# EECS 6893: Big Data Analytics
# HW4
# Yutao Zhou UNI: yz4359

**Task 1 Helloworld**

**Q1.1**

**(1)**

**(2)**



Airflow — DAGs | Security | Browse | Admin | Docs — 21:12 UTC — YZ

Do not use **SQLite** as metadata DB in production – it should only be used for dev/testing We recommend using Postgres or MySQL. **Click here** for more information.

Do not use **SequentialExecutor** in production. **Click here** for more information.

## DAGs

| | DAG | Owner | Runs | Schedule | Last Run | Next Run | Recent Tasks |
|---|---|---|---|---|---|---|---|
| | example_bash_operator example example2 | airflow | ○○○○○ | 0 0 * * * | | 2022-11-21, 00:00:00 | ○○○○○○○○○○○ |
| | example_branch_datetime_operator_2 example | airflow | ○○○○○ | @daily | | 2022-11-21, 00:00:00 | ○○○○○○○○○○○ |
| | example_branch_dop_operator_v3 example | airflow | ○○○○○ | */1 * * * * | | 2022-11-22, 21:10:00 | ○○○○○○○○○○○ |
| | example_branch_labels | airflow | ○○○○○ | @daily | | 2022-11-21, 00:00:00 | ○○○○○○○○○○○ |

All **33**  Active **0**  Paused **33**

# Q1.2
# (1)SequentialExecutor



Airflow — DAGs | Security | Browse | Admin | Docs — 22:28 UTC — YZ

DAG: **helloworld** A simple toy DAG

Schedule: 1 day, 0:00:00  Next Run: 2022-11-22, 21:21:36

📍 Tree | 🔳 Graph | 📅 Calendar | ⏱ Task Duration | 🔀 Task Tries | 📈 Landing Times | ☰ Gantt | ⚠ Details | <> Code

2022-11-22T22:23:33Z  Runs  25  Update

○ BashOperator  ○ PythonOperator    ■ queued ■ running ■ success ■ failed ■ up_for_retry ■ up_for_reschedule ■ upstream_failed ■ skipped ■ scheduled ■ deferred □ no_status

Auto-refresh

DAG: **helloworld** A simple toy DAG

🟣 Tree    🔲 Graph    📅 Calendar    ⏳ Task Duration    ⇄ Task Tries    🛬 Landing Times    ≡ Gantt    ⚠ Details    <> Code

| 📅 | 2022-11-22T22:23:34Z | Runs | 25 ⌄ | Run | manual__2022-11-22T22:23:33.963878+00:00 ⌄ | Layout | Left > Right ⌄ | Update |

BashOperator    PythonOperator

Schedule: 1 day, 0:00:00    Next Run: 2022-11-22, 21:21:36

DAG: **helloworld** A simple toy DAG

🟣 Tree    🔲 Graph    📅 Calendar    ⏳ Task Duration    ⇄ Task Tries    🛬 Landing Times    ≡ Gantt    ⚠ Details    <> Code    ▶ 🗑

| 📅 | 2022-11-22T22:23:34Z | Runs | 25 ⌄ | Run | manual__2022-11-22T22:23:33.963878+00:00 ⌄ | Update |

**(1)LocalExecutor**

DAG: **helloworld** A simple toy DAG                    Schedule: 1 day, 0:00:00    Next Run: 2022-11-23, 22:45:36

💡 Tree    📊 Graph    📅 Calendar    ⏳ Task Duration    ⇄ Task Tries    📐 Landing Times    ☰ Gantt    ⚠ Details    <> Code    ▶    🗑

📅 | 2022-11-21T22:45:36Z    Runs    25 ⌄    Update

○ BashOperator   ○ PythonOperator          ■ queued  ■ running  ■ success  ■ failed  ■ up_for_retry  ■ up_for_reschedule  ■ upstream_failed  ■ skipped  ■ scheduled  ■ deferred  ☐ no_status

Auto-refresh    ⟳

Nov 21, 17:45

○ [DAG]
○ t1
  ○ t2_2
    ○ t3_2
      ○ t4_1
  ○ t2_1
    ○ t3_1
      ○ t4_1
  ○ t2_3
    ○ t4_1

DAG: **helloworld** A simple toy DAG            success    Schedule: 1 day, 0:00:00    Next Run: 2022-11-22, 22:45:36

💡 Tree    📊 Graph    📅 Calendar    ⏳ Task Duration    ⇄ Task Tries    📐 Landing Times    ☰ Gantt    ⚠ Details    <> Code    ▶    🗑

📅 | 2022-11-21T22:45:37Z    Runs    25 ⌄    Run    scheduled__2022-11-21T22:45:36.163039+00:00 ⌄    Layout    Left > Right ⌄    Find Task...

