

Práctica 0: Despliegue del entorno y primeros pasos

Procesamiento Masivo de Datos y Visualización (PMDV)

Curso 2025-2026

1. Introducción

Con el objetivo de poner en práctica los conceptos aprendidos en teoría, se dispone de un entorno que incluye algunas de las principales herramientas.

- **Almacenamiento distribuido:** MinIO
- **Procesamiento Distribuido:** Apache Spark, para su uso con Python (PySpark). Además, se incluyen las librerías necesarias para usar los formatos Delta e Iceberg.
- **Orquestación:** Apache Airflow
- **Calidad de Datos:** Great Expectations
- **Exploración y Visualización:** Jupyter Lab. Además, se incluyen las librerías Plotly para Python.

Las anteriores herramientas conforman una arquitectura Big Data moderna, de código abierto y ampliamente utilizadas en los proyectos de datos de las empresas y organizaciones.

Para facilitar su despliegue, el entorno se ha definido con la herramienta Docker, la cual permite empaquetar y distribuir las distintas herramientas que conforman la arquitectura. Las herramientas ya se encuentran pre-configuradas para uso conjunto.

El aprendizaje de Docker queda fuera del alcance de esta asignatura. No obstante, se proporcionan los comandos básicos para su uso, incluyendo el arranque y parada del entorno.

En cuanto al entorno de desarrollo, se usará principalmente de Visual Studio Code/PyCharm. No obstante, también se usará la herramienta de Notebooks, Jupyter Lab. Esta herramienta permitirá ejecutar el código Spark (PySpark) de forma interactiva, facilitando el aprendizaje de Spark, la exploración interactiva de datos y la generación de visualizaciones sobre los datos procesados y almacenados en nuestro Data Lake.

Tanto Docker Desktop como Visual Studio Code, son herramientas gratuitas y se encuentran disponibles en los ordenadores de los laboratorios. Se usará el SO Windows para las prácticas y solo se garantiza el funcionamiento en este SO.

No obstante, el entorno y herramientas son compatibles con Linux (ej. Ubuntu), aunque podrían ser necesarios algunos ajustes en los scripts de Docker que definen el entorno.

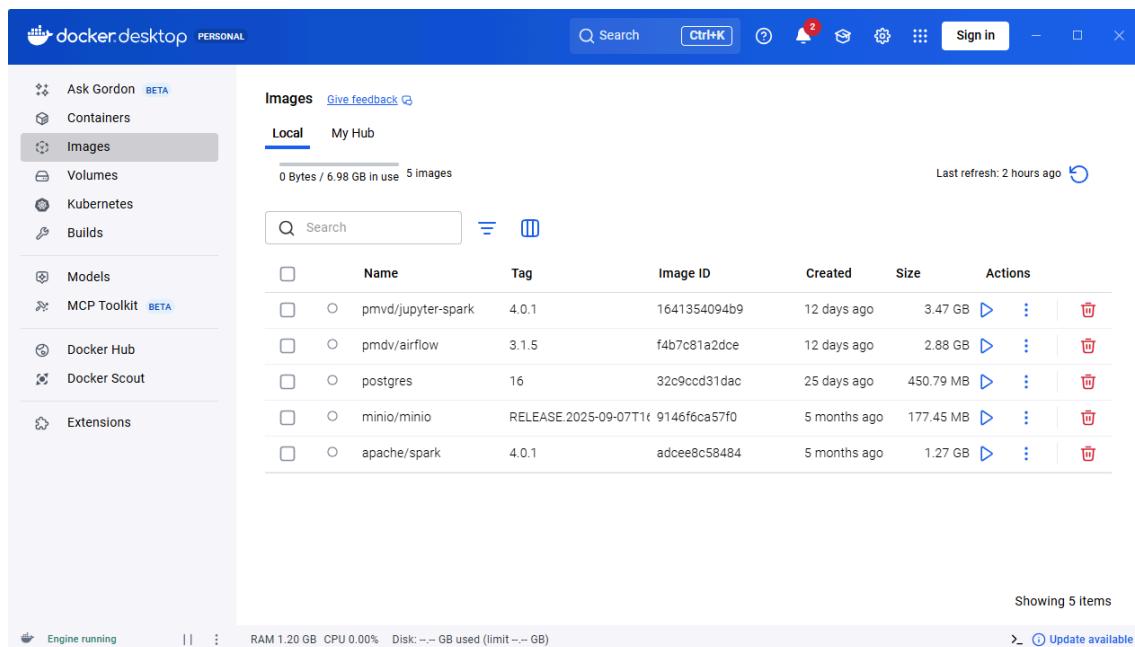
2. Descripción de la práctica

El objetivo principal de esta primera práctica es aprender a arrancar el entorno y probar el acceso a las distintas herramientas, como punto de partida fundamental para poder llevar a cabo las siguientes prácticas de la asignatura.

En los siguientes apartados se detalla cada una las tareas a llevar a cabo, así como la salida de cada una de ellas que forma parte de los archivos a entregar para la evaluación de esta práctica. En el apartado “3. Entrega de la Práctica”, se enumeran los archivos a entregar y las condiciones para dicha entrega.

2.1 Arranque del entorno

En primer lugar, arrancaremos Docker Desktop desde el menú de inicio. Será necesario para poder ejecutar los comandos de Docker más adelante, en la terminal que abriremos desde Visual Studio Code.



En los laboratorios, dado que se solicitó la creación previa del entorno con Docker en cada puesto, debería aparecer el listado de imágenes anteriores. En otro caso se han revisar los anexos para alternativas de uso mediante USB o instalación en tu equipo personal.

Tras esto, se ha de localizar o descargar la carpeta “infra”. En los ordenadores de los laboratorios debería encontrarse en “C:\lsw\infra”. En otro caso puede descargarse desde el Moodle de la asignatura “Software de la asignatura sobre Docker”.

Dentro de esta carpeta, tendremos a su vez dos carpetas:

- **datalake**: Contiene los archivos imprescindibles para desplegar el entorno con Docker.
 - docker-compose.yml
 - requirements.txt
 - .env
 - spark
 - docker-airflow
 - docker-jupyterlab
- **tests**: Contiene los archivos con los ejemplos de código a ejecutar como parte de esta Práctica 0.
 - airflow
 - data
 - jobs
 - notebooks

Se ha de mover o copiar las 4 carpetas dentro de *tests* en la carpeta *datalake*.

