

```
In [1]: import pandas as pd
```

```
In [23]: import pandas as pd
```

```
# Use 'latin1' or 'ISO-8859-1' to fix encoding error
df = pd.read_csv(r"C:\Users\sunny\SampleSuperstore.csv", encoding='latin1')
df.head()
```

```
Out[23]:
```

| | Row ID | Order ID | Order Date | Ship Date | Ship Mode | Customer ID | Customer Name | Segment | Co |
|--|--------|----------|------------|-----------|-----------|-------------|---------------|---------|----|
|--|--------|----------|------------|-----------|-----------|-------------|---------------|---------|----|

| | | | | | | | | | |
|---|---|----------------|-----------|------------|--------------|----------|-------------|----------|---|
| 0 | 1 | CA-2016-152156 | 11/8/2016 | 11/11/2016 | Second Class | CG-12520 | Claire Gute | Consumer | l |
|---|---|----------------|-----------|------------|--------------|----------|-------------|----------|---|

| | | | | | | | | | |
|---|---|----------------|-----------|------------|--------------|----------|-------------|----------|---|
| 1 | 2 | CA-2016-152156 | 11/8/2016 | 11/11/2016 | Second Class | CG-12520 | Claire Gute | Consumer | l |
|---|---|----------------|-----------|------------|--------------|----------|-------------|----------|---|

| | | | | | | | | | |
|---|---|----------------|-----------|-----------|--------------|----------|-----------------|-----------|---|
| 2 | 3 | CA-2016-138688 | 6/12/2016 | 6/16/2016 | Second Class | DV-13045 | Darrin Van Huff | Corporate | l |
|---|---|----------------|-----------|-----------|--------------|----------|-----------------|-----------|---|

| | | | | | | | | | |
|---|---|----------------|------------|------------|----------------|----------|----------------|----------|---|
| 3 | 4 | US-2015-108966 | 10/11/2015 | 10/18/2015 | Standard Class | SO-20335 | Sean O'Donnell | Consumer | l |
|---|---|----------------|------------|------------|----------------|----------|----------------|----------|---|

| | | | | | | | | | |
|---|---|----------------|------------|------------|----------------|----------|----------------|----------|---|
| 4 | 5 | US-2015-108966 | 10/11/2015 | 10/18/2015 | Standard Class | SO-20335 | Sean O'Donnell | Consumer | l |
|---|---|----------------|------------|------------|----------------|----------|----------------|----------|---|

5 rows × 21 columns

```
In [25]: print(df.info())
print(df.describe())
print(df.isnull().sum())
```

```
<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
```

```
None
```

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

| | Profit |
|-------|--------------|
| count | 9994.000000 |
| mean | 28.656896 |
| std | 234.260108 |
| min | -6599.978000 |
| 25% | 1.728750 |
| 50% | 8.666500 |
| 75% | 29.364000 |
| max | 8399.976000 |

| | |
|---------------|---|
| Row ID | 0 |
| Order ID | 0 |
| Order Date | 0 |
| Ship Date | 0 |
| Ship Mode | 0 |
| Customer ID | 0 |
| Customer Name | 0 |
| Segment | 0 |
| Country | 0 |
| City | 0 |
| State | 0 |
| Postal Code | 0 |

```
Region ..... 0
Product ID ..... 0
Category ..... 0
Sub-Category ..... 0
Product Name ..... 0
Sales ..... 0
Quantity ..... 0
Discount ..... 0
Profit ..... 0
dtype: int64
```

```
In [27]: df = df.drop_duplicates()
```

```
In [29]: df
```

