

Superstore Sales Analysis using Python

In [1]:

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import plotly.express as px
6 import plotly.graph_objects as go
7 import plotly.io as pio
8 import plotly.colors as colors
9
```

In [2]:

```
1 #Load the data from CSV file
2 # encoding = 'latin-1' this is particularly useful if the file contains special characters not represented
3 # in the default encoding(usually UTF-8).
4 # df.head() Display the first few rows of the dataset
5 df = pd.read_csv(r"F:\kailash\Project\python superstore sale\superstore sale dataset.csv", encoding = 'latin-1')
6 print(df.head())
7
```

	Row ID	Order ID	Order Date	Ship Date	Ship Mode	Customer ID	\
0	1	CA-2016-152156	11-08-2016	11-11-2016	Second Class	CG-12520	
1	2	CA-2016-152156	11-08-2016	11-11-2016	Second Class	CG-12520	
2	3	CA-2016-138688	06-12-2016	6/16/2016	Second Class	DV-13045	
3	4	US-2015-108966	10-11-2015	10/18/2015	Standard Class	SO-20335	
4	5	US-2015-108966	10-11-2015	10/18/2015	Standard Class	SO-20335	

	Customer Name	Segment	Country	City	...	\
0	Claire Gute	Consumer	United States	Henderson	...	
1	Claire Gute	Consumer	United States	Henderson	...	
2	Darrin Van Huff	Corporate	United States	Los Angeles	...	
3	Sean O'Donnell	Consumer	United States	Fort Lauderdale	...	
4	Sean O'Donnell	Consumer	United States	Fort Lauderdale	...	

	Postal Code	Region	Product ID	Category	Sub-Category	\
0	42420	South	FUR-BO-10001798	Furniture	Bookcases	
1	42420	South	FUR-CH-10000454	Furniture	Chairs	
2	90036	West	OFF-LA-10000240	Office Supplies	Labels	
3	33311	South	FUR-TA-10000577	Furniture	Tables	
4	33311	South	OFF-ST-10000760	Office Supplies	Storage	

	Product Name	Sales	Quantity	\
0	Bush Somerset Collection Bookcase	261.9600	2	
1	Hon Deluxe Fabric Upholstered Stacking Chairs,...	731.9400	3	
2	Self-Adhesive Address Labels for Typewriters b...	14.6200	2	
3	Bretford CR4500 Series Slim Rectangular Table	957.5775	5	
4	Eldon Fold 'N Roll Cart System	22.3680	2	

	Discount	Profit
0	0.00	41.9136
1	0.00	219.5820
2	0.00	6.8714
3	0.45	-383.0310
4	0.20	2.5164

[5 rows x 21 columns]

In [3]:

```
1 # Display Summary statistics
2 df.describe()
3
```

Out[3]:

	Row ID	Postal Code	Sales	Quantity	Discount	Profit
count	9994.000000	9994.000000	9994.000000	9994.000000	9994.000000	9994.000000
mean	4997.500000	55190.379428	229.858001	3.789574	0.156203	28.656896
std	2885.163629	32063.693350	623.245101	2.225110	0.206452	234.260108
min	1.000000	1040.000000	0.444000	1.000000	0.000000	-6599.978000
25%	2499.250000	23223.000000	17.280000	2.000000	0.000000	1.728750
50%	4997.500000	56430.500000	54.490000	3.000000	0.200000	8.666500
75%	7495.750000	90008.000000	209.940000	5.000000	0.200000	29.364000
max	9994.000000	99301.000000	22638.480000	14.000000	0.800000	8399.976000

In [4]:

```
1 # Display ALL Columns name
2 df.columns
3
```

Out[4]:

```
Index(['Row ID', 'Order ID', 'Order Date', 'Ship Date', 'Ship Mode',
      'Customer ID', 'Customer Name', 'Segment', 'Country', 'City', 'State',
      'Postal Code', 'Region', 'Product ID', 'Category', 'Sub-Category',
      'Product Name', 'Sales', 'Quantity', 'Discount', 'Profit'],
      dtype='object')
```

```
In [5]: 1 # Display basic information about the dataset
2 df.info()
3

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 9994 entries, 0 to 9993
Data columns (total 21 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Row ID                 9994 non-null   int64
1   Order ID               9994 non-null   object
2   Order Date             9994 non-null   object
3   Ship Date              9994 non-null   object
4   Ship Mode              9994 non-null   object
5   Customer ID            9994 non-null   object
6   Customer Name          9994 non-null   object
7   Segment                9994 non-null   object
8   Country                9994 non-null   object
9   City                   9994 non-null   object
10  State                  9994 non-null   object
11  Postal Code            9994 non-null   int64
12  Region                 9994 non-null   object
13  Product ID             9994 non-null   object
14  Category               9994 non-null   object
15  Sub-Category           9994 non-null   object
16  Product Name           9994 non-null   object
17  Sales                  9994 non-null   float64
18  Quantity               9994 non-null   int64
19  Discount               9994 non-null   float64
20  Profit                 9994 non-null   float64
dtypes: float64(3), int64(3), object(15)
memory usage: 1.6+ MB
```

```
In [6]: 1 # Display the number of unique values in each column
2 df.nunique()
3
```

Out[6]:

Row ID	9994
Order ID	5009
Order Date	1237
Ship Date	1334
Ship Mode	4
Customer ID	793
Customer Name	793
Segment	3
Country	1
City	531
State	49
Postal Code	631
Region	4
Product ID	1862
Category	3
Sub-Category	17
Product Name	1850
Sales	5825
Quantity	14
Discount	12
Profit	7287

dtype: int64

```
In [7]: 1 #Create new column Like order month,order year,order day
2
3 # Convert the 'Order Date' column to datetime format
4 df['Order Date'] = pd.to_datetime(df['Order Date'])
5 df['Ship Date'] = pd.to_datetime(df['Ship Date'])
6 # Create new columns for order month, year, and day
7 df['Order Day'] = df['Order Date'].dt.day
8 df['Order Month'] = df['Order Date'].dt.month
9 df['Order Year'] = df['Order Date'].dt.year
10
11
12 df['Ship Day'] = df['Ship Date'].dt.day
13 df['Ship Month'] = df['Ship Date'].dt.month
14 df['Ship Year'] = df['Ship Date'].dt.year
15
16
17 # Display the updated DataFrame
18 print(df[['Order Date', 'Order Day', 'Order Month', 'Order Year',
19           'Ship Date', 'Ship Day', 'Ship Month', 'Ship Year']])
20
21
```

	Order Date	Order Day	Order Month	Order Year	Ship Date	Ship Day	\
0	2016-11-08	8	11	2016	2016-11-11	11	
1	2016-11-08	8	11	2016	2016-11-11	11	
2	2016-06-12	12	6	2016	2016-06-16	16	
3	2015-10-11	11	10	2015	2015-10-18	18	
4	2015-10-11	11	10	2015	2015-10-18	18	
...	
9989	2014-01-21	21	1	2014	2014-01-23	23	
9990	2017-02-26	26	2	2017	2017-03-03	3	
9991	2017-02-26	26	2	2017	2017-03-03	3	
9992	2017-02-26	26	2	2017	2017-03-03	3	
9993	2017-05-04	4	5	2017	2017-05-09	9	

	Ship Month	Ship Year
0	11	2016
1	11	2016
2	6	2016
3	10	2015
4	10	2015
...
9989	1	2014
9990	3	2017
9991	3	2017
9992	3	2017
9993	5	2017

[9994 rows x 8 columns]

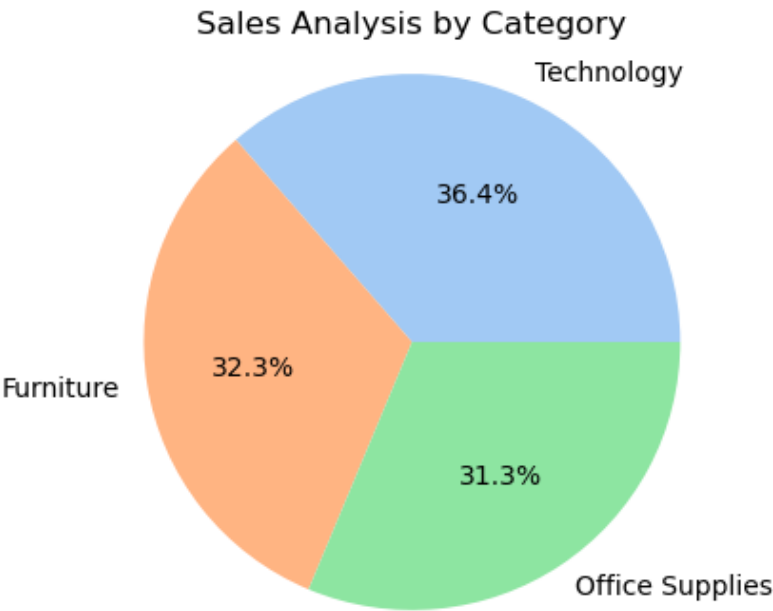
.