En este momento, se iniciará Visual Studio Code, que será nuestro entorno de desarrollo principal. En Visual Studio, hacemos clic en *File->Open Folder*. Se debe buscar la carpeta *datalake* y seleccionarla para abrirla.

No es necesario abrir ningún archivo de la carpeta *datalake* para arrancar el entorno. No obstante, el archivo raíz con la definición del entorno se encuentra en *datalake/docker-compose.yml*. Este archivo es el que usará por defecto el comando de arranque de Docker que se ejecutará a continuación.

Para ejecutar ese comando y cualquier otro comando Docker, se usará una terminal en Visual Studio (zona de abajo a la derecha en la captura anterior). Si no apareciese una terminal, es posible abrir una nuevo haciendo clic en el menú superior, *Terminal->New Terminal*.

```

version: '3.7'
services:
  airflow:
    image: pmdv/airflow:3.1.5
    pull_policy: never
    build:
      context: .
      dockerfile: docker-airflow/Dockerfile
    volumes:
      - ./airflow/dags:/opt/airflow/dags
      - ./airflow/logs:/opt/airflow/logs
      - ./airflow/config:/opt/airflow/config
      - ./airflow/plugins:/opt/airflow/plugins
      - ./jobs:/opt/airflow/jobs
      - ./data:/data
      - ./var/run/docker.sock:/var/run/docker.sock
    env_file: .env
    user: "${AIRFLOW_UID:-50000}:0"
    depends_on:
      - airflow-postgres
      - condition: service_healthy
    networks:
      - datalake
    restart: unless-stopped
  airflow-runtime:
    image: pmdv/airflow:3.1.5
    pull_policy: never
    env_file: .env
    volumes:
      - ./airflow/dags:/opt/airflow/dags
      - ./airflow/logs:/opt/airflow/logs
      - ./airflow/config:/opt/airflow/config
      - ./airflow/plugins:/opt/airflow/plugins
      - ./jobs:/opt/airflow/jobs
      - ./data:/data
      - ./var/run/docker.sock:/var/run/docker.sock
    user: "${AIRFLOW_UID:-50000}:0"
    depends_on:
      - airflow-postgres
      - condition: service_healthy
    networks:
      - datalake
    restart: unless-stopped
  minio:
    image: minio/minio:RELEASE.2025-09-07T16-13-09Z-cpu1
    container_name: minio
    command: server /data --console-address ":9001"
    environment:
      MINIO_ROOT_USER: minio
      MINIO_ROOT_PASSWORD: minio123
    volumes:
      - ./data:/data
    networks:
      - datalake
    restart: unless-stopped
  airflow-postgres:
    image: postgres:15
    environment:
      POSTGRES_DB: airflow
      POSTGRES_USER: airflow
      POSTGRES_PASSWORD: airflow
    networks:
      - datalake
    restart: unless-stopped
  spark-master:
    image: spark:3.2.1
    environment:
      SPARK_MASTER_IP: $(minio_ip)
      SPARK_MASTER_PORT_HTTP: 7077
    networks:
      - datalake
    restart: unless-stopped
  spark-worker:
    image: spark:3.2.1
    environment:
      SPARK_MASTER_URL: http://$(spark_master_ip):7077
    networks:
      - datalake
    restart: unless-stopped
  spark-dagprocessor:
    image: spark:3.2.1
    environment:
      SPARK_MASTER_URL: http://$(spark_master_ip):7077
    networks:
      - datalake
    restart: unless-stopped
  spark-scheduler:
    image: spark:3.2.1
    environment:
      SPARK_MASTER_URL: http://$(spark_master_ip):7077
    networks:
      - datalake
    restart: unless-stopped
  spark-apiserver:
    image: spark:3.2.1
    environment:
      SPARK_MASTER_URL: http://$(spark_master_ip):7077
    networks:
      - datalake
    restart: unless-stopped
  jupyter:
    image: jupyter/tensorflow-notebook:latest
    environment:
      JUPYTER_TENSORFLOW_NOTEBOOK_KERNEL_NAME: tensorflow
    networks:
      - datalake
    restart: unless-stopped
  airflow-triggerer:
    image: pmdv/airflow-triggerer:3.1.5
    environment:
      AIRFLOW_CONN_AIRFLOW_DEFAULT: "mysql://airflow:airflow@$(minio_ip)/airflow"
    networks:
      - datalake
    restart: unless-stopped
  airflow-dagprocessor:
    image: pmdv/airflow-dagprocessor:3.1.5
    environment:
      AIRFLOW_CONN_AIRFLOW_DEFAULT: "mysql://airflow:airflow@$(minio_ip)/airflow"
    networks:
      - datalake
    restart: unless-stopped
  airflow-scheduler:
    image: pmdv/airflow-scheduler:3.1.5
    environment:
      AIRFLOW_CONN_AIRFLOW_DEFAULT: "mysql://airflow:airflow@$(minio_ip)/airflow"
    networks:
      - datalake
    restart: unless-stopped
  airflow-apiserver:
    image: pmdv/airflow-apiserver:3.1.5
    environment:
      AIRFLOW_CONN_AIRFLOW_DEFAULT: "mysql://airflow:airflow@$(minio_ip)/airflow"
    networks:
      - datalake
    restart: unless-stopped

```

Una vez en el terminal se ejecutará el siguiente comando para arrancar el entorno:

- **docker compose up -d**

Si todo va bien, se mostrarán y listarán los distintos contenedores que contienen los servicios que conforman nuestro entorno (minio, spark-master, spark-worker,...), informado sobre el estado del arranque. Si el arranque concluye con éxito obtendremos el siguiente resultado.

```

PS C:\pmvd-lab\infra\datalake> docker compose up -d
[+] Running 11/11
  ✓ Network datalake
  ✓ Container airflow-postgres          Created                         0.1s
  ✓ Container minio                     Healthy                        7.8s
  ✓ Container spark-master             Started                        2.6s
  ✓ Container spark-worker             Started                        2.6s
  ✓ Container datalake-airflow-init-1 Exited                         41.5s
  ✓ Container spark-worker             Started                        2.7s
  ✓ Container jupyter                  Started                        2.9s
  ✓ Container airflow-triggerer        Started                        41.8s
  ✓ Container airflow-dagprocessor     Started                        41.7s
  ✓ Container airflow-scheduler       Started                        41.7s
  ✓ Container airflow-apiserver       Started                        41.8s
PS C:\pmvd-lab\infra\datalake>

```

En este momento ya se puede avanzar al siguiente apartado “2.2 Acceso a MinIO y carga manual de datos”.

