

- IBM AICTE MAJOR PROJECT FINAL CODE
- student id- STU62a5e631eb60e1655039537

```
dataset = pd.read_csv('/content/SampleSuperstore.csv')
```

	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	261,860.0	2	0.00	41,913.6
1	Second Class	Consumer	United States	Henderson	Kentucky	42420	South	Furniture	Chairs	731,940.0	3	0.00	219,582.0
2	Second Class	Corporate	United States	Los Angeles	California	90036	West	Office Supplies	Labels	14,620.0	2	0.00	6,871.4
3	Standard Class	Consumer	United States	Fort Lauderdale	Florida	33311	South	Furniture	Tables	957,577.5	6	0.45	-383,310.0
4	Standard Class	Consumer	United States	Fort Lauderdale	Florida	33311	South	Office Supplies	Storage	22,368.0	2	0.20	2,164.0

```
Index(['Ship Mode', 'Segment', 'Country', 'City', 'State', 'Postal Code',
       'Region', 'Category', 'Sub-Category', 'Sales', 'Quantity', 'Discount',
       'Profit'],
      dtype='object')
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
Data columns (total 13 columns):
#   Column              Non-Null Count  Dtype
---  -
```

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

9	Sales	9994 non-null	float64
10	Quantity	9994 non-null	int64

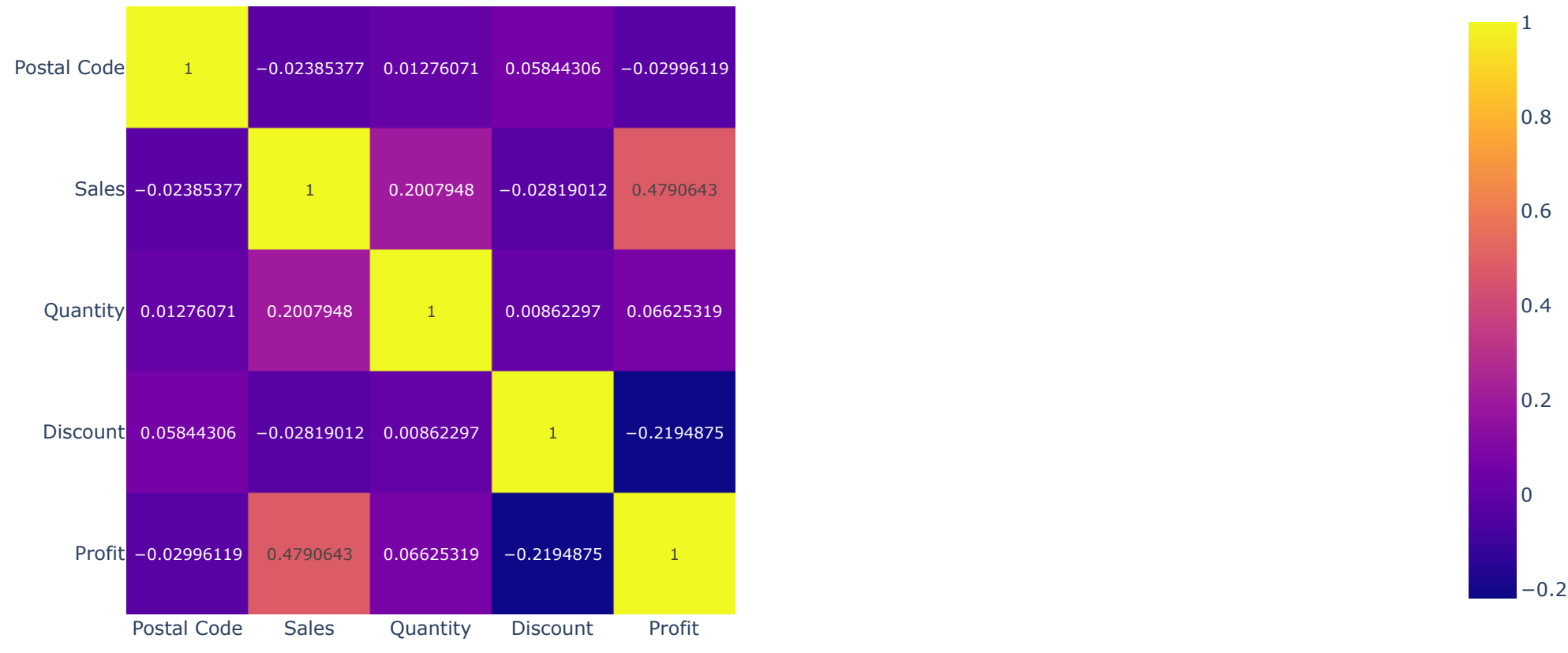
11	Discount	9994	non-null	float64
12	Profit	9994	non-null	float64

```
memory usage: 1015.1+ KB
```

