Práctica PADRON

Práctica Hive + Impala + HDFS + Spark

- A partir de los datos (CSV) de Padrón de Madrid (https://datos.madrid.es/egob/catalogo/200076-1-padron.csv) llevar a cabo lo siguiente:
- 1- Creación de tablas en formato texto.
 - 1.1) Crear Base de datos "datos_padron" .

```
hive> CREATE DATABASE datos_padron
> ;
OK
Time taken: 1.627 seconds
```

• 1.2) Crear la tabla de datos padron_txt con todos los campos del fichero CSV y cargar los datos mediante el comando LOAD DATA LOCAL INPATH. La tabla tendrá formato texto y tendrá como delimitador de campo el caracter ';' y los campos que en el documento original están encerrados en comillas dobles "" no deben estar envueltos en estos caracteres en la tabla de Hive (es importante indicar esto utilizando el serde de OpenCSV, si no la importación de las variables que hemos indicado como numéricas fracasará ya que al estar envueltos en comillas los toma como strings) y se deberá omitir la cabecera del fichero de datos al crear la tabla.

Creación de la tabla con los parámetros que nos dicen

```
hive> CREATE TABLE padron txt (
   > COD_DISTRITO INT,
> DESC_DISTRITO STRING,
        COD DIST BARRIO INT,
        DESC BARRIO STRING,
        COD_BARRIO INT,
        COD_DIST_SECCION INT,
        COD_SECCION INT,
COD_EDAD_INT INT,
        ESPANOLESHOMBRES INT,
        ESPANOLESMUJERES INT,
        EXTRANJEROSHOMBRES INT,
         EXTRANJEROSMUJERES INT, FX CARGA TIMESTAMP,
         FX DATOS INI DATE,
         FX DATOS FIN DATE
    > )
    > ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde'
    > WITH SERDEPROPERTIES (
    > "separatorChar" = "\;",
          "quoteChar" = "\""
          "escapeChar" = "\\"
    > )
    > STORED AS TEXTFILE
    > TBLPROPERTIES ("skip.header.line.count"="1");
Time taken: 0.058 seconds
```

Cargamos los datos desde el csv a la tabla

```
hive> LOAD DATA LOCAL INPATH '/home/cloudera/Desktop/padron.csv' INTO TABLE padron_txt;
Loading data to table datos_padron.padron_txt
Table datos_padron.padron_txt stats: [numFiles=1, totalSize=34057160]
OK
Time taken: 0.485 seconds
```

• 1.3) Hacer trim sobre los datos para eliminar los espacios innecesarios guardando la tabla resultado como padron_txt_2. (Este apartado se puede hacer creando la tabla con una sentencia CTAS.)

Create Table As Select (CTAS)

```
hive> CREATE TABLE padron txt 2 AS
    > SELECT
           COD DISTRITO,
    >
           TRIM(DESC DISTRITO) AS DESC DISTRITO,
    >
           COD DIST BARRIO,
           TRIM(DESC BARRIO) AS DESC BARRIO,
           COD BARRIO,
    >
           COD DIST SECCION,
           COD SECCION,
    >
           COD EDAD INT,
           ESPANOLESHOMBRES,
    >
           ESPANOLESMUJERES.
    >
           EXTRANJEROSHOMBRES,
    >
           EXTRANJEROSMUJERES,
           FX CARGA,
    >
    >
           FX DATOS INI,
           FX DATOS FIN
    > FROM padron txt;
Query ID = cloudera 20240730115757 fd3ba867-d836-45eb-9e32-43d03be635f4
Total jobs = 3
Stage-Stage-1: Map: 1 Cumulative CPU: 4.31 sec HDFS Read: 34062895 HDFS Write: 26422298 SUCCESS
Total MapReduce CPU Time Spent: 4 seconds 310 msec
            desc distrito cod dist barrio desc barrio cod barrio
                                                                cod dist seccion
      fx datos ini
                 fx datos fin
rga
Time taken: 19.981 seconds
```

• 1.4) Investigar y entender la diferencia de incluir la palabra LOCAL en el comando LOAD DATA.

La palabra clave LOCAL en el comando LOAD DATA en Hive tiene una importancia significativa en cuanto a la ubicación de los datos que estás cargando en la tabla de Hive.