Aunque no lo ejecutes aún si vas a continuar con la práctica, a continuación, se indica el comando para detener el entorno:

- `docker compose down`

```
PS C:\pmvd-lab\infra\datalake> docker compose down
[+] Running 11/11
✓ Container airflow-scheduler      Removed   2.4s
✓ Container airflow-dagprocessor   Removed   2.8s
✓ Container spark-worker          Removed   1.2s
✓ Container airflow-triggerer    Removed   2.3s
✓ Container airflow-jupyter      Removed   2.1s
✓ Container airflow-apiserver    Removed   2.6s
✓ Container minio                 Removed   1.2s
✓ Container spark-master         Removed   1.6s
✓ Container datalake-init-1     Removed   0.1s
✓ Container airflow-postgres     Removed   0.9s
✓ Network datalake              Removed   0.6s
PS C:\pmvd-lab\infra\datalake>
```

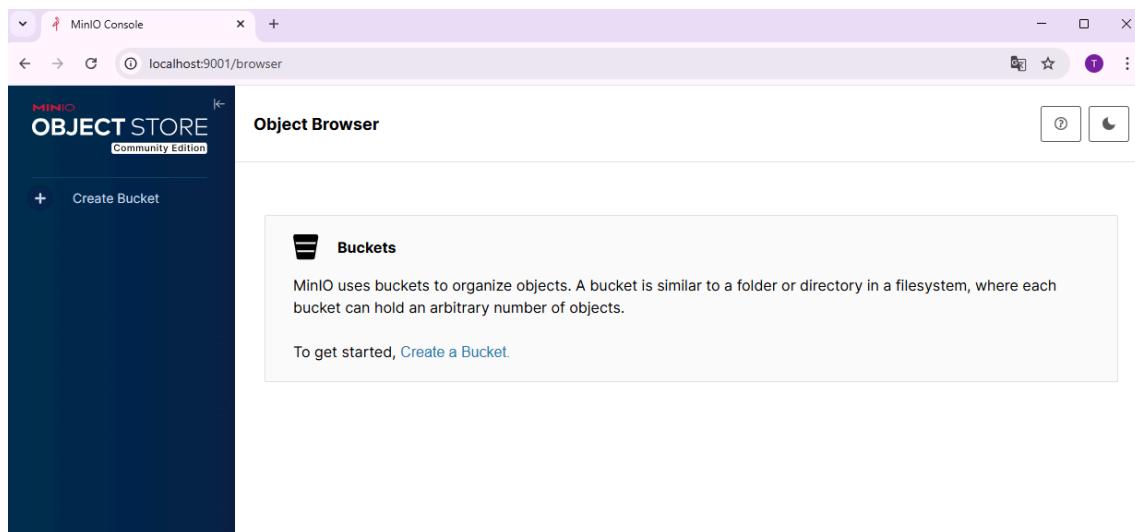
2.2 Acceso a MinIO y carga manual de datos

MinIO es un sistema de almacenamiento distribuido de código abierto. Está basado la tecnología de almacenamiento de objetos y se ha diseñado para su uso en la nube, siendo su interfaz compatible con la de AWS S3, servicio de almacenamiento cloud ampliamente usado y compatible con las herramientas de Big Data.

Una vez arrancado el entorno con Docker, es posible acceder a MinIO a través de tu navegador:

- **url:** <http://localhost:9001>
- **usuario:** minioadmin
- **pass:** minioadmin123

Si no hemos usado el entorno previamente el almacenamiento en MinIO estará vacío, como se muestra en la siguiente captura.

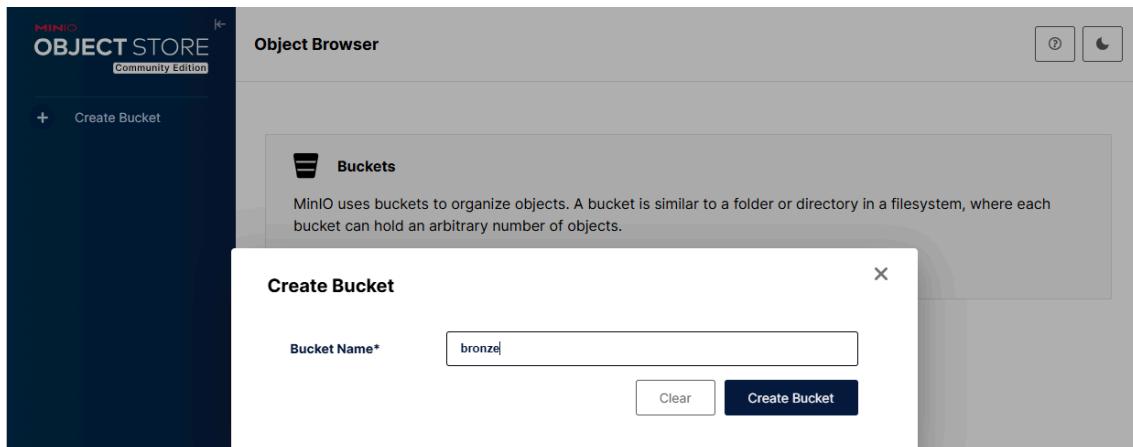


En MinIO y, en general, en sistemas basados en almacenamiento de objetos, las distintas zonas de datos se aíslan en buckets. Un bucket es un contenedor de almacenamiento donde se pueden subir archivos y organizarlos en carpetas.

