

Build ETL Data Pipelines with PythonOperator using Apache Airflow



Estimated time needed: **90** minutes.

Project Scenario

You are a data engineer at a data analytics consulting company. You have been assigned a project to de-congest the national highways by analyzing the road traffic data from different toll plazas. Each highway is operated by a different toll operator with a different IT setup that uses different file formats. Your job is to collect data available in different formats and consolidate it into a single file.

Objectives

In this assignment, you will develop an Apache Airflow DAG that will:

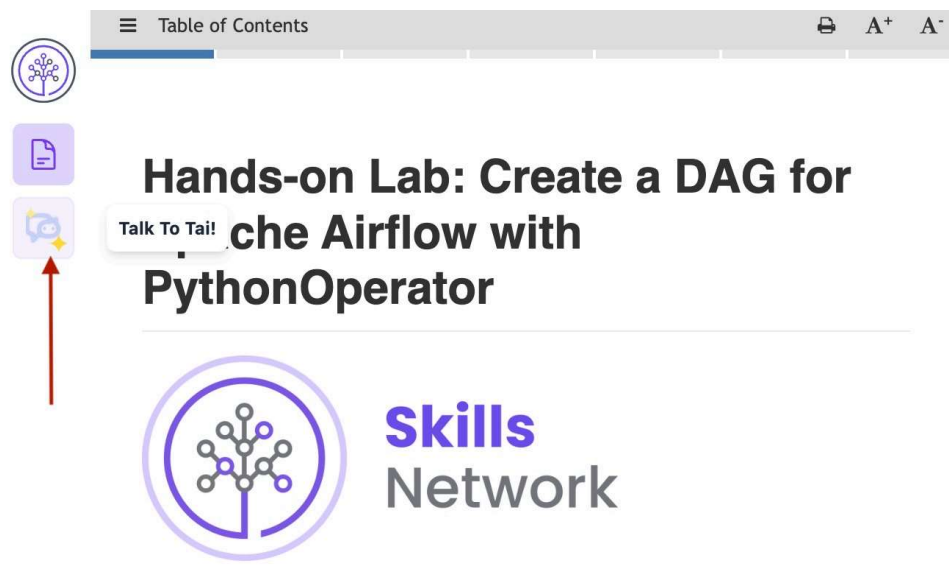
- Extract data from a csv file
- Extract data from a tsv file
- Extract data from a fixed-width file
- Transform the data
- Load the transformed data into the staging area

About Skills Network Cloud IDE

Skills Network Cloud IDE (based on Theia and Docker) provides an environment for hands-on labs for course and project-related labs. Theia is an open-source IDE (Integrated Development Environment) that can be run on a desktop or on the cloud. To complete this lab, you will be using the Cloud IDE based on Theia, running in a Docker container.

Important notice about this lab environment

Please be aware that sessions for this lab environment are not persistent. A new environment is created for you every time you connect to this lab. Any data you may have saved in an earlier session will get lost. To avoid losing your data, please plan to complete these labs in a single session. You can use the Tai AI assistant to complete this task.



Exercise 1: Prepare the lab environment

1. Start Apache Airflow.

Open Apache Airflow in IDE

Please wait until Airflow starts up fully and is active before you proceed further. If there is an error starting Airflow, please restart it.

2. Open a terminal and create a directory structure for staging area as follows:
/home/project/airflow/dags/python_etl/staging.

1. 1

1. `sudo mkdir -p /home/project/airflow/dags/python_etl/staging`

Copied! Executed!

3. Execute the following commands to avoid any permission issues in writing to the directories.

1. 1

1. `sudo chmod -R 777 /home/project/airflow/dags/python_etl`

Copied! Executed!

Exercise 2: Add imports, define DAG arguments, and define DAG

1. Create a file named `ETL_toll_data.py` in `/home/project` directory and add the necessary imports and DAG arguments to it.

| Parameter | Value |
|-------------|-------------------------------|
| owner | <You may use any dummy name> |
| start_date | today |
| email | <You may use any dummy email> |
| retries | 1 |
| retry_delay | 5 minutes |

2. Create a DAG as per the following details.

| Parameter | Value |
|--------------|--|
| DAG id | ETL_toll_data |
| Schedule | Daily once |
| default_args | as you have defined in the previous step |
| description | Apache Airflow Final Assignment |

Exercise 3: Create Python functions

1. Create a Python function named `download_dataset` to download the data set from the source to the destination. You will call this function from the task.

Source: <https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-DB0250EN-SkillsNetwork/labs/Final%20Assignment/tolldata.tgz>

Destination: `/home/project/airflow/dags/python_etl/staging`

2. Create a Python function named `untar_dataset` to untar the downloaded data set.

3. Create a function named `extract_data_from_csv` to extract the fields Rowid, Timestamp, Anonymized Vehicle number, and Vehicle type from the `vehicle-data.csv` file and save them into a file named `csv_data.csv`.
4. Create a function named `extract_data_from_tsv` to extract the fields Number of axles, Tollplaza id, and Tollplaza code from the `tollplaza-data.tsv` file and save it into a file named `tsv_data.csv`.
5. Create a function named `extract_data_from_fixed_width` to extract the fields Type of Payment code and Vehicle Code from the fixed width file `payment-data.txt` and save it into a file named `fixed_width_data.csv`.
6. Create a function named `consolidate_data` to create a single csv file named `extracted_data.csv` by combining data from the following files:
 - `csv_data.csv`
 - `tsv_data.csv`
 - `fixed_width_data.csv`

The final csv file should use the fields in the order given below:

Rowid, Timestamp, Anonymized Vehicle number, Vehicle type, Number of axles, Tollplaza id, Tollplaza code, Type of Payment code, and Vehicle Code

7. Create a function named `transform_data` to transform the `vehicle_type` field in `extracted_data.csv` into capital letters and save it into a file named `transformed_data.csv` in the staging directory.

Exercise 4: Create a tasks using PythonOperators and define pipeline

1. Create 7 tasks using Python operators that does the following using the Python functions created in Task 2.

1. `download_dataset`
2. `untar_dataset`
3. `extract_data_from_csv`
4. `extract_data_from_tsv`
5. `extract_data_from_fixed_width`
6. `consolidate_data`
7. `transform_data`

2. Define the task pipeline based on the details given below:

| Task | Functionality |
|-------------|----------------------------|
| First task | <code>download_data</code> |
| Second task | <code>unzip_data</code> |

| Task | Functionality |
|--------------|--|
| Third task | <code>extract_data_from_csv</code> |
| Fourth task | <code>extract_data_from_tsv</code> |
| Fifth task | <code>extract_data_from_fixed_width</code> |
| Sixth task | <code>consolidate_data</code> |
| Seventh task | <code>transform_data</code> |

Exercise 5: Save, submit, and run DAG

1. Save the DAG you defined.
2. Submit the DAG by copying it into `$AIRFLOW_HOME/dags` directory.
 - [Click here if your DAG does not get submitted properly.](#)
3. Use CLI or Web UI to unpause the task.
4. Observe the outcome of the tasks in DAG on the Airflow console.

Solution

- [Click here for the solution](#)

Authors

[Lavanya T S](#)

Ramesh Sannareddy

Other Contributors

Rav Ahuja

© IBM Corporation. All rights reserved.