



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

**Assignment – 3  
Date – 17<sup>th</sup> November, 2023**

**Submitted By**  
Sony Shrestha  
Aayush Paudel

## Question 1

### Definition

A Directed Acyclic Graph (DAG) is a graph structure that consists of nodes (vertices) connected by directed edges (arcs or arrows). The "directed" aspect implies that each edge has a specific direction, indicating a relationship between two nodes where one node is the source, and the other is the destination. "Acyclic" means that there are no cycles or loops in the graph, ensuring a clear and deterministic order of traversal.

### DAGs in Apache Airflow

In the realm of Apache Airflow, a DAG is used to define and orchestrate a workflow - a sequence of tasks that need to be executed in a particular order. Each task within a DAG is represented by an operator, which is a class that defines what the task does.

DAG has been prepared to analyze a web server log file, which extracts IP addresses from the given log file, performs transformation to filter given IP addresses and generates a tar file out of it. Finally, the user is notified about the successful execution of DAG via an email.

The details of the DAG created is provided below:

### DAG Details

**DAG ID:** process\_web\_log

**Schedule Interval:** daily

**Start Date:** datetime (2023,11,12)

### Tasks

Above DAG consists of following tasks.

#### Task 1: scan\_for\_log

**Task ID:** scan\_for\_log

#### Task Description:

1. This task checks whether a log file named **log.txt** is present in the input directory.
2. **BashOperator** has been used to check file availability.
3. Process continues if the file is found, otherwise it fails with an error indicating "Log File Not Found".

#### Dependencies:

1. None

#### Parameters:

1. Input Directory: This is the directory where our process will check for the log file.

## Task 2: extract\_data

Task ID: extract\_data

### Task Description:

1. This task extracts IP addresses from the input log file.
2. **PythonOperator**, which makes use of regular expressions to identify IP addresses, has been used.
3. **Extracted IP** addresses are then dumped into a file named **extracted\_data.txt**.

### Dependencies:

1. This task executes upon successful execution of **scan\_for\_log** (if **log.txt** file is present in input directory).

### Parameters:

1. Input Directory: This is the directory where the log file is present.
2. Output Directory: This is the directory where **extracted\_data.txt** file is generated by the task.

## Task 3: transform\_data

Task ID: transform\_data

### Task Description:

1. This task filters out all the occurrences of IP address **198.46.149.143** from the "extracted\_data.txt" file and saves the output to a file named "**transformed\_data.txt**".
2. **PythonOperator** is used for this purpose.

### Dependencies:

1. This task executes upon successful execution of **extract\_data** (if **extracted\_data.txt** file has been generated by the previous task).

### Parameters:

1. Input Directory: This is the directory where the **extracted\_data.txt** file is present.
2. Output Directory: This is the directory where **transformed\_data.txt** file is generated by the task.
3. IP address: This is the IP address which needs to be filtered out.

## Task 4: load\_data

Task ID: load\_data

**Task Description:**

1. This task archives the file “**transformed\_data.txt**” into a tar file named “**weblog.tar**”.
2. **BashOperator** has been used for this purpose.

**Dependencies:**

1. This task executes upon successful execution of transform\_data (if **transformed\_data.txt** file has been generated by the previous task).

**Parameters:**

1. Input Directory: This is the directory where the **transformed\_data.txt** file is present.
2. Output Directory: This is the directory where the **weblog.tar** file is generated by the task.

**Task 5:** send\_email

**Task ID:** send\_email

**Task Description:**

1. This task sends an email to the user, notifying him/her about the successful execution of the workflow.
2. **EmailOperator** has been used for this purpose.

**Dependencies:**

1. This task executes upon successful execution of **load\_data** (if weblog.tar file has been generated by the previous task).

**Parameters:**

1. Input Directory: This is the directory where the **weblog.tar** file is present.
2. Recipient Email: Email address of the user to be notified about the successful execution of workflow is provided.
3. Subject: Subject of the email
4. Content: Body of the email

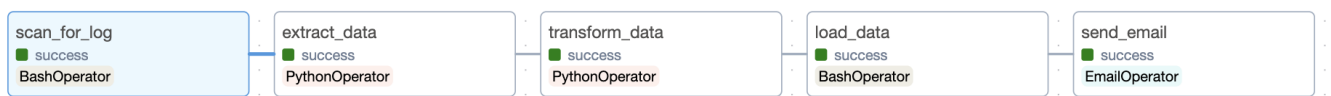
**Overall Dependencies:**

The tasks present in workflow are executed in the following sequence.

