

# BUSINESS CHALLENGE: DATA CLEANING AND EDA

Presented by: Natalia Domínguez Jiménez





# OVERVIEW

01

Problem  
Introduction

02

Dataset  
Overview

03

Project  
Focus

04

Data Cleaning

05

EDA  
(Exploratory Data  
Analysis)

06

Conclusions

01

# PROBLEM INTRODUCTION

Sales performance is a critical metric for any business. Understanding the factors that influence it enables informed decision-making, strategy optimization, and ultimately, business growth.

This project aims to identify the key drivers of sales performance by analyzing a comprehensive sales dataset.

## 02 DATASET OVERVIEW

The dataset contains over 2,800 sales records documenting various transactions of classic and specialty vehicles.

It includes information such as product quantity and price, total sales amount, and the date of each transaction. Additionally, it provides details about customers and their geographical locations.

Order ID	Customer ID	Product ID	Quantity	Unit Price	Total Price	Shipped Date	Status	QTR_ID	MONTH_ID	YEAR_ID	Category
1000000000	1	1	2	2003	Motorcycles	2003-01-01 00:00:00	Shipped	1	2	2003	Motorcycles
1000000001	2	2	5	2003	Motorcycles	2003-02-01 00:00:00	Shipped	2	5	2003	Motorcycles
1000000002	3	3	7	2003	Motorcycles	2003-03-01 00:00:00	Shipped	3	7	2003	Motorcycles
1000000003	4	4	8	2003	Motorcycles	2003-04-01 00:00:00	Shipped	4	8	2003	Motorcycles
1000000004	5	5	10	2003	Motorcycles	2003-05-01 00:00:00	Shipped	4	10	2003	Motorcycles
1000000005	6	6	10	2003	Motorcycles	2003-06-01 00:00:00	Shipped	4	10	2003	Motorcycles
1000000006	7	7	11	2003	Motorcycles	2003-07-01 00:00:00	Shipped	4	11	2003	Motorcycles
1000000007	8	8	11	2003	Motorcycles	2003-08-01 00:00:00	Shipped	4	11	2003	Motorcycles
1000000008	9	9	12	2003	Motorcycles	2003-09-01 00:00:00	Shipped	4	12	2003	Motorcycles
1000000009	10	10	1	2003	Motorcycles	2003-10-01 00:00:00	Shipped	4	1	2003	Motorcycles
1000000010	11	11	1	2003	Motorcycles	2003-11-01 00:00:00	Shipped	4	1	2003	Motorcycles
1000000011	12	12	1	2003	Motorcycles	2003-12-01 00:00:00	Shipped	4	1	2003	Motorcycles
1000000012	13	13	1	2004	Motorcycles	2004-01-01 00:00:00	Shipped	1	1	2004	Motorcycles
1000000013	14	14	1	2004	Motorcycles	2004-02-01 00:00:00	Shipped	2	6	2004	Motorcycles
1000000014	15	15	1	2004	Motorcycles	2004-03-01 00:00:00	Shipped	3	7	2004	Motorcycles
1000000015	16	16	1	2004	Motorcycles	2004-04-01 00:00:00	Shipped	4	8	2004	Motorcycles
1000000016	17	17	1	2004	Motorcycles	2004-05-01 00:00:00	Shipped	4	9	2004	Motorcycles
1000000017	18	18	1	2004	Motorcycles	2004-06-01 00:00:00	Shipped	4	10	2004	Motorcycles
1000000018	19	19	1	2004	Motorcycles	2004-07-01 00:00:00	Shipped	4	11	2004	Motorcycles
1000000019	20	20	1	2004	Motorcycles	2004-08-01 00:00:00	Shipped	4	11	2004	Motorcycles
1000000020	21	21	1	2004	Motorcycles	2004-09-01 00:00:00	Shipped	4	12	2004	Motorcycles
1000000021	22	22	1	2004	Motorcycles	2004-10-01 00:00:00	Shipped	4	1	2004	Motorcycles
1000000022	23	23	1	2004	Motorcycles	2004-11-01 00:00:00	Shipped	4	1	2004	Motorcycles
1000000023	24	24	1	2004	Motorcycles	2004-12-01 00:00:00	Shipped	4	1	2004	Motorcycles
1000000024	25	25	1	2005	Motorcycles	2005-01-01 00:00:00	Shipped	1	2	2005	Motorcycles
1000000025	26	26	1	2005	Motorcycles	2005-02-01 00:00:00	Shipped	1	3	2005	Motorcycles
