

Innobyte Services Internship Task

Name : PALLAVI PATIL

Amazon Sales Analysis

Objective : To Analyze and Provide Insights on Amazon Sales Report

Libraries

```
In [47]: 1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import calendar
```

Dataset

```
In [3]: 1 df = pd.read_csv("C:\\Users\\Dell\\Downloads\\Amazon_sales.csv")
```

1.Shape

```
In [4]: 1 df.shape
```

Out[4]: (128976, 21)

```
In [5]: 1 df.head(3)
```

Out[5]:

	index	Order ID	Date	Status	Fulfilment	Sales Channel	ship-service-level	Category	Size	Courier Status	...	currency	Amount
0	0	405-8078784-5731545	04-30-22	Cancelled	Merchant	Amazon.in	Standard	T-shirt	S	On the Way	...	INR	647.62
1	1	171-9198151-1101146	04-30-22	Shipped - Delivered to Buyer	Merchant	Amazon.in	Standard	Shirt	3XL	Shipped	...	INR	406.00
2	2	404-0687676-7273146	04-30-22	Shipped	Amazon	Amazon.in	Expedited	Shirt	XL	Shipped	...	INR	329.00

3 rows × 21 columns

In [6]:

1df.tail(3)

Out[6]:

	index	Order ID	Date	Status	Fulfilment	Sales Channel	ship-service-level	Category	Size	Courier Status	...	currency	Amo
128973	128972	407-9547469-3152358	05-31-22	Shipped	Amazon	Amazon.in	Expedited	Blazzer	XXL	Shipped	...	INR	69
128974	128973	402-6184140-0545956	05-31-22	Shipped	Amazon	Amazon.in	Expedited	T-shirt	XS	Shipped	...	INR	119
128975	128974	408-7436540-8728312	05-31-22	Shipped	Amazon	Amazon.in	Expedited	T-shirt	S	Shipped	...	INR	69

3 rows × 21 columns

2. Info

In [7]:

1df.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
```

3. Describe

In [8]: 1 df.describe()

Out[8]:

	index	Qty	Amount	ship-postal-code	New	PendingS
count	128976.000000	128976.000000	121176.000000	128941.000000	0.0	0.0
mean	64486.130427	0.904401	648.562176	463945.677744	NaN	NaN
std	37232.897832	0.313368	281.185041	191458.488954	NaN	NaN
min	0.000000	0.000000	0.000000	110001.000000	NaN	NaN
25%	32242.750000	1.000000	449.000000	382421.000000	NaN	NaN
50%	64486.500000	1.000000	605.000000	500033.000000	NaN	NaN
75%	96730.250000	1.000000	788.000000	600024.000000	NaN	NaN
max	128974.000000	15.000000	5584.000000	989898.000000	NaN	NaN

4. Missing Values

In [9]: 1 df.isna().sum()

Out[9]:

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
New	128976
PendingS	128976

dtype: int64

Handling missing values

In [10]: 1 df1=df.drop(columns=["New", "PendingS", "index"])

In [11]: 1 df1[df1.duplicated()].shape

Out[11]: (959, 18)

In [12]: 1 df1.drop_duplicates(inplace=True)

In [13]: 1 df1 = df1.loc[df1["Qty"] != 0]

```
In [14]: 1 df1.isna().sum()
```

```
Out[14]: 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      124
Amount        124
ship-city     27
ship-state    27
ship-postal-code 27
ship-country  27
B2B           0
fulfilled-by  83341
dtype: int64
```

```
In [15]: 1 df1.dropna(subset=["Amount", "currency"], inplace=True)
```

```
In [16]: 1 df1.dropna(subset=["ship-city", "ship-state", "ship-postal-code", "ship-country"], inplace=True)
```

```
In [17]: 1 df1.isna().sum()
```

```
Out[17]: 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      0
Amount        0
ship-city     0
ship-state    0
ship-postal-code 0
ship-country  0
B2B           0
fulfilled-by  83200
dtype: int64
```

```
In [18]: 1 df1['fulfilled-by'].fillna('Unknown', inplace=True)
```

```
In [19]: 1 df1.isna().sum()
```

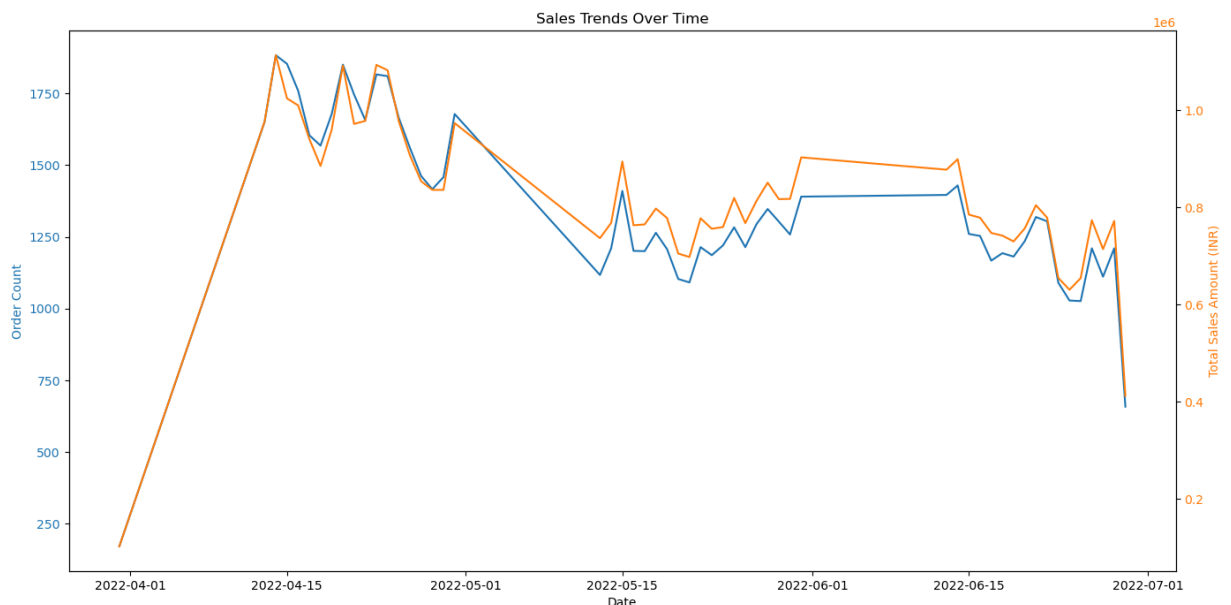
```
Out[19]: Order ID      0
Date          0
Status        0
Fulfilment    0
Sales Channel 0
ship-service-level 0
Category      0
Size          0
Courier Status 0
Qty           0
currency      0
Amount        0
ship-city     0
ship-state    0
ship-postal-code 0
ship-country  0
B2B           0
fulfilled-by  0
dtype: int64
```

1. Sales Overview: To Understand the overall sales performance, trends, and patterns over time.

```
In [21]: 1 df['Date'] = pd.to_datetime(df['Date'], format='%m-%d-%y', errors="coerce")
```

```
In [22]: 1 df = df.dropna(subset=['Date'])
```

```
In [23]: 1 sales_by_date = df.groupby('Date').agg({
2         'Order ID': 'count',
3         'Amount': 'sum'
4     }).rename(columns={'Order ID': 'Order Count', 'Amount': 'Total Sales Amount'})
5 fig, ax1 = plt.subplots(figsize=(14, 7))
6 ax1.set_xlabel('Date')
7 ax1.set_ylabel('Order Count', color='tab:blue')
8 ax1.plot(sales_by_date.index, sales_by_date['Order Count'], color='tab:blue', label='Order Count')
9 ax1.tick_params(axis='y', labelcolor='tab:blue')
10 ax2 = ax1.twinx()
11 ax2.set_ylabel('Total Sales Amount (INR)', color='tab:orange')
12 ax2.plot(sales_by_date.index, sales_by_date['Total Sales Amount'], color='tab:orange', label='Total Sales Amount')
13 ax2.tick_params(axis='y', labelcolor='tab:orange')
14 fig.tight_layout()
15 plt.title('Sales Trends Over Time')
16 plt.show()
```



This shows Peak Sales in the month of April.

```
In [24]: 1 df1['Date'] = pd.to_datetime(df1['Date'], errors='coerce')
          2 print(df1['Date'].isna().sum())
```

0

C:\Users\De11\AppData\Local\Temp\ipykernel_31296\2111675325.py:1: UserWarning: Could not infer format, so each element will be parsed individually, falling back to `dateutil`. To ensure parsing is consistent and as-expected, please specify a format.

