

# Pipeline Documentation:

The data pipeline consists of three tasks: extract, transform, and load.

The extract task extracts data from three CSV files stored in the local file system.

The transform task cleans and aggregates the extracted data to produce a new DataFrame that contains the customer lifetime value (CLV) for each customer.

The load task loads the transformed data into a PostgreSQL database table called "customer\_ltv."

The pipeline is designed to run daily using Airflow, a popular open-source platform for creating, scheduling, and monitoring data pipelines. The pipeline is written in Python and uses several libraries, including Pandas, Psycopg2, and Airflow.

## Best Practices:

**Modularity:** The pipeline is modular, which makes it easier to maintain and update. Each task has a well-defined input and output, making it easy to test and debug individual tasks independently.

**Error Handling:** The pipeline includes error handling for each task, which ensures that the pipeline continues to run even if a task fails. Each task logs any errors or exceptions that occur, making it easier to troubleshoot and debug.

**Documentation:** The pipeline includes detailed documentation, including comments in the code and this document. The documentation makes it easier for other developers to understand how the pipeline works and how to modify or extend it.

## Recommendations:

The pipeline can be deployed to a cloud-based provider, such as Amazon Web Services (AWS), Microsoft Azure, or Google Cloud Platform (GCP). Some of the recommendations for deploying and running the pipeline in a cloud-based environment are:

- **Use a managed PostgreSQL database service:** Instead of running a PostgreSQL database on a virtual machine, use a managed PostgreSQL database service provided by the cloud provider. This will simplify database management and reduce operational overhead.

- **Use object storage:** Store the input CSV files in object storage (e.g., AWS S3, Azure Blob Storage, or Google Cloud Storage) instead of storing them in the local file system. This will make it easier to scale the pipeline and reduce the risk of data loss.
- **Use a containerized approach:** Use Docker to containerize the pipeline and deploy it to a container orchestration platform, such as Kubernetes. This will make it easier to manage and scale the pipeline and reduce the risk of environment-related issues.
- **Use a monitoring and alerting service:** Use a monitoring and alerting service, such as AWS CloudWatch, Azure Monitor, or GCP Stack driver, to monitor the pipeline and receive alerts if any issues occur. This will help ensure that the pipeline is running smoothly and reduce the risk of data loss or downtime.