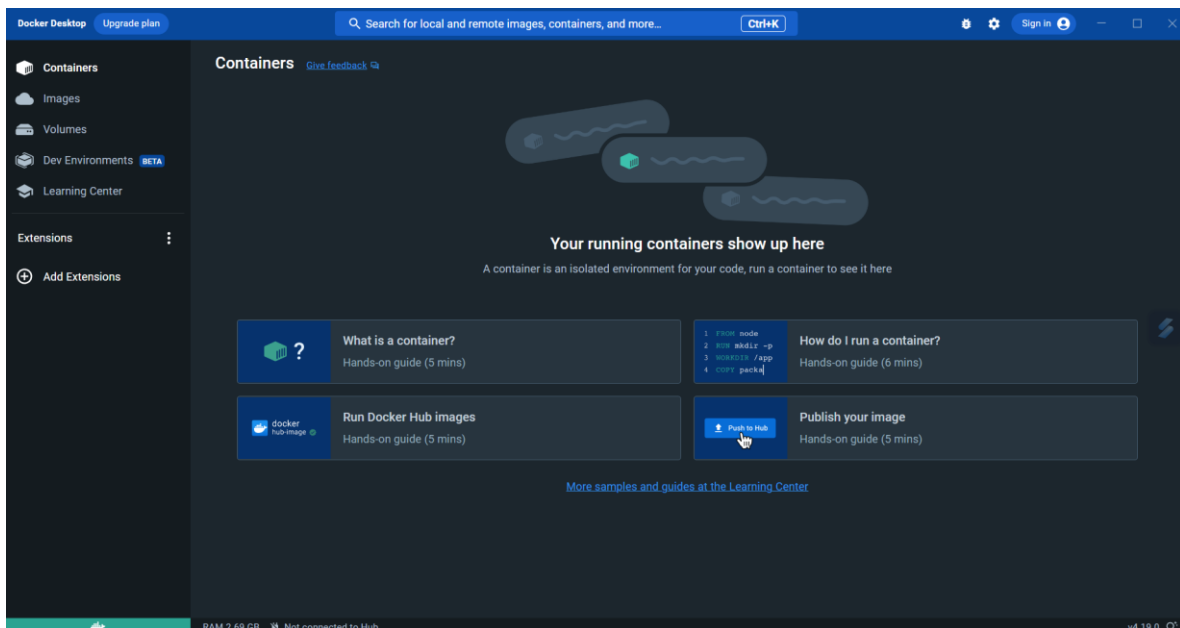


Solución desafío técnico: construir un pipeline con Docker y Apache Ariflow

Objetivo: Crear un pipeline de datos simple usando Apache Airflow que extraiga datos de un archivo CSV, los transforme y los exporte a otro archivo CSV, todo ejecutándose en un contenedor Docker. Este desafío te ayudará a demostrar tus habilidades en Python, procesos ETL, trabajo con DAG's en Apache Airflow y contenerización con Docker. Cualquier duda enviar un correo ageracollante95@gmail.com.

1. Configurar Apache Airflow con Docker
 - a. Instalar Docker en tu máquina local.

R/ Se descarga docker y se instala en una maquina local.



