

# **Data Warehouse (DWH) and ETL Pipeline Project**

## **Project Overview**

This document outlines the design, implementation, and analysis of a Data Warehouse (DWH) project. The project focuses on integrating sales, customer behavior, and delivery datasets into a comprehensive DWH, supported by an ETL pipeline built with Python, and visualizing the data using Power BI.

## **Business Objectives**

The main objectives of this project are:

- To create a centralized Data Warehouse for sales, customer behavior, and delivery data.
- To implement an ETL pipeline for data extraction, transformation, and loading into the DWH.
- To develop interactive and insightful Power BI dashboards for data analysis and visualization.
- To support business decision-making with comprehensive data insights.

## **Project Architecture**

The project architecture consists of the following components:

- Data Sources: Various CSV files and databases containing sales, customer behavior, and delivery data.
- ETL Pipeline: A Python-based ETL process to extract, transform, and load data into the DWH.
- Data Warehouse: A PostgreSQL database designed to store integrated and processed data.

- Data Visualization: Power BI dashboards for analyzing and visualizing the data.

### **Tools and Technologies:**

- ETL Tool: Python
- Database: PostgreSQL
- Visualization Tool: Power BI

### **Data cleaning :**

The point of this process is to drop missing values or fill in empty cells to get a more complete dataset with fewer errors

Code :

```
import pandas as pd
import os

#customer behavior
file_path= ''
df=pd.read_csv(file_path)

df.dropna( inplace=True)

df.drop(columns=['Search_Result_Exploration',
'Customer_Reviews_Importance', 'Add_to_Cart_Browsing',
'Cart_Completion_Frequency', 'Cart_Abandonment_Factors',
'Saveforlater_Frequency',
'Review_Left', 'Review_Reliability',
'Review_Helpfulness',
'Personalized_Recommendation_Frequency',
'Recommendation_Helpfulness',
```

```

        'Shopping_Satisfaction', 'Service_Appreciation',
        'Improvement_Areas'], inplace=True)

customer_behavior_dtype_conversions = {
    'index': 'int64',
    'Order ID': 'int64',
    'Date': 'datetime64[ns]',
    'Status': 'category',
    'Fulfilment': 'category',
    'Sales Channel': 'category',
    'SKU': 'category',
    'Category': 'category',
    'Size': 'category',
    'ASIN': 'category',
    'Qty': 'int64',
    'currency': 'category',
    'Amount': 'float64'
}

for column, dtype in customer_behavior_dtype_conversions.items():
    if column in df.columns:
        try:
            df[column] = df[column].astype(dtype)
        except ValueError as e:
            print(f"Error converting column {column} to {dtype}: {e}")

print(df.dtypes)

print(df)

os.makedirs('', exist_ok=True)
cleaned_file_path = ''
df.to_csv(cleaned_file_path, index=False)
print(f"CleaneD dataset saved to {cleaned_file_path}")