```
df1['Date'] = pd.to_datetime(df1['Date'], errors='coerce')
```

```
In [25]: 1 df1['Year'] = df1['Date'].dt.year
          2 df1['Month'] = df1['Date'].dt.month
          3 df1['Day'] = df1['Date'].dt.day
          4 df1['Week'] = df1['Date'].dt.isocalendar().week
          5 df1['Month-Year'] = df1['Date'].dt.to_period('M')
```


In [26]:

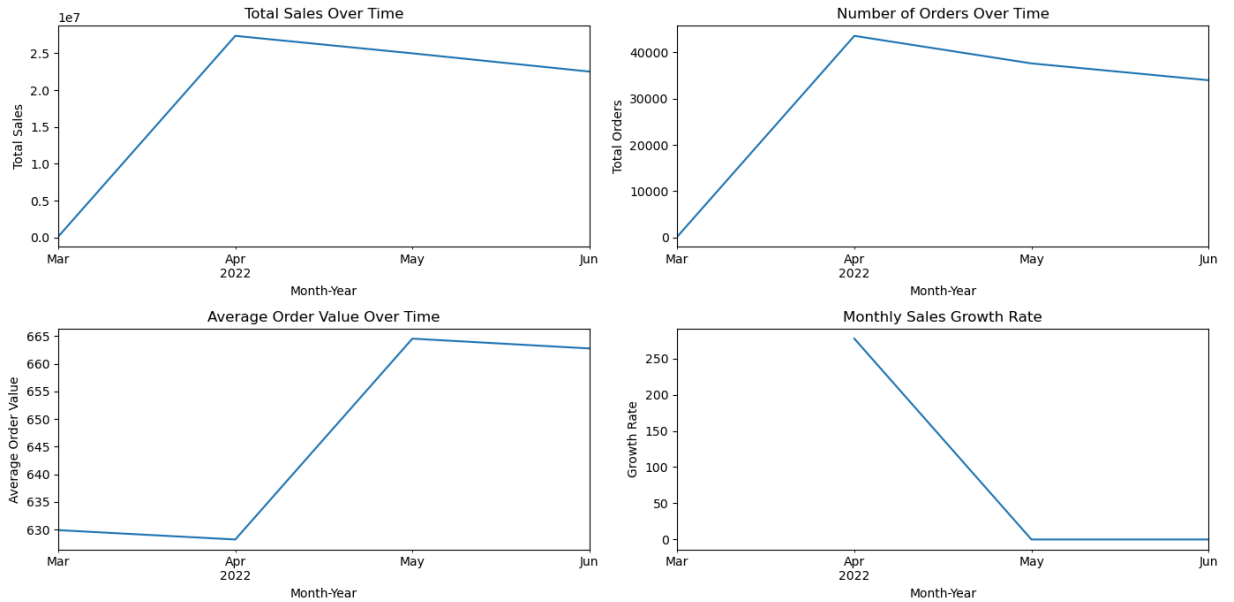
```

1 print(df1['Date'].dtype)
2
3 df1['Date'] = pd.to_datetime(df1['Date'], errors='coerce')
4
5 num_nat = df1['Date'].isna().sum()
6 print(f"Number of NaT values: {num_nat}")
7
8 df1.dropna(subset=['Date'], inplace=True)
9
10 df1['Year'] = df1['Date'].dt.year
11 df1['Month'] = df1['Date'].dt.month
12 df1['Day'] = df1['Date'].dt.day
13 df1['Week'] = df1['Date'].dt.isocalendar().week
14 df1['Month-Year'] = df1['Date'].dt.to_period('M')
15
16 monthly_sales = df1.groupby('Month-Year').agg({'Amount': 'sum', 'Order ID': 'count'}).rename(columns={'Amount': 'Total Sales', 'Order ID': 'Total Orders'})
17 monthly_sales['Average Order Value'] = monthly_sales['Total Sales'] / monthly_sales['Total Orders']
18
19 monthly_sales['Monthly Sales Growth'] = monthly_sales['Total Sales'].pct_change()
20
21 plt.figure(figsize=(14, 7))
22
23 # Total Sales Over Time
24 plt.subplot(2, 2, 1)
25 monthly_sales['Total Sales'].plot()
26 plt.title('Total Sales Over Time')
27 plt.xlabel('Month-Year')
28 plt.ylabel('Total Sales')
29
30 # Number of Orders Over Time
31 plt.subplot(2, 2, 2)
32 monthly_sales['Total Orders'].plot()
33 plt.title('Number of Orders Over Time')
34 plt.xlabel('Month-Year')
35 plt.ylabel('Total Orders')
36
37 # Average Order Value Over Time
38 plt.subplot(2, 2, 3)
39 monthly_sales['Average Order Value'].plot()
40 plt.title('Average Order Value Over Time')
41 plt.xlabel('Month-Year')
42 plt.ylabel('Average Order Value')
43
44 # Sales Growth Rate Over Time
45 plt.subplot(2, 2, 4)
46 monthly_sales['Monthly Sales Growth'].plot()
47 plt.title('Monthly Sales Growth Rate')
48 plt.xlabel('Month-Year')
49 plt.ylabel('Growth Rate')
50
51 plt.tight_layout()
52 plt.show()
53
54 # Calculate specific insights
55 peak_sales_month = monthly_sales['Total Sales'].idxmax()
56 peak_sales_value = monthly_sales['Total Sales'].max()
57
58 average_sales_per_month = monthly_sales['Total Sales'].mean()
59 average_orders_per_month = monthly_sales['Total Orders'].mean()
60
61 print(f"Peak sales month: {peak_sales_month}")
62 print(f"Peak sales value: {peak_sales_value:.2f}")
63 print(f"Average sales per month: {average_sales_per_month:.2f}")
64 print(f"Average number of orders per month: {average_orders_per_month:.2f}")
65
66 high_growth_months = monthly_sales[monthly_sales['Monthly Sales Growth'] > 0.20]
67 print("Months with more than 20% sales growth:")
68 print(high_growth_months[['Total Sales', 'Monthly Sales Growth']])

```

datetime64[ns]

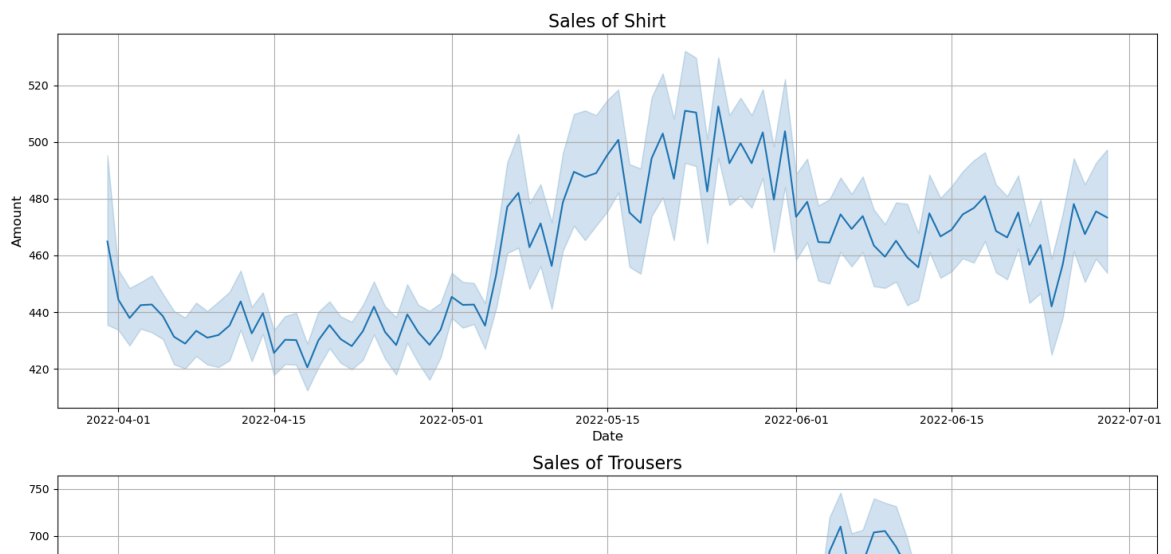
Number of NaT values: 0



Peak sales month: 2022-04
 Peak sales value: 27371152.00
 Average sales per month: 18746710.50
 Average number of orders per month: 28830.75
 Months with more than 20% sales growth:
 Total Sales Monthly Sales Growth
 Month-Year
 2022-04 27371152.0 277.555602

