



**2023/2024**

# **DEPARTMENT STORE ANALYSIS REPORT**

**BUSINESS INTELLIGENCE AND  
DATABASE MANAGEMENT SYSTEMS**

**Work done by :**

Sarah Ibn El hadj - Wadie Tliche  
Becher Zribi - Nour Attia

**Professors:**

Prof.Manel Abdelkader  
Prof.Ameni Azzouz

# TABLE OF CONTENTS

INTRODUCTION	1
--------------	---

IMPLEMENTATION	2
----------------	---

DATA GATHERING	2.1
----------------	-----

DATA PREPARATION	2.2
------------------	-----

DATA STORAGE	2.3
--------------	-----

STORAGE	2.3.1
---------	-------

FACT	2.3.2
------	-------

DIMENSIONS	2.3.3
------------	-------

DATA VISUALIZATION	2.4
--------------------	-----

CONCLUSION	3
------------	---

# INTRODUCTION

In the dynamic retail landscape, our research focuses on understanding **shopping trends in department stores in the United States**, providing insights for strategic planning. Utilizing a **comprehensive dataset covering customer demographics, purchasing patterns, product preferences, and inventory levels**, we aim to enhance decision-making for businesses.

By Examining various information from the gathered datasets, our primary objective is to distil key insights to offer a nuanced understanding of consumer preferences and industry trends, aiding businesses in navigating the complexities of the local United States retail market.

# INTRODUCTION

Concretely, we focused on 4 important key performance indicators(KPI):

1. **Demographic Analysis:** By understanding the relationships between marital status, income and location distribution of customers to identify our potential target markets and tailor customized marketing strategies to increase the demand.
2. **Product Preferences:** By examining the types items customers purchase, their preferred category (furniture, office suppliers, technology), their preferred sub-category, seasonal inclinations, we can optimize our inventory management and product assortment.
3. **Purchase Behaviour:** By studying the frequency of purchases and shipping mode to optimize the overall shopping experience.
4. **Inventory Optimization:** By examining the quantity on hand, the total costs and the inventory turnovers to update our management strategies and reduce expenses.

By achieving these goals, businesses can attain a more profound insight into their customer base. This, in turn, empowers them to **fine-tune marketing strategies, optimize inventory management, and elevate overall customer satisfaction**. Ultimately, these efforts contribute to sustained business growth in the competitive retail landscape.

## 2.IMPLEMENTATION

### 2.1.DATA GATHERING:

We obtained the dataset on Shopping Trends ratings in the USA from Kaggle.

This is a link to the dataset.

### 2.2.DATA PREPARATION:

Preparing the data , we utilized Python to streamline the manipulation and configuration of our dataset for integration into the data warehouse. The dataset existed in two formats: JSON and CSV.

Upon careful examination, a critical observation revealed the existence of duplicated data and Non-Available data. To tackle this issue, we utilized Python to identify and rectify duplicate rows within the dataset. We also converted data into the right data types.

