

DATA 226- DATAWAREHOUSE

Homework 5

Name : Vimalanandhan Sivanandham

SJSU ID: 017596436

Porting homework #4 to Airflow (13 pts)

- (+2) Create tasks using @task decorator (refer to [GitHub link](#)[Links to an external site.](#))
 - You can use as many tasks as you want
 - Schedule the tasks properly (task dependency)

```
from airflow import DAG
from airflow.models import Variable
from airflow.decorators import task
from airflow.providers.snowflake.hooks.snowflake import SnowflakeHook
from datetime import datetime
import requests
import logging

def return_snowflake_conn():
    """Initialize Snowflake connection using SnowflakeHook."""
    # Initialize the SnowflakeHook using the connection ID stored in Airflow
    hook = SnowflakeHook(snowflake_conn_id='snowflake_conn')

    # Return the cursor object
    return hook.get_conn().cursor()

@task
def extract():
    """Extract AAPL stock data from Alpha Vantage API"""
    api_key = Variable.get("ALPHA_VANTAGE_API_KEY") # Get API key from Airflow Variables
    symbol = "AAPL"
    url = f"https://www.alphavantage.co/query?function=TIME_SERIES_DAILY&symbol={symbol}&apikey={api_key}&outputsize=compact"

    try:
        response = requests.get(url)
        response.raise_for_status()
        data = response.json().get("Time Series (Daily)", {})
        logging.info(f"Extracted {len(data)} records")
        return data # XCom push
    except Exception as e:
        logging.error(f"Error in extract: {str(e)}")
        raise
```

```

@task
def transform(data):
    """Transform extracted stock data into structured format"""
    records = []
    for date, values in data.items():
        records.append([
            date, float(values["1. open"]), float(values["2. high"]),
            float(values["3. low"]), float(values["4. close"]),
            int(values["5. volume"])
        ])
    logging.info(f"Transformed {len(records)} records")
    return records # XCom push

@task
def load(records):
    """Load transformed data into Snowflake using SnowflakeHook"""
    cur = return_snowflake_conn()
    try:
        cur.execute("BEGIN;")
        cur.execute("""
            CREATE TABLE IF NOT EXISTS stock_data.raw.stock_data (
                date DATE PRIMARY KEY,
                open FLOAT,
                high FLOAT,
                low FLOAT,
                close FLOAT,
                volume INT
            );
        """)
        cur.execute("DELETE FROM stock_data.raw.stock_data;") # Full refresh

        for record in records:
            # Debug: Print each record before inserting
            print(f"Inserting record: {record}")

```

```

def load(records):
    for record in records:
        # Debug: Print each record before inserting
        print(f"Inserting record: {record}")

        # Check for NULL values and replace them
        cleaned_record = [
            value if value is not None else 0 # Replace NULLs with 0
            for value in record
        ]

        sql = f"""
            INSERT INTO stock_data.raw.stock_data (date, open, high, low, close, volume)
            VALUES ('{cleaned_record[0]}', {cleaned_record[1]}, {cleaned_record[2]}, {cleaned_record[3]}, {cleaned_record[4]}, {cleaned_r
            """
        cur.execute(sql)

        cur.execute("COMMIT;")
        logging.info("Data successfully loaded into Snowflake")
    except Exception as e:
        cur.execute("ROLLBACK;")
        logging.error(f"Error in load: {str(e)}")
        raise

# Define the Airflow DAG
with DAG(
    dag_id='AAPL_Stock_ETL',
    start_date=datetime(2024, 9, 21),
    catchup=False,
    tags=['ETL', 'Stock Data'],
    schedule_interval='30 2 * * *' # Run daily at 2:30 AM UTC
) as dag:
    data = extract()
    transformed_data = transform(data)
    load(transformed_data)

```

- (+1) Set up a variable for Alpha Vantage API key
 - Use the variable in your code (Variable.get)
 - Capture the Admin -> Variables screenshot (an example will be provided ①)

```
from airflow import DAG
from airflow.models import Variable
from airflow.decorators import task
from airflow.providers.snowflake.hooks.snowflake import SnowflakeHook
from datetime import datetime
import requests
import logging

def return_snowflake_conn():
    """Initialize Snowflake connection using SnowflakeHook."""

    # Initialize the SnowflakeHook using the connection ID stored in Airflow
    hook = SnowflakeHook(snowflake_conn_id='snowflake_conn')

    # Return the cursor object
    return hook.get_conn().cursor()

@task
def extract():
    """Extract AAPL stock data from Alpha Vantage API"""
    api_key = Variable.get("ALPHA_VANTAGE_API_KEY") # Get API key from Airflow Variables
    symbol = "AAPL"
    url = f"https://www.alphavantage.co/query?function=TIME_SERIES_DAILY&symbol={symbol}&apikey={api_key}&outputsize=compact"

    try:
        response = requests.get(url)
        response.raise_for_status()
        data = response.json().get("Time Series (Daily)", {})
        logging.info(f"Extracted {len(data)} records")
        return data # XCom push
    except Exception as e:
        logging.error(f"Error in extract: {str(e)}")
        raise
```

