

Exploratory_Data_Analysis_Retail

```
In [45]: import numpy as np
import pandas as pd
import seaborn as sns
from plotnine import *
import warnings
warnings.filterwarnings('ignore')
import matplotlib.pyplot as plt
```

Reading the dataset

```
In [46]: sample = pd.read_csv("SampleSuperstore.csv")
```

```
In [47]: sample.head()
```

Out[47]:

	Ship Mode	Segment	Country	City	State	Postal Code	Region	Category	Sub-Category	Sales	Quantity	Discount	Profit
0	Second Class	Consumer	United States	Henderson	Kentucky	42420	South	Furniture	Bookcases	261.9600	2	0.00	41.9136
1	Second Class	Consumer	United States	Henderson	Kentucky	42420	South	Furniture	Chairs	731.9400	3	0.00	219.5820
2	Second Class	Corporate	United States	Los Angeles	California	90036	West	Office Supplies	Labels	14.6200	2	0.00	6.8714
3	Standard Class	Consumer	United States	Fort Lauderdale	Florida	33311	South	Furniture	Tables	957.5775	5	0.45	-383.0310
4	Standard Class	Consumer	United States	Fort Lauderdale	Florida	33311	South	Office Supplies	Storage	22.3680	2	0.20	2.5164

```
In [48]: sample.tail()
```

Out[48]:

	Ship Mode	Segment	Country	City	State	Postal Code	Region	Category	Sub-Category	Sales	Quantity	Discount	Profit
9989	Second Class	Consumer	United States	Miami	Florida	33180	South	Furniture	Furnishings	25.248	3	0.2	4.1028
9990	Standard Class	Consumer	United States	Costa Mesa	California	92627	West	Furniture	Furnishings	91.960	2	0.0	15.6332
9991	Standard Class	Consumer	United States	Costa Mesa	California	92627	West	Technology	Phones	258.576	2	0.2	19.3932
9992	Standard Class	Consumer	United States	Costa Mesa	California	92627	West	Office Supplies	Paper	29.600	4	0.0	13.3200
9993	Second Class	Consumer	United States	Westminster	California	92683	West	Office Supplies	Appliances	243.160	2	0.0	72.9480

```
In [49]: sample.shape
```

Out[49]: (9994, 13)

Checking for the data's information, i.e type

```
In [50]: sample.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 9994 entries, 0 to 9993
Data columns (total 13 columns):
#   Column          Non-Null Count  Dtype
---  ---
0   Ship Mode       9994 non-null   object
1   Segment         9994 non-null   object
2   Country         9994 non-null   object
3   City            9994 non-null   object
4   State           9994 non-null   object
5   Postal Code     9994 non-null   int64
6   Region          9994 non-null   object
7   Category        9994 non-null   object
8   Sub-Category    9994 non-null   object
9   Sales           9994 non-null   float64
10  Quantity        9994 non-null   int64
11  Discount        9994 non-null   float64
12  Profit          9994 non-null   float64
dtypes: float64(3), int64(2), object(8)
memory usage: 1015.1+ KB
```

```
In [51]: sample.describe()
```

Out[51]:

	Postal Code	Sales	Quantity	Discount	Profit
count	9994.000000	9994.000000	9994.000000	9994.000000	9994.000000
mean	55190.379428	229.858001	3.789574	0.156203	28.656896
std	32063.693350	623.245101	2.225110	0.206452	234.260108
min	1040.000000	0.444000	1.000000	0.000000	-6599.978000
25%	23223.000000	17.280000	2.000000	0.000000	1.728750
50%	56430.500000	54.490000	3.000000	0.200000	8.666500
75%	90008.000000	209.940000	5.000000	0.200000	29.364000
max	99301.000000	22638.480000	14.000000	0.800000	8399.976000

```
In [52]: sample.isnull().sum()
```

Out[52]:

```
Ship Mode      0
Segment        0
Country        0
City           0
State          0
Postal Code    0
Region         0
Category       0
Sub-Category   0
Sales          0
Quantity       0
Discount       0
Profit         0
dtype: int64
```

Checking for the duplicate data. If yes, then dropping those data

```
In [53]: sample.duplicated().sum()
```

Out[53]: 17

```
In [54]: sample.drop_duplicates()
```

Out[54]:

	Ship Mode	Segment	Country	City	State	Postal Code	Region	Category	Sub-Category	Sales	Quantity	Discount	Profit
0	Second Class	Consumer	United States	Henderson	Kentucky	42420	South	Furniture	Bookcases	261.9600	2	0.00	41.9136
1	Second Class	Consumer	United States	Henderson	Kentucky	42420	South	Furniture	Chairs	731.9400	3	0.00	219.5820
2	Second Class	Corporate	United States	Los Angeles	California	90036	West	Office Supplies	Labels	14.6200	2	0.00	6.8714
3	Standard Class	Consumer	United States	Fort Lauderdale	Florida	33311	South	Furniture	Tables	957.5775	5	0.45	-383.0310
4	Standard Class	Consumer	United States	Fort Lauderdale	Florida	33311	South	Office Supplies	Storage	22.3680	2	0.20	2.5164
...
9989	Second Class	Consumer	United States	Miami	Florida	33180	South	Furniture	Furnishings	25.2480	3	0.20	4.1028
9990	Standard Class	Consumer	United States	Costa Mesa	California	92627	West	Furniture	Furnishings	91.9600	2	0.00	15.6332
9991	Standard Class	Consumer	United States	Costa Mesa	California	92627	West	Technology	Phones	258.5760	2	0.20	19.3932
9992	Standard Class	Consumer	United States	Costa Mesa	California	92627	West	Office Supplies	Paper	29.6000	4	0.00	13.3200
9993	Second Class	Consumer	United States	Westminster	California	92683	West	Office Supplies	Appliances	243.1600	2	0.00	72.9480