En nuestro caso vamos a crear **cuatro buckets**:

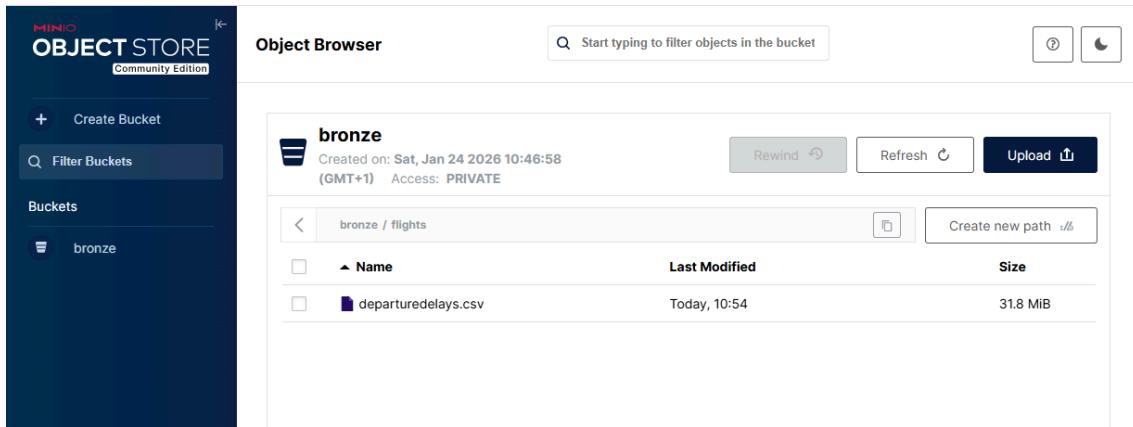
- **bronze:** Contendrá los conjuntos de datos que subamos manualmente, en bruto y sin transformar, y que se usarán como fuentes de datos de las prácticas. Los enlaces para descargar los conjuntos de datos necesarios para cada práctica se irán publicando en el Moodle de la asignatura (sección “*Datasets necesarios para las prácticas*”).
- **silver:** Se usará como destino de los conjuntos de datos procesados a partir de bronze, tras aplicar operaciones básicas de limpieza y estructuración a nivel de conjunto de datos, en formatos optimizados (por ejemplo, Parquet) para su uso con herramientas como Spark.
- **lakehouse:** Representa la capa gold, que almacenará los datos procesados tras aplicar uniones, agregaciones y modelado, de forma que queden listos para su consulta, visualización (por ejemplo, cuadros de mando) o IA / Machine Learning.

La forma más sencilla de crear Buckets en MinIO es usar la interfaz web mediante la opción “+ Create Bucket” (menú de la izquierda).



Una vez creado el primer bucket *bronze* se pueden empezar a subir archivos dentro del mismo. A continuación, se debe cargar el archivo "departuredelays.csv" dentro de una carpeta /flights en dicho bucket. Dentro del bucket *bronze* se hará clic en "Create new path", se introduce el nombre flights y hacemos clic en "Create". Para subir archivos en la nueva carpeta se hará clic en "Upload->Upload File" y se selecciona en nuestro equipo local el archivo "departuredelays.csv".

Si se han seguido correctamente los pasos se visualizará el archivo subido a la carpeta, que a su vez está dentro del bucket *bronze*.



The screenshot shows the MinIO Object Store interface. On the left, there's a sidebar with 'Create Bucket', 'Filter Buckets', and a list of 'Buckets' containing 'bronze'. The main area is titled 'Object Browser' with a search bar 'Start typing to filter objects in the bucket'. It shows a list of objects in the 'bronze' bucket:

Name	Last Modified	Size
departuredelays.csv	Today, 10:54	31.8 MiB

Para esta práctica 0 solo es necesario usar este conjunto de datos en el bucket *bronze* y, además, usaremos el bucket *silver* para la salida de las pruebas. No obstante, se deben crear los otros dos buckets (gold y lakehouse) para dejar el entorno preparado para el resto de las prácticas.

2.2 Acceso a Jupyter Lab y prueba de un Notebook con código Spark

JupyterLab es una herramienta que permite la creación y ejecución de *notebooks*, combinando código, texto y resultados de forma interactiva.

En las prácticas de la asignatura se utilizará JupyterLab como apoyo para:

- La exploración inicial de los conjuntos de datos.
- El desarrollo de una primera versión del código de los procesos de datos con Spark + Python, antes de su orquestación.
- Facilitar el aprendizaje y la depuración del código Spark, gracias a su ejecución interactiva y a la visualización inmediata de resultados.

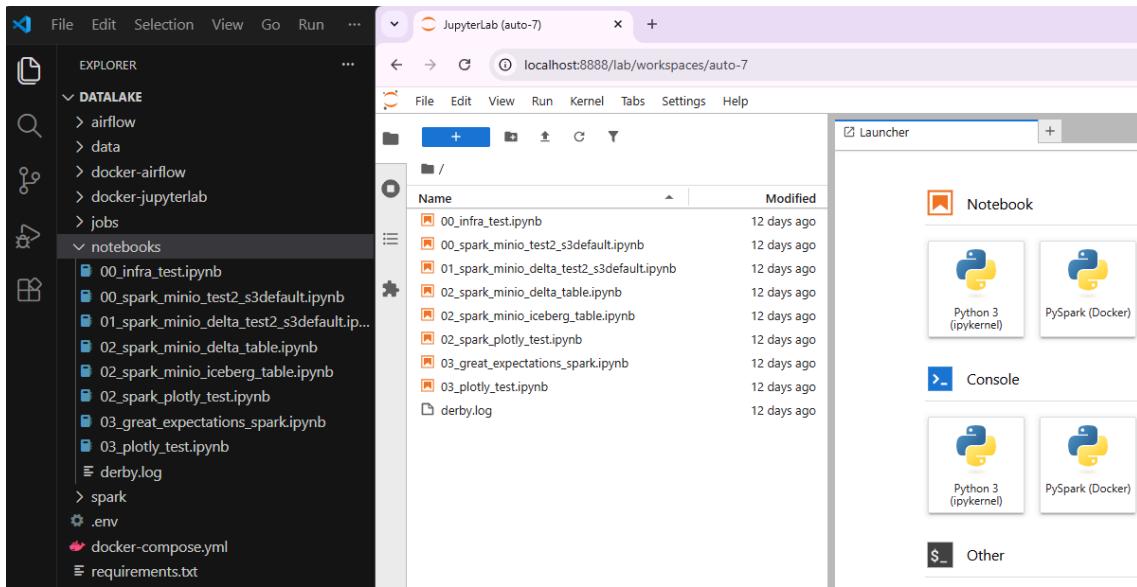
Además, en la Práctica 4, JupyterLab se empleará para la generación de visualizaciones a partir de los datos ya procesados y almacenados en el Data Lake.