Before I start,I want to tell you.I have Created Two graph of each analysis,first using Matplotlib.pyplot library and second using plotly.express library

Analysis of sales by Category

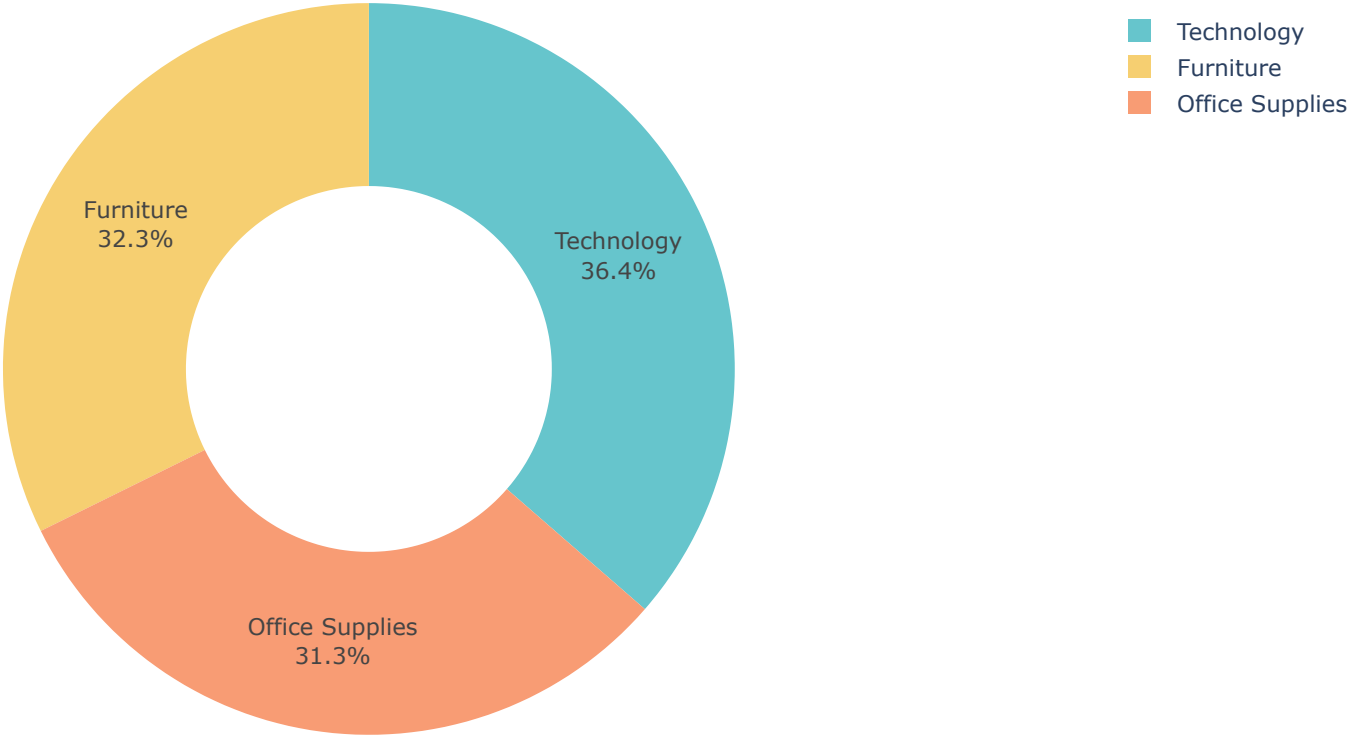
```
In [8]: 1 print("Using Seaborn and Matplotlib library")
2 # Group by category and sum the sales
3 category_sales = df.groupby('Category')['Sales'].sum().reset_index()
4
5 # Sort the results by sales in descending order
6 category_sales = category_sales.sort_values(by='Sales', ascending=False)
7
8 # Set the color palette
9 sns.set_palette("pastel")
10
11 # Create a pie chart
12 plt.figure(figsize=(4, 4))
13 plt.pie(category_sales['Sales'], labels=category_sales['Category'], autopct='%1.1f%%')
14 plt.title('Sales Analysis by Category')
15 plt.axis('equal') # Equal aspect ratio ensures that pie is drawn as a circle.
16
17 # Show the plot
18 plt.show()
19
20 #####
21
22 print("Using plotly.express library")
23
24 category_sales = df.groupby('Category')['Sales'].sum().reset_index()
25
26 fig = px.pie(category_sales,values='Sales',names='Category',hole=0.5,
27             color_discrete_sequence=px.colors.qualitative.Pastel)
28 fig.update_traces(textposition='inside', textinfo='percent+label')
29 fig.update_layout(title_text='Sales Analysis by Category', title_font=dict(size=14))
30 fig.show()
31
```

Using Seaborn and Matplotlib library



Using plotly.express library

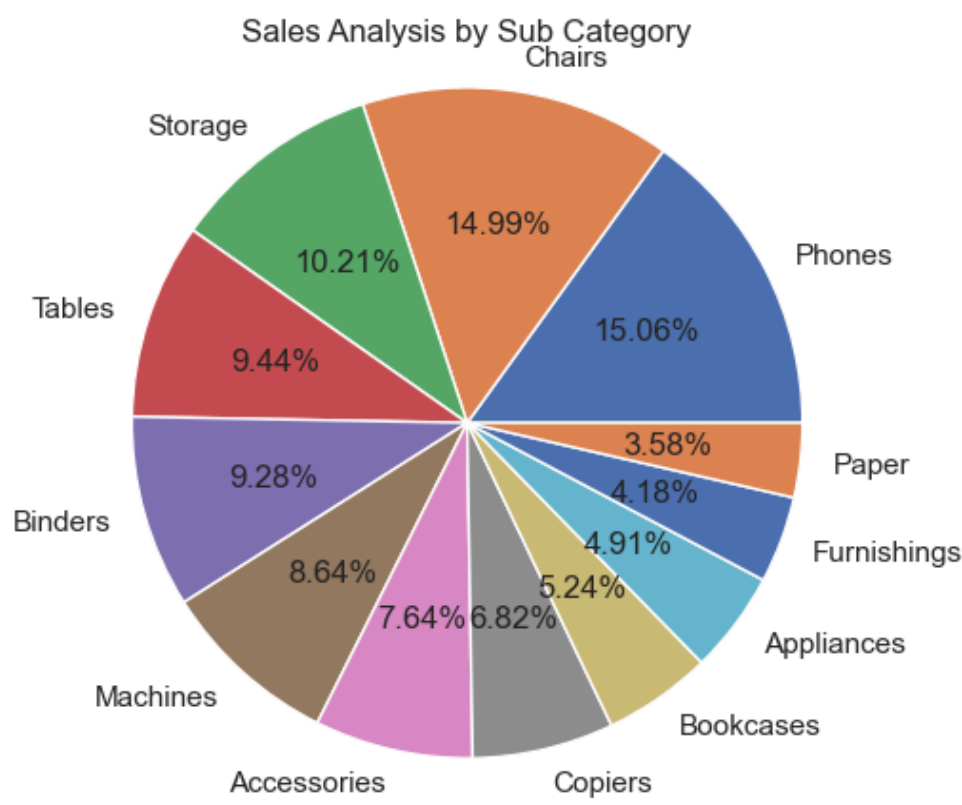
Sales Analysis by Category



Analysis of sales by SubCategory

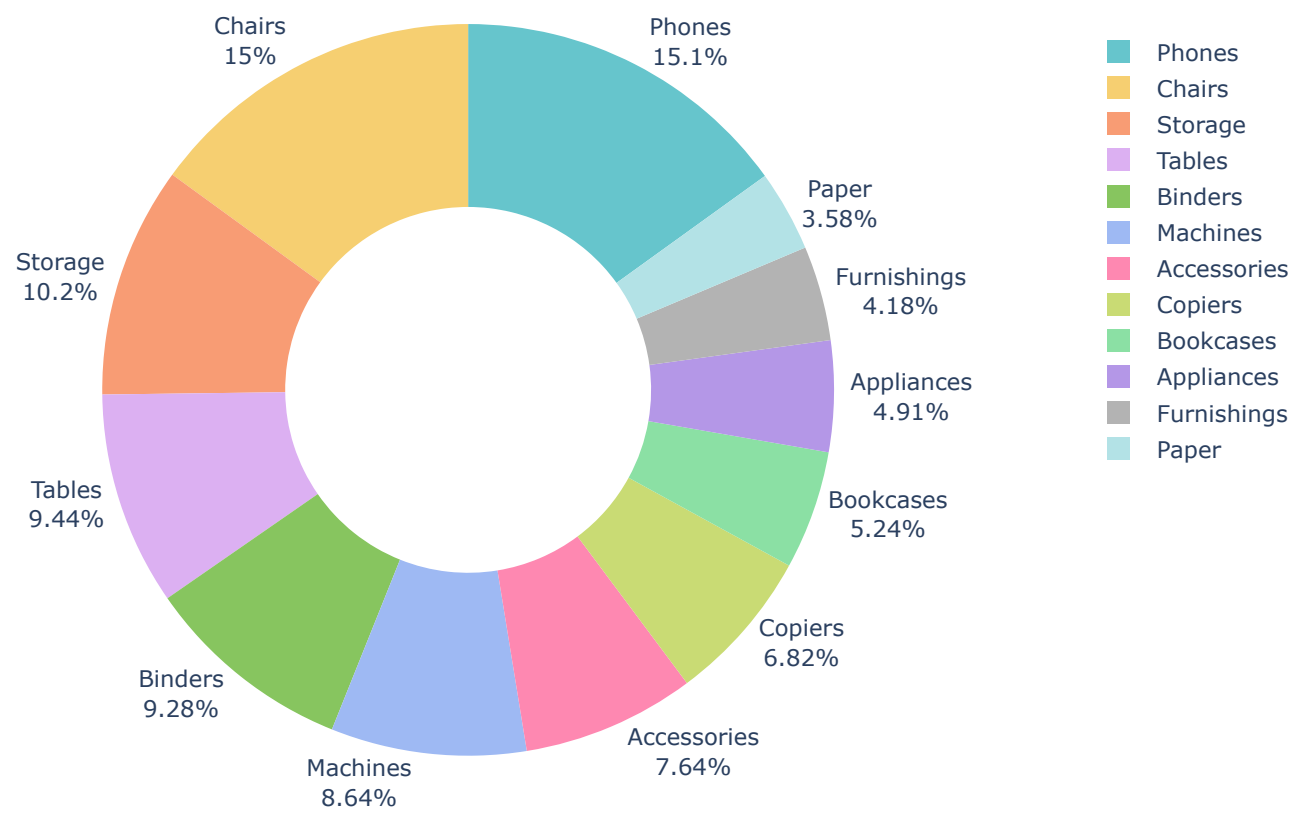
```
In [9]: 1 print("Using Seaborn and Matplotlib library")
2 # Group by Sub-Category and sum the sales
3 subcategory_sales = df.groupby('Sub-Category')['Sales'].sum().reset_index()
4
5 # Sort the results by sales in descending order
6 subcategory_sales = subcategory_sales.sort_values(by='Sales', ascending=False).head(12)
7
8 # Set the color palette
9 sns.set(style='whitegrid')
10
11 # Create a pie chart
12 plt.figure(figsize=(5, 5))
13 plt.pie(subcategory_sales['Sales'], labels=subcategory_sales['Sub-Category'], autopct='%0.2f%%')
