Debre Birhan University

Collage Of Computing

Department OF Software engineering

# Fundamental OF Big Data Analytics And BI

## End To End Pipeline

Prepared by - Surafel G/selassie
Id- 0430/13

Github Link -
https://github.com/0001sura/bigdata.git

# Table of content

# ETL Pipeline Documentation

## Introduction

An ETL (Extract, Transform, Load) pipeline is essential for processing and managing large datasets efficiently. This documentation provides an in-depth guide to an ETL pipeline designed to process an e-commerce dataset. The pipeline automates data extraction, transformation, and loading into a PostgreSQL database.

## Objective

The primary objective of this ETL pipeline is to:

1. Extract raw data from CSV files.
2. Clean and transform data to ensure consistency and accuracy.
3. Load processed data into a PostgreSQL database for storage and further analysis.
4. Automate the ETL process for efficiency and scalability.

## Prerequisites

Before running the ETL pipeline, ensure that you have the following dependencies installed:

pip install pandas sqlalchemy psycopg2

## System Requirements

- **Operating System**: Windows, macOS, or Linux
- **Python Version**: 3.7 or higher
- **Database**: PostgreSQL 12 or higher
- **Memory Requirement**: At least 4GB RAM for large datasets

**ETL Process**

**1. Data Extraction and Transformatio(clean_and_transform())**

**a) Load Raw Data**

The raw dataset is read from:

data/raw/Pakistan Largest Ecommerce Dataset.csv

The dataset contains date columns which are parsed during loading.

**b) Handling Missing Values**

- **Categorical Columns:**

    o sales_commission_code: Missing values are replaced with 'Unknown'.
    o BI Status: '#REF!' values are replaced with 'Unknown', and missing values are also set to 'Unknown'.

- **Numerical Columns:**

    o Columns: price, grand_total, discount_amount, qty_ordered
    o Converted to numeric types; invalid values are coerced to NaN.
    o Missing values are replaced with the median for robustness.

**c) Data Cleaning**

- **Duplicate Removal:** All duplicate rows are dropped.
- **Date Columns Handling:** created_at, Working Date, and M-Y are converted to datetime format.

**d) Dropping Unnecessary Columns**

The following columns are removed:

MV, Year, Month, Customer Since, Unnamed: 19-25

## 2. Load Data to PostgreSQL (load_to_postgres(cleaned_df))

## a) Database Connection

The function connects to a PostgreSQL database using:

postgresql://postgres:1221@localhost:5432/bigdata

## b) Table Creation

- A table named ecommerce_data is created.
- Only the schema (column names and types) is written initially.

## c) Data Insertion

- Data is inserted in chunks of 10,000 rows for efficiency.
- Existing data is preserved, and new data is appended.

## 3. Execution

The script is executed by running:

```
if __name__ == "__main__":
    df = clean_and_transform()
    load_to_postgres(df)
```

## Error Handling

- If the CSV file is missing or corrupted, an appropriate error message is displayed.
- If database credentials are incorrect, the script terminates with a connection error message.
- If column names in the dataset change, the script will log an error and halt execution.

# Implementatoin (Code)

## 1. Import the necessary files

```
import pandas as pd
from sqlalchemy import create_engine
```

## 2. Load the data

```
def clean_and_transform():
    # Load raw data
    raw_df = pd.read_csv(
        'data/raw/Pakistan Largest Ecommerce Dataset.csv',
        low_memory=False,
        parse_dates=['created_at', 'Working Date', 'M-Y'],
        infer_datetime_format=True
    )
```