Una vez arrancado el entorno con Docker, es posible acceder a Jupyter Lab a través de tu navegador:

- url: <http://localhost:8888/>

En la siguiente captura se muestra el entorno de Jupyter a la derecha y, a la izquierda, la estructura de carpetas desde Visual Studio Code. Se puede observar que los archivos (notebooks) que están disponibles en Jupyter se almacenan en la carpeta local “*/infra/notebooks*”, por tanto, los podemos

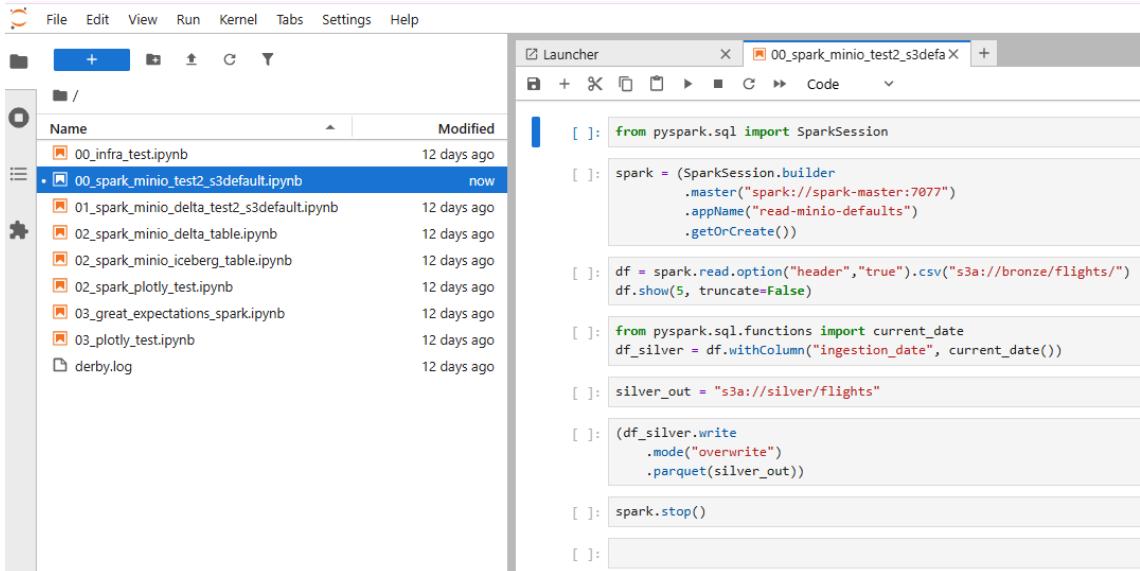
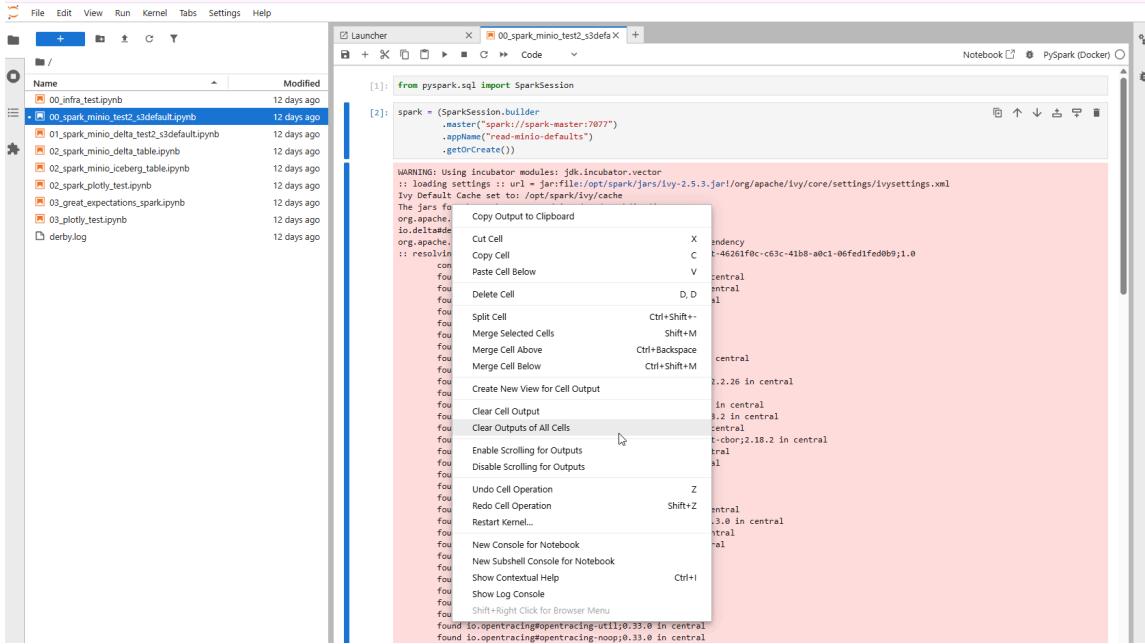
exportar para futuras ejecuciones y evitar que los borren (especialmente si usamos un equipo del laboratorio).



Es posible crear un nuevo notebook seleccionando el icono de “PySpark (Docker)”. Esto permitirá que se use el clúster de Spark desplegado con Docker y no una nueva instancia de Spark dentro del mismo contenedor de Jupyter.