Update

BashOperator   PythonOperator          queued  running  success  failed  up_for_retry  up_for_reschedule  upstream_failed  skipped  scheduled  deferred  no_status
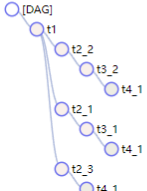
Auto-refresh    ⟳

**(2)**

**Calender:** This is a great features because it allows me to see the history states of job in one place. If we have a large project and fatch data for many days, this function make it clear on which day have successfully executed and which day fail. So, we could use thai data to monitor the pipeline.

**Task Duration:** This is a great features because it shows the duration of each sub task in our tree or graph. So, we will know sipecifically which task cost what amount of time. This is great for debugging. For example, when we accidentally wrote a infinite loop or wrtote some bugs that takes a really long time. This would make it clear to identify which task is wrong. Also, this could be used for optimize the runtime.

# Task 2 Build workflows

# Q2.1

# (1)

## DAG: Q2.1 Task 2 impliment the DAG below

🌲 Tree    🔧 Graph    📅 Calendar    ⏳ Task Duration    ⇄ Task Tries

📅   2022-11-23T00:03:47Z    Runs   25 ▾   Update

⚪ BashOperator   ⚪ PythonOperator

Nov 21, 19:03

```
○ [DAG]
 ○ t1
  ○ t2
   ○ t6
  ○ t5
   ○ t8
    ○ t15
     ○ t18
      ○ t19
    ○ t10
     ○ t14
      ○ t17
       ● t18
      ○ t16
       ○ t19
  ○ t9
   ○ t12
    ● t14
   ○ t11
    ● t14
 ○ t3
  ○ t7
   ● t18
   ● t14
   ○ t13
    ● t18
  ● t12
 ○ t4
```

◐ DAG: **Q2.1** Task 2 impliment the DAG below

📍 Tree   🔳 Graph   📅 Calendar   ⧖ Task Duration   ⇄ Task Tries   ⤓ Landing Times   ☰ Gantt   ⚠ Details   <> Code

| 📅 | 2022-11-23T00:03:48Z | Runs | 25 ⌄ | Run | manual__2022-11-23T00:03:47.602492+00:00 ⌄ | Layout | Left > Right ⌄ | Update |

BashOperator  PythonOperator



**(2)**

📍 Tree   🔳 Graph   📅 Calendar   ⧖ Task Duration   ⇄ Task Tries   ⤓ Landing Times   ☰ Gantt   ⚠ Details   <> Code

| 📅 | 2022-11-23T00:03:48Z | Runs | 25 ⌄ | Run | manual__2022-11-23T00:03:47.602492+00:00 ⌄ | Update |

**(3)**

| | | State | Dag Id | Execution Date | Run Id | Run Type | Queued At | Start Date | End Date | External T |
|---|---|---|---|---|---|---|---|---|---|---|
| ☐ | ✎ 🗑 | success | Q2.1 | 2022-11-22, 20:01:38 | scheduled__2022-11-23T01:01:38.780484+00:00 | scheduled | 2022-11-22, 20:31:38 | 2022-11-22, 20:31:38 | 2022-11-22, 20:32:02 | False |
| ☐ | ✎ 🗑 | success | Q2.1 | 2022-11-22, 19:31:38 | scheduled__2022-11-23T00:31:38.780484+00:00 | scheduled | 2022-11-22, 20:01:39 | 2022-11-22, 20:01:39 | 2022-11-22, 20:02:00 | False |
| ☐ | ✎ 🗑 | success | Q2.1 | 2022-11-22, 19:01:38 | scheduled__2022-11-23T00:01:38.780484+00:00 | scheduled | 2022-11-22, 19:32:05 | 2022-11-22, 19:32:05 | 2022-11-22, 19:32:28 | False |

