

## Project: Data pipelines with Apache Airflow

### Documentation of the Pipeline

The DAG pipeline is an ETL (extract, transform, load) process that extracts customer, order, and payment data from CSV files, transforms it by merging, grouping, and aggregating the data, and then loads the transformed data into a PostgreSQL database.

### Best Practices

Here are some of the best practices used during the implementation of the pipeline

- i. **Modularization/ Separation of concerns:-**  
The pipeline has been broken down into three distinct stages extracting data, transforming data, and loading data. Each stage is encapsulated in a separate function, which makes it easier to test and debug each stage independently. Moreover, this modularization approach also makes it easier to maintain and update the pipeline in the future.  
This makes the code easier to understand, test, and maintain.
- ii. **Error Handling:-**  
The pipeline has been designed to handle errors that might occur during data extraction, transformation, or loading. For example, if there is an error in extracting data from the CSV files, an exception will be raised, and the pipeline will stop processing. Similarly, if there is an error in transforming or loading data, the pipeline will log the error and retry the operation after a specified interval.
- iii. **Logging and Monitoring:-**  
The pipeline logs every step of the process, including the start and end time of each task, the input and output of each task, and any errors or exceptions that occur during the process. This information is stored in a centralized log database, which can be used to monitor and troubleshoot the pipeline.
- iv. **Parameterization:-**  
The use of parameterization allows for flexibility in the DAG data pipeline. The parameters can be easily modified to handle different data sources, file formats, and other configurations. This reduces the need for hard-coded values and makes the pipeline more adaptable to changes.
- v. **Code versioning:-**  
The code is stored in a Git repository, which allows for version control and collaboration among team members. Each change to the code is tracked, making it easier to revert to a previous version if necessary.

### Recommendations for Deployment and Running the Pipeline in a Cloud-based Provider

By following these recommendations, you can ensure that your DAG pipeline is deployed and run efficiently and securely on a cloud-based provider.

When deploying this data pipeline, it is recommended to:

- a. **Use Managed Services:-**

Instead of setting up your own infrastructure, consider using managed services provided by your cloud provider. For example, you can use Google Cloud SQL to host your PostgreSQL database, which eliminates the need to manage the database infrastructure yourself. Similarly, you can use Google Cloud Storage to store your CSV files, which eliminates the need to manage the storage infrastructure yourself.

**b. Utilize scalable resources**

If you expect a large amount of data to be processed by the pipeline, you should consider scaling the infrastructure horizontally to handle the increased load. For example, you can create multiple instances of the pipeline running in parallel, or use a load balancer to distribute the load across multiple instances.

**c. Security:-**

Make sure to follow the security best practices recommended by your cloud provider. For example, you can use Google Cloud IAM to manage access to your pipeline resources, and use Google Cloud KMS to encrypt sensitive data such as database passwords. You should also make sure to regularly update your pipeline components to address any security vulnerabilities that may arise.

**d. Backup and Recovery:-**

Make sure to configure backup and recovery procedures for your pipeline data. For example, you can use Google Cloud Storage to store backups of your PostgreSQL database, and use Google Cloud IAM to manage access to the backups. You should also regularly test your backup and recovery procedures to ensure they are working as expected.

**e. Cost Optimization:-**

Make sure to optimize the cost of running your pipeline. For example, you can use Google Cloud Cost Management to monitor your pipeline usage and optimize your resource allocation. You should also make sure to regularly review your resource usage and adjust your infrastructure as needed to minimize costs.