



UA

**Universidad
de Alicante**

ESCUELA POLITÉCNICA SUPERIOR

Ejercicio de Middleware

Curso académico 2024 / 2025

Semestre 2

Grado en Ingeniería en Inteligencia Artificial

Grupo 1

Profesora: Ricardo Moreno Rodríguez

Alumno: Santiago Álvarez Geanta



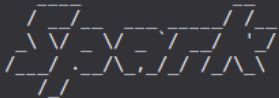
ÍNDICE

Introducción	2
Desarrollo	2
1. Pasos seguidos para construir la solución	2
2. Dificultades encontradas y resolución	3
3. Análisis de la ejecución en Spark UI	4
Conclusión	5

Durante el desarrollo de esta práctica de Computación de Alto Rendimiento, mi objetivo fue construir un flujo de análisis distribuido de eventos técnicos registrados por múltiples brazos robóticos. Para ello, empleé Apache Spark ejecutándose sobre un contenedor Docker, siguiendo los principios de procesamiento distribuido de datos.

Desarrollo

En primer lugar, diseñé cuatro archivos de registros simulados (`eventos_robot1.txt`, `eventos_robot2.txt`, `eventos_robot3.txt`, `eventos_robot4.txt`) dentro de una carpeta llamada `registros/`. Estos archivos contenían eventos realistas tales como **ALERTA_SOBRECALENTAMIENTO**, **REINICIO_SISTEMA**, **PAUSA_NO_PROGRAMADA**, entre otros. Cada archivo representa los sucesos de un brazo robótico distinto.

```
o sant_vz6@santvz6-IdeaPad-Gaming-3-15ACH6:~/Escritorio/UA-2/2 CUATRI/CAR/MIDDLEWARE/app$ doc  
ker run -it --rm -p 4040:4040 -v "$PWD/registros:/app/registros" bitnami/spark:latest /opt/  
bitnami/spark/bin/spark-shell  
spark 19:01:52.74 INFO ==>  
spark 19:01:52.75 INFO ==> Welcome to the Bitnami spark container  
spark 19:01:52.75 INFO ==> Subscribe to project updates by watching https://github.com/bit  
nami/containers  
spark 19:01:52.75 INFO ==> Did you know there are enterprise versions of the Bitnami catal  
og? For enhanced secure software supply chain features, unlimited pulls from Docker, LTS su  
pport, or application customization, see Bitnami Premium or Tanzu Application Catalog. See  
https://www.arrow.com/globalecs/na/vendors/bitnami/ for more information.  
spark 19:01:52.75 INFO ==>  
  
Setting default log level to "WARN".  
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel)  
.  
25/04/25 19:01:57 WARN NativeCodeLoader: Unable to load native-hadoop library for your plat  
form... using builtin-java classes where applicable  
Spark context Web UI available at http://f7f35c23a287:4040  
Spark context available as 'sc' (master = local[*, app id = local-1745607717942).  
Spark session available as 'spark'.  
Welcome to  
  
 version 3.5.5  
  
Using Scala version 2.12.18 (OpenJDK 64-Bit Server VM, Java 17.0.14)  
Type in expressions to have them evaluated.  
Type :help for more information.
```

Dentro del spark-shell, procedí a realizar la lectura de todos los archivos .txt mediante `textFile`, filtré los eventos relevantes con `filter`, transformé cada evento en un par clave-valor (evento, 1) usando `map`, y finalmente apliqué una reducción distribuida con `reduceByKey` para obtener la frecuencia de cada evento.

El código principal fue el siguiente:

```
scala> val registros = spark.sparkContext.textFile("/app/registros/*.txt")
registros: org.apache.spark.rdd.RDD[String] = /app/registros/*.txt MapPartitionsRDD[1] at textFile at <console>:22

scala>

scala> val eventosClave = List(
  | "ALERTA_SOBRECALENTAMIENTO",
  | "REINICIO_SISTEMA",
  | "CALIBRACION_AUTOMATICA",
  | "CARGA_EXCESIVA",
  | "PAUSA_NO_PROGRAMADA"
  | )
eventosClave: List[String] = List(ALERTA_SOBRECALENTAMIENTO, REINICIO_SISTEMA, CALIBRACION_AUTOMATICA, CARGA_EXCESIVA, PAUSA_NO_PROGRAMADA)

scala>

scala> val eventosFiltrados = registros
eventosFiltrados: org.apache.spark.rdd.RDD[String] = /app/registros/*.txt MapPartitionsRDD[1] at textFile at <console>:22

scala> .filter(linea => eventosClave.exists(evento => linea.contains(evento)))
res0: org.apache.spark.rdd.RDD[String] = MapPartitionsRDD[2] at filter at <console>:25

scala> .map(evento => (evento, 1))
res1: org.apache.spark.rdd.RDD[(String, Int)] = MapPartitionsRDD[3] at map at <console>:24

scala> .reduceByKey(_ + _)
res2: org.apache.spark.rdd.RDD[(String, Int)] = ShuffledRDD[4] at reduceByKey at <console>:24
```

Este enfoque permitió procesar de manera distribuida los datos provenientes de los registros y obtener un resumen de la frecuencia de los eventos críticos.

