

Documentation of the Data Pipeline

This pipeline is designed to extract, transform, and load data from a CSV file containing customer call logs into a PostgreSQL database. The pipeline consists of three main functions:

- i. **extract_data():**
This function uses pandas to extract the data from the CSV file and cache it in Redis for faster retrieval.
- ii. **transform_data():**
This function retrieves the data from Redis cache, transforms it by cleaning, structuring, and formatting it, and returns the transformed data.
- iii. **load_data(transformed_data):** This function connects to the PostgreSQL database, creates a table to store the data, and inserts the transformed data into the table.

The main function **data_pipeline()** calls these functions in the order listed above to execute the complete pipeline.

Best Practices during implementation

The following best practices were used for this pipeline.

- a) Separation of Concerns: Each function in the pipeline has a specific responsibility, which enhances the readability and maintainability of the code.
- b) Caching: Data was cached in Redis for faster retrieval, reducing the data extraction time and increasing the efficiency of the pipeline.

Recommendations

For Deployment and Running the Pipeline with a Cloud-Based Provider, the below are the recommendations:

- **Use a containerization tool:** To deploy the pipeline to a cloud-based provider, it is recommended to use a containerization tool such as Docker. This allows the pipeline to be packaged as a container and deployed to various cloud platforms without compatibility issues.
- **Utilize a managed database service:** Instead of deploying a PostgreSQL database on a virtual machine, it is recommended to utilize a managed database service provided by cloud platforms such as AWS RDS or Azure Database. This reduces the administrative burden and provides scalability, availability, and durability for the database.
- **Implement Continuous Integration and Deployment:** Implementing a CI/CD pipeline automates the process of building, testing, and deploying the pipeline to the cloud. This provides faster feedback on code changes, reduces the risk of deployment errors, and increases the agility of the pipeline.
- **Set up monitoring and alerting:** To monitor the pipeline's health and performance, it is recommended to set up monitoring and alerting using tools such as Prometheus and Grafana. This allows for proactive maintenance and troubleshooting of the pipeline.
- **Implement security measures:** When deploying the pipeline to the cloud, it is recommended to implement security measures such as SSL encryption for database connections and access control to prevent unauthorized access and data breaches.