I set the start date be the date and time right now(I first change the time zone on my vm to New York) so that it would tart immediately. I set the schedule interval to be 30minutes.

```python
with DAG(
    'Q2.1',
    default_args=default_args,
    description='Task 2 impliment the DAG below',
    schedule_interval=timedelta(minutes=30),
    start_date=datetime(2022, 11, 22, 19, 35),
    catchup=False,
    tags=['homework'],
) as dag:
```

## Q2.2

I build this workflow by separating our goal to three parts. First, I need to fetch data and push these data with xcom. Second, I need to train the model and get prediction for next day and save it to a file called 'predict.csv'. Third, I need to calculate the error between yesterday's prediction and today's price. I then saved all errors to another file called 'errors.csv'.

I manage cross task communication by using xcom_push and xcom_pull. These command will push key and value pairs to dictionaries in each task. I could access values in each dictionary with task id and keys.

I setup scheduler by specifying 'start date' and 'schedule interval' when setting up DAG. I let the schedule interval be 7 am UTC each day. Here is my DAG and code.



```
1  from datetime import datetime, timedelta, date
2  from textwrap import dedent
3  import time
4  import yfinance as yf
5  import numpy as np
6  from sklearn.linear_model import LinearRegression
7  import pandas as pd
8  from csv import writer, reader
9  from airflow import DAG
10 from airflow.operators.bash import BashOperator
11 from airflow.operators.python import PythonOperator
12
13 def fatch_data(**context):
14     symbles = ["AAPL", "GOOGL", "META", "MSFT", "AMZN"]
15     for s in symbles:
16         tickers = yf.Tickers(s)
17         data = tickers.tickers[s].info
18         open_price = data["open"]
19         high_price = data["dayHigh"]
20         low_price = data["dayLow"]
21         close_price = data["previousClose"]
22         volume = data["volume"]
23         context['ti'].xcom_push(key=s, value=[open_price, high_price, low_price, close_price, volume])
24
25 def train_model(**context):
26     symbles = ["AAPL", "GOOGL", "META", "MSFT", "AMZN"]
27     all_predict = []
28     for s in symbles:
29         todayPrice = context['ti'].xcom_pull(key=s,task_ids="fatch_data")
30         data = yf.download(tickers = s, period = "11d", interval = "1d")
31         data = data.drop(columns=["Adj Close"])
32         x = data[:-1]
33         y = data["High"][1:]
34         reg = LinearRegression().fit(x, y)
35         print(f"Current model accuracy: {reg.score(x, y)}")
36         features = pd.DataFrame([todayPrice])
37         prediction = reg.predict(features)[0]
38         print(f"Prediction for next day's high: {prediction}")
39         all_predict.append(prediction)
```

```python
            context['ti'].xcom_push(key=s, value=prediction)
    today = date.today()
    d = today.strftime("%m/%d/%y")
    write = [d] + all_predict
    with open(f'predict.csv','a') as file:
        writer_object = writer(file)
        writer_object.writerow(write)
        file.close()

def caculate_error(**context):
    symbles = ["AAPL", "GOOGL", "META", "MSFT", "AMZN"]
    errors = []
    predictions = []
    today = date.today()
    d = today.strftime("%m/%d/%y")
    yesterday = today - timedelta(days = 1)
    yesterday = yesterday.strftime("%m/%d/%y")
    with open('predict.csv', 'r') as file:
        reader_object = reader(file)
        for row in reader_object:
            if row and row[0] == yesterday:
                predictions = row[1:]
                break
    print(f"Yesterday's prediction: {predictions}")