2. Dificultades encontradas y resolución

Durante el montaje del contenedor Docker, inicialmente cometí un error al situarme dentro de la carpeta `registros` y ejecutar el comando de arranque. Esto provocó que Docker no montase los archivos correctamente, resultando en el error **Input Pattern file:/app/registros/*.txt matches 0 files**.

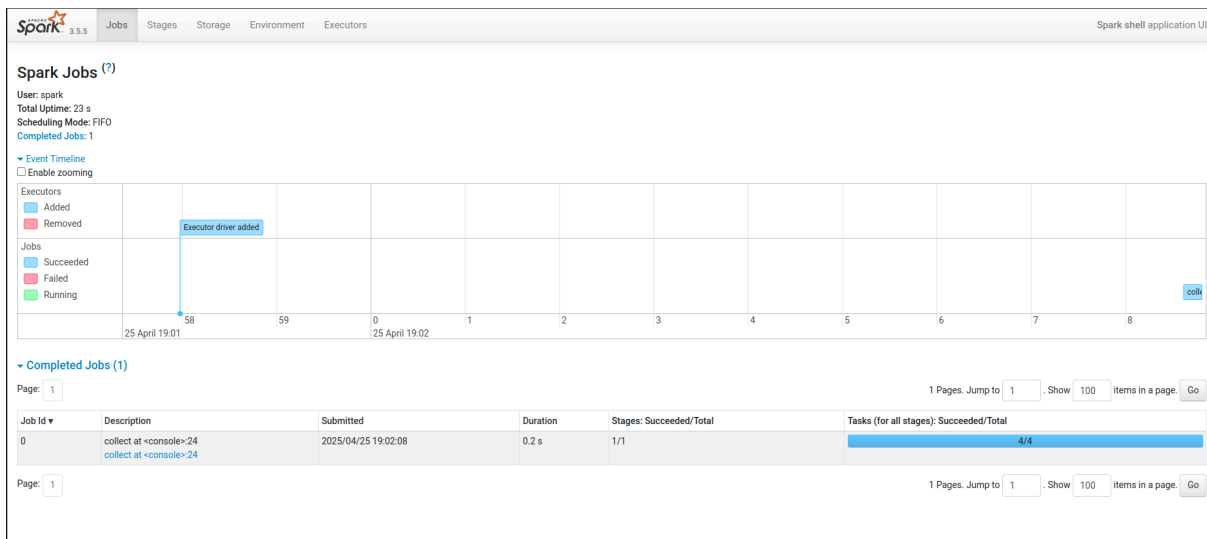
La solución fue volver a la carpeta `app/` (el directorio padre) y lanzar nuevamente el contenedor, corrigiendo así la ruta del volumen. Además, verifiqué dentro del contenedor que los archivos estuvieran correctamente montados ejecutando un pequeño script en Scala que listaba el contenido de `/app/registros`.

Esta experiencia me ayudó a comprender la importancia de la correcta gestión de rutas en entornos virtualizados y la necesidad de validar que los volúmenes estén accesibles.

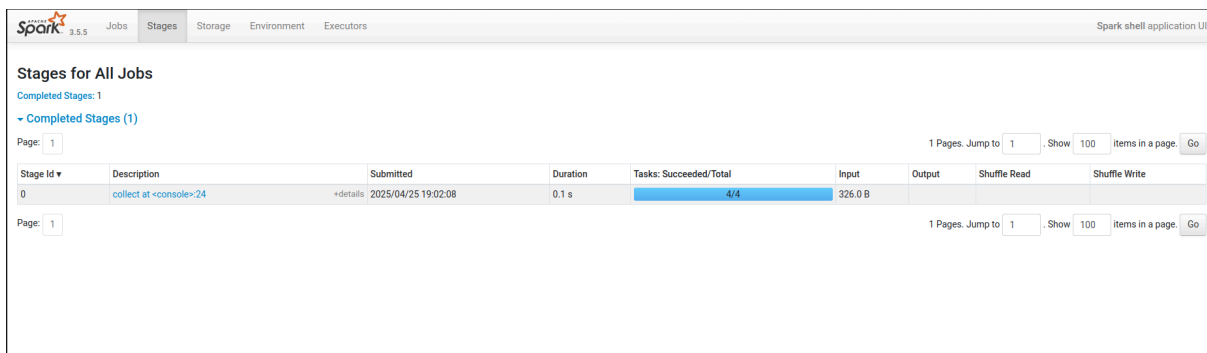
3. Análisis de la ejecución en Spark UI

Accedí a la interfaz web de Spark en **http://localhost:4040** para observar cómo se distribuyó el procesamiento.

