

## Analyzing Consumer Shopping Behavior:

- Insights from Demographics, Purchases, and Geospatial Trends

### ✓ Introduction

In this project, I aim to analyze customer shopping behavior by utilizing a dataset that includes valuable details about customer demographics, items purchased, payment methods, and more. By combining demographic, purchasing, and geographic data, I plan to identify patterns in consumer behavior, understand how different payment methods influence purchasing decisions, and explore how geographic factors like location and seasonality play a role in shopping trends. With this comprehensive approach, the project will provide insights into the interrelationships between various factors that influence consumer choices.

Additionally, I will explore the spatial dimension of shopping behavior using GeoJSON data to create geographical visualizations. These visualizations will map the purchasing patterns across different US states, allowing me to identify regional variations in customer behavior. By examining these trends through the lens of data visualization and statistical analysis, I aim to deliver a deeper understanding of the drivers of consumer behavior and provide actionable insights for businesses to enhance their marketing strategies.

### Datasets Used

1. Customer Shopping Trends Dataset (from Kaggle):

- [https://www.kaggle.com/datasets/iamsouravbanerjee/customer-shopping-trends-dataset?select=shopping\\_trends\\_updated.csv](https://www.kaggle.com/datasets/iamsouravbanerjee/customer-shopping-trends-dataset?select=shopping_trends_updated.csv)
- This dataset contains valuable information about customer purchasing behavior, including various attributes such as customer demographics, items purchased, payment methods, and more. It can be used for understanding shopping trends, segmenting customers, and analyzing purchasing patterns. You can explore correlations between different variables such as payment method, category, and customer profile.

2. HD Pulse Data (from NIMHD):

- <https://hdpulse.nimhd.nih.gov/data-portal/social/table>
- The HD Pulse Data provides insights into social and economic factors that can affect health outcomes. This dataset offers a range of social determinants of health data across different regions and populations. It's useful for exploring the relationship between social factors (like income, education, etc.) and health disparities, which can contribute to understanding the broader context of well-being and health interventions.

3. US States GeoJSON Data (from ERIC):

- <https://eric.clst.org/tech/usgeojson/>
- This dataset provides geographical boundaries of US states in GeoJSON format. It is typically used for creating visualizations such as choropleth maps, where geographic boundaries are essential to display spatial data like population, income levels, or other regional metrics. The dataset allows for visual exploration of US states and is ideal for geographic visualizations when combined with relevant socio-economic or demographic data.

```
1 import pandas as pd
2 import geopandas as gpd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import plotly.express as px
6 import json
```

### ✓ Load the datasets


```
1 shopping_trends = pd.read_csv('shopping_trends_updated.csv')
```

```
1 Income_by_states = pd.read_json("Income by states.json")
```


```
1 with open('us-states geodata.json', 'r') as f:
2     geojson = json.load(f)
```

### ✓ Data transformations

```
1 shopping_trends.head()
```



	Customer ID	Age	Gender	Item Purchased	Category	Purchase Amount (USD)	Location	Size	Color	Season	Review Rating	Subscription Status	Shipping Type	Discount Applied
0	1	55	Male	Blouse	Clothing	53	Kentucky	L	Gray	Winter	3.1	Yes	Express	Y
1	2	19	Male	Sweater	Clothing	64	Maine	L	Maroon	Winter	3.1	Yes	Express	Y
2	3	50	Male	Jeans	Clothing	73	Massachusetts	S	Maroon	Spring	3.1	Yes	Free Shipping	Y
3	4	21	Male	Sandals	Footwear	90	Rhode Island	M	Maroon	Spring	3.5	Yes	Next Day Air	Y
4	5	45	Male	Blouse	Clothing	49	Oregon	M	Turquoise	Spring	2.7	Yes	Free Shipping	Y




Next steps:

[Generate code with shopping\\_trends](#)

 [View recommended plots](#)


