**amazon**

**Business Intelligence & Database Management Project**

**Paper written by** *Mouna Saadaoui & Houssine Khlif*

**Major:** Marketing/BA - **Minor:** IT

## Table of contents

## 1. Introduction and requirement gathering

### 1.1 About the project

This project focuses on optimizing the performance of Amazon's operations by creating a data warehouse (DWH) and leveraging analytical tools to derive insights. The primary goal is to enhance decision-making processes, improve resource allocation, and support strategic planning

## 1.2 Requirements gathering

Key questions identified for this project:

1. Which product categories and subcategories perform best?
2. What regions exhibit the highest sales and customer engagement?
3. How can marketing campaigns be optimized for better ROI?
4. What are the seasonal trends and customer purchasing behaviors?
5. Which operational inefficiencies can be mitigated to improve performance?

---

## 2. Data and resources gathering
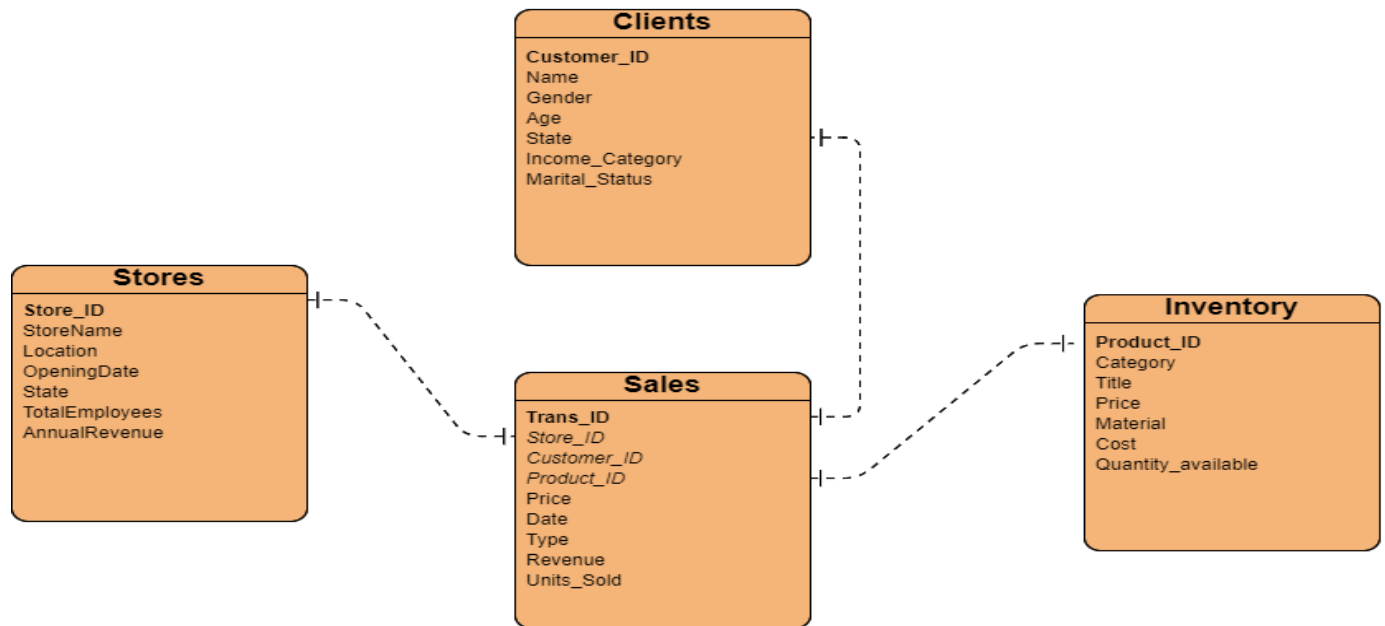
### 2.1 Data sources

- ***Internal database:***

The company's internal database is a well-organized relational system designed to store and manage critical business information. It serves as the backbone for key operations such as inventory management, sales monitoring, and client relationship management.

This database, named **Company_Performance_DWH**, consists of four interconnected tables: **SALES**, **INVENTORY**, **CLIENTS**, and **STORES**. These tables are linked through defined relationships, including Primary and Foreign Keys. Below is an overview of its

Relational Schema to demonstrate the database structure:

**Internal Database Relational Schema**



➢ *The data in this database originates from multiple sources. The **CLIENTS** and **SALES** datasets were sourced from the Kaggle website, containing information aligned with the attributes outlined in the schema. On the other hand, the **STORES** and **INVENTORY** datasets were manually generated using Microsoft Excel, with simulated and random data tailored to fit the structure of the company's database. Finally, these four datasets were combined to build the SQL database in SQL Workbench, where the files were imported, and their relationships were established.*

● ***External data sources:***

Amazon's core competency lies in efficient logistics, customer service, and offering a vast range of products, which includes managing its inventory budget strategically. To optimize operations, the company requires up-to-date information about consumer behavior and market trends, as well as historical data on regional and demographic insights. This information helps Amazon allocate its resources effectively across its supply chain and inventory management.