- El comando LOAD DATA en Hive se utiliza para mover datos desde un archivo en el sistema de archivos al directorio de la tabla en el sistema de archivos de Hadoop (HDFS).

```
LOAD DATA [LOCAL] INPATH 'path' [OVERWRITE] INTO TABLE table_name;
```

 Cuando usas LOAD DATA sin la palabra clave LOCAL, Hive asume que el archivo o directorio especificado está ubicado en el sistema de archivos distribuido de Hadoop (HDFS). En este caso, el archivo debe estar en HDFS, no en el sistema de archivos local del nodo donde se ejecuta Hive.

LOAD DATA INPATH '/user/hive/warehouse/data.csv' INTO TABLE padron_txt;

• 1.5) En este momento te habrás dado cuenta de un aspecto importante, los datos nulos de nuestras tablas vienen representados por un espacio vacío y no por un identificador de nulos comprensible para la tabla. Esto puede ser un problema para el tratamiento posterior de los datos. Podrías solucionar esto creando una nueva tabla utiliando sentencias case when que sustituyan espacios en blanco por 0. Para esto primero comprobaremos que solo hay espacios en blanco en las variables numéricas correspondientes a las últimas 4 variables de nuestra tabla (podemos hacerlo con alguna sentencia de HiveQL) y luego aplicaremos las sentencias case when para sustituir por 0 los espacios en blanco. (Pista: es útil darse cuenta de que un espacio vacío es un campo con longitud 0). Haz esto solo para la tabla padron_txt.

Como se puede ver, por la manera en la que se han tratado, no hay espacios en blanco.

• 1.6) Una manera tremendamente potente de solucionar todos los problemas previos (tanto las comillas como los campos vacíos que no son catalogados como null y los espacios innecesarios) es utilizar expresiones regulares (regex) que nos proporciona OpenCSV.

Para ello utilizamos:

ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.RegexSerDe' WITH SERDEPROPERTIES ('input.regex'='XXXXXXX')

Donde XXXXXX representa una expresión regular que debes completar y que identifique el formato exacto con el que debemos interpretar cada una de las filas de nuestro CSV de entrada. Para ello puede ser útil el portal "regex101". Utiliza este método para crear de nuevo la tabla padron_txt_2.

La inserción de datos, los metemos omitiendo el header, par que no de NULL los campos del header, y luego lo volvemos a añadir los campos manual:

```
hive> LOAD DATA LOCAL INPATH '/home/cloudera/Desktop/padron.csv' INTO TABLE padron_txt_2;
Loading data to table default.padron_txt_2
Table default.padron_txt_2 stats: [numFiles=1, totalSize=34056941]
OK
Time taken: 0.475 seconds
```

```
| Padron txt 2.cod distrito | Padron txt 2.desc distrito | Padron txt 2.cod dist barrio | Padron txt 2.cod dist seccion | Padron txt 2.cod seccion | Padron txt 2.cod edad int | Padron txt 2.cod edad
```

Una vez finalizados todos estos apartados deberíamos tener una tabla padron_txt que conserve los espacios innecesarios, no tenga comillas envolviendo los campos y los campos nulos sean tratados como valor 0 y otra tabla padron_txt_2 y padron2 sin espacios innecesarios, sin comillas envolviendo los campos y con los campos nulos como valor 0. Idealmente esta tabla ha sido creada con las regex de OpenCSV.

A priori y por lo que he comprobado, lo tengo listo para trabajar

2- Investigamos el formato columnar parquet.

• 2.1) ¿Qué es CTAS?

CTAS significa "Create Table As Select" y es una técnica en SQL (y otras bases de datos) utilizada para crear una nueva tabla basada en los resultados de una consulta SELECT. Este método es muy útil para copiar datos de una tabla a otra o para crear una nueva tabla con una estructura y datos específicos sin tener que definir manualmente cada columna y tipo de datos.

 2.2) Crear tabla Hive padron_parquet (cuyos datos serán almacenados en el formato columnar parquet) a través de la tabla padron_txt mediante un CTAS.

```
hive> CREATE TABLE padron_parquet
        STORED AS PAROUET
        SELECT
        COD DISTRITO,
       DESC DISTRITO,
         COD DIST BARRIO,
   >
   >
         DESC BARRIO,
         COD BARRIO,
         COD DIST SECCION,
        COD SECCION,
       COD EDAD INT,
       ESPANOLESHOMBRES,
       ESPANOLESMUJERES,
       EXTRANJEROSHOMBRES,
       EXTRANJEROSMUJERES,
       FX CARGA,
        FX DATOS INI,
         FX DATOS FIN
         FROM padron txt;
Query ID = cloudera 20240731091414 728a4d71-9bdf-4efa-9361-59e51971002b
Total jobs = 3
```

