# E-Commerce Order Data Analysis with Missing Value Handling

**SUBMITTED BY:**

      **NAME : MANOHARI.M**

      **USN : 4GW23CI029**

      **EMAIL :manoharim987@gmail.com**

      **DATE : 02/09/2005**

# Problem Statement

- The dataset contains customer purchase records with missing values, duplicate entries, and inconsistent formats. The goal is to clean the dataset, handle missing values, and perform order-level analysis to extract business insights.

# TABLE OF CONTENTS

# Objectives

❑ **Clean and preprocess raw data**

– Remove errors, duplicates, and inconsistencies from raw e-commerce orders.

❑ **Handle missing values effectively**

– Apply imputation or removal techniques to ensure data reliability.

❑ **Remove duplicate & inconsistent records**

– Ensure each order is unique and data remains accurate.

❑ **Standardize formats for uniformity**

– Convert dates, numeric fields, and categories into a consistent format.

❑ **Perform customer & order-level analysis**

– Derive insights on customer behavior, order frequency, and spending.

❑ **Generate business insights through visualization**

– Use charts and graphs to identify sales trends, top customers, and product categories.

# Dataset Overview

## ❑ Dataset Columns:

- **OrderID** → Unique identifier for each customer order.
- **CustomerID** → Unique identifier for each customer; some entries missing.
- **Product** → Name/category of product purchased; inconsistent naming observed.
- **Quantity** → Number of items ordered per transaction.
- **Price** → Cost per unit of product; used to calculate total revenue.
- **OrderDate** → Date on which the order was placed; multiple formats present.

## ❑ Characteristics of the Dataset:

- Contains a **large volume of customer orders** collected over time.
- Covers **multiple product categories**, giving wide insights into sales.
- Data suffers from **quality issues**: missing values, duplicates, and inconsistent formatting.
- Rich enough for analysis **once properly cleaned and standardized**.

# Data cleaning steps

❑ **Duplicate Handling**
- •Checked for repeated OrderID values.
- •Removed duplicates to ensure each order is counted only once.

❑ **Missing Value Treatment**
- •Filled missing CustomerID using available patterns or frequent values.
- •Dropped records only when critical fields were unusable.

❑ **Standardization**
- •Converted all OrderDate entries into **YYYY-MM-DD** format.
- •Corrected invalid numeric values (e.g., negative or zero Quantity/Price).

❑ **Data Uniformity**
- •Standardized product names to avoid duplicates (e.g., "Laptop" vs. "laptop").
- •Ensured consistent naming across categories for reliable grouping and analysis.

# Handling Missing Values

❑ **Approach Applied**

•**Imputation:** Replaced minor missing values using mean/mode substitution.

•**Forward/Backward Fill:** Applied where sequential data (e.g., time-series orders) allowed logical filling.

•**Record Dropping:** Removed entries with missing critical fields (OrderID, Price) that could not be recovered.

❑ **Outcome**

•Achieved a **clean dataset with over 95% usable records**.

•Reduced noise from incomplete data.

•Improved **data reliability**, ensuring accurate customer and order-level analysis.

# Standardization and Formating

❑ **Standardization**
- **Date Format:** Converted all OrderDate entries to **YYYY-MM-DD** for consistency.
- **Numeric Fields:** Verified Quantity and Price → corrected invalid entries (e.g., negative or zero values).
- **Data Types:** Ensured proper types (integer for Quantity, float for Price).

❑ **Formatting for Uniformity**
- **Product Categories:** Standardized inconsistent product names (e.g., *"Laptop", "laptop", "LAPTOP" → "Laptop"*).
- **Consistent Units:** Ensured price and quantity follow uniform measurement standards.
- **Readable Structure:** Created a clean, structured dataset ready for analysis and visualization.

# Data Uniformity

•**Product Names:**
   •Standardized capitalization and spelling.
   •Merged similar entries (e.g., *"Mobile Phone"*, *"Mobiles"*, *"mobile phone"* → *"Mobile Phone"*).
•**Customer Records:**
   •Checked for duplicate CustomerID entries.
   •Consolidated information to avoid multiple profiles for the same customer.
•**Order Records:**
   •Verified each OrderID linked correctly to a unique customer and product.
   •Removed mismatched or incomplete references.
❑  **Benefits Achieved**
•Eliminated confusion caused by inconsistent data entry.
•Improved **grouping, filtering, and aggregation** for sales and customer analysis.
•Enhanced **data quality** → more reliable insights at customer, product, and order levels.

# Exploratory Data Analysis (EDA)

❑ **Key Analysis**

- **Customers:** Active vs. inactive users, repeat purchases, loyalty patterns.

- **Products:** Top-selling categories, low-performing items.

- **Orders:** Volume trends (daily, monthly, seasonal), revenue contribution.

- **Missing Data:** Checked concentration of nulls and business impact.

❑ **Insights**

- Clear view of **sales distribution** across customers & products.

- Identified **seasonal spikes** and sales patterns.

- Detected **customer behavior trends** for targeted marketing.

# Results & Insights

❑ **Key Findings**

• **Customer Behavior:**

   – A small group of loyal customers contributed to a large share of total revenue.

• **Product Performance:**

   – Electronics and fashion emerged as **top-selling categories**.

   – Certain low-demand products added little value and increased inventory costs.

• **Order Trends:**

   – Seasonal peaks observed (e.g., festival sales).

   – Revenue spikes linked to discount periods and promotions.

❑ **Business Impact**

• Enabled focus on **high-value customers** for retention strategies.

• Provided data-driven insights for **stock management** and **pricing decisions**.

• Improved clarity on **when to run promotions** for maximum impact.

# Conclusion and future scope

❑ **Conclusion**

• Successfully cleaned and preprocessed raw e-commerce dataset.

• Handled missing values, duplicates, and inconsistent formats to achieve a **95%+ usable dataset**.

• Performed detailed order-level and customer-level analysis.

• Extracted key insights on **customer behavior, product performance, and seasonal trends**.

• Improved dataset quality, enabling **reliable and data-driven decision making**.

❑ **Future Scope**

• Incorporate **predictive analytics** (e.g., forecasting demand, churn prediction).

• Expand analysis to include **customer demographics and regional trends**.

• Build a **dashboard/BI tool** for real-time monitoring of sales and customer activity.

# SCREENSHOTS

## Orders Dashboard

Summary Stats    Category Summary    Top Customers

### Summary Stats

| Metric | Value |
|--------|-------|
| TotalOrders | 51.000000 |
| TotalRevenue | 37749.692308 |
| AverageOrderValue | 754.993846 |
| UniqueCustomers | 47.000000 |

## Orders Dashboard

Summary Stats    Category Summary    Top Customers

### Category Summary

| ProductCategory | TotalRevenue | AverageOrderValue | OrderCount |
|-----------------|-------------|-------------------|------------|
| Clothing | 3807.000000 | 317.250000 | 12 |
| Electronics | 9172.692308 | 611.512821 | 15 |
| Furniture | 23450.000000 | 2605.555556 | 9 |
| Grocery | 1320.000000 | 94.285714 | 14 |

# Orders Dashboard

## Top 5 Customers

| CustomerID | TotalAmount |
| --- | --- |
| C114 | 4400.0 |
| C133 | 4200.0 |
| C123 | 2300.0 |
| C145 | 2250.0 |
| C139 | 2150.0 |

## Total Revenue by Product Category