To support this, the following external data is utilized, though it has been simulated for analytical purposes:

1. **Geographical data**:
   ○ **File Type**: CSV
   ○ Data representing various regions, including population density, infrastructure, and market potential. This data is used to help optimize warehouse locations and delivery routes.
2. **Demographic data**:
   ○ **File Type**: CSV
   ○ Data provides insights into customer segments, such as age, income, and preferences. This information enables Amazon to tailor its marketing strategies and product offerings for different markets.

This external data, combined with internal databases, ensures that Amazon can analyze trends, strategize effectively, and maintain its competitive edge while managing resources efficiently.
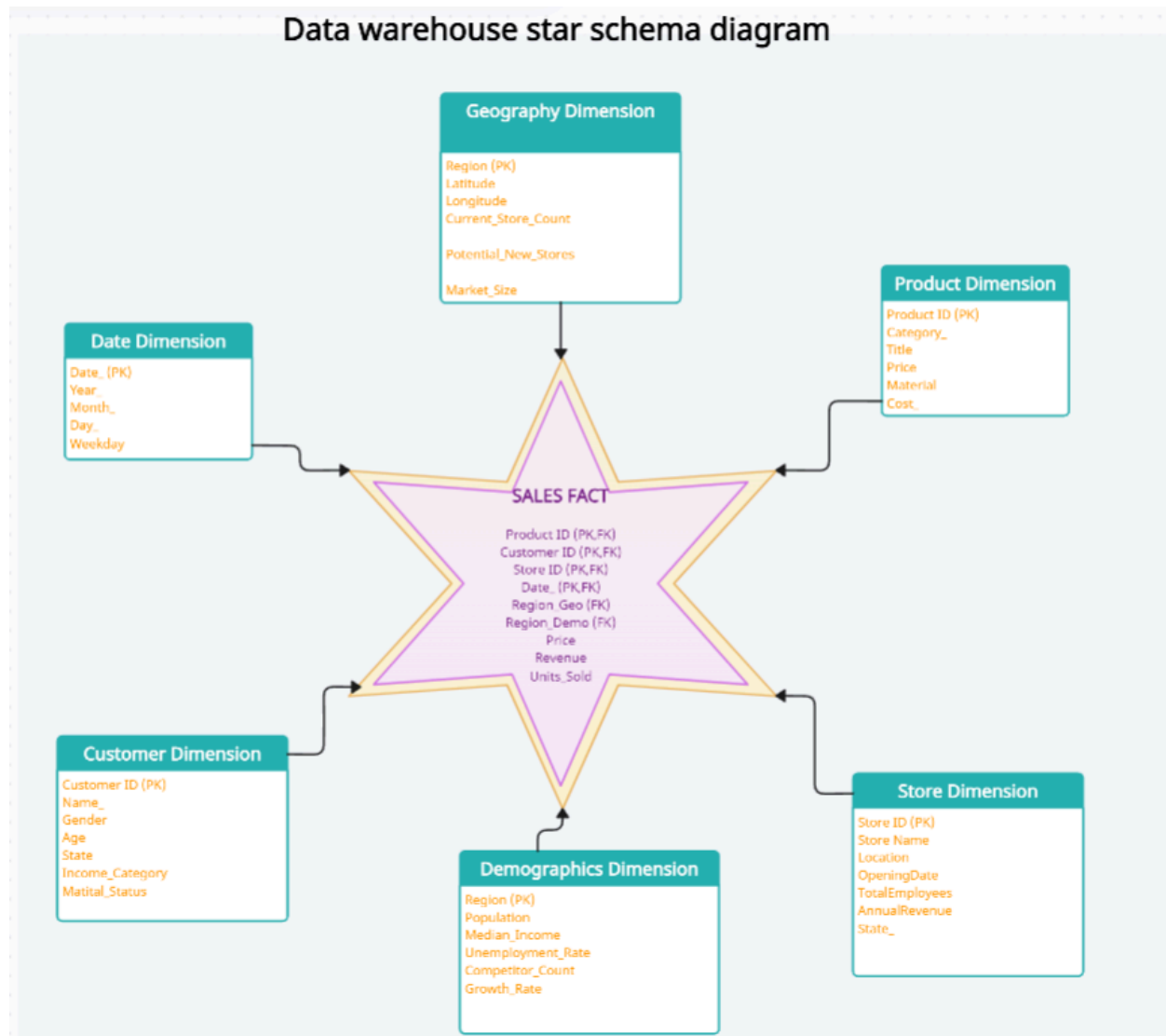
> ➢ **Throughout the project, we will be working with 2 different data sources types: SQL database and CSV file.**

---

## 3. Multidimensional modeling

As the Data Warehouse is yet to be established, it is more practical to design and develop its dimensional model before initiating the ETL process. The transformation step in ETL involves adapting the extracted data to align with the target system's structure, which cannot proceed without a clearly defined model.