```

In [27]: 1 fig, ax = plt.subplots(9, 1, figsize=(15, 50))
          2
          3 cat = df1['Category'].unique().tolist()
          4 cat = iter(cat)
          5
          6 for i in range(9):
          7     a = next(cat)
          8     df = df1[df1['Category'] == a].copy()
          9     sns.lineplot(x='Date', y='Amount', data=df, ax=ax[i])
         10     ax[i].set_title(f'Sales of {a}', fontsize=16)
         11     ax[i].set_xlabel('Date', fontsize=12)
         12     ax[i].set_ylabel('Amount', fontsize=12)
         13     ax[i].grid(True)
         14
         15 plt.tight_layout()
         16 plt.show()
  
```

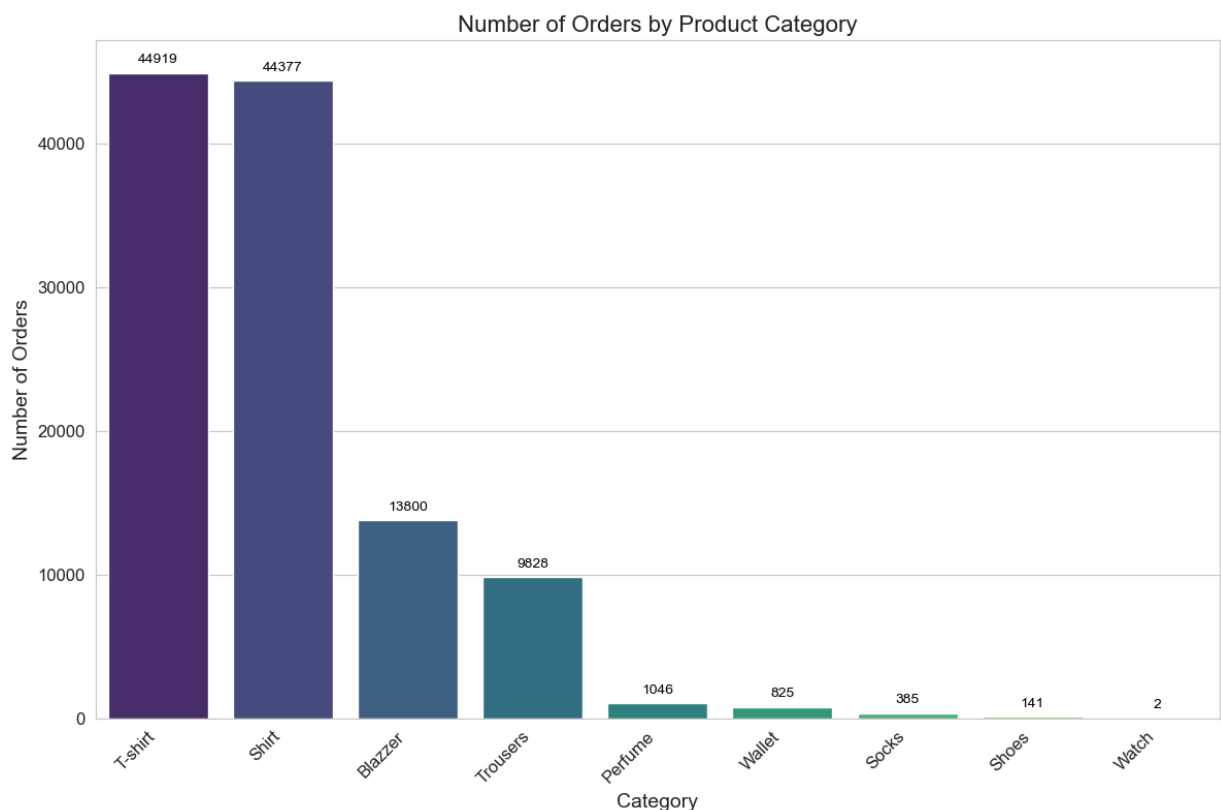


1 . The average sales of each category of product is good in between 1-05-2022 to 1-06-2022.

2 . The watch has no sales as it is the least bought product.

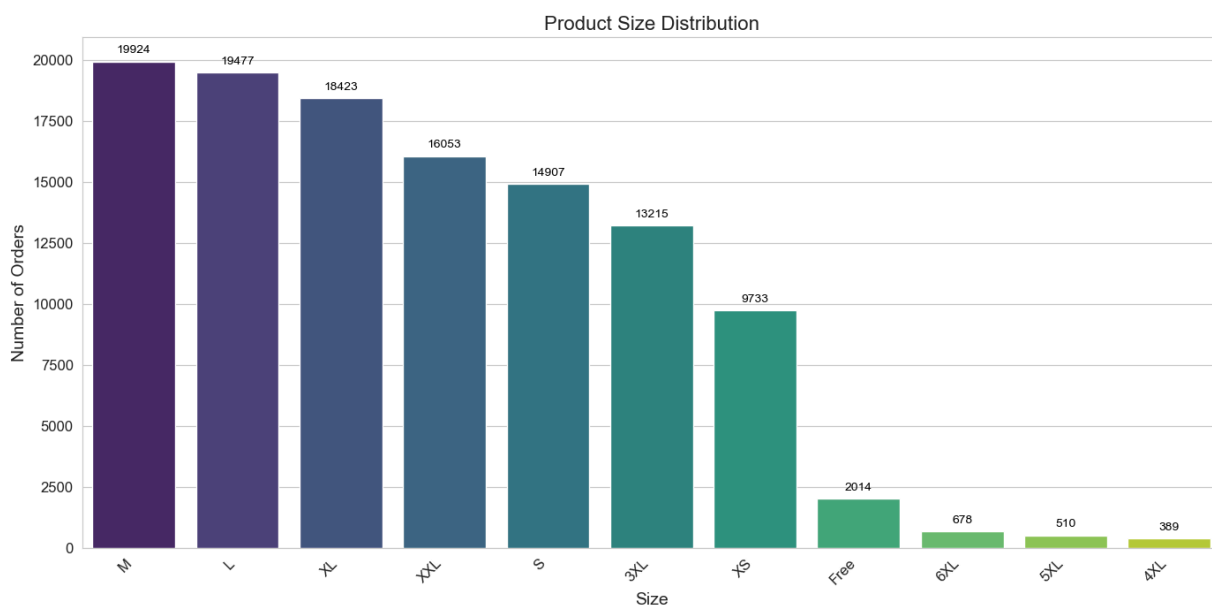
2. Product Analysis: Analyze the distribution of product categories, sizes, and quantities sold to identify popular products.

```
In [28]: 1 category_orders = df1['Category'].value_counts()
2 plt.figure(figsize=(12, 8))
3 sns.set_style("whitegrid")
4
5 ax = sns.barplot(x=category_orders.index, y=category_orders.values, palette="viridis")
6 for p in ax.patches:
7     ax.annotate(format(p.get_height(), '.0f'),
8                 (p.get_x() + p.get_width() / 2., p.get_height()),
9                 ha='center', va='center',
10                xytext=(0, 9),
11                textcoords='offset points',
12                fontsize=10, color='black')
13 plt.title('Number of Orders by Product Category', fontsize=16)
14 plt.xlabel('Category', fontsize=14)
15 plt.ylabel('Number of Orders', fontsize=14)
16 plt.xticks(rotation=45, ha='right', fontsize=12)
17 plt.yticks(fontsize=12)
18 plt.tight_layout()
19 plt.show()
```



This graph shows the number of orders for different categories, where T-shirt has the maximum number of orders and Watch has the minimum number of orders.

```
In [29]: 1 size_distribution = df1['Size'].value_counts()
2 plt.figure(figsize=(14, 7))
3 sns.set_style("whitegrid")
4 ax = sns.barplot(x=size_distribution.index, y=size_distribution.values, palette="viridis")
5 for p in ax.patches:
6     ax.annotate(format(p.get_height(), '.0f'),
7                 (p.get_x() + p.get_width() / 2., p.get_height()),
8                 ha='center', va='center',
9                 xytext=(0, 9),
10                textcoords='offset points',
11                fontsize=10, color='black')
12
13 plt.title('Product Size Distribution', fontsize=16)
14 plt.xlabel('Size', fontsize=14)
15 plt.ylabel('Number of Orders', fontsize=14)
16 plt.xticks(rotation=45, ha='right', fontsize=12)
17 plt.yticks(fontsize=12)
18 plt.tight_layout()
19 plt.show()
```

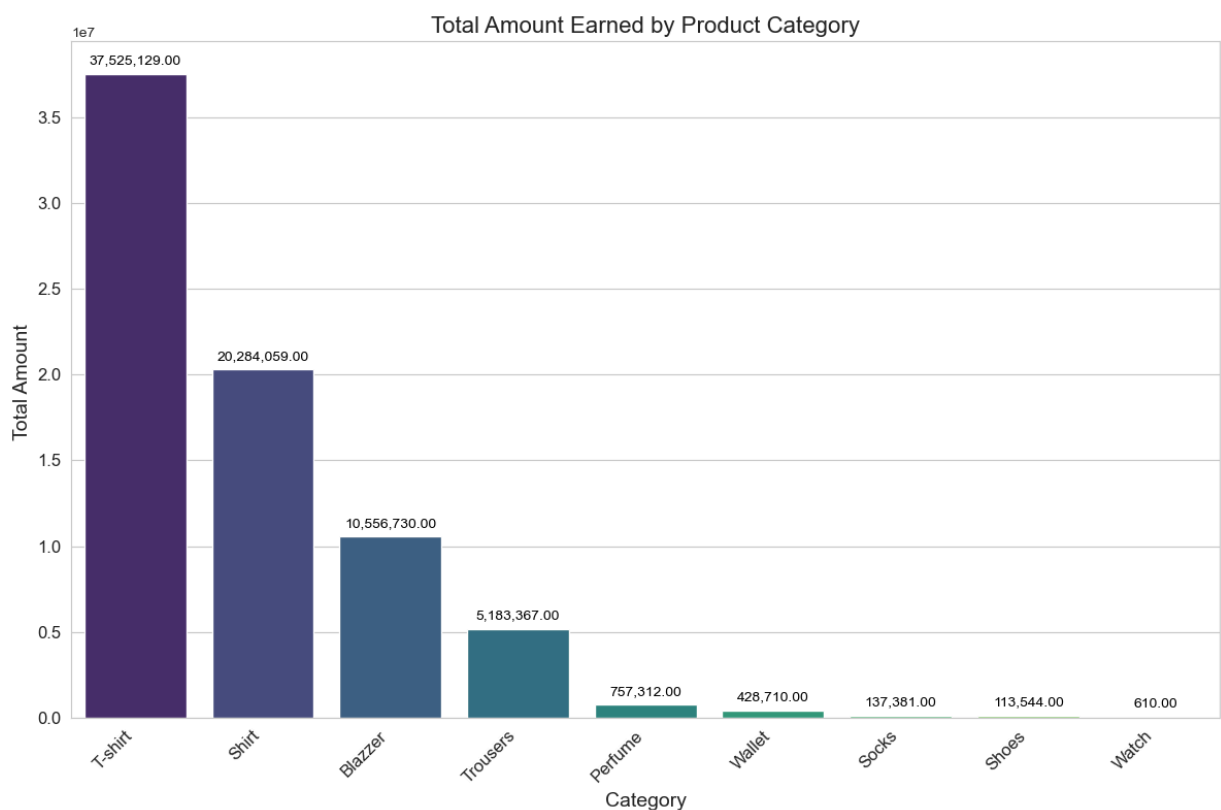


This graph shows the distribution of product size, where M has maximum number of orders and 4XL has minimum.