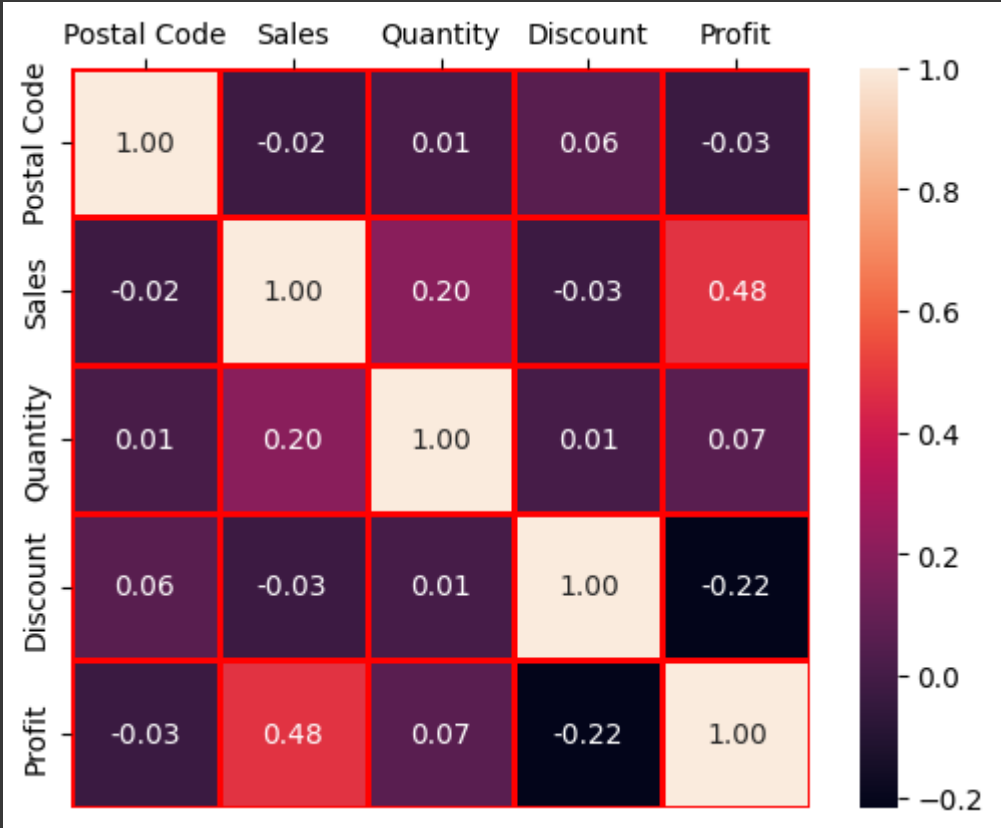
	Postal code	Sales	Quantity	Discount	Profit
count	9994.000000	9994.000000	9994.000000	9994.000000	9994.000000
mean	55190.379428	229.858001	3.789574	0.156203	28.656896
std	32063.693350	623.245101	2.225110	0.206452	234.261001
min	1940.000000	0.444300	1.000000	0.000000	-8589.978000
25%	12323.000000	17.288000	2.000000	0.000000	1.728750
50%	56430.500000	54.490000	3.000000	0.200000	8.666800
75%	90308.000000	209.940000	5.000000	0.200000	29.364000
max	99031.000000	22638.480000	14.000000	0.800000	8399.976000

	Postal Code	Sales	Quantity	Discount	Profit
Postal Code	1.000000	-0.023854	0.012761	0.058443	-0.029961
Sales	-0.023854	1.000000	0.200795	-0.026190	0.479064
Quantity	0.012761	0.200795	1.000000	0.008623	0.066253
Discount	0.058443	-0.026190	0.008623	1.000000	-0.219487
Profit	-0.029961	0.479064	0.066253	-0.219487	1.000000

```
import plotly.express as px
```

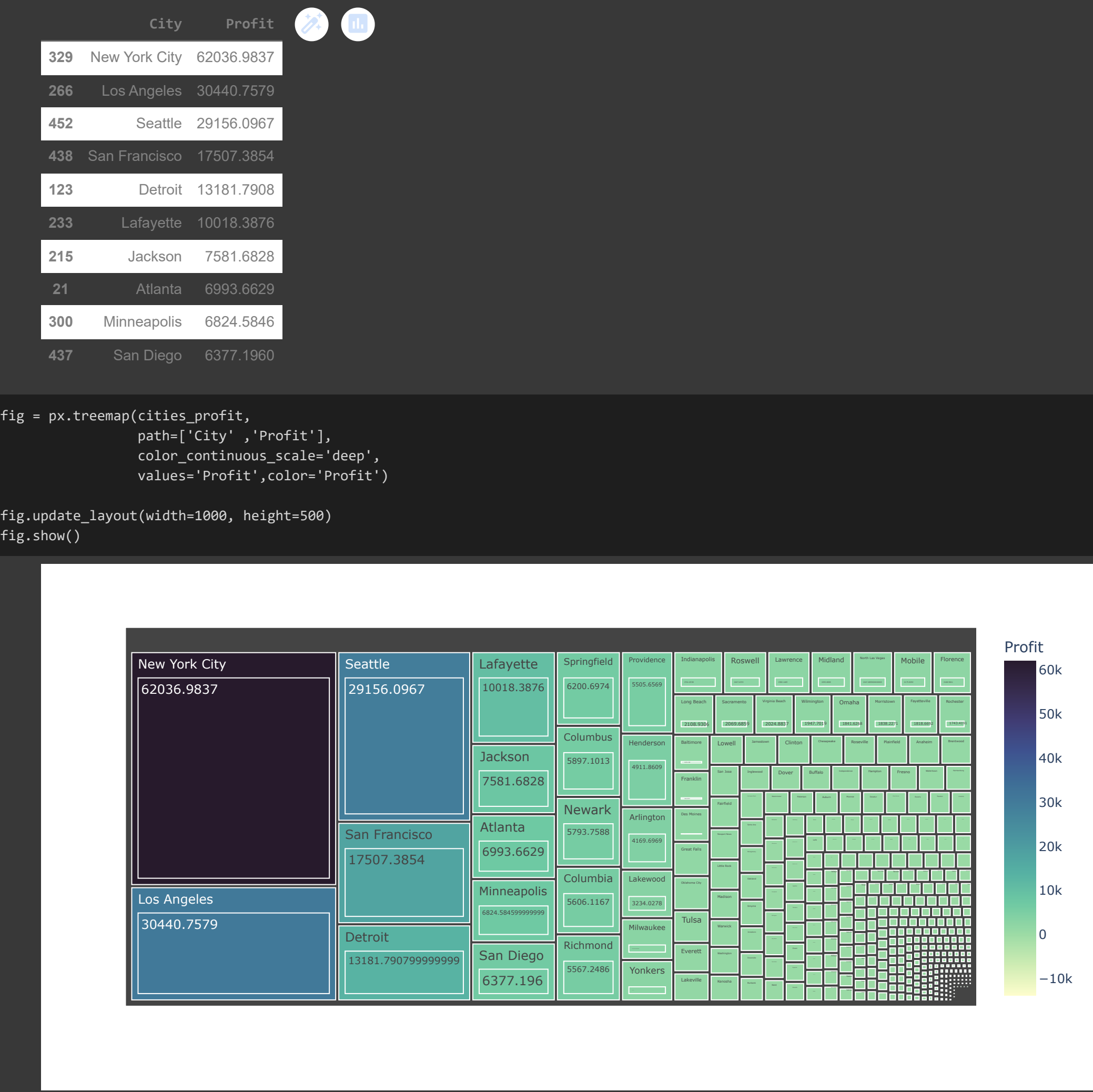


```
ax=sns.heatmap(dataset.corr(),annot=True,cbar=True,square=True,linewidth="red",linewidths=1,xticklabels="auto",yticklabels="auto",fmt=".2f")
ax.set(xlabel="", ylabel="")
ax.xaxis.tick_top()
```