9977 rows × 13 columns

```
In [55]: sample.nunique()
```

Out[55]:

Ship Mode	4
Segment	3
Country	1
City	531
State	49
Postal Code	631
Region	4
Category	3
Sub-Category	17
Sales	5825
Quantity	14
Discount	12
Profit	7287
dtype:	int64

Dropping irrelevant columns

```
In [56]: col = ['Postal Code']
sample1 = sample.drop(columns=col, axis=1)
sample1
```

Out[56]:

	Ship Mode	Segment	Country	City	State	Region	Category	Sub-Category	Sales	Quantity	Discount	Profit
0	Second Class	Consumer	United States	Henderson	Kentucky	South	Furniture	Bookcases	261.9600	2	0.00	41.9136
1	Second Class	Consumer	United States	Henderson	Kentucky	South	Furniture	Chairs	731.9400	3	0.00	219.5820
2	Second Class	Corporate	United States	Los Angeles	California	West	Office Supplies	Labels	14.6200	2	0.00	6.8714
3	Standard Class	Consumer	United States	Fort Lauderdale	Florida	South	Furniture	Tables	957.5775	5	0.45	-383.0310
4	Standard Class	Consumer	United States	Fort Lauderdale	Florida	South	Office Supplies	Storage	22.3680	2	0.20	2.5164
...
9989	Second Class	Consumer	United States	Miami	Florida	South	Furniture	Furnishings	25.2480	3	0.20	4.1028
9990	Standard Class	Consumer	United States	Costa Mesa	California	West	Furniture	Furnishings	91.9600	2	0.00	15.6332
9991	Standard Class	Consumer	United States	Costa Mesa	California	West	Technology	Phones	258.5760	2	0.20	19.3932
9992	Standard Class	Consumer	United States	Costa Mesa	California	West	Office Supplies	Paper	29.6000	4	0.00	13.3200
9993	Second Class	Consumer	United States	Westminster	California	West	Office Supplies	Appliances	243.1600	2	0.00	72.9480

9994 rows × 12 columns

Checking statistical relation between the various rows & columns

```
In [57]: # Correlation between variables
sample1.corr()
```

Out[57]:

	Sales	Quantity	Discount	Profit
Sales	1.000000	0.200795	-0.028190	0.479064
Quantity	0.200795	1.000000	0.008623	0.066253
Discount	-0.028190	0.008623	1.000000	-0.219487
Profit	0.479064	0.066253	-0.219487	1.000000

```
In [58]: # Covariance of columns
sample1.cov()
```

Out[58]:

	Sales	Quantity	Discount	Profit
Sales	388434.455308	278.459923	-3.627228	69944.096586
Quantity	278.459923	4.951113	0.003961	34.534769
Discount	-3.627228	0.003961	0.042622	-10.615173
Profit	69944.096586	34.534769	-10.615173	54877.798055

```
In [59]: sample1.head() # Loads the first five rows
```

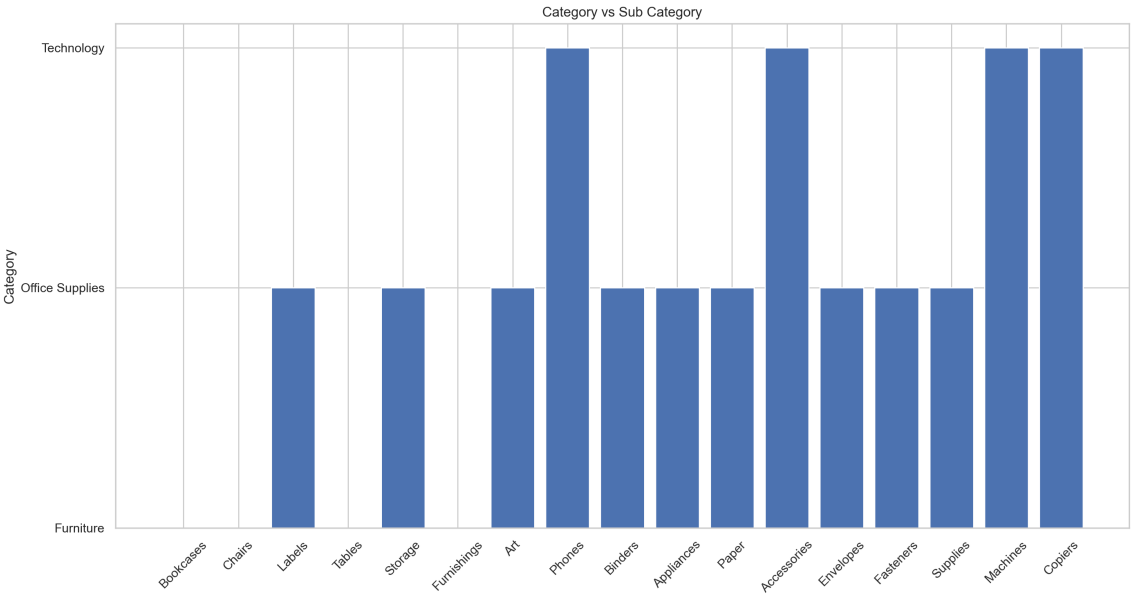
Out[59]:

	Ship Mode	Segment	Country	City	State	Region	Category	Sub-Category	Sales	Quantity	Discount	Profit
0	Second Class	Consumer	United States	Henderson	Kentucky	South	Furniture	Bookcases	261.9600	2	0.00	41.9136
1	Second Class	Consumer	United States	Henderson	Kentucky	South	Furniture	Chairs	731.9400	3	0.00	219.5820
2	Second Class	Corporate	United States	Los Angeles	California	West	Office Supplies	Labels	14.6200	2	0.00	6.8714
3	Standard Class	Consumer	United States	Fort Lauderdale	Florida	South	Furniture	Tables	957.5775	5	0.45	-383.0310
4	Standard Class	Consumer	United States	Fort Lauderdale	Florida	South	Office Supplies	Storage	22.3680	2	0.20	2.5164

Data Visualization

In [60]:

```
plt.figure(figsize=(16,8))
plt.bar('Sub-Category', 'Category', data=sample1)
plt.title('Category vs Sub Category')
plt.xlabel('Sub-Category')
plt.ylabel('Category')
plt.xticks(rotation=45)
plt.show()
```



In [61]:

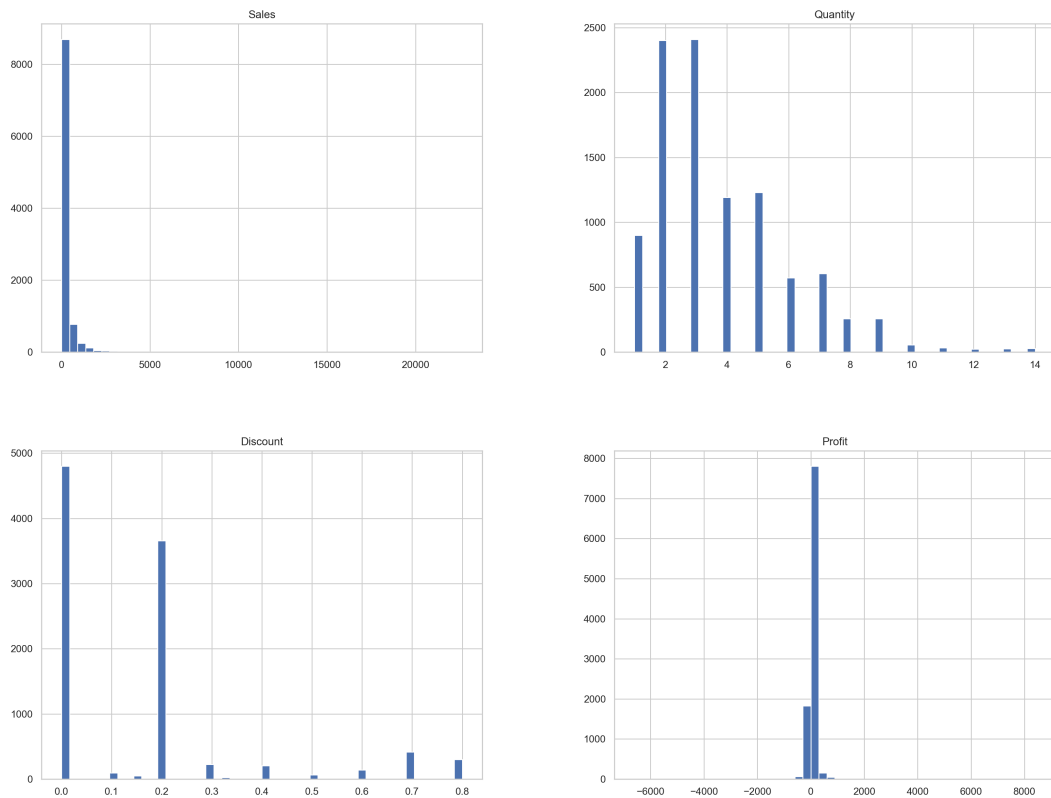
```
sample1.corr() # Checking the correlation
```

Out[61]:

	Sales	Quantity	Discount	Profit
Sales	1.000000	0.200795	-0.028190	0.479064
Quantity	0.200795	1.000000	0.008623	0.066253
Discount	-0.028190	0.008623	1.000000	-0.219487
Profit	0.479064	0.066253	-0.219487	1.000000

