

# Comprehensive Documentation of End-to-End Data Processing and Visualization

## I. Summary :

The objective of this assignment was to demonstrate the ability to manage an end-to-end data pipeline and deliver insights in a professional and structured manner. The task involved leveraging public datasets available on [Open Data Paris](https://opendata.paris.fr/) to showcase skills in data engineering, analysis, and visualization.

### Theme and Subject

- **Theme:** *Administration et Finances Publiques*
- **Subject:** *Compte de Résultat* (Financial Results)
- **URL=** ["https://opendata.paris.fr/api/explore/v2.1/catalog/datasets/compte-de-resultat/records"](https://opendata.paris.fr/api/explore/v2.1/catalog/datasets/compte-de-resultat/records)

The focus was on analyzing the financial results dataset to extract meaningful insights and present them using intuitive and interactive visualizations.

### Tools and Technologies

- **Programming Language:** Python
- **Database:** PostgreSQL
- **Visualization Tool:** Metabase
- **Deployment:** Docker

### Assignment Workflow

1. **Data Retrieval:** Accessing the dataset through the provided API.
2. **Database Integration:** Storing and organizing the data in a PostgreSQL database.
3. **Data Analysis:** Developing algorithms to clean, transform, and analyze the data.
4. **Visualization:** Creating interactive dashboards in Metabase to display the results.
5. **Deployment:** Dockerizing the solution for ease of replication and sharing via GitHub.

This project highlighted expertise in data extraction, transformation, analysis, and visualization while adhering to best practices in software deployment and documentation.

## II. Accomplished Assignment :

### 1. Data Retrieval

Data retrieval marked the starting point of our journey, where we focused on accessing and gathering the necessary information to drive our analysis. Through a consistent and methodical approach, we ensured that all data was collected efficiently and accurately, forming a robust foundation for the next steps. Each step of this process is documented, offering insights into the under-the-hood operations for those interested in the technical details.

### 2.Database Integration

The integration of our retrieved data into a structured database allowed us to centralize and organize the information effectively. This step ensured seamless accessibility and reliability for subsequent tasks, laying the groundwork for meaningful insights. All database-related operations are fully documented for further review if needed.

### 3. Data Analysis

Our analysis utilized several tailored algorithms, each designed to extract meaningful insights from the data. Here is an overview of the algorithms used, along with their purpose and what we aimed to accomplish:

- **Profit Calculation:**  
This algorithm calculated the yearly profit by analyzing the difference between revenue (products) and expenses (charges). The goal was to determine the overall financial health and trends over time.
- **Exceptional Charges Analysis:**  
Focused on identifying and summing exceptional charges grouped by year, this algorithm provided insights into irregular expenses that could impact financial performance.
- **Profitability Analysis:**  
This algorithm analyzed core profits and core costs associated with operations to evaluate the efficiency and profitability of the business's core activities.
- **Expense Categorization:**  
Designed to group expenses by categories (e.g., operational, administrative) and analyze them by year, this algorithm helped identify major cost drivers and their trends over time.
- **Return on Investment (ROI) Analysis:**  
Calculated yearly returns on financial investments to evaluate how effectively the resources were being utilized to generate profits.

Each algorithm played a crucial role in uncovering trends and supporting data-driven decision-making. The detailed documentation ensures that anyone interested can explore the inner workings of these algorithms.

## 4. Visualization

**metabase** brought our analysis to life, transforming raw data and insights into intuitive and compelling visuals. By connecting to the database via a metabase, leveraging SQL scripting, and integrating the analytical results that we already developed in the data analysis section, we created comprehensive dashboards. These were then exported to PDF format and embedded into HTML pages, ensuring accessibility and clarity for all stakeholders.

## 5. Unit Testing

To ensure the reliability and accuracy of our work, we implemented a unit testing process. This step involved meticulously validating every function and algorithm to guarantee that they met the expected outcomes. The documentation provides a clear view of this process for those interested in understanding how we maintained high standards throughout the project.

## 6. Deployment

The final stage of deployment brought all components together, ensuring a smooth transition from development to operational use. This comprehensive process included finalizing all functionalities, organizing the project for ease of use, and delivering a cohesive and functional system. Every aspect of deployment has been documented to provide a roadmap for future enhancements or replication.

# III. Project Workflow Overview :

This project was divided into two main containers, each serving a specific role in the data pipeline and visualization workflow:

## 1. ETL (Extract, Transform, Load)

The ETL container was designed to handle the data ingestion and storage data. It includes the following steps:

This container encapsulates the entire ETL pipeline .