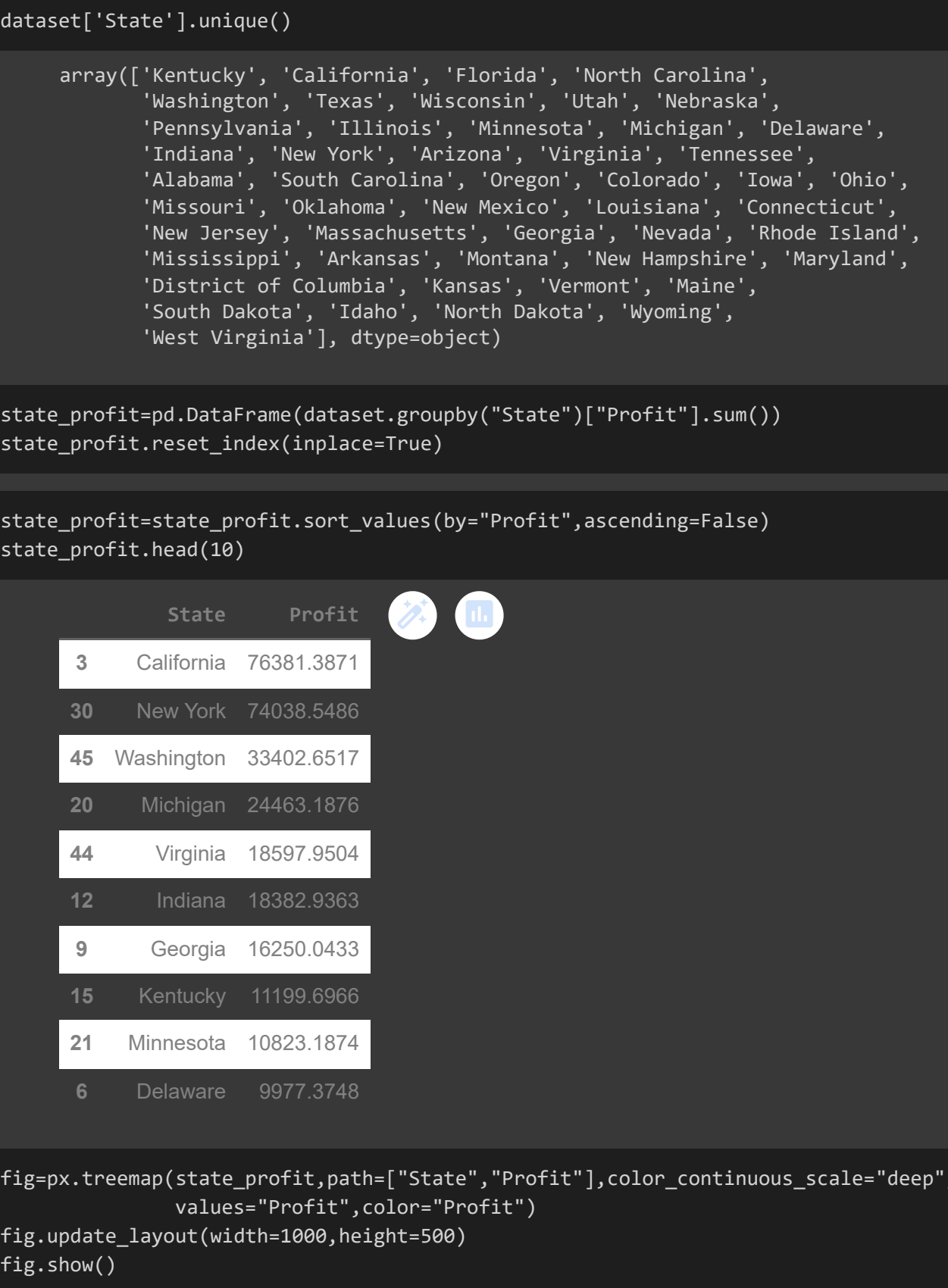


'Miskola', 'Kello', 'Port Orange', 'Medford', 'Charlotteville',
 'Missula', 'Appoka', 'Reading', 'Broomfield', 'Pateron',
 'Oshawa', 'Hawthorne', 'Hawthorne', 'Hawthorne',
 'San Bernardino', 'Leconstein', 'Bozeman', 'Parth Amby',
 'Ontario', 'Rancho Cucamonga', 'Moonhead', 'Mesquite', 'Stockton',
 'Ormond beach', 'Sunnyvale', 'York', 'College Station',
 'San Antonio', 'San Antonio', 'San Antonio',
 'Knoxville', 'Little Rock', 'Lincoln Park', 'Marion', 'Littletown',
 'Rogers', 'Southaven', 'New Castle', 'Midland', 'Sloat Falls',
 'Fond du Lac', 'Fond du Lac', 'Fond du Lac',
 'Malden', 'Holyoke', 'Albuquerque', 'Sparks', 'Coachella',
 'Elmhurst', 'Passaic', 'North Charleston', 'Newport News',
 'Jamestown', 'Mishawaka', 'La Quinta', 'Tallahassee', 'Nashville',
 'San Francisco', 'London', 'London',
 'Summerville', 'Hot Springs', 'Englewood', 'Las Cruces', 'Hooover',
 'Frisco', 'Vacaville', 'Mackinaw', 'Bakersfield', 'Pompano Beach',
 'Fond du Lac', 'Fond du Lac', 'Fond du Lac',
 'Lake Charles', 'Highland Park', 'Hempstead', 'Noblesville',
 'Apple Valley', 'Mount Pleasant', 'Sterling Heights', 'Eau Claire',
 'Pahr', 'Billings', 'Gresham', 'Chattanooga', 'Meridian',
 'Pearland', 'San Mateo', 'Grand Rapids', 'Visalia',
 'Overland Park', 'Temecula', 'Yucapita', 'Revere', 'Conroe',
 'Tinley Park', 'Dubuque', 'Barron Heights', 'Santa Fe',
 'Elk River', 'Caracas', 'San Antonio',
 'Plantation', 'Port Saint Lucie', 'Rochester', 'Altamira',
 'West Falls', 'Chula Vista', 'Manhattan', 'Holliston', 'Thornton',
 'Channahon', 'Channahon', 'Channahon',
 'Woonsocket', 'Superior', 'Bedford', 'Covington', 'Broken Arrow',
 'Miramar', 'Hollywood', 'Deer Park', 'Nichita', 'Mallin',
 'Towson', 'Towson', 'Towson',