1000000026	27	27	1	2005	Motorcycles	2005-03-01 00:00:00	Shipped	2	4	2005	Motorcycles
1000000027	28	28	1	2005	Motorcycles	2005-04-01 00:00:00	Shipped	2	5	2005	Motorcycles
1000000028	29	29	1	2005	Motorcycles	2005-05-01 00:00:00	Shipped	2	6	2005	Motorcycles
1000000029	30	30	1	2005	Motorcycles	2005-06-01 00:00:00	Shipped	2	7	2005	Motorcycles
1000000030	31	31	1	2005	Motorcycles	2005-07-01 00:00:00	Shipped	2	8	2005	Motorcycles
1000000031	32	32	1	2005	Motorcycles	2005-08-01 00:00:00	Shipped	2	9	2005	Motorcycles
1000000032	33	33	1	2005	Motorcycles	2005-09-01 00:00:00	Shipped	2	10	2005	Motorcycles
1000000033	34	34	1	2005	Motorcycles	2005-10-01 00:00:00	Shipped	2	11	2005	Motorcycles
1000000034	35	35	1	2005	Motorcycles	2005-11-01 00:00:00	Shipped	2	11	2005	Motorcycles
1000000035	36	36	1	2005	Motorcycles	2005-12-01 00:00:00	Shipped	2	1	2005	Motorcycles
1000000036	37	37	1	2006	Classic Cars	2006-01-01 00:00:00	Shipped	1	1	2006	Classic Cars
1000000037	38	38	1	2006	Classic Cars	2006-02-01 00:00:00	Shipped	1	3	2006	Classic Cars
1000000038	39	39	1	2006	Classic Cars	2006-03-01 00:00:00	Shipped	2	5	2006	Classic Cars
1000000039	40	40	1	2006	Classic Cars	2006-04-01 00:00:00	Shipped	3	7	2006	Classic Cars
1000000040	41	41	1	2006	Classic Cars	2006-05-01 00:00:00	Shipped	3	9	2006	Classic Cars
1000000041	42	42	1	2006	Classic Cars	2006-06-01 00:00:00	Shipped	3	11	2006	Classic Cars
1000000042	43	43	1	2006	Classic Cars	2006-07-01 00:00:00	Shipped	3	11	2006	Classic Cars
1000000043	44	44	1	2006	Classic Cars	2006-08-01 00:00:00	Shipped	3	12	2006	Classic Cars
1000000044	45	45	1	2006	Classic Cars	2006-09-01 00:00:00	Shipped	3	1	2006	Classic Cars
1000000045	46	46	1	2006	Classic Cars	2006-10-01 00:00:00	Shipped	3	3	2006	Classic Cars
1000000046	47	47	1	2006	Classic Cars	2006-11-01 00:00:00	Shipped	3	7	2006	Classic Cars
1000000047	48	48	1	2006	Classic Cars	2006-12-01 00:00:00	Shipped	3	11	2006	Classic Cars
1000000048	49	49	1	2007	Classic Cars	2007-01-01 00:00:00	Shipped	1	1	2007	Classic Cars
1000000049	50	50	1	2007	Classic Cars	2007-02-01 00:00:00	Shipped	1	3	2007	Classic Cars
1000000050	51	51	1	2007	Classic Cars	2007-03-01 00:00:00	Shipped	2	5	2007	Classic Cars
1000000051	52	52	1	2007	Classic Cars	2007-04-01 00:00:00	Shipped	3	7	2007	Classic Cars
1000000052	53	53	1	2007	Classic Cars	2007-05-01 00:00:00	Shipped	3	9	2007	Classic Cars
1000000053	54	54	1	2007	Classic Cars	2007-06-01 00:00:00	Shipped	3	11	2007	Classic Cars
1000000054	55	55	1	2007	Classic Cars	2007-07-01 00:00:00	Shipped	3	11	2007	Classic Cars
1000000055	56	56	1	2007	Classic Cars	2007-08-01 00:00:00	Shipped	3	12	2007	Classic Cars
1000000056	57	57	1	2007	Classic Cars	2007-09-01 00:00:00	Shipped	3	1	2007	Classic Cars
1000000057	58	58	1	2007	Classic Cars	2007-10-01 00:00:00	Shipped	3	3	2007	Classic Cars
1000000058	59	59	1	2007	Classic Cars	2007-11-01 00:00:00	Shipped	3	7	2007	Classic Cars
1000000059	60	60	1	2007	Classic Cars	2007-12-01 00:00:00	Shipped	3	11	2007	Classic Cars
1000000060	61	61	1	2008	Classic Cars	2008-01-01 00:00:00	Shipped	1	1	2008	Classic Cars
1000000061	62	62	1	2008	Classic Cars	2008-02-01 00:00:00	Shipped	1	3	2008	Classic Cars
1000000062	63	63	1	2008	Classic Cars	2008-03-01 00:00:00	Shipped	2	5	2008	Classic Cars
1000000063	64	64									