No obstante, en esta práctica 0 el objetivo es abrir el notebook de test “*00_spark_minio_test2_s3default.ipynb*” y ejecutarlo paso a paso. Para abrirlo se hará clic en dicho archivo en el menú de la izquierda en Jupyter y como resultado se verá el código de este a la derecha.

Además, para facilitar una ejecución limpia, se hará clic en con el botón derecho del ratón en cualquier zona del notebook abierto y se hará clic en “Clear Output of All Cells”.



Ahora, se ejecutará el notebook bloque a bloque (también llamados celdas). Para ello, colocamos el cursor para garantizar la selección de la primera celda y ejecutamos pulsando Shift+Enter (o haciendo clic en el ícono ➤). Esto ejecutará el código de la celda y moverá el foco a la siguiente, que podrá ser ejecutada de la misma forma. Repetiremos hasta llegar al final del Notebook.

Si alguna celda genera algún error al ejecutar el código Spark, debes revisarla para que se ejecute con éxito. Si se genera algún warning se debe ignorar para esta práctica.

A continuación, se muestra la salida correcta.

```
[7]: 1 df = spark.read.option("header","true").csv("s3a://bronze/flights/")
2 df.show(5, truncate=False)

26/01/24 10:30:24 WARN MetricsConfig: Cannot locate configuration: tried hadoop-metrics2-s3a-file-system.properties,hadoop-metrics2.properties
SLF4J: Failed to load class "org.slf4j.impl.StaticLoggerBinder".
SLF4J: Defaulting to no-operation (NOP) logger implementation
SLF4J: See http://www.slf4j.org/codes.html#StaticLoggerBinder for further details.

+-----+-----+-----+
|date |delay|distance|origin|destination|
+-----+-----+-----+-----+
|01011245|6 |602 |ABE |ATL |
|01020600|-8 |369 |ABE |DTW |
|01021245|-2 |602 |ABE |ATL |
|01020605|-4 |602 |ABE |ATL |
|01031245|-4 |602 |ABE |ATL |
+-----+-----+-----+-----+
only showing top 5 rows

[8]: 1 from pyspark.sql.functions import current_date
2 df_silver = df.withColumn("ingestion_date", current_date())

[9]: 1 silver_out = "s3a://silver/flights"

[11]: 1 (df_silver.write
2 .mode("overwrite")
3 .parquet(silver_out))

[12]: 1 spark.stop()
```

Para guardar el notebook ejecutado junto con la salida de sus celdas (no las borres como hicimos al inicio), puedes pulsar Ctrl+S o hacer clic en . **Como parte de la entrega de la práctica se ha de incluir una copia del notebook ejecutado**, incluyendo el resultado de cada celda. De esta forma podremos conocer la fecha de ejecución y verificar que es una ejecución única.

2.3 Acceso a Airflow y prueba de un proceso de orquestación

Apache Airflow es la herramienta que usaremos para la orquestación o automatización de la ejecución de los procesos de Spark desarrollados.

Una vez arrancado el entorno con Docker, es posible acceder a Airflow a través de tu navegador:

- **url:** <http://localhost:8080/>
- **usuario:** airflow
- **pass:** airflow

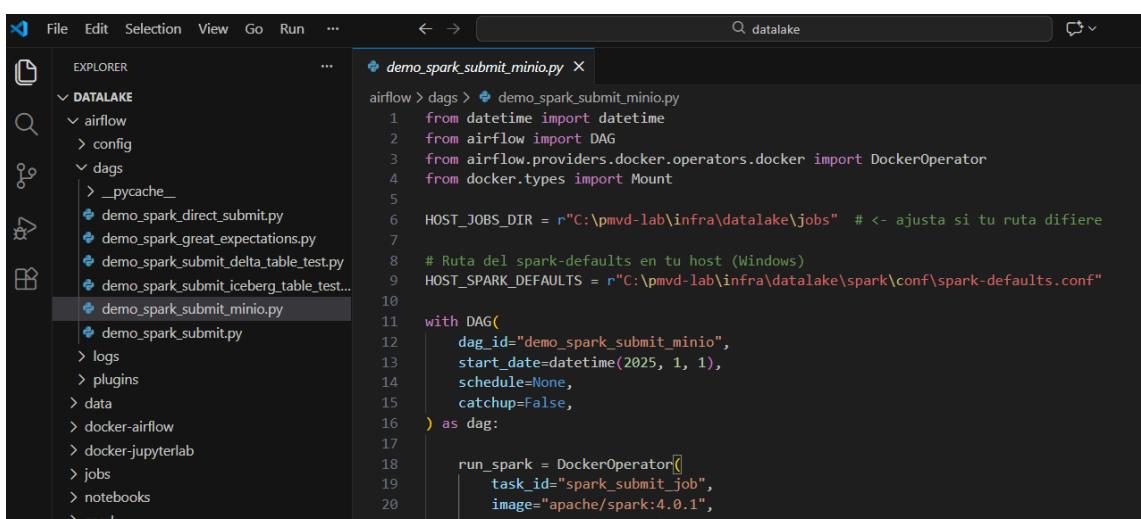
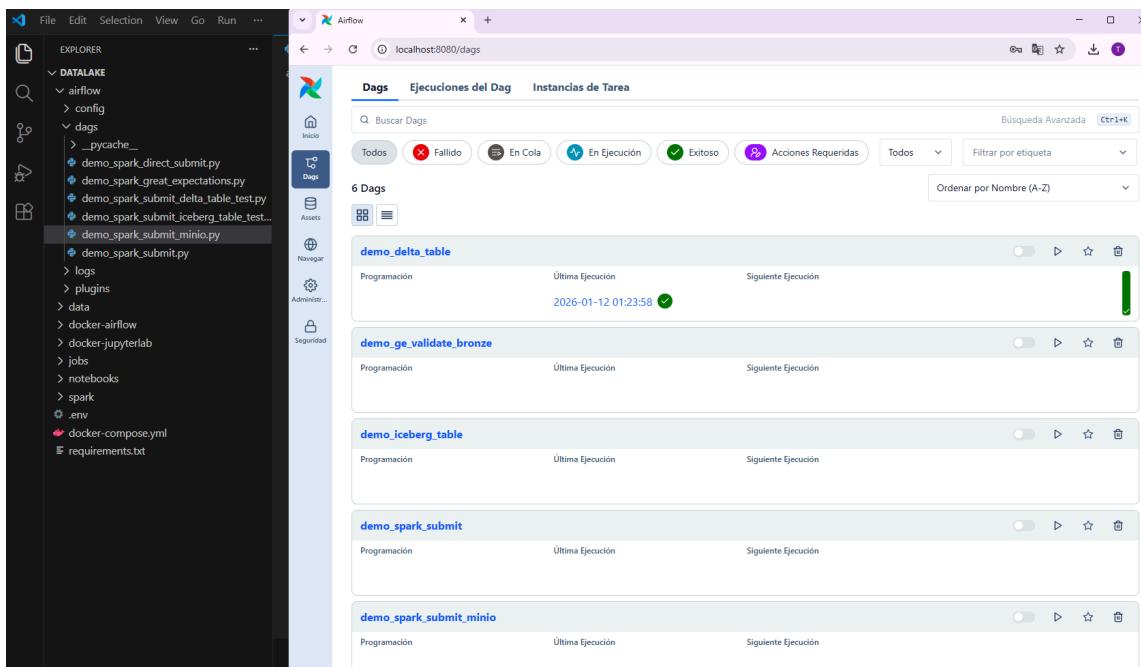
The screenshot shows the Airflow homepage. On the left is a sidebar with icons for Inicio, Dags, Assets, Navegar, Administración, and Seguridad. The main area has a title "Te damos la bienvenida". Below it are three buttons: "Dags Fallidos" (0), "Dags en Ejecución" (0), and "Dags Activos" (1). A section titled "Primeros 10 Dags Favoritos" is shown with a note: "Aún no hay favoritos. Haz clic en el ícono de estrella junto a un Dag en la lista para agregarlo a tus favoritos." To the right is a "Slots del Pool" counter at 128 and a "Gestionar Pools" button. Below these are sections for "Salud" (Base de datos de la metadata, Programador, Triggerer, Procesador de Dags) and "Historial" (Últimas 24 Horas from 2026-01-23 11:52:54 to 2026-01-24 11:52:54). A large green bar indicates 128 active DAGs. On the right, there's a "Eventos de Asset" section with a note: "No se encontraron Eventos de Asset."