[New interactive sheet](#)

```
1 shopping_trends.shape
```




```
(3900, 18)
```

```
1 shopping_trends.columns
```



```
Index(['Customer ID', 'Age', 'Gender', 'Item Purchased', 'Category',
      'Purchase Amount (USD)', 'Location', 'Size', 'Color', 'Season',
      'Review Rating', 'Subscription Status', 'Shipping Type',
      'Discount Applied', 'Promo Code Used', 'Previous Purchases',
      'Payment Method', 'Frequency of Purchases'],
      dtype='object')
```

```
1 shopping_trends.info()
```



```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3900 entries, 0 to 3899
Data columns (total 18 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Customer ID                          3900 non-null  int64
1   Age                                  3900 non-null  int64
2   Gender                               3900 non-null  object
3   Item Purchased                       3900 non-null  object
4   Category                             3900 non-null  object
5   Purchase Amount (USD)                3900 non-null  int64
6   Location                             3900 non-null  object
7   Size                                 3900 non-null  object
8   Color                                3900 non-null  object
9   Season                               3900 non-null  object
10  Review Rating                         3900 non-null  float64
11  Subscription Status                   3900 non-null  object
12  Shipping Type                         3900 non-null  object
13  Discount Applied                     3900 non-null  object
14  Promo Code Used                      3900 non-null  object
15  Previous Purchases                   3900 non-null  int64
16  Payment Method                       3900 non-null  object
17  Frequency of Purchases                3900 non-null  object
dtypes: float64(1), int64(4), object(13)
memory usage: 548.6+ KB
```

```
1 # Convert specified columns into category type
2 categorical_columns = [
3     'Gender',
4     'Category',
5     'Size',
6     'Color',
7     'Season',
8     'Subscription Status',
9     'Shipping Type',
10    'Discount Applied',
11    'Promo Code Used',
12    'Payment Method',
13    'Frequency of Purchases'
14 ]
15
16 # Apply category dtype to these columns
17 for col in categorical_columns:
18     shopping_trends[col] = shopping_trends[col].astype('category')
19
```

```
20 # Verify the changes
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3900 entries, 0 to 3899
Data columns (total 18 columns):
#   Column                                Non-Null Count  Dtype
---  ---                                -
0   Customer ID                          3900 non-null   int64
1   Age                                  3900 non-null   int64
2   Gender                              3900 non-null   category
3   Item Purchased                       3900 non-null   object
4   Category                             3900 non-null   category
5   Purchase Amount (USD)                3900 non-null   int64
6   Location                             3900 non-null   object
7   Size                                 3900 non-null   category
8   Color                                3900 non-null   category
9   Season                               3900 non-null   category
10  Review Rating                        3900 non-null   float64
11  Subscription Status                  3900 non-null   category
12  Shipping Type                       3900 non-null   category
13  Discount Applied                    3900 non-null   category
14  Promo Code Used                     3900 non-null   category
15  Previous Purchases                   3900 non-null   int64
16  Payment Method                      3900 non-null   category
17  Frequency of Purchases               3900 non-null   category
dtypes: category(11), float64(1), int64(4), object(2)
memory usage: 257.9+ KB
None
```

```
1 # Convert 'Item Purchased' and 'Location' columns to string type
2 shopping_trends['Item Purchased'] = shopping_trends['Item Purchased'].astype('string')
3 shopping_trends['Location'] = shopping_trends['Location'].astype('string')
```

```
1 Income_by_states.head()
```

	State	FIPS	Value (Dollars)	Rank within US (of 52 states)
0	Puerto Rico	72001	24,002	52
1	Mississippi	28000	52,985	51
2	West Virginia	54000	55,217	50
3	Arkansas	5000	56,335	49
4	Louisiana	22000	57,852	48

Next steps:

[Generate code with Income\\_by\\_states](#)

[View recommended plots](#)

[New interactive sheet](#)

```
1 Income_by_states.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 52 entries, 0 to 51
Data columns (total 4 columns):
#   Column                                Non-Null Count  Dtype
---  ---                                -
0   State                                  52 non-null     object
1   FIPS                                  52 non-null     int64
2   Value (Dollars)                       52 non-null     object
3   Rank within US (of 52 states)         52 non-null     int64
dtypes: int64(2), object(2)
memory usage: 1.8+ KB
```

✓ Pivot tables

Average Purchase Amount by Gender and Category

```
1 pivot_gender = shopping_trends.pivot_table(
2     values='Purchase Amount (USD)',
3     index='Gender',
4     columns='Category',
5     aggfunc='mean',
6     observed=True
7 )
8
9 pivot_gender
```

Category	Accessories	Clothing	Footwear	Outerwear
Gender				
Female	60.762755	60.496403	59.472362	58.425743
Male	59.411557	59.803556	60.645000	56.605381

Next steps:

[Generate code with pivot\\_gender](#)[View recommended plots](#)[New interactive sheet](#)