In [62]:

```
sample1.hist(bins=50,figsize=(20,15))
plt.show();
```

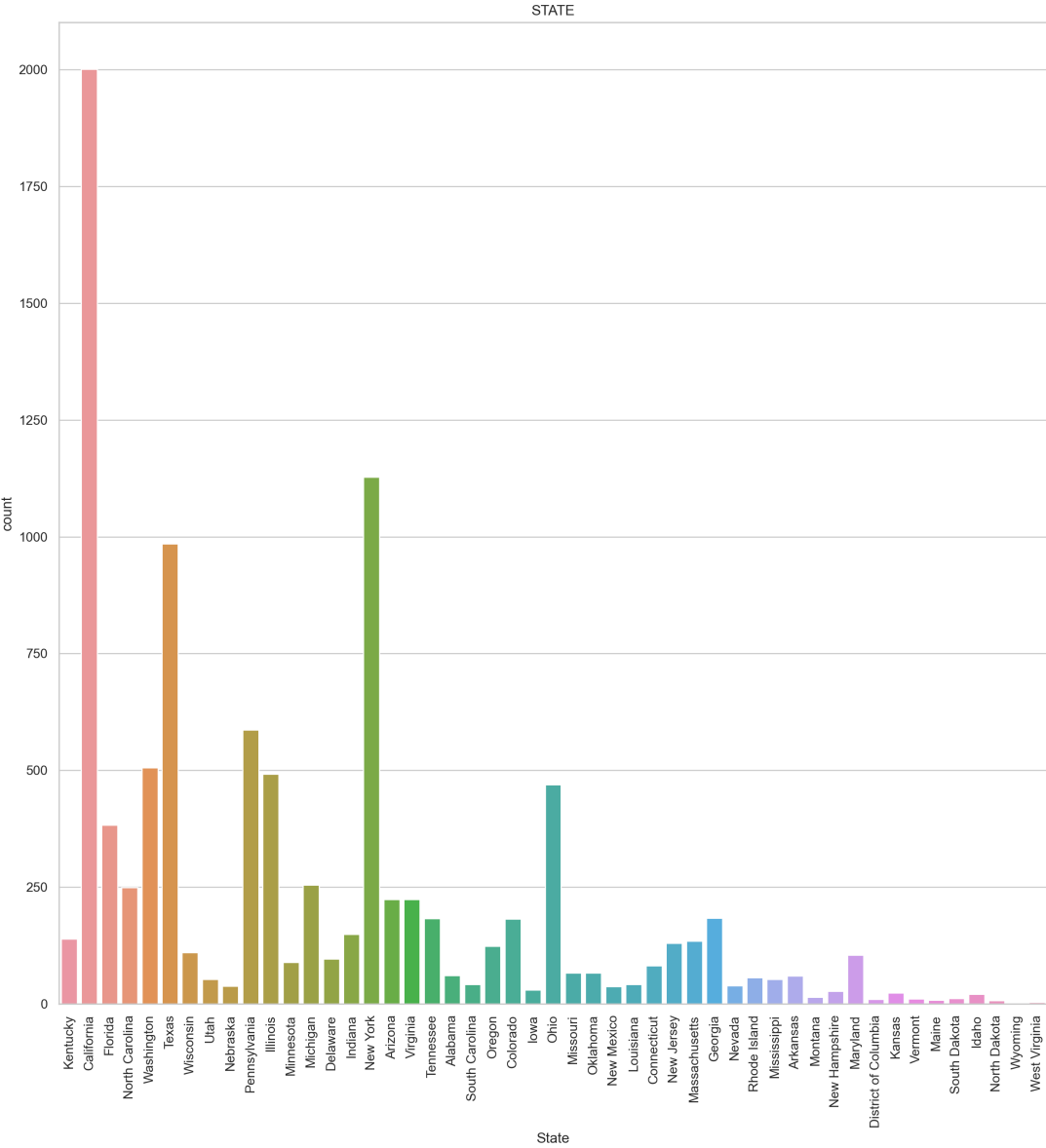


So, the data here is not normal as revealed by this histogram graph.

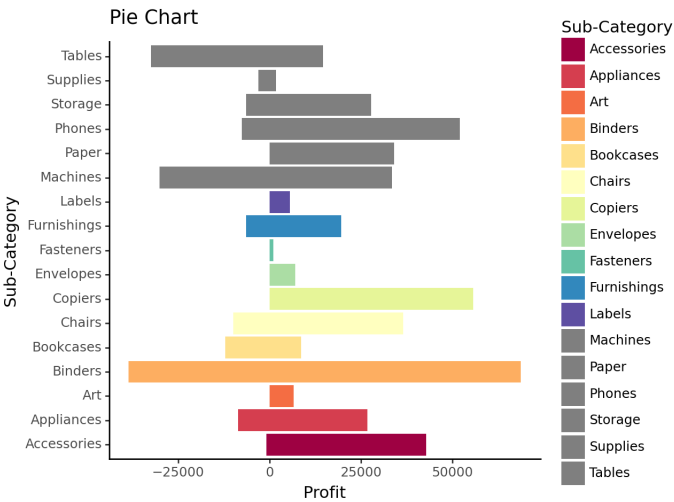
```
In [63]: # Count the total repeatable states
sample1['State'].value_counts()
```

```
Out[63]: California      2001
New York      1128
Texas         985
Pennsylvania  587
Washington    506
Illinois      492
Ohio          469
Florida       383
Michigan       255
North Carolina 249
Arizona       224
Virginia      224
Georgia       184
Tennessee     183
Colorado      182
Indiana       149
Kentucky      139
Massachusetts 135
New Jersey    130
Oregon        124
Wisconsin     110
Maryland      105
Delaware      96
Minnesota     89
Connecticut   82
Oklahoma      66
Missouri      66
Alabama       61
Arkansas      60
Rhode Island  56
Utah          53
Mississippi   53
Louisiana     42
South Carolina 42
Nevada        39
Nebraska      38
New Mexico    37
Iowa          30
New Hampshire 27
Kansas        24
Idaho         21
Montana       15
South Dakota  12
Vermont       11
District of Columbia 10
Maine         8
North Dakota  7
West Virginia 4
Wyoming       1
Name: State, dtype: int64
```

```
In [66]: plt.figure(figsize=(15,15))
sns.countplot(x=sample1['State'])
plt.xticks(rotation=90)
plt.title("STATE")
plt.show()
```



```
In [25]: Profit_plot = (ggplot(sample, aes(x='Sub-Category', y='Profit', fill='Sub-Category')) + geom_col() + coord_flip()
+ scale_fill_brewer(type='div', palette="Spectral") + theme_classic() + ggtitle('Pie Chart'))
display(Profit_plot)
```



<Figure Size: (640 x 480)>

The above pie chart shows the profit and loss of each and every subcategories.

