

# Amazon Sales Rapport

Tekup data science team

28-07-2024

## Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Project Objectives</b>	<b>3</b>
<b>3</b>	<b>Problem Statement</b>	<b>3</b>
<b>4</b>	<b>Existing Solutions</b>	<b>3</b>
<b>5</b>	<b>Our Proposed Solution</b>	<b>4</b>
<b>6</b>	<b>Justification for Our Solution</b>	<b>4</b>
<b>7</b>	<b>Technologies and Tools Used</b>	<b>5</b>
7.1	Technology Comparisons . . . . .	5
7.1.1	Database Management Systems . . . . .	5
7.1.2	Data Visualization Tools . . . . .	5
<b>8</b>	<b>Methodology - Scrum</b>	<b>6</b>
8.1	Scrum Roles . . . . .	6
8.2	Scrum Events . . . . .	6
8.3	Milestones . . . . .	6
<b>9</b>	<b>Results</b>	<b>7</b>
9.1	Data Collection and Cleaning . . . . .	7
9.2	ETL Process . . . . .	7
9.3	Data Visualization and Reporting . . . . .	8
9.4	Sales Forecasting . . . . .	9
9.5	GUI Development . . . . .	9
<b>10</b>	<b>Conclusion</b>	<b>10</b>

## 1 Introduction

The Amazon Sales Rapport Data Warehouse (DW) project was undertaken to establish a robust and comprehensive data infrastructure for analyzing and reporting on sales data from Amazon. This data warehouse aims to provide valuable insights into customer behavior, product performance, market trends, and future sales forecasts.

## 2 Project Objectives

The project objectives were defined as follows:

1. **Build a robust data warehouse architecture:** This involved designing and implementing a data warehouse system capable of storing and managing large volumes of sales data efficiently.
2. **Develop an ETL process:** An Extract, Transform, Load (ETL) process was designed to retrieve data from various sources, transform it for consistency and integrity, and load it into the data warehouse.
3. **Enable visualization and reporting:** Tools like Power BI were integrated to provide interactive dashboards and reports for analyzing sales data.
4. **Implement sales forecasting capabilities:** By leveraging time series analysis and machine learning algorithms in R, the project aimed to predict future sales trends.

## 3 Problem Statement

Amazon, with its vast and diverse sales data, lacked a centralized and structured system for effectively analyzing and leveraging this information. The existing data sources were fragmented, lacked consistency, and presented challenges for extracting meaningful insights. This project aimed to address this problem by creating a dedicated data warehouse for comprehensive sales data analysis.

## 4 Existing Solutions

Several existing solutions were considered, including:

1. **Traditional data warehousing solutions:** While these solutions offered mature infrastructure and tools, they were often complex and expensive to implement.
2. **Cloud-based data warehousing services:** These services provided scalability and flexibility but could be limited in terms of customization options.

3. **Open-source data warehouse solutions:** These solutions offered cost-effective alternatives but required more technical expertise to manage.

## 5 Our Proposed Solution

Our proposed solution encompassed a comprehensive approach that combined the advantages of traditional and modern data warehousing methodologies. We opted for a hybrid solution that leveraged the power of Python, MySQL, and R for data extraction, transformation, loading, visualization, and forecasting. The key aspects of our solution included:

1. **A robust data warehouse architecture:** We implemented a relational data warehouse architecture with optimized indexing and normalization techniques to ensure efficient data storage and retrieval.
2. **A comprehensive ETL process:** We designed an ETL pipeline using Python scripts to extract data from various sources, transform it using data cleaning and transformation libraries, and load it into the MySQL database.
3. **Visualization and reporting capabilities:** We integrated Power BI to provide interactive dashboards and reports for analyzing sales trends and customer insights.
4. **Advanced sales forecasting with R:** We leveraged time series analysis and machine learning algorithms in R to develop a sales forecasting model, providing predictions for future sales based on historical data.