A dimensional model serves as a blueprint for a particular business process. In line with the managers' requirements, this project will concentrate on analyzing and addressing

the **SALES BUSINESS PROCESS.**



## Data warehouse star schema diagram

**Geography Dimension**
Region (PK)
Latitude
Longitude
Current_Store_Count

Potential_New_Stores

Market_Size

**Product Dimension**
Product ID (PK)
Category_
Title
Price
Material
Cost_

**Date Dimension**
Date_ (PK)
Year_
Month_
Day_
Weekday

**SALES FACT**
Product ID (PK,FK)
Customer ID (PK,FK)
Store ID (PK,FK)
Date_ (PK,FK)
Region_Geo (FK)
Region_Demo (FK)
Price
Revenue
Units_Sold

**Customer Dimension**
Customer ID (PK)
Name_
Gender
Age
State
Income_Category
Matital_Status

**Demographics Dimension**
Region (PK)
Population
Median_Income
Unemployment_Rate
Competitor_Count
Growth_Rate

**Store Dimension**
Store ID (PK)
Store Name
Location
OpeningDate
TotalEmployees
AnnualRevenue
State_

## 03.01 FACT/DIMENSION IDENTIFICATION

This dimensional model, combined with these files, lays the groundwork for constructing the Data Warehouse and performing the ETL process effectively.

**Fact table: sales fact table**

- **Attributes:**
  - Product_ID
  - Customer_ID
  - StoreID

- Date
- Price
- Revenue
- Units_Sold
- Region-Geo
- Region_Demo

---

## Dimensions and their attributes

1. **Store dimension**
   - StoreID
   - StoreName
   - Location
   - OpeningDate
   - TotalEmployees
   - AnnualRevenue
   - State
2. **Product dimension**
   - Product_ID
   - Category
   - Title
   - Price
   - Material
   - Cost
3. **Geography dimension**
   - Region
   - Latitude
   - Longitude
   - Current_Store_Count
   - Potential_New_Stores
   - Market_Size
4. **Demographics dimension**
   - Region
   - Population
   - Median_Income
   - Unemployment_Rate
   - Competitor_Count
   - Growth_Rate
5. **Date dimension**
   - Date

- Year
      - Month
      - Day
      - Weekday
  6. **Customer dimension**
      - Customer_ID
      - Name
      - Gender
      - Age
      - State
      - Income_Category
      - Marital_Status

### 3.2 Star schema

Due to the straightforward structure of the dimensions and their direct connection to the fact table, we chose to use a star schema. Furthermore, a star schema is generally preferred over a snowflake schema for OLAP because it offers simpler, more intuitive query structures, leading to faster query performance and easier optimization. The denormalized design of the star schema minimizes the need for complex joins, making it more efficient for large-scale aggregations and analytical queries. Which is our main objective in this project.

## 4. ETL Process design and development

### 04.01 Data extraction

Before constructing the Data Warehouse, we first carried out the ETL process using **Python**. This phase involved extracting data from various available sources, both internal and external.

### External files extraction

Using Python, we extracted data from external files and ensured its quality by applying cleaning procedures. The extraction process included two main types:

  1. *Logical extraction:*
     A full extraction was performed, where all the data from source files was retrieved completely.

2. ***Physical extraction:***
    - **Offline extraction**: Data from external and internal sources was copied into CSV files for offline processing and fetched for further transformations.

The cleaning was performed using a custom Python function, which:

- Removed duplicates.
- Handled missing values.
- Reformatted column names (example: replacing spaces with underscores).
- Ensured numeric columns were properly typed.
- Added columns that were necessary for the functionality of the star schema.