```

```

#sales
sales_path= ''
df=pd.read_csv(sales_path)
print(df.dtypes)

df.dropna( inplace=True)

df.drop (columns=['ship-service-level','Style','Courier
Status','ship-city','ship-state','ship-postal-code',
            'ship-country','promotion-ids','B2B','fulfilled-by','Unnamed:
22'], inplace=True)

os.makedirs('' , exist_ok=True)
cleaned_file = ''
df.to_csv(cleaned_file, index=False)
print(f"Cleanded dataset saved to {cleaned_file}")

#delivery
delivery_path= ''
df=pd.read_csv(delivery_path)
print(df.dtypes)

df.dropna( inplace=True)

df.drop(columns=
['Store_Latitude','Store_Longitude','Drop_Latitude','Drop_Longitude','Pick
up_Time','Weather'
    , 'Traffic','Vehicle','Delivery_Time'] , inplace=True)

print(df)

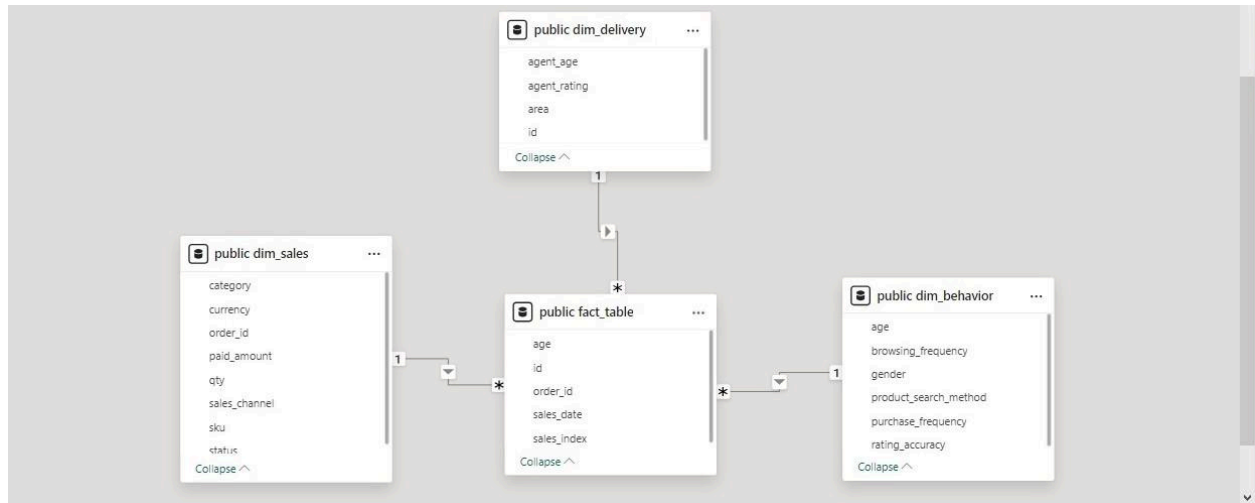
os.makedirs('' , exist_ok=True)
cleaned_file = ''
df.to_csv(cleaned_file, index=False)
print(f"Cleanded dataset saved to {cleaned_file}")

```

## Data Warehouse Design:

### Schema Design:

The DWH follows a star schema design, with fact tables connected to dimension tables.



### Staging tables sql script :

```
CREATE TABLE staging_sales (  
  index TEXT,  
  Order_ID TEXT,  
  Date timesatmp,  
  Status varchar(200),  
  Fulfilment varchar(200),  
  Sales_Channel varchar(200),  
  SKU varchar(200),  
  Category varchar(200),  
  Size varchar(200),
```

```
    ASIN varchar(200),  
    Qty int,  
    Currency varchar(200),  
    Amount numeric  
);
```

```
CREATE TABLE staging_delivery (  
    Order_ID varchar(200),  
    Agent_Age int ,  
    Agent_Rating numeric,  
    Order_Date date,  
    Order_Time TEXT,  
    Area varchar(200),  
    Category varchar(200)  
);
```

```
CREATE TABLE staging_behavior (  
    Timestamp Timestamp,  
    Age int ,  
    Gender varchar(200),  
    Purchase_Frequency varchar(200),  
    Purchase_Categories varchar(200),  
    Browsing_Frequency varchar(200),  
    Product_Search_Method varchar(200),  
    Personalized_Recommendation_Frequency varchar(200),  
    Rating_Accuracy TEXT  
);
```

