Create and execute a Shell script from Airflow

Create a Basic Shell Script



Foillowing creating a basic shell script and starting an instance of airflow through docker

Explore Anatomy of DAG

```
An Apache Airflow DAG is a python program. It consists of these logical blocks.

    Imports

    DAG Arguments

    DAG Definition

    Task Definitions

    Task Pipeline

A typical imports block looks like this.
         # import the libraries
  1
  2
        from datetime import timedelta
  4
         from airflow import DAG
  6
         from airflow.operators.bash_operator import BashOperator
         # This makes scheduling easy
  8
                                                                                              4
         from airflow.utils.dates import days_ago
  9
A typical DAG Arguments block looks like this.
   1
          #defining DAG arguments
   2
          # You can override them on a per-task basis during operator initialization
          default_args = {
              'owner': 'Ramesh Sannareddy',
              'start_date': days_ago(0),
   6
    7
              'email': ['ramesh@somemail.com'],
              'email_on_failure': True,
   8
   9
              'email_on_retry': True,
              'retries': 1,
  10
              'retry_delay': timedelta(minutes=5),
  11
                                                                                               @
  12
          }
DAG arguments are like settings for the DAG.
```

The above settings mention

- · the owner name,
- · when this DAG should run from: days_age(0) means today,
- · the email address where the alerts are sent to,
- · whether alert must be sent on failure,
- · whether alert must be sent on retry,
- · the number of retries in case of failure, and
- · the time delay between retries.

A typical DAG definition block looks like this.

```
# define the DAG
dag = DAG(
dag_id='sample-etl-dag',
default_args=default_args,
description='Sample ETL DAG using Bash',
schedule_interval=timedelta(days=1),
)
```

Here we are creating a variable named dag by instantiating the DAG class with the following parameters.

sample-etl-dag is the ID of the DAG. This is what you see on the web console.

We are passing the dictionary default_args, in which all the defaults are defined.

description helps us in understanding what this DAG does.

schedule_interval tells us how frequently this DAG runs. In this case every day. (days=1).

A typical task definitions block looks like this:

```
# define the tasks

# define the task named extract_transform_and_load to call the shell script

extract_transform_and_load = BashOperator(

task_id='extract_transform_and_load',

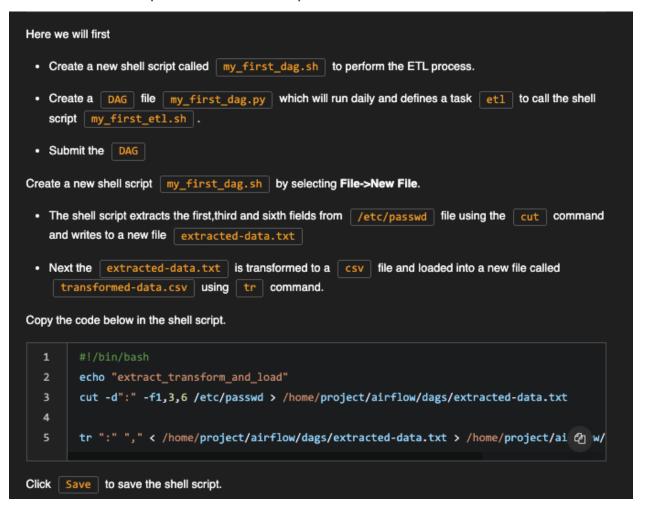
bash_command='/home/project/airflow/dags/extract_transform_load.sh "',

dag=dag,

)
```



Exercise 4 - ETL process on a /etc/passwd file

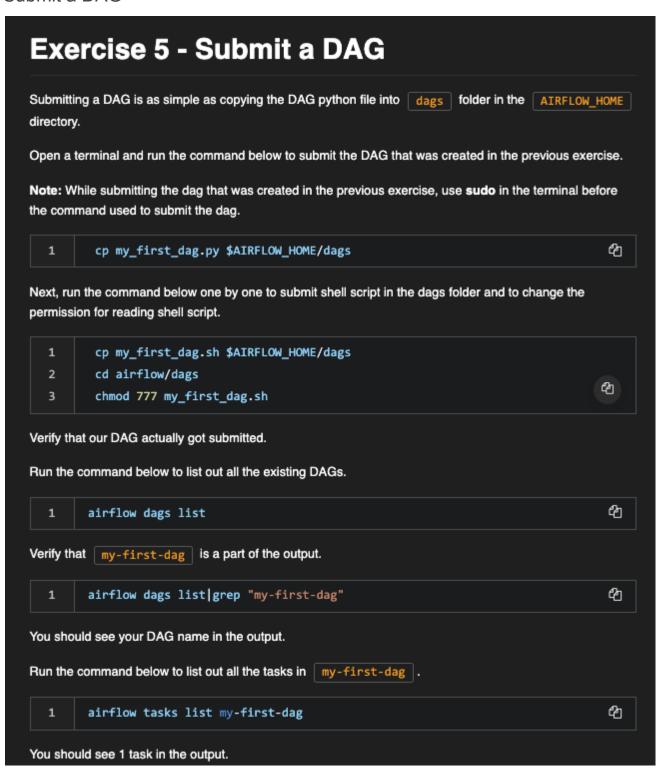


```
Create a new DAG file my_first_dag.py by selecting File->New File.

This DAG has one task etl that calls the shell script my_first_dag.sh
```

```
1
      # import the libraries
 2
 3
      from datetime import timedelta
 4
      from airflow import DAG
 5
      # Operators; we need this to write tasks!
 6
 7
      from airflow.operators.bash_operator import BashOperator
      # This makes scheduling easy
 8
 9
      from airflow.utils.dates import days_ago
10
      #defining DAG arguments
11
12
      # You can override them on a per-task basis during operator initialization
13
14
      default_args = {
           'owner': 'Ramesh Sannareddy',
15
           'start_date': days_ago(0),
16
17
           'email': ['ramesh@somemail.com'],
           'email_on_failure': False,
18
          'email_on_retry': False,
19
          'retries': 1,
20
21
           'retry_delay': timedelta(minutes=5),
22
      }
23
      # defining the DAG
24
25
```

```
# define the DAG
26
      dag = DAG(
27
          'my-first-dag',
28
          default_args=default_args,
29
30
          description='My first DAG',
          schedule_interval=timedelta(days=1),
31
      )
32
33
      # define the task **extract_transform_and_load** to call shell script
34
35
      #calling the shell script
36
      extract_transform_load = BashOperator(
37
          task_id="extract_transform_load",
38
          bash_command="/home/project/airflow/dags/my_first_dag.sh ",
39
          dag=dag,
40
41
      )
42
43
      # task pipeline
                                                                                     @
44
      extract_transform_load
```



```
theia@theiadocker-craigtrupp8:/home/project$ cp my first dag.py
$AIRFLOW_HOME/dags
theia@theiadocker-craigtrupp8:/home/project$ cp my_first_dag_1.py
cp: missing destination file operand after 'my first dag 1.py'
Try 'cp --help' for more information.
theia@theiadocker-craigtrupp8:/home/project$ cp my_first_dag_1.py
$AIRFLOW HOME/dags
theia@theiadocker-craigtrupp8:/home/project$ cp my first dag.sh
$AIRFLOW HOME/dags
theia@theiadocker-craigtrupp8:/home/project$ cd airflow/dags
theia@theiadocker-craigtrupp8:/home/project/airflow/dags$ chmod 777
my first dag.sh
theia@theiadocker-craigtrupp8:/home/project/airflow/dags$ airflow dags list
/home/airflow/.local/lib/python3.7/site-packages/airflow/configuration.py:5
28: DeprecationWarning: The sql_alchemy_conn option in [core] has been
moved to the sql_alchemy_conn option in [database] - the old setting has
been used, but please update your config.
 option = self._get_environment_variables(deprecated_key,
deprecated_section, key, section)
/home/airflow/.local/lib/python3.7/site-packages/airflow/configuration.py:5
28: DeprecationWarning: The auth_backend option in [api] has been renamed
to auth_backends - the old setting has been used, but please update your
config.
 option = self. get environment variables(deprecated key,
deprecated_section, key, section)
/home/airflow/.local/lib/python3.7/site-packages/airflow/configuration.py:3
60: FutureWarning: The auth backends setting in [api] has had
airflow.api.auth.backend.session added in the running config, which is
needed by the UI. Please update your config before Apache Airflow 3.0.
  FutureWarning,
Error: Failed to load all files. For details, run `airflow dags
list-import-errors`
                                        filepath
dag id
owner
               paused
=========++=====++=====++=====++=====+
example bash operator
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
xample_dags/example_bash_operator.py
example_branch_datetime_operator
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
```

```
True
xample_dags/example_branch_datetime_operator.py
example branch datetime operator 2
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
True
xample dags/example branch datetime operator.py
example branch dop operator v3
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
True
xample_dags/example_branch_python_dop_operator_3.py
example_branch_labels
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
True
xample_dags/example_branch_labels.py
example_branch_operator
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
True
xample_dags/example_branch_operator.py
example branch python operator decorator
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
True
xample_dags/example_branch_operator_decorator.py
example complex
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
True
                                          xample_dags/example_complex.py
example_dag_decorator
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
True
xample_dags/example_dag_decorator.py
example_external_task_marker_child
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
True
xample_dags/example_external_task_marker_dag.py
```

```
example external task marker parent
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
True
xample_dags/example_external_task_marker_dag.py
example_kubernetes executor
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
True
xample_dags/example_kubernetes_executor.py
example local kubernetes executor
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
True
xample_dags/example_local_kubernetes_executor.py
example nested branch dag
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
True
xample dags/example nested branch dag.py
example passing params via test command
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
True
xample_dags/example_passing_params_via_test_command.py
example_python_operator
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
True
xample_dags/example_python_operator.py
example_short_circuit_operator
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
True
xample_dags/example_short_circuit_operator.py
example skip dag
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
True
                                          xample_dags/example_skip_dag.py
example_sla_dag
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
```

```
xample dags/example sla dag.py
example_subdag_operator
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
True
xample_dags/example_subdag_operator.py
example subdag operator.section-1
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
True
xample_dags/example_subdag_operator.py
example_subdag_operator.section-2
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
True
xample dags/example subdag operator.py
example_task_group
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
True
xample_dags/example_task_group.py
example_task_group_decorator
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
True
xample dags/example task group decorator.py
example_time_delta_sensor_async
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
True
xample_dags/example_time_delta_sensor_async.py
example trigger controller dag
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
True
xample dags/example trigger controller dag.py
example_trigger_target_dag
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
True
xample_dags/example_trigger_target_dag.py
example_weekday_branch_operator
```

```
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
True
xample dags/example branch day of week operator.py
example xcom
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
True
                                          xample dags/example xcom.py
example_xcom_args
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
True
                                          xample_dags/example_xcomargs.py
example_xcom_args_with_operators
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
True
                                          xample_dags/example_xcomargs.py
latest only
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
True
xample dags/example latest only.py
latest_only_with_trigger
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
True
xample dags/example_latest_only_with_trigger.py
my-first-dag
                                          my_first_dag_1.py
| Darth Vader | True
tutorial
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
True
                                          xample_dags/tutorial.py
tutorial etl dag
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
True
                                          xample dags/tutorial etl dag.py
tutorial_taskflow_api_etl
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
```

```
True
xample dags/tutorial taskflow api etl.py
tutorial taskflow api etl virtualenv
/home/airflow/.local/lib/python3.7/site-packages/airflow/e | airflow
True
xample dags/tutorial taskflow api etl virtualenv.py
theia@theiadocker-craigtrupp8:/home/project/airflow/dags$ airflow dags list
grep 'my-first-dag-1'
theia@theiadocker-craigtrupp8:/home/project/airflow/dags$ airflow dags list
grep 'my-first-dag'
                                         | my_first_dag_1.py
my-first-dag
| Darth Vader | True
theia@theiadocker-craigtrupp8:/home/project/airflow/dags$ airflow tasks
lits my-first-dag
/home/airflow/.local/lib/python3.7/site-packages/airflow/configuration.py:5
28: DeprecationWarning: The sql_alchemy_conn option in [core] has been
moved to the sql alchemy conn option in [database] - the old setting has
been used, but please update your config.
  option = self._get_environment_variables(deprecated_key,
deprecated section, key, section)
/home/airflow/.local/lib/python3.7/site-packages/airflow/configuration.py:5
28: DeprecationWarning: The auth_backend option in [api] has been renamed
to auth_backends - the old setting has been used, but please update your
config.
 option = self. get environment variables(deprecated key,
deprecated_section, key, section)
/home/airflow/.local/lib/python3.7/site-packages/airflow/configuration.py:3
60: FutureWarning: The auth_backends setting in [api] has had
airflow.api.auth.backend.session added in the running config, which is
needed by the UI. Please update your config before Apache Airflow 3.0.
  FutureWarning,
usage: airflow tasks [-h] COMMAND ...
Manage tasks
positional arguments:
  COMMAND
                      Clear a set of task instance, as if they never ran
    clear
                      Returns the unmet dependencies for a task instance
    failed-deps
```

```
list
                      List the tasks within a DAG
                      Render a task instance's template(s)
    render
                      Run a single task instance
    run
                      Get the status of a task instance
    state
    states-for-dag-run
                      Get the status of all task instances in a dag run
                     Test a task instance
    test
optional arguments:
  -h, --help
                      show this help message and exit
airflow tasks command error: argument COMMAND: invalid choice: 'lits'
(choose from 'clear', 'failed-deps', 'list', 'render', 'run', 'state',
'states-for-dag-run', 'test'), see help above.
theia@theiadocker-craigtrupp8:/home/project/airflow/dags$ airflow tasks
list
/home/airflow/.local/lib/python3.7/site-packages/airflow/configuration.py:5
28: DeprecationWarning: The sql_alchemy_conn option in [core] has been
moved to the sql_alchemy_conn option in [database] - the old setting has
been used, but please update your config.
  option = self. get environment variables(deprecated key,
deprecated_section, key, section)
/home/airflow/.local/lib/python3.7/site-packages/airflow/configuration.py:5
28: DeprecationWarning: The auth backend option in [api] has been renamed
to auth_backends - the old setting has been used, but please update your
config.
  option = self._get_environment_variables(deprecated_key,
deprecated_section, key, section)
/home/airflow/.local/lib/python3.7/site-packages/airflow/configuration.py:3
60: FutureWarning: The auth_backends setting in [api] has had
airflow.api.auth.backend.session added in the running config, which is
needed by the UI. Please update your config before Apache Airflow 3.0.
 FutureWarning,
usage: airflow tasks list [-h] [-S SUBDIR] [-t] [-v] dag_id
List the tasks within a DAG
positional arguments:
  dag_id
                       The id of the dag
optional arguments:
  -h, --help
                        show this help message and exit
  -S SUBDIR, --subdir SUBDIR
```

```
File location or directory from which to look for
the dag. Defaults to '[AIRFLOW_HOME]/dags' where [AIRFLOW_HOME] is the
value you set for 'AIRFLOW_HOME' config you set in 'airflow.cfg'
  -t, --tree
                        Tree view
                        Make logging output more verbose
  -v, --verbose
airflow tasks list command error: the following arguments are required:
dag id, see help above.
theia@theiadocker-craigtrupp8:/home/project/airflow/dags$ airflow tasks
list my-first-dag
/home/airflow/.local/lib/python3.7/site-packages/airflow/configuration.py:5
28: DeprecationWarning: The sql_alchemy_conn option in [core] has been
moved to the sql_alchemy_conn option in [database] - the old setting has
been used, but please update your config.
  option = self._get_environment_variables(deprecated_key,
deprecated_section, key, section)
/home/airflow/.local/lib/python3.7/site-packages/airflow/configuration.py:5
28: DeprecationWarning: The auth_backend option in [api] has been renamed
to auth_backends - the old setting has been used, but please update your
config.
  option = self. get environment variables(deprecated key,
deprecated_section, key, section)
/home/airflow/.local/lib/python3.7/site-packages/airflow/configuration.py:3
60: FutureWarning: The auth backends setting in [api] has had
airflow.api.auth.backend.session added in the running config, which is
needed by the UI. Please update your config before Apache Airflow 3.0.
  FutureWarning,
extract_transform_load
theia@theiadocker-craigtrupp8:/home/project/airflow/dags$
```