En primer lugar, se creará una conexión con clúster de Spark desplegado en otros contenedores de Docker como parte del entorno. Para ello, el menú lateral hacemos clic en “Administración” y dentro, arriba a la derecha, en “Agregar Conexión”. Se hará clic en Tipo de Conexión (puede tardar unos segundos en listarlas la primera vez) y seleccionaremos Spark. En la siguiente captura se muestran el resto de los parámetros a llenar y los valores que se deben introducir.

The screenshot shows the "Agregar Conexión" (Add Connection) dialog box. The left sidebar is identical to the previous one. The dialog box has fields for "ID de la Conexión" (spark_connection) and "Tipo de Conexión" (Spark). A note says: "¿Falta el Tipo de conexión? Asegúrate de haber instalado el paquete de proveedores de Airflow." The "Campos Estándar" section contains fields for "Descripción" (empty), "Host" (spark://spark-master), and "Puerto" (7077). Below these are sections for "Campos Extra" and "Campos Extra (tipo JSON)". At the bottom is a "Guardar" (Save) button.

Para terminar de crearla se hará clic en Guardar, tras lo que se debe mostrar la nueva conexión añadida en el listado de conexiones.

A continuación, se revisará la sección Dags (clic en el menú lateral izquierdo). En esta sección se muestran los procesos de orquestación disponibles y listos para ejecutar, de forma manual programada. Estos procesos se denominan DAG y se crean mediante código Python. Los archivos que contienen el código fuente se almacenan en la carpeta *infra/airflow/dags* y se pueden editar para su desarrollo con Visual Studio Code. Al estar en esa carpeta se sincronizan automáticamente con el servidor de Airflow, donde se ejecutan y monitorizan de forma visual.



```
airflow > dags > demo_spark_submit_minio.py
1  from datetime import datetime
2  from airflow import DAG
3  from airflow.providers.docker.operators.docker import DockerOperator
4  from docker.types import Mount
5
6  HOST_JOBS_DIR = r"C:\pmvd-lab\infra\datalake\jobs" # <- ajusta si tu ruta difiere
7
8  # Ruta del spark-defaults en tu host (Windows)
9  HOST_SPARK_DEFAULTS = r"C:\pmvd-lab\infra\datalake\spark\conf\spark-defaults.conf"
10
11 with DAG(
12     dag_id="demo_spark_submit_minio",
13     start_date=datetime(2025, 1, 1),
14     schedule=None,
15     catchup=False,
16 ) as dag:
17
18     run_spark = DockerOperator(
19         task_id="spark_submit_job",
20         image="apache/spark:4.0.1",
21         command="spark-submit --class com.yoursparkjob.YourSparkJob /path/to/your/spark/app.jar"
22     )
```

En esta práctica 0 se ha de ejecutar el DAG ***demo_spark_submit_minio***. Si se explora el código de este en el correspondiente archivo .py , podemos observar como hace referencia a un proceso de Spark desarrollado con Python "*jobs/spark_job_read_minio_demo.py*".

Este proceso es equivalente al código del Notebook que ejecutamos con Jupyter en el apartado anterior. En las siguientes prácticas veremos como convertir un Notebook de Jupyter a un script de Python .py, siendo necesario para su orquestación desde Airflow.

Estos archivos que definen los proceso de Spark con Python se almacenarán en la carpeta *infra/jobs*. Además, en los procesos de Airflow podrá ser necesario ajustar la siguiente variable en función de dónde tengamos ubicada la carpeta *infra* en nuestro equipo local.

```
# Ajusta si tu ruta difiere  
HOST_JOBS_DIR = r"C:\pmvd-Lab\infra\dataLake\jobs"
```

Tras revisar lo anterior (y ajustar la ruta en caso necesario), desde la interfaz web de Airflow se ejecutará el DAG “*demo_spark_submit_minio*”. Para ello, se hará clic en el botón Play (Activar Dag).