14 plt.title('Sales Analysis by Sub Category')
15 plt.axis('equal') # Equal aspect ratio ensures that pie is drawn as a circle.
16
17 # Show the plot
18 plt.show()
19
20 #####
21 print("Using plotly.express library")
22
23 subcategory_sales = df.groupby('Sub-Category')['Sales'].sum().reset_index().nlargest(12,'Sales')
24
25 fig = px.pie(subcategory_sales,values='Sales',names='Sub-Category',hole=0.5,
26             color_discrete_sequence=px.colors.qualitative.Pastel)
27 fig.update_traces(textposition='outside', textinfo='percent+label')
28 fig.update_layout(title_text='Sales Analysis by SubCategory', title_font=dict(size=18))
29 fig.show()
30
```

Using Seaborn and Matplotlib library



Using plotly.express library

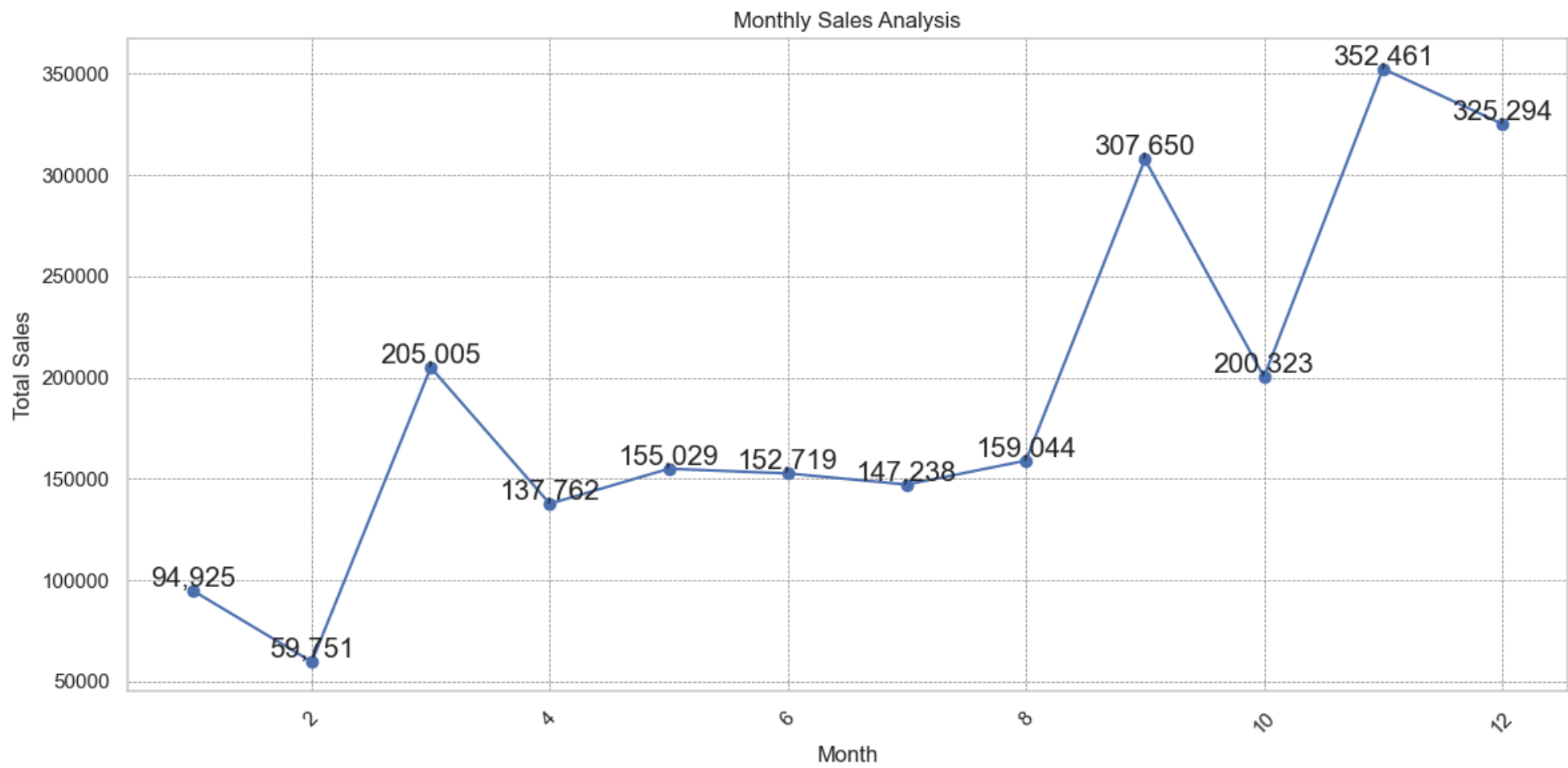
Sales Analysis by SubCategory



Monthly Sales Analysis

```
In [10]: 1 print("Using matplotlib.pyplot")
2
3 # Group by month and sum the sales
4 monthly_sales = df.groupby(['Order Month'])['Sales'].sum().reset_index()
5
6 # Plotting the monthly sales
7 plt.figure(figsize=(12, 6))
8 plt.plot(monthly_sales['Order Month'], monthly_sales['Sales'], marker='o')
9
10 # annotate each data point with its value
11 for i in range(len(monthly_sales)):
12     plt.text(monthly_sales['Order Month'][i],monthly_sales['Sales'][i],
13             f'{monthly_sales["Sales"][i]:,.0f}',fontsize=15, ha='center',va='bottom')
14
15 plt.title('Monthly Sales Analysis')
16 plt.xlabel('Month')
17 plt.ylabel('Total Sales')
18 plt.xticks(rotation=45)
19 plt.grid(color='gray',linestyle='--',linewidth=0.5)
20 plt.tight_layout()
21 plt.show()
22
23
24 #####
25 print("Using plotly.express")
26
27
28 monthly_sales=df.groupby('Order Month')['Sales'].sum().reset_index()
29 fig = px.line(monthly_sales, x = 'Order Month', y = 'Sales',
30             title = 'Monthly Sales Anallysis',markers=True)
31
32 #add data labels to each point
33 for i,row in monthly_sales.iterrows():
34     fig.add_annotation( x = row['Order Month'], y = row['Sales'],
35                     text = f'{row["Sales"]:,}.0f',ax = 0, ay = -10)
