

Azure Meal- CD

(Code Documentation)

Overview:

The **Meal Project** is a cloud-based ETL solution hosted on Microsoft Azure. It processes raw data stored in Azure Blob Storage, performs transformations using Azure Function App, and ingests the results into an Azure SQL Database. The processed data is visualized in Power BI to generate actionable insights.

Azure Blob Etl:

Azure Function TimerTrigger ETL Process:

This Python script implements an ETL (Extract, Transform, Load) pipeline using Azure Function. It processes blob files from Azure Blob Storage, extracts data, transforms it, and loads the results into a database.

The script uses the following libraries and modules:

- 1- **logging**: For logging messages.
- 2- **azure.functions**: To define Azure Functions.
- 3- **azure.storage.blob**: For interaction with Azure Blob Storage.
- 4- **pandas**: For data manipulation.
- 5- **pymssql and pyodbc**: For database connections.
- 6- **dotenv**: To load environment variables from config.py file.
- 7- **re, os, traceback**: For regex operations, file manipulation, and error handling.

Key Constants and Variables

Config.py file having Name of the Azure Blob Storage container, Connection string for Azure Blob Storage., Folder paths for input, output, and faulty files.

Below the code screenshot is provided, **Config.py**

Config.py

```
config.py > ...
1  # config.py
2  import os
3  from dotenv import load_dotenv
4
5  # Load environment variables from .env file
6  load_dotenv()
7
8  # Database connection details
9  server_name = os.environ.get('SERVER')
10 user_name = os.environ.get('USER_NAME')
11 pwd = os.environ.get('PASSWORD')
12 db_name = os.environ.get('DATABASE')
13
14 # Other configurations, such as storage connection string
15 STORAGE_CONNECTION_STRING = os.environ.get('AzureWebJobsStorage')
16 CONTAINER_NAME = 'csv-files'
17 INPUT_FOLDER = 'InputFiles/'
18 OUTPUT_FOLDER = 'Archive/'
19 FAULTY_FOLDER = 'Faulty/'
```

Databases:

database.py file having insert_metdata and insert_weights function.

In insert_metdata we are inserting data group_id, file_code, filename, variety, number_of_samples and formulations.

In insert_weights we are inserting data group_id, file_code, wld, weight and datedifference.

Below the code screenshot is provided, **Database.py**

Database.py

```
1 import os
2 import logging
3 import traceback
4 import pymysql
5
6 def insert_metadata(file_metadata_records, conn):
7     # Ensure connection is valid
8     if conn is None:
9         logging.error("Database connection is None. Aborting insert.")
10        return
11
12    # Test the connection by running a lightweight query
13    try:
14        cursor = conn.cursor()
15        cursor.execute("SELECT 1") # Lightweight query to validate the connection
16    except pymysql.OperationalError as e:
17        logging.error(f"Database connection is not active. Error: {e}")
18        return
19    try:
20        for _, row in file_metadata_records.iterrows():
21            try:
22                # Check if the record exists using the condition
23                cursor.execute("""
24                    SELECT COUNT(1) FROM dbo.File_Metadata
25                    WHERE group_id = %s
26                    AND file_code = %s
27                    """, (row['group_id'], row['file_code']))
28
29                record_exists = cursor.fetchone()[0]
30
31                if record_exists == 0: # No record found, perform the insert
32                    cursor.execute("""
33                        INSERT INTO dbo.File_Metadata (Filename, Variety, Formulation, Number_of_Samples, UploadDate, group_id,
34                        file_code)
35                        VALUES (%s, %s, %s, %s, %s, %s, %s)
36                        """, (row['Filename'], row['Variety'], row['Formulation'], row['Number_of_Samples'], row['UploadDate'], row
37                        ['group_id'], row['file_code']))
38                    logging.info(f"Inserted File_Metadata record: {row.to_dict()}")
39                else:
40                    logging.info(f"Skipped File_Metadata record (duplicate): {row.to_dict()}")
41            except Exception as e:
42                logging.error(f"Error during database insertion: {traceback.format_exc()}")
43            except pymysql.Error as e:
44                logging.error(f"Error during insertion: {e}")
45            except pymysql.Error as db_error:
46                logging.error(f"Database error: {db_error}")
47        finally:
48            conn.commit()
49            cursor.close()
50
51 def insert_weights(weight_records, conn):
52     # Check if the connection is None
53     if conn is None:
54         logging.error("Database connection is None. Aborting insert.")
55         return
56
57     # Test the connection by running a lightweight query
58     try:
59         cursor = conn.cursor()
60         cursor.execute("SELECT 1") # Lightweight query to validate the connection
61     except pymysql.OperationalError as e:
62         logging.error(f"Database connection is not active. Error: {e}")
63         return
64     try:
65         for _, row in weight_records.iterrows():
66             try:
67                 # Check if the record exists using the condition
68                 cursor.execute("""
69                     SELECT COUNT(1) FROM dbo.Weightloss
70                     WHERE group_id = %s
71                     AND file_code = %s
72                     AND WLD = %s
73                     AND Weight = %s
74                     AND DateDifference = %s
75                     """, (row['group_id'], row['file_code'], row['WLD'], row['Weight'], row['DateDifference']))
76
77                 record_exists = cursor.fetchone()[0]
78
79                 if record_exists == 0: # No record found, perform the insert
80                     cursor.execute("""
81                         INSERT INTO dbo.Weightloss (group_id, file_code, WLD, Weight, DateDifference)
82                         VALUES (%s, %s, %s, %s, %s)
83                         """, (row['group_id'], row['file_code'], row['WLD'], row['Weight'], row['DateDifference']))
84                     logging.info(f"Inserted weight record: {row.to_dict()}")
85                 else:
86                     logging.info(f"Skipped weight record (duplicate): {row.to_dict()}")
87             except Exception as e:
88                 logging.error(f"Error processing weight record {row.to_dict()}: {e}")
89            except pymysql.Error as db_error:
90                logging.error(f"Database error: {db_error}")
91        finally:
92            conn.commit()
93            cursor.close()
94
95
```