Tras esto seleccionaremos “ejecución única” y haremos clic en “Trigger”.



En este momento, comenzará la ejecución del DAG. Se hará clic en el nombre de este para entrar al detalle las ejecuciones, actuales e histórico. Cada DAG se compone de 1 o más Tasks, que se ejecutarán por separado. Podremos acceder al logs las task para revisar en detalla su ejecución en caso necesario.

The screenshot shows the Airflow web interface. On the left, the sidebar has 'Dags' selected. In the main area, the DAG 'demo_spark_submit_minio' is shown with one task, 'spark_submit_job'. The task card indicates it was last run on 2026-01-24 at 12:35:28, with a duration of 00:01:17. Below the task card is a chart titled 'Última Ejecución de Dag' showing the duration from 00:00:40 to 00:01:20.

Para acceder a los logs de una ejecución, dentro del detalle de ese DAG, hemos de:

- Clic en Ejecuciones para ver el listado de ejecuciones. Haremos clic en la ejecución que nos interese para acceder al detalle de monitorización y logs, en este caso la ejecución en curso o terminada.
- Una vez en la ejecución, seleccionaremos la Task que nos interese. En nuestro caso se llama “spark_submit_job”. Si hacemos clic en la misma, se nos mostrará su log, el cual se actualizará en vivo durante su ejecución.

Como parte de la entrega de la práctica se ha de incluir una copia del log de ejecución generado al término de una ejecución exitosa. Se descargará en formato .txt al hacer clic en .

The screenshot shows the execution details for the task 'spark_submit_job'. The 'Logs' tab is selected, displaying the terminal output of the task's execution. The log output includes several INFO messages related to Docker and Airflow, such as 'Docker container created' and 'Starting docker container from image apache/spark:4.0.1'. A tooltip 'Pulsar d para descargar los registros' is visible over the log area.

2.4 Apagado del entorno

Para terminar la sesión de forma segura, es necesario detener el entorno con el siguiente comando de Docker:

- **docker compose down**

```
PS C:\pmvd-lab\infra\datalake> docker compose down
[+] Running 11/11
✓ Container airflow-scheduler      Removed          2.4s
✓ Container airflow-dagprocessor   Removed          2.8s
✓ Container spark-worker          Removed          1.2s
✓ Container airflow-triggerer    Removed          2.3s
✓ Container jupyter              Removed          2.1s
✓ Container airflow-apiserver    Removed          2.6s
✓ Container minio                Removed          1.2s
✓ Container spark-master         Removed          1.6s
✓ Container datalake-airflow-init-1 Removed        0.1s
✓ Container airflow-postgres     Removed          0.9s
✓ Network datalake             Removed          0.6s
PS C:\pmvd-lab\infra\datalake>
```

3. Entrega de la Práctica

La fecha límite de entrega para esta práctica es el **viernes 13 de febrero**, hasta las **23:59**.

Debes entregar un fichero llamado **practica0.zip** que incluya:

- **00_spark_minio_test2_s3default.ipynb**: Copia ejecutada con éxito por el alumno, incluyendo obligatoriamente la salida de las celdas (incluyen log).
- **[logs_airflow].txt**: Archivo de logs descargado de Airflow. El nombre logs_airflow se sustituirá por el autogenerado por Airflow, por ejemplo, *logs_demo_spark_submit_minio_manual_2026-01-24T11_35_28+00_00_spark_submit_job_-1_1.txt*.
- **alumno.txt**: Archivo de texto incluyendo:
 - Nombre y apellidos
 - DNI
 - Grupo de prácticas

La entrega se realizará **exclusivamente a través de la página de Moodle de la asignatura**.

4. Evaluación

Dado que esta práctica es guiada, se evaluará de la siguiente forma:

- 0 puntos por no entregar en fecha
- 4 puntos por entrega a la que le faltan partes
- 6 puntos por entrega conforme a lo especificado, pero con errores
- 10 puntos por práctica entregada conforme a lo especificado (sin errores)

Anexo 1: Instalación en tu equipo personal

Es posible la instalación del entorno en tu equipo personal, siempre y cumpla los siguientes requisitos:

- [Docker Desktop](#) instalado en la máquina
- Mínimo estimado de 12GB de RAM
- Aunque no es obligatorio, se recomienda también instalar [Visual Studio Code](#) como entorno de desarrollo.

Para ello, se descargará el archivo “Software de la asignatura sobre Docker” desde Moodle. A continuación, descomprimir el archivo zip descargado *infra-PMDV-UA.zip* en tu equipo, obteniendo de esta forma las carpetas ***datalake*** y ***tests***.

A partir de este punto se deben seguir los pasos indicados en el apartado **2.1 Arranque del entorno** y siguientes:

- Mover o copiar las 4 carpetas dentro de *tests* en la carpeta *datalake*.
- Abrir una terminal (ej. en Visual Studio), acceder a la ruta donde se encuentre *datalake* y ejecutar el comando `docker compose up -d`

Cuando se ejecuta el comando anterior la primera vez, se van a descargar y compilar algunas imágenes en Docker, pudiendo llevar hasta 10-20 minutos esa instalación, según las características de tu equipo. Esto solo ocurrirá en la primera ejecución, aunque es posible evitarlo usando la [copia de las imágenes Docker](#) ya generadas que se proporcionan para en versión portable, usando el comando `docker load -i pmvd-images.tar` para cargarlas antes ejecutar el comando `compose`.

Anexo 2: Versión portable USB

Alternativa para casos en los que el PC de laboratorio no tenga instalado correctamente el entorno o se vaya a usar otro equipo cualquiera de forma puntual. Las diferencias principales con la instalación completa previamente descrita en el Anexo 1 son:

- Se cargarán las imágenes de Docker previamente compiladas desde un USB o disco duro externo ([pmvd-images.zip](#) – 3,48 GB comprimido)
- Se usará una carpeta project, que contiene una versión especial de la definición de Docker para que sea más fácil copiarla desde el USB al equipo y viceversa. Ya incluye los ejemplos de test copiados en las correspondientes carpetas ([project.zip](#))
- En el zip anterior se incluye unas instrucciones (.txt) detalladas para su uso.