```

In [30]: 1 category_amount = df1.groupby("Category")["Amount"].sum().sort_values(ascending=False)
2
3
4 plt.figure(figsize=(12, 8))
5 sns.set_style("whitegrid")
6 ax = sns.barplot(x=category_amount.index, y=category_amount.values, palette="viridis")
7 for p in ax.patches:
8     ax.annotate("{:,.2f}".format(p.get_height()),
9                 (p.get_x() + p.get_width() / 2., p.get_height()),
10                ha='center', va='center',
11                xytext=(0, 9),
12                textcoords='offset points',
13                fontsize=10, color='black')
14
15 plt.title('Total Amount Earned by Product Category', fontsize=16)
16 plt.xlabel('Category', fontsize=14)
17 plt.ylabel('Total Amount ', fontsize=14)
18 plt.xticks(rotation=45, ha='right', fontsize=12)
19 plt.yticks(fontsize=12)
20 plt.tight_layout()
21 plt.show()

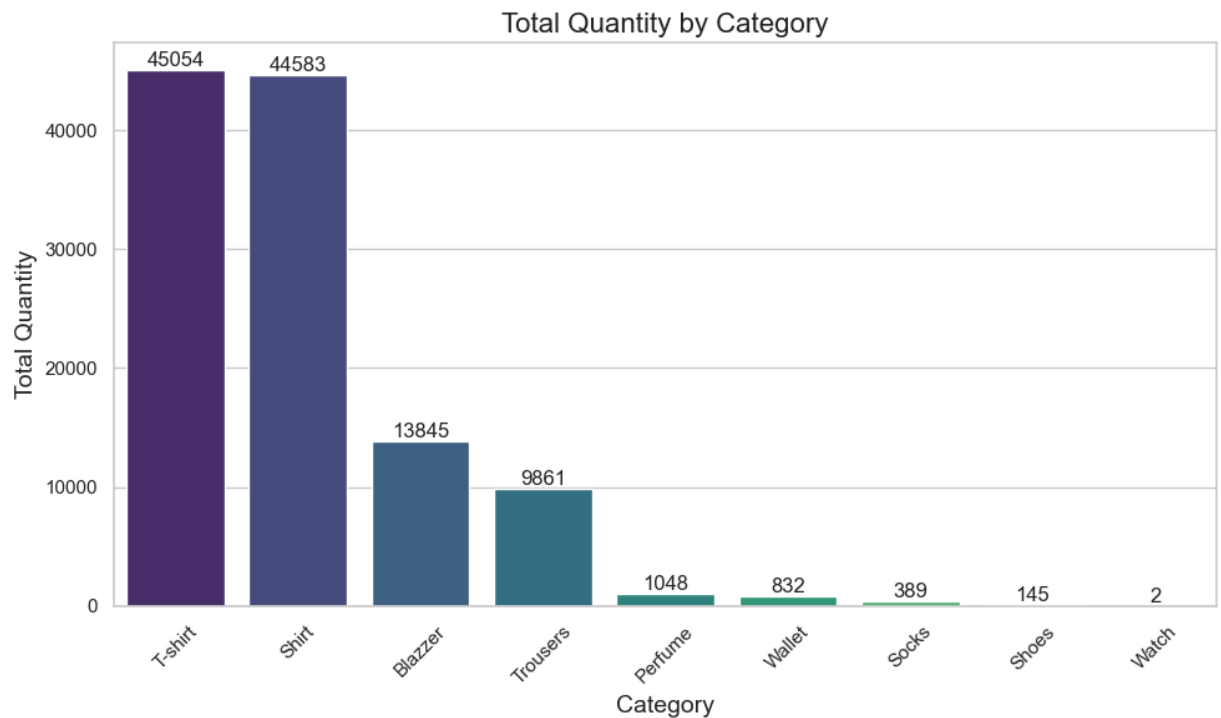
```



This shows the total amount earned by different categories where T-shirt has maximum Sales Amount.

```
In [54]: 1 category_amount = df1.groupby("Category")["Qty"].sum().reset_index()
2 category_amount_sorted = category_amount.sort_values(by='Qty', ascending=False)
3 print(category_amount_sorted)
4 sns.set_style("whitegrid")
5 plt.figure(figsize=(10, 6))
6 bar_plot = sns.barplot(x="Category", y="Qty", data=category_amount_sorted, palette="viridis")
7 bar_plot.set_title('Total Quantity by Category', fontsize=16)
8 bar_plot.set_xlabel('Category', fontsize=14)
9 bar_plot.set_ylabel('Total Quantity', fontsize=14)
10 plt.xticks(rotation=45)
11 for container in bar_plot.containers:
12     bar_plot.bar_label(container)
13 plt.tight_layout()
14 plt.show()
```

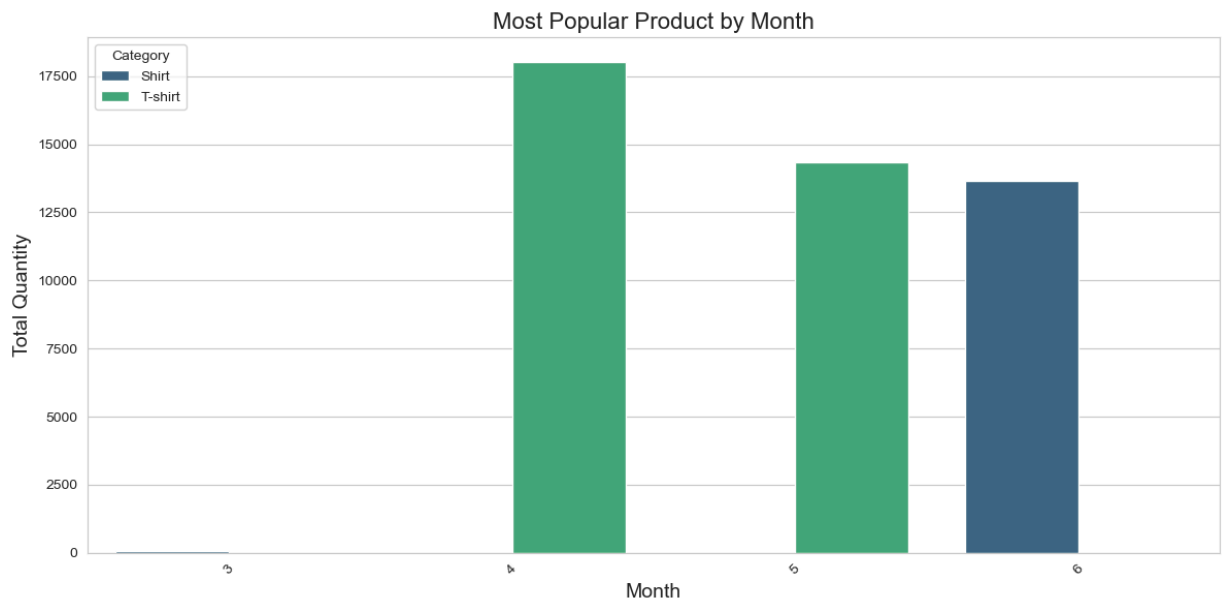
	Category	Qty
5	T-shirt	45054
2	Shirt	44583
0	Blazzer	13845
6	Trousers	9861
1	Perfume	1048
7	Wallet	832
4	Socks	389
3	Shoes	145
8	Watch	2



This graph shows T-shirt has maximum quantity So we can focus on the Marketing strategies and Product Diversification for Tshirt

```
In [32]: 1 monthly_category_amount = df1.groupby(['Month', 'Category'])['Qty'].sum().reset_index()
2 most_popular_by_month = monthly_category_amount.loc[monthly_category_amount.groupby('Month')
3
4 print(most_popular_by_month)
5 sns.set_style("whitegrid")
6 plt.figure(figsize=(12, 6))
7 bar_plot = sns.barplot(x="Month", y="Qty", hue="Category", data=most_popular_by_month, palette="magma")
8 bar_plot.set_title('Most Popular Product by Month', fontsize=16)
9 bar_plot.set_xlabel('Month', fontsize=14)
10 bar_plot.set_ylabel('Total Quantity', fontsize=14)
11 plt.xticks(rotation=45)
12 plt.tight_layout()
13 plt.show()
```

	Month	Category	Qty
2	3	Shirt	71
11	4	T-shirt	18015
19	5	T-shirt	14354
24	6	Shirt	13665



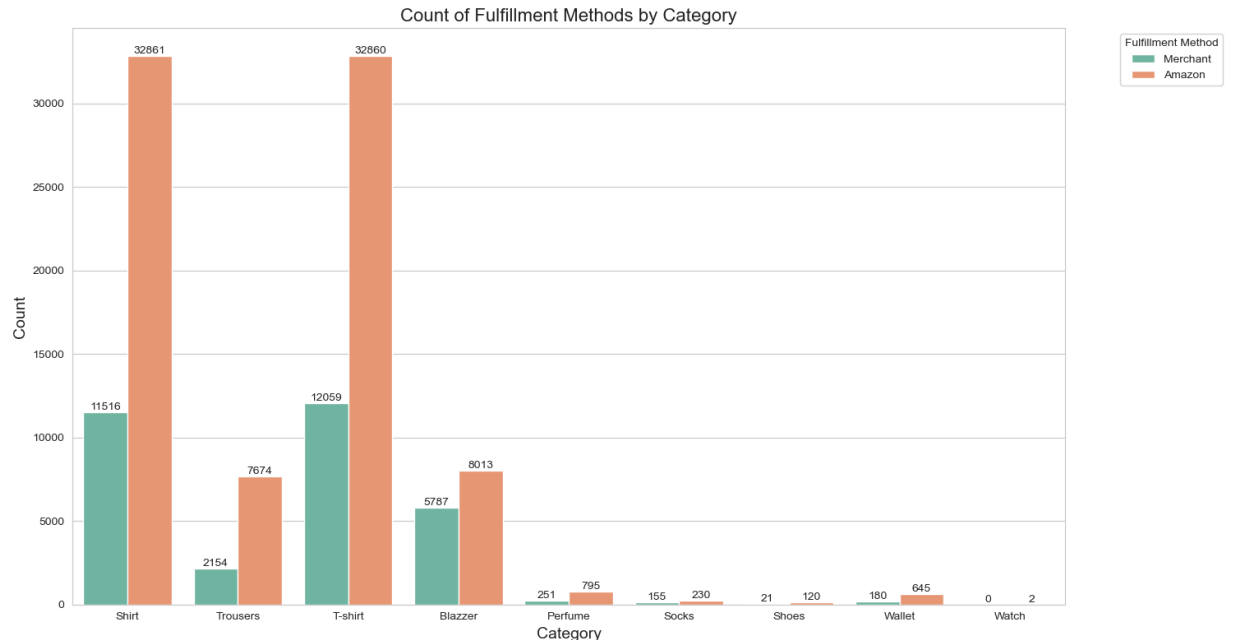
This shows Famous Product by month Where T-shirt has maximum sales in April and May but in June the Sales of Shirt is Maximum

3. Fulfillment Analysis: Investigate the fulfillment methods used and their effectiveness in delivering orders.