 'Elk River', 'Daytona Beach', 'Willard City', 'Portage', 'Fargo',
 'Plantation', 'San Gabriel', 'Margate', 'Sandy Springs', 'Mentor',
 'Lampton', 'Hampton', 'Rome', 'La Crosse', 'Leviston',
 'Huntington', 'Huntington', 'Huntington',
 'Aurouville', 'Marietta', 'Yuma', 'Mausau', 'Pasco', 'Oak Park',
 'Pensacola', 'League City', 'Gaitersburg', 'Lehi', 'Tuscaloosa',
 'Moreno Valley', 'Georgetown', 'Lowland', 'Chandler', 'Helena',
 'Huntington', 'Huntington', 'Huntington',
 'Woodbury', 'Rogers', 'Clovis', 'Duplister', 'Santa Barbara',
 'Cedar Hill', 'Norfolk', 'Draper', 'Ann Arbor', 'La Mesa',
 'Poughkeepsie', 'Poughkeepsie', 'Poughkeepsie',
 'Redding', 'Chico', 'Utica', 'Conway', 'Cheyenne', 'Owensboro',
 'Caldwell', 'Kenner', 'Nashua', 'Bartlett', 'Redwood City',
 'Lebanon', 'Santa Maria', 'Des Plaines', 'Logansville',
 'Huntington', 'Huntington', 'Huntington',
 'Eggen', 'Omard', 'Renton', 'Glenview', 'Delray Beach',
 'Commerce City', 'Texas City', 'Wilson', 'Rio Rancho', 'Goldsboro',
 'Huntington', 'Huntington', 'Huntington',
 'Normal', 'Saint Charles', 'Camarillo', 'Hillsboro', 'Burnsaw',
 'Modesto', 'Garden City', 'Atlantic City', 'Longmont', 'Davis',
 'Morgan Hill', 'Clifton', 'Shenoyagan', 'East Point', 'Rapid City',
 'Huntington', 'Huntington', 'Huntington',
 'San Marcos', 'Greely', 'Mansfield', 'Elyria', 'Twain Falls',
 'Coral Gables', 'Romeoville', 'Marlborough', 'Lynch', 'Bryan',
 'Pahr', 'Pahr', 'Pahr',
 'Arlington Heights', 'Owensboro', 'Con Rapids', 'San Clemente',
 'San Luis Obispo', 'Springdale', 'Lodi', 'Mason', 'dtype=object'

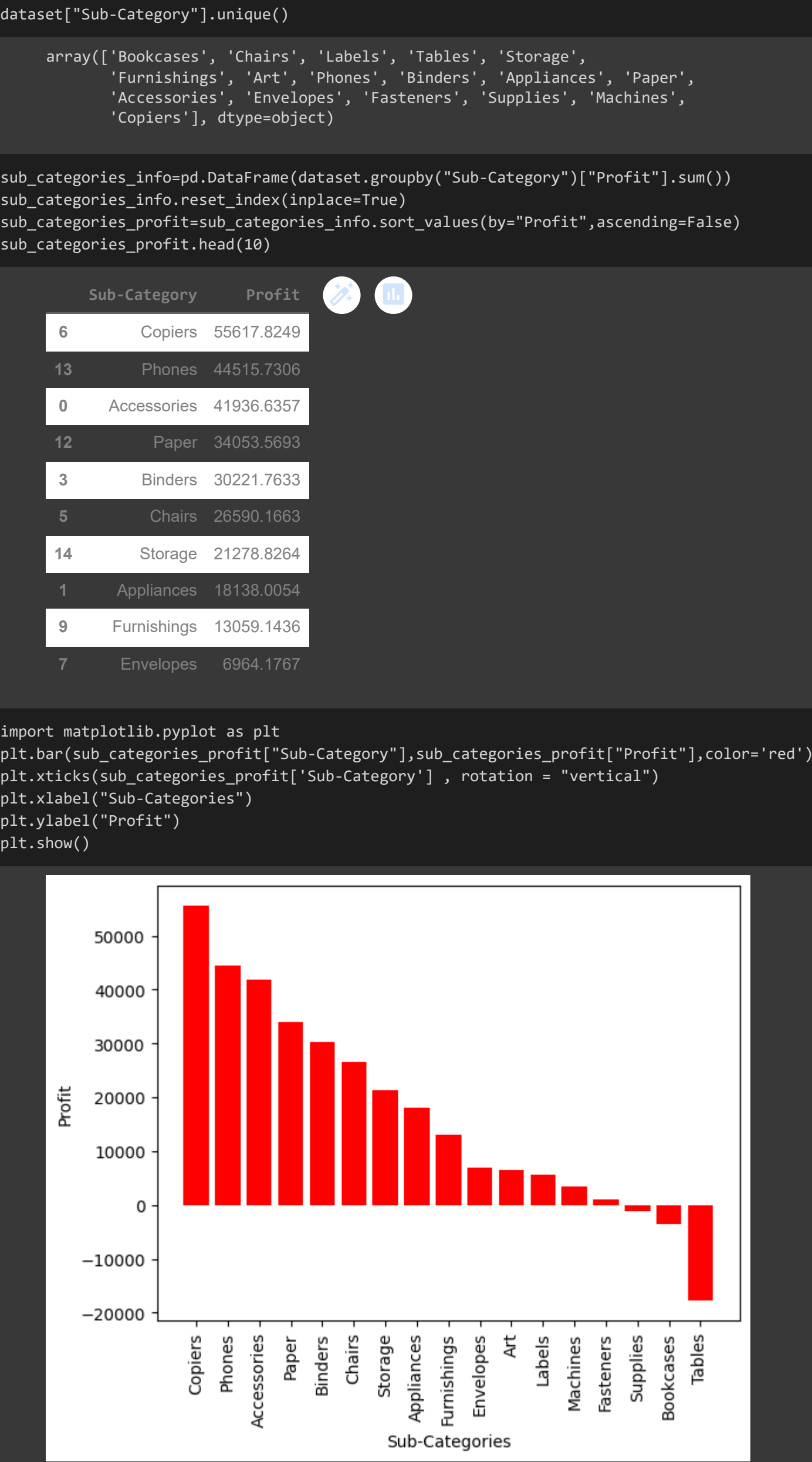
```
cities_profit=pd.DataFrame(dataset.groupby("City")["Profit"].sum())
cities_profit.reset_index(inplace=True)
cities_profit=cities_profit.sort_values(by="Profit",ascending=False)
cities_profit.head(10)
```