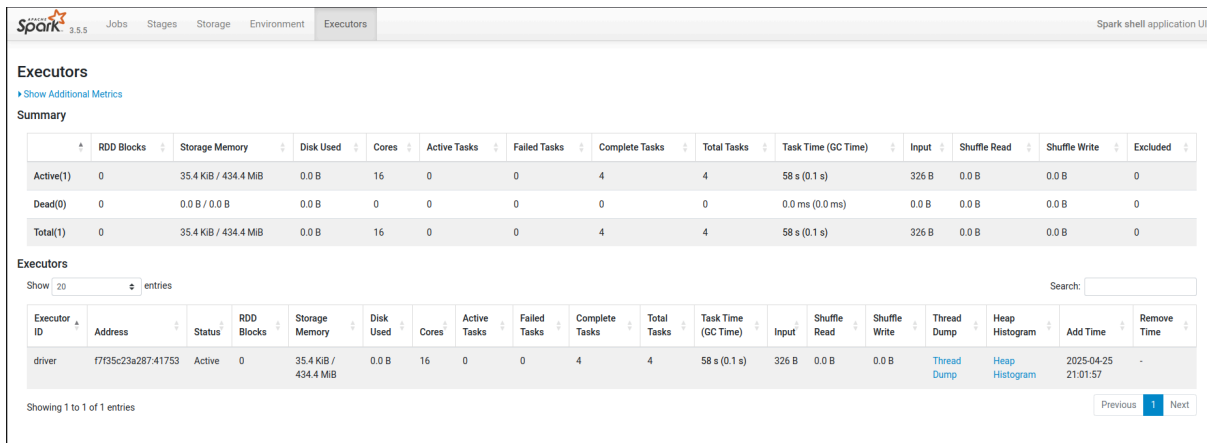
En la pestaña **Jobs**, pude identificar claramente el Job correspondiente al procesamiento de los registros. El job incluía las tareas de lectura, filtrado, mapeo y reducción, agrupadas en una única ejecución debido al encadenamiento de operaciones transformadoras.



En la pestaña **Stages**, se evidenciaron las etapas generadas por Spark. Cada stage correspondía a una transformación importante en el flujo de datos. Observé que se produjo un **shuffle** durante la operación **reduceByKey**, lo cual es lógico ya que esta operación requiere re-agrupamiento de datos por clave entre particiones.



Opcionalmente, también revisé la pestaña **Executors**, donde se mostraba el uso de recursos. Aunque en mi caso sólo había un ejecutor principal (ya que la ejecución fue local), me permitió visualizar estadísticas como número de tareas completadas y cantidad de memoria usada.



The screenshot shows the Spark UI 'Executors' tab. It includes a 'Summary' section with a table of executor statistics and a main 'Executors' table listing individual executors.

Summary

	RDD Blocks	Storage Memory	Disk Used	Cores	Active Tasks	Failed Tasks	Complete Tasks	Total Tasks	Task Time (GC Time)	Input	Shuffle Read	Shuffle Write	Excluded
Active(1)	0	35.4 KiB / 434.4 MiB	0.0 B	16	0	0	4	4	58 s (0.1 s)	326 B	0.0 B	0.0 B	0
Dead(0)	0	0.0 B / 0.0 B	0.0 B	0	0	0	0	0	0.0 ms (0.0 ms)	0.0 B	0.0 B	0.0 B	0
Total(1)	0	35.4 KiB / 434.4 MiB	0.0 B	16	0	0	4	4	58 s (0.1 s)	326 B	0.0 B	0.0 B	0

Executors

Show 20 entries

Executor ID	Address	Status	RDD Blocks	Storage Memory	Disk Used	Cores	Active Tasks	Failed Tasks	Complete Tasks	Total Tasks	Task Time (GC Time)	Input	Shuffle Read	Shuffle Write	Thread Dump	Heap Histogram	Add Time	Remove Time
driver	f7f35c23a28741753	Active	0	35.4 KiB / 434.4 MiB	0.0 B	16	0	0	4	4	58 s (0.1 s)	326 B	0.0 B	0.0 B	Thread Dump	Heap Histogram	2025-04-25 21:01:57	-

Showing 1 to 1 of 1 entries

La interfaz gráfica de Spark fue esencial para entender visualmente la secuencia de operaciones y cómo Spark gestiona los recursos en segundo plano.

Conclusión

Realizar esta práctica me permitió consolidar mis conocimientos sobre procesamiento distribuido y el middleware Spark.

Aprendí cómo Spark divide los datos en particiones y ejecuta operaciones de transformación de forma perezosa, lanzándolas efectivamente sólo tras una acción como `collect()`. También comprendí mejor la diferencia entre transformaciones que requieren shuffle y aquellas que no lo requieren, lo cual tiene un impacto directo en el rendimiento del sistema.

Además, la experiencia me permitió familiarizarme con aspectos prácticos de Docker y el manejo de volúmenes, algo que en contextos reales de computación distribuida es crucial para garantizar que los datos estén correctamente accesibles por los nodos de procesamiento.

En conclusión, esta práctica no sólo reforzó los conceptos de alto rendimiento y procesamiento paralelo, sino que también me dio herramientas prácticas para desplegar y analizar aplicaciones distribuidas en entornos controlados y reproducibles.