

BMW Sales Analysis Data 2020 – 2025

Comprehensive Data Cleaning and Exploration Report

1. Data Import and Setup

- The analysis begins by importing essential Python libraries such as Pandas for data manipulation, NumPy for numerical operations, Matplotlib and Seaborn for visualization, and Random for generating synthetic values.
- The primary dataset, `BMW_Sales_2020_2025.csv`, was loaded into a Pandas DataFrame named `df` for further analysis.
- The data includes global BMW sales records covering multiple years and countries.
- The loading process ensures that all required dependencies are installed and verified before execution.
- Data integrity was checked at the time of import to confirm there were no issues with encoding or file corruption.
- Consistent naming conventions were followed to maintain readability and traceability of code. Proper file path verification was implemented to handle potential errors during loading.
- This stage sets the foundation for all subsequent steps by ensuring a clean and reproducible environment for analysis.

2. Initial Data Overview

- The dataset comprises 5,760 rows and 18 columns, capturing BMW's global sales transactions from 2020 to 2025.
- Key columns include Year, Month, Model, Type, Power, Region, Country, Units_Sold, Avg_Price_USD, Revenue_USD, Cost_USD, Profit_USD, Profit_Margin, Gross_Margin, Operating_Expenses_USD, Net_Profit_USD, ROI, and EBITDA_Margin. Each row represents a unique sales record with associated financial and operational metrics.
- This wide variety of fields provides a holistic view of BMW's business performance across markets and years.

- Basic exploration using `head()`, `tail()`, and `info()` functions confirmed the presence of meaningful data types, mostly
- numerical and categorical.
- The presence of both continuous financial metrics and categorical variables makes the dataset suitable for descriptive statistics and predictive modeling.
- Understanding the breadth of data is critical before moving to cleaning or transformation.
- This overview highlights the need for careful treatment of missing values and outliers to maintain analytical accuracy.

3. Descriptive Statistics

- Descriptive analysis revealed that the dataset spans the years 2020 through 2025, with Units Sold ranging from a minimum of 201 to a maximum of 2,500, and an average around 1,230.
- The average price of vehicles ranges between \$35,012 and \$119,995, with a mean near \$77,405.
- Revenue, cost, and profit variables follow expected distributions but show variability across regions and models.
- Profit margins fall between 15% and 34.99%, with a mean close to 24.79%, indicating overall healthy profitability.
- Return on Investment (ROI) spans 2.73 to 34.32, with a mean around 16.85, revealing significant differences in investment efficiency.
- The spread of values across key metrics suggests opportunities to identify high-performing product lines and markets.
- These statistics are vital to identify potential outliers and guide the cleaning process.
- They also provide initial signals of business trends before deeper exploratory analysis.

4. Data Quality Check

- A comprehensive quality check exposed several missing values across important columns.
- Critical fields such as `Avg_Price_USD` had 283 missing values, `Revenue_USD` had 326, `Cost_USD` had 283, `Profit_USD` had 289,

Operating_Expenses_USD had 326, Net_Profit_USD had 599, and ROI had 852.

- Despite these gaps, no duplicate rows were found, ensuring unique transactional records.
- This assessment emphasized the necessity of strategic imputation methods to prevent bias.
- Special attention was required for financial columns where inaccurate replacement could distort profitability metrics.
- Categorical fields such as Model, Type, and Region were verified to
- contain valid categories with no structural inconsistencies.
- Identifying these quality issues early allowed for targeted solutions in the cleaning phase.
- The absence of duplicates simplified the cleaning process but reinforced the need for accurate missing value handling.

5. Data Cleaning and Missing Value Treatment

- Cleaning involved systematic handling of missing data and preparation for analysis.
- Duplicate records were dropped even though none were detected during the initial check, ensuring consistency in methodology. Numerical columns were imputed using either median values to maintain central tendencies or random integers within observed ranges to simulate realistic variability.
- Categorical columns were filled with random selections from existing categories to preserve natural distributions.
- This dual approach ensured that imputed values neither skewed summary statistics nor introduced unrealistic patterns.
- Care was taken to replicate natural relationships between cost, revenue, and profit metrics during imputation.
- The cleaned dataset was saved as BMW_Sales_Filled_Replace.csv for future reference and reproducibility.
- All transformations were documented to maintain transparency and facilitate future auditing or updates.

6. Exploration and Group Insights

- Exploratory grouping of the cleaned dataset produced valuable business insights.
- Powertrain analysis revealed Electric vehicles as the leading category with 2,277 average units sold, followed by Petrol with 1,744, Diesel with 875, and Hybrid with 864.
- Regional comparisons highlighted Japan as a top-performing market with an average of 1,299 units sold, followed by the UK with 1,252, Spain with 1,251, and the USA with 1,236. Group-wise analysis by Model and Type identified high sales for Convertibles in South Africa and Sedans in Japan.
- Cost and profit groupings revealed substantial differences in profitability between models and regions.
- Insights from grouped statistics guide BMW in optimizing production and marketing strategies.
- The exploration also confirmed correlations between units sold and revenue, validating expected economic relationships.
- These group-level patterns provide a foundation for more advanced predictive analytics.
- They also suggest markets and product lines that may benefit from focused investment.

7. Key Observations

- Several key observations emerged from the combined cleaning and exploration process.
- Electric Vehicles dominate the dataset, reflecting BMW's strategic pivot toward sustainable mobility and electric technology.
- Japan consistently appears as a high-demand market across multiple vehicle categories, signaling strong brand penetration.
- Profit margins remain around 25% on average, indicating solid operational efficiency despite regional variability.

- ROI variability shows inconsistent cost structures, suggesting opportunities to improve capital allocation.
- The combination of high revenue but modest profit in some models indicates potential pricing or cost inefficiencies.
- Regional EBITDA margins display differences that may be attributed to taxation, logistics, and economic factors. The cleaned dataset is now suitable for advanced machine learning and visualization tasks.
- These insights collectively inform BMW's strategic decision-making for future growth and profitability.

8. Conclusion

- This comprehensive cleaning and exploratory analysis ensures a robust and reliable dataset for deeper statistical modeling and predictive analytics.
- The process demonstrated the importance of handling missing data carefully to preserve key business relationships.
- Exploratory grouping provided actionable insights into product performance, regional demand, and cost efficiency.
- BMW can leverage these findings to refine marketing strategies and production planning.
- The detailed imputation methodology supports reproducibility and future scalability of analysis.
- Strong patterns in electric vehicle sales reinforce the company's sustainability strategy.
- Regional profitability differences highlight markets requiring further operational optimization.
- Overall, this report provides a solid analytical foundation for BMW's future decision-making and advanced forecasting initiatives.