10 States by most Profit



Most Profitable Sub Category



Category thaty sales the most

```
dataset["Category"].value_counts()

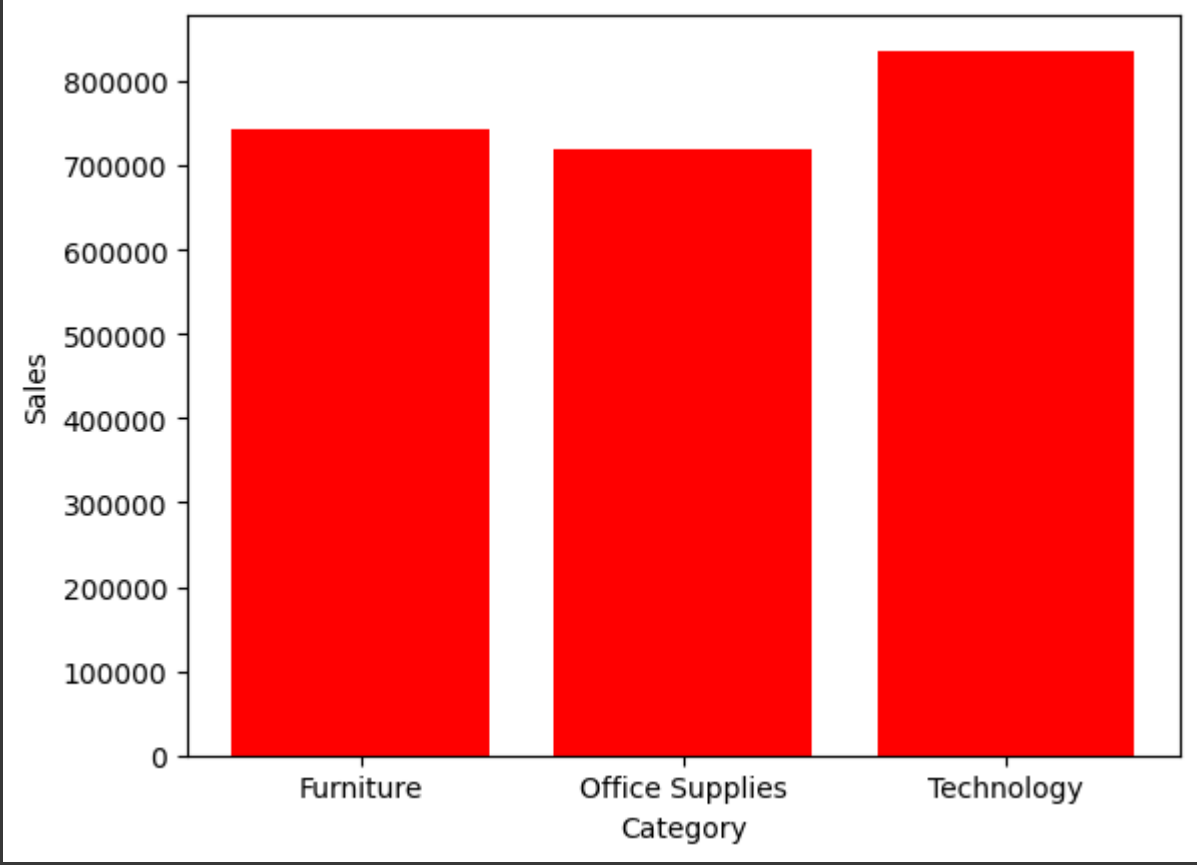
Office Supplies    6026
Furniture          2321
Technology         1847
Name: Category, dtype: int64

categories_sales=dataset.groupby("Category").sum()
categories_sales.reset_index(inplace=True)
categories_sales.head(10)
```



```
<ipython-input-54-e07ea8d23804d>11: FutureWarning:
The default value of numeric_only in DataFrameGroupBy.sum is deprecated. In a future version, numeric_only will default to False. Either specify numeric_only or select only columns which should be valid for the function.
```