## 3. Data Cleaning And Transformation

```
    # --- Null Handling ---
    # Categorical columns
    raw_df['sales_commission_code'] =
raw_df['sales_commission_code'].fillna('Unknown')
    raw_df['BI Status'] = raw_df['BI Status'].replace('#REF!',
'Unknown').fillna('Unknown')

    # Numerical columns
    num_cols = ['price', 'grand_total', 'discount_amount', 'qty_ordered']
    for col in num_cols:
        raw_df[col] = pd.to_numeric(raw_df[col], errors='coerce')
        raw_df[col] = raw_df[col].fillna(raw_df[col].median())

    # Remove duplicates
    raw_df = raw_df.drop_duplicates()

    date_cols = ['created_at', 'Working Date', 'M-Y']
    for col in date_cols:
        raw_df[col] = pd.to_datetime(raw_df[col], errors='coerce')

    # --- Column Removal ---
    cols_to_drop = [
        'MV', 'Year', 'Month', 'Customer Since',
        'Unnamed: 19', 'Unnamed: 20', 'Unnamed: 21',
        'Unnamed: 22', 'Unnamed: 23', 'Unnamed: 24', 'Unnamed: 25'
    ]
    return raw_df.drop(columns=[c for c in cols_to_drop if c in raw_df.columns])
```
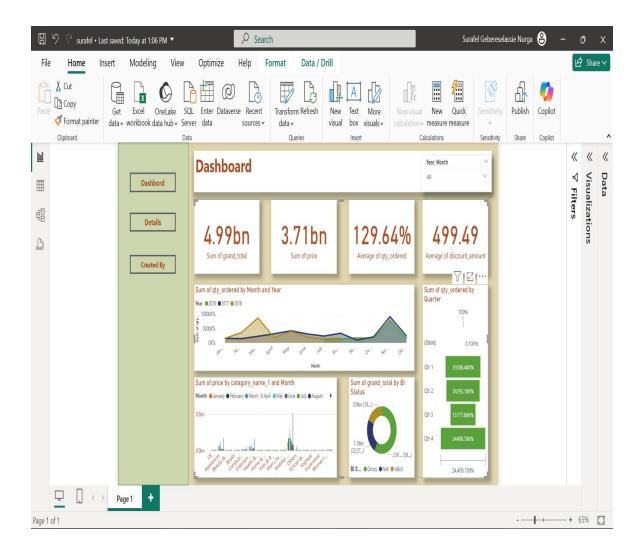
## 4. Load To POStgreSQL

```python
def load_to_postgres(cleaned_df):
    # Database connection
    engine = create_engine('postgresql://postgres:1221@localhost:5432/bigdata')

    cleaned_df.head(0).to_sql('ecommerce_data', engine, if_exists='replace', index=False)

    # Load in chunks (for large datasets)
    cleaned_df.to_sql('ecommerce_data', engine, if_exists='append', index=False, chunksize=10000)
    print("Data successfully loaded to PostgreSQL!")
    print(cleaned_df.shape)


if name == "main":
    df = clean_and_transform()
    load_to_postgres(df)
```

### Creating Table

```sql
CREATE TABLE ecommerce_data (
    item_id BIGINT,
    status VARCHAR(50),
    created_at TIMESTAMP,
    sku VARCHAR(100),
    price FLOAT,
    qty_ordered INT,
    grand_total FLOAT,
    increment_id VARCHAR(50),
    category_name_1 VARCHAR(100),
    sales_commission_code VARCHAR(100),
    discount_amount FLOAT,
    payment_method VARCHAR(50),
    working_date TIMESTAMP,
    bi_status VARCHAR(50),
    m_y TIMESTAMP
);
```
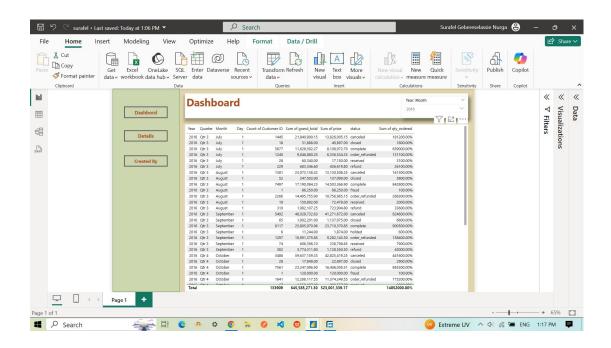
**Visualization (Dashboard)**

**Link to the visualization(powerbi)**

**https://app.powerbi.com/links/6W9FCQRh4j?ctid=1695066a-e388-40d1-8ed5-5d0b28ba9f80&pbi_source=linkShare**
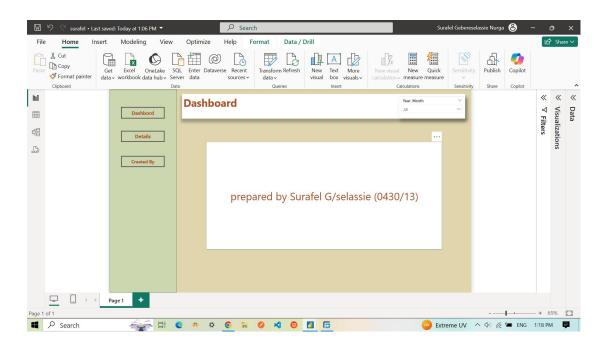
**1. Dashboard page**

## 2. Details Page



## 3. Prepaerd by page

## Performance Optimization

- **Batch Processing**: Data is inserted into PostgreSQL in chunks to handle large datasets efficiently.
- **Indexing**: Indexing key columns in the database improves query performance.
- **Logging**: Implementing logging for monitoring ETL operations.

## Security Considerations

- **Database Credentials**: Store credentials securely using environment variables.
- **Data Privacy**: Ensure sensitive customer data is anonymized if necessary.
- **Access Control**: Restrict database access to authorized users only.

## Future Enhancements

- **Automated Scheduling**: Implementing cron jobs or Airflow for scheduled ETL runs.
- **Data Validation**: Adding data validation checks before inserting into the database.

## Notes

- Ensure PostgreSQL is running and accessible at localhost:5432.
- Update database credentials as needed.
- Modify cols_to_drop if additional columns should be removed.

## Conclusion

This ETL pipeline automates data cleaning and loading into a PostgreSQL database, making it efficient for processing large e-commerce datasets. With improvements in performance optimization and security, the pipeline can scale for larger datasets and enterprise applications.