### Average Review Rating by Category

```

1 pivot_rating = shopping_trends.pivot_table(
2     values='Review Rating',
3     index='Category',
4     aggfunc='mean',
5     observed=False
6 )
7
8
9 pivot_rating

```

Category	Review Rating
Accessories	3.768629
Clothing	3.723143
Footwear	3.790651
Outerwear	3.746914

Next steps:

[Generate code with pivot\\_rating](#)[View recommended plots](#)[New interactive sheet](#)

### Total Previous Purchases by Payment Method

```

1 pivot_payment_methods = shopping_trends.pivot_table(
2     values='Previous Purchases',
3     index='Payment Method',
4     aggfunc='sum',
5     observed=True
6 )
7
8 pivot_payment_methods
9

```

Payment Method	Previous Purchases
Bank Transfer	14995
Cash	16920
Credit Card	17170
Debit Card	16257
PayPal	17270
Venmo	16259

Next steps:

[Generate code with pivot\\_payment\\_methods](#)[View recommended plots](#)[New interactive sheet](#)

### Total Purchase Amount by Subscription Status

```

1 pivot_subscription = shopping_trends.pivot_table(
2     values='Purchase Amount (USD)',
3     index='Subscription Status',
4     aggfunc='sum',
5     observed=False
6 )
7
8 pivot_subscription
9

```

	Purchase Amount (USD)	
Subscription Status		
No	170436	
Yes	62645	

Next steps:

[Generate code with pivot\\_subscription](#)[View recommended plots](#)[New interactive sheet](#)

## ✓ Data visualizations

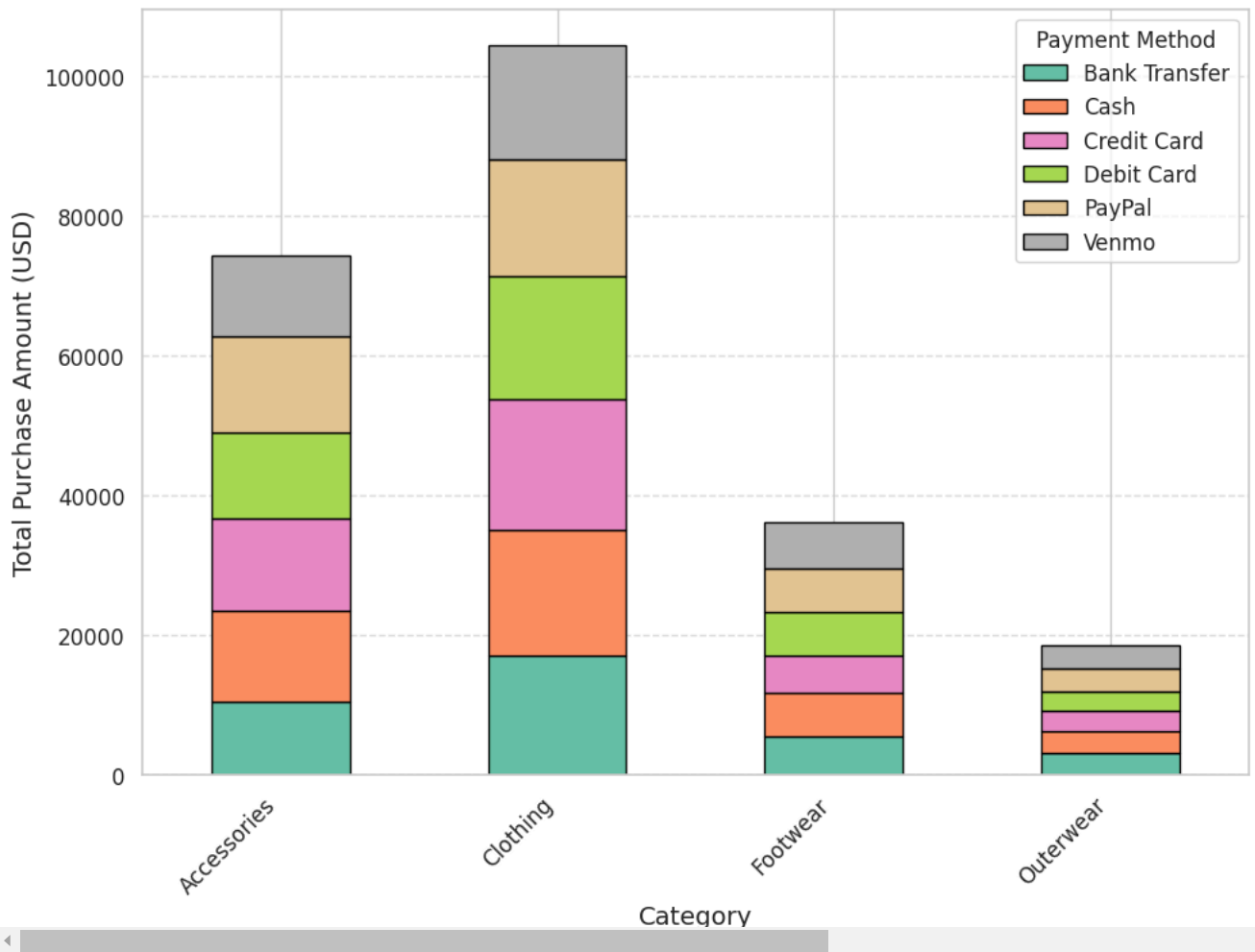
```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3 import pandas as pd
```