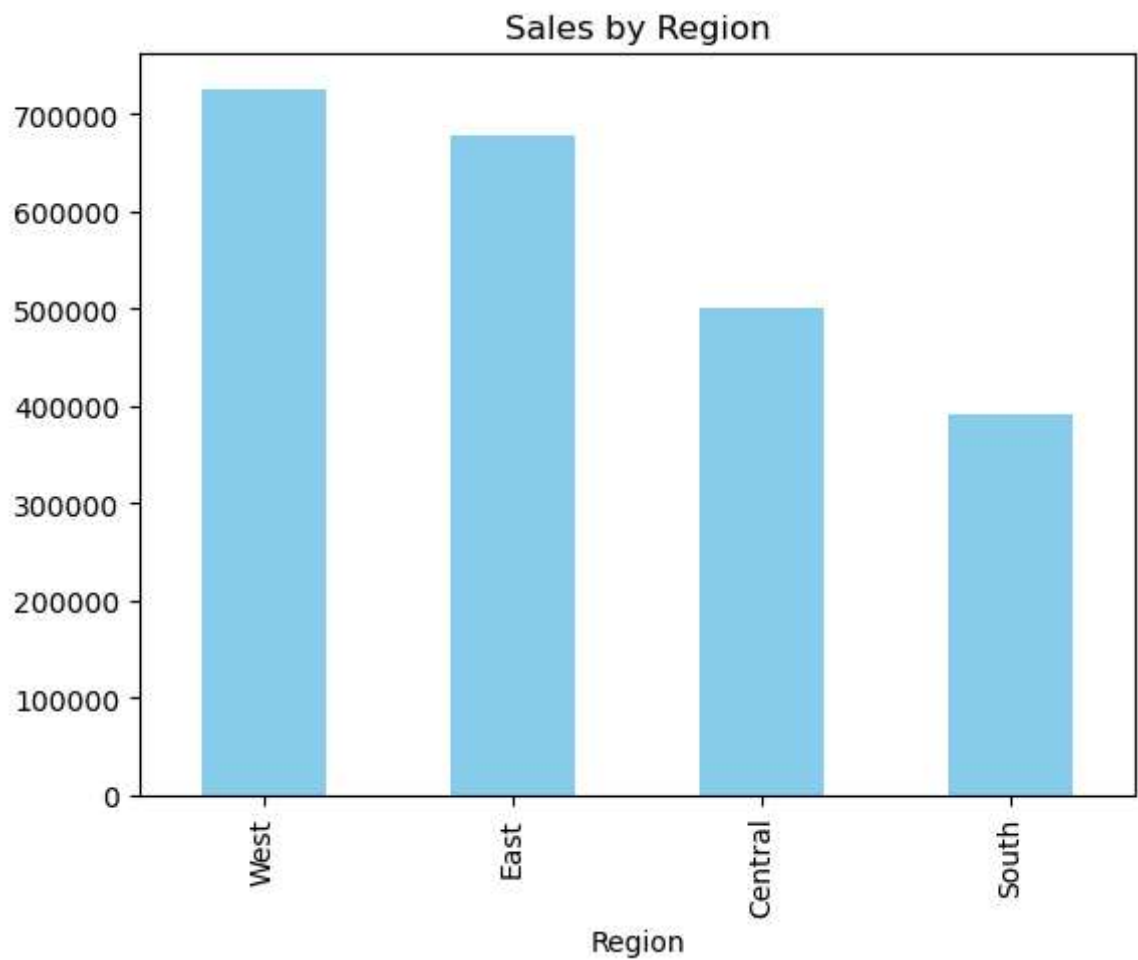
Out[29]:

| | Row ID | Order ID | Order Date | Ship Date | Ship Mode | Customer ID | Customer Name | Segment |
|-------------|--------|----------------|------------|------------|----------------|-------------|------------------|-----------|
| 0 | 1 | CA-2016-152156 | 11/8/2016 | 11/11/2016 | Second Class | CG-12520 | Claire Gute | Consumer |
| 1 | 2 | CA-2016-152156 | 11/8/2016 | 11/11/2016 | Second Class | CG-12520 | Claire Gute | Consumer |
| 2 | 3 | CA-2016-138688 | 6/12/2016 | 6/16/2016 | Second Class | DV-13045 | Darrin Van Huff | Corporate |
| 3 | 4 | US-2015-108966 | 10/11/2015 | 10/18/2015 | Standard Class | SO-20335 | Sean O'Donnell | Consumer |
| 4 | 5 | US-2015-108966 | 10/11/2015 | 10/18/2015 | Standard Class | SO-20335 | Sean O'Donnell | Consumer |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 9989 | 9990 | CA-2014-110422 | 1/21/2014 | 1/23/2014 | Second Class | TB-21400 | Tom Boeckenhauer | Consumer |
| 9990 | 9991 | CA-2017-121258 | 2/26/2017 | 3/3/2017 | Standard Class | DB-13060 | Dave Brooks | Consumer |
| 9991 | 9992 | CA-2017-121258 | 2/26/2017 | 3/3/2017 | Standard Class | DB-13060 | Dave Brooks | Consumer |
| 9992 | 9993 | CA-2017-121258 | 2/26/2017 | 3/3/2017 | Standard Class | DB-13060 | Dave Brooks | Consumer |
| 9993 | 9994 | CA-2017-119914 | 5/4/2017 | 5/9/2017 | Second Class | CC-12220 | Chris Cortes | Consumer |