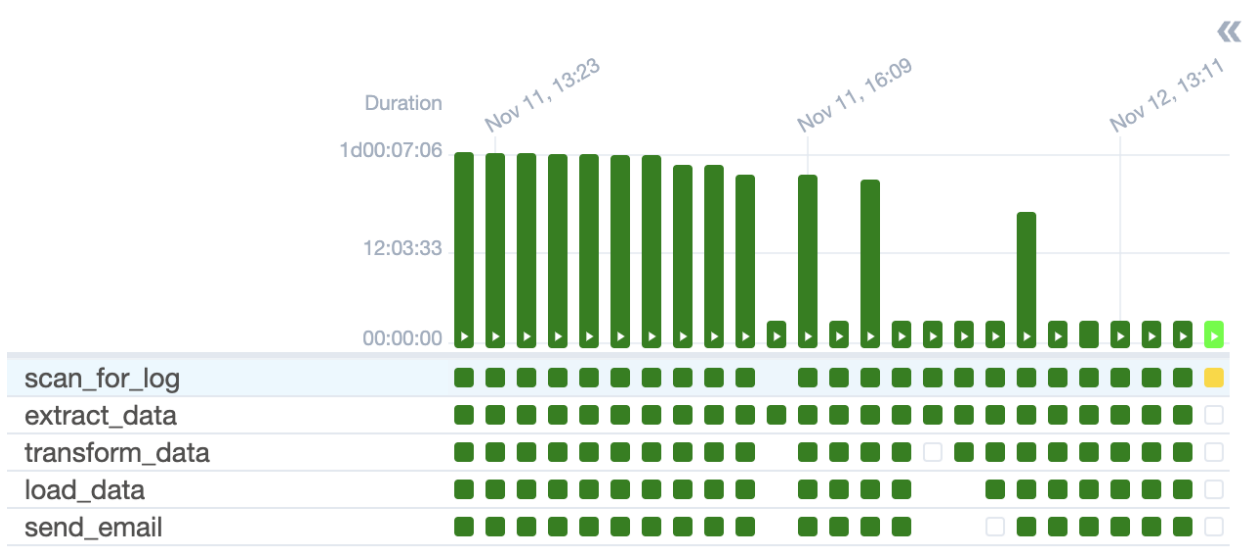
**scan\_for\_log\_task >> extract\_data\_task >> transform\_data\_task >> load\_data\_task >> email\_task**

## Question 2

The following image provides the sequential flow of the task.



Number of executions were done. Corresponding execution chart has been provided below.



On average, it took 9 seconds for entire workflow to be executed successfully. Details for each task instance has been provided in diagram below.

### Details for scan\_for\_log task which was executed successfully

### Task Instance Details

Status	<span>■</span> success
Task ID	scan_for_log <a href="#">🔗</a>
Run ID	manual__2023-11-12T13:59:29.943405+00:00 <a href="#">🔗</a>
Operator	BashOperator
Trigger Rule	all_success
Duration	00:00:00
Started	2023-11-12, 13:59:32 UTC

Details for extract\_data which was executed successfully

Task Instance Details

Status	<div><div></div> success</div>
Task ID	extract_data <a href="#">🔗</a>
Run ID	manual__2023-11-11T17:15:31.691286+00:00 <a href="#">🔗</a>
Operator	PythonOperator
Trigger Rule	all_success
Duration	00:00:00
Started	2023-11-11, 17:15:35 UTC

Details for transform\_data task which was executed successfully

Task Instance Details

Status	<div><div></div> success</div>
Task ID	transform_data <a href="#">🔗</a>
Run ID	manual__2023-11-11T17:15:31.691286+00:00 <a href="#">🔗</a>
Operator	PythonOperator
Trigger Rule	all_success
Duration	00:00:00
Started	2023-11-11, 17:15:36 UTC