### Total Purchase Amount by Category and Payment Method

```
1 sns.set(style="whitegrid")
2
3
4 # Group by 'Category' and 'Payment Method', and sum the 'Purchase Amount (USD)'
5 purchase_by_category_payment = shopping_trends.groupby(['Category', 'Payment Method'], observed=False)['Purchase Amount (USD)'].sum()
6
7 # Plot with customized colors and styling
8 purchase_by_category_payment.plot(kind='bar', stacked=True, figsize=(10, 8), colormap='Set2', edgecolor='black')
9
10 # Add title and labels with improved formatting
11 plt.title('Total Purchase Amount by Category and Payment Method', fontsize=18, fontweight='bold', pad=20)
12 plt.xlabel('Category', fontsize=14)
13 plt.ylabel('Total Purchase Amount (USD)', fontsize=14)
14 plt.xticks(rotation=45, ha='right', fontsize=12)
15 plt.yticks(fontsize=12)
16
17 plt.legend(title='Payment Method', fontsize=12, loc='upper right')
18
19 # Add gridlines for better readability
20 plt.grid(True, axis='y', linestyle='--', alpha=0.6)
21
22 # Adjust layout to prevent overlap and improve spacing
23 plt.tight_layout()
24
25 # Show the plot
26 plt.show()
```



## Total Purchase Amount by Category and Payment Method



- This stacked bar chart visualizes the total purchase amount by category and payment method. Each bar represents a product category, with the total purchase amount split by the different payment methods used, including Bank Transfer, Cash, Credit Card, Debit Card, PayPal, and Venmo.
- The colors in the bars distinguish between these payment methods. The chart shows that Clothing has the highest total purchase amount, with a significant portion paid via Debit Card and Credit Card, while categories like Footwear and Outerwear have comparatively lower totals.
- The legend on the right clarifies which color corresponds to each payment method.

```

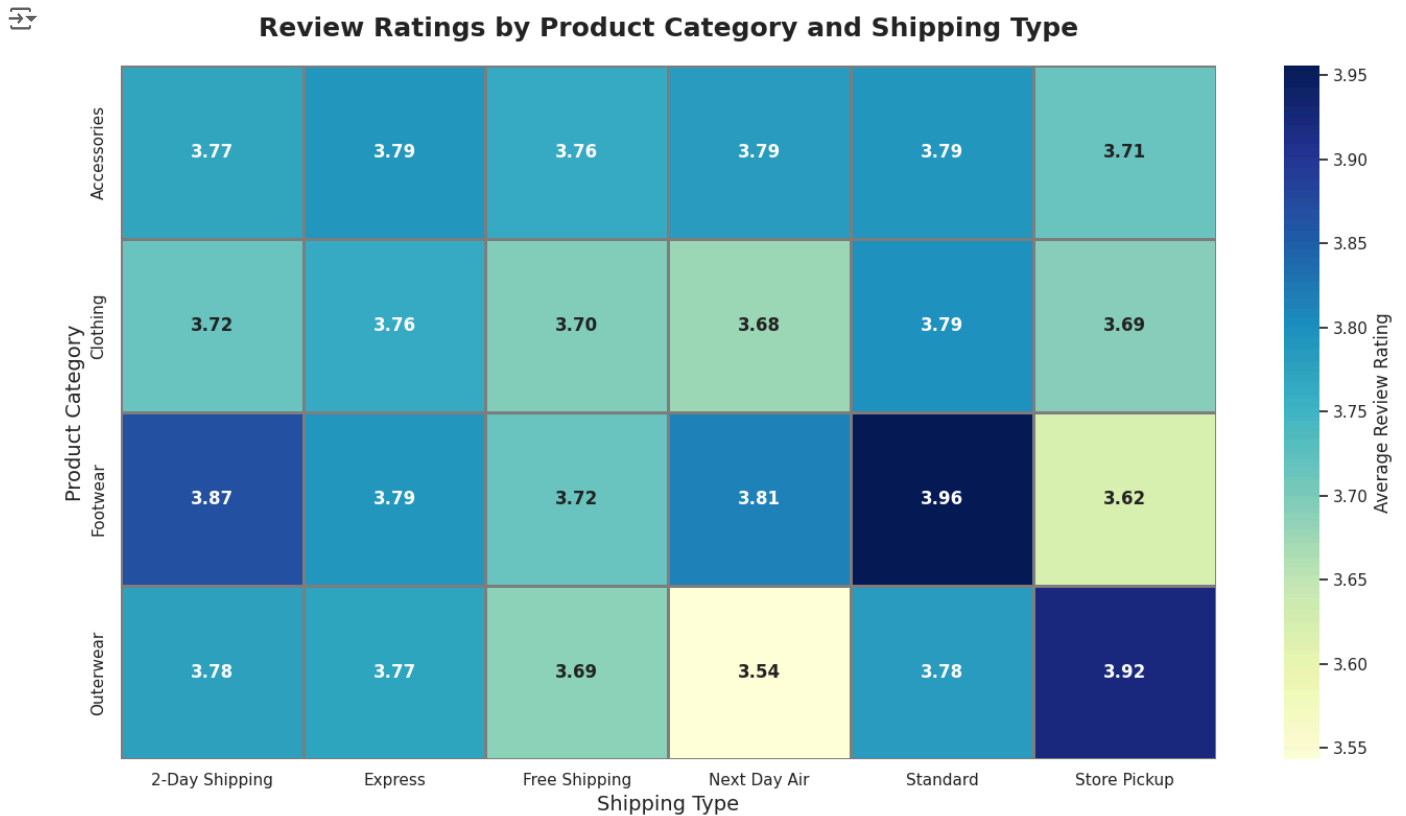
1 avg_review_by_category_shipping = shopping_trends.groupby(['Category', 'Shipping Type'], observed=False)['Review Rating'].mean().unstack()
2
3 # Create the heatmap with customized colors and annotations
4 plt.figure(figsize=(14, 8))
5
6 sns.heatmap(
7     avg_review_by_category_shipping,
8     annot=True,
9     fmt='.2f',
10    cmap='YlGnBu', # Changed to a more visually appealing color palette
11    cbar_kws={'label': 'Average Review Rating'},
12    linewidths=1, # Add separation between cells for better readability
13    linecolor='gray', # Color of the lines between cells
14    annot_kws={'size': 12, 'weight': 'bold'} # Customize font size and weight for annotations
15 )
16
17 # Adding title and labels with improved formatting
18 plt.title('Review Ratings by Product Category and Shipping Type', fontsize=18, fontweight='bold', pad=20)
19 plt.xlabel('Shipping Type', fontsize=14)
20 plt.ylabel('Product Category', fontsize=14)
21
22
23 # Adding gridlines for better readability
24 plt.grid(True, axis='y', linestyle='--', alpha=0.6)
25
26 # Adjust the layout to prevent overlapping elements
27 plt.tight_layout()
28

```

```

29 # Show the plot
30 plt.show()
31

```



- This heatmap visualizes average review ratings for different product categories and shipping types.
- Each cell represents the average review rating for a specific combination of product category (Accessories, Clothing, Footwear, and Outerwear) and shipping type (2-Day Shipping, Express, Free Shipping, Next Day Air, Standard, and Store Pickup).
- The color gradient reflects the average ratings, with darker shades indicating higher ratings and lighter shades indicating lower ratings. For example, Footwear with Next Day Air has the highest average review rating of 3.96, while Outerwear with Free Shipping has the lowest at 3.54.
- This heatmap provides an easy way to compare how different shipping methods influence customer satisfaction across various product categories.

```

1 # Grouping by 'Location' and calculating the average of 'Age', 'Purchase Amount (USD)', and 'Previous Purchases'
2 shopping_trends_final = shopping_trends.groupby('Location').agg({
3     'Age': 'mean',
4     'Purchase Amount (USD)': 'mean',
5     'Previous Purchases': 'mean'
6 }).reset_index()
7
8 shopping_trends_final.head()

```

	Location	Age	Purchase Amount (USD)	Previous Purchases
0	Alabama	44.314607	59.112360	27.449438
1	Alaska	43.000000	67.597222	28.097222
2	Arizona	45.276923	66.553846	28.369231
3	Arkansas	44.101266	61.113924	27.063291
4	California	42.663158	59.000000	24.494737

Next steps:

[Generate code with shopping\\_trends\\_final](#)[View recommended plots](#)[New interactive sheet](#)

```
1 shopping_trends_final = shopping_trends_final.rename(columns={'Location': 'State'})
```

```
1 shopping_trends_final.head()
```

	State	Age	Purchase Amount (USD)	Previous Purchases	
0	Alabama	44.314607	59.112360	27.449438	
1	Alaska	43.000000	67.597222	28.097222	
2	Arizona	45.276923	66.553846	28.369231	
3	Arkansas	44.101266	61.113924	27.063291	
4	California	42.663158	59.000000	24.494737	

Next steps:

[Generate code with shopping\\_trends\\_final](#)[View recommended plots](#)[New interactive sheet](#)

Merge shopping\_trends\_final and Income\_by\_states to to combine the average purchase behavior data (age, purchase amount, previous purchases) from the shopping trends dataset with income-related information

```
1 # Merge the two datasets on the 'State' column
2 merged_data = pd.merge(shopping_trends_final, Income_by_states, how='inner',
3 on='State')
4
5 # Show the first few rows of the merged dataset
6 merged_data.head()
```

	State	Age	Purchase Amount (USD)	Previous Purchases	FIPS	Value (Dollars)	Rank within US (of 52 states)	
0	Alabama	44.314607	59.112360	27.449438	1000	59,609	46	
1	Alaska	43.000000	67.597222	28.097222	2900	86,370	13	
2	Arizona	45.276923	66.553846	28.369231	4000	72,581	24	
3	Arkansas	44.101266	61.113924	27.063291	5000	56,335	49	
4	California	42.663158	59.000000	24.494737	6000	91,905	6	

Next steps:

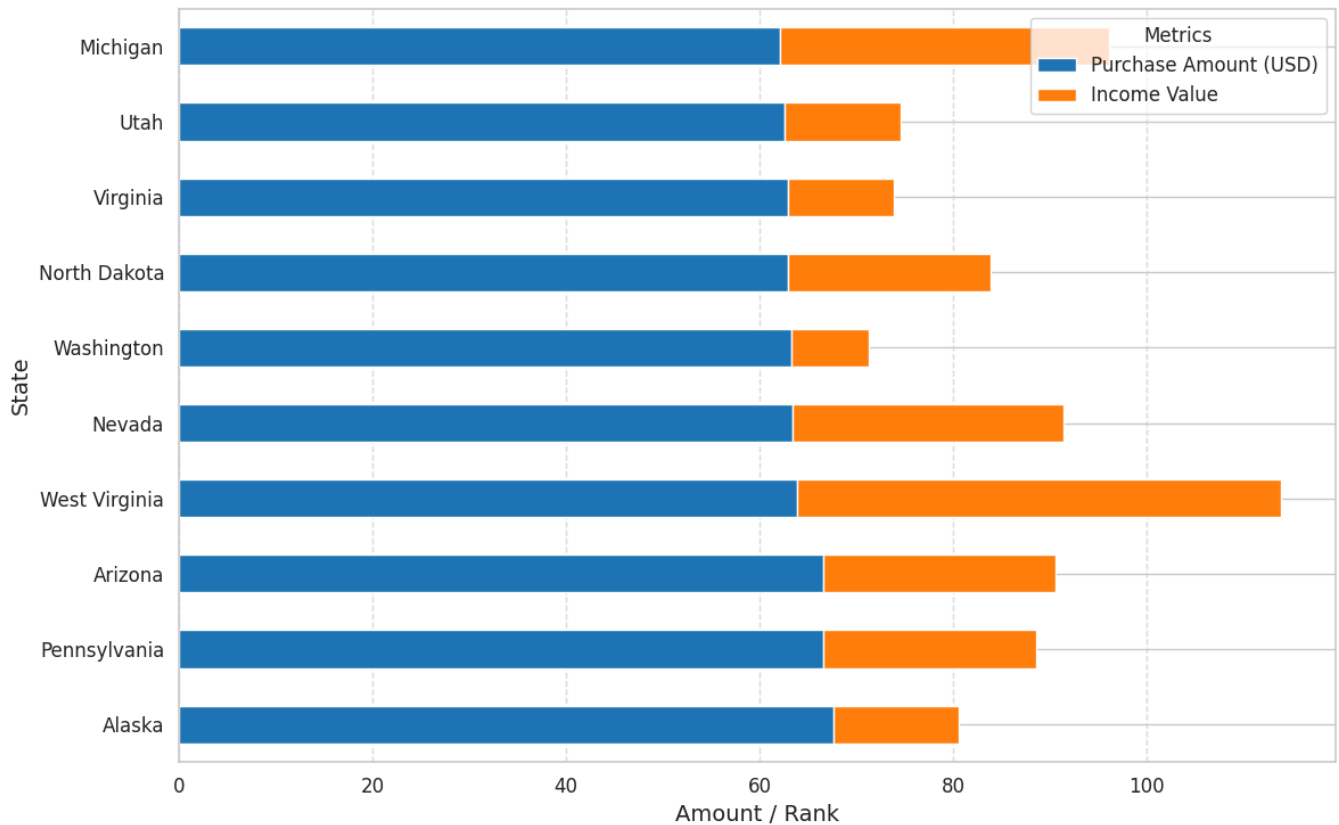
[Generate code with merged\\_data](#)[View recommended plots](#)[New interactive sheet](#)

```
1 # Sort the merged data by 'Purchase Amount (USD)' and select the top 10 states
2 top_10_states = merged_data.sort_values(by='Purchase Amount (USD)', ascending=False).head(10)
3
4 # Plotting the Stacked Bar Chart for Purchase Amount, Income, and Rank
5 top_10_states.set_index('State')[['Purchase Amount (USD)', 'Value (Dollars)', 'Rank within US (of 52 states)']].plot(kind='barh', stacked=True)
6
7 # Adding title and labels with improved formatting
8 plt.title('Top 10 States by Purchase Amount, Income, and Rank', fontsize=18, fontweight='bold', pad=20)
9 plt.xlabel('Amount / Rank', fontsize=14)
10 plt.ylabel('State', fontsize=14)
11 plt.xticks(fontsize=12)
12 plt.yticks(fontsize=12)
13
14 plt.legend(title='Metrics', labels=['Purchase Amount (USD)', 'Income Value', 'Rank'], loc='upper right', fontsize=12)
15
16 # Adding gridlines for better readability
17 plt.grid(True, axis='x', linestyle='--', alpha=0.6)
18
19 # Adding some padding to the layout to avoid congestion
20 plt.tight_layout()
21
22 # Display the chart
23 plt.show()
```





## Top 10 States by Purchase Amount, Income, and Rank



- This bar chart visualizes the Top 10 States by Purchase Amount (USD) and Income Value (Dollars).
- The blue bars represent the Purchase Amount (USD), while the orange bars correspond to the Income Value of each state. The chart allows us to compare the purchasing behavior of these states against their income levels, highlighting which states spend more relative to their income.
- From the chart, it's clear that some states with higher income (such as Michigan and Virginia) also show significant purchase amounts, suggesting a correlation between wealth and spending.
- The legend on the right differentiates between the two metrics for easy comparison.

```

1  geojson_df = pd.json_normalize(geojson['features'])
2
3  merged_data = pd.merge(geojson_df, shopping_trends_final, left_on='properties.
   name', right_on='State', how='left')
4
5  # Create the choropleth map
6  fig = px.choropleth(merged_data,
7                      geojson=geojson,
8                      locations='properties.name',
9                      color='Purchase Amount (USD)',
10                     featureidkey="properties.name",
11                     hover_name='properties.name',
12                     hover_data=['Purchase Amount (USD)'],
13                     color_continuous_scale="bluered")
14
15  fig.update_geos(fitbounds="locations", visible=False)
16
17  # Show the map
18  fig.show()

```