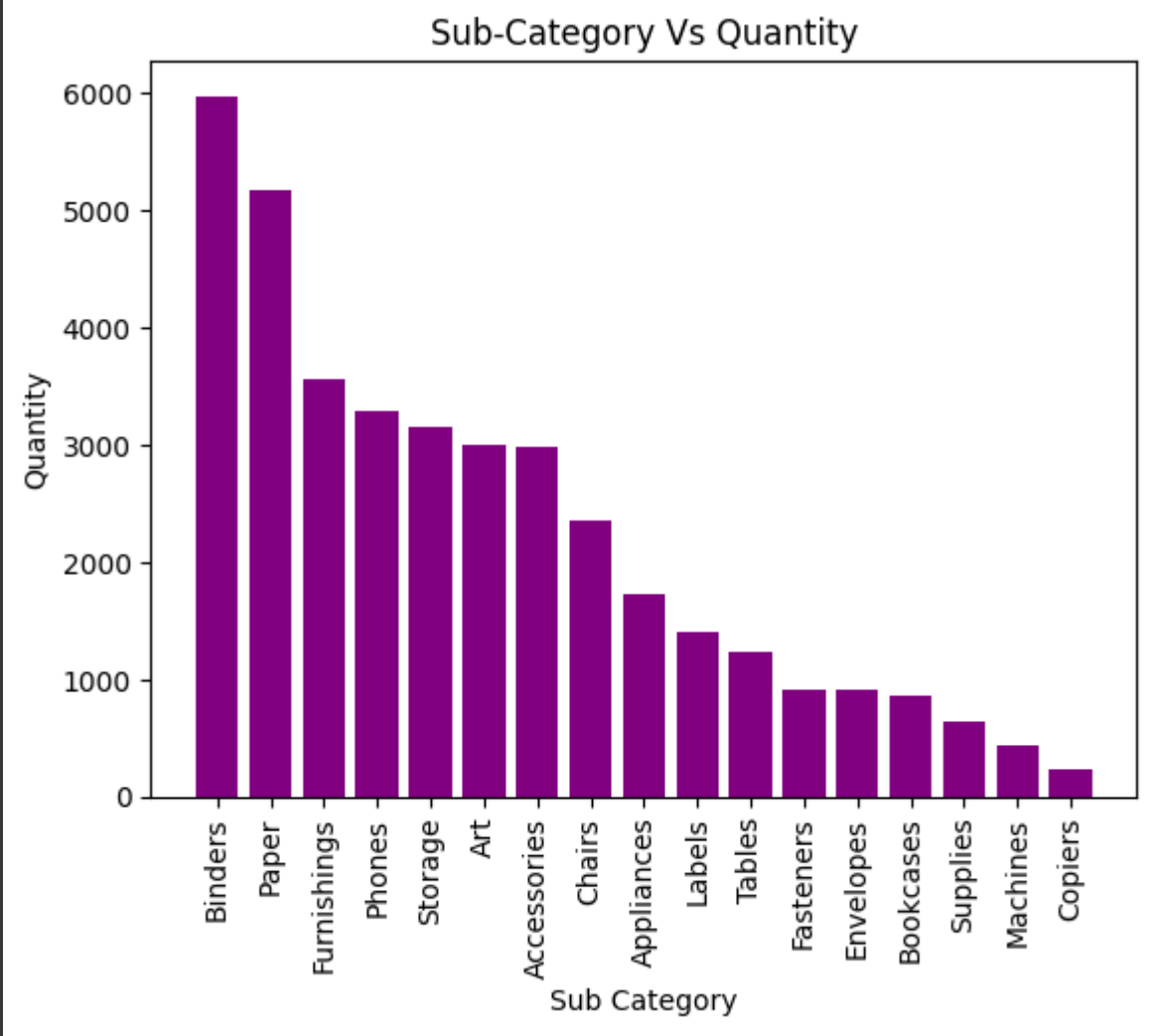
```
Category Postal Code Sales Quantity Discount Profit
0
1
2
plt.bar(categories_sales["Category"], categories_sales['Sales'],color = 'red')
plt.xticks(categories_sales.index)
plt.xlabel('Category')
plt.ylabel('Sales')
plt.show()
```



Most popular Subcategory

```
popular_sub_category=pd.DataFrame(dataset.groupby("Sub-Category")["Quantity"].sum())
popular_sub_category.reset_index(inplace=True)
popular_sub_category=popular_sub_category.sort_values(by="Quantity",ascending=False)
```

```
plt.bar(popular_sub_category["Sub-Category"],popular_sub_category["Quantity"],color="purple")
plt.xticks(popular_sub_category["Sub-Category"],rotation="vertical")
plt.title("Sub-Category Vs Quantity")
plt.xlabel("Sub Category")
plt.ylabel("Quantity")
plt.show()
```



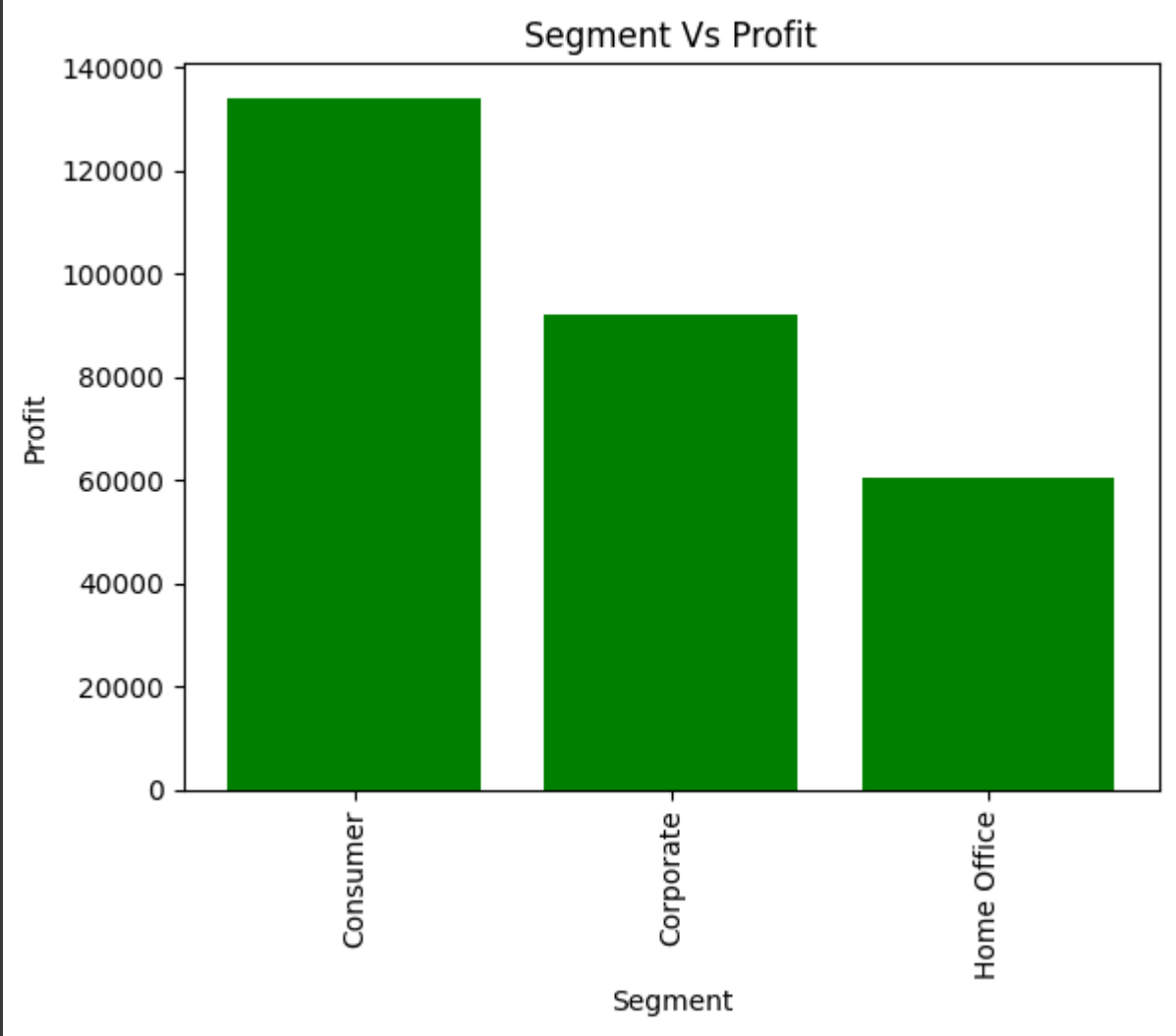
Most Profitable customer segment

```
customer_segment=dataset.groupby("Segment").sum()
customer_segment.reset_index(inplace=True)
customer_segment
```

```
<ipython-input-58-a0be9968528>11: FutureWarning:
The default value of numeric_only in DataFrameGroupBy.sum is deprecated. In a future version, numeric_only will default to False. Either specify numeric_only or select only columns which should be valid for the function.
```