Comprobamos que la creación ha ido bien:

```
hive> select * from padron_parquet limit 3;
OK
padron parquet.cod_distrito
padron_parquet.cod_distrito
padron_parquet.cod_distrito
padron_parquet.cod_dist_seccion
padron_parquet.cod_dist_seccion
padron_parquet.cod_dist_seccion
padron_parquet.cod_dist_seccion
padron_parquet.cod_dist_seccion
padron_parquet.cod_dist_seccion
padron_parquet.cod_dist_seccion
padron_parquet.cod_dist_seccion
padron_parquet.expanolesnwij
padron_parquet.fx_carga padron_parquet.fx_datos_ini
padron_parquet.fx_carga padron_parquet.fx_
```

• 2.3) Crear tabla Hive padron_parquet_2 a través de la tabla padron_txt_2 mediante un CTAS. En este punto deberíamos tener 4 tablas, 2 en txt (padron_txt y padron_txt_2, la primera con espacios innecesarios y la segunda sin espacios innecesarios) y otras dos tablas en formato parquet (padron_parquet y padron_parquet_2, la primera con espacios y la segunda sin ellos).

Parquet en Hive no soporta "date", por lo que hay que hacer algún cambio: En vez de date, lo metemos como timestamp el dato.

```
hive> CREATE TABLE padron_parquet_2 (
           COD_DISTRITO INT,
DESC_DISTRITO STRING,
           COD_DIST_BARRIO INT,
DESC BARRIO STRING.
           COD_BARRIO INT,
           COD_DIST_SECCION INT, COD SECCION INT,
           COD_EDAD_INT INT
           ESPANOLESHOMBRES INT.
           ESPANOLESMUJERES INT
           EXTRANJEROSHOMBRES INT, EXTRANJEROSMUJERES INT,
           FX_CARGA TIMESTAMP,
FX_DATOS_INI STRING,
           FX_DATOS_FIN STRING
    > STORED AS PARQUET;
Time taken: 0.087 seconds
    > INSERT INTO padron_parquet_2
    > SELECT
            COD DISTRITO
           TRIM(DESC_DISTRITO) AS DESC_DISTRITO, COD_DIST_BARRIO,
            TRIM(DESC_BARRIO) AS DESC_BARRIO,
            COD BARRIO,
           COD_DIST_SECCION,
           COD SECCION,
            CASE WHEN TRIM(COD EDAD INT) = '' THEN 0 ELSE CAST(TRIM(COD EDAD INT) AS INT) END AS COD EDAD INT,
           ESPANOLESHOMBRES.
           ESPANOLESMUJERES,
            EXTRANJEROSHOMBRES
           EXTRANJEROSMUJERES,
           FX_CARGA,
           FX DATOS INI.
            FX DATOS FIN
     > FROM padron_txt_2;
```

Comprobamos que ha ido bien la creación:

```
hive> select * from padron_parquet_2 limit 3;
OK
padron_parquet_2.cod_distrito padron_parquet_2.desc_distrito padron_parquet_2.cod_dist_barrio
d_barrio padron_parquet_2.cod_dist_seccion padron_parquet_2.cod_seccion padron_pa
adron_parquet_2.espanolesmujeres padron_parquet_2.extranjeroshombres padron_parquet_2.
ron_parquet_2.fx_datos_ini padron_parquet_2.fx_datos_fin
2 ARGANZUELA 202 ACACIAS 2 2094 94 47 6 16 0 0
                                                                                                                                                             padron_parquet_2.desc_barrio
                                                                                                                       eccion padron_parquet_2.cod_edad_int
padron_parquet_2.extranjerosmujeres
                                                                                                                                                                                      padron parquet 2.espanoleshombres p
                                                                                                                                                                                   padron_parquet_2.fx_carga
                                                                                                                                                                                                                                      pad
4-07-01
            VILLA DE VALLECAS 1802 SANTA EUGENIA 2
                                                                                                                         77
                                                                                                                                                             0
                                                                                                                                                                         2
                                                                                                                                                                                      2024-07-07 00:00:01.8166596
4-07-01 2024-07-01
            ARGANZUELA
                               205 DELICIAS
                                                                                                                                                                         2024-07-07 00:00:01.8166596
024-07-01
Time taken: 0.074 seconds, Fetched: 3 row(s)
```

2.4) Opcionalmente también se pueden crear las tablas directamente desde 0 (en lugar de mediante CTAS) en formato parquet igual que lo hicimos para el formato txt incluyendo la sentencia STORED AS PARQUET. Es importante para comparaciones posteriores que la tabla padron_parquet conserve los espacios innecesarios y la tabla padron parquet 2 no los tenga. Dejo a tu elección cómo hacerlo.

