**INFO-H420 Management of Data Science and Business Workflows**

# Assignment 3 : Workflows with Apache Airflow

Min Zhang 000586970

Yutao Chen 000585954

Prof. Dimitris SACHARIDIS

Nov 24 2023

# Table des matières

# Listings

# Table des figures

# Exercise 1 (10 points))

— Define the DAG. Create a DAG named `process_web_log` that runs daily.

— Create a task to scan for a log. Create a task named `scan_for_log` that scans a folder `the_logs` for a `log.txt` file and triggers the rest of the workflow.

— Create a task to extract data. Create a task named `extract_data`. This task should extract the `ipaddress` field from the web server log file and save it into a file named `extracted_data.txt`.

— Create a task to transform data. Create a task named `transform_data`. This task should filter out all the occurrences of ipaddress `198.46.149.143` from `extracted_data.txt` and save the output to a file named `transformed_data.txt`.

— Create a task to load the data. Create a task named `load_data`. This task should archive the file `transformed_data.txt` into a tar file named `weblog.tar`.

— Define the workflow that executes the aforementioned tasks in sequence.

Save the DAG you defined into a file named `process_web_log.py`. In the report, include code snippets describing how you achieved this.

## 1.1 Assumptions

1.The DAG does not depend on the success of its previous run.

2.The DAG runs on a calendar every day, including weekdays and weekends.

## 1.2 Solution

### 1.2.1 DAG definition

Listing 1.1: DAG definition with default arguments.

```python
default_args = {
    'owner': 'airflow',
    'depends_on_past': False,
    'start_date': datetime(2023, 1, 1),
    'email_on_failure': False,
    'email_on_retry': False,
    'retries': 1,
    'retry_delay': timedelta(seconds=1),
}

dag = DAG('process_web_log',
          default_args=default_args,
          description='DAG for processing web log',
          schedule_interval='@daily',
          catchup=False,
          tags=['DSBW'])
```

### 1.2.2 File path definition

Listing 1.2: Definition of file paths for logs and processed data.

```python
log_dir = '/usr/local/airflow'
log_file = f'{log_dir}/log.txt'
extracted_data_file = f'{log_dir}/extracted_data.txt'
transformed_data_file = f'{log_dir}/transformed_data.txt'
tar_file = f'{log_dir}/weblog.tar'
```

### 1.2.3 Function definition

Listing 1.3: Functions for scanning, extracting, transforming, and loading data.

```python
def scan_for_log(**context):
    if os.path.isfile(log_file):
        return log_file
```

```
4        else:
5            raise ValueError("log.txt not found")
6
7  def extract_data(**context):
8      task_instance = context['ti']
9      log_path = task_instance.xcom_pull(task_ids='scan_for_log')
10
11     with open(log_path, 'r') as file, open(extracted_data_file, '
           w') as out_file:
12         for line in file:
13             ip_address = line.split()[0]  # Assuming IP address
                   is the first element in the log line
14             out_file.write(ip_address + '\n')
15
16 def transform_data(**context):
17     with open(extracted_data_file, 'r') as file, open(
           transformed_data_file, 'w') as out_file:
18         for line in file:
19             if '198.46.149.143' not in line:
20                 out_file.write(line)
21
22 def load_data(**context):
23     with tarfile.open(tar_file, 'w') as tar:
24         tar.add(transformed_data_file, arcname='transformed_data.
               txt')
```

### 1.2.4 Task definition

Listing 1.4: Task definitions in the Airflow DAG.

```
1  scan_task = PythonOperator(
2      task_id='scan_for_log',
3      python_callable=scan_for_log,
4      provide_context=True,
```

```
 5        dag=dag)

 6

 7  extract_task = PythonOperator(

 8        task_id='extract_data',

 9        python_callable=extract_data,

10        provide_context=True,

11        dag=dag)

12

13  transform_task = PythonOperator(

14        task_id='transform_data',

15        python_callable=transform_data,

16        provide_context=True,

17        dag=dag)

18

19  load_task = PythonOperator(

20        task_id='load_data',

21        python_callable=load_data,

22        provide_context=True,

23        dag=dag)
```

# Exercise 2 (2 points)

Do a test run for each of the tasks you defined. Once all works as expected, do a test run for the workflow. Finally, trigger/run the workflow and monitor a few runs. In the report, document the test runs, and include any findings or observations that you may have from the runs.

## 2.1 Assumptions

1. These tasks run independently from other tasks or resources on a daily basis.

2. The system is equipped with the necessary dependencies and environment for the tasks to run effectively.

3. The data required for each task is available and accessible at the time of execution.

## 2.2 CLI Test

To ensure the robustness and reliability of the `process_web_log` DAG, each individual task was subjected to a rigorous testing procedure using the Airflow CLI. This facilitated a granular evaluation of task logic and functionality in isolation from the workflow. The CLI command used for the test was :

```
airflow tasks test <dag_id> <task_id> <execution_date>
```

The figures below illustrate the outcome of each task test, with a focus on confirming the expected behavior of the task logic and its successful interaction with the required data and resources.