## 04.02 Data transformation

The transformation phase was executed in Python to ensure the raw data met the quality, consistency, and compatibility standards required for the star schema. Each table underwent specific transformations, described below:

1. **Date dimension**:
    - Extracted unique dates from the Sales dataset.
    - Added columns for year, month, day, and weekday using Python's datetime module.
2. **Product dimension**:
    - Created by extracting product-related attributes (Product_ID, Category, Title, Price, etc.) from the Inventory dataset.
3. **Customer dimension**:
    - Derived from the Clients dataset by extracting customer-related fields (Customer_ID, Name, Gender, etc.).
    - Duplicates were removed based on Customer_ID.
4. **Store dimension:**
    - Created using store-related data (StoreID, StoreName, Location, etc.) from the Stores dataset.
5. **Geography dimension:**
    - Cleaned and processed data from the Geographical_Data.csv file, keeping relevant columns like Region, Latitude, and Market_Size.
6. **Demographics dimension:**
    - Cleaned and processed data from the Demographic_Data.csv file, focusing on fields such as Population and Median_Income.
7. **Sales fact table:**
    - Extracted from the Sales dataset.

○ Renamed and linked Region to both Geography and Demographics dimensions as Region-Geo and Region-Demo.

### *Data cleaning example*

- **State data transformation**: Underscores in state names (e.g., New_Jersey) were replaced with spaces (New Jersey) using a Python replace function.
- **Column consistency**: Column names were reformatted to avoid spaces and ensure consistency.

## 04.03 Schema creation

After transformations, the star schema was created with Python:

- Each dimension and the fact table were saved as separate CSV files.
- SQL files were generated for each table, including CREATE TABLE statements and INSERT INTO commands for loading data into the database.