```
In [33]: 1 df1.groupby("Category")["Fulfilment"].value_counts()
```

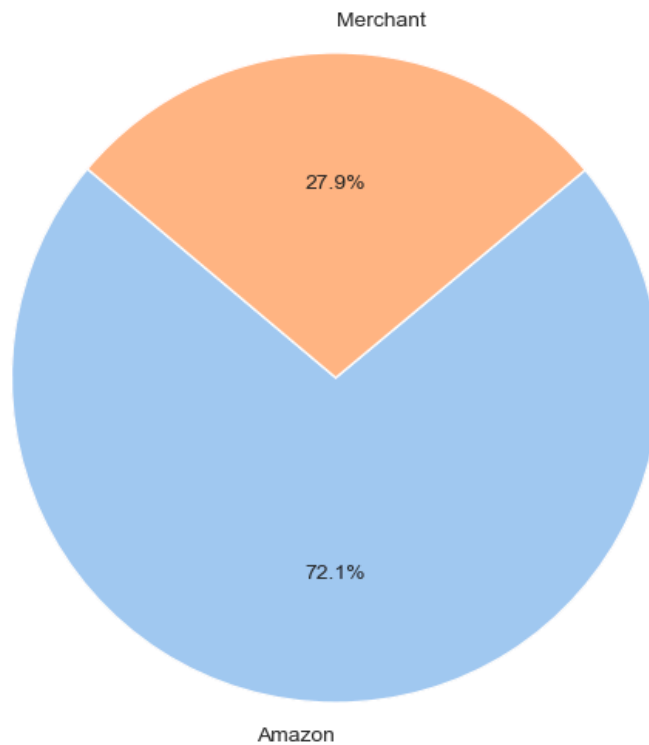
```
Out[33]: Category Fulfilment
Blazzer Amazon      8013
          Merchant    5787
Perfume Amazon       795
          Merchant    251
Shirt Amazon    32861
          Merchant  11516
Shoes Amazon      120
          Merchant    21
Socks Amazon      230
          Merchant   155
T-shirt Amazon   32860
          Merchant 12059
Trousers Amazon   7674
          Merchant  2154
Wallet Amazon     645
          Merchant   180
Watch Amazon        2
Name: count, dtype: int64
```

```
In [34]: 1 sns.set_style("whitegrid")
2 palette = sns.color_palette("Set2")
3 plt.figure(figsize=(15, 8))
4 ax = sns.countplot(data=df1, x="Category", hue="Fulfilment", palette=palette)
5 for container in ax.containers:
6     ax.bar_label(container)
7 plt.title('Count of Fulfilment Methods by Category', fontsize=16)
8 plt.xlabel('Category', fontsize=14)
9 plt.ylabel('Count', fontsize=14)
10 plt.legend(title='Fulfilment Method', bbox_to_anchor=(1.05, 1), loc='upper left')
11 plt.tight_layout()
12 plt.show()
```



This bar chart shows the distribution of fulfilment methods used for Products

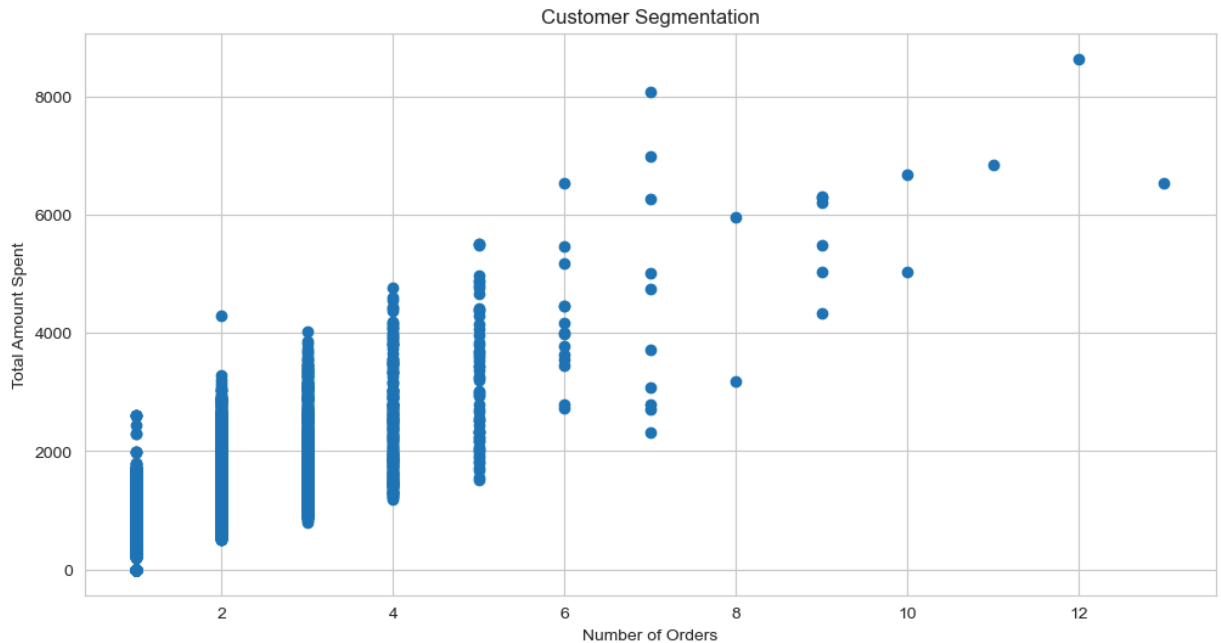
```
In [35]: 1 fulfillment_counts = df1['Fulfilment'].value_counts(normalize=True) * 100
2 sns.set_style("whitegrid")
3 colors = sns.color_palette("pastel")
4 plt.figure(figsize=(10, 6))
5 plt.pie(fulfillment_counts, labels=fulfillment_counts.index, autopct='%1.1f%%', startangle=0)
6 plt.axis('equal')
7 plt.show()
```



Fulfilment by Amazon is the most efficient method.

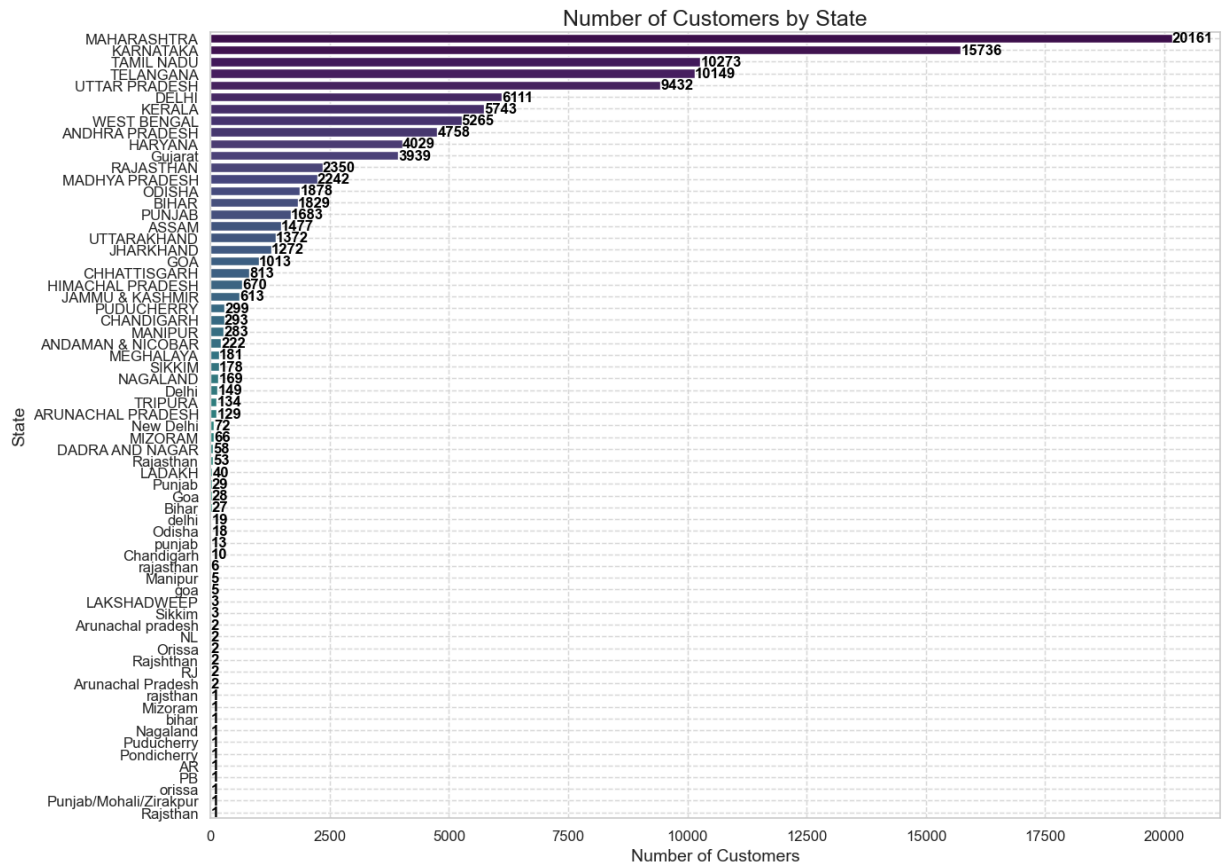
4. Customer Segmentation: Segment customers based on their buying behaviour, location, and other relevant factors.


```
In [39]: 1 customer_segments = df1.groupby('Order ID').agg({
2         'Amount': 'sum',
3         'Qty': 'sum'
4     }).rename(columns={'Amount': 'total_spent', 'Qty': 'order_count'})
5
6     # Visualize customer segmentation
7     plt.figure(figsize=(12, 6))
8     plt.scatter(customer_segments['order_count'], customer_segments['total_spent'])
9     plt.title('Customer Segmentation')
10    plt.xlabel('Number of Orders')
11    plt.ylabel('Total Amount Spent')
12    plt.grid(True)
13    plt.show()
```



- 1 This scatter plot illustrates customer segmentation based on the number of orders and the total amount spent.amount they spent.
- 2 Positive Correlation: There's a clear trend showing that as the number of orders increases, the total amount spent also increases. This suggests that frequent buyers tend to spend more overall.
- 3 Customer Distribution: Most customers have placed between 1 to 6 orders, with spending generally ranging from ₹1,000 to ₹4,000. Outliers are seen where customers with a higher number of orders have spent significantly more, some exceeding ₹8,000.
- 4 Insight: The data indicates a potential opportunity to target customers with fewer orders through personalized marketing to encourage repeat purchases, thereby increasing their total spend.