### 2.2.1 Test of Scan for Log Task

The `scan_for_log` task was tested to validate its capability to locate and identify the log file. As depicted in Figure 2.1, the task executed successfully, confirming the presence of the log file.

FIGURE 2.1: test_scan_for_log : The successful detection of the log file necessary for further processing.

### 2.2.2 Test of Extract Data Task

Following the successful log scanning, the `extract_data` task was assessed. This task's responsibility is to extract the requisite data fields from the log file. Figure 2.2 showcases a successful extraction, laying the groundwork for the subsequent transformation step.



FIGURE 2.2: test_extract_data : Extraction of data fields verified, ensuring data integrity for the next phase.

### 2.2.3 Test of Transform Data Task

The `transform_data` task's purpose is to apply the necessary transformations to the extracted data. The test, as shown in Figure 2.3, confirmed that the task performed the transformations

correctly and the output was as intended.



FIGURE 2.3: test_transform_data : Validation of data transformation logic and output formatting.

### 2.2.4  Test of Load Data Task

Lastly, the load_data task was evaluated to ensure it could successfully load the transformed data to the target destination. The success of this task, as observed in Figure 2.4, signifies the readiness of the data for any downstream processes or storage solutions.



FIGURE 2.4: test_load_data : Confirmation of the loading mechanism's operational status and data readiness.

### 2.2.5 Triggering the Whole Workflow

Having each task confirmed for operational readiness, the entire workflow was triggered to observe the cohesive execution of the tasks in a sequence that mimics the production environment. Figure 2.5 reflects the successful execution of the entire workflow as a unified process, demonstrating the DAG's effectiveness in orchestrating the defined tasks.



FIGURE 2.5: trigger_whole_process : The culmination of individual task tests in the successful execution of the full workflow.

## 2.3 GUI Test

A comprehensive graphical user interface (GUI) testing was conducted to visually monitor the execution of the defined tasks within the `process_web_log` DAG. This involved observing and recording the run logs, execution times, and statuses of each task to ensure operational transparency and identify any potential bottlenecks or errors.

### 2.3.1 Overview of Task Execution

Figure 2-1 presents a consolidated view of the execution of the four tasks. It provides a snapshot of each task's running status, execution time, and duration, offering insights into the workflow's overall efficiency and task coordination.

### 2.3.2 Detailed Run Logs

The run logs for each task were meticulously documented to trace the execution flow and debug any issues that arose during the tasks' lifecycle. Figures 2-2 to 2-5 capture the detailed logs, providing a granular look into the internal operations of each task.

FIGURE 2.6: Consolidated results overview of the four tasks, illustrating their execution status and timing.



FIGURE 2.7: Detailed run log for Task 1, highlighting the step-by-step execution process.



FIGURE 2.8: Execution log for Task 2, displaying the task's interactions and data handling.

FIGURE 2.9: Log entries for Task 3, showing the transformation steps and outcomes.



FIGURE 2.10: Operational log for Task 4, detailing the data loading and finalization stages.

### 2.3.3 Workflow Scheduling and Results

Figure 2-6 depicts the scheduling sequence and execution results of the workflow comprising the four tasks. It chronicles the tasks' interdependencies and execution order, which are pivotal for understanding the DAG's orchestration and pinpointing the critical path of the data pipeline.



FIGURE 2.11: Schematic representation of the workflow execution, showcasing the scheduling and results of the four-task sequence.

## 2.4 Conclusion

Working with Airflow for the `process_web_log` DAG has been a real learning curve. Running it on Linux meant grappling with Docker, which was new to me. I spent hours troubleshooting why Airflow couldn't find my log files, until I realized they needed to be inside the Docker container, not on my local machine.

I also learned how Airflow relies on a database to keep track of everything. When I made changes in Docker and didn't update the database, it led to some confusing errors. It's clear now how important that database is for keeping the workflow running smoothly.

All this has given me a much better grip on how Airflow works, especially the bits about Docker and databases that you don't really think about at first. I'm looking forward to getting even more familiar with Airflow and building more complex and reliable data workflows down the line.

## Exercise 3 (3 points)

The objective of this exercise was to augment the existing `process_web_log` DAG with an additional task. This new task sends out a notification upon the successful execution of the workflow. The flexibility of the implementation meant we could choose from various messaging platforms. We opted to send the notification to a Slack channel, which is widely used for team communication and can easily integrate with Airflow through webhooks.