## Folder Schema :

Each file and folder has a distinct purpose:

- **common/**: Centralized configuration for database connections and shared utilities.
- **models/**: Maps database tables to Python objects for seamless data interaction using SQLAlchemy.
- **main.py**: Core logic of the ETL pipeline; orchestrates data extraction, transformation, and loading.
- **tools.py**: Utility functions to support transformation, cleaning, and handling of edge cases.
- **Dockerfile**: Defines the environment for running the ETL pipeline as a container.

- **requirements.txt**: Ensures consistent dependencies for the Python environment.

## 2. LTE (Load , Transform , Export)

The second container focuses on data analysis and visualization. It was built to perform the following:

1. **Website Creation:**
  - Developed a simple web application using Flask to display analysis results and dashboards.
2. **Analysis Algorithms:**
  - Created algorithms to explore and derive insights from our dataset, such as detecting trends .
3. **Dashboard Integration:**
  - Interactive dashboards were designed using Metabase, which provided a user-friendly interface for visualizing the data.
4. **PDF Deployment:**
  - Exported Metabase dashboards as PDFs and integrated them into the Flask web application for easy access.

### Folder Schema :

- **static/ Dashboard.pdf** : Exported visualization in PDF format
- **templates**

    /**dashboard.html** : HTML template for dashboard view

    /**index.html** : Homepage template

    /**table.html** : Template for displaying tables

- **app.py** : Flask app for serving the visualization and routes.
- **Processing\_funs.py** : Data analysis functions for extracting insights.
- **Dockerfile** : Docker image configuration for Flask app
- **requirements.txt** : Dependencies for the project

## 3. Orchestration

To manage and deploy these containers, docker-compose was used. This ensured seamless communication between the ETL and visualization components while maintaining isolated environments for each part of the project.

This overview highlights the modularity and efficiency of the workflow, enabling end-to-end data processing and visualization. The next sections will dive into each container, detailing the technologies, processes, and challenges encountered.

## IV. Future Optimizations :

we identify critical areas for improvement and propose targeted optimizations to enhance the overall performance, security, and reliability of the system:

1. **Strengthened Backend Security:**

The backend system can be made more secure by implementing robust security measures, such as token-based authentication. This approach would ensure secure access to APIs and enhance protection against unauthorized access and potential breaches.

2. **Improved Database Security:**

PostgreSQL can benefit from better security configurations to safeguard sensitive data. Measures such as enabling SSL connections, setting up strict role-based access control (RBAC), and enforcing data encryption both at rest and in transit would bolster database security, aligning it with best practices.

3. **Enhance Architectural Abstraction:**

Transform the ETL process from a single-use program into a multi-purpose solution. This improved architecture would feature an automated mapping mechanism capable of detecting data types from JSON files and dynamically generating the required tables and columns in the PostgreSQL database. With this approach, data ingestion would be streamlined, requiring only a list of target URLs to be provided for seamless storage in the database.

## V. Process Overview :

To begin, clone the Git repository to your local machine. Next, navigate to the project folder using the `cd` command. Once inside the directory, execute **`docker-compose up`** to initiate the setup process.

The **`docker-compose`** command will build the necessary containers and establish a default network, which will be automatically removed when the containers are shut down. This ensures a clean environment for each deployment.

After the setup, verify that the data has been successfully ingested by using a database viewer, such as **DBeaver**, to inspect the PostgreSQL database. Alternatively, visit the website hosted at **`localhost:8080`**, ensuring all specified ports are accessible. From there, click on **`sample_data`** to view the ingested data directly.

Once you've verified that the data has been successfully ingested, you can explore the additional analysis options provided by the platform. Navigate to the website hosted at **`localhost:8080`**. Beyond viewing the **`sample_data`**, you will find various tools and functionalities designed to facilitate deeper insights into the data.

These options include a Dashboard, and advanced analysis results generated from the algorithms applied to the ingested data. By clicking through the different sections of the website, you can seamlessly access and interpret these insights, making it easier to derive actionable conclusions and support decision-making processes.

## **VI. Conclusion**

This project successfully delivered a robust data pipeline capable of extracting and loading data into a PostgreSQL database while maintaining flexibility for future scalability. We implemented an efficient ETL process, designed an intuitive API, and provided visualization tools to enable insightful data analysis. By integrating modern tools and technologies, we ensured seamless deployment and usability, empowering users to access and manage data effortlessly.