

# Auto Sales Data Visualization: Project Documentation

## Project Overview

The **Auto Sales Data Visualization** project aims to analyze automotive sales data to uncover key insights, trends, and patterns. By using advanced data analytics and visualization tools, this project will enable stakeholders to understand the underlying factors affecting automotive sales, identify market trends, and make informed business decisions based on data-driven insights.

## Objectives

The main objectives of this project are:

- **Analyze Various Dimensions:** Investigate key aspects of the dataset, such as monthly trends, sales distribution by model, brand popularity, and geographical variations in sales.
- **Create Interactive Visualizations:** Build interactive charts and graphs that allow stakeholders to explore the data and uncover insights easily.
- **Identify Sales Influencing Factors:** Examine potential factors affecting sales performance, such as product pricing, vehicle type, and regional market preferences.

## Dataset Overview

The dataset used in this project is sourced from Kaggle: Auto Sales Data. It contains detailed records of automobile sales, offering various attributes that describe each transaction.

## Dataset Summary

The dataset contains **2,747 entries** and **20 columns**. Here is a breakdown of each column in the dataset:

1. **ORDERNUMBER**: A unique identifier for each order.
2. **QUANTITYORDERED**: The quantity of the product ordered.
3. **PRICEEACH**: The price per individual item.
4. **ORDERLINENUMBER**: A line number within the order indicating the sequence of items.
5. **SALES**: The total sales amount for the order line.
6. **ORDERDATE**: The date the order was placed.
7. **DAYS\_SINCE\_LASTORDER**: The number of days since the customer's last order.
8. **STATUS**: The status of the order (e.g., Shipped, On Hold, Canceled).
9. **PRODUCTLINE**: The category of the product (e.g., Classic Cars, Motorcycles).
10. **MSRP**: Manufacturer's suggested retail price for the product.
11. **PRODUCTCODE**: A unique code identifying the product.
12. **CUSTOMERNAME**: Name of the customer placing the order.
13. **PHONE**: Contact phone number of the customer.
14. **ADDRESSLINE1**: Primary address line of the customer.
15. **CITY**: The city where the customer is located.
16. **POSTALCODE**: Postal code of the customer's location.
17. **COUNTRY**: The country where the customer is located.
18. **CONTACTLASTNAME**: Last name of the contact person for the customer.
19. **CONTACTFIRSTNAME**: First name of the contact person for the customer.

## Libraries Used

The following Python libraries were employed in this project to perform data analysis and create visualizations:

- **Pandas:** Used for data manipulation and analysis, allowing efficient handling of data frames and performing operations like filtering, grouping, and aggregation.
- **Matplotlib:** A library for creating static data visualizations, which helped in generating bar charts, line plots, and histograms.
- **Seaborn:** Built on top of Matplotlib, Seaborn facilitates statistical data visualization and is particularly useful for visualizing relationships and distributions.
- **Plotly:** A tool for creating interactive visualizations that can be embedded or shared online, providing a more engaging experience for users.
- **Jupyter Notebook:** Used as the environment for conducting the analysis interactively, allowing for real-time visualizations and documenting the entire workflow.

## Data Preprocessing

Data preprocessing is a crucial step to ensure the dataset is clean, consistent, and ready for analysis. The following preprocessing steps were applied:

### 1. Data Type Adjustments

To facilitate efficient analysis, certain columns were converted to appropriate data types:

- **ORDERDATE:** Initially stored as a string, the ORDERDATE column was converted to a datetime type. This allows for time-based operations, such as calculating trends over time and grouping data by year or month.
- **PRICEEACH and SALES:** Both columns were initially stored as strings but were converted to numeric types to allow for accurate calculations like total sales and averages.

## 2. Column Renaming and Standardization

To ensure consistency and readability, the following changes were made:

- **Renaming for Clarity:** Columns with spaces were renamed for consistency. For instance, "PRICE EACH" was changed to PRICEEACH and "ORDER DATE" was renamed to ORDERDATE.
- **Case Standardization:** All column names were converted to upper-case and formatted without spaces to reduce errors and ensure uniformity throughout the dataset.

## 3. Missing Value Handling

Upon inspection, no missing values were found in the dataset. This ensured that the data was complete and ready for analysis without the need for additional handling like imputation or removal of rows.

## 4. Outlier Detection

Outliers were identified and handled to ensure the integrity of the dataset:

- **Outlier Identification:** The SALES and PRICEEACH columns were checked for values that deviated significantly from the norm, such as negative prices or extremely high values.
- **Handling Outliers:** Invalid data points (e.g., negative prices) were removed, while valid high-priced transactions (e.g., luxury cars) were retained.

## Exploratory Data Analysis (EDA) Summary

Exploratory data analysis (EDA) was conducted to uncover important trends and relationships in the dataset. Below are the key findings:

- **Price and MSRP Relationship:** A strong positive correlation exists between PRICEEACH and MSRP. This suggests that products with higher prices tend to have higher manufacturer-recommended prices.
- **Product Line Sales Trends:** "Classic Cars" dominated in sales, particularly in 2019. Other categories like "Vintage Cars" saw an increase in sales during the same period, while "Trains" consistently had low sales.

- **Seasonal and Annual Sales Trends:** Sales showed clear seasonal patterns, with notable peaks in November of 2019. 2020 experienced an early surge in sales, followed by a steep decline.
- **Regional Sales Distribution:** North America and Europe accounted for the majority of sales, while Africa and Australia had the lowest shares.
- **Quarterly Revenue Patterns:** Revenue consistently increased each quarter, with Q4 showing the highest sales, especially in 2019.
- **Country-wise Sales Performance:** The USA, Canada, and Germany were the top-performing countries in terms of sales volume.