Untitled-1

```
1 # %%
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import json
5 import numpy as np
6 from datetime import datetime
7
8 # %% [markdown]
9 # Exploring the first Dataset
10
11 # %%
12 df= pd.read_csv('sales_data.csv')
13
14 # %%
15 df.info()
16
17 # %%
18 df.describe()
19
20 # %%
21 #checking for null values
22 df.isnull().sum()
23
24 # %% [markdown]
25 # Treating the first dataset
26 #
27
28 # %%
29 #tracking the null values
30 null_rows = df.loc[df['postal_code'].isnull(), ['postal_code', 'region','city','state']]
31 print(null_rows)
32
33 # %%
34 print(df.loc[df['city']=='Burlington',['postal_code','state','city']])
35
36 # %%
37 #correcting the null values
38 df['postal_code'] = df['postal_code'].fillna('05402')
```

```
39 df.isnull().sum()
40
41 # %%
42 #converting dates to the date datatype
43 df['order_date']= pd.to_datetime(df['order_date'], format='%d/%m/%Y')
44 df['ship_date']= pd.to_datetime(df['ship_date'], format='%d/%m/%Y')
45 df.info()
46
47 # %% [markdown]
48 # Treating the second dataset
49
50 # %%
51 with open('customer_ds.json', 'r') as json_file:
52     df2 = json.load(json_file)
53
54 # %%
55 #exploring the structure
56 # Check if data is a list
57 if isinstance(df2, list):
58     # Access the first element in the list
59     first_element = df2[0]
60
61     # Display keys of the first dictionary in the list
62     print("Keys of the first dictionary in the list:", first_element.keys())
63 else:
64     print("The data is not a list of dictionaries.")
65
66 # %%
67 print(pd.json_normalize(df2).describe())
68
69 # %%
70 print(pd.json_normalize(df2).isnull().sum())
71
72 # %%
73 # Extract non-null incomes
74 incomes = [item.get('Income', 0) for item in df2]
75
76 # Calculate the mean of non-null incomes
77 mean_income = np.mean([income for income in incomes if income is not None])
78
79 # Replace null incomes with the mean
```

```
80 for item in df2:
81     if item.get('Income') is None:
82         item['Income'] = mean_income
83
84 # %%
85 print(pd.json_normalize(df2).isnull().sum())
86
87 # %%
88 #convertin to the right data types
89 for item in df2:
90     if 'Year_Birth' in item:
91         try:
92             year_as_int = int(item['Year_Birth'])
93             item['Year_Birth'] = datetime(year_as_int, 1, 1) # Assuming January 1st of the given year
94         except ValueError:
95             print("Invalid year format for item: {item}")
96
97 # %% [markdown]
98 # Building Dimensions
99 #
100
101 # %%
102 #product dimension
103
104 # %%
105 pdm = df.loc[:, ['product_id','product_name','sub_category','category']].copy()
106 pdm.head()
107
108 # %%
109 #eliminating the duplicates and adding the quantities :
110
111 # Group by 'id' and count occurrences
112 pdm['quantity'] = pdm.groupby('product_id')['product_id'].transform('count')
113
114 # Drop duplicates from the 'id' column
115 pdm.drop_duplicates(subset='product_id', inplace=True)
116
117 # Reset index after dropping duplicates
118 pdm.reset_index(drop=True, inplace=True)
119
120
```

# 2.IMPLEMENTATION

## 2.3 DATA STORAGE :

### 2.3.1.STORAGE:

After obtaining the refined version of the data we chose csv files to serve as a staging area.

```
1 | # %%  
2 | pd.json_normalize(df2).to_csv('cleaned_cus')  
3 |  
4 | # %%  
5 | dp.to_csv('cleaned_data.csv')  
6 |
```

We used python to create the Desired Fact and dimensions , and using the same tool we exported the output to our data warehouse which is MySQL.

## ESTABLISHING CONNECTION

```
255 # %%
256 import mysql.connector
257
258 # %%
259 conn = mysql.connector.connect(
260     host='localhost',
261     user='root',
262     password='Wdie.66.99',
263     database='warehouse2'
264 )
265
266 # %%
267 cursor = conn.cursor()
268
```

## DEFINED FUNCTION TO CREATE TABLES

```
270 #creating a function to insert the data
271 def create_table_from_excel(cursor, excel_file_path, table_name):
272     # Load Excel data into a Pandas DataFrame
273     df = pd.read_excel(excel_file_path)
274
275     # Create the table
276     create_table_query = f"CREATE TABLE {table_name} ({', '.join([f'{col} VARCHAR(255)' for col in df.columns])})"
277     cursor.execute(create_table_query)
278
279
```

