Northeastern University

**From: Rohit Raman**

**Course: Business Intelligence and Decision support.**

**To: Prof Marcus Ellis**

**Date: 22/02/2024**

**Project: Amazon retail services.**

For this project, I acquired three distinct datasets from Amazon: one encompassing overall sale, another detailing international sale, and a third named "sales report." By amalgamating these datasets through an inner join process, we consolidated them into a comprehensive dataset comprising 1000 rows and 33 columns. To analyse this dataset effectively, I employed various tools, including MySQL for data manipulation, Tableau for advanced visualization, and Python for detailed exploratory data analysis.

Leveraging this extensive dataset, I addressed pertinent inquiries such as: Which apparel items exhibit significant popularity across diverse cities in India, and what corresponding revenue streams do they generate? Additionally, I delved into consumer preferences regarding apparel types, sizes, and colours across different regions. Furthermore, I conducted analyses to ascertain the average pricing of apparel items across various states and to discern the distribution of Amazon's delivery service utilization among its customer base.

In my Python-based analysis, I conducted comprehensive examinations to uncover insights such as the distribution of clothing sales across different seasons, the identification of potential correlations between disparate data points, and the identification of high-volume purchasing states along with their preferred shipping methods.

Using Tableau, I crafted interactive dashboards to address additional business queries, including but not limited to: Trends in average apparel pricing over successive years, identification of the top 7 states by order volume on Amazon, and the examination of monthly apparel purchasing trends. Moreover, I explored metrics such as the popularity of specific apparel items based on order frequency, the monthly revenue generated by Amazon from diverse ethnic groups, and the monthly performance of Amazon's shipping operations.
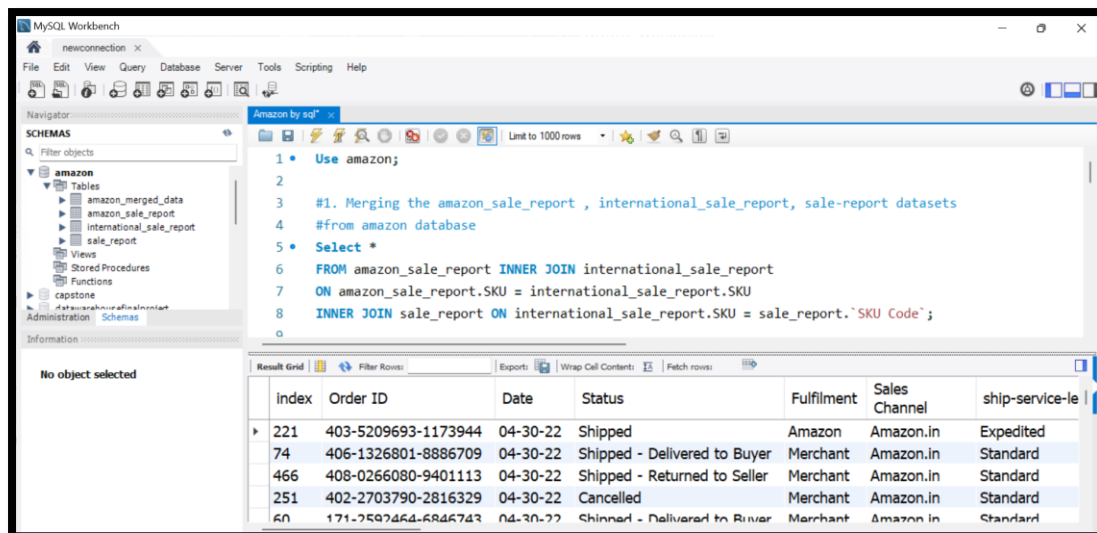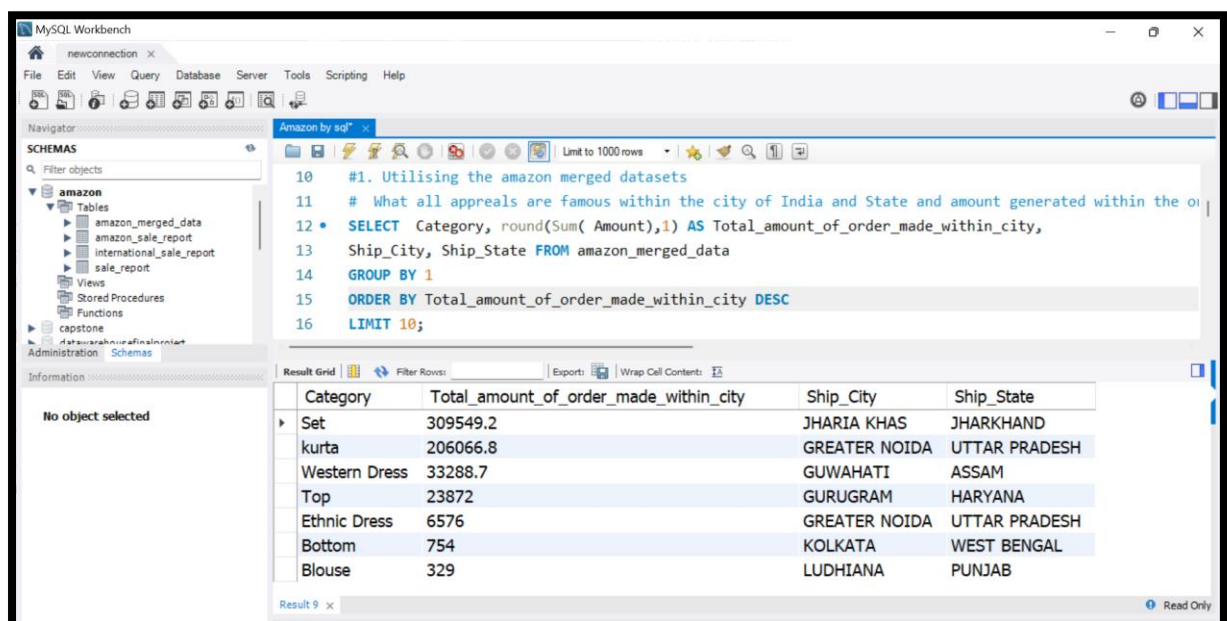
**SQL**