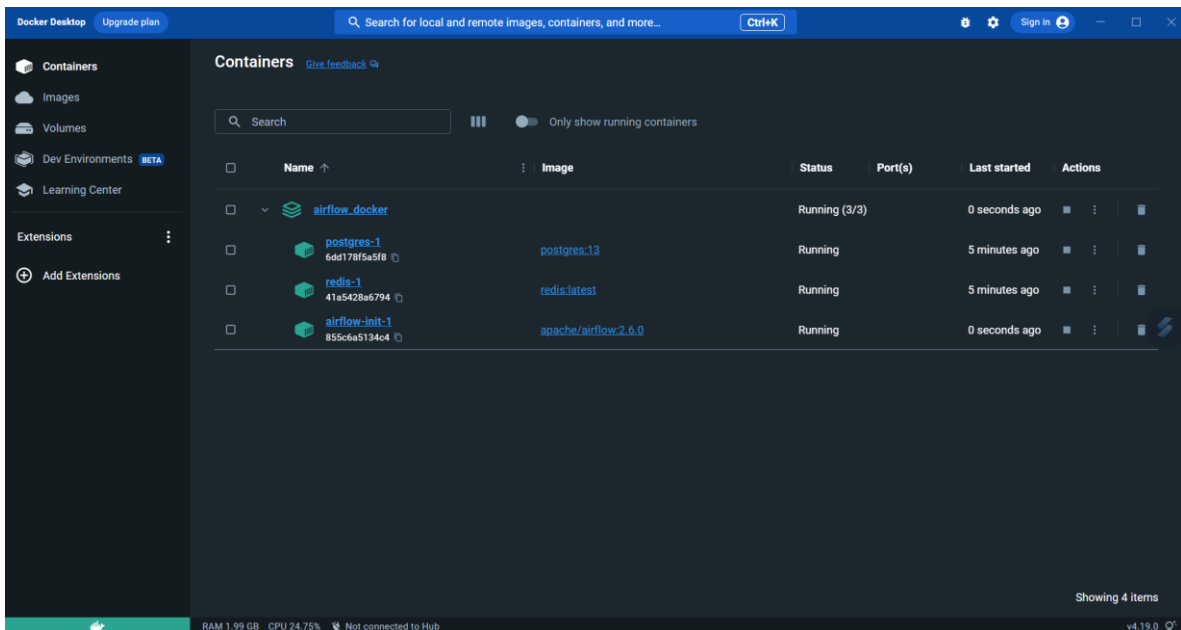
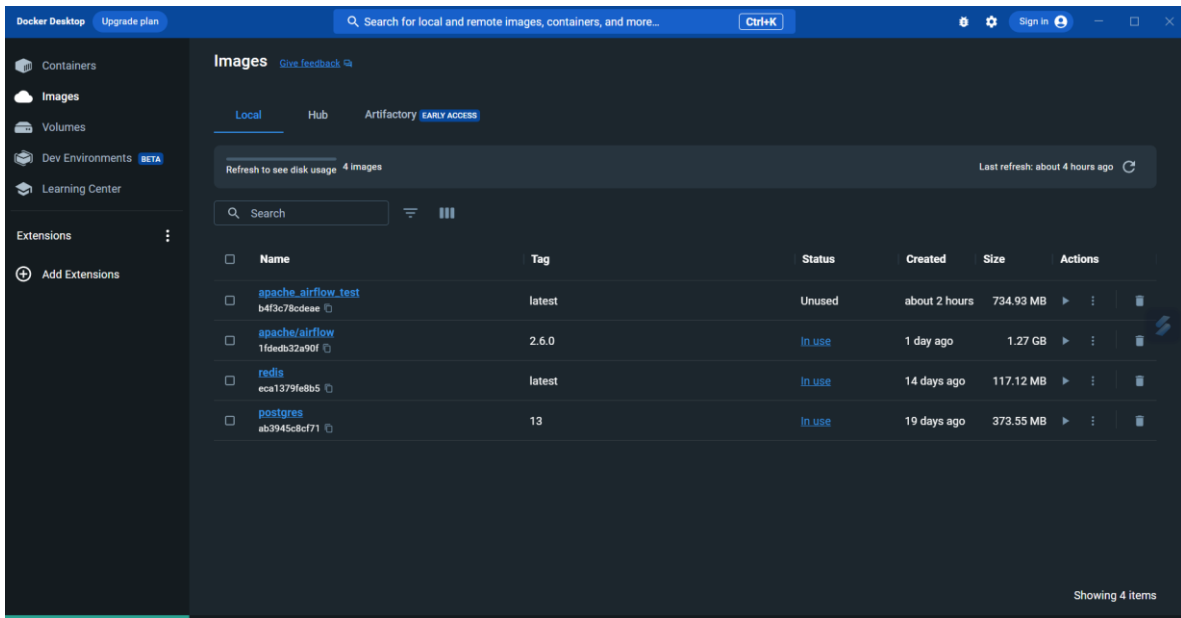
Nota: tener en cuenta la actualización para ejecutar subsistemas Linux (WSL2).