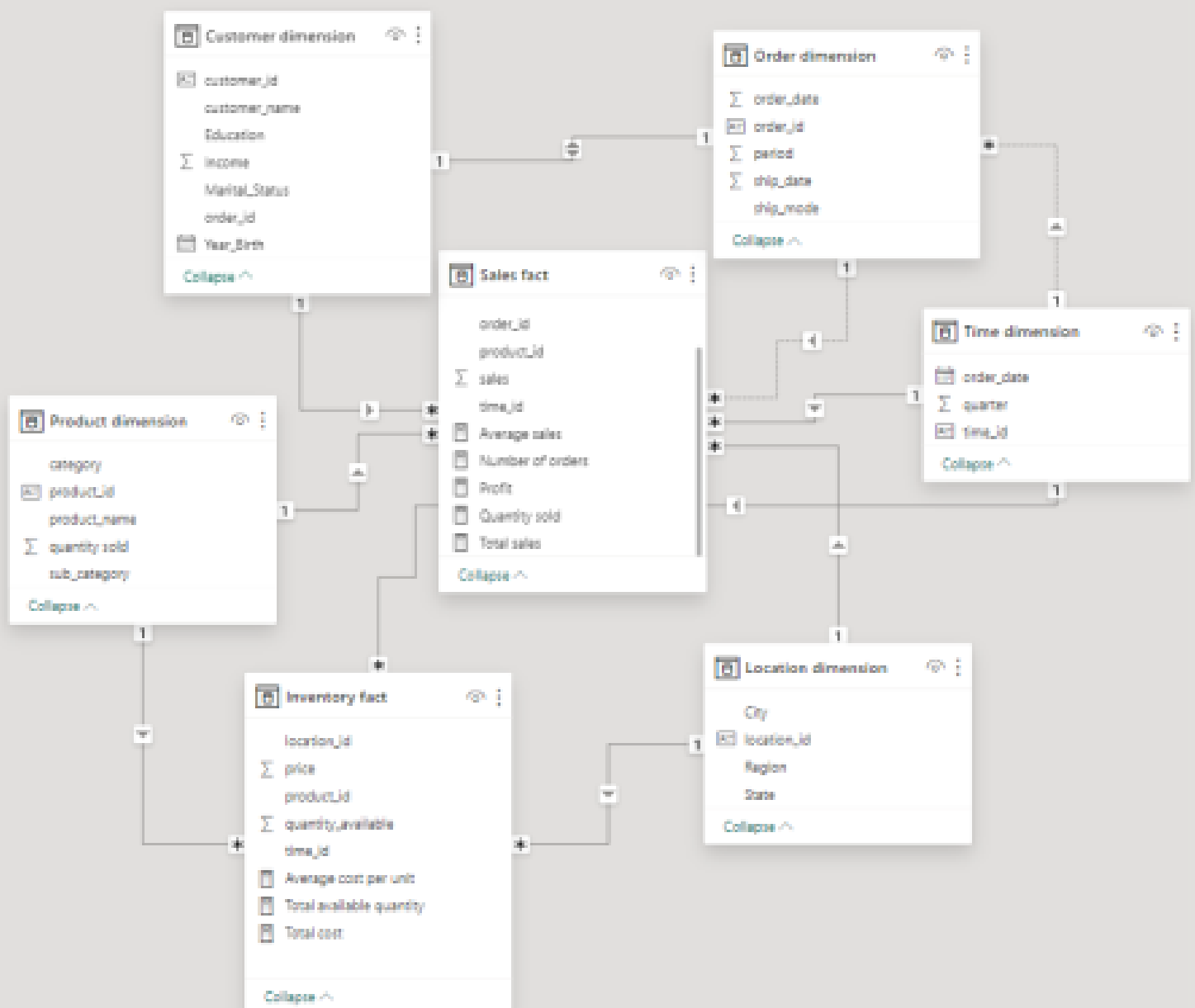
## DEFINED FUNCTION TO INSERT THE DATA

```
302 # %%
303 def insert_data_from_excel(cursor, excel_file_path, table_name):
304     # Load Excel data into a Pandas DataFrame
305     df = pd.read_excel(excel_file_path)
306
307     # Create the INSERT query
308     columns = ', '.join(df.columns)
309     values_placeholder = ', '.join(['%s' for _ in df.columns])
310     insert_query = f"INSERT INTO {table_name} ({columns}) VALUES ({values_placeholder})"
311
312     # Iterate through rows and insert data
313     for index, row in df.iterrows():
314         cursor.execute(insert_query, tuple(row))
315
```



# Schema choice: Snowflake

We decided to implement the Snowflake Schema to optimize data organization and minimize redundancy in our database. Breaking down dimension tables into interconnected tables facilitates more efficient updates and modifications. While this may result in increased query complexity due to necessary joins, the overall benefits encompass improved scalability, simplified maintenance, and enhanced performance in specific scenarios.



## 2.IMPLEMENTATION

### 2.3 DATA STORAGE :

#### 2.3.2.FACT:

We identified two Fact tables :

- **Sales fact** that holds the informations about the order, the product, the customer, the time, the location, and contained measures like average sales, number of orders and profit.
- **Inventory fact** that holds the details about the product, time, locations, quantity available and price.

The two fact tables allowed us to gain insight about important keys related to sales.