**Python code**

```python
import pandas as pd

def clean_data(file_path):
    # Load the data
    data = pd.read_csv(file_path)

    # Drop duplicate rows
    data = data.drop_duplicates()

    # Handle missing values (e.g., fill with median for numeric columns)
    for column in data.columns:
        if data[column].dtype in ['int64', 'float64']:
            data[column] = data[column].fillna(data[column].median())
        else:
            data[column] = data[column].fillna('Unknown')

    # Strip whitespace from column names
    data.columns = data.columns.str.strip()

    # Rename columns for consistency (optional: customize as needed)
    data.rename(columns=lambda x: x.strip().replace(" ", "_"),
inplace=True)
```

```python
22.
23.    # Ensure numeric columns are correctly typed
24.    for column in data.select_dtypes(include=['object']).columns:
25.        try:
26.            data[column] = pd.to_numeric(data[column])
27.        except ValueError:
28.            pass
29.
30.    # Return the cleaned DataFrame
31.    return data
32.
33. # File paths
34. file_path_geographical =
    'C:/Users/21650/Downloads/New_script/Geographical_Data.csv'
35. file_path_demographic =
    'C:/Users/21650/Downloads/New_script/Demographic_Data.csv'
36.
37. # Clean the datasets
38. cleaned_geographical_data = clean_data(file_path_geographical)
39. cleaned_demographic_data = clean_data(file_path_demographic)
40.
41. # Save cleaned data to new files
42. cleaned_geographical_data.to_csv('C:/Users/21650/Downloads/New_script_out
    put/Cleaned_Geographical_Data.csv', index=False)
43. cleaned_demographic_data.to_csv('C:/Users/21650/Downloads/New_script_outp
    ut/Cleaned_Demographic_Data.csv', index=False)
44.
45. print("Datasets have been cleaned and saved.")
46.
47. # Star schema creation
48. # Load datasets
49. sales_data = pd.read_csv("C:/Users/21650/Downloads/New_script/Sales.csv")
50. inventory_data =
    pd.read_csv("C:/Users/21650/Downloads/New_script/Inventory.csv")
51. clients_data =
    pd.read_csv("C:/Users/21650/Downloads/New_script/Clients.csv")
52. stores_data =
    pd.read_csv("C:/Users/21650/Downloads/New_script/Stores.csv")
53. cleaned_geographical_data =
    pd.read_csv("C:/Users/21650/Downloads/New_script_output/Cleaned_Geographi
    cal_Data.csv")
54. cleaned_demographic_data =
    pd.read_csv("C:/Users/21650/Downloads/New_script_output/Cleaned_Demograph
    ic_Data.csv")
```

```python
55.
56. # Remove duplicate Customer_ID rows in clients_data
57. clients_data = clients_data.drop_duplicates(subset=['Customer_ID'])
58.
59. # Create Date Dimension
60. sales_data['Date'] = pd.to_datetime(sales_data['Date'])
61. date_dimension = sales_data[['Date']].drop_duplicates()
62. date_dimension['Year'] = date_dimension['Date'].dt.year
63. date_dimension['Month'] = date_dimension['Date'].dt.month
64. date_dimension['Day'] = date_dimension['Date'].dt.day
65. date_dimension['Weekday'] = date_dimension['Date'].dt.day_name()
66.
67. # Create Product Dimension
68. product_dimension = inventory_data[['Product_ID', 'Category', 'Title',
    'Price', 'Material', 'Cost']].drop_duplicates()
69.
70. # Create Customer Dimension
71. customer_dimension = clients_data[['Customer_ID', 'Name', 'Gender',
    'Age', 'State', 'Income_Category', 'Marital_Status']].drop_duplicates()
72.
73. # Create Store Dimension
74. store_dimension = stores_data[['StoreID', 'StoreName', 'Location',
    'OpeningDate', 'TotalEmployees', 'AnnualRevenue',
    'State']].drop_duplicates()
75.
76. # Create Geography Dimension
77. geography_dimension = cleaned_geographical_data[['Region', 'Latitude',
    'Longitude', 'Current_Store_Count', 'Potential_New_Stores',
    'Market_Size']].drop_duplicates()
78.
79. # Create Demographics Dimension
80. demographics_dimension = cleaned_demographic_data[['Region',
    'Population', 'Median_Income', 'Unemployment_Rate', 'Competitor_Count',
    'Growth_Rate']].drop_duplicates()
81.
82. # Create Sales Fact Table
83. sales_fact = sales_data[['Product_ID', 'Customer_ID', 'StoreID', 'Date',
    'Price', 'Revenue', 'Units_Sold', 'Region']].drop_duplicates()
84.
85. # Rename 'Region' column to 'Region-Geo' and add 'Region_Demo'
86. sales_fact.rename(columns={'Region': 'Region-Geo'}, inplace=True)
87. sales_fact['Region_Demo'] = sales_fact['Region-Geo']
88.
89. # Save all dimensions and the fact table as CSV files
```

```python
90.  output_paths = {
91.      "Date Dimension":
     "C:/Users/21650/Downloads/New_script_output/Date_Dimension.csv",
92.      "Product Dimension":
     "C:/Users/21650/Downloads/New_script_output/Product_Dimension.csv",
93.      "Customer Dimension":
     "C:/Users/21650/Downloads/New_script_output/Customer_Dimension.csv",
94.      "Store Dimension":
     "C:/Users/21650/Downloads/New_script_output/Store_Dimension.csv",
95.      "Geography Dimension":
     "C:/Users/21650/Downloads/New_script_output/Geography_Dimension.csv",
96.      "Demographics Dimension":
     "C:/Users/21650/Downloads/New_script_output/Demographics_Dimension.csv",
97.      "Sales Fact Table":
     "C:/Users/21650/Downloads/New_script_output/Sales_Fact_Table.csv"
98.  }
99.
100.     date_dimension.to_csv(output_paths["Date Dimension"], index=False)
101.     product_dimension.to_csv(output_paths["Product Dimension"],
     index=False)
102.     customer_dimension.to_csv(output_paths["Customer Dimension"],
     index=False)
103.     store_dimension.to_csv(output_paths["Store Dimension"], index=False)
104.     geography_dimension.to_csv(output_paths["Geography Dimension"],
     index=False)
105.     demographics_dimension.to_csv(output_paths["Demographics Dimension"],
     index=False)
106.     sales_fact.to_csv(output_paths["Sales Fact Table"], index=False)
107.
108.     print("Star schema created and saved:")
109.     for name, path in output_paths.items():
110.         print(f"- {name}: {path}")
111.
112.     # Save each table as an SQL file
113.     output_paths_sql = {
114.         "Date_Dimension.sql": date_dimension,
115.         "Product_Dimension.sql": product_dimension,
116.         "Customer_Dimension.sql": customer_dimension,
117.         "Store_Dimension.sql": store_dimension,
118.         "Geography_Dimension.sql": geography_dimension,
119.         "Demographics_Dimension.sql": demographics_dimension,
120.         "Sales_Fact_Table.sql": sales_fact
121.     }
122.
```

```python
123.    for filename, dataframe in output_paths_sql.items():
124.        table_name = filename.split('.')[0]
125.        sql_file_path = f"C:/Users/21650/Downloads/New_script/{filename}"
126.
127.        with open(sql_file_path, 'w', encoding='utf-8') as sql_file:
128.            # Write the CREATE TABLE statement
129.            columns = ', '.join([f'[{col}] TEXT' for col in
    dataframe.columns])  # Simplified as TEXT
130.            sql_file.write(f"CREATE TABLE {table_name} ({columns});\n")
131.
132.            # Write the INSERT INTO statements
133.            for _, row in dataframe.iterrows():
134.                values = []
135.                for value in row.values:
136.                    if isinstance(value, str):
137.                        value = value.replace("'", "''")  # Escape single
    quotes for SQL
138.                        values.append(f"'{value}'")  # Wrap strings in
    single quotes
139.                    elif pd.isna(value):
140.                        values.append("NULL")  # Handle missing values
141.                    else:
142.                        values.append(str(value))  # Convert other values
    to string
143.
144.                values_str = ', '.join(values)
145.                sql_file.write(f"INSERT INTO {table_name} VALUES
    ({values_str});\n")
```