## 6 Justification for Our Solution

Our chosen solution offered several advantages over existing solutions, including:

1. **Flexibility and scalability:** The use of open-source tools and technologies provided flexibility and scalability for handling evolving data requirements.
2. **Cost-effectiveness:** The use of open-source tools significantly reduced implementation costs compared to proprietary solutions.
3. **Customization and control:** We had greater control over the data warehouse architecture and ETL processes, allowing for tailored solutions to meet specific business needs.
4. **Integration with advanced analytics:** The integration of R provided a powerful platform for advanced time series analysis and machine learning, enabling robust sales forecasting capabilities.

## 7 Technologies and Tools Used

We employed a combination of technologies and tools to build our data warehouse solution. These included:

- (i) **Python:** The core language for developing ETL scripts, data cleaning and transformation processes, and building the user interface.
- (ii) **MySQL:** A relational database management system chosen for its scalability, reliability, and compatibility with various tools.
- (iii) **Data Transformation Libraries (Python):** We utilized popular Python libraries like Pandas, NumPy, and Scikit-learn for data manipulation, cleaning, and feature engineering.
- (iv) **R:** Employed for time series analysis, forecasting models, and data visualization.
- (v) **Power BI:** A powerful business intelligence tool integrated with the data warehouse for interactive dashboards and reporting.
- (vi) **GUI Interface (Python):** We developed a user-friendly GUI interface using Python libraries like Tkinter to facilitate data updates and create new data warehouse instances.

### 7.1 Technology Comparisons

#### 7.1.1 Database Management Systems

Table 1: Comparison of Database Management Systems

MySQL	Open-source, high scalability, cost-effective, strong community support
PostgreSQL	Open-source, high scalability, cost-effective, strong community support
Oracle	Proprietary, very high scalability, expensive, excellent support

#### 7.1.2 Data Visualization Tools

Table 2: Comparison of Data Visualization Tools

Power BI	Microsoft product, strong integration with Microsoft ecosystem, user-friendly interface, cost-effective
Tableau	Independent vendor, extensive data connectors, user-friendly interface, more expensive
Qlik Sense	Independent vendor, excellent data exploration features, advanced analytics features, more expensive

## 8 Methodology - Scrum

To ensure efficient project management and timely delivery, we adopted the Scrum agile methodology. This framework facilitated collaboration, iterative development, and continuous improvement.

### 8.1 Scrum Roles

1. **Product Owner:** Responsible for defining the project vision, prioritizing tasks, and ensuring alignment with business goals.
2. **Scrum Master:** Facilitates the Scrum process, removes impediments, and ensures team efficiency.
3. **Development Team:** A cross-functional team of developers responsible for implementing the data warehouse solution.

### 8.2 Scrum Events

1. **Sprint Planning:** The team plans and commits to a set of tasks for the upcoming sprint.
2. **Daily Scrum:** A short daily meeting to discuss progress, identify road-blocks, and coordinate efforts.
3. **Sprint Review:** The team demonstrates the completed work to the product owner and stakeholders.
4. **Sprint Retrospective:** The team reflects on the sprint, identifies areas for improvement, and plans for the next sprint.

### 8.3 Milestones

We defined key milestones throughout the project to track progress and ensure timely delivery:

1. **Data Collection and Cleaning:** The initial phase involved collecting data from various sources, cleaning it for consistency and accuracy, and preparing it for loading into the data warehouse.
2. **Data Warehousing and ETL Development:** The team designed and implemented the data warehouse architecture, developed the ETL process, and ensured data integrity.
3. **Visualization and Reporting:** The team integrated Power BI with the data warehouse, creating dashboards and reports for analyzing sales data.
4. **Sales Forecasting Model Development:** The team developed a sales forecasting model using R, applying time series analysis and machine learning techniques.

5. **GUI Development and Deployment:** The team built a user-friendly GUI for data updates and new instance creation, ensuring ease of use and scalability.

## 9 Results

### 9.1 Data Collection and Cleaning

1. **Source Identification:** We identified various data sources, including Amazon sales reports, customer feedback, and market trends.
2. **Data Extraction:** Python scripts were developed to extract data from identified sources, handling different formats and structures.
3. **Data Cleaning:** Data was cleaned to remove duplicates, handle missing values, and ensure consistency across different datasets.