La tabla "mala"

> STORED AS PARQUET;

Time taken: 0.049 seconds

```
hive> CREATE TABLE padron parquet (
                      COD DISTRITO STRING,
                      DESC DISTRITO STRING,
                      COD DIST BARRIO STRING,
                      DESC BARRIO STRING,
        > COD_BARRIO STRING,
> COD_BARRIO STRING,
> COD_DIST_SECCION STRING,
> COD_SECCION STRING,
> COD_EDAD_INT STRING,
> ESPANOLESHOMBRES STRING,
> ESPANOLESMUJERES STRING,
> EXTRANJEROSHOMBRES STRING,
> EXTRANJEROSHUJERES STRING,
> EXTRANJEROSMUJERES STRING,
                  FX CARGA STRING,
                      FX DATOS INI STRING,
                      FX DATOS FIN STRING
         >
         > )
         > STORED AS PARQUET;
 0K
 Time taken: 0.085 seconds
hive> INSERT INTO padron parquet SELECT * FROM padron txt;
Comprobación:
hive> select * from padron parquet limit 3;
     _parquet.cod_distrito padron_parquet.desc_distrito padron_parquet.cod_dist_barrio_padron_parquet.desc_barrio padron_parquet.cod_barrio
edad_int_padron_parquet.espanoleshombres_padron_parquet.espanolesmujeres_padron_parquet.extranjeroshombres_padron_parquet.extranjerosmujeres
                                                                                                                       padron_parquet.cod_dist_seccion padron_parquet.cod_sec
padron_parquet.fx_carga padron_parquet.fx_datos_ini
05.11 DESC DISTRITIO NULL DESC BARRIO NULL NULL NULL COD EDAD INT NULL NULL NULL NULL NULL HX DATOS INT FX DATOS FIN 2 "ARGANZUELA " 202 "ACACIAS " 2 2994 54 " 47" 6 16 0 0 NULL "2024-07-01" 18 "VILLA DE VALLECAS " 1802 "SANTA EUGENIA " 2 18029 29 " 77" 7 3 0 2 NULL "2024-07-01" Time taken: 0.047 seconds, Fetched: 3 row(s)
Ahora la tabla "buena"
hive> CREATE TABLE padron parquet 2 (
      > COD DISTRITO STRING,
               DESC_DISTRITO STRING,
               COD_DIST_BARRIO STRING, DESC_BARRIO STRING,
                 COD BARRIO STRING,
                COD_DIST_SECCION_STRING,
                COD SECCION STRING,
                COD_EDAD_INT STRING,
               ESPANOLESHOMBRES STRING,
               ESPANOLESMUJERES STRING,
               EXTRANJEROSHOMBRES STRING,
               EXTRANJEROSMUJERES STRING,
               FX CARGA STRING,
               FX DATOS INI STRING,
                 FX DATOS FIN STRING
       > )
```

hive> INSERT INTO padron parquet 2 SELECT * FROM padron txt 2;

Comprobación:

```
hive select * from padron parquet 2 cloud distrito obraquet 2 limit 3;

OK padron parquet 2 cloud distrito adron parquet 2 cloud distrito padron parquet 2 cloud barro padron padron parquet 2 cloud barro padron padro
```

• 2.5) Investigar en qué consiste el formato columnar parquet y las ventajas de trabajar con este tipo de formatos.

Parquet es un formato de almacenamiento columnar que organiza los datos por columnas en lugar de filas. Esto permite una compresión eficiente y reduce el tamaño del almacenamiento. Su diseño facilita la lectura rápida de solo las columnas necesarias para consultas, mejorando el rendimiento en el procesamiento de datos. Además, Parquet es compatible con herramientas de Big Data como Apache Spark y Hive, y permite la evolución del esquema sin afectar los datos existentes.