### 4.3 Loading

After transforming internal and external raw data into a structured and stable format, we load it into the RDBMS in CSV format:

- Loading type : Full Loading

---

## 5. Data warehouse creation and data storage

After completing the ETL process, we proceeded to create and populate the data warehouse, a critical step in transforming raw data into a format that supports analytical queries and decision-making. For this purpose, we utilized SQL Developer as our RDBMS, leveraging its robust set of features to streamline the process and optimize performance.
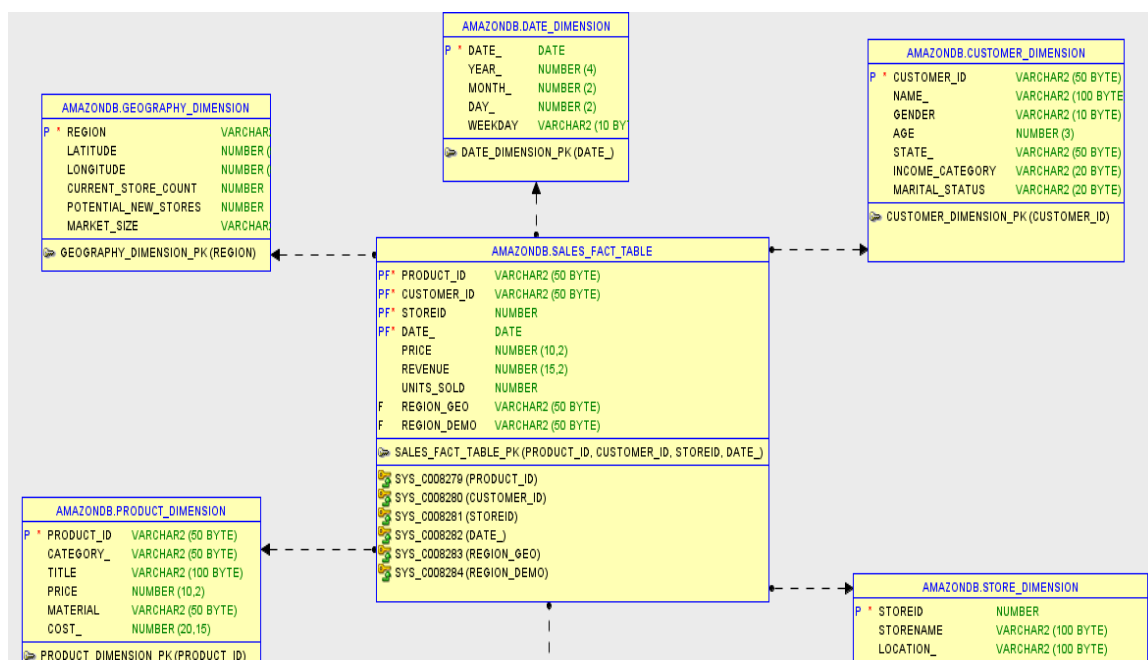
The data warehouse design was based on the **star schema**, a widely recognized and effective model for organizing data. This design consisted of dimension tables and a central fact table, all of which we loaded into the database using CSV files. The use of the star schema ensured a simple yet powerful structure, enabling us to efficiently store data while facilitating intuitive querying and reporting.

We chose SQL Developer for several key reasons. First, it supports **multidimensional modeling**, acting as a ROLAP (Relational Online Analytical Processing) server. This allowed us to build a relational foundation for analyzing large datasets without compromising flexibility. One of the standout features we utilized was its support for **materialized views**, which played a crucial role in enhancing query performance. By precomputing and storing complex aggregations, materialized views significantly reduced the time required to run analytical queries, enabling faster insights.

Throughout this process, we saved and documented our progress to maintain transparency and traceability. The ERD (Entity-Relationship Diagram) of the data model served as a visual representation of the relationships between our tables, providing a clear blueprint for the underlying structure of the data warehouse. Having this documentation not only supported our implementation but also provided a useful reference for future development or troubleshooting.

What made this stage particularly interesting was the tangible shift from handling raw, unorganized data to seeing it transformed into a well-structured system capable of delivering insights. The structured, stable nature of the star schema, combined with the power of SQL Developer, brought our data warehouse to life.