## Fact and dim tables sql script :

```
CREATE TABLE DIM_SALES (  
    order_id text PRIMARY KEY,  
    status varchar(200) ,  
    sales_channel varchar(200),  
    sku varchar(200),  
    category varchar(200),  
    qty int,  
    currency varchar(200),  
    paid_amount numeric  
);
```

```
INSERT INTO DIM_SALES (order_id, status, sales_channel, sku, category  
,qty , currency , paid_amount)  
SELECT order_id, status, sales_channel, sku, category ,qty , currency ,  
paid_amount  
FROM staging_sales  
ON CONFLICT (order_id) DO NOTHING;
```

```
CREATE TABLE DIM_BEHAVIOR (  
    customer_age int PRIMARY KEY,  
    gender varchar(50),  
    purchase_frequency varchar(500),  
    browsing_frequency varchar(500),  
    product_search_method varchar(500),  
    rating_accuracy text  
);
```

```
INSERT INTO DIM_BEHAVIOR (age ,  
    gender ,  
    purchase_frequency ,  
    browsing_frequency ,
```

```
    product_search_method ,  
    rating_accuracy )  
SELECT age ,  
    gender ,  
    purchase_frequency ,  
    browsing_frequency ,  
    product_search_method ,  
    rating_accuracy  
FROM staging_behavior  
ON CONFLICT (age) DO NOTHING;
```

```
select * from dim_behavior
```

```
CREATE TABLE DIM_DELIVERY (  
    delivery_id text PRIMARY KEY,  
    agent_age int,  
    agent_rating numeric ,  
    area varchar(200)  
);
```

```
INSERT INTO DIM_DELIVERY (  
    id ,  
    agent_age ,  
    agent_rating ,  
    area )  
SELECT id ,  
    agent_age ,  
    agent_rating ,  
    area  
FROM staging_delivery
```

```
select * from DIM_DELIVERY ;
```



```

CREATE TABLE fact_table (
  sales_index text PRIMARY KEY,
  order_id text REFERENCES DIM_SALES(order_id),
  sales_date text,
  id text REFERENCES DIM_DELIVERY(id),
  age int REFERENCES DIM_BEHAVIOR(age)
);

```

--columns names are the same from the raw tables to the staging tables while inserting data  
 --making data type of columns in staging tables and dim and act tables the same would make insertion of data more smooth

```

INSERT INTO fact_table (
  sales_index,
  order_id,
  sales_date,
  id,
  age
)
-- Replace NULL with 'NA' (or adjust the value as needed)
SELECT
  COALESCE(staging_sales.sales_index, 'NA'),
  staging_sales.order_id,
  staging_sales.sales_date,
  staging_delivery.order_id,
  staging_behavior.age
FROM staging_sales
FULL OUTER JOIN staging_delivery
  ON staging_sales.order_id = staging_delivery.order_id
FULL OUTER JOIN staging_behavior
  ON staging_sales.category = staging_behavior.category
-- Ensure no duplicate insertions based on primary key (if necessary)

```

```
WHERE staging_sales.sales_index IS NOT NULL  
ON CONFLICT (sales_index) DO NOTHING;
```

```
SELECT * FROM fact_table
```

## ETL Pipeline

The ETL process involves three main steps:

1. Extract: Data is extracted from various sources.
2. Transform: Data is cleaned, transformed, and normalized.
3. Load: Data is loaded into the PostgreSQL DWH.

Code :

```
import pandas as pd  
from sqlalchemy import create_engine  
import logging  
  
logging.basicConfig(level=logging.INFO)  
logger = logging.getLogger(__name__)  
  
# Database info  
DB_USER = ''  
DB_PASSWORD = ''  
DB_HOST = ''  
DB_PORT = ''  
DB_NAME = ''  
  
# 1 - Extracting process
```

```

file_path = ''

def extract_data(file_path):
    logger.info("Extracting data from %s", file_path)
    try:
        data = pd.read_csv(file_path)
        logger.info("Data extracted successfully.")
        return data
    except Exception as e:
        logger.error("Error extracting data: %s", e)
        raise

# 2 - Transformation process
# Transforming CSV file or editing data types or cleaning

def transform_data(data):
    logger.info("Transforming data.")
    try:
        data = data.dropna()          # Remove rows with missing values
        logger.info("Data transformed successfully.")
        return data
    except Exception as e:
        logger.error("Error transforming data: %s", e)
        raise

# 3 - Loading process
# Loading data into the database

def load_data(data, table_name):
    logger.info("Loading data into %s table.", table_name)
    try:
        engine =
create_engine(f'postgresql://{DB_USER}:{DB_PASSWORD}@{DB_HOST}:{DB_PORT}/{
DB_NAME}')
        with engine.connect() as connection:
            data.to_sql(table_name, connection, if_exists='replace',
index=False)
        logger.info("Data loaded successfully.")
    except Exception as e:
        logger.error("Error loading data: %s", e)

```

```
        raise

def main():
    # ETL process
    file_path = ''
    table_name = 'sales_report'

    try:
        # Extract
        data = extract_data(file_path)

        # Transform
        data = transform_data(data)

        # Load
        load_data(data, table_name)

        logger.info("ETL process completed successfully.")
    except Exception as e:
        logger.error("ETL process failed: %s", e)

if __name__ == "__main__":
    main()
```