Segment	Postal Code	Sales	Quantity	Discount	Profit
0	Consumer	288878609	1.161401e+06	19521	820.91
1	Corporate	164036230	7.061404e+05	11608	477.85
2	Home Office	98157713	4.296531e+05	6744	262.33

```
plt.bar(customer_segment["Segment"],customer_segment["Profit"],color="green")
plt.xticks(customer_segment["Segment"],rotation="vertical")
plt.title("Segment Vs Profit")
plt.xlabel("Segment")
plt.ylabel("Profit")
plt.show()
```



```
dataset["Region"].value_counts()
West 3283
East 2848
Central 2323
South 1628
Name: Region, dtype: int64
```

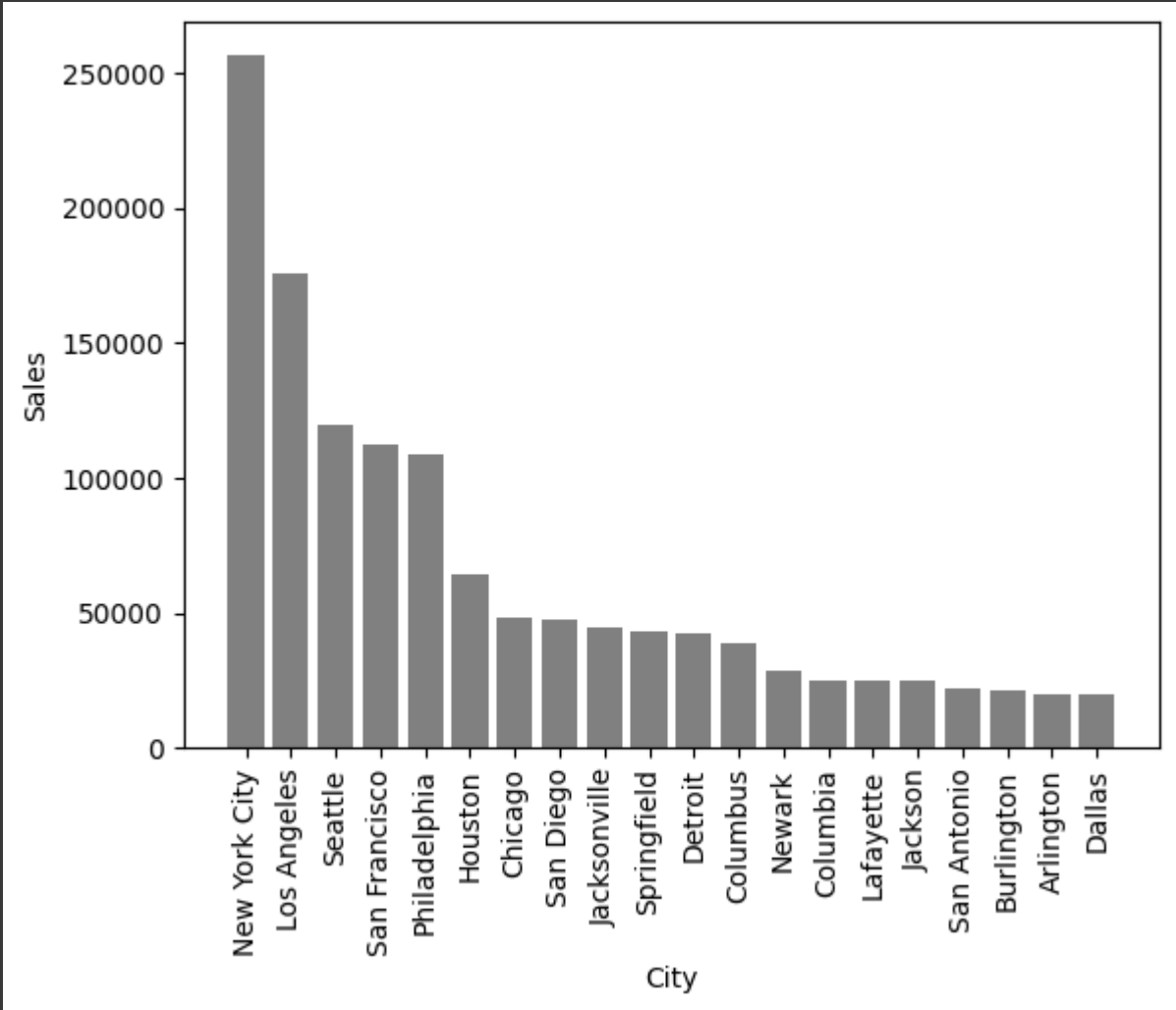
```
profitable_region=pd.DataFrame(dataset.groupby("Region")["Profit"].sum())
profitable_region.reset_index(inplace=True)
profitable_region
```