Figure1: SQL code for merging three different datasets of Amazon.

## 1: What all apparels are famous within the city of India and amount generated within the order?



The picture above shows how many orders each city made for their favourite clothes. Sets are really popular in Jharia Khas from Jharkhand, kurtas are the trend in Greater Noida, Uttar Pradesh and lots of people in Guwahati, Assam like Western dresses, ordering them the most.

## 2: What type of clothing are popular in the region, wand what sizes and colour are customers interested in?

In the pictures above, blue emerges as the top-selling color in Uttar Pradesh, with the highest number of units sold. In Telangana, sets in turquoise are the most favoured color among customers. Meanwhile, pink kurtas reign as the preferred choice in Kerala etc.

## 3: What is the average price of apparels in the states?

The pricing of apparel differs from state to state in India. In Rajasthan, ethnic clothing costs Rs 50 more compared to Karnataka. Likewise, sets are priced over Rs 40 higher in West Bengal than in Maharashtra and Telangana, and the pattern continues.

## 4: What percentage of customer have utilised each service model?



Approximately 80% of customers have utilized the expedited Ship Service Level, while the rest have opted for the Standard service level.

**Python**

## Exploratory Data Analysis using Python:



Figure1: Importing the file from Sql to python

```
In [4]: merged_df.shape

Out[4]: (1000, 33)

In [5]: #Type of varaibles in Amazon dataset
        print(merged_df.info())

        <class 'pandas.core.frame.DataFrame'>
        RangeIndex: 1000 entries, 0 to 999
        Data columns (total 33 columns):
         #   Column              Non-Null Count  Dtype
        ---  ------              --------------  -----
         0   index               1000 non-null   int64
         1   Order_Id            1000 non-null   object
         2   Date                1000 non-null   object
         3   Status              1000 non-null   object
         4   Fulfilment          1000 non-null   object
         5   Sales_Channel       1000 non-null   object
         6   Ship-Service-Level  1000 non-null   object
         7   Style               1000 non-null   object
         8   SKU                 1000 non-null   object
         9   Category            1000 non-null   object
         10  Size                1000 non-null   object
         11  ASIN                1000 non-null   object
         12  Courier_Status      981 non-null    object
         13  Qty                 1000 non-null   int64
         14  Currency            1000 non-null   object
         15  Amount              1000 non-null   float64
         16  Ship_City           1000 non-null   object
         17  Ship_State          1000 non-null   object
         18  Ship_Postal_Code    1000 non-null   int64
         19  Ship_Country        1000 non-null   object
         20  Promotion-Ids       608 non-null    object
         21  B2B                 1000 non-null   bool
         22  Fulfilled_By        180 non-null    object
         23  DATE                1000 non-null   object
         24  Months              1000 non-null   object
         25  CUSTOMER            1000 non-null   object
         26  PCS                 1000 non-null   int64
         27  RATE                1000 non-null   float64
         28  GROSS AMT           1000 non-null   int64
         29  Design No.          1000 non-null   object
         30  Stock               1000 non-null   int64
         31  Size.1              1000 non-null   object
         32  Color               1000 non-null   object
        dtypes: bool(1), float64(2), int64(6), object(24)
        memory usage: 251.1+ KB
        None
```
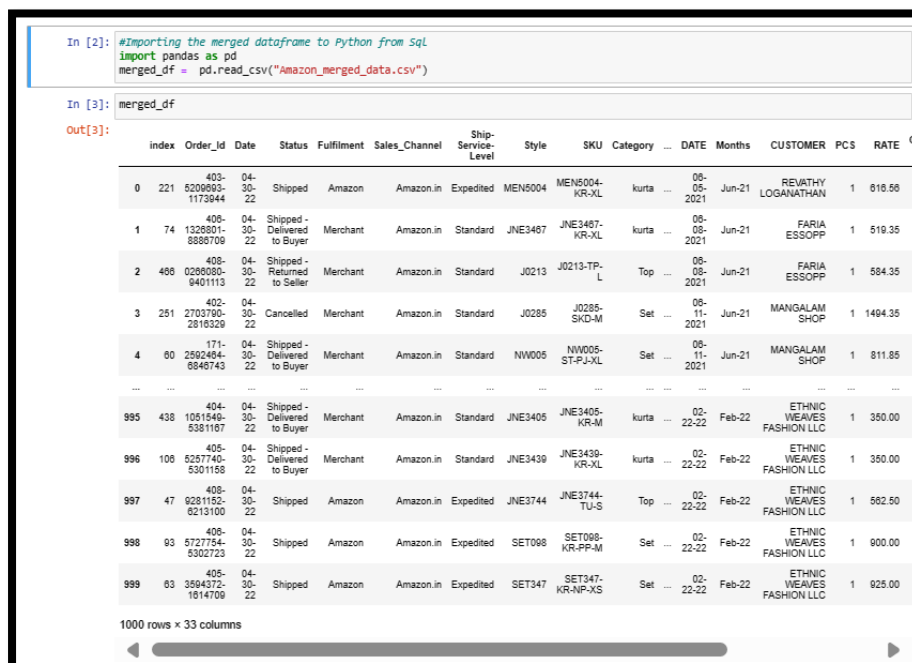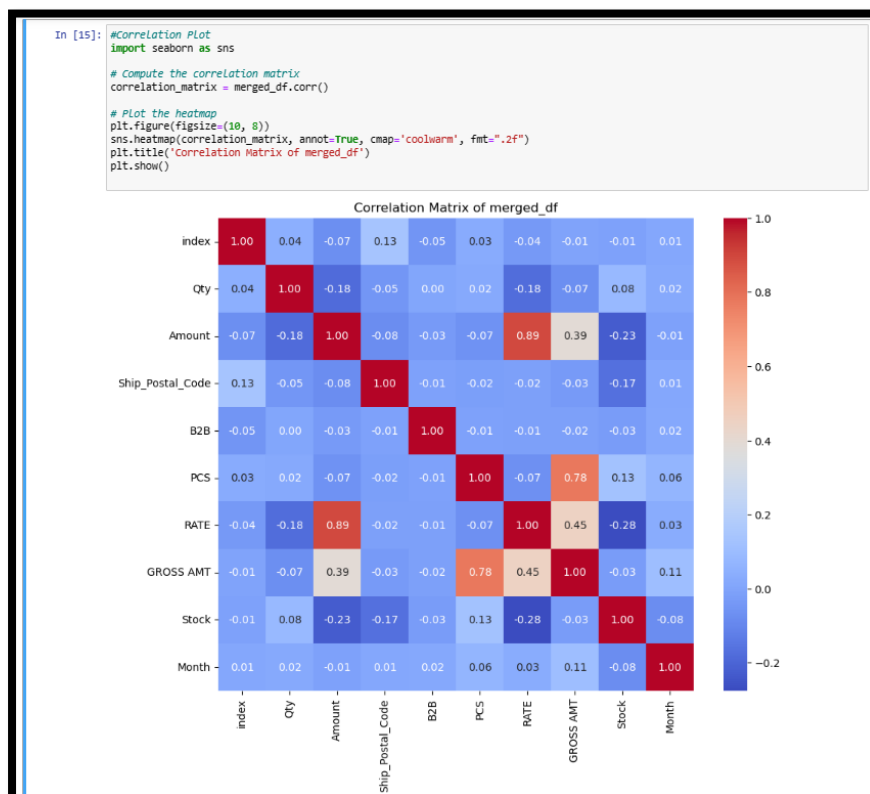
Figure 2: Dimension and type of dataset.

The dataset contains 1000 observations and 33 columns. The majority of these columns represent categorical variables, while others, such as Qty, Ship_postal_code, PCs, and Gross Amt, are of integer type.

## Correlation Plot:

```
In [15]: #Correlation Plot
         import seaborn as sns

         # Compute the correlation matrix
         correlation_matrix = merged_df.corr()

         # Plot the heatmap
         plt.figure(figsize=(10, 8))
         sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt=".2f")
         plt.title('Correlation Matrix of merged_df')
         plt.show()
```



Correlation Matrix of merged_df

From the correlation plot, The Gross amount and PCs exhibit a strong positive correlation, with a coefficient of 0.78. Conversely, Qty and Postal code show a negative correlation, as do Qty and Gross amount. Additionally, Ship_postal_code demonstrates negative correlations with Qty, amount, PCs, rate, and postal code.

## Seasonality Trends:

## 6: what are the sales of different categories of apparels across various season?

To understand the sale of each category of apparels I have performed the seasonality trends to know how the sale performed over the three seasons in India.

```python
#Seasonality trends
import pandas as pd

# Convert 'Date' column to datetime
merged_df['DATE'] = pd.to_datetime(merged_df['DATE'])

# Define a function to map each month to a season
def get_season(month):
    if month in [12, 1, 2]:
        return 'Winter'
    elif month in [3, 4, 5]:
        return 'Spring'
    elif month in [6, 7, 8]:
        return 'Summer'
    else:
        return 'Autumn'

# Extract month and season from the 'Date' column
merged_df['Month'] = merged_df['DATE'].dt.month
merged_df['Season'] = merged_df['Month'].apply(get_season)

# Group by 'Season', 'Month', and 'Category', then sum the 'PCS' column
season_month_category_sales = merged_df.groupby(['Season', 'Month', 'Category'])['PCS'].sum().reset_index()

# Display the table
print(season_month_category_sales)
```

```
      Season  Month        Category  PCS
0     Autumn      9             Set  106
1     Autumn      9             Top   11
2     Autumn      9           kurta  110
3     Autumn     10          Bottom    8
4     Autumn     10   Ethnic Dress    3
5     Autumn     10             Set  160
6     Autumn     10  Western Dress   12
7     Autumn     10           kurta   77
8     Autumn     11          Blouse    1
9     Autumn     11          Bottom    2
10    Autumn     11             Set   38
11    Autumn     11             Top    1
12    Autumn     11  Western Dress    3
13    Autumn     11           kurta   55
14    Summer      6             Set   43
15    Summer      6             Top    9
16    Summer      6  Western Dress   21
17    Summer      6           kurta  155
18    Summer      7   Ethnic Dress    1
19    Summer      7             Set   33
20    Summer      7             Top    6
21    Summer      7  Western Dress   13
22    Summer      7           kurta   57
23    Summer      8             Set   37
24    Summer      8             Top    9
25    Summer      8  Western Dress    3
26    Summer      8           kurta  190
```

```
In [24]: #Seasonality trends
         import pandas as pd

         # Convert 'Date' column to datetime
         merged_df['DATE'] = pd.to_datetime(merged_df['DATE'])

         # Define a function to map each month to a season
         def get_season(month):
             if month in [12, 1, 2]:
                 return 'Winter'
             elif month in [3, 4, 5]:
                 return 'Spring'
             elif month in [6, 7, 8]:
                 return 'Summer'
             else:
                 return 'Autumn'

         # Extract month and season from the 'Date' column
         merged_df['Month'] = merged_df['DATE'].dt.month
         merged_df['Season'] = merged_df['Month'].apply(get_season)

         # Group by 'Season', 'Month', and 'Category', then sum the 'PCS' column
         season_month_category_sales = merged_df.groupby(['Season', 'Month', 'Category'])['PCS'].sum().reset_index()

         # Display the table
         print(season_month_category_sales)
```

```
In [12]: import matplotlib.pyplot as plt

         # Drop the 'Total' column
         season_aggregate_without_total = season_aggregate.drop(columns=['Total'])

         # Transpose the DataFrame twice for easy plotting
         season_aggregate_transposed = season_aggregate_without_total.transpose()
         season_category_transposed = season_aggregate_transposed.transpose()

         # Plot the line graph
         season_category_transposed.plot(kind='line', figsize=(10, 6))
         plt.title('Total PCS Sold by Season for Each Category')
         plt.xlabel('Season')
         plt.ylabel('Total PCS Sold')
         plt.xticks(rotation=45, ha='right')
         plt.legend(title='Category')
         plt.grid(axis='y')

         # Show the plot
         plt.tight_layout()
         plt.show()
```



Total PCS Sold by Season for Each Category

Based on the seasonal trends, it's clear that apparel sales decrease during the winter months. However, kurta sales peak during the summer, recording the highest number of pieces sold. Conversely, sets experience a declining trend, beginning with 300 pieces sold in autumn and dropping to approximately 90 pieces in winter.

## 5: What are the trending ship services on Amazon across different states, and how many pieces were delivered using these services?

```python
# Group by 'Ship_State' and 'Ship-Service-Level', then count unique categories
ship_service_level_counts = merged_df.groupby(['Ship_State', 'Ship-Service-Level'])['Category'].nunique().reset_index()

# Display the result
print(ship_service_level_counts)
```

```
        Ship_State Ship-Service-Level  Category
0    ANDHRA PRADESH          Expedited         4
1    ANDHRA PRADESH           Standard         1
2             ASSAM          Expedited         3
3             BIHAR          Expedited         1
4       CHHATTISGARH         Expedited         1
5       CHHATTISGARH          Standard         2
6             DELHI          Expedited         2
7              GOA          Expedited         2
8              Goa          Expedited         1
9           Gujarat          Expedited         3
10          HARYANA          Expedited         3
11          HARYANA           Standard         4
12 HIMACHAL PRADESH          Expedited         1
13   JAMMU & KASHMIR         Expedited         1
14        JHARKHAND           Standard         1
15        KARNATAKA          Expedited         5
16        KARNATAKA           Standard         2
17           KERALA          Expedited         2
18           KERALA           Standard         2
19   MADHYA PRADESH          Expedited         2
20   MADHYA PRADESH           Standard         1
21      MAHARASHTRA          Expedited         3
22      MAHARASHTRA           Standard         3
23        MEGHALAYA           Standard         2
24         NAGALAND          Expedited         1
25           ODISHA          Expedited         3
26           ODISHA           Standard         1
27       PUDUCHERRY          Expedited         1
28       PUDUCHERRY           Standard         1
29           PUNJAB          Expedited         1
30           PUNJAB           Standard         1
31        RAJASTHAN          Expedited         2
32        RAJASTHAN           Standard         1
33       TAMIL NADU          Expedited         4
34       TAMIL NADU           Standard         2
35        TELANGANA          Expedited         4
36        TELANGANA           Standard         2
37    UTTAR PRADESH          Expedited         5
38    UTTAR PRADESH           Standard         2
39      UTTARAKHAND          Expedited         3
40      UTTARAKHAND           Standard         1
41      WEST BENGAL          Expedited         3
42      WEST BENGAL           Standard         3
```
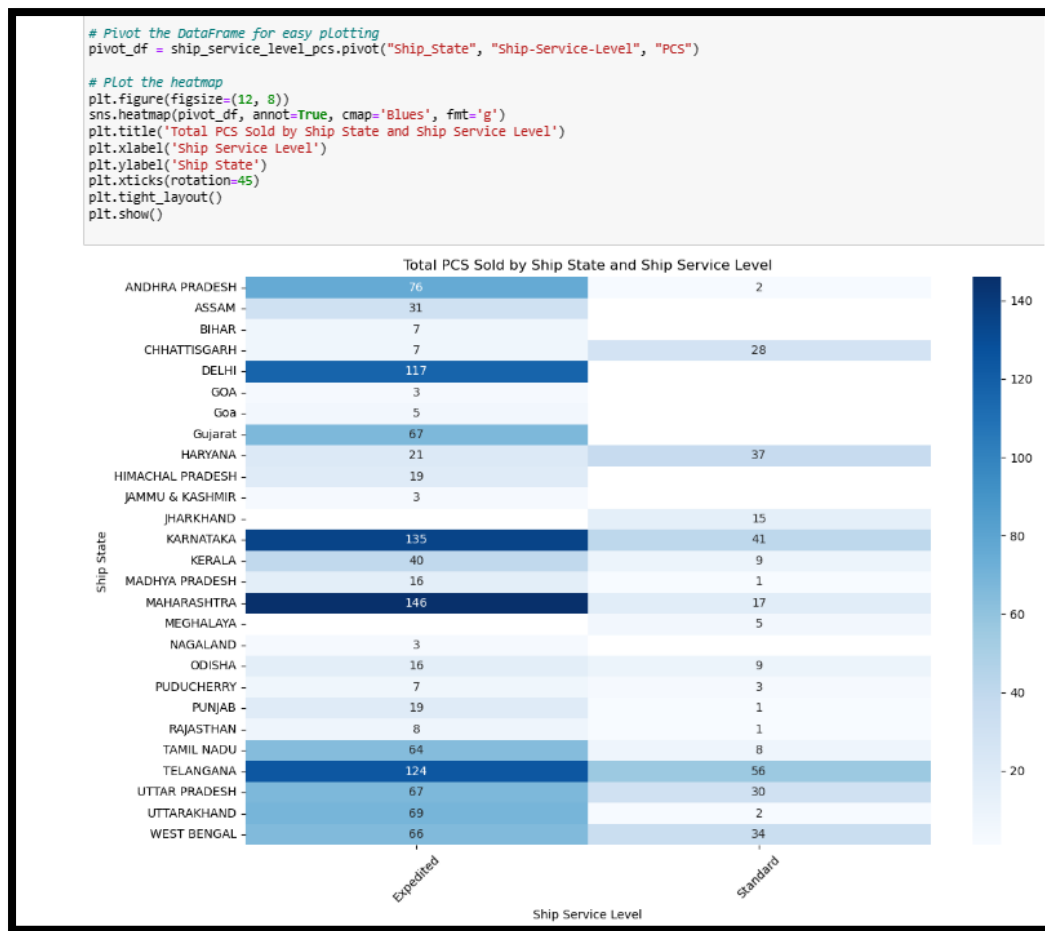
```python
# Group by 'Ship_State' and 'Ship-Service-Level', then sum the PCS
ship_service_level_pcs = merged_df.groupby(['Ship_State', 'Ship-Service-Level'])['PCS'].sum().reset_index()

# Display the result
print(ship_service_level_pcs)
```

```
        Ship_State Ship-Service-Level  PCS
0    ANDHRA PRADESH          Expedited   76
1    ANDHRA PRADESH           Standard    2
2             ASSAM          Expedited   31
3             BIHAR          Expedited    7
4       CHHATTISGARH         Expedited    7
5       CHHATTISGARH          Standard   28
6             DELHI          Expedited  117
7              GOA          Expedited    3
8              Goa          Expedited    5
9           Gujarat          Expedited   67
10          HARYANA          Expedited   21
11          HARYANA           Standard   37
12 HIMACHAL PRADESH          Expedited   19
13   JAMMU & KASHMIR         Expedited    3
14        JHARKHAND           Standard   15
15        KARNATAKA          Expedited  135
16        KARNATAKA           Standard   41
17           KERALA          Expedited   40
18           KERALA           Standard    9
19   MADHYA PRADESH          Expedited   16
20   MADHYA PRADESH           Standard    1
21      MAHARASHTRA          Expedited  146
22      MAHARASHTRA           Standard   17
23        MEGHALAYA           Standard    5
24         NAGALAND          Expedited    3
25           ODISHA          Expedited   16
26           ODISHA           Standard    9
27       PUDUCHERRY          Expedited    7
28       PUDUCHERRY           Standard    3
29           PUNJAB          Expedited   19
30           PUNJAB           Standard    1
31        RAJASTHAN          Expedited    8
32        RAJASTHAN           Standard    1
33       TAMIL NADU          Expedited   64
34       TAMIL NADU           Standard    8
35        TELANGANA          Expedited  124
36        TELANGANA           Standard   56
37    UTTAR PRADESH          Expedited   67
38    UTTAR PRADESH           Standard   30
39      UTTARAKHAND          Expedited   69
40      UTTARAKHAND           Standard    2
41      WEST BENGAL          Expedited   66
42      WEST BENGAL           Standard   34
```
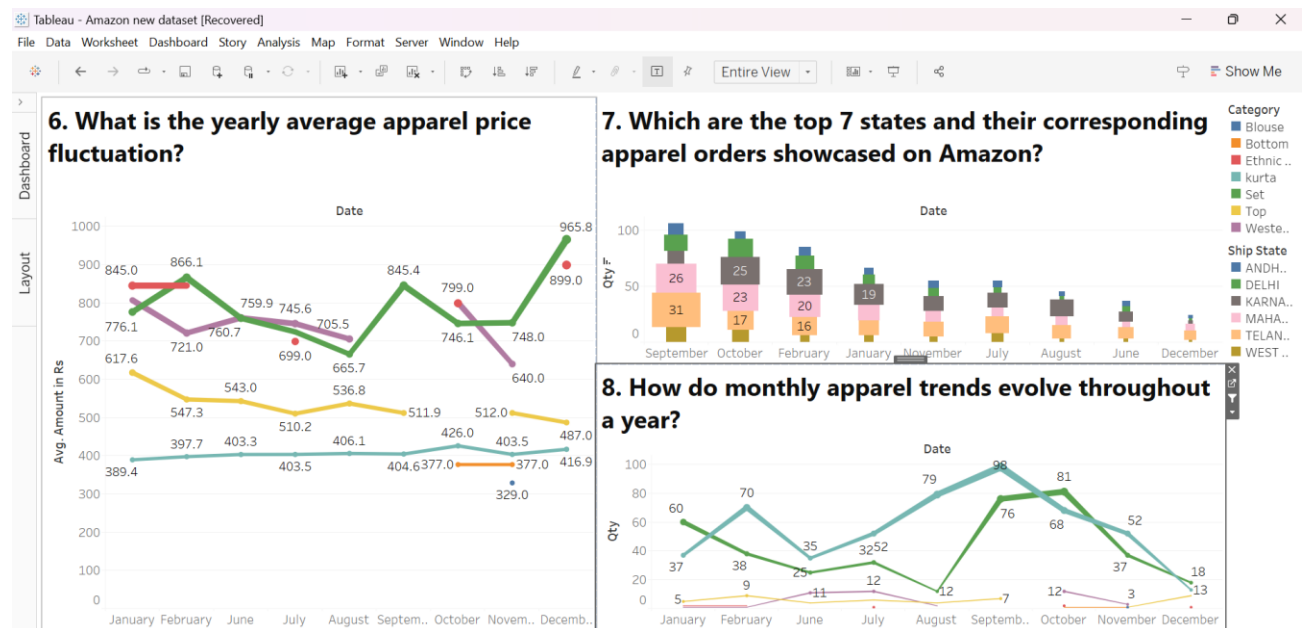
```
# Pivot the DataFrame for easy plotting
pivot_df = ship_service_level_pcs.pivot("Ship_State", "Ship-Service-Level", "PCS")

# Plot the heatmap
plt.figure(figsize=(12, 8))
sns.heatmap(pivot_df, annot=True, cmap='Blues', fmt='g')
plt.title('Total PCS Sold by Ship State and Ship Service Level')
plt.xlabel('Ship Service Level')
plt.ylabel('Ship State')
plt.xticks(rotation=45)
plt.tight_layout()
plt.show()
```



Total PCS Sold by Ship State and Ship Service Level

Based on the analysis above, Telangana, Maharashtra, Delhi, Karnataka, and Andhra Pradesh are the states utilizing the highest number of expedited services, while states like Goa, Himachal Pradesh, Nagaland, Gujarat, Assam, and Bihar have not opted for the standard ship service. Investigating the reasons behind their choices is necessary.

## Tableau: Data Visualisation



### Dashboard1:

From the dashboard provided above, some questions were addressed using visualization, such as:

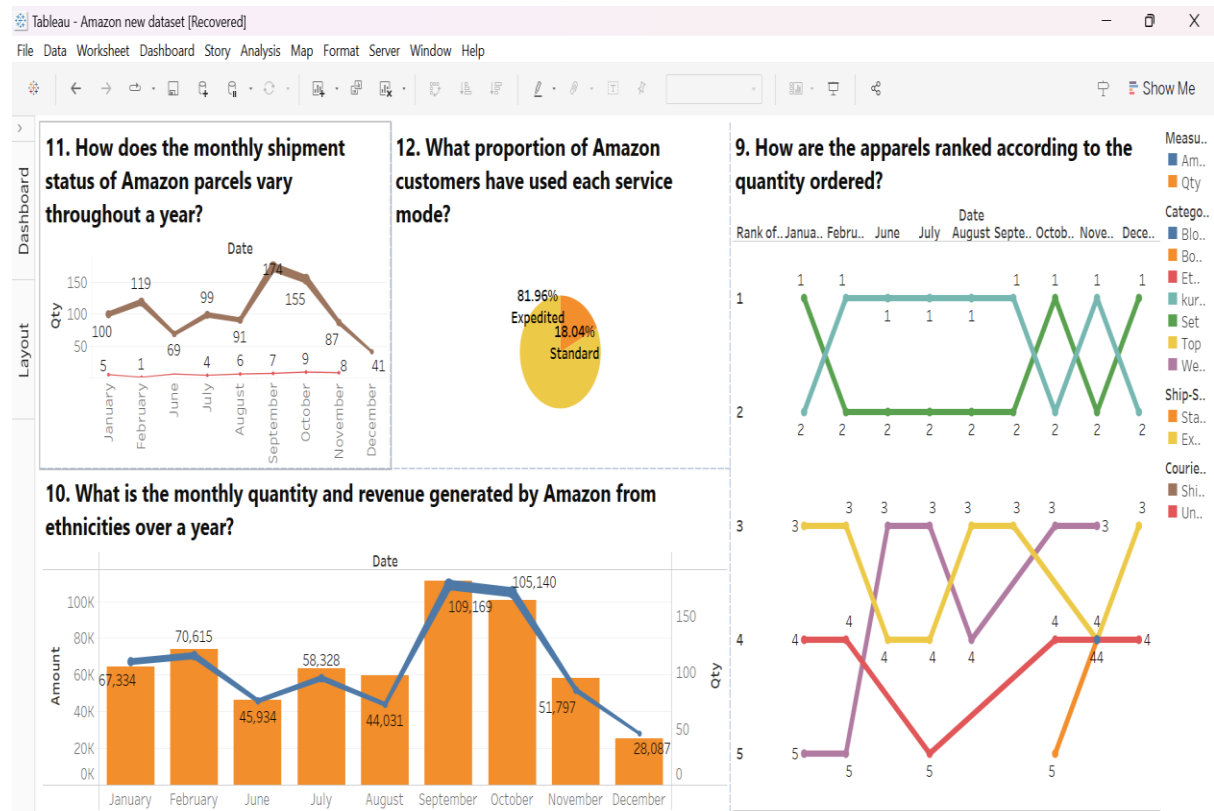### 6. What is the yearly average apparel price fluctuation?