In addition to the technical benefits, working through this process also gave us a deeper appreciation of the role a well-designed data warehouse plays in the overall data pipeline.
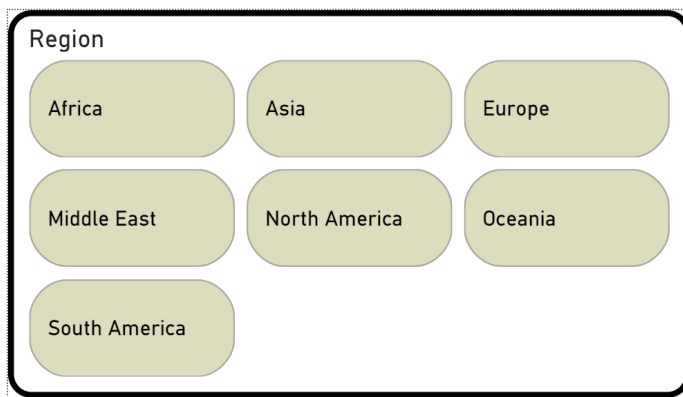
# 6. Data visualization and analysis

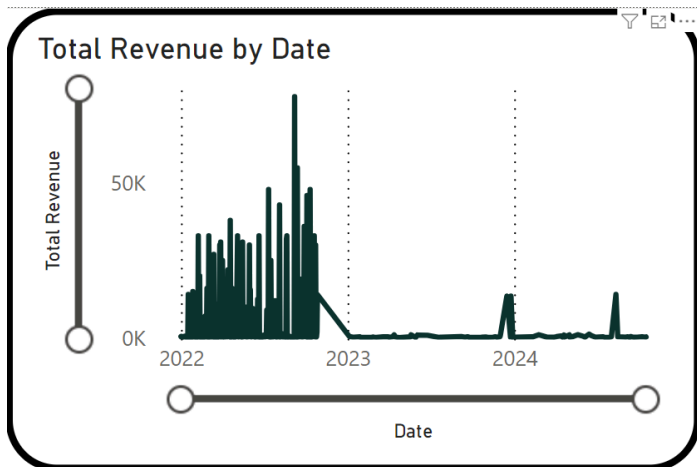Power BI was employed for creating visualizations and dashboards.

## 1. Regions overview :

- The dashboard has clickable regions (Africa, Asia, Europe, Middle East, North America, Oceania, and South America).
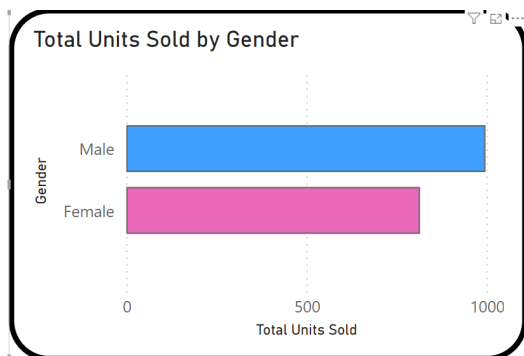- This likely serves as a filter to view data specific to each region.



## 2. Total revenue by date :

- A time-series graph shows revenue trends over time (2022 to 2024).
- There is a notable spike in revenue during certain periods, indicating high-performing months or campaigns.
- There are also dips in revenue, possibly related to seasonality or operational issues.

Total Revenue by Date

## 3. Total units sold by gender :

- Males appear to have purchased more units compared to females, indicating higher participation or demand from male customers.
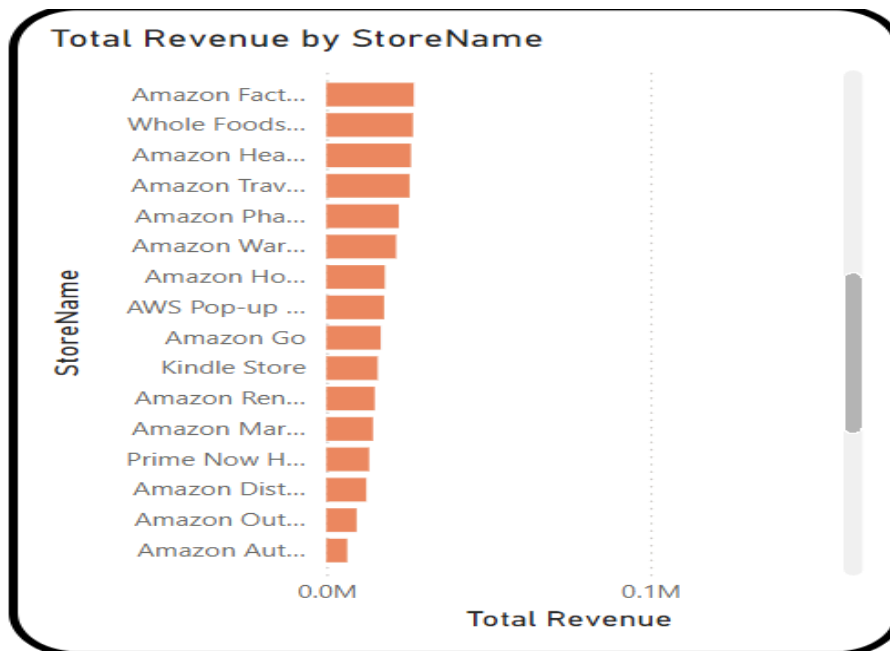


Total Units Sold by Gender

## 4. Total units sold and revenue by category of products:

- USB cables dominate the sales and revenue share with 86.12% of units sold.
- Other categories contribute minimally (e.g., HDMI cables, remote controls, etc.), suggesting product portfolio imbalance.