36 fig.show()
37
```

Using matplotlib.pyplot



Using plotly.express

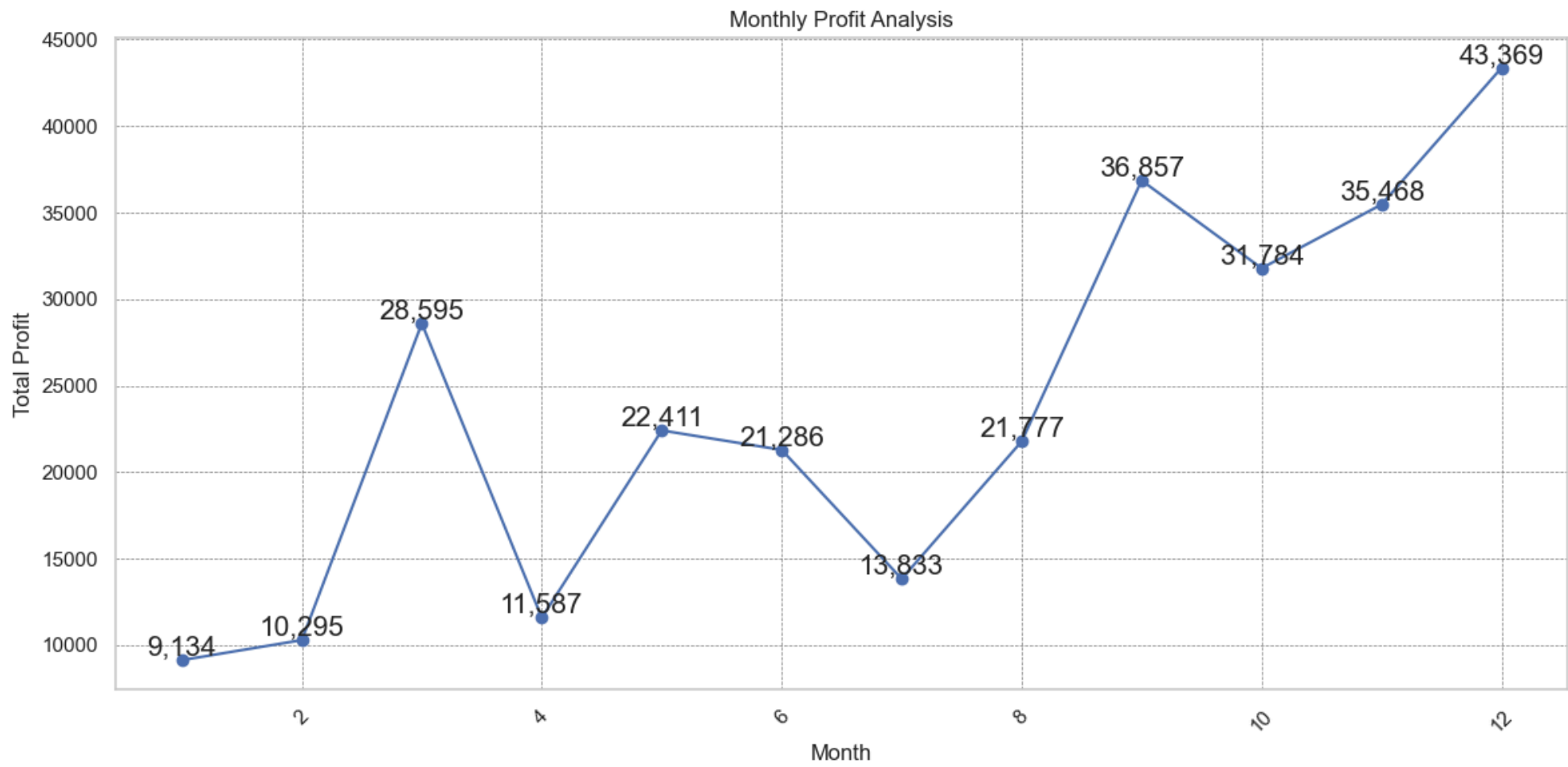
Monthly Sales Analysis



Monthly Profit Analysis

```
In [11]: 1 print("Using matplotlib.pyplot")
2
3 # Group by month and sum the Profit
4 monthly_profit = df.groupby(['Order Month'])['Profit'].sum().reset_index()
5
6 # Plotting the monthly Profit
7 plt.figure(figsize=(12, 6))
8 plt.plot(monthly_profit['Order Month'], monthly_profit['Profit'], marker='o')
9
10 # annotate each data point with its value
11 for i in range(len(monthly_profit)):
12     plt.text(monthly_profit['Order Month'][i],monthly_profit['Profit'][i],
13             f'{monthly_profit["Profit"][i]:,.0f}',fontsize=15, ha='center',va='bottom')
14
15
16 plt.title('Monthly Profit Analysis')
17 plt.xlabel('Month')
18 plt.ylabel('Total Profit')
19 plt.xticks(rotation=45)
20 plt.grid(color='gray',linestyle='--',linewidth=0.5)
21 plt.tight_layout()
22 plt.show()
23
24 #####
25 print("Using plotly.express")
26
27 monthly_profit=df.groupby('Order Month')['Profit'].sum().reset_index()
28 fig = px.line(monthly_profit, x = 'Order Month', y = 'Profit',
29             title = 'Monthly Profit Anallysis',markers=True)
30
31 #add data Labels to each point
32 for i,row in monthly_profit.iterrows():
33     fig.add_annotation( x = row['Order Month'], y = row['Profit'],
34                     text = f'{row["Profit"]:,.0f}',ax = 0, ay = -10)
35 fig.show()
36
```

Using matplotlib.pyplot



Using plotly.express

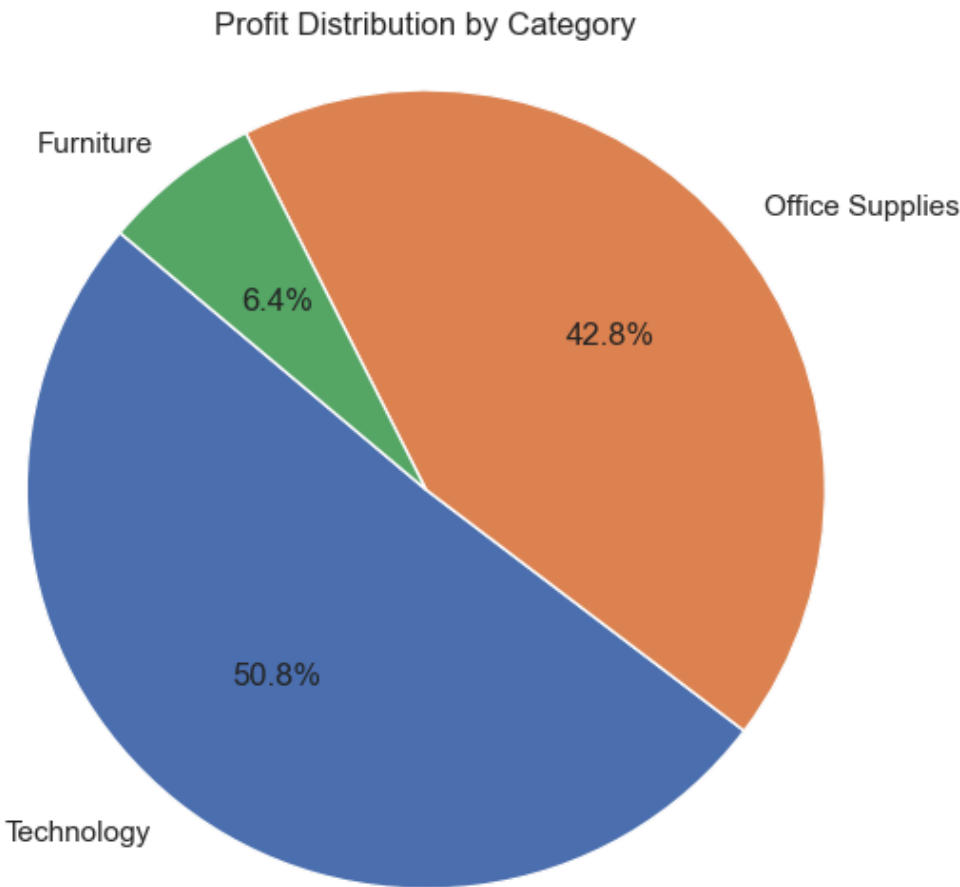
Monthly Profit Anallysis



Profit Analysis by Category

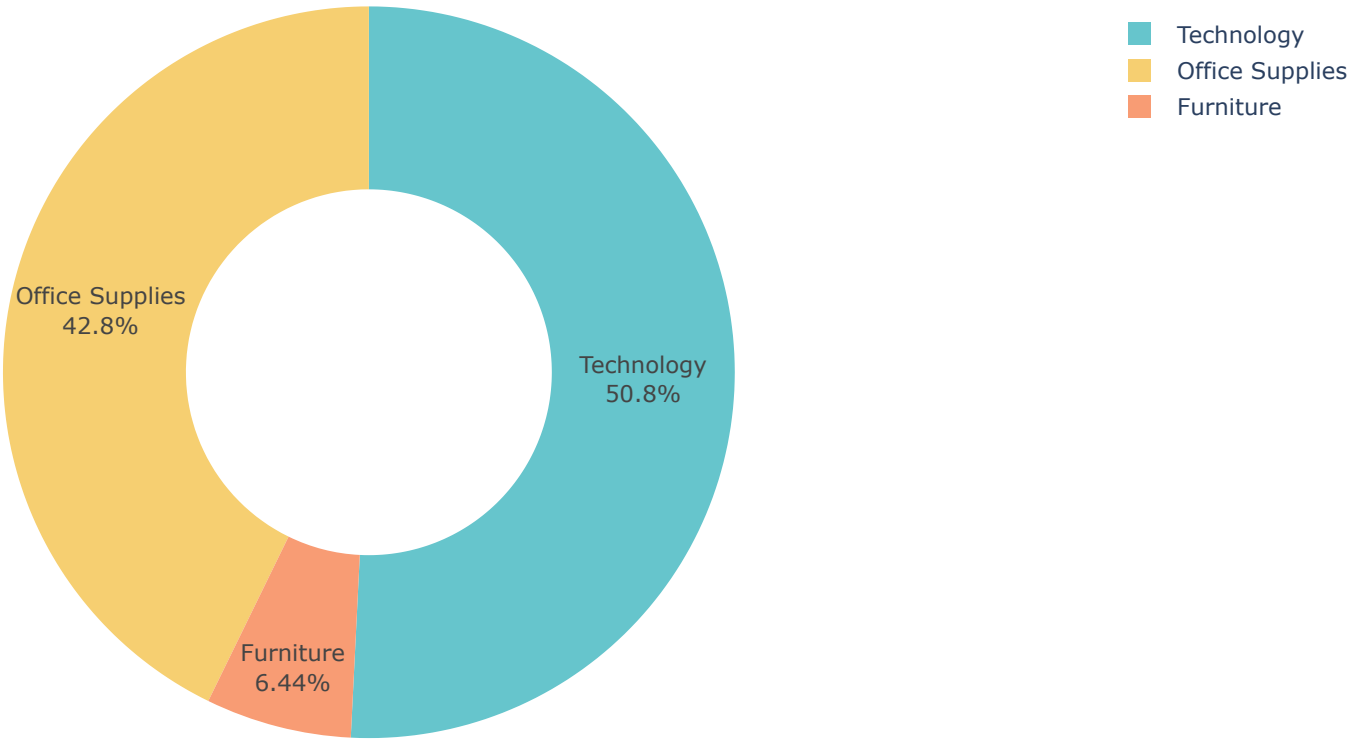
```
In [12]: 1 print("using matplotlib library")
2
3 # Group by category and sum the profits
4 category_profit = df.groupby('Category')['Profit'].sum().reset_index()
5
6 # Sort the results by profit in descending order
7 category_profit = category_profit.sort_values(by='Profit', ascending=False)
8
9 # Create a pie chart
10 plt.figure(figsize=(7, 6))
11 plt.pie(category_profit['Profit'], labels=category_profit['Category'], autopct='%1.1f%%', startangle=140)