- b. Crear un Dockerfile que incluya las dependencias requeridas y la configuración para Apache Airflow

R/ La instalación de airflow se realiza por un archivo “.yaml”, ya que la instalación por este medio es más manejable que por una serie de instrucciones armadas por medio de un dockerfile, el “.yaml” fue descargado desde la página de Apache Airflow ([LINK](#)), tener en cuenta que se debe instalar docker-compose, el cual también se explica [AQUI](#)

```
PROBLEMS  OUTPUT  DEBUG CONSOLE  TERMINAL
PS D:\Airflow-Docker> docker-compose up airflow-init
time="2023-05-01T11:55:59-05:00" level=warning msg="The \"AIRFLOW_UID\" variable is not set. Defaulting to a blank string."
time="2023-05-01T11:55:59-05:00" level=warning msg="The \"AIRFLOW_UID\" variable is not set. Defaulting to a blank string."
[+] Running 7/4
 - airflow-init 19 layers [#####] 53.72MB/315.8MB Pulling 17.0s
 - postgres 12 layers [#####] 0B/0B Pulling 17.0s
 - redis 5 layers [#####] 0B/0B Pulling 17.0s
```

```
PS D:\Airflow_docker> docker-compose up airflow-init
time="2023-05-01T11:55:59-05:00" level=warning msg="The \"AIRFLOW_UID\" variable is not set. Defaulting to a blank string."
time="2023-05-01T11:55:59-05:00" level=warning msg="The \"AIRFLOW_UID\" variable is not set. Defaulting to a blank string."
[+] Running 39/4
  ✓ airflow-init 19 layers [#####] 00/00 Pulled 344.9%
  ✓ postgres 12 layers [#####] 00/00 Pulled 281.9%
  ✓ redis 5 layers [#####] 00/00 Pulled 66.4%
[+] Running 5/5
  ✓ Network airflow_docker_default Created 0.5s
  ✓ Volume "airflow_docker_postgres-db-volume" Created 0.1s
  ✓ Container airflow_docker-redis-1 Created 3.3s
  ✓ Container airflow_docker-postgres-1 Created 3.3s
  ✓ Container airflow_docker-airflow-init-1 Created 1.0s
Attaching to airflow_docker-airflow-init-1
```



c. Construir y ejecutar el contenedor de Airflow utilizando el Dockerfile

R/ Dado que se usó el “.yaml” el proceso se ejecuta automáticamente una vez se inicia docker.

2. Extracción de Datos:

- a. Descargar el archivo CSV de ejemplo que contiene datos abiertos de accidentes de tráfico en la ciudad de Seattle.

R/ Solución explicada al finalizar las preguntas en general

- b. Crear una tarea PythonOperator de Airflow para leer el archivo CSV.

R/ Solución explicada al finalizar las preguntas en general

3. Transformación de Datos:

- a. Crear una tarea PythonOperator de Airflow para limpiar y transformar los datos eliminando cualquier fila con valores faltantes.

R/ Se identificó que el archivo NO tiene información nula o vacía, sin embargo, se incorporó al proceso una validación para suplir este requerimiento explicada al finalizar las preguntas en general

- b. Calcular el número total de accidentes por tipo de clima.

R/ No se encontró información del clima, en su lugar se realizó un conteo de accidentalidad por año

4. Exportación de Datos:

- a. Crear una tarea PythonOperator de Airflow para exportar los datos transformados a un nuevo archivo CSV.

R/ Solución explicada al finalizar las preguntas en general

5. Crear el DAG:

- a. Definir un DAG que incluya las tareas de extracción, transformación y exportación de datos creadas en los pasos 2-4.

