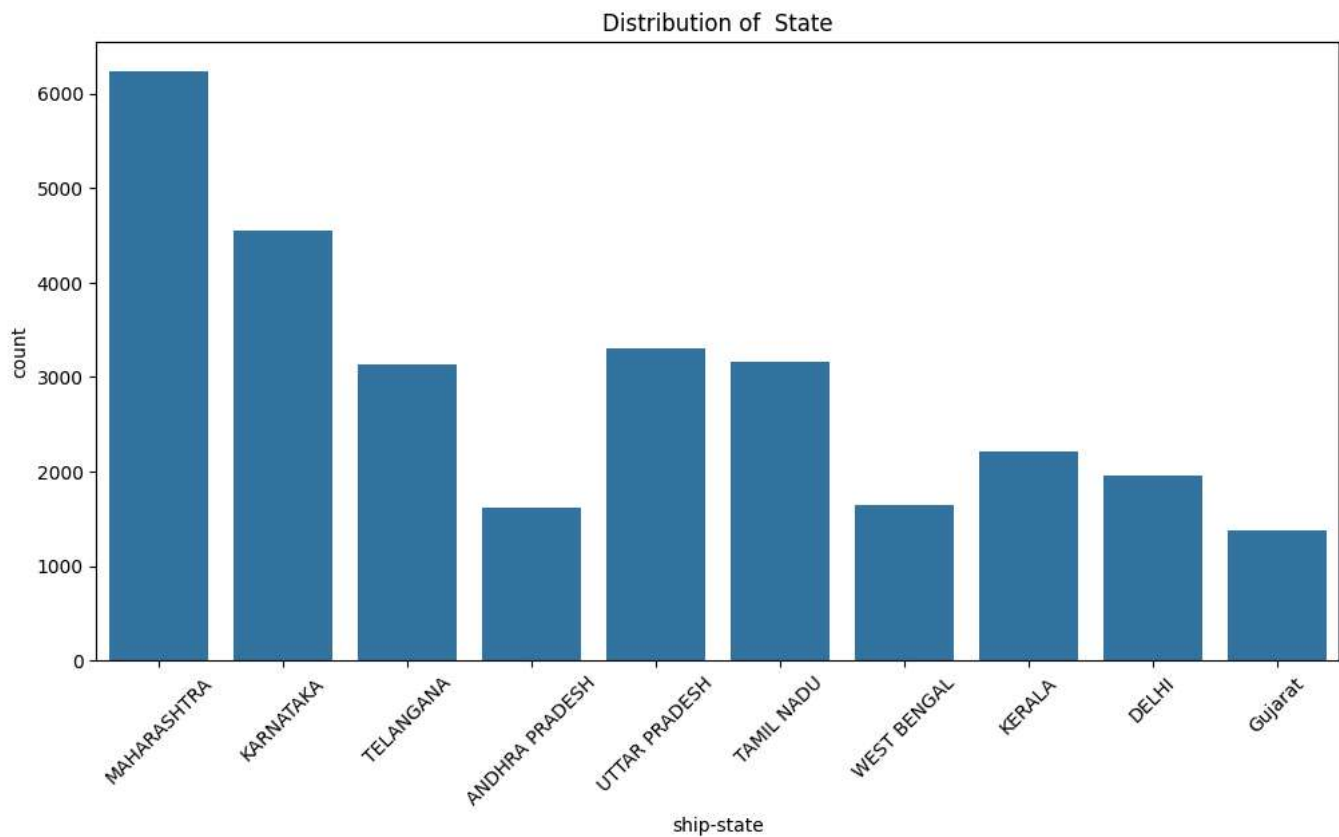
```
#histogram  
df['Size'].hist()
```



```
df['Category'] = df['Category'].astype(str)
```



```
sns.countplot(data=df[df['ship-state'].isin(top_10_state.index)], x='ship-state')
plt.xlabel('ship-state')
plt.ylabel('count')
plt.title('Distribution of State')
plt.xticks(rotation=45)
plt.show()
```



Note: From above Graph you can see that most of the 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 sees M-Size as the preferred choice among buyers.