## Descriptive Statistics

- **QUANTITYORDERED:** Mean of 35.1 and a standard deviation of 9.76.
- **SALES:** The average total sales amounted to \$3,553, with a maximum value reaching \$14,082.

## Key Visualizations and Insights

### 4.1 Regression Plot of Price vs. MSRP

**Graph Type:** Scatter Plot with Regression Line

**Description:** Visualizes the positive correlation between PRICEEACH and MSRP, indicating that higher-priced products generally align with higher MSRP.

**Insight:** This relationship is crucial for ensuring pricing strategies are consistent with manufacturer guidelines.

### 4.2 Total Sales by Product Line (by Year)

**Graph Type:** Bar Chart

**Description:** Displays the total sales for each product line from 2018 to 2020.

**Insight:** "Classic Cars" had the highest sales in 2019, while "Trains" had the lowest consistently.

### 4.3 Pair Plot of Key Variables

**Graph Type:** Pair Plot

**Description:** A matrix of scatter plots and histograms, visualizing relationships between PRICEEACH, QUANTITYORDERED, SALES, and MSRP.

**Insight:** Strong correlations between price, quantity, and sales are evident.

### 4.4 Monthly Sales Trend by Year

**Graph Type:** Line Graph

**Description:** The line graph tracks total sales across the months of 2018, 2019, and 2020.

**Insights:**

- Stable sales were recorded in 2018 with a peak in November.
- Gradual sales increase in 2019, peaking again in November.
- A significant sales surge occurred in early 2020, peaking in May, but sales dropped to zero from June onwards.

### 4.5 Monthly Trends in Deal Size

**Graph Type:** Line Graph

**Description:** This graph categorizes deals into Small, Medium, and Large sizes, showing their trends over time.

**Insights:**

- Significant spikes in Small and Medium deals in October of 2018 and 2019.
- Large deals remained low and stable over time.
- A recurring pattern of spikes in October, with sharp drops in November for Small and Medium deals.

### 4.6 Distribution of Sales Across Regions

**Graph Type:** Pie Chart

**Description:** The pie chart shows the sales distribution across different geographical regions: North America, Europe, Asia, South America, Africa, and Australia.

**Insights:**

- North America and Europe accounted for the highest percentages of sales.
- Africa and Australia recorded the lowest sales percentages.

## 4.7 Customer Satisfaction Ratings by Product Line

**Graph Type:** Bar Chart

**Description:** The bar chart shows customer satisfaction ratings for various product lines, including Classic Cars, Vintage Cars, Motorcycles, Planes, Ships, Trucks and Buses, and Trains.

**Insights:**

- Classic Cars and Vintage Cars received the highest customer satisfaction ratings.
- Trains and Trucks and Buses had the lowest satisfaction ratings.

## 4.8 Quarterly Revenue Comparison

**Graph Type:** Stacked Bar Chart

**Description:** The stacked bar chart compares quarterly revenues for the years 2018, 2019, and 2020.

**Insights:**

- Revenue increased consistently from Q1 to Q4 in each year.
- The highest revenue was recorded in Q4 each year, with a notable peak in Q4 2019.

## 4.9 Heatmap of Correlation Between Variables

**Graph Type:** Heatmap

**Description:** The heatmap shows correlation coefficients between various variables in the dataset.

**Insights:**

- There is a strong positive correlation between `PRICE_EACH` and `MSRP`, as well as between `QUANTITY_ORDERED` and `TOTAL_PRICE`.
- Some variable pairs showed weak or no significant correlation.

## 4.10 Sales Performance by Country

**Graph Type:** Geographical Heat Map

**Description:** The heat map illustrates the sales performance across different countries.

**Insights:**

- The highest sales were recorded in the USA, Canada, and Germany.
- Darker colors indicate higher sales volumes, representing geographic hotspots for sales.

## Challenges and Limitations

### 1. Data Completeness

While no missing values were detected, the dataset may still contain incomplete records that impact the accuracy of certain analyses. This may include incorrect or inconsistent data entries that could influence some of the insights.

### 2. Outliers

Outliers, particularly in pricing and quantities, can significantly affect statistical analysis. While efforts were made to remove any invalid data points (such as negative prices), some extreme values may still exist. These outliers can distort the understanding of overall trends and should be carefully handled in future analyses.

### 3. Scalability

The dataset used in this project is relatively small and manageable, allowing for quick analysis. However, for larger and more complex datasets, optimizations and scalable methods would be required to ensure efficient data processing and visualization. Methods such as distributed computing or sampling may be necessary for handling large-scale datasets.

## Conclusion

This analysis of the auto sales dataset has provided critical insights into the automotive industry, identifying key trends, correlations, and geographical performance. Notable findings include:



- A strong positive correlation between PRICEEACH and MSRP.
- Seasonal sales trends, with notable peaks in the last quarter of each year.
- Regional performance, with the USA, Canada, and Germany leading in total sales.

Through interactive visualizations and a thorough exploratory data analysis, stakeholders are equipped with a clearer understanding of the factors driving automotive sales. However, some challenges remain, such as outlier handling and scalability concerns, which could impact the analysis if more extensive datasets are used in the future.

Overall, this project serves as a foundation for deeper analysis and decision-making in the automotive industry, providing valuable data-driven insights to help businesses optimize their sales strategies.

## References

**Pandas Documentation:** <https://pandas.pydata.org/pandas-docs/stable/>

**Matplotlib Documentation:** <https://matplotlib.org/stable/contents.html>

**Seaborn Documentation:** <https://seaborn.pydata.org/>

**Plotly Documentation:** <https://plotly.com/python/>