## Data Visualization with Power BI.

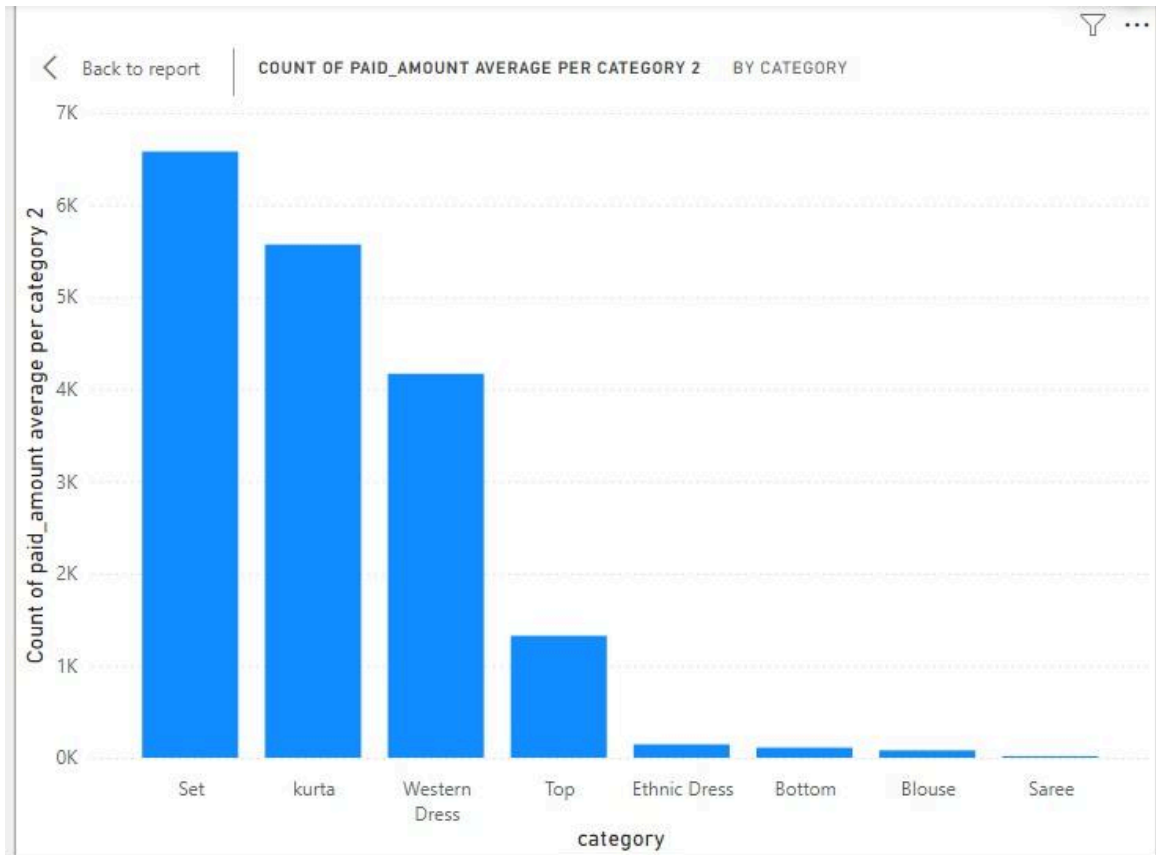
### Power BI Implementation