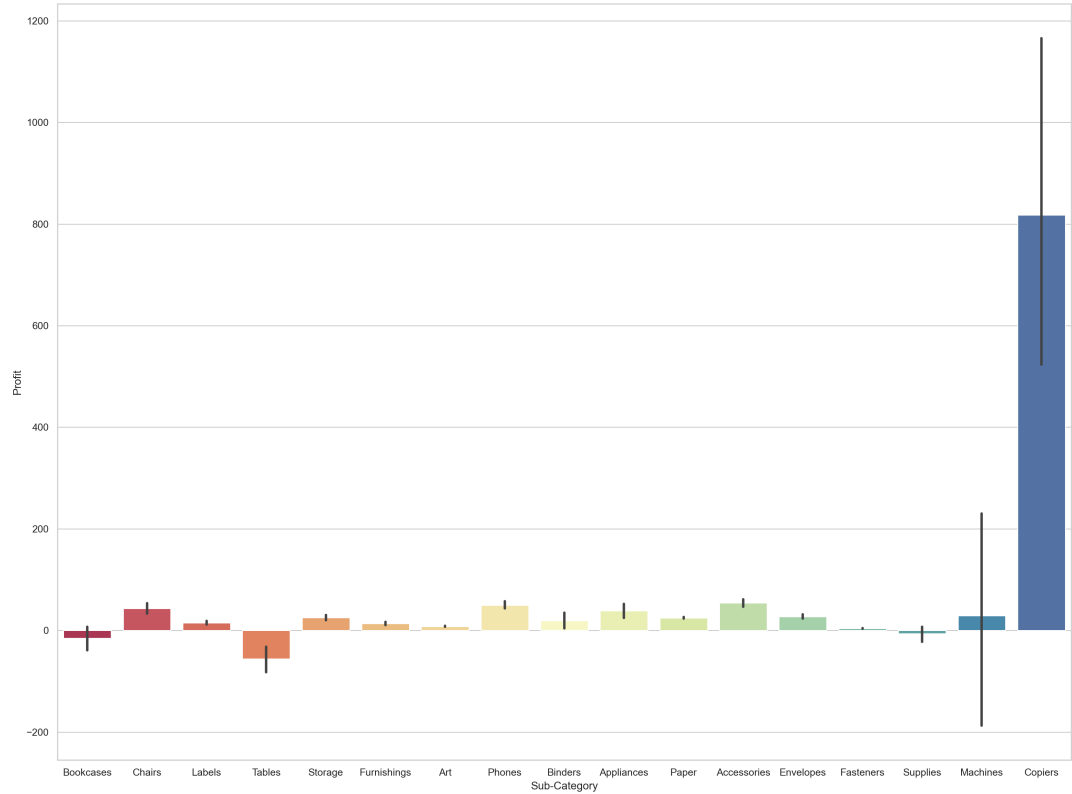
Here from the graph we can visualize that "binders" sub-category has suffered the highest amount of loss and also profit amongst all other sub-Categories (For now we can't say that what is the reason it may be because of discounts given on binders subcategory).

Next,"Copiers" sub-category has gained highest amount of profit with no loss.There are other sub-categories too haven't faced any kind of losses but their profit margins are also low.

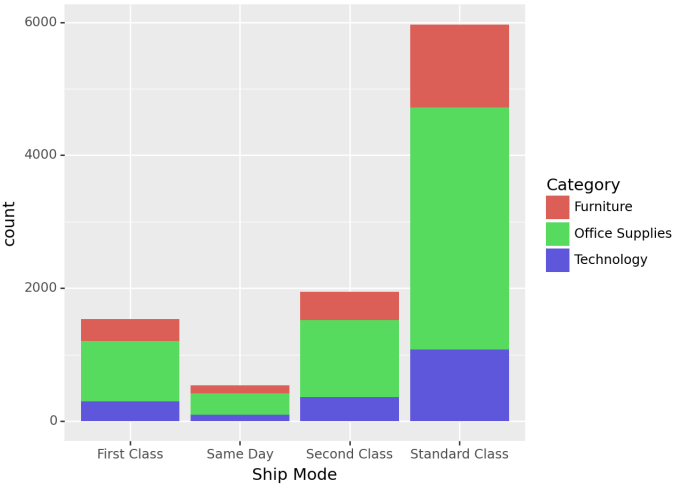
Next, suffering from highest loss is machines.

```
In [26]: sns.set(style="whitegrid")
plt.figure(2, figsize=(20,15))
sns.barplot(x='Sub-Category',y='Profit', data=sample, palette='Spectral')
plt.suptitle('Pie Consumption Patterns in the United States', fontsize=16)
plt.show()
```