**Total Units Sold and Sum of Revenue by Category**

**Category**
- USBCables
- SmartTelevisions
- RemoteControls
- HDMICables
- WirelessUSBAda...
- StandardTelevisi...
- TVWall&Ceiling...
- Projectors
- RCACables
- OpticalCables
- AVReceivers&A...
- DVICables
- Adapters
- SpeakerCables

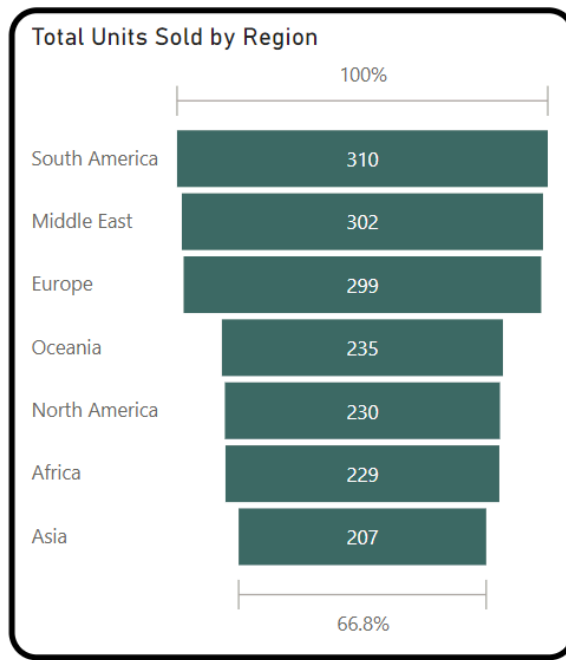0M (0.01%)  0M (0%)  0M (0.07%)

1.33M (86.12%)

## 5. Total revenue by store name :

- Amazon Pet Store generates the highest revenue, followed by Amazon Logistics and Amazon Publications.
- Revenue seems diverse across stores, though there is a clear top performer.

**Total Revenue by StoreName**

| StoreName |
| --- |
| Amazon Fact... |
| Whole Foods... |
| Amazon Hea... |
| Amazon Trav... |
| Amazon Pha... |
| Amazon War... |
| Amazon Ho... |
| AWS Pop-up ... |
| Amazon Go |
| Kindle Store |
| Amazon Ren... |
| Amazon Mar... |
| Prime Now H... |
| Amazon Dist... |
| Amazon Out... |
| Amazon Aut... |

0.0M          0.1M
**Total Revenue**

## 6. Total units sold by region :

- South America has the highest units sold (310), followed by the Middle East (302) and Europe (299).
- Asia has the lowest units sold (207), possibly highlighting a lack of market penetration or lower demand.

**Total Units Sold by Region**

| Region | 100% |
|---|---|
| South America | 310 |
| Middle East | 302 |
| Europe | 299 |
| Oceania | 235 |
| North America | 230 |
| Africa | 229 |
| Asia | 207 |

66.8%

## ➢ Key insights:

1. **Top product**: USB cables are a strong driver of both sales and revenue, suggesting their critical importance in the business.
2. **High-performing region**: South America leads in units sold, which could indicate a strong market presence or demand there.
3. **Demographic insights**: Male customers purchase more, indicating potential for targeting female customers through campaigns.
4. **Revenue concentration**: A small number of store names contribute significantly to total revenue; diversification or optimization in lower-performing stores could help.
5. **Seasonality or event impact**: The time-series chart highlights spikes in revenue, which could align with events, promotions, or holidays.

## References

1. GitHub repository: https://github.com/houssinekhlif/Amazon_DWH_Project
2. External sources:

https://www.kaggle.com/code/mehakiftikhar/amazon-sales-dataset-eda
https://www.youtube.com/watch?v=0BKlUySopU4&list=PLwIcJx1aSL1SeTJgPbFgf1V-5CfsV4l1l
https://youtu.be/I7DZP4rVQOU?list=PLTsu3dft3CWhOUPyXdLw8DGy_1l2oK1yy

---

**Emails:**

**mouna.saadaoui202@gmail.com**

**houssine.khlif.7@gmail.com**

***Thank you for reading***