R/ Se agrupa respuesta con las preguntas anteriores, el dag se encuentra en el [GITHUB](#) con todo lo relacionado

- b. Establecer las dependencias de las tareas para garantizar que se ejecuten en el orden correcto.

R/ igual que el punto anterior

- c. Programar el DAG para que se ejecute en un intervalo específico (por ejemplo, todos los días a la media noche).

R/ igual que el punto anterior

Como fue mencionado en las respuestas relacionadas, se creó un DAG, donde hay 3 funciones:

- ✓ Función input_file:

La función descarga el archivo <https://data.seattle.gov/api/views/38vd-gyvt/rows.csv?accessType=DOWNLOAD> y lo almacena en la carpeta input

```
# ----- Function definition -----
def input_file(url):
    """Descarga archivo a local de acuerdo a url determinada."""
    response = requests.get(url)
    input_file_name = os.path.join(in_path, file_name)
    print(f'\n\n {response.status_code} \n\n')
    print(f'\n\n {input_file_name} \n\n')
    if response.status_code == 200:
        with open(input_file_name, "wb") as file:
            file.write(response.content)
        return '\n\n Done input file \n\n'
    else:
        return '\n\n Don't input file \n\n'
```

se crea una variable en airflow para setear la url mencionada y así dar oportunidad a modificar el link de obtención de la información



Edit Variable	
Key *	url_test
Val	https://data.seattle.gov/api/views/38vd-gyvt/rows.csv?accessType=DOWNLOAD
Description	
<div>Save ↩</div>	

✓ Función output_file:

La función captura el archivo descargado en la carpeta input, realiza borrado de filas si encuentra alguna nula y lo deja con las 4 primeras columnas, el archivo se deja en la carpeta output

```
def output_file(in_p, out_p):
    """Transforma el archivo importado."""
    path_file = os.path.join(in_p, file_name)
    out_path_file = os.path.join(out_p, output_file_name)
    file_df = pd.read_csv(path_file, sep=',', header=0)
    print(f'\n\n {file_df.head(10)} \n\n')
    file_df = file_df.dropna(how='any').reset_index(drop=True)
    file_df = file_df.iloc[:, [0, 1, 2, 3]]
    print(f'\n\n {file_df.head(10)} \n\n')
    file_df.to_csv(out_path_file, index=False, sep=';')
    return '\n\n Done output file \n\n'
```

✓ Función reader_file:

La función captura el archivo alojado en la carpeta output y hace una operación simple, realiza un conteo de accidentes por año y los muestra en el log de la tarea en airflow.

```
def reader_file(out_p):
    """Realiza operaciones sobre el archivo transformado."""
    out_path_file = os.path.join(out_p, output_file_name)
    file_df = pd.read_csv(out_path_file, sep=';', header=0)
    conteo = file_df.groupby(['YEAR']).size().reset_index(name='COUNT')
    print(f'\n\n {file_df.info()} \n\n')
    return f'\n\n {conteo} \n\n'
```

```
localhost:8080/log?dag_id=airflow_test&task_id=reader_file_task&execution_date=2023-05-03T22:05:00%3A00%3A00%2B00%3A00

Airflow DAGs Security Browse Admin Docs Astronomer

1

*** Reading local file: /usr/local/airflow/logs/dag_id=airflow_test/run_id=scheduled__2023-05-03T22:05:00+00:00/task_id=reader_file_task/attempt=1.log
[2023-05-03, 22:10:09 UTC] {taskinstance.py:1173} INFO - Dependencies all met for <TaskInstance: airflow_test.reader_file_task scheduled__2023-05-03T22:05:00+00:00 [queued]
[2023-05-03, 22:10:09 UTC] {taskinstance.py:1173} INFO - Dependencies all met for <TaskInstance: airflow_test.reader_file_task scheduled__2023-05-03T22:05:00+00:00 [queued]
[2023-05-03, 22:10:09 UTC] {taskinstance.py:1370} INFO -

[2023-05-03, 22:10:09 UTC] {taskinstance.py:1371} INFO - Starting attempt 1 of 1
[2023-05-03, 22:10:09 UTC] {taskinstance.py:1372} INFO -

[2023-05-03, 22:10:09 UTC] {taskinstance.py:1391} INFO - Executing <Task(PythonOperator): reader_file_task> on 2023-05-03 22:05:00+00:00
[2023-05-03, 22:10:09 UTC] {standard_task_runner.py:52} INFO - Started process 13145 to run task
[2023-05-03, 22:10:09 UTC] {standard_task_runner.py:79} INFO - Running: ['airflow', 'tasks', 'run', 'airflow_test', 'reader_file_task', 'scheduled__2023-05-03T22:05:00+00:00:00']
[2023-05-03, 22:10:09 UTC] {standard_task_runner.py:80} INFO - Job 2711: Subtask reader_file_task
[2023-05-03, 22:10:09 UTC] {task_command.py:370} INFO - Running <TaskInstance: airflow_test.reader_file_task scheduled__2023-05-03T22:05:00+00:00 [running]> on host 926a63f
[2023-05-03, 22:10:10 UTC] {taskinstance.py:1583} INFO - Exporting the following env vars:
AIRFLOW_CTX_DAG_OWNER=gustavo8704@hotmail.com
AIRFLOW_CTX_DAG_ID=airflow_test
AIRFLOW_CTX_TASK_ID=reader_file_task
AIRFLOW_CTX_EXECUTION_DATE=2023-05-03T22:05:00+00:00
AIRFLOW_CTX_TRY_NUMBER=1
AIRFLOW_CTX_DAG_RUN_ID=scheduled__2023-05-03T22:05:00+00:00
[2023-05-03, 22:10:10 UTC] {logging_mixin.py:115} INFO - <class 'pandas.core.frame.DataFrame'>
RangeIndex: 437 entries, 0 to 436
Data columns (total 4 columns):
#   Column          Non-Null Count  Dtype
---  --
0   OBJECTID        437 non-null    int64
1   STNAME          437 non-null    object
2   COUNT_LOCATION  437 non-null    object
3   YEAR            437 non-null    int64
dtypes: int64(2), object(2)
memory usage: 13.8+ KB
[2023-05-03, 22:10:10 UTC] {logging_mixin.py:115} INFO -

None
[2023-05-03, 22:10:10 UTC] {python.py:173} INFO - Done. Returned value was:

   YEAR  COUNT
0  2017    278
1  2018    159

[2023-05-03, 22:10:10 UTC] {taskinstance.py:1409} INFO - Marking task as SUCCESS. dag_id=airflow_test, task_id=reader_file_task, execution_date=20230503T220500, start_date=
[2023-05-03, 22:10:10 UTC] {local_task_job.py:156} INFO - Task exited with return code 0
[2023-05-03, 22:10:10 UTC] {local_task_job.py:273} INFO - 0 downstream tasks scheduled from follow-on schedule check
```