• 2.6) Comparar el tamaño de los ficheros de los datos de las tablas padron_txt (txt), padron_txt_2 (txt pero no incluye los espacios innecesarios), padron_parquet y padron_parquet_2 (alojados en hdfs cuya ruta se puede obtener de la propiedad location de cada tabla por ejemplo haciendo "show create table").

```
Para ello se pillará de cada tabla su LOCATION, accediendo a ellas con el siguente comando:

hive> SHOW CREATE TABLE PADRON_TXT;

LOCATION
  'hdfs://quickstart.cloudera:8020/user/hive/warehouse/padron_txt'

LOCATION
  'hdfs://quickstart.cloudera:8020/user/hive/warehouse/padron_txt_2'

LOCATION
  'hdfs://quickstart.cloudera:8020/user/hive/warehouse/padron_parquet'

LOCATION
  'hdfs://quickstart.cloudera:8020/user/hive/warehouse/padron_parquet'

LOCATION
  'hdfs://quickstart.cloudera:8020/user/hive/warehouse/padron_parquet_2'
```

```
[cloudera@quickstart ~]$ hdfs dfs -du -h /user/hive/warehouse/padron_txt;
32.5 M 32.5 M /user/hive/warehouse/padron_txt/padron.csv
[cloudera@quickstart ~]$ hdfs dfs -du -h /user/hive/warehouse/padron_txt_2;
32.5 M 32.5 M /user/hive/warehouse/padron_txt_2/padron.csv
[cloudera@quickstart ~]$ hdfs dfs -du -h /user/hive/warehouse/padron_parquet;
2.4 M 2.4 M /user/hive/warehouse/padron_parquet/000000_0
[cloudera@quickstart ~]$ hdfs dfs -du -h /user/hive/warehouse/padron_parquet_2;
2.4 M 2.4 M /user/hive/warehouse/padron_parquet_2/0000000_0
```

Se puede ver que Parquet ocupa 30 M menos que txt, lo que es sensacional para trabajar con datos grandes.

3- Juguemos con Impala.

• 3.1) ¿Qué es Impala?

Impala es un motor de consultas SQL de código abierto para Hadoop, diseñado para ofrecer consultas rápidas y analíticas en grandes volúmenes de datos. Lee datos directamente desde HDFS y Apache HBase, usa procesamiento en memoria para mejorar la velocidad, y es compatible con SQL estándar.

Es ideal para análisis en tiempo real, escalabilidad horizontal y se integra con herramientas de BI para crear dashboards y reportes.

• 3.2) ¿En qué se diferencia de Hive?

Impala y Hive son motores SQL para Hadoop con diferencias notables. Impala es más rápido y ofrece consultas interactivas en tiempo real al procesar datos en memoria, mientras que Hive, que usa MapReduce o Tez, puede ser más lento debido a la escritura en disco. Impala también se integra bien con HBase, permitiendo consultas sobre datos NoSQL, mientras que Hive se usa principalmente con HDFS. Impala está optimizado para consultas rápidas y complejas, mientras que Hive es más adecuado para procesamiento por lotes y tareas ETL.

• 3.3) Comando INVALIDATE METADATA, ¿en qué consiste?

Se utiliza para forzar la actualización de la caché de metadatos de la base de datos. Esto es necesario cuando los metadatos de las tablas o particiones han cambiado fuera de Impala, como al agregar, eliminar o modificar archivos de datos. Ejecutar este comando garantiza que Impala lea los metadatos actualizados y refleje los cambios en las consultas posteriores.

• 3.4) Hacer invalidate metadata en Impala de la base de datos datos_padron.

Entro en Impala:

```
[cloudera@quickstart ~]$ impala-shell
Starting Impala Shell without Kerberos authentication
Connected to quickstart.cloudera:21000
Server version: impalad version 2.10.0-cdh5.13.0 RELEASE (build 2511805fleaa991df1460276c7e9f19d819cd4e4)
Welcome to the Impala shell.
(Impala Shell v2.10.0-cdh5.13.0 (2511805) built on Wed Oct 4 10:55:37 PDT 2017)
Want to know what version of Impala you're connected to? Run the VERSION command to
find out!
              [quickstart.cloudera:21000] > INVALIDATE METADATA datos_padron.padron_txt;
Query: invalidate METADATA datos_padron.padron_txt
Query submitted at: 2024-07-31 12:51:39 (Coordinator: http://quickstart.cloudera:25000)
Query progress can be monitored at: http://quickstart.cloudera:25000/query_plan?query_id=e94f811570d74aff:87a7489c00000000
Fetched 0 row(s) in 0.17s
[quickstart.cloudera:21000] > INVALIDATE METADATA datos padron.padron txt 2;
Query: invalidate METADATA datos_padron.padron_txt_2
Query submitted at: 2024-07-31 12:52:03 (Coordinator: http://quickstart.cloudera:25000)
Query progress can be monitored at: http://quickstart.cloudera:25000/query_plan?query_id=4245a405523548f1:5500b64b00000000
Fetched 0 row(s) in 0.12s
```

• 3.5) 3.6) Calcular el total de EspanolesHombres, espanolesMujeres, ExtranjerosHombres y ExtranjerosMujeres agrupado por DESC_DISTRITO y DESC_BARRIO.

Llevar a cabo las consultas en Hive en las tablas padron_txt_2 y padron_parquet_2 (No deberían incluir espacios innecesarios). ¿Alguna conclusión?

```
hive> SELECT

DESC_DISTRITO,

DESC_BARRIO,

SUM(ESPANOLESHOMBRES) AS Total_EspanolesHombres,

SUM(ESPANOLESHOMBRES) AS Total_EspanolesMujeres,

SUM(EXTRANJEROSHOMBRES) AS Total_ExtranjerosHombres,

SUM(EXTRANJEROSMUJERES) AS Total_ExtranjerosMujeres

FROM

Addos_padron.padron_txt_2

GROUP BY

DESC_DISTRITO,

DESC_BARRIO;
```

```
desc_distrito desc_barrio total_espanoleshombres total_espanolesmujeres total_extranjeroshombres
                                                                                                           total extranierosmuieres
                               0.0 0.0
 ARGANZUELA
                                              1.0
                                                     0.0
                       ACACIAS
                                              15297.0 17871.0 1636.0 1895.0
 ARGANZUELA
                                              899.0 921.0 129.0 152.0
7786.0 9357.0 1635.0 1806.0
 ARGANZUELA
                       ATOCHA
                       CHOPERA
ARGANZUELA
Ahora en parquet...
hive> SELECT
              DESC DISTRITO,
              DESC BARRIO,
      >
              SUM(ESPANOLESHOMBRES) AS Total EspanolesHombres,
              SUM(ESPANOLESMUJERES) AS Total EspanolesMujeres,
              SUM(EXTRANJEROSHOMBRES) AS Total ExtranjerosHombres,
              SUM(EXTRANJEROSMUJERES) AS Total ExtranjerosMujeres
      > FROM
              padron parquet 2
      > GROUP BY
              DESC DISTRITO,
              DESC BARRIO;
desc_distrito desc_barrio total_espanoleshombres total_espanolesmujeres total_extranjeroshombres total_extranjerosmujeres
NULL NULL NULL NULL NULL
ARGANZUELA ACACIAS 15297.0 17871.0 1636.0 1895.0
                       NULL NULL
ACACIAS
ARGANZUELA
                                              899.0 921.0 129.0 152.0
7786.0 9357.0 1635.0 1806.0
                       ATOCHA
ARGANZUELA
                       CHOPERA
                                             11928.0 13417.0 2098.0 2331.0
ARGANZUELA
                       DELICIAS
```

• 3.7) 3.8) Llevar a cabo la misma consulta sobre las mismas tablas en Impala. ¿Alguna conclusión?

¿Se percibe alguna diferencia de rendimiento entre Hive e Impala?

```
[quickstart.cloudera:21000] > SELECT
                         DESC DISTRITO,
                     >
                         DESC BARRIO,
                         SUM(CAST(ESPANOLESHOMBRES AS INT)) AS Total Es
                     >
panolesHombres.
                         SUM(CAST(ESPANOLESMUJERES AS INT)) AS Total Es
panolesMujeres,
                         SUM(CAST(EXTRANJEROSHOMBRES AS INT)) AS Total
                     >
ExtranjerosHombres,
                         SUM(CAST(EXTRANJEROSMUJERES AS INT)) AS Total
ExtranjerosMujeres
                     > FROM
                         default.padron parquet 2
                    >
                    > GROUP BY
                         DESC DISTRITO,
                         DESC BARRIO;
spanolesmujeres | total_extranjeroshombres | total_extranjerosmujeres |
| LEGAZPI
ARGANZUELA
                                  8956
                                                      9338
                                950
            859
                 | RIOS ROSAS
CHAMBERI
                                  10525
                                                      13407
            1915
                                2397
                                                     П
                                                      | 3182
                 CASCO H.BARAJAS
BARAJAS
                                 2983
            832
                                955
                                                      | 10971
| USERA
                 | ORCASITAS
                                9758
            1480
                                | 1650
                                   1 5276
I DADATAC
                 I TTMON
                                                      1 5005
```

En impala he percibido que la ejecución se ha realizado mucho más rápido que en hive, me ha sorprendido lo rápido que ha extraído los datos. Impala es más rápido que Hive porque ejecuta consultas en memoria y usa un motor de procesamiento distribuido optimizado, mientras que Hive utiliza MapReduce, que es más lento debido a su enfoque en disco.

- 4- Sobre tablas particionadas.
- 4.1) Crear tabla (Hive) padron_particionado particionada por campos DESC_DISTRITO y DESC_BARRIO cuyos datos estén en formato parquet.

```
hive> CREATE TABLE padron particionado (
          COD DISTRITO INT,
         COD DIST BARRIO INT.
         COD BARRIO INT,
         COD DIST SECCION INT,
         COD SECCION INT,
        COD EDAD INT STRING,
        ESPANOLESHOMBRES INT,
         ESPANOLESMUJERES INT,
         EXTRANJEROSHOMBRES INT,
        EXTRANJEROSMUJERES INT,
        FX CARGA STRING,
         FX DATOS INI STRING,
         FX_DATOS_FIN STRING
   > PARTITIONED BY (DESC DISTRITO STRING, DESC BARRIO STRING)
    > STORED AS PARQUET;
oĸ
```

• 4.2) Insertar datos (en cada partición) dinámicamente (con Hive) en la tabla recién creada a partir de un select de la tabla padron_parquet_2.

Me ha dado un error en la inserción de datos ya que me indicaba que Hive intentaba crear demasiadas particiones dinámicas. Este problema es común cuando se intenta particionar por columnas con una alta cardinalidad sin ajustar los límites de particiones dinámicas.

Para resolverlo, he incrementado los límites de particion:

```
hive> SET hive.exec.max.dynamic.partitions = 1000;
hive> SET hive.exec.max.dynamic.partitions.pernode = 1000;
hive> INSERT INTO TABLE padron particionado PARTITION (DESC DISTRITO, DESC BARRIO)
        COD DISTRITO,
        COD_DIST_BARRIO,
        COD BARRIO,
        COD DIST SECCION,
        COD SECCION,
        COD EDAD INT.
        ESPANOLESHOMBRES,
        ESPANOLESMUJERES,
        EXTRANJEROSHOMBRES.
        EXTRANJEROSMUJERES,
        FX CARGA,
        FX DATOS INI,
        FX DATOS FIN,
        DESC DISTRITO,
        DESC BARRIO
    > FROM padron parquet 2;
```

• 4.3) Hacer invalidate metadata en Impala de la base de datos padron_particionado.

[[quickstart.cloudera:21000] > INVALIDATE METADATA default.padron particionado;

• 4.4) • 4.5) Calcular el total de EspanolesHombres, EspanolesMujeres, ExtranjerosHombres y ExtranjerosMujeres agrupado por DESC_DISTRITO y DESC_BARRIO para los distritos CENTRO, LATINA, CHAMARTIN, TETUAN, VICALVARO y BARAJAS.

Llevar a cabo la consulta en Hive en las tablas padron_parquet y padron_partitionado. ¿Alguna conclusión?

```
hive> SELECT
          DESC DISTRITO,
          DESC BARRIO,
          SUM(ESPANOLESHOMBRES) AS Total_EspanolesHombres,
    >
          SUM(ESPANOLESMUJERES) AS Total_EspanolesMujeres,
          SUM(EXTRANJEROSHOMBRES) AS Total ExtranjerosHombres,
          SUM(EXTRANJEROSMUJERES) AS Total ExtranjerosMujeres
          default.padron parquet
    > WHERE
          DESC DISTRITO IN ('CENTRO', 'LATINA', 'CHAMARTIN', 'TETUAN', 'VICALVARO', 'BARAJAS')
    > GROUP BY
         DESC DISTRITO,
          DESC BARRIO;
hive> SELECT
          DESC DISTRITO,
          DESC BARRIO,
          SUM(ESPANOLESHOMBRES) AS Total EspanolesHombres,
          SUM(ESPANOLESMUJERES) AS Total EspanolesMujeres,
          SUM(EXTRANJEROSHOMBRES) AS Total ExtranjerosHombres,
          SUM(EXTRANJEROSMUJERES) AS Total ExtranjerosMujeres
    > FROM
          default.padron particionado
    > WHERE
          DESC DISTRITO IN ('CENTRO', 'LATINA', 'CHAMARTIN', 'TETUAN', 'VICALVARO', 'BARAJAS')
    > GROUP BY
          DESC DISTRITO.
          DESC BARRIO;
```

La particionada se realiza a mayor velocidad, ya que al estar particonada, no tiene que leer toda la tabla completa, por lo que es más eficiente, por lo que la tabla particionada, mejora el rendimiento de la consulta.

4.6) Llevar a cabo la consulta en Impala en las tablas padron_parquet y padron_particionado. ¿Alguna conclusión?

En Impala, la ejecución de consultas es mucho más rápida que en Hive, ya que Impala procesa datos en memoria con un motor distribuido optimizado, mientras que Hive usa MapReduce, más lento por su enfoque en disco. Además, las tablas particionadas en Hive mejoran el rendimiento porque solo leen las particiones relevantes en lugar de toda la tabla.

• 4.7) Hacer consultas de agregación (Max, Min, Avg, Count) tal cual el ejemplo anterior con las 3 tablas (padron_txt_2, padron_parquet_2 y padron_particionado) y comparar rendimientos tanto en Hive como en Impala y sacar conclusiones.

En impala:

En hive:

```
hive> SELECT

DESC_DISTRITO,
DESC_BARRIO,
SUM(CAST(ESPANOLESHOMBRES AS INT)) AS Total_EspanolesHombres,
SUM(CAST(ESPANOLESMUJERES AS INT)) AS Total_EspanolesMujeres,
SUM(CAST(EXTRANJEROSHOMBRES AS INT)) AS Total_ExtranjerosHombres,
SUM(CAST(EXTRANJEROSMUJERES AS INT)) AS Total_ExtranjerosMujeres
FROM
default.padron_txt_2
WHERE
DESC_DISTRITO IN ('CENTRO', 'LATINA', 'CHAMARTIN', 'TETUAN', 'VICALVARO', 'BARAJAS')
GROUP BY
DESC_DISTRITO,
DESC_DISTRITO,
```

En memoria con un motor distribuido optimizado, mientras que Hive usa MapReduce, más lento por su enfoque en disco. Además, las tablas particionadas en Hive mejoran el rendimiento porque solo leen las particiones relevantes en lugar de toda la tabla.

5- Trabajando con tablas en HDFS.

A continuación vamos a hacer una inspección de las tablas, tanto externas (no gestionadas) como internas (gestionadas). Este apartado se hará si se tiene acceso y conocimiento previo sobre cómo insertar datos en HDFS.

• 5.1) Crear un documento de texto en el almacenamiento local que contenga una secuencia de números distribuidos en filas y separados por columnas, llámalo datos1 y que sea por ejemplo: 1,2,3 4,5,6 7,8,9



datos1.txt

 5.2) Crear un segundo documento (datos2) con otros números pero la misma estructura.



datos2.txt

• 5.3) Crear un directorio en HDFS con un nombre a placer, por ejemplo, /test. Si estás en una máquina Cloudera tienes que asegurarte de que el servicio HDFS está activo ya que puede no iniciarse al encender la máquina (puedes hacerlo desde el Cloudera Manager). A su vez, en las máquinas Cloudera es posible (dependiendo de si usamos Hive desde consola o desde Hue) que no tengamos permisos para crear directorios en HDFS salvo en el directorio /user/cloudera.

```
[cloudera@quickstart ~]$ hdfs dfs -mkdir /user/cloudera/test
[cloudera@quickstart ~]$ hdfs dfs -ls /user/cloudera
Found 5 items
drwxrwxrwx - cloudera cloudera 0 2024-07-10 11:42 /user/cloudera/ejercicios
drwxrwxrwx - cloudera cloudera 0 2024-07-11 14:05 /user/cloudera/ejercicios_master
drwxrwxrwx - cloudera cloudera 0 2024-07-10 13:37 /user/cloudera/terasort-input
drwxrwxrwx - cloudera cloudera 0 2024-07-10 13:38 /user/cloudera/terasort-output
drwxr-xr-x - cloudera cloudera 0 2024-07-31 14:31 /user/cloudera/test
```

• 5.4) Mueve tu fichero datos1 al directorio que has creado en HDFS con un comando desde consola.

• 5.5) Desde Hive, crea una nueva database por ejemplo con el nombre numeros. Crea una tabla que no sea externa y sin argumento location con tres columnas numéricas, campos separados por coma y delimitada por filas. La llamaremos por ejemplo numeros_tbl.

```
hive> CREATE DATABASE numeros;
0K
Time taken: 0.343 seconds
hive> USE numeros;
Time taken: 0.023 seconds
hive> CREATE TABLE numeros tbl (
          num1 INT,
    >
          num2 INT,
          num3 INT
    > )
    > ROW FORMAT DELIMITED
    > FIELDS TERMINATED BY ','
    > LINES TERMINATED BY '\n';
0K
Time taken: 0.244 seconds
```

• 5.6) Carga los datos de nuestro fichero de texto datos1 almacenado en HDFS en la tabla de Hive. Consulta la localización donde estaban anteriormente los datos almacenados