9994 rows × 21 columns

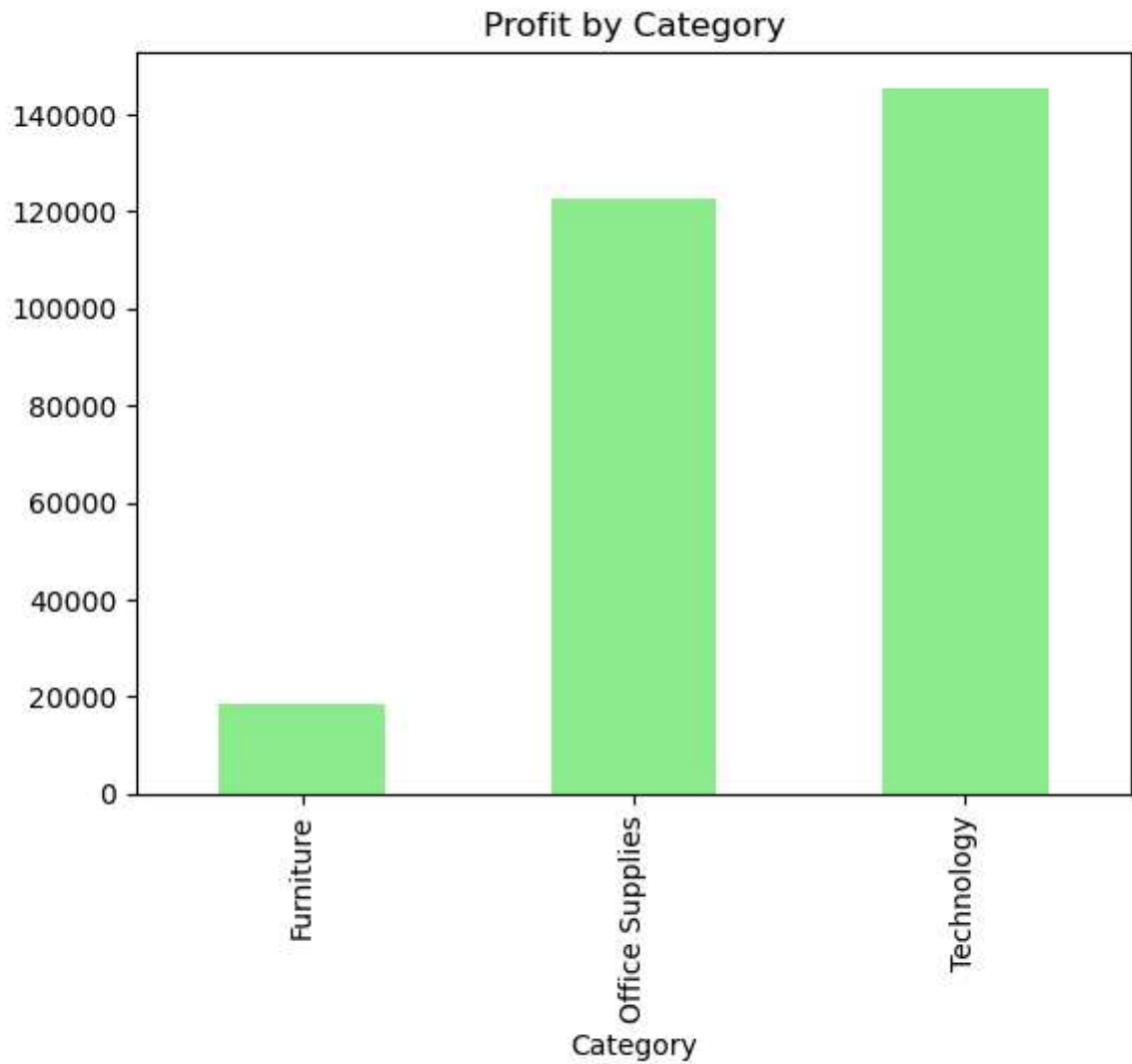
```
In [31]: region_sales = df.groupby('Region')['Sales'].sum().sort_values(ascending=False)
region_sales.plot(kind='bar', title='Sales by Region', color='skyblue')
```