```
In [40]: 1 state_customer_count = df1['ship-state'].value_counts().reset_index()
2 state_customer_count.columns = ['Ship State', 'Customer Count']
3 plt.figure(figsize=(14, 10))
4 sns.set(style="whitegrid")
5 barplot = sns.barplot(x='Customer Count', y='Ship State', data=state_customer_count, palette=
6 plt.title('Number of Customers by State', fontsize=18)
7 plt.xlabel('Number of Customers', fontsize=14)
8 plt.ylabel('State', fontsize=14)
9 for index, value in enumerate(state_customer_count['Customer Count']):
10     barplot.text(value + 0.1, index, str(value), va='center', ha='left', color='black', font
11 plt.grid(True, linestyle='--', alpha=0.7)
12 plt.xticks(fontsize=12)
13 plt.yticks(fontsize=12)
14 plt.tight_layout()
15 plt.show()
```

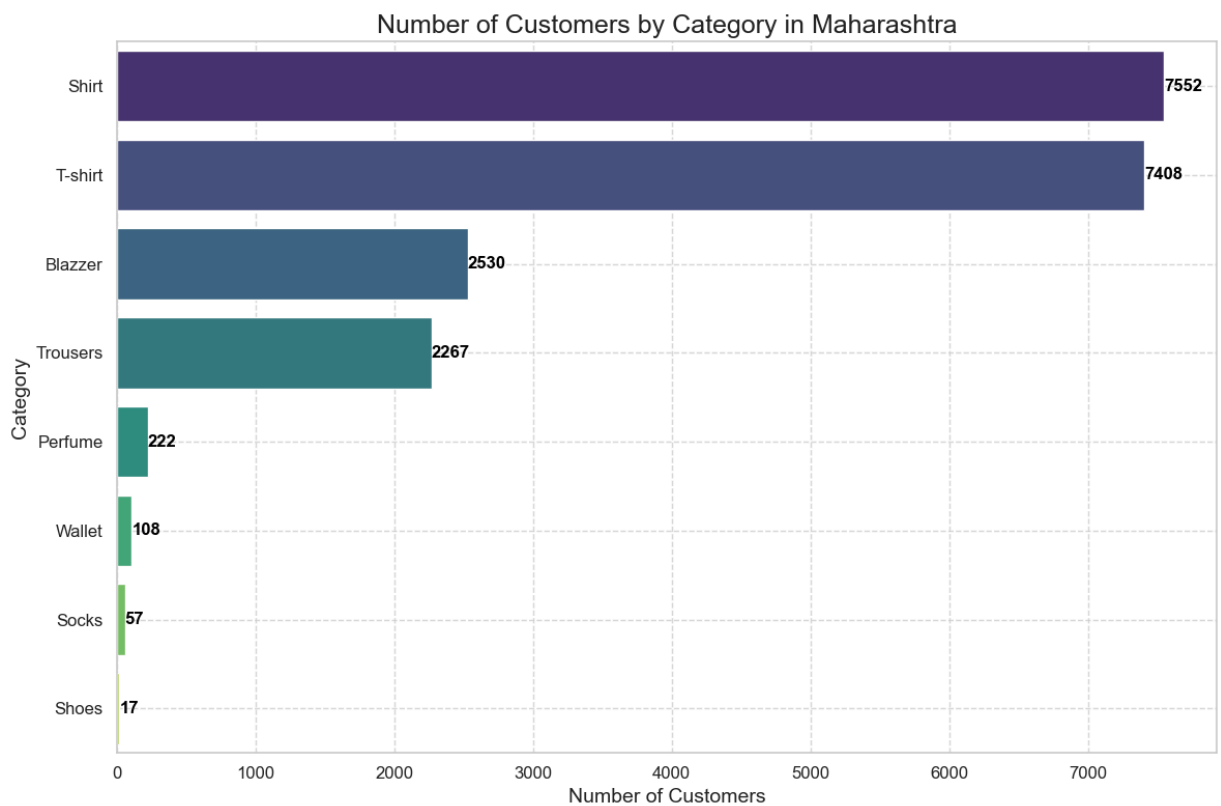


- 1 Top States:
- 2 Maharashtra leads with the highest number of customers (20,161).
- 3 Karnataka follows with 15,736 customers.
- 4 Tamil Nadu and Telangana also have significant customer bases with 10,273 and 10,149 customers.
- 5 Mid-Tier States:
- 6 Uttar Pradesh, Kerala, and West Bengal have moderate customer numbers ranging from approximately 5,000 to 9,000.
- 7 Lower-Tier States:
- 8 Several states, including Haryana, Gujarat, Rajasthan, and Madhya Pradesh, have fewer than 5,000 customers.
- 9 States like Sikkim, Nagaland, and Tripura have customer numbers below 200.
- 10 Least Populated States:
- 11 Lakshadweep, Sikkim, and Mizoram have the smallest customer bases, each with fewer than 100 customers.
- 12 Recommendations:
- 13 Focus on High Customer States: Develop targeted marketing strategies and promotions in states like Maharashtra, Karnataka, Tamil Nadu, and Telangana to further increase sales and customer engagement.
- 14 Expand in Mid-Tier States: Strengthen presence and explore growth opportunities in states with moderate customer bases such as Uttar Pradesh, Kerala, and West Bengal.

- 15 Address Low Customer States: Investigate the low customer numbers in states with minimal representation to identify potential barriers and tailor specific strategies to attract new customers.

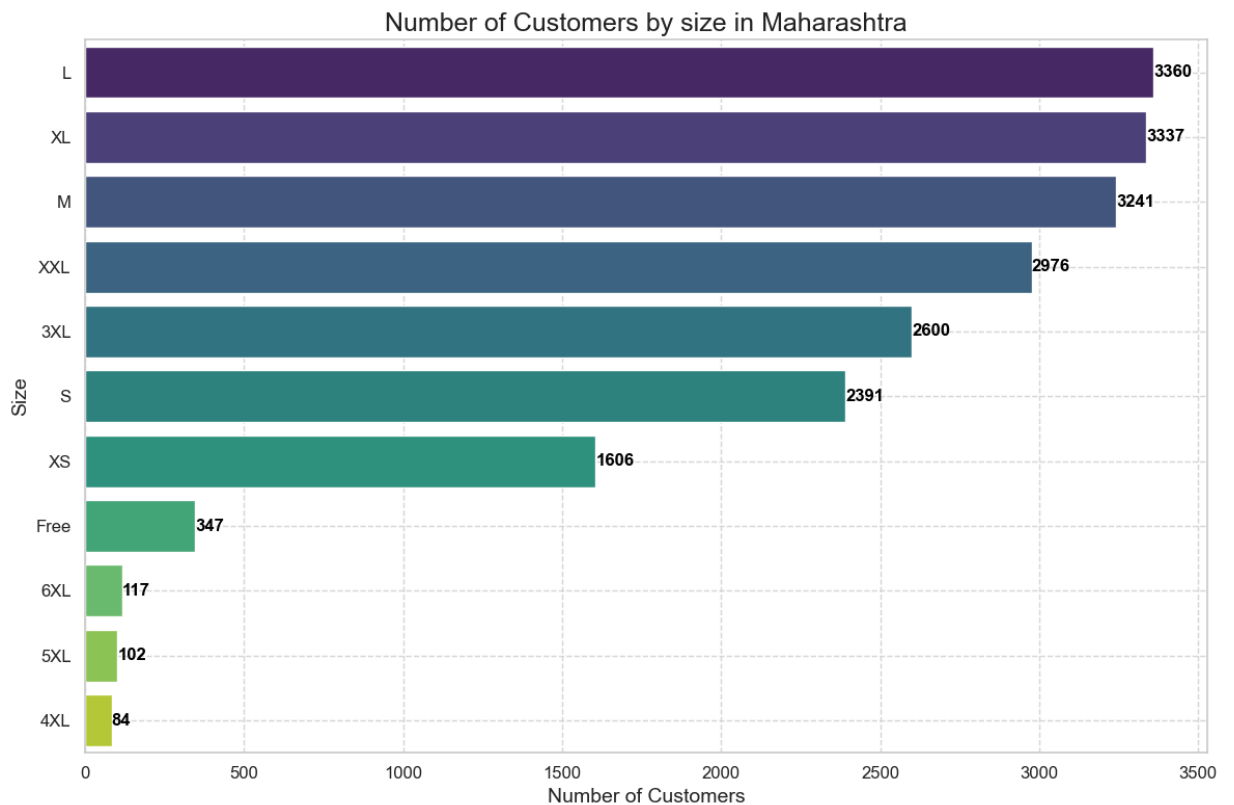
In [41]:

```
1 maharashtra_orders = df1[df1['ship-state'] == 'MAHARASHTRA']
2 category_customer_count = maharashtra_orders['Category'].value_counts().reset_index()
3 category_customer_count.columns = ['Category', 'Customer Count']
4 plt.figure(figsize=(12, 8))
5 sns.set(style="whitegrid")
6 barplot = sns.barplot(x='Customer Count', y='Category', data=category_customer_count, palette=
7 plt.title('Number of Customers by Category in Maharashtra', fontsize=18)
8 plt.xlabel('Number of Customers', fontsize=14)
9 plt.ylabel('Category', fontsize=14)
10 for index, value in enumerate(category_customer_count['Customer Count']):
11     barplot.text(value + 0.1, index, str(value), va='center', ha='left', color='black', font
12 plt.grid(True, linestyle='--', alpha=0.7)
13 plt.xticks(fontsize=12)
14 plt.yticks(fontsize=12)
15 plt.tight_layout()
16 plt.show()
```



- 1 To focus on the state Maharashtra with maximum sales and orders, here the Shirt has maximum number of customers.