Based on the line plot above, it's clear that **Set** has experienced the most significant fluctuation over the years. Regarding apparel, the average price of the **Top** has declined consistently over the years, possibly influenced by seasonal factors affecting sales prices. However, the price of **Kurta** remained stable throughout the years.

### 7. Which are the top 7 states and their corresponding apparel orders showcased on Amazon?

Based on the chart, it can be inferred that **Maharashtra, Telangana, Karnataka**, and **Delhi** are the top states in terms of the number of orders placed across all months.

### 8. How do monthly apparel trends evolve throughout a year?

**Kurta** consistently outperforms other apparel types in terms of the number **of units sold** each month, with significantly higher **sales figures**. **Set** ranks second in popularity among trending **apparel,** with consistent **purchases** throughout the year, peaking notably in **September, October, and November**.

## Dashboard2

**From dashboard 2**

### 9. How are the apparels ranked according to the quantity ordered?

Based on the preceding line graph, both sets and kurtas consistently maintained the top two positions in preference over the years. Conversely, tops and ethnic dresses consistently held the third and fourth positions, respectively, in terms of total quantity ordered.

### 10. What is the monthly quantity and revenue generated by Amazon from ethnicities over multiple years?

Based on the dual chart, it's clear that Amazon achieves its peak revenue in February, September, October, and November, likely due to increased sales volume during these months. Conversely, December experiences the lowest sales and orders for Amazon.

## 11. How does the monthly shipment status of Amazon parcels vary throughout a year?

According to the line chart, the highest volume of parcels was shipped during **February, September, and October**, corresponding to increased order numbers. Conversely, **these months** also exhibit the highest **count of unshipped orders**.

## 12. What proportion of Amazon customers have used each service mode?

The preferred delivery service option among customers is **Expedited service**, with a significant adoption rate of **81.96%,** followed by the less popular **Standard option**, chosen by only **18.04%** of customers

## Summary:

Amazon can enhance its operational efficiency and sales by implementing the following recommendations:

**Regional Targeting: Amazon** should focus on **stocking and promoting apparel** types that are popular in specific regions. For example, in **Jharia Khas, Jharkhand,** where **sets** are highly **favored**, Amazon can ensure a robust supply of sets to meet local demand. Similarly, in **Greater Noida, Uttar Pradesh**, where **kurtas** are trending, Amazon can highlight and promote **kurta** collections to attract more customers.

**Customized Offerings:** Understanding customer preferences regarding sizes and colors is crucial. Amazon should tailor its offerings based on regional preferences. For instance, in Uttar Pradesh, blue-colored apparel is popular, so Amazon can prioritize stocking more blue-colored items in that region.

**Dynamic Pricing Strategies: Since** apparel prices vary across states, Amazon can employ dynamic pricing strategies to remain competitive. Offering competitive prices in regions where prices are comparatively higher can attract more customers.

**Service Mode Optimization:** Amazon should focus on optimizing its delivery service modes based on customer preferences. Since **expedited service** is **highly favoured**, **Amazon can invest more resources in improving and expanding its expedited delivery network** to provide faster and more reliable service.

**Seasonal Marketing:** Leveraging seasonal trends in apparel sales can be beneficial. For instance, since **kurta** sales peak during the **summer**, <span style="color:red">**Amazon can run marketing campaigns specifically targeting kurta enthusiasts during this period**</span>.

By implementing these recommendations, Amazon can improve its operational efficiency, cater to customer preferences more effectively, and ultimately increase sales and revenue.