Out[31]: <Axes: title={'center': 'Sales by Region'}, xlabel='Region'>



```
In [33]: category_profit = df.groupby('Category')['Profit'].sum()
category_profit.plot(kind='bar', title='Profit by Category', color='lightgreen')
```

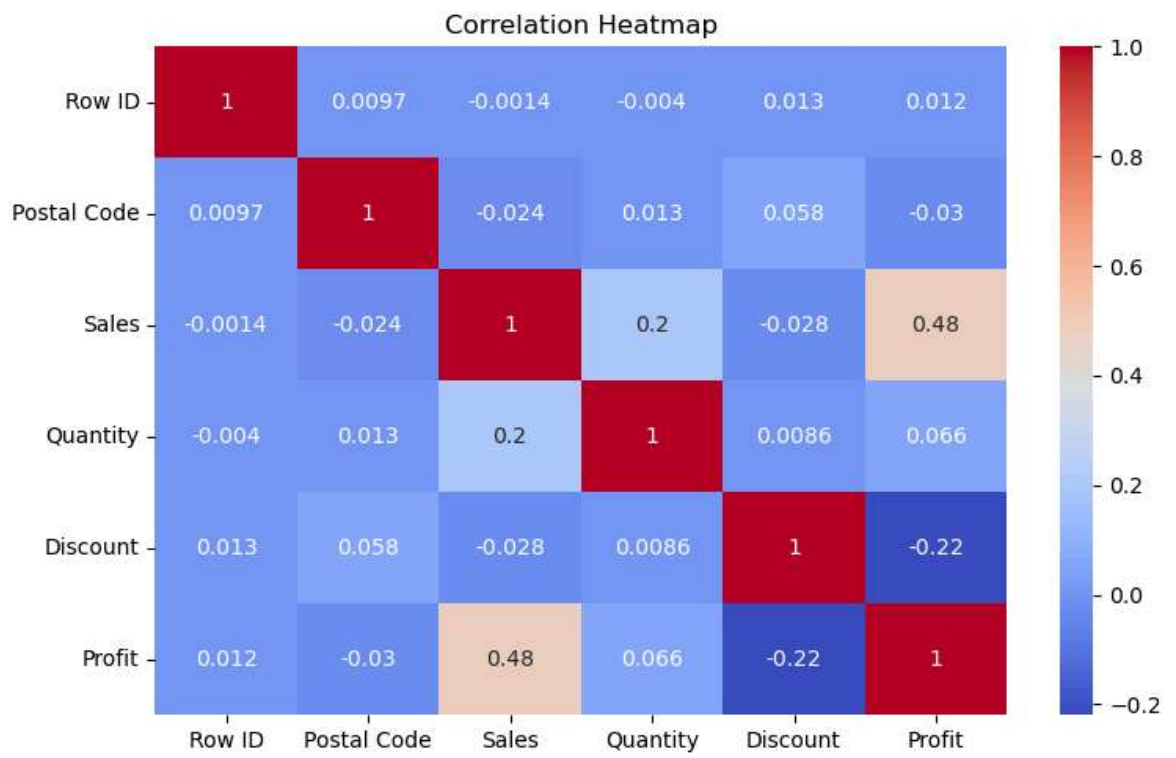
Out[33]: <Axes: title={'center': 'Profit by Category'}, xlabel='Category'>



```
In [37]: import seaborn as sns
import matplotlib.pyplot as plt

# ☒ Select only numerical columns for correlation
numeric_df = df.select_dtypes(include=['number'])

plt.figure(figsize=(8,5))
sns.heatmap(numeric_df.corr(), annot=True, cmap='coolwarm')
plt.title("Correlation Heatmap")
plt.tight_layout()
plt.savefig("correlation_heatmap.png")
plt.show()
```



In []: