

Tarea evaluable

CE_5075 8.1

Apartado 1

Big data aplicado



Índice

APARTADO 1	2
1 - CONFIGURACIÓN DEL CLÚSTER	2
2 - QUERIES	7

APARTADO 1

Vas a trabajar con un dataset que contiene todos los partidos de fútbol entre selecciones nacionales desde 1872 hasta la actualidad.

De los cuatro archivos, solo te interesan dos:

1. **results.csv**: contiene la información general de todos los partidos (equipos, marcador, campeonato, sede, etc.).
2. **goalscorers.csv**: contiene detalles de los goles anotados (partido, equipo, jugador, minuto, si fue penal o en propia portería).

Debes:

- Usar **Azure Databricks**.
- Crear un **cuaderno (notebook)** que responda a **las mismas preguntas** que en la tarea anterior (con Hive).
- Puedes usar **Python con Spark SQL** o directamente **SQL**.
- Asegúrate de **cargar, preparar y analizar los datos correctamente** para responder a cada pregunta.

1 - CONFIGURACIÓN DEL CLÚSTER

Empezamos configurando el clúster en Databricks con las instrucciones aprendidas en los apuntes..

Primero accedemos al Azure Portal y vamos a buscar “databricks” en el marketplace. Ajustamos y rellenamos los datos necesarios y creamos el recurso. Luego pulsamos “acceder al recurso” y podemos iniciar el área de trabajo de databricks.

Microsoft Azure

Inicio > Crear un recurso >

Marketplace

...

Comenzar

Proveedores de servicios

Administración

Marketplace privado

Administración de ofertas privadas

Mi Marketplace

Favoritos

Mis soluciones

Creado recientemente

Planes privados

Categorías

databricks

☐ Solo servicios de Azure

Mostrando 1/20 de un total de 77 res

Azure Databricks

Microsoft

Azure Service

Azure Databricks is the fast, easy and collaborative Apache Spark-based analytics platform.

Crear

Microsoft Azure

Buscar recursos, servicios y documentos (G+/)

Inicio > Crear un recurso > Marketplace >

Creación de un área de trabajo de Azure Databricks

...

Datos básicos

Redes

Cifrado

Security & compliance

Etiquetas

Revisar y crear

Detalles del proyecto

Seleccione la suscripción para administrar recursos implementados y los costes. Use los grupos de recursos como carpetas para organizar y administrar todos los recursos.

Suscripción *

Azure for Students

Grupo de recursos *

(Nuevo) apartado_1_databricks

Crear nuevo

Detalles de instancia

Nombre del área de trabajo *

apartado_1_databricks

Región *

West Europe

Plan de tarifa *

Standard (Apache Spark, Secure with Microsoft Entra ID)

Nombre del grupo de recursos administrados

apartado_1_databricks_group

3

Lo siguiente va a ser crear un clúster de Spark. El clúster solamente dispondrá de un nodo ya que es lo que nos permite nuestra suscripción de Azure for Students. Pulsando en “Nuevo” se abren una serie de opciones, y seleccionamos “Clúster”.

Al seguir la configuración de los apuntes me he encontrado con el siguiente problema:

Parece que no es posible usar “Standard_DS3_v2” con localización en West Europe. Por lo tanto he decidido usar “Standard_D4s_v3” con 16 Gb de memoria y 4 núcleos.

The screenshot shows the Databricks cluster configuration interface. At the top, it says 'Cómputo > Formulario simple: DESACTIVADO'. The cluster name is 'David Ramírez's Cluster'. Below the name, there are tabs for 'Configuración', 'Cuadernos (0)', 'Bibliotecas', 'Log de eventos', 'IU de Spark', 'Logs del driver', 'Métricas', 'Aplicaciones', and 'IU de cómputo Spark - Máster'. The 'Configuración' tab is active. On the left, there are sections for 'Multi-nodo' and 'Nodo único' (selected), 'Modo de acceso' (set to 'Usuario único'), 'Rendimiento' (set to '15.4 LTS (includes Apache Spark 3.5.0, Scala 2.12)' and 'Utilizar aceleración Photon' checked), 'Tipo de nodo' (set to 'Standard_D4s_v3' with '16 GB de memoria y 4 núcleos'), and 'Etiquetas'. On the right, a 'Resumen' box shows: '1 driver', '16 GB de memoria y 4 núcleos', 'Runtime: 15.4.x-scala2.12', and 'Photon Standard_D4s_v3 1.5 DBU/h'. At the bottom, there are 'Terminar' and 'Editar' buttons.

The screenshot shows the 'Cómputo' (Compute) page in Databricks. It has tabs for 'Cómputo interactivo', 'Cómputo de trabajos', 'Pools', and 'Aplicaciones'. Below the tabs, there is a search bar 'Filtre los cómputos a los que tiene acceso', a 'Creador' dropdown, and a 'Solo los anclados' checkbox. A 'Crear cómputo' button is on the right. Below this is a table with the following columns: 'Estado', 'Nombre', 'Runtime', 'Memoria acti...', 'Núcleos acti...', 'DBU/h activo', 'Fuente', 'Creador', and 'Cuadernos'. The table contains one entry: 'David Ramírez's Cluster' with a green status dot, runtime '15.4', memory '16 GB', 4 cores, 1.5 DBU/h, and source 'UI'. The creator is 'David Ramírez'.

Ahora sí, el clúster está preparado, vamos a añadir los datos.

The screenshot shows the 'Añadir datos' (Add data) page in Databricks. It has a sidebar with navigation links: 'Nuevo', 'Workspace', 'Recientes', 'Catálogo', 'Flujos de trabajo', 'Cómputo', 'Ingeniería de datos', 'Ejecuciones', 'Machine Learning', 'Zona de pruebas', 'Experimentos', 'Características', 'Modelos', and 'Servicio'. The main content area has a search bar 'Buscar datos, cuadernos, recientes y más...'. Below it, there is a section 'Añadir datos' with the text 'Dé sus primeros pasos conectándose a una fuente de datos o cargando un archivo local.' and a 'Buscar fuentes de datos' dropdown. Below this is a section 'Conectores de Databricks' with a 'Archivos' subsection containing a 'Crear o modificar tabla' button and the text 'Cargar archivos de datos tabulares para crear una nueva tabla o reemplazar una existente'. At the bottom, there is a section 'Conectores de Fivetran' with the text 'Ver todos los socios de ingesta disponibles en Partner Connect'. Below this are four boxes for 'OneDrive', 'Google Drive', 'Jira', and 'GitHub'. At the very bottom, there is a section 'Productos heredados'.

2 - QUERIES

Una vez tenemos todo configurado creamos un cuaderno donde ejecutaremos las consultas de los datos. Pulsamos “Nuevo” y después “Notebook”



Ahora vamos a cargar los datos en los data frames usando Spark

```
▶ ✓ Ahora mismo (5 s) 1
df_results = spark.read.table('hive_metastore.default.results')
df_goalscorers = spark.read.table('hive_metastore.default.goalscorers')
▶ df_goalscorers: pyspark.sql.dataframe.DataFrame = [date: date, home_team: string ... 6 campos adicionales]
▶ df_results: pyspark.sql.dataframe.DataFrame = [date: date, home_team: string ... 7 campos adicionales]

▶ ✓ Ahora mismo (<1 s) 2
df_results.head()
▶ (1) trabajos de Spark
Row(date=datetime.date(1872, 11, 30), home_team='Scotland', away_team='England', home_score=0, away_score=0, tournament='Friendly', ci
ty='Glasgow', country='Scotland', neutral=False)

▶ ✓ Ahora mismo (<1 s) 3 Python
df_goalscorers.head()
▶ (1) trabajos de Spark
Row(date=datetime.date(1916, 7, 2), home_team='Chile', away_team='Uruguay', team='Uruguay', scorer='José Piendibene', minute='44', own
_goal=False, penalty=False)
```


1. Número de goles que ha marcado Lionel Messi (sin contar goles en propia puerta).

```
Ahora mismo (2 s) 5

df_goalscorers.createOrReplaceTempView('goalscorers')
spark.sql('SELECT COUNT(*) AS goals_from_messi FROM goalscorers WHERE scorer = "Lionel Messi" AND own_goal = "FALSE"').show()

(2) trabajos de Spark

+-----+
|goals_from_messi|
+-----+
|                55|
+-----+
```

2. Listado de los 5 partidos más recientes que ha jugado la selección española.

```
Hace 3 minutos (1 s) 7

df_results.createOrReplaceTempView("goalscorers")
spark.sql("SELECT * FROM goalscorers WHERE home_team = 'Spain' OR away_team = 'Spain' ORDER BY date DESC LIMIT 5").show()

(1) trabajos de Spark

+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|date|home_team|away_team|home_score|away_score|tournament|city|country|neutral|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|2025-03-23|Spain|Netherlands|3|3|UEFA Nations League|Valencia|Spain|false|
|2025-03-20|Netherlands|Spain|2|2|UEFA Nations League|Rotterdam|Netherlands|false|
|2024-11-18|Spain|Switzerland|3|2|UEFA Nations League|Santa Cruz de Ten...|Spain|false|
|2024-11-15|Denmark|Spain|1|2|UEFA Nations League|Copenhagen|Denmark|false|
|2024-10-15|Spain|Serbia|3|0|UEFA Nations League|Cordoba|Spain|false|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+
```

3. Número de goles que ha marcado España en toda su historia. Esta información debe extraerse del archivo results.csv, ya que goalscorers.csv no contiene todos los goles.

```
Ahora mismo (1 s) 9

df_results.createOrReplaceTempView('results')
spark.sql(''
SELECT
SUM(CASE WHEN home_team = 'Spain' THEN home_score ELSE 0 END) +
SUM(CASE WHEN away_team = 'Spain' THEN away_score ELSE 0 END) AS total_goals_from_spain
FROM results
'').show()

(2) trabajos de Spark

+-----+
|total_goals_from_spain|
+-----+
|                1567|
+-----+
```

4. Listado de los 5 máximos goleadores con la selección española (sin contar goles en propia puerta).

```
Ahora mismo (1 s)

df_goalscorers.createOrReplaceTempView('goalscorers')
spark.sql('''
    SELECT scorer, COUNT(*) AS goals
    FROM goalscorers
    WHERE team = 'Spain' AND own_goal = 'FALSE'
    GROUP BY scorer
    ORDER BY goals DESC
    LIMIT 5
''').show()

(2) trabajos de Spark

+-----+-----+
| scorer|goals|
+-----+-----+
| David Villa| 41|
| Raúl| 32|
| Álvaro Morata| 29|
| Fernando Torres| 28|
| Fernando Hierro| 25|
+-----+-----+
```

5. Listado de los jugadores españoles que han marcado algún gol de penalti en una Eurocopa (UEFA Euro), ordenados alfabéticamente.

```
Ahora mismo (2 s) 13

df_goalscorers.createOrReplaceTempView('goalscorers')
df_results.createOrReplaceTempView('results')
spark.sql('''
    SELECT DISTINCT scorer
    FROM results r
    JOIN goalscorers g
    ON r.date = g.date AND r.home_team = g.home_team AND r.away_team = g.away_team
    WHERE g.team = 'Spain' AND g.penalty = 'TRUE' AND r.tournament LIKE '%Euro%'
    ORDER BY scorer
''').show()

(3) trabajos de Spark

+-----+
| scorer|
+-----+
| Andrés Iniesta|
| Daniel Ruiz|
| David Villa|
| Fernando Hierro|
| Francisco José Ca...|
| Gaizka Mendieta|
| José Claramunt|
| Juan Antonio Señor|
| Michel|
| Pirri|
| Sergio Ramos|
| Xabi Alonso|
| Álvaro Morata|
+-----+
```

6. Listado de los 5 máximos goleadores de las fases finales de los mundiales (FIFA World Cup) (sin contar goles en propia puerta).

```
▶ ✓ Ahora mismo (1 s) 15

df_goalscorers.createOrReplaceTempView('goalscorers')
df_results.createOrReplaceTempView('results')
spark.sql('''
    SELECT scorer, COUNT(*) AS goals
    FROM results r
    JOIN goalscorers g
    ON r.date = g.date AND r.home_team = g.home_team AND r.away_team = g.away_team
    WHERE r.tournament = 'FIFA World Cup' AND g.own_goal = 'FALSE'
    GROUP BY scorer
    ORDER BY goals DESC
    LIMIT 5
''').show()

▶ (3) trabajos de Spark

+-----+-----+
|      scorer|goals|
+-----+-----+
|Miroslav Klose|   16|
|      Ronaldo|   15|
|   Gerd Müller|   14|
|  Lionel Messi|   13|
|   Just Fontaine|  13|
+-----+-----+
```