Pie Consumption Patterns in the United States



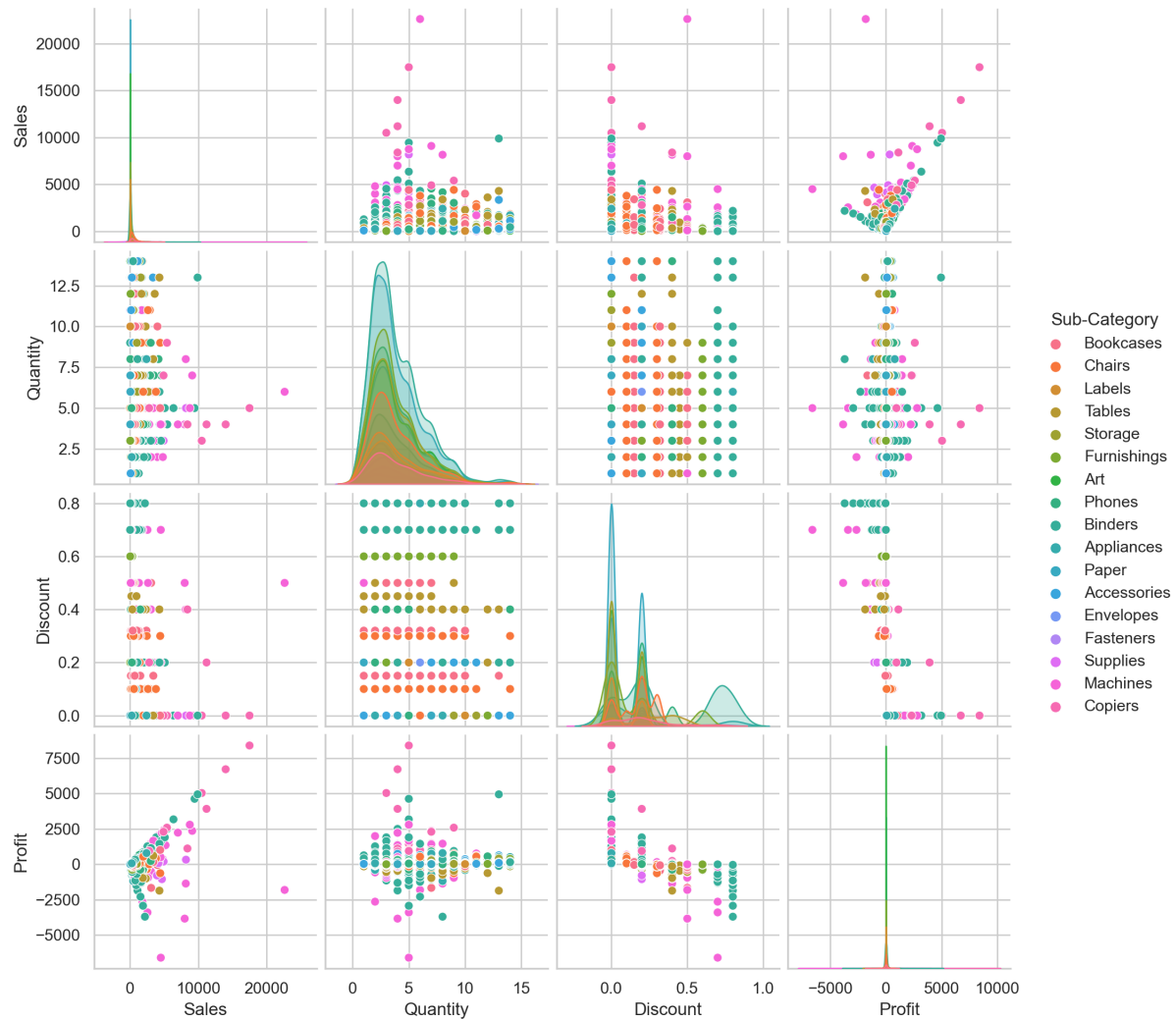
```
In [31]: ggplot(sample, aes(x='Ship Mode', fill = 'Category')) + geom_bar(stat = 'count')
```



Out[31]: <Figure Size: (640 x 480)>

```
In [28]: figsize=(15,10)
sns.pairplot(sample1,hue='Sub-Category')
plt.show
```

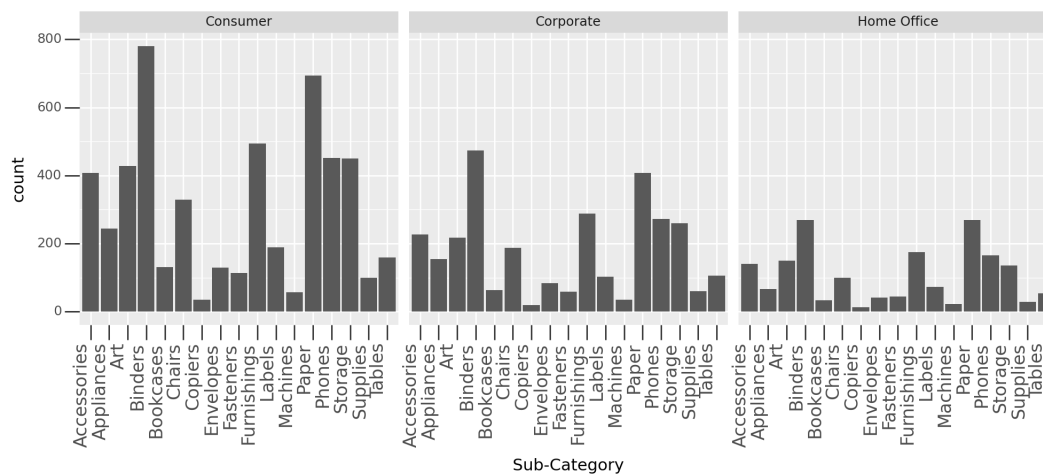
Out[28]: <function matplotlib.pyplot.show(close=None, block=None)>



From the above plot we can say that our data is not normal and it has some amount of outliers too. Let's explore more about these outliers by using boxplots. 1st we'll check sales from every segments of whole data.

```
In [32]: flip_xlabels = theme(axis_text_x = element_text(angle=90, hjust=1), figure_size=(10,5),
axis_ticks_length_major=10, axis_ticks_length_minor=5)
(ggplot(sample, aes(x='Sub-Category', fill='Sales')) + geom_bar() + facet_wrap(['Segment']))
+ flip_xlabels + theme(axis_text_x = element_text(size=12)) + ggtitle("Sales From Every Segment Of United States of Whole Data")
```