Region	Profit
0 Central	39706.3625
1 East	91622.7890
2 South	46749.4303
3 West	108418.4489

```
city_volume = pd.DataFrame(dataset.groupby('City')['Sales'].sum())
city_volume =city_volume.sort_values(by="Sales",ascending=False)
city_volume.head()
```

City	Sales
New York City	256368.161
Los Angeles	175851.341
Seattle	119540.742
San Francisco	112669.092
Philadelphia	109077.013

```
plt.bar(city_volume.index[0:20], city_volume['Sales'][0:20],color="gray")
plt.xticks(city_volume.index[0:20],rotation = 'vertical')
plt.xlabel('City')
plt.ylabel('Sales')
plt.show()
```



MACHINE LEARNING MODEL

dataset.head()												
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	261.9800	2	0.00
1	Second Class	Consumer	United States	Henderson	Kentucky	42420	South	Furniture	Chairs	731.9400	3	0.00
2	Second Class	Corporate	United States	Los Angeles	California	90036	West	Office Supplies	Labels	14.6200	2	0.00
3	Standard Class	Consumer	United States	Fort Lauderdale	Florida	33311	South	Furniture	Tables	957.5775	5	0.45
4	Standard Class	Consumer	United States	Fort Lauderdale	Florida	33311	South	Office Supplies	Storage	22.3680	2	0.20

```
dataset1=dataset.drop(columns=["Country","City","State","Sub-Category","Postal Code"])
dataset1.head()
```

```
dataset1.head()
```

	Ship Mode	Segment	Region	Category	Sales	Quantity	Discount	Profit
0	1	0	3	1	261.9600	2	0.00	41.9136
1	1	0	3	1	731.9400	3	0.00	219.5820
2	1	1	0	0	14.6200	2	0.00	6.8714
3	0	0	3	1	957.5775	5	0.45	-383.0310
4	0	0	3	0	22.3680	2	0.20	2.5164

```
dataset1.shape
```

(9994, 8)

```
dataset1["Ship Mode"].value_counts()
```

0	5968
1	1946
2	1538
3	543

Name: Ship Mode, dtype: int64

```
dataset["Category"].value_counts()
```

Office Supplies	6026
Furniture	2121
Technology	1847

Name: Category, dtype: int64

```
dataset1["Region"].value_counts()
```

0	3283
1	2848
2	2323
3	1620

Name: Region, dtype: int64

```
dataset1.replace({"Ship Mode":{"Standard Class":0,"Second Class":1,"First Class":2,"Same Day":3},"Segment":{"Consumer":0,"Corporate":1,"Home Office":2},  
                  "Category":{"Office Supplies":0,"Furniture":1,"Technology":2},"Region":{"West":0,"East":1,"Central":2,"South":3}},inplace=True)
```

```
dataset1.head()
```

	Ship Mode	Segment	Region	Category	Sales	Quantity	Discount	Profit
0	1	0	3	1	261.9600	2	0.00	41.9136
1	1	0	3	1	731.9400	3	0.00	219.5820
2	1	1	0	0	14.6200	2	0.00	6.8714
3	0	0	3	1	957.5775	5	0.45	-383.0310
4	0	0	3	0	22.3680	2	0.20	2.5164

```
x=dataset1.iloc[:,:-1]  
print(x)
```

	Ship Mode	Segment	Region	Category	Sales	Quantity	Discount
0	1	0	3	1	261.9600	2	0.00
1	1	0	3	1	731.9400	3	0.00
2	1	1	0	0	14.6200	2	0.00
3	0	0	3	1	957.5775	5	0.45
4	0	0	3	0	22.3680	2	0.20
...
9989	1	0	3	1	25.2480	3	0.20
9990	0	0	0	1	91.9600	2	0.00
9991	0	0	0	2	256.5760	2	0.20
9992	0	0	0	0	29.6000	4	0.00
9993	1	0	0	0	243.1600	2	0.00

[9994 rows x 7 columns]

```
y=dataset1.iloc[:,1]  
print(y)
```

0	41.9136
1	219.5820
2	6.8714
3	-383.0310
4	2.5164
...	...
9989	4.1028
9990	15.0332
9991	19.3812
9992	13.3200
9993	72.9480

Name: Profit, Length: 9994, dtype: float64

```
from sklearn.preprocessing import MinMaxScaler  
scaler=MinMaxScaler()  
  
num_vars=["Sales","Quantity"]  
x[num_vars]=scaler.fit_transform(x[num_vars])  
print(x)
```

	Ship Mode	Segment	Region	Category	Sales	Quantity	Discount
0	1	0	3	1	0.011552	0.076923	0.00
1	1	0	3	1	0.032313	0.153846	0.00
2	1	1	0	0	0.000626	0.076923	0.00
3	0	0	3	1	0.042200	0.307692	0.45
4	0	0	3	0	0.000968	0.076923	0.20
...
9989	1	0	3	1	0.001806	0.153846	0.20
9990	0	0	0	1	0.004843	0.076923	0.00
9991	0	0	0	2	0.011403	0.076923	0.20
9992	0	0	0	0	0.001288	0.230769	0.00
9993	1	0	0	0	0.010722	0.076923	0.00

[9994 rows x 7 columns]

```
from sklearn.model_selection import train_test_split  
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.25,random_state=100)  
print(x_train.shape,x_test.shape,y_train.shape,y_test.shape)  
from sklearn.ensemble import RandomForestRegressor  
model=RandomForestRegressor(random_state=0,n_estimators=117)
```

(7495, 7) (2499, 7) (7495,) (2499,)

```
model=RandomForestRegressor(random_state=0,n_estimators=117)  
model.fit(x_train,y_train)
```

```
RandomForestRegressor  
RandomForestRegressor(n_estimators=117, random_state=0)
```

```
y_pred=model.predict(x_test)
```

```
from sklearn.metrics import r2_score  
print("R2 Score=",r2_score(y_pred,y_test)*100)
```

R2 Score= 74.08625172960967

```
fig = plt.figure()  
plt.scatter(y_test,y_pred)
```

plt.xlabel('y_test', fontsize=18)
plt.ylabel('Y_pred', fontsize=16)