Function_App

TimerTrigger Function

The function runs on a schedule (every 10 minutes) using the TimerTrigger decorator. It triggers the ETL process by processing blobs in a specified container.

The core ETL process:

1. Lists blobs from the input folder.
2. Downloads and processes each blob.
3. Transforms the data using `pandas`.
4. Inserts transformed data into the database.
5. Archives successful blobs or moves faulty blobs.

Extract:

- List blobs from the input folder using `container_client.list_blobs()`.
- Download blob content and read it into a pandas DataFrame.

Transform:

- Validate required columns (`Sample Number`, `Variety`, `Formulation`, `Number of Samples`).
- Reshape data using `pandas.melt()`.
- Calculate additional metrics (e.g., average weight loss).

Load:

- Insert metadata and weight data into the database using `insert_metadata` and `insert_weights` functions.
- Handle errors with rollback and proper logging.

Archive and Faulty Handling:

- Successful blobs are moved to the archive folder.
- Faulty blobs are moved to the faulty folder.

Below the code link is provided

<https://drive.google.com/drive/folders/1Fiv4TglDPFCA9nickwoiumRogaEMfbPU?usp=sharing>

Commands to Run App:

To run your Azure Function commands in the project environment using **CMD** (Command Prompt) or PowerShell, you need to ensure you're working in the correct environment and using the Azure Functions Core Tools.

Here's a step-by-step guide to run:

1- Install Azure Functions Core Tools (if not already installed):

```
npm install -g azure-functions-core-tools@4 --unsafe-perm true
```

2- Navigate to the project directory, open the cmd in the project directory and run the command. Use the **cd** command to navigate to the directory where your Azure Functions project is located. Example:

```
cd:path-to-your-project\MyAzureFunctionApp\.venv\Scripts\activate
```

3- Set Up Your Environment (If applicable):

- Ensure any virtual environments or dependencies (like Python environments) are activated. If you're using Python, you can activate the environment like this:

```
cd:path-to-your-project\MyAzureFunctionApp
```

4- Run the Azure Function Commands:

- Once you're in the correct environment, you can use the following commands:
- **Publish (Deploy) Changes:** If you made any changes to your function code and need to deploy it to Azure, use the command mention below :

```
func azure functionapp publish FunctionApp12144
```

5- Start the Function Locally or using portal: using portal you can simply click on If you want to test the function locally (instead of using Azure), use the command mention below :

```
func start --verbose
```