12 plt.title('Profit Distribution by Category')
13 plt.axis('equal') # Equal aspect ratio ensures that pie is drawn as a circle.
14
15 # Show the plot
16 plt.show()
17
18
19
20 #####
21 print("using plotly.express library")
22
23 # Group by category and sum the profits
24 category_profit = df.groupby('Category')['Profit'].sum().reset_index()
25
26 # Create a pie chart
27 fig= px.pie(category_profit,values='Profit', names='Category',hole=0.5,
28             color_discrete_sequence=px.colors.qualitative.Pastel)
29 fig.update_traces(textposition='inside', textinfo='percent+label')
30 fig.update_layout(title_text='Profit Analysis by Category', title_font=dict(size=18))
31
32 # Show the plot
33 fig.show()
34
```

using matplotlib library



using plotly.express library

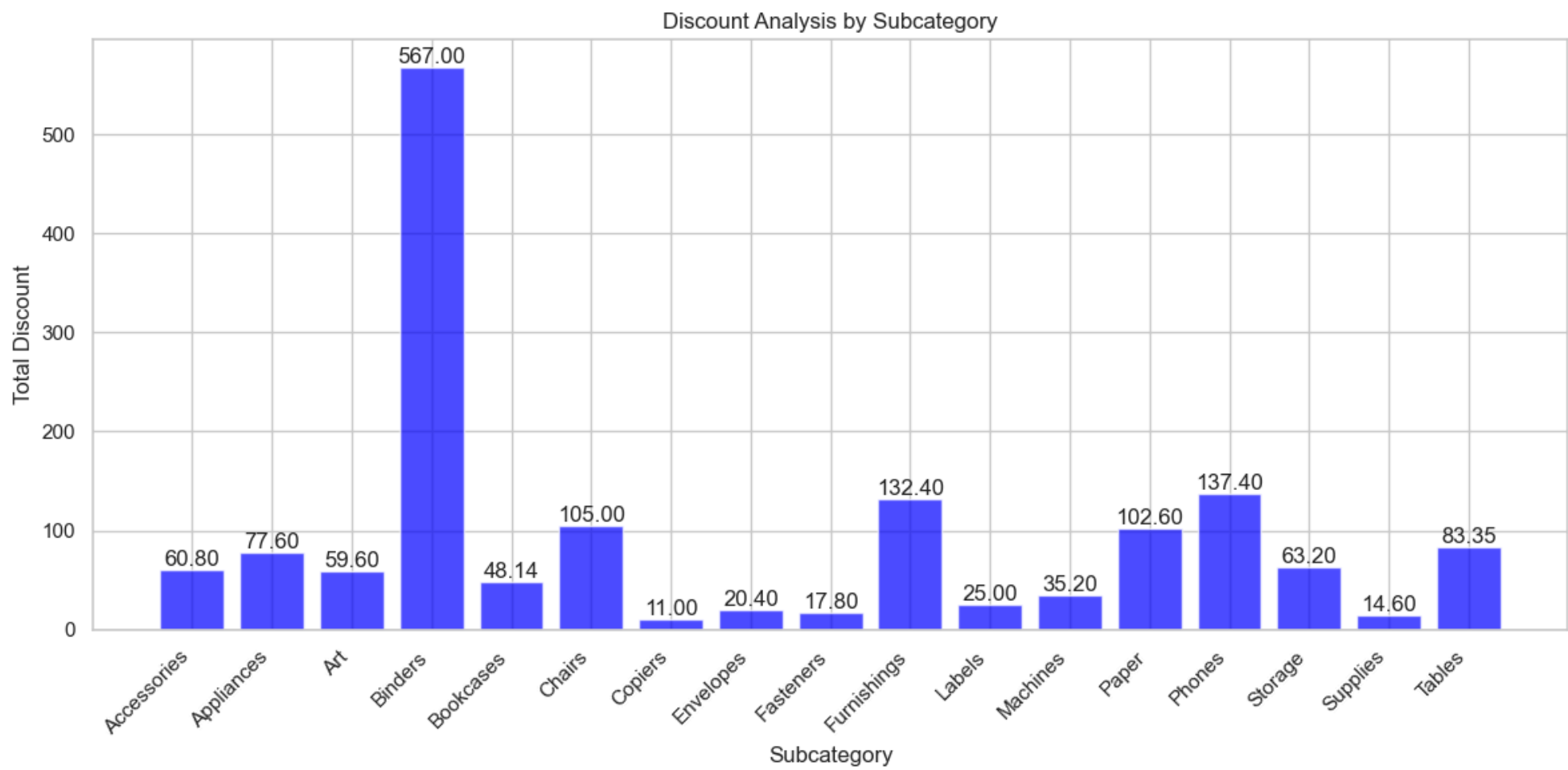
Profit Analysis by Category



Discount Analysis by Subcategory

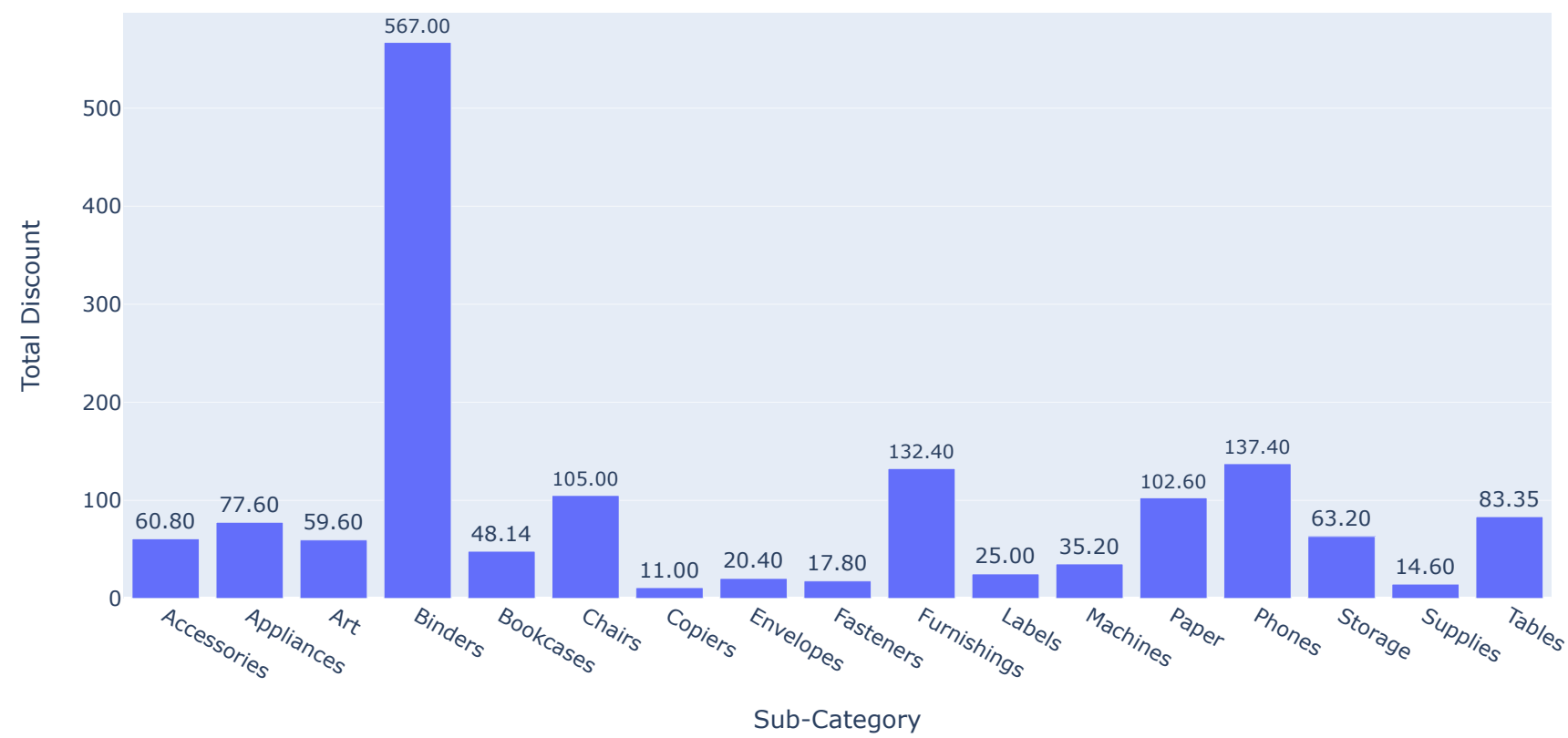
```
In [13]: 1 print("Using Matplotlib library")
2
3 # Group by Subcategory and sum the Discounts
4 subcategory_analysis = df.groupby('Sub-Category')['Discount'].sum().reset_index()
5
6 # Set the figure size
7 plt.figure(figsize=(12, 6))
8
9 # Create a bar chart
10 bars = plt.bar(subcategory_analysis['Sub-Category'], subcategory_analysis['Discount'], color='blue', alpha=0.7)
11
12 # Add titles and Labels
13 plt.title('Discount Analysis by Subcategory')
14 plt.xlabel('Subcategory')
15 plt.ylabel('Total Discount')
16
17 # Annotate the bars with their values without a loop
18 plt.bar_label(bars, label_type='edge', fmt='%.2f')
19
20 # Rotate x-ticks for better readability
21 plt.xticks(rotation=45, ha='right')
22
23 # Show the plot
24 plt.tight_layout()
25 plt.show()
26
27
28 #####
29 print("Using Plotly.Express library")
30
31 # Group by Subcategory and sum the Discounts
32 subcategory_analysis = df.groupby('Sub-Category')['Discount'].sum().reset_index()
33
34 # Create a bar chart using Plotly Express
35 fig = px.bar(subcategory_analysis, x='Sub-Category', y='Discount',
36              title='Discount Analysis by Subcategory',