```
In [42]: 1 maharashtra_orders = df1[df1['ship-state'] == 'MAHARASHTRA']
2 category_customer_count = maharashtra_orders['Size'].value_counts().reset_index()
3 category_customer_count.columns = ['Size', 'Customer Count']
4 plt.figure(figsize=(12, 8))
5 sns.set(style="whitegrid")
6 barplot = sns.barplot(x='Customer Count', y='Size', data=category_customer_count, palette='magma')
7 plt.title('Number of Customers by size in Maharashtra', fontsize=18)
8 plt.xlabel('Number of Customers', fontsize=14)
9 plt.ylabel('Size', fontsize=14)
10 for index, value in enumerate(category_customer_count['Customer Count']):
11     barplot.text(value + 0.1, index, str(value), va='center', ha='left', color='black', fontweight='bold')
12 plt.grid(True, linestyle='--', alpha=0.7)
13 plt.xticks(fontsize=12)
14 plt.yticks(fontsize=12)
15 plt.tight_layout()
16 plt.show()
```



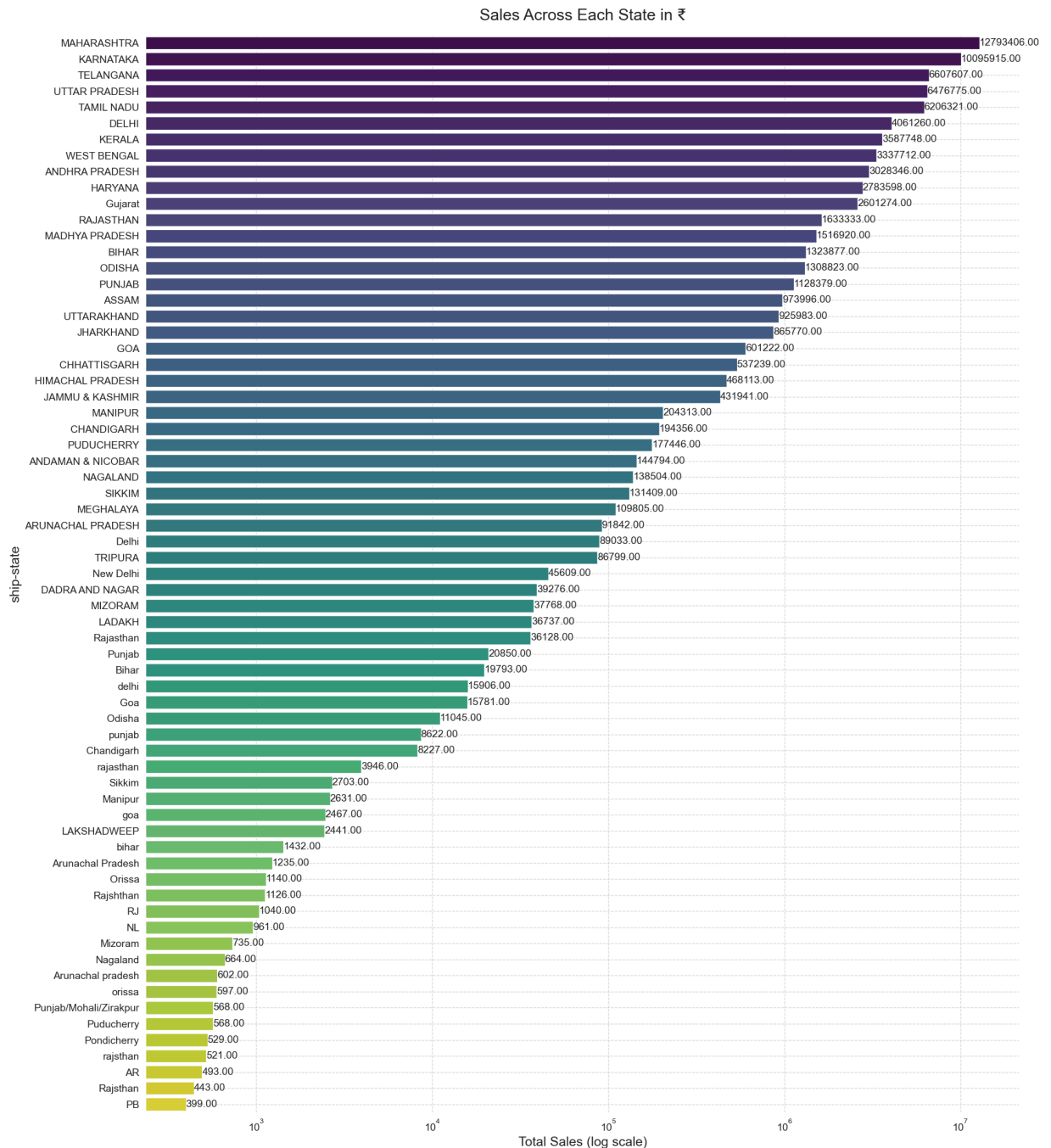
1 This shows in Maharashtra L size has maximum number of customers.

In [43]:

```

1 plt.figure(figsize=(20, 22))
2 df2 = pd.DataFrame(df1.groupby("ship-state")["Amount"].sum())
3 df2 = df2.sort_values(by="Amount", ascending=False)
4 sns.set(style="whitegrid")
5 axs = sns.barplot(y=df2.index, x="Amount", data=df2, palette="viridis", log=True)
6 for container in axs.containers:
7     axs.bar_label(container, fmt="%0.2f", label_type='edge', fontsize=14)
8 plt.ylabel("ship-state", fontsize=18)
9 plt.xlabel("Total Sales (log scale)", fontsize=18)
10 plt.title("Sales Across Each State in ₹", fontsize=22, pad=20)
11 plt.xticks(fontsize=14)
12 plt.yticks(fontsize=14)
13 plt.grid(True, linestyle='--', alpha=0.7)
14 sns.despine(left=True, bottom=True)
15 plt.tight_layout()
16 plt.show()

```



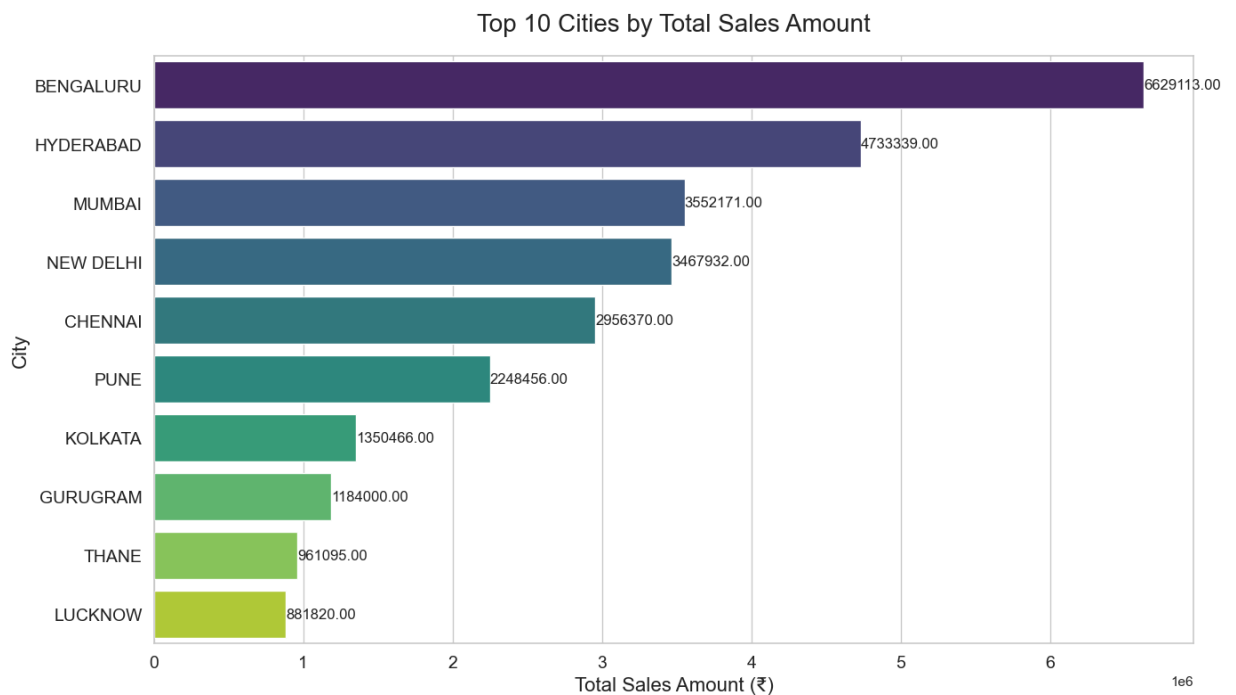
In []:

```

1 This shows Maharashtra has maximum sales.

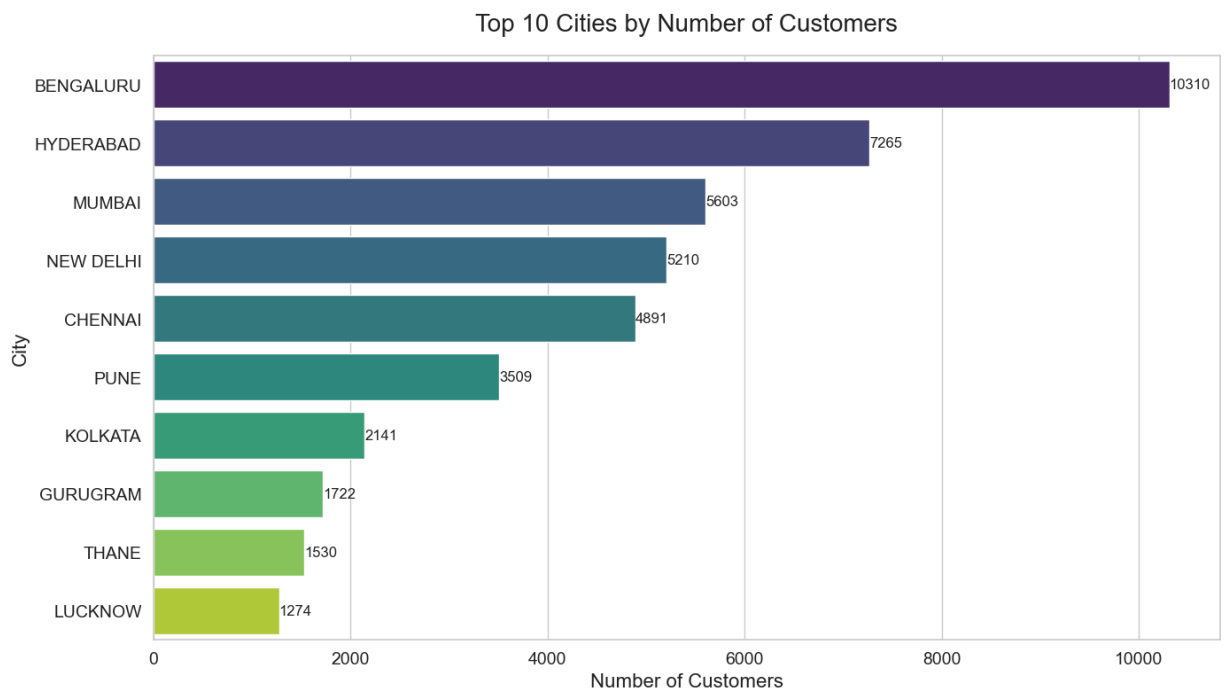
```

```
In [44]: 1 city_sales = df1.groupby('ship-city')['Amount'].sum().reset_index()
2 top_cities = city_sales.sort_values(by='Amount', ascending=False).head(10)
3 plt.figure(figsize=(14, 8))
4 sns.set(style="whitegrid")
5 barplot = sns.barplot(y='ship-city', x='Amount', data=top_cities, palette='viridis')
6 for container in barplot.containers:
7     barplot.bar_label(container, fmt='%.2f', label_type='edge', fontsize=12)
8 plt.ylabel("City", fontsize=16)
9 plt.xlabel("Total Sales Amount (₹)", fontsize=16)
10 plt.title("Top 10 Cities by Total Sales Amount", fontsize=20, pad=20)
11 plt.xticks(fontsize=14)
12 plt.yticks(fontsize=14)
13 plt.tight_layout()
14 plt.show()
```

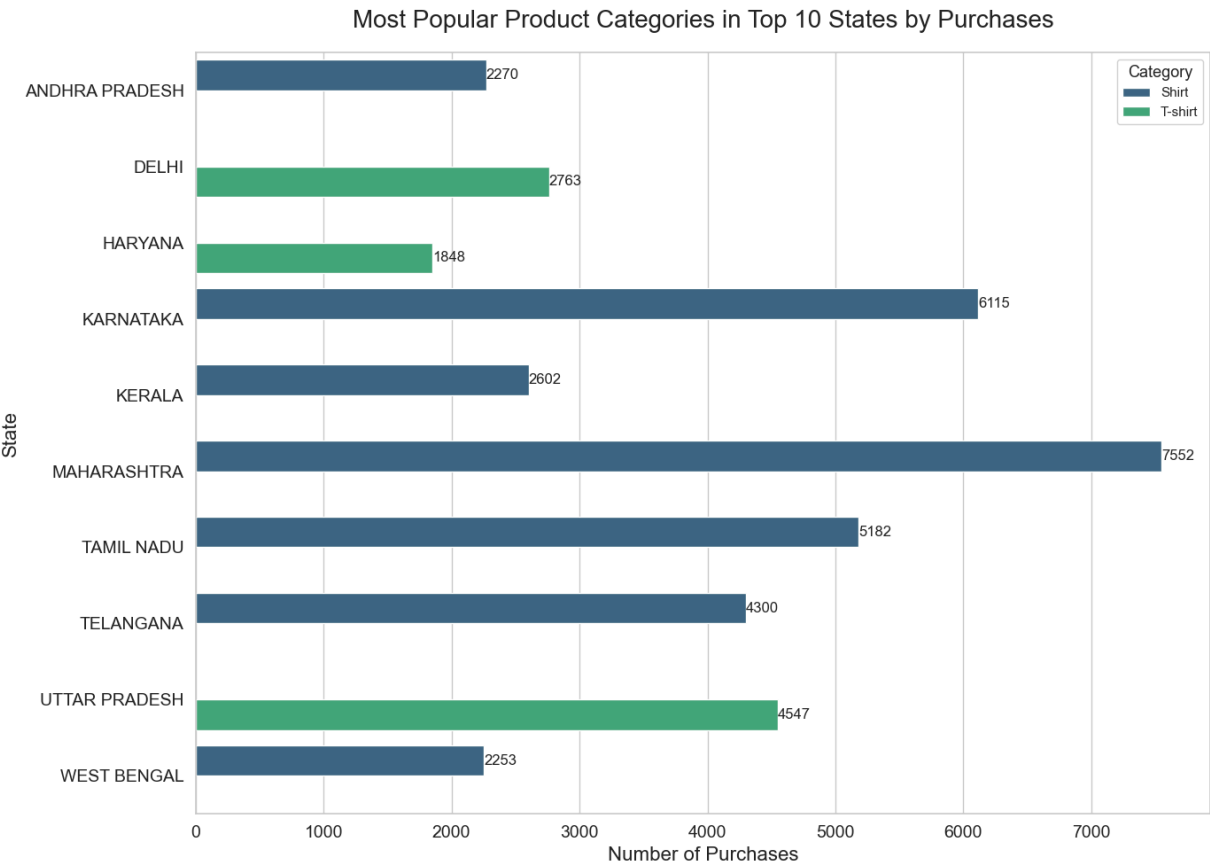


1 This shows the City with maximum sales amount, here Bengaluru has maximum sales amount

```
In [45]: 1 city_customers = df1.groupby('ship-city').size().reset_index(name='Customer Count')
2 top_cities_customers = city_customers.sort_values(by='Customer Count', ascending=False).head(10)
3 plt.figure(figsize=(14, 8))
4 sns.set(style="whitegrid")
5 barplot = sns.barplot(y='ship-city', x='Customer Count', data=top_cities_customers, palette="magma")
6 for container in barplot.containers:
7     barplot.bar_label(container, fmt='%d', label_type='edge', fontsize=12)
8 plt.ylabel("City", fontsize=16)
9 plt.xlabel("Number of Customers", fontsize=16)
10 plt.title("Top 10 Cities by Number of Customers", fontsize=20, pad=20)
11 plt.xticks(fontsize=14)
12 plt.yticks(fontsize=14)
13 plt.tight_layout()
14 plt.show()
```



```
In [46]: 1 state_purchase_counts = df1.groupby('ship-state').size().reset_index(name='Total Purchases')
2 top_10_states = state_purchase_counts.sort_values(by='Total Purchases', ascending=False).head(10)
3 top_10_states_df = df1[df1['ship-state'].isin(top_10_states['ship-state'])]
4 popular_categories = top_10_states_df.groupby(['ship-state', 'Category']).size().reset_index(name='Counts')
5 most_popular_categories = popular_categories.loc[popular_categories.groupby('ship-state')['Counts'].rank(1)]
6 plt.figure(figsize=(14, 10))
7 sns.set(style="whitegrid")
8 barplot = sns.barplot(y='ship-state', x='Counts', hue='Category', data=most_popular_categories)
9 for container in barplot.containers:
10     barplot.bar_label(container, fmt='%d', label_type='edge', fontsize=12)
11 plt.ylabel("State", fontsize=16)
12 plt.xlabel("Number of Purchases", fontsize=16)
13 plt.title("Most Popular Product Categories in Top 10 States by Purchases", fontsize=20, pad=10)
14 plt.xticks(fontsize=14)
15 plt.yticks(fontsize=14)
16 plt.legend(title='Category', title_fontsize='13', fontsize='11', loc='upper right')
17 plt.tight_layout()
18 plt.show()
```



- 1 Maharashtra:
- 2 The state has the highest number of purchases in the Shirt category with 7,552 purchases.
- 3 Karnataka:
- 4 Follows Maharashtra with 6,115 purchases in the Shirt category.
- 5 Delhi:
- 6 Notably, T-shirts are more popular with 2,763 purchases compared to Shirts.
- 7 Uttar Pradesh:
- 8 Similar to Delhi, T-shirts lead with 4,547 purchases.
- 9 Tamil Nadu and Telangana:
- 10 These states have significant purchases in the Shirt category with 5,182 and 4,300.
- 11 Kerala, Andhra Pradesh, West Bengal, and Haryana:
- 12 All these states show a higher preference for Shirts over T-shirts.
- 13 Actionable Recommendations:
- 14 Maharashtra and Karnataka: Focus on maintaining the inventory and introducing new styles in the Shirt category to sustain and grow the market.
- 15 Delhi and Uttar Pradesh: Increase marketing efforts and promotions for T-shirts given their popularity in these states.

16	Other States: Tailor marketing strategies to emphasize the popular product category in each state, potentially boosting sales by aligning with local preferences.
----	---

6. Business Insights: Provide actionable insights and recommendations based on the analysis to optimize sales strategies, improve customer satisfaction, and enhance overall business performance.

1. Peak sales periods are observed in April. Plan inventory and marketing campaigns accordingly.
2. Product categories T-shirt, Shirt, Blazzer are most popular. Consider expanding these categories.
3. High-value customers are those with more than 2 orders and total spending above 828.00. Focus on loyalty programs for these customers.
4. States MAHARASHTRA, KARNATAKA, TELANGANA and cities BENGALURU, HYDERABAD, MUMBAI show high sales. Consider increasing marketing efforts and presence in these areas.
5. Maharashtra has the maximum number of Customers with maximum Sales, the famous Product in Maharashtra is Shirt and the most borrowed size is "L"
6. The Fulfilment by Amazon is the most efficient fulfilment method.