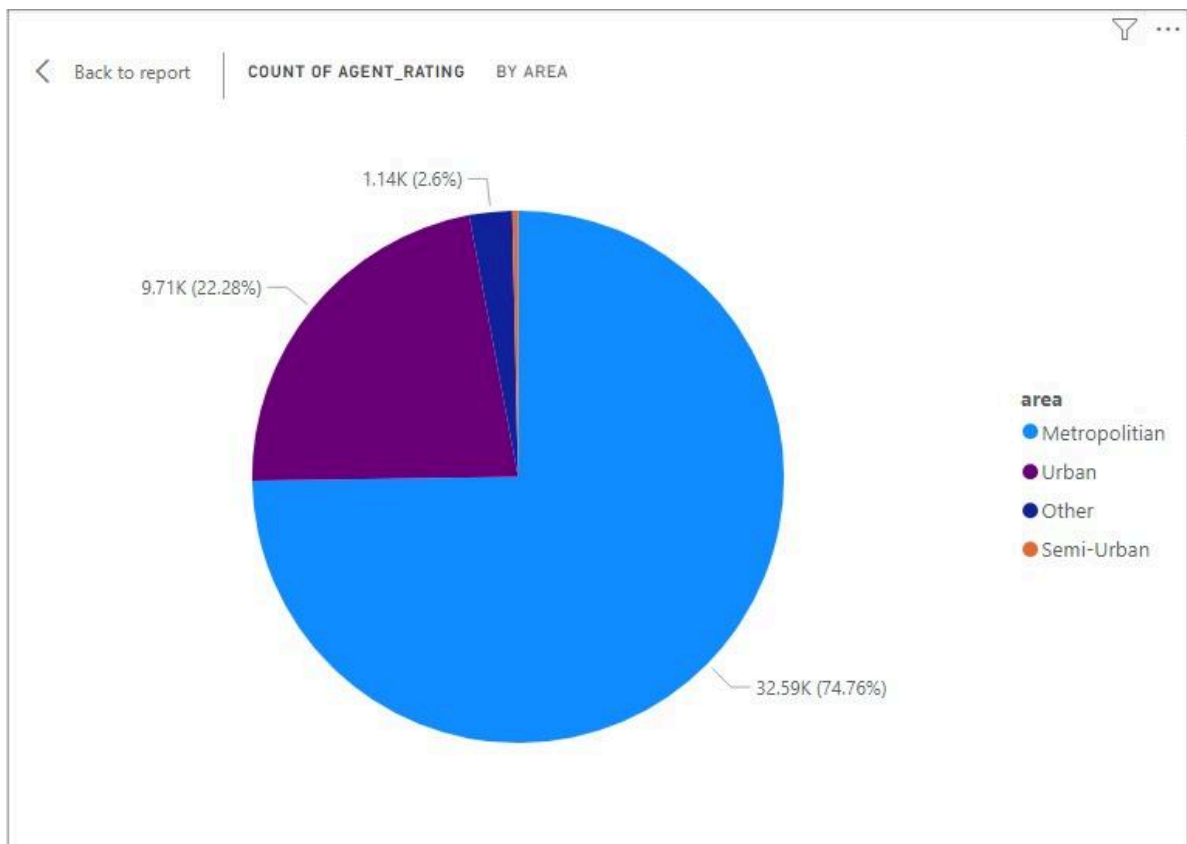
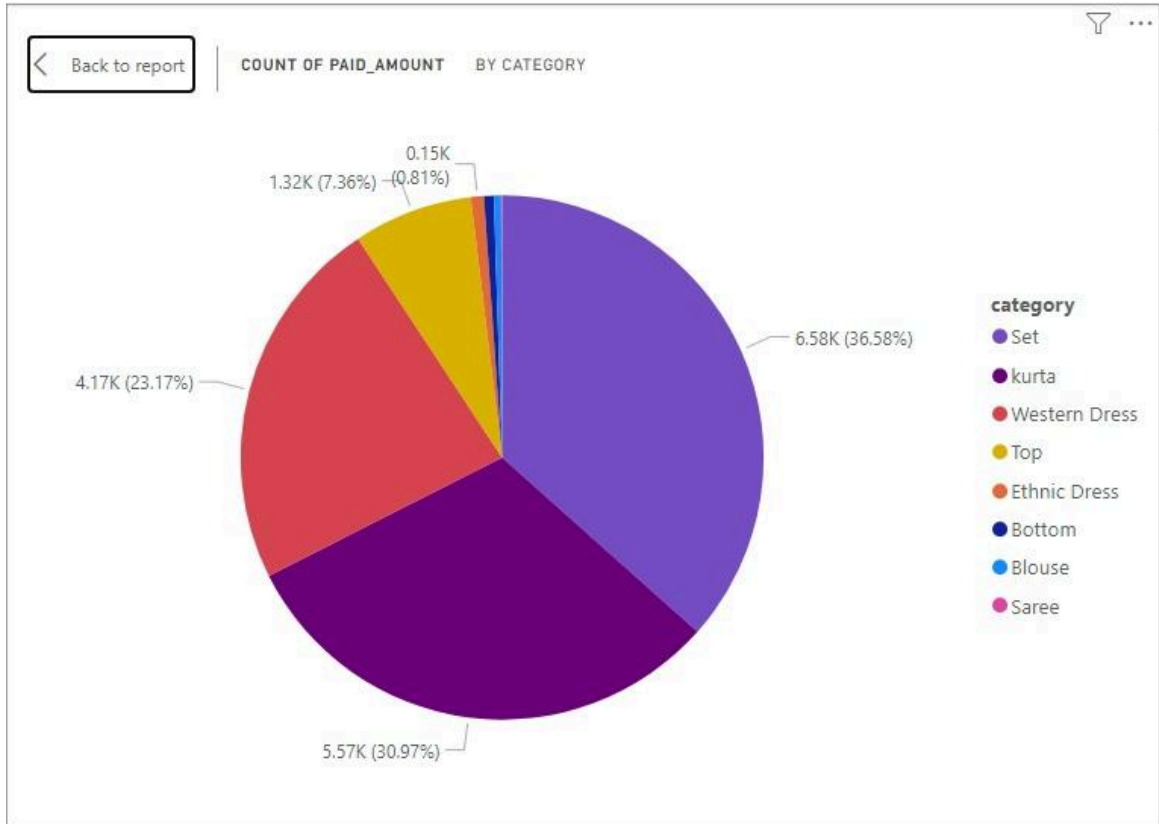
connected Power BI to the PostgreSQL DWH.