Added Row

Choose File No file chosen ☒ Overwrite if exists ☐ Fail if exists ☐ Skip if exists

List Variable

Search-

+ Actions -

Record Count: 5

	Key	Val	Description	Is Encrypted
<input type="checkbox"/>	ALPHA_VANTAGE_API_KEY	*****		False
<input type="checkbox"/>	ALPHA_VANTAGE_SYMBOL	AAPL		False
<input type="checkbox"/>	SNOWFLAKE_ACCOUNT	10657309		False
<input type="checkbox"/>	SNOWFLAKE_PASSWORD	*****		False
<input type="checkbox"/>	SNOWFLAKE_USER	vimalnandham		False

Version: v2.10.1
 OR Version: _release:854173176f272660980ed446286c32cb75945e

- (+2) Set up Snowflake Connection (refer to [GitHub link](#)[Links to an external site.](#))
 - Use the connection in your code
 - Capture the Connection detail page screenshot (an example will be provided ②)

```
from airflow import DAG
from airflow.models import Variable
from airflow.decorators import task
from airflow.providers.snowflake.hooks.snowflake import SnowflakeHook
from datetime import datetime
import requests
import logging

def return_snowflake_conn():
    """Initialize Snowflake connection using SnowflakeHook."""

    # Initialize the SnowflakeHook using the connection ID stored in Airflow
    hook = SnowflakeHook(snowflake_conn_id='snowflake_conn')

    # Return the cursor object
    return hook.get_conn().cursor()

@task
def extract():
    """Extract AAPL stock data from Alpha Vantage API"""
    api_key = Variable.get("ALPHA_VANTAGE_API_KEY") # Get API key from Airflow Variables
    symbol = "AAPL"
    url = f"https://www.alphavantage.co/query?function=TIME_SERIES_DAILY&symbol={symbol}&apikey={api_key}&outputsize=compact"

    try:
        response = requests.get(url)
        response.raise_for_status()
        data = response.json().get("Time Series (Daily)", {})
        logging.info(f"Extracted {len(data)} records")
        return data # XCom push
    except Exception as e:
        logging.error(f"Error in extract: {str(e)}")
        raise
```

Connection successfully tested

21:49 UTC

Edit Connection

Connection id * snowflake_conn

Connection Type * Snowflake
Connection Type missing? Make sure you've installed the corresponding Airflow Provider Package.

Description

Schema snowflake schema

Login vrm@stanandham

Password

Extra

```
{
  "account": "usdt500",
  "warehouse": "compute-wh",
  "database": "dev",
  "schema": "dev",
  "role": "snowflake_role",
  "private_key_path": null
}
```

Account usdt500

Warehouse compute-wh

Database dev

Region snowflake hosted region

Role snowflake role

Private key (Path) Path of snowflake private key (PEM Format)

Private key (Text)

- (+5) Ensure the overall DAG is implemented properly and runs successfully
 - A github link with the entire code needs to be submitted (2 pts)
 - Implement the same full refresh using SQL transaction (3 pts)