Sales From Every Segment Of United States of Whole Data

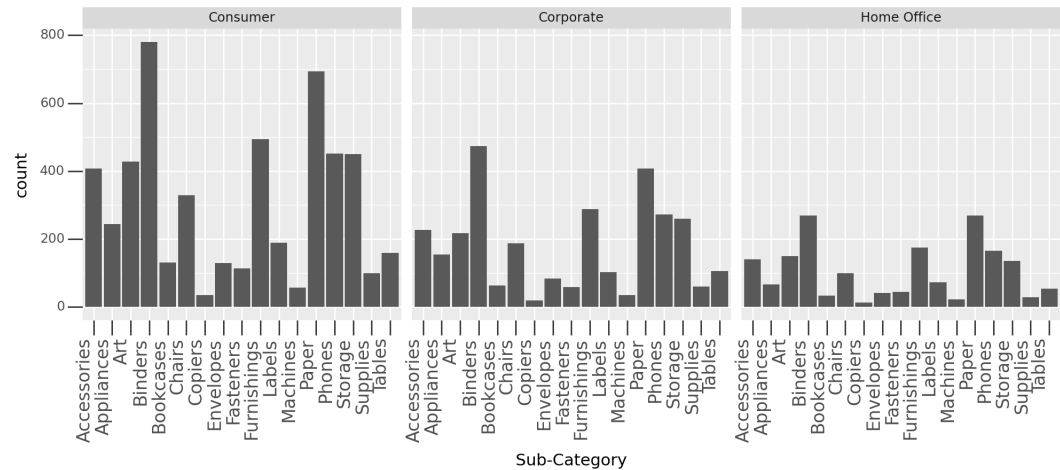


Out[32]: <Figure Size: (1000 x 500)>

From the above Graph we can say that "Home Office" segment has less purchased sub-categories and in that "Tables", "Supplies", "Machines", "Copiers", "Bookcases" has the lowest sales. "Consumer" has purchased more sub-categories as compared to other segments.

```
In [33]: flip_xlabels = theme(axis_text_x = element_text(angle=90, hjust=1), figure_size=(10,5),
axis_ticks_length_major=10, axis_ticks_length_minor=5)
(ggplot(sample, aes(x='Sub-Category', fill='Discount')) + geom_bar() + facet_wrap(['Segment']))
+ flip_xlabels + theme(axis_text_x = element_text(size=12)) + ggtitle("Discount on Categories From Every Segment Of United States of Whole Data")
```

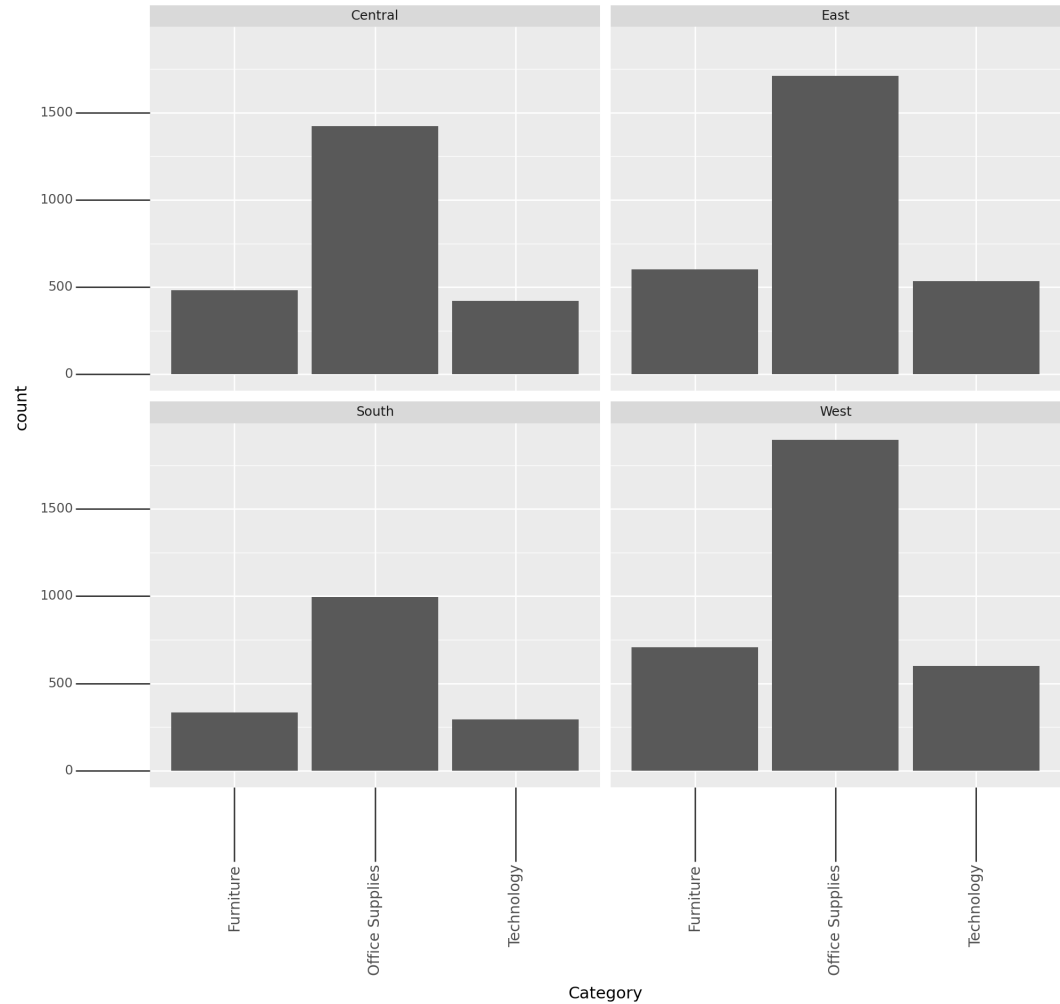
Discount on Categories From Every Segment Of United States of Whole Data