Cada función se ejecuta por medio de PythonOperator y se determinó el orden de ejecución como se refleja a continuación:

```
# ----- Dag definition -----
default_args = {
    'owner': 'gustavo8704@hotmail.com',
    'depends_on_past': False,
    'retries': 0
}

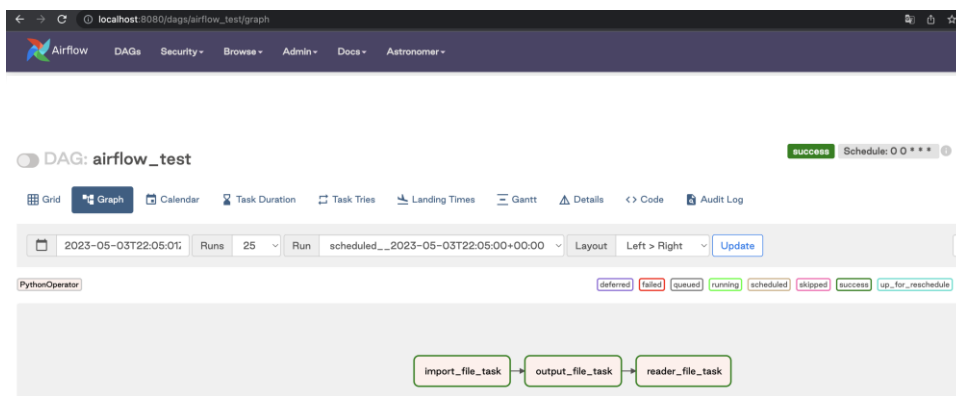
dag = DAG('airflow_test',
          max_active_runs=1,
          schedule_interval='0 0 * * *',
          start_date=datetime(2022, 1, 13),
          tags=['1.input_file', '2.output_file', '3.reader_file'],
          default_args=default_args,
          catchup=False
          )
```

```

with dag:
    import_file_task = PythonOperator(
        task_id='import_file_task',
        python_callable=input_file,
        op_args=[url]
    )
    output_file_task = PythonOperator(
        task_id='output_file_task',
        python_callable=output_file,
        op_args=[in_path, out_path]
    )
    reader_file_task = PythonOperator(
        task_id='reader_file_task',
        python_callable=reader_file,
        op_args=[out_path]
    )
    import_file_task
    output_file_task.set_upstream(import_file_task)
    reader_file_task.set_upstream(output_file_task)

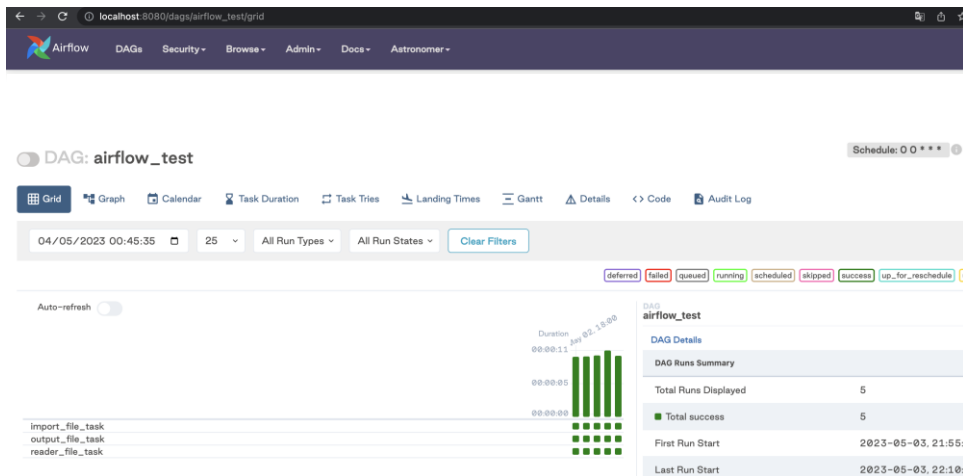
```

Se ilustra la dependencia de las tareas en la siguiente imagen



También se evidencia como se ve el dag desde airflow

DAG	Owner	Runs	Schedule	Last Run	Next Run	Recent Tasks
airflow_test	gustavo8704@hotmail.com	5	0 0 * * *	2023-05-03, 22:05:00	2023-05-03, 22:10:00	1:input_file, 2:output_file, 3:reader_file



Entregables:

- Dockerfile para el contenedor de Airflow.
- Scripts de Python para las tareas de Airflow, incluidos extracción, transformación y exportación de datos.
- El archivo de definición del DAG de Airflow.
- Un archivo README que contenga instrucciones sobre cómo configurar y ejecutar el proyecto localmente utilizando Docker, así como cualquier dependencia requerida.

Porfavor, crea un repositorio en GitHub que contenga todos los archivos requeridos y comparte el enlace del repositorio para enviar tu solución. ¡Buenasuerte!

Los entregables se encuentran en el siguiente link de [GITHUB](https://github.com/gagarb/airflow_tesst/tree/main/airflow_test)
https://github.com/gagarb/airflow_tesst/tree/main/airflow_test