## 3.1 Implementation

To implement the notification feature, a Python function was written to post a message to a specified Slack channel using the channel's webhook URL. This function was then used as a callable in an Airflow PythonOperator task. Below is the code snippet that outlines the implementation of the Slack notification task :

Listing 3.1: Python code for sending a Slack notification.

```python
def send_slack_message(**context):
    webhook_url = "https://hooks.slack.com/services/..."
    message = "Workflow executed successfully on " + \
              datetime.now().strftime("%Y-%m-%d %H:%M:%S")
    data = {"text": message}
    response = requests.post(webhook_url, json=data)
    print("Message sent to Slack: Status code", response.
        status_code)

success_message_task = PythonOperator(
    task_id='send_success_message',
```

```
11      python_callable=send_slack_message,
12      provide_context=True,
13      dag=dag)
```

This task was placed at the end of the DAG to ensure that the message would only be sent after the successful completion of all preceding tasks.

## 3.2   Results and Evidence of Execution

Upon triggering the workflow, all tasks, including the new Slack notification task, executed as expected. The message delivered to the Slack channel served as both a confirmation of successful execution and a real-time update to the team monitoring the workflow. The screenshot below captures the message as it appeared in the Slack channel, providing evidence of the task's successful execution.



FIGURE 3.1: Slack message confirming successful DAG execution.

## 3.3   Conclusion

The addition of the Slack notification task to the `process_web_log` DAG represents a significant step toward operationalizing our workflow. It not only ensures immediate feedback on the status of the DAG but also enhances the transparency of the data pipeline's performance. This exercise has reinforced the importance of integrating communication tools within data workflows, paving the way for more sophisticated alerting mechanisms in the future.

# 4

Listing 4.1: DAG definition with default arguments.

```python
from airflow import DAG
from airflow.operators.python_operator import PythonOperator
from datetime import datetime, timedelta
import os
import tarfile
import requests

default_args = {
    'owner': 'airflow',
    'depends_on_past': False,
    'start_date': datetime(2023, 1, 1),
    'email_on_failure': False,
    'email_on_retry': False,
    'retries': 1,
    'retry_delay': timedelta(seconds=1),
}

dag = DAG('process_web_log',
          default_args=default_args,
          description='DAG for processing web log',
          schedule_interval='@daily',
          catchup=False,
          tags=['DSBW'])
```

```
24
25  # Update the path to the new location of your log file inside the
        Docker container
26  log_dir = '/usr/local/airflow'
27  log_file = f'{log_dir}/log.txt'
28  extracted_data_file = f'{log_dir}/extracted_data.txt'
29  transformed_data_file = f'{log_dir}/transformed_data.txt'
30  tar_file = f'{log_dir}/weblog.tar'
31
32  def scan_for_log(**context):
33      if os.path.isfile(log_file):
34          return log_file
35      else:
36          raise ValueError("log.txt not found")
37
38  def extract_data(**context):
39      task_instance = context['ti']
40      log_path = task_instance.xcom_pull(task_ids='scan_for_log')
41
42      with open(log_path, 'r') as file, open(extracted_data_file, '
          w') as out_file:
43          for line in file:
44              ip_address = line.split()[0]   # Assuming IP address
                  is the first element in the log line
45              out_file.write(ip_address + '\n')
46
47  def transform_data(**context):
48      with open(extracted_data_file, 'r') as file, open(
          transformed_data_file, 'w') as out_file:
49          for line in file:
50              if '198.46.149.143' not in line:
51                  out_file.write(line)
52
```

```python
53  def load_data(**context):
54      with tarfile.open(tar_file, 'w') as tar:
55          tar.add(transformed_data_file, arcname='transformed_data.
                txt')
56
57  def send_slack_message(**context):
58      webhook_url = "https://hooks.slack.com/services/T067ARZ6K18/
            B06747785P0/ixNxGWM33fovA0uIc0mf0FjV"
59      message = "Workflow executed successfully on " + datetime.now
            ().strftime("%Y-%m-%d %H:%M:%S")
60      data = {"text": message}
61      response = requests.post(webhook_url, json=data)
62      print("Message sent to Slack: Status code", response.
            status_code)
63
64
65  # Define tasks
66  scan_task = PythonOperator(
67      task_id='scan_for_log',
68      python_callable=scan_for_log,
69      provide_context=True,
70      dag=dag)
71
72  extract_task = PythonOperator(
73      task_id='extract_data',
74      python_callable=extract_data,
75      provide_context=True,
76      dag=dag)
77
78  transform_task = PythonOperator(
79      task_id='transform_data',
80      python_callable=transform_data,
81      provide_context=True,
```

```
82        dag=dag)
83
84  load_task = PythonOperator(
85        task_id='load_data',
86        python_callable=load_data,
87        provide_context=True,
88        dag=dag)
89
90  success_message_task = PythonOperator(
91        task_id='send_success_message',
92        python_callable=send_slack_message,
93        provide_context=True,
94        dag=dag)
95
96
97  # Set up the workflow
98  scan_task >> extract_task >> transform_task >> load_task >>
        success_message_task
```