### 9.2 ETL Process

1. **Extract:** Python scripts extracted data from various sources, including CSV files, APIs, and databases.
2. **Transform:** Data was transformed using Pandas and NumPy for cleaning, normalization, and feature engineering.
3. **Load:** The transformed data was loaded into the MySQL data warehouse, ensuring data integrity and consistency.

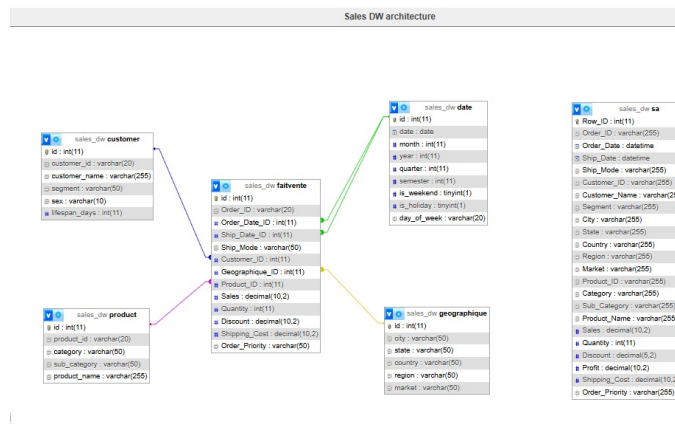


Figure 1: Data Warehouse Architecture Capture

### 9.3 Data Visualization and Reporting

Power BI was integrated with the data warehouse to provide interactive dashboards and reports. Key features included:

1. **Sales Trends:** Visualization of sales trends over time, highlighting peak sales periods and patterns.

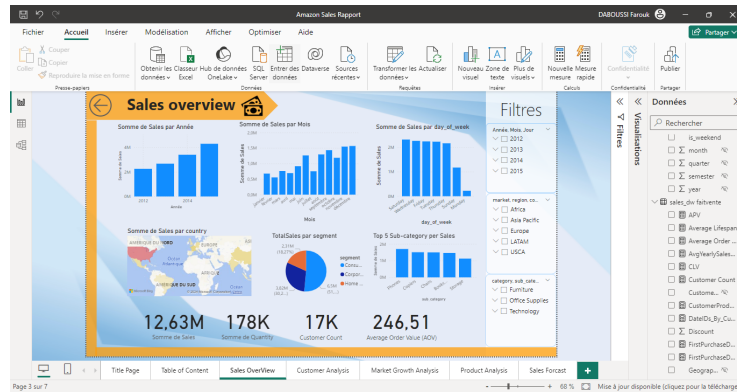


Figure 2: Power BI Report 1: Sales Trends Visualization

2. **Customer Insights:** Analysis of customer behavior, preferences, and demographics.

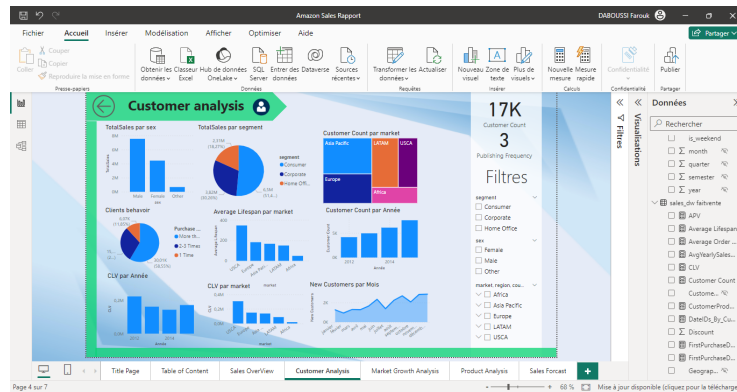


Figure 3: Power BI Report 2: Customer Insights Analysis

3. **Product Performance:** Evaluation of product performance, identifying top-selling products and underperformers.