03

## PROJECT FOCUS



What impact sales?

# DATA CLEANING AND INITIAL ANALYSIS

## Dataset Observations

During the data cleaning and exploration phase, we identified several key points that impact the analysis and structure of the dataset:

- No Unique Key: The dataset lacks a primary key to uniquely identify each row.
- Order-to-Product Relationship: A single **ORDERNUMBER** can be associated with multiple **PRODUCTCODE** entries. This indicates that each row represents a line item within an order, not a complete order itself.
- Null Values: Several columns, such as **ADDRESSLINE2**, **STATE**, **POSTALCODE**, and **TERRITORY**, contain a significant number of null values.
- Sales Column Discrepancy: Initially, the meaning of the **SALES** column was not immediately clear, but it was determined to represent the multiplication of **QUANTITYORDERED \* PRICEEACH**. Upon validation, it was discovered that some rows contained an incorrect value of 100, which was an error in data entry.

04

# DATA CLEANING AND INITIAL ANALYSIS

## Creating a Subset Table

After performing this preliminary analysis, I decided to create a subset table to focus on the study of three primary factors: geographical, temporal, and product line.

The table was created by selecting the following columns from the original dataset:

OrderNumber	ProductCode	QuantityOrdered	Sales	City	Country	Month_ID	QTR_ID	ProductLine	
0	10107	S10_1678	30	2871.00	NYC	USA	2	1	Motorcycles
1	10121	S10_1678	34	2765.90	Reims	France	5	2	Motorcycles
2	10134	S10_1678	41	3884.34	Paris	France	7	3	Motorcycles
3	10145	S10_1678	45	3746.70	Pasadena	USA	8	3	Motorcycles
4	10159	S10_1678	49	5205.27	San Francisco	USA	10	4	Motorcycles
--	--	--	--	--	--	--	--	--	
2818	10350	S72_3212	20	2244.40	Madrid	Spain	12	4	Ships
2819	10373	S72_3212	29	3978.51	Oulu	Finland	1	1	Ships
2820	10386	S72_3212	43	5417.57	Madrid	Spain	3	1	Ships
2821	10397	S72_3212	34	2116.16	Toulouse	France	3	1	Ships
2822	10414	S72_3212	47	3079.44	Boston	USA	5	2	Ships

# EDA (EXPLORATORY DATA ANALYSIS)

## 🔍 Exploratory Data Analysis (EDA)

Following the data cleaning phase, I conducted an exploratory data analysis to better understand the key business characteristics and the relationships between variables. I based my analysis on the three factors I identified and followed these steps:

- **Variable Types Check:** I first examined the data types of each column to ensure they were appropriate for my analysis. In this case, no changes were necessary.
- **Univariate EDA:** I analyzed individual variables to understand their distributions, central tendencies, and unique values (e.g., looking at sales distribution or the count of each product line).
- **Bivariate EDA:** Finally, I explored the relationships between pairs of variables (e.g., how sales relate to the month or city).

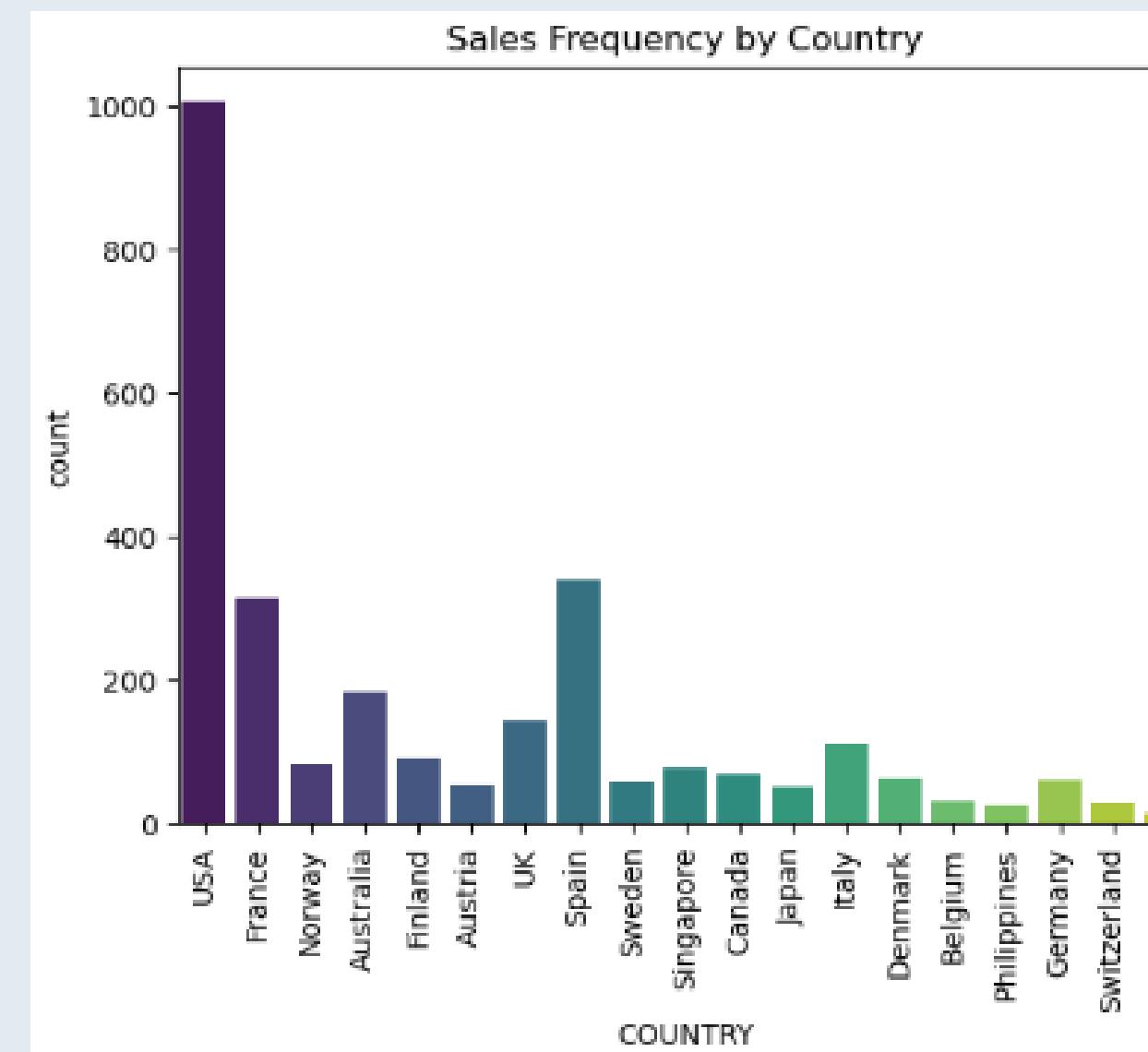
Here are some of the findings from my EDA:

# EDA (EXPLORATORY DATA ANALYSIS)



## Geographical Factors

- Sales by Country: I identified the United States (USA) and Spain (Spain) as the countries with the highest sales volume.

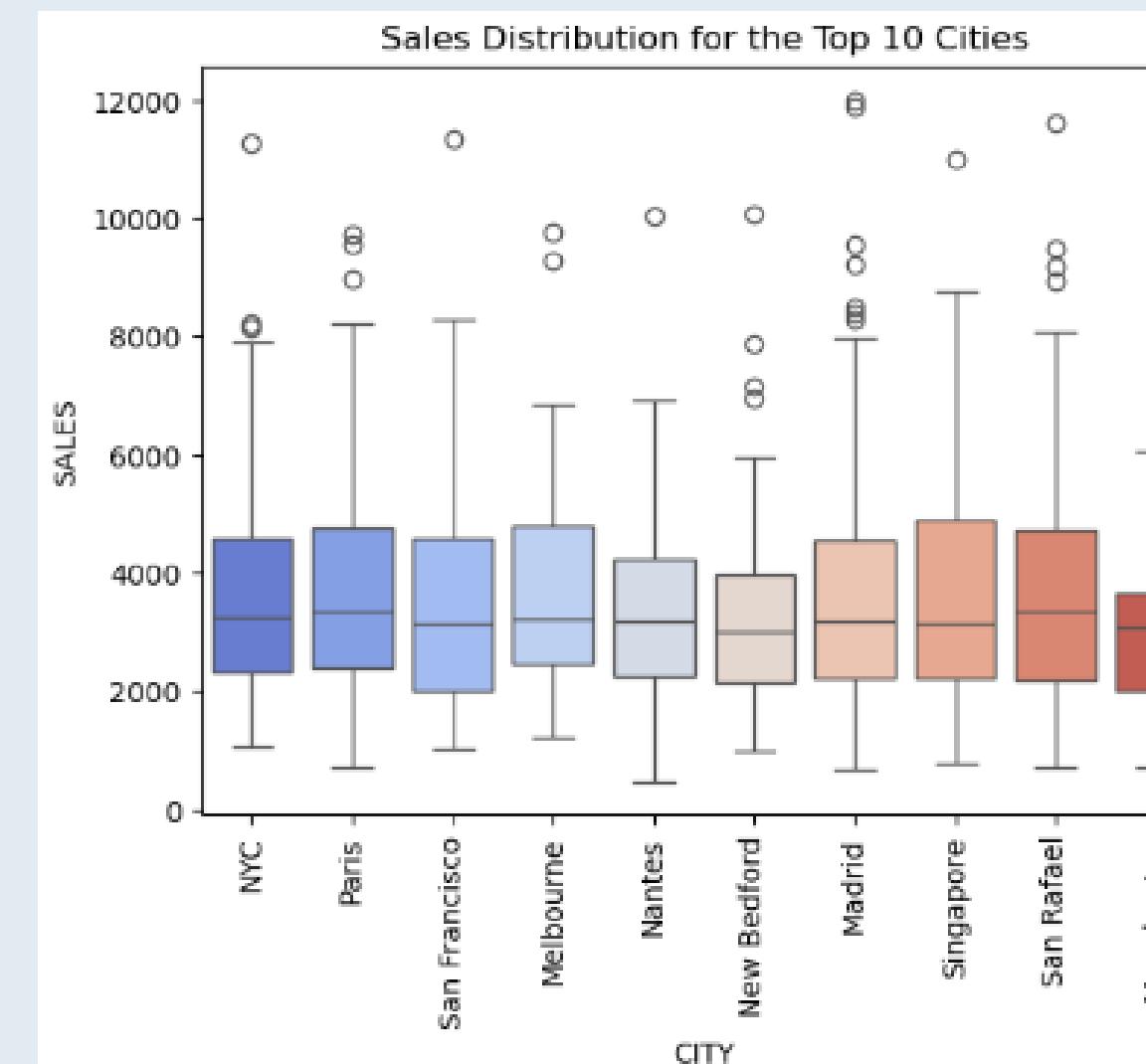


# EDA (EXPLORATORY DATA ANALYSIS)



## Geographical Factors

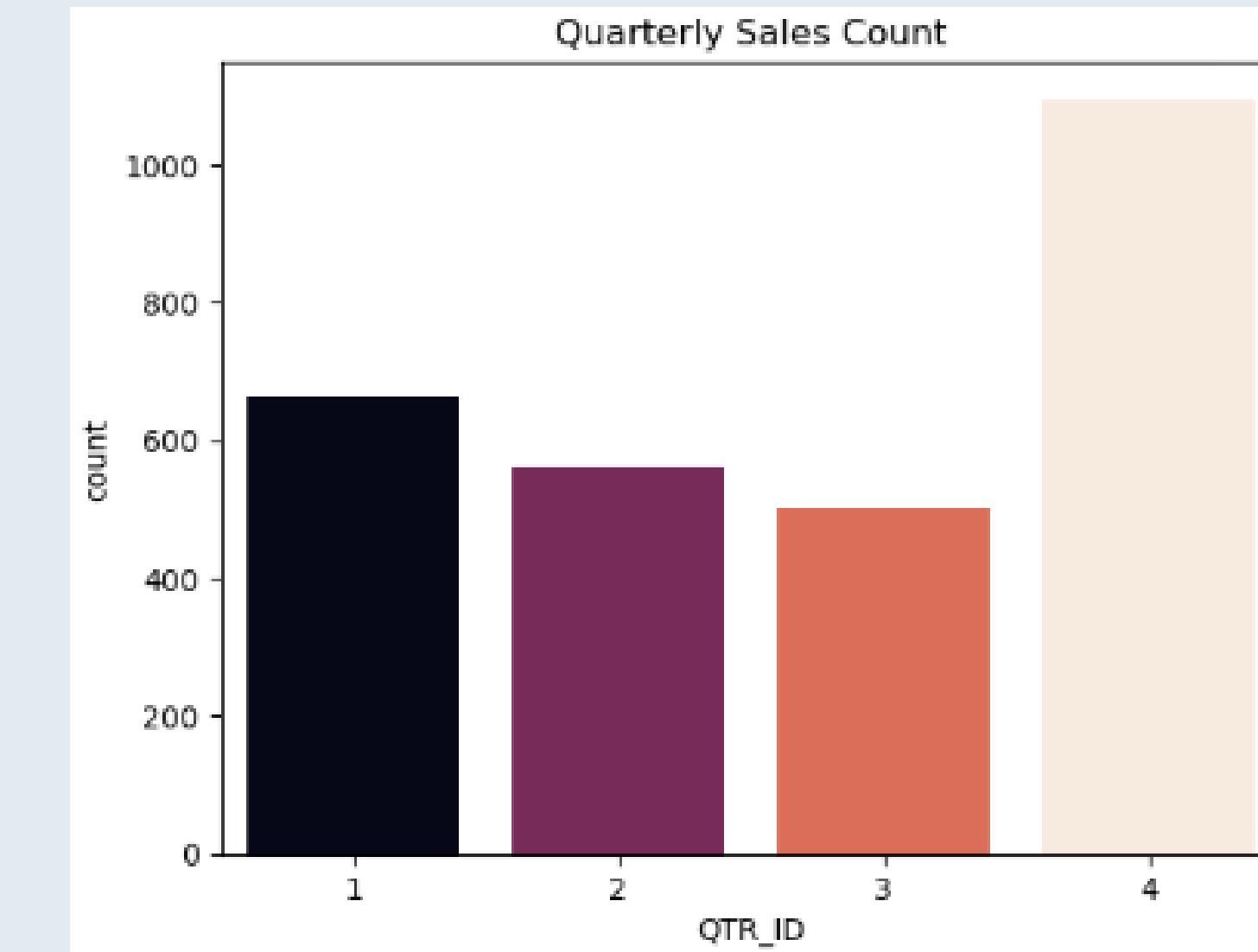
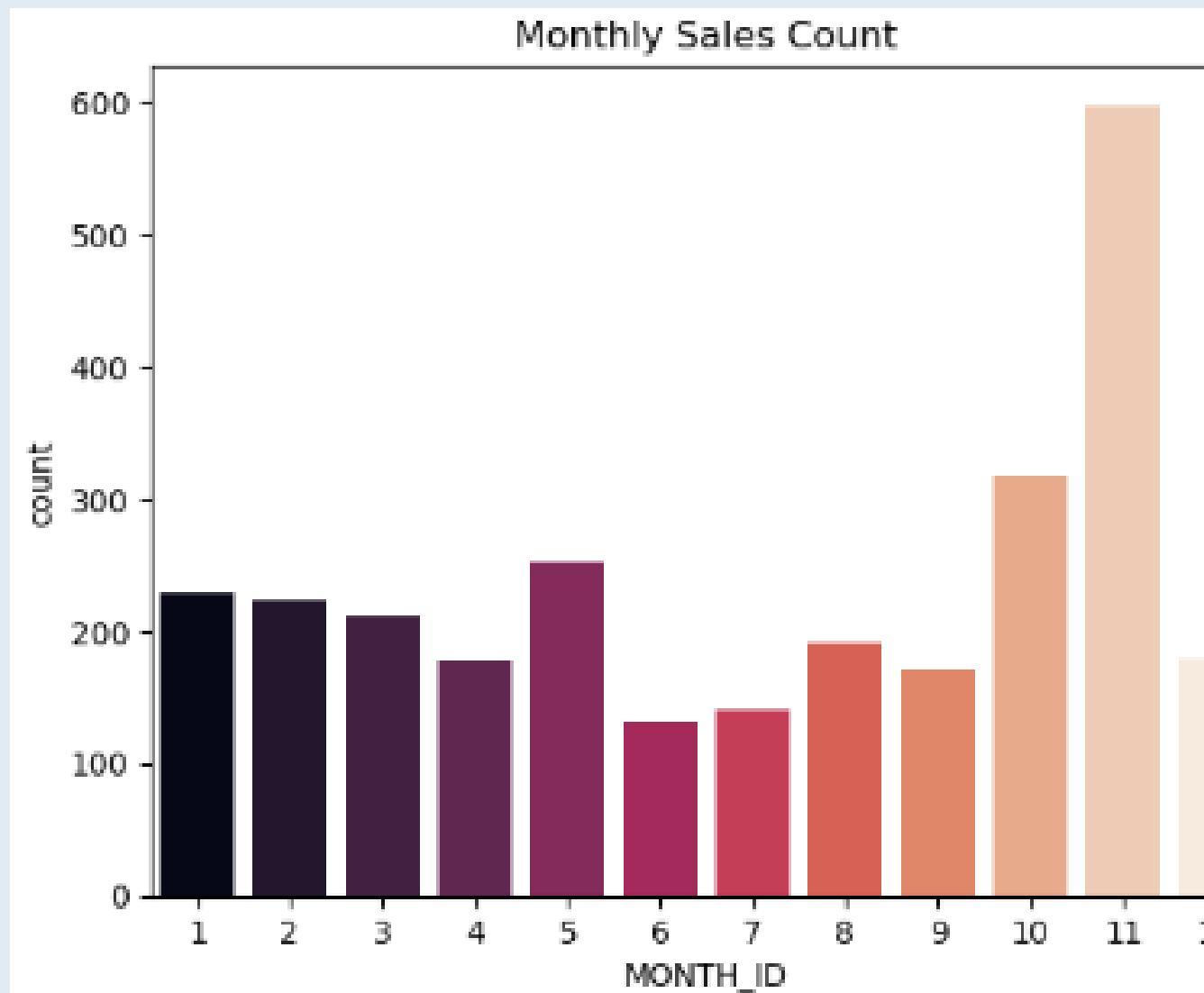
- Geographical Distribution: The boxplot shows that cities like NYC, Paris, and San Francisco generally have higher median sales and a wider distribution of high-value transactions. In contrast, cities such as San Rafael and Manchester exhibit a more concentrated distribution of lower sales values, with fewer high-value outliers.



# EDA (EXPLORATORY DATA ANALYSIS)

## (Temporal Factors)

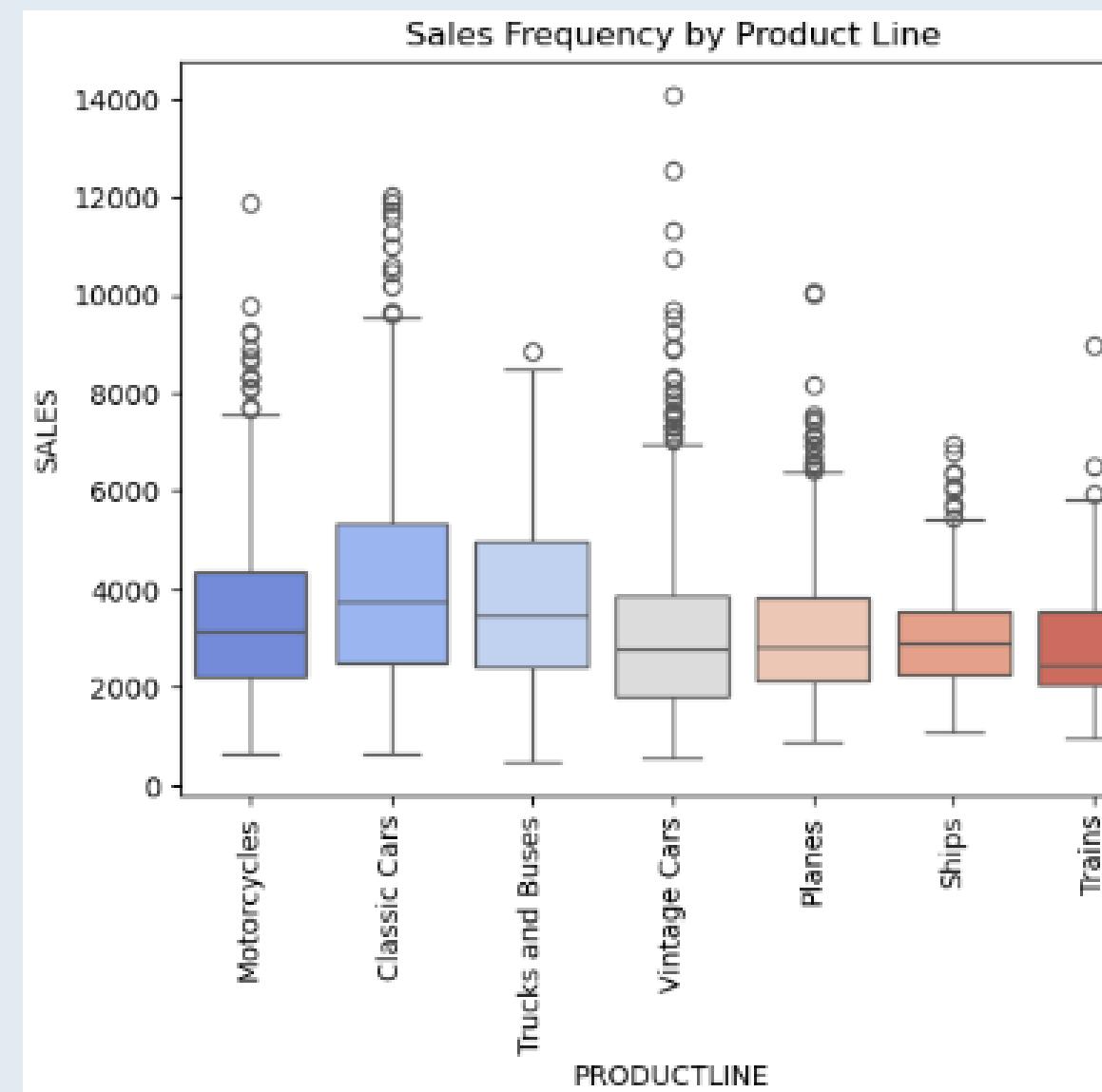
- Annual Cyclical Pattern: I observed a clear cyclical pattern in sales. The fourth quarter (QTR\_ID 4), particularly the month of November, shows a significant sales spike, indicating strong seasonal behavior.



# EDA (EXPLORATORY DATA ANALYSIS)

## 🛍️ Product Line Factors

- Top-Selling Products: My analysis of average sales by product line reveals a clear hierarchy. Classic cars have the highest average amount of sales, whereas trains have the least average amount of sales.



# CONCLUSIONS

06

This project gave me a comprehensive understanding of the sales dataset. The main challenge I faced was during the initial analysis, as the dataset lacked clear definitions and I had to determine the meaning of the SALES column.

This project gave me useful insights into sales behavior, where customers are located, and how products are performing. These results will help in creating future predictive models and making better strategic decisions.

For the complete code and a detailed breakdown of each step, please refer to my notebook.

# THANK YOU