Out[33]: <Figure Size: (1000 x 500)>

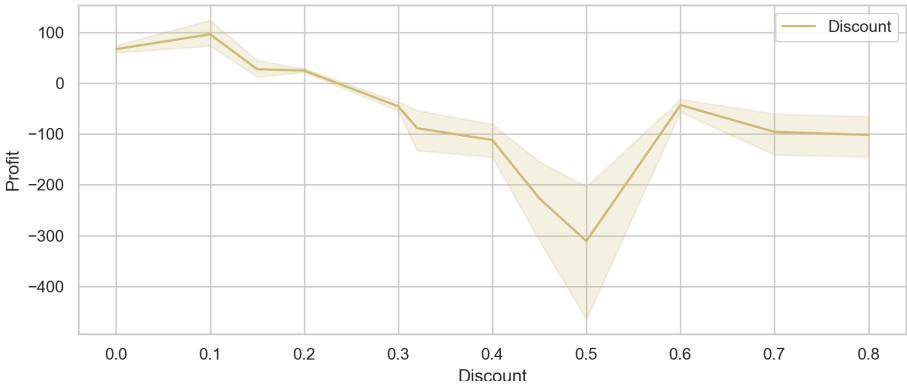
```
In [34]: flip_xlabels = theme(axis_text_x = element_text(angle=90, hjust=10),figure_size=(10,10),
axis_ticks_length_major=50,axis_ticks_length_minor=50)
(ggplot(sample1, aes(x='Category', fill='Sales')) + geom_bar() + theme(axis_text_x = element_text(size=10))
+ facet_wrap(['Region']) + flip_xlabels+ ggtitle("Sales From Every Region Of United States of Whole Data"))
```

Sales From Every Region Of United States of Whole Data



Out[34]: <Figure Size: (1000 x 1000)>

```
In [36]: plt.figure(figsize=(10,4))
sns.lineplot(x='Discount', y='Profit', data=sample1, color='y', label='Discount')
plt.legend()
plt.show()
```

```
In [37]: import plotly.express as px
import plotly.graph_objects as go
from plotly.subplots import make_subplots

In [38]: state_code = {'Alabama': 'AL','Alaska': 'AK','Arizona': 'AZ','Arkansas': 'AR','California': 'CA','Colorado': 'CO','Connecticut': 'CT','Delaware': 'DE','Florida': 'FL','Georgia': 'GA','Hawaii': 'HI','Idaho': 'ID','Illinois': 'IL','Indiana': 'IN','Iowa': 'IA','Kansas': 'KS','Kentucky': 'KY','Louisiana': 'LA','Maine': 'ME','Maryland': 'MD','Massachusetts': 'MA','Michigan': 'MI','Minnesota': 'MN','Mississippi': 'MS','Missouri': 'MO','Montana': 'MT','Nebraska': 'NE','Nevada': 'NV','New Hampshire': 'NH','New Jersey': 'NJ','New Mexico': 'NM','New York': 'NY','North Carolina': 'NC','North Dakota': 'ND','Ohio': 'OH','Oklahoma': 'OK','Oregon': 'OR','Pennsylvania': 'PA','Rhode Island': 'RI','South Carolina': 'SC','South Dakota': 'SD','Tennessee': 'TN','Texas': 'TX','Utah': 'UT','Vermont': 'VT','Virginia': 'VA','Washington': 'WA','West Virginia': 'WV','Wisconsin': 'WI','Wyoming': 'WY'}
sample1['state_code'] = sample1.State.apply(lambda x: state_code[x])

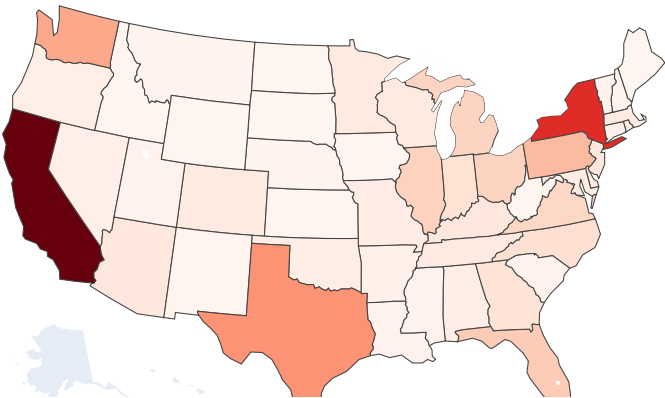
In [39]: state_data = sample1[['Sales', 'Profit', 'state_code']].groupby(['state_code']).sum()

fig = go.Figure(data=go.Choropleth(
    locations=state_data.index,
    z = state_data.Sales,
    locationmode = 'USA-states',
    colorscale = 'Reds',
    colorbar_title = 'Sales in USD',
))

fig.update_layout(
    title_text = 'Total State-Wise Sales',
    geo_scope='usa',
    height=800,
)

fig.show()
```

Total State-Wise Sales



Now, let us analyze the sales of a few random states from each profit bracket (high profit, medium profit, low profit, low loss and high loss) and try to observe some crucial trends which might help us in increasing the sales.

We have a few **questions** to answer here.

- 1. What products do the most profit making states buy?
- 2. What products do the loss bearing states buy?
- 3. What product segment needs to be improved in order to drive the profits higher?

```
In [40]: def state_data_viewer(states):
    """Plots the turnover generated by different product categories and sub-categories for the list of given states.
    Args:
```

```
states- List of all the states you want the plots for
Returns:
None
"""
product_data = sample1.groupby(['State'])
for state in states:
    data = product_data.get_group(state).groupby(['Category'])
    fig, ax = plt.subplots(1, 3, figsize = (28,5))
    fig.suptitle(state, fontsize=14)
    ax_index = 0
    for cat in ['Furniture', 'Office Supplies', 'Technology']:
        cat_data = data.get_group(cat).groupby(['Sub-Category']).sum()
        sns.barplot(x = cat_data.Profit, y = cat_data.index, ax = ax[ax_index])
        ax[ax_index].set_ylabel(cat)
        ax_index +=1
    fig.show()
```

In [41]: states = ['California', 'Washington', 'Mississippi', 'Arizona', 'Texas']
state_data_viewer(states)



Thank You!

GitHub: <https://github.com/anujtiwari21?tab=repositories>