## 2.IMPLEMENTATION

### 2.3 DATA STORAGE :

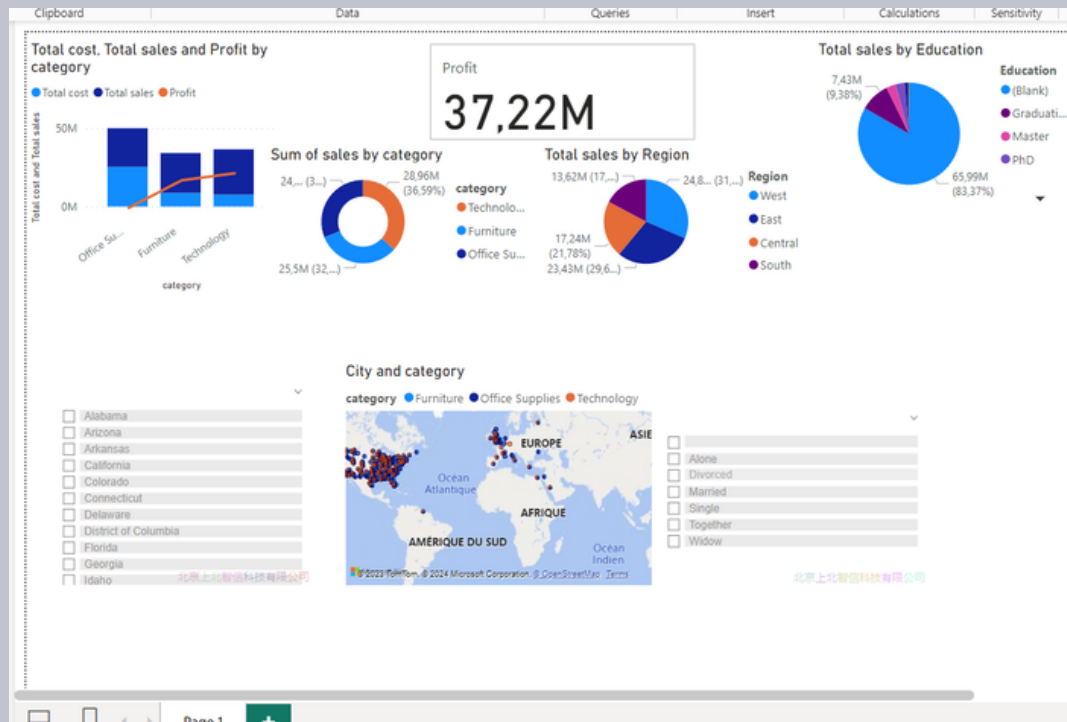
#### 2.3.3.DIMENSIONS:

To enrich our analysis, we implied various dimension tables each providing a distinct perspective on the data. These dimensions include:

- **Customer dimension:** represents the customer information.
- **Product Dimension:** represents the product information.
- **Location dimension:** represents the different locations of transactions.
- **Order Dimension:** Represents different characteristics of the order.
- **Time Dimension:** Contains information about the order time, shipment...

# 2.IMPLEMENTATION

## 2.4 DATA VISUALIZATION:



- The total profit amount is 37.22 M
- Technology prove to be the most profitable item generating 21.29M while office supplies are the least profitable with lost of 855.26 K
- California is our top selling state with 15.62M in sales followed by new york with 10.72
- Most of sales are coming from married people with 941 K and together generating 527 K

## INTERPRETATION:

- the average quantity available is half about 5 M and it ranges from 0 to 10 M.
- "Total Sales by Product Category" provides a comprehensive view of our sales distribution across different product categories.
- "Total Sales by Product Category" provides a comprehensive view of our sales distribution across different regions
- The Product Profitability Analysis, visualized through a stacked bar chart, serves as a powerful tool for businesses seeking to optimize their product offerings, pricing strategies, and overall profitability.

### 3. CONCLUSION

In summary Our research focuses on understanding how people shop in the ever-changing retail world of the USA. We use a big dataset covering customer demographics, buying habits, product preferences, and inventory levels to help businesses make smart decisions. We also explore global trends in office supplies, furniture, and technology retail, studying both big and new players to provide key insights for businesses navigating the dynamic global market.

# CHALLENGES FACED

We Faced several Obstacles throughout our project journey:

- **Challenges in data gathering phase:**  
The research was not easy, we faced difficulties in tracking the ideal dataset for our project.
- **challenges in data quality:** managing numerous duplicates and navigating through misleading data introduced complexity to the project. insuring the accuracy and reliability of the data became a crucial task.
- **Challenges in using talend:** Our group encountered challenges that prevented us from utilizing talend effectively.
- **Challenges in Power BI:** one oversight has the potential to affect the entire of our progress. To overcome these challenges, innovative problem-solving and iterative methods were employed to ensure the successful execution of our data management and visualization project.