# Clear your schedule: Running Data Science with python and airflow

●●●

Jonathan Stott

mago

@namelessjon

Production?

@namelessjon

Schedule?

```
* * * * *

1 * * * *

5 8 * * *
```

@namelessjon

# More Schedules

# Enter: Airflow



`$ pip install apache-airflow`

# Directed Acyclic Graphs

Operators

@namelessjon

Sensors

@namelessjon

# UI / Dashboard

# Building a workflow

```python
from airflow import DAG
from airflow.sensors.s3_key_sensor import S3KeySensor
from airflow.operators.docker_operator import DockerOperator
import datetime as dt

default_args = {
    'start_date': dt.datetime(2020, 4, 29),
    'retries': 1,
    'retry_delay': dt.timedelta(minutes=5),
}

with DAG('example', default_args=default_args, schedule_interval="13 5 * * *") as dag:

    # Look for s3://example-bucket/Data/2020-04-29/data.csv
    t1 = S3KeySensor(
        task_id='check_for_data',
        bucket='example-bucket'
        bash_command='Data/{{ ds }}/data.csv',)

    t2 = DockerOperator(
        task_id='process_1',
        container='some_container',
        command=["python3", "process_1.py",
                 "--input-file", "s3://example-bucket/Data/{{ ds }}/data.csv",
                 "--output-file",  "s3://example-bucket/Output/{{ ds }}/output.csv"],
        retries=3,)

    t1 >> t2
    # more tasks ...
```

# Other helpful features

- Task logging
- Connections
- Alerting
- Celery
- SLAs
- ... and more!

Use operators where you can

@namelessjon

# Use containers for scripts

@namelessjon

Keep business logic out of the DAG file

Photo by Sarah Shaffer on Unsplash

Remember DAGs are python

@namelessjon

# Summary

- Use airflow to manage workflows with dependencies between tasks
- Docker is great for packing up tasks in workflow
- Start small and scale up!

# Thank you

# Links

https://airflow.apache.org/

https://towardsdatascience.com/best-practices-for-airflow-developers-990c8a04f7c6

https://medium.com/wbaa/datas-inferno-7-circles-of-data-testing-hell-with-airflow-cef4adff58d8