Details for load\_data task which was executed successfully

Task Instance Details

Status	<div><div></div> success</div>
Task ID	load_data <a href="#">🔗</a>
Run ID	manual__2023-11-11T17:15:31.691286+00:00 <a href="#">🔗</a>
Operator	BashOperator
Trigger Rule	all_success
Duration	00:00:00
Started	2023-11-11, 17:15:38 UTC

Details for send\_mail task which was executed successfully

Task Instance Details

Status	<div><div></div> success</div>
Task ID	send_email <a href="#">🔗</a>
Run ID	manual__2023-11-12T13:59:29.943405+00:00 <a href="#">🔗</a>
Operator	EmailOperator
Trigger Rule	all_success
Duration	00:00:01
Started	2023-11-12, 13:59:38 UTC

Details for scan\_for\_log task which failed because of log file not found scenario

Task Instance Details

Status	<div><div></div>failed</div>
Task ID	scan_for_log <a href="#">🔗</a>
Run ID	manual__2023-11-12T14:03:39.909381+00:00 <a href="#">🔗</a>
Operator	BashOperator
Trigger Rule	all_success
Duration	00:00:00
Started	2023-11-12, 14:08:42 UTC

Details for scan\_for\_log task which was scanning for the input log file

Task Instance Details

Status	<div><div></div>up_for_retry</div>
Task ID	scan_for_log <a href="#">🔗</a>
Run ID	manual__2023-11-12T14:03:39.909381+00:00 <a href="#">🔗</a>
Operator	BashOperator
Trigger Rule	all_success
Duration	00:00:00
Started	2023-11-12, 14:03:41 UTC



### Question 3

Additional task has been created in airflow that is responsible for sending an email to the user, notifying him/her about the successful execution of the workflow.

Task Detail for **send\_email** has been provided below.

#### Task 5: send\_email

**Task ID:** send\_email

#### Task Description:

1. This task sends an email to the user, notifying him/her about the successful execution of the workflow.
2. **EmailOperator** has been used for this purpose.

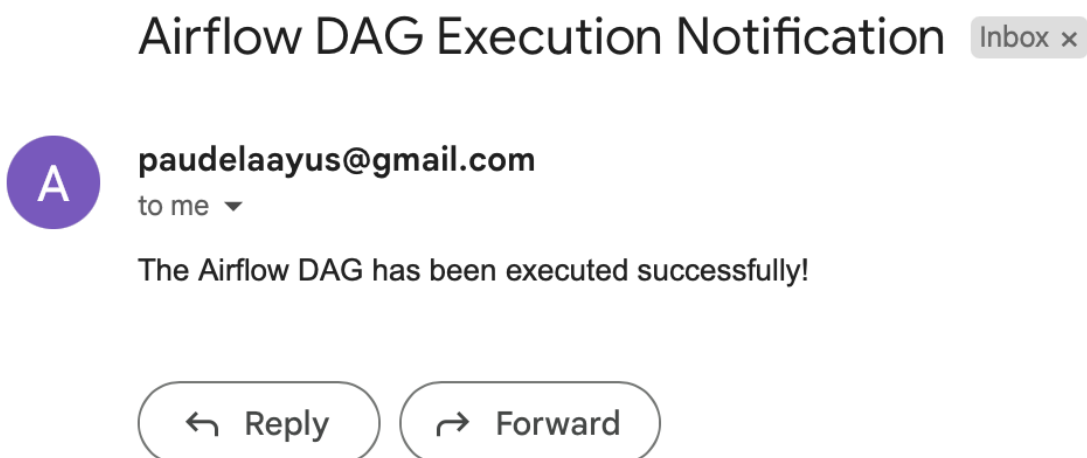
#### Dependencies:

1. This task executes upon successful execution of load\_data (if **weblog.tar** file has been generated by the previous task).

#### Parameters:

1. Input Directory: This is the directory where the **weblog.tar** file is present.
2. Recipient Email: Email address of the user to be notified about the successful execution of workflow is provided.
3. Subject: Subject of the email
4. Content: Body of the email

Email is received as follows.



In order to make sure that Email operator works as expected, airflow.cfg file is configured for SMTP as follows:

```
[smtp]
smtp_host = smtp.gmail.com
smtp_starttls = True
smtp_ssl = False
smtp_user = your_gmail_username@gmail.com
smtp_password = your_gmail_password
smtp_port = 587
smtp_mail_from = your_airflow_email@gmail.com
```