The screenshot displays the Power BI Desktop interface. On the left, the 'Data' pane shows a list of tables under the 'localhost: dwh project and etl [7]' connection. The table 'public.dim\_delivery' is selected. The central area shows a preview of the 'public.dim\_delivery' table, which contains 20 rows of data. The columns are 'delivery\_id', 'agent\_age', 'agent\_rating', 'area', and 'public.fact\_table'. The 'area' column contains values like 'Metropolitan' and 'Urban'. The 'public.fact\_table' column contains the value 'Table' for all rows. On the right, the 'Visualizations' pane is visible, showing various chart types and a 'Build visual' section. The 'Values' section has a text box 'Add data fields here'. The 'Drill through' section has a 'Cross-report' toggle set to 'Off' and a 'Keep all filters' toggle set to 'On'. The 'Add drill-through fields here' section is also visible.

delivery_id	agent_age	agent_rating	area	public.fact_table
aaar186826409	38	5.0	Metropolitan	Table
aabo036315307	37	4.5	Urban	Table
aabo151954572	39	4.7	Metropolitan	Table
aabt975198794	20	4.5	Metropolitan	Table
aabx020665819	38	5.0	Metropolitan	Table
aaca221183558	23	4.5	Metropolitan	Table
aack303984035	25	4.9	Metropolitan	Table
aacn882394574	26	4.8	Urban	Table
aacr165500044	21	4.6	Metropolitan	Table
aacu540600648	20	4.9	Metropolitan	Table
aaeg065438532	24	5.0	Urban	Table
aaem314845836	35	4.8	Other	Table
aaer219698621	26	4.9	Metropolitan	Table
aaff291190267	39	5.0	Metropolitan	Table
aافر395199003	22	4.8	Urban	Table
aagk068425859	25	5.0	Metropolitan	Table
aagp176989136	22	4.8	Metropolitan	Table
aaha266438982	26	5.0	Other	Table
aahf239075270	38	4.0	Metropolitan	Table
aahj033176282	37	4.8	Metropolitan	Table
aahu068587422	38	4.4	Metropolitan	Table
aaif308012919	38	4.6	Urban	Table









## **Conclusion**

The DWH and ETL pipeline project successfully integrates sales, customer behavior, and delivery datasets, providing valuable insights through interactive Power BI charts. This implementation supports data-driven decision-making and enhances business analysis capabilities.