<https://github.com/Vimalanandhan/DATA-226---DATAWAREHOUSE/tree/main/Homework/Homework5>

```
@task
def extract():
    """Extract AAPL stock data from Alpha Vantage API"""
    api_key = Variable.get("ALPHA_VANTAGE_API_KEY") # Get API key from Airflow Variables
    symbol = "AAPL"
    url = f"https://www.alphavantage.co/query?function=TIME_SERIES_DAILY&symbol={symbol}&apikey={api_key}&outputsize=compact"

    try:
        response = requests.get(url)
        response.raise_for_status()
        data = response.json().get("Time Series (Daily)", {})
        logging.info(f"Extracted {len(data)} records")
        return data # XCom push
    except Exception as e:
        logging.error(f"Error in extract: {str(e)}")
        raise

@task
def transform(data):
    """Transform extracted stock data into structured format"""
    records = []
    for date, values in data.items():
        records.append([
            date, float(values["1. open"]), float(values["2. high"]),
            float(values["3. low"]), float(values["4. close"]),
            int(values["5. volume"])
        ])
    logging.info(f"Transformed {len(records)} records")
    return records # XCom push

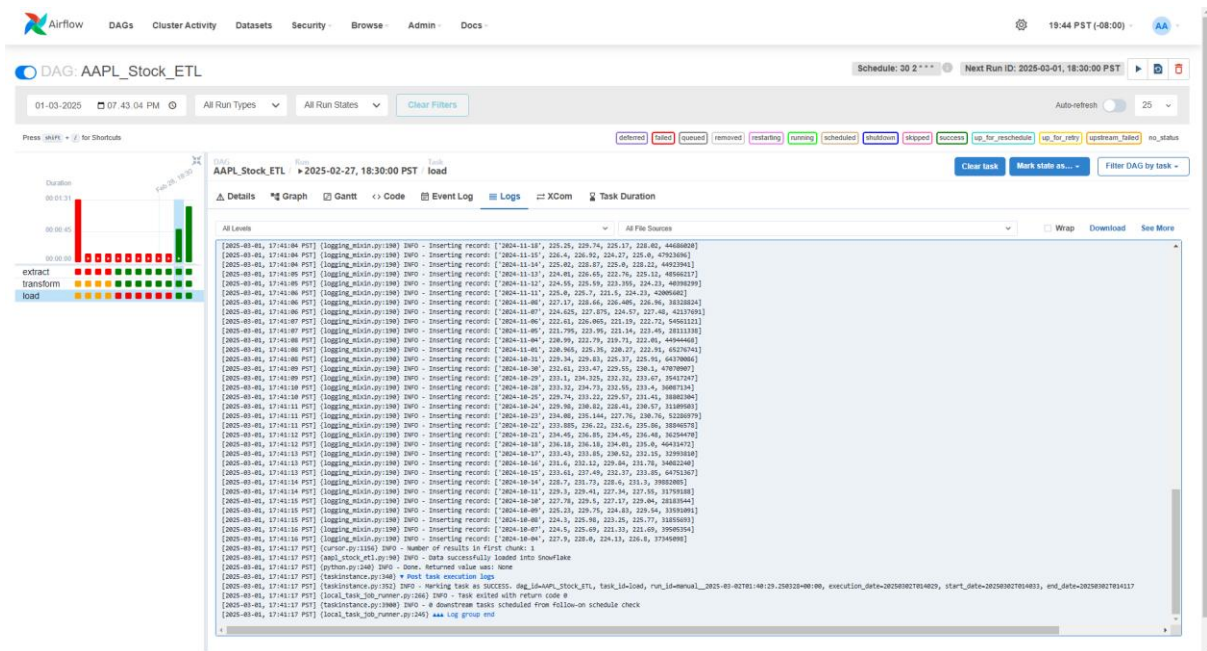
@task
def load(records):
    """Load transformed data into Snowflake using SnowflakeHook"""
    cur = return_snowflake_conn()
    try:
        cur.execute("BEGIN;")
        cur.execute("""
            CREATE TABLE IF NOT EXISTS stock_data.raw.stock_data (
                date DATE PRIMARY KEY,
                open FLOAT,
                high FLOAT,
                low FLOAT,
                close FLOAT,
                volume INT
            );
        """)
        cur.execute("DELETE FROM stock_data.raw.stock_data;") # Full refresh

        for record in records:
            # Debug: Print each record before inserting
            print(f"Inserting record: {record}")
```

- (+2) Capture two screenshot of your Airflow Web UI (examples to follow)
 - One with the Airflow homepage showing the DAG (③)
 - The other with the log screen of the DAG (④)

This screenshot shows the Airflow DAGs homepage. At the top, there are navigation links for DAGs, Cluster Activity, Datasets, Security, Browse, Admin, and Docs. The main section displays a table of DAGs with columns for DAG name, Owner, Runs, Schedule, Last Run, Next Run, Recent Tasks, and Actions. Two DAGs are visible: 'AAPL_Stock_ETL' (owned by 'airflow') and 'HelloWorld' (owned by 'user'). The 'AAPL_Stock_ETL' DAG has a status of 'Running' and a schedule of '30 2 * * *'. The 'HelloWorld' DAG has a status of 'Failed' and a schedule of '0 2 * * *'. Below the table, there is a pagination bar showing 'Showing 1-2 of 2 DAGs'.

This screenshot shows the Airflow DAG log screen for the 'AAPL_Stock_ETL' DAG. The top section displays the DAG name, schedule, and next run. Below this, there is a 'Press ctrl + / for shortcuts' bar and a 'Task Duration' chart. The main section shows the task logs for the 'extract' task, which is a 'DagRun' task. The logs show a series of 'Inserting record' messages for various stock prices. The task is currently in a 'Running' state. The bottom section shows the task logs for the 'transform' task, which is a 'DagRun' task. The logs show a series of 'Inserting record' messages for various stock prices. The task is currently in a 'Running' state.



- (+1) Overall formatting

4 screenshot examples are in the lecture notes (from slides 62 to 64)