    for i in range(len(symbles)):
        high_price = context['ti'].xcom_pull(key=symbles[i],task_ids="fatch_data")[1]
        error = (float(predictions[i]) - float(high_price)) / float(high_price)
        errors.append(error)
    write = [d] + errors
    with open(f'errors.csv','a') as file:
        writer_object = writer(file)
        writer_object.writerow(write)
        file.close()
    print(f"Today's error: {write}")

default_args = {
    'owner': 'yutao',
    'depends_on_past': False,
    'email': ['yz4359@columbia.edu'],
```

```python
        'email_on_failure': True,
        'email_on_retry': False,
        'retries': 1,
        'retry_delay': timedelta(seconds=3),
        # 'queue': 'bash_queue',
        # 'pool': 'backfill',
        # 'priority_weight': 10,
        # 'end_date': datetime(2016, 1, 1),
        # 'wait_for_downstream': False,
        # 'dag': dag,
        # 'sla': timedelta(hours=2),
        # 'execution_timeout': timedelta(seconds=300),
        # 'on_failure_callback': some_function,
        # 'on_success_callback': some_other_function,
        # 'on_retry_callback': another_function,
        # 'sla_miss_callback': yet_another_function,
        # 'trigger_rule': 'all_success'
}

with DAG(
    'Q2.2',
    default_args=default_args,
    description='Stock price fetching, prediction, and storage every day.',
    schedule_interval='0 7 * * *',
    start_date=datetime(2022, 11, 24),
    catchup=False,
    tags=['homework'],
) as dag:

    fatch_data = PythonOperator(
        task_id='fatch_data',
        python_callable=fatch_data,
        retries=3,
        provide_context=True
    )

    caculate_error = PythonOperator(
        task_id='caculate_error',
        python_callable=caculate_error,
```

```
118              retries=3,
119              provide_context=True
120         )
121
122     train_model = PythonOperator(
123              task_id='train_model',
124              python_callable=train_model,
125              retries=3,
126              provide_context=True
127         )
128
129     fatch_data >> train_model
130     train_model >> caculate_error
```
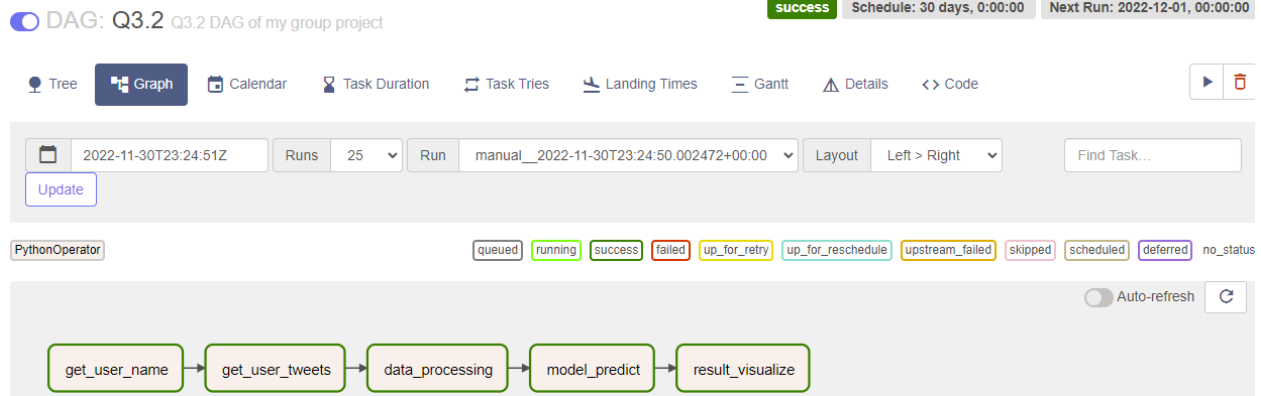
## Task 3 Written parts
## Q3.1
## (1)

| Executor | Pros | Cons |
|---|---|---|
| **SequentialExecutor** | Every task is in sequential. It is very easy to debug and write code since there can be only one task executing at any given time(even branch task from same root). | When there are many task or some task have many dependency, SequentialExecutorwould be really slow. A task will wait for all of its depended task finish executing one by one to continue. |
| **LocalExecutor** | Execute task in parallel. It would be a | This would be harder for debugging and |

| | | |
|---|---|---|
| | lot faster and utilizing more computational resource when executing task. | programing. We do not know which task will finish executing before hand. As a result we have to carefully write dependency and make sure task would not collide with each other(read and write same file at same time). We have to use lock when necessary to avoid collision. |
| **CeleryExecutor** | Most mature option because it is the oldest adoption. Many resource online since a lot of people and company are using it. | It require infrastructure support to work. Need Celery and Celery's backend. |
| **KubernetesExecutor** | Able to use different docker images for different task. More flexibility. Works well with Kubernetes ecosystem. | Not as robest as other executor because it is newer. |

**Q3.2**

**(3)**

We will schedule our tasks to run every month(every 30 days as shown in the screen shoot).