37              labels={'Discount': 'Total Discount'})
38
39 # Add data labels to each bar
40 fig.update_traces(texttemplate='%{y:,.2f}', textposition='outside')
41
42 # Show the plot
43 fig.show()
44
```

Using Matplotlib library



Using Plotly.Express library

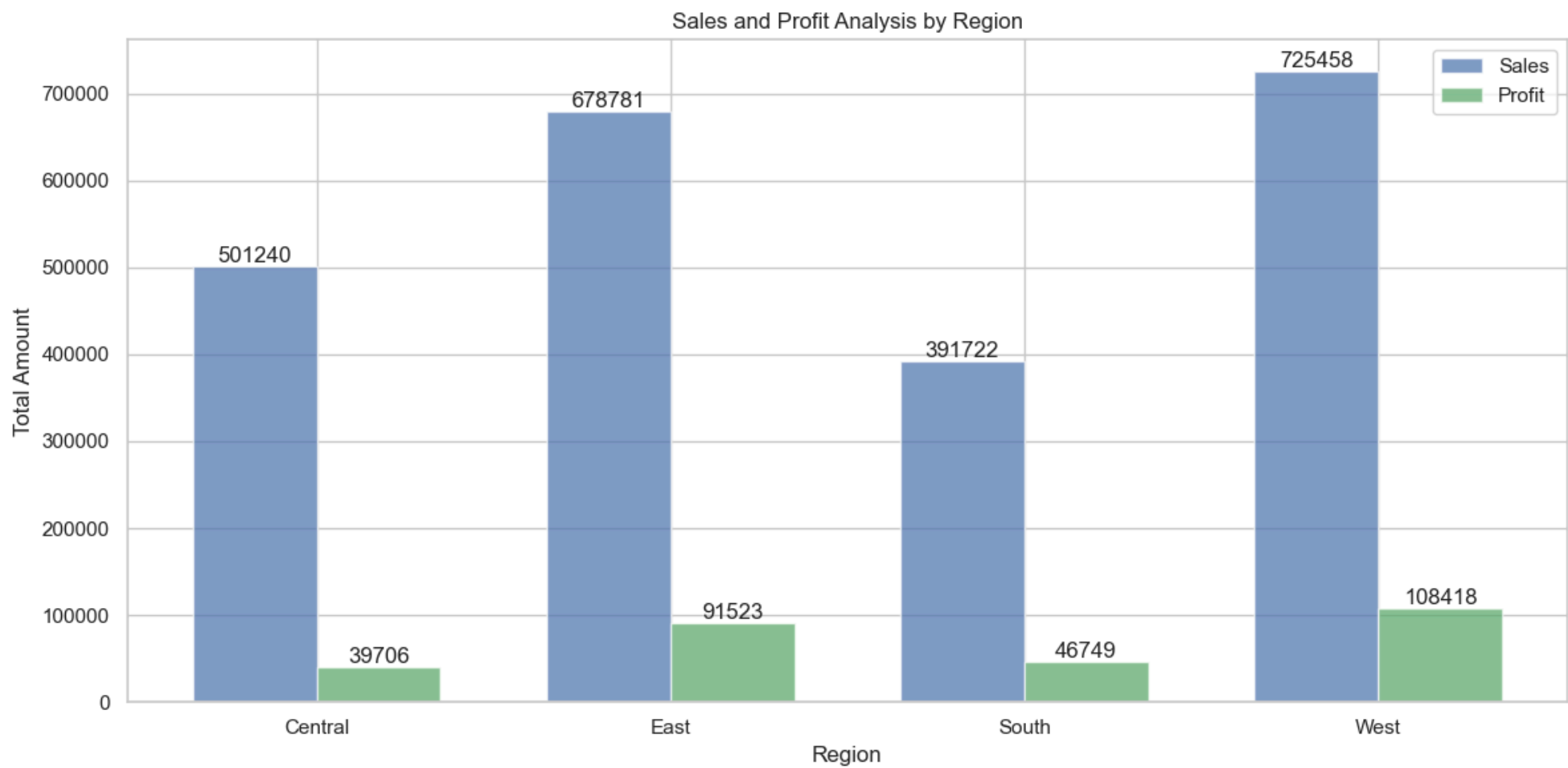
Discount Analysis by Subcategory



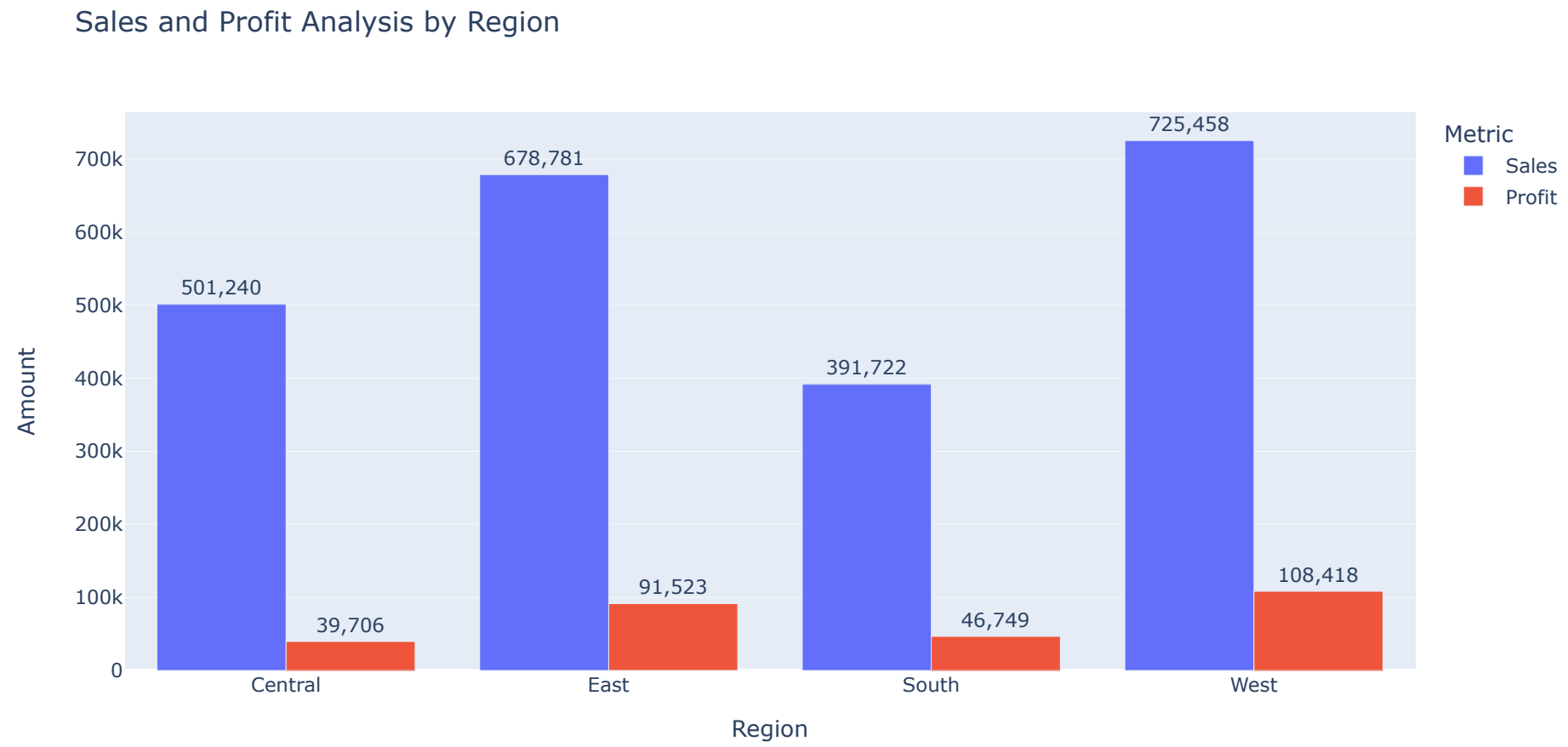
Sales and Profit Analysis by Region

```
In [14]: 1 print("Using Matplotlib library")
2
3 # Group by Region and sum the Sales and Profit
4 region_analysis = df.groupby('Region')[['Sales', 'Profit']].sum().reset_index()
5
6 # Set the figure size
7 plt.figure(figsize=(12, 6))
8
9 # Create bar plots for Sales and Profit
10 bar_width = 0.35
11 x = range(len(region_analysis))
12
13 # Create bars for Sales
14 sales_bars = plt.bar(x, region_analysis['Sales'], width=bar_width, label='Sales',
15                     color='b', alpha=0.7)
16
17 # Create bars for Profit
18 profit_bars = plt.bar([p + bar_width for p in x], region_analysis['Profit'],
19                      width=bar_width, label='Profit', color='g', alpha=0.7)
20
21 # Set the x-ticks and Labels
22 plt.xticks([p + bar_width / 2 for p in x], region_analysis['Region'])
23
24 # Add titles and Labels
25 plt.title('Sales and Profit Analysis by Region')
26 plt.xlabel('Region')
27 plt.ylabel('Total Amount')
28 plt.legend()
29
30 #Annotate the bars with values
31 plt.bar_label(sales_bars,label_type='edge',fmt='%.0f')
32 plt.bar_label(profit_bars,label_type='edge',fmt='%.0f')
33 # Show the plot
34 plt.tight_layout()
35 plt.show()
36
37 #####
38 print("Using Plotly.Express library")
39
40 # Group by Region and sum the Sales and Profit
41 region_analysis = df.groupby('Region')[['Sales', 'Profit']].sum().reset_index()
42
43 # Melt the DataFrame for Plotly
44 region_melted = region_analysis.melt(id_vars='Region', value_vars=['Sales', 'Profit'],
45                                   var_name='Metric', value_name='Amount')
46
47 # Create a bar chart using Plotly Express
48 fig = px.bar(region_melted, x='Region', y='Amount', color='Metric',
49             title='Sales and Profit Analysis by Region',
50             barmode='group')
51
52 #Add data Labels to each bar
53 fig.update_traces(texttemplate='%{y:,.0f}',textposition='outside')
54 # Show the plot
55 fig.show()
56
```

Using Matplotlib library



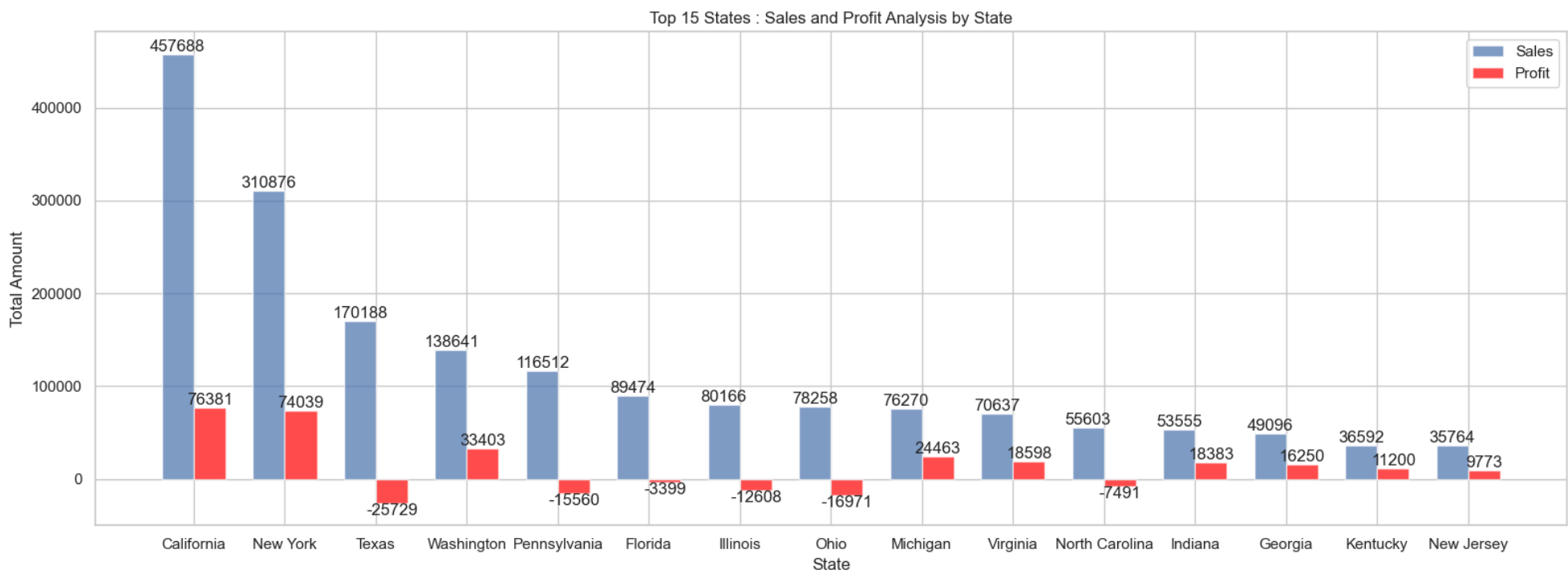
Using Plotly.Express library



Sales and Profit Analysis by Statewise

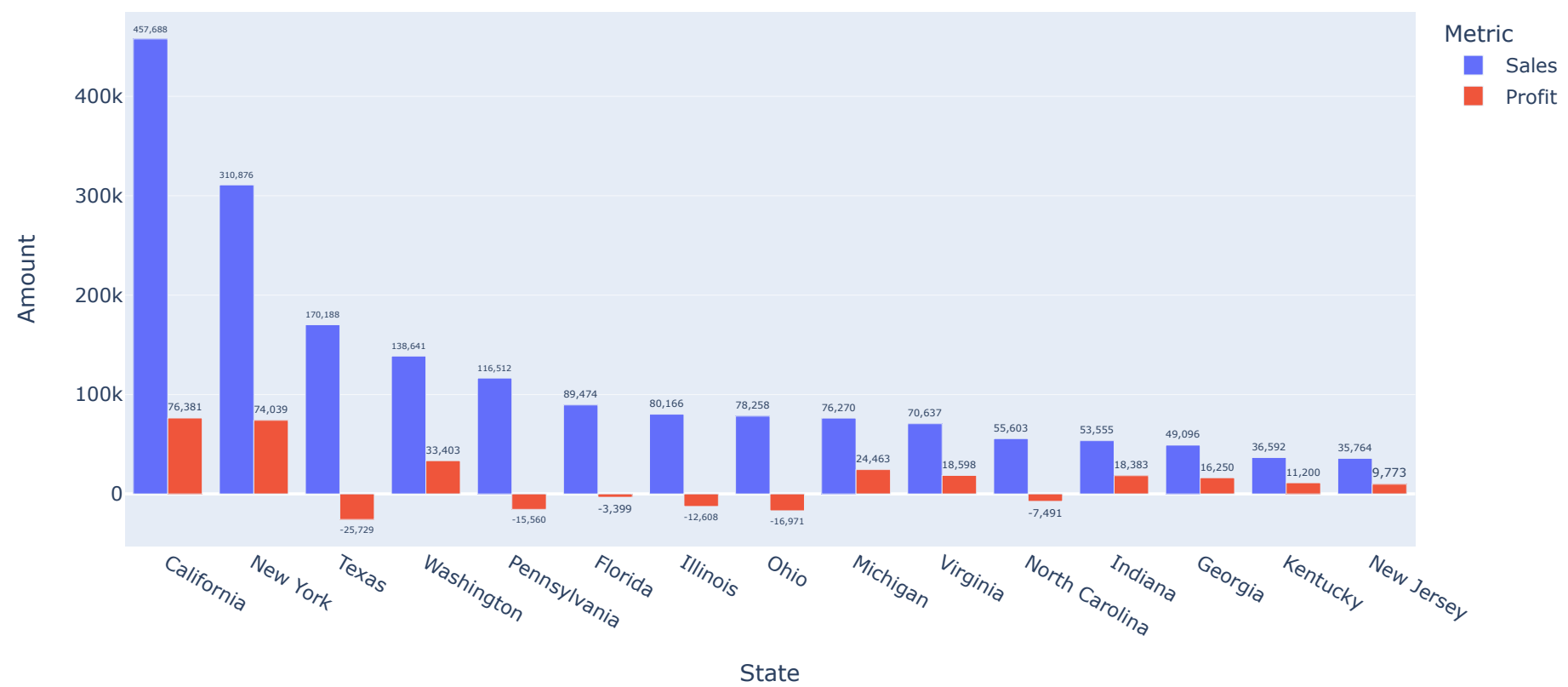
```
In [15]: 1 print("Using Matplotlib library")
2
3 # Group by State and sum the Sales and Profit
4 state_analysis = df.groupby('State')[['Sales', 'Profit']].sum().reset_index()
5
6 #Sort by sales to get the top 15 States
7 top_states = state_analysis.sort_values(by='Sales', ascending = False).head(15)
8
9 # Set the figure size
10 plt.figure(figsize=(16, 6))
11
12 # Set bar width and positions
13 bar_width = 0.35
14 x = range(len(top_states))
15
16 # Create bars for Sales
17 sales_bars = plt.bar(x, top_states['Sales'], width=bar_width, label='Sales', color='b', alpha=0.7)
18
19 # Create bars for Profit
20 profit_bars = plt.bar([p + bar_width for p in x], top_states['Profit'], width=bar_width, label='Profit',
21                       color='red', alpha=0.7)
22
23 # Set the x-ticks and Labels
24 plt.xticks([p + bar_width / 2 for p in x], top_states['State'])
25
26 # Add titles and Labels
27 plt.title('Top 15 States : Sales and Profit Analysis by State')
28 plt.xlabel('State')
29 plt.ylabel('Total Amount')
30 plt.legend()
31
32 # Annotate the bars with values
33 plt.bar_label(sales_bars, label_type='edge', fmt='%.0f')
34 plt.bar_label(profit_bars, label_type='edge', fmt='%.0f')
35
36 # Show the plot
37 plt.tight_layout()
38 plt.show()
39
40 #####
41 print("Using plotly.Express library")
42
43 # Group by State and sum the Sales and Profit
44 state_analysis = df.groupby('State')[['Sales', 'Profit']].sum().reset_index()
45
46 #Sort by sales to get the top 15 States
47 top_states = state_analysis.sort_values(by='Sales', ascending = False).head(15)
48
49 # Melt the DataFrame for Plotly
50 top_states_melted = top_states.melt(id_vars='State', value_vars=['Sales', 'Profit'],
51                                   var_name='Metric', value_name='Amount')
52
53 # Create a bar chart
54 fig = px.bar(top_states_melted, x='State', y='Amount', color='Metric',
55             title='Top 15 States : Sales and Profit Analysis by State',
56             barmode='group')
57
58 # Add data labels to each bar
59 fig.update_traces(texttemplate='%{y:,.0f}', textposition='outside')
60
61 # Show the plot
62 fig.show()
63
```

Using Matplotlib library



Using plotly.Express library

Top 15 States : Sales and Profit Analysis by State



summarizing key insights from a sales dataset

1. Category Sales Contribution:
- Technology has the highest sales at 36.4%.
 - The highest sales contribution in the subcategory is from Phones, accounting for 15.06% of sales.
2. Monthly Sales Performance:
- The months of November and December show the highest sales contributions.
3. Profit Insights:
- The highest profit occurs from September to December.
 - The peak profit of ₹43,369 is recorded in December.
 - Technology leads in profit contribution at 50.8%, followed by Office Supplies at 42.8%, and Furniture at 6.4%.
4. Discount Contributions:
- Binders have the maximum discount contribution at ₹567.
5. Regional Sales and Profit:
- The West Region contributes the most to sales and profit, totaling ₹725,458, followed by the East, Central, and South Regions (least at ₹391,722).
6. State Contributions:
- Among the top 15 states, California tops with sales of ₹457,688, followed closely by New York and Texas.
 - Profit Contribution are as follows: California: ₹76,381 New York: ₹74,039 Texas shows a loss of -₹25,729.

Conclusion

This analysis highlights significant trends in sales and profit across various categories, months, regions, and states. Technology consistently outperforms in both sales and profit, while the West Region dominates overall contributions. Understanding these dynamics can guide strategic decisions in inventory, marketing, and regional focus for improved performance.