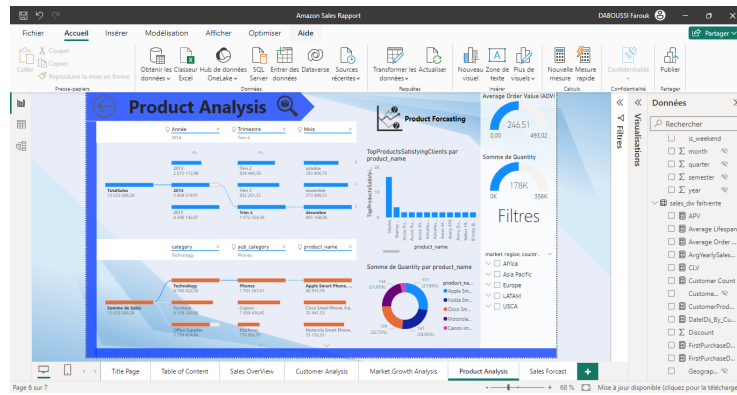


Figure 4: Power BI Report 3: Product Performance Evaluation

## 9.4 Sales Forecasting

We developed a sales forecasting model using R, applying time series analysis and machine learning techniques. The model provided predictions for future sales based on historical data, enabling informed decision-making.

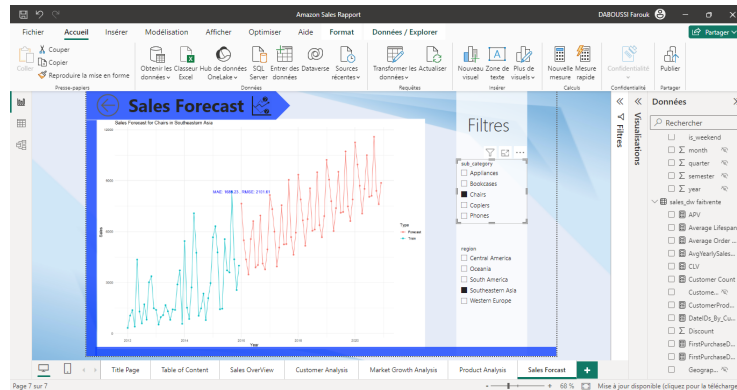


Figure 5: Power BI Report 4: Sales forecast

## 9.5 GUI Development

A user-friendly GUI interface was developed using Python libraries like Tkinter. Key features included:

1. **Data Updates:** Users could update the data warehouse with new data sources, ensuring the system remained up-to-date.
2. **Instance Creation:** Users could create new instances of the data warehouse, facilitating scalability and flexibility.

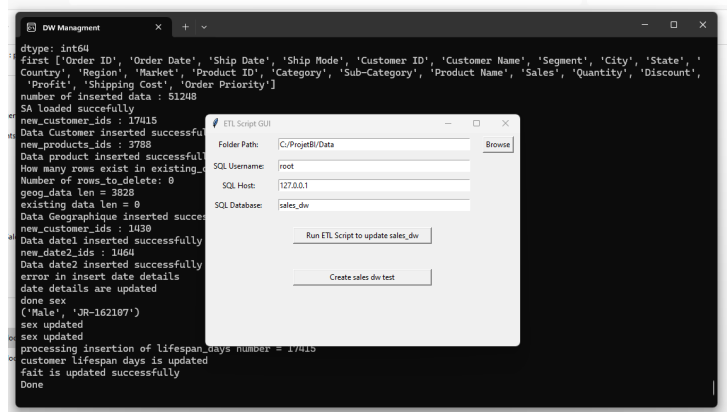


Figure 6: GUI Capture: User Interface for Data Updates and Instance Creation

## 10 Conclusion

The Amazon Sales Rapport Data Warehouse project successfully established a comprehensive data infrastructure for analyzing and reporting on sales data. The project achieved its objectives, providing valuable insights into customer behavior, product performance, market trends, and future sales forecasts. Our solution, leveraging Python, MySQL, R, and Power BI, offered flexibility, scalability, and advanced analytics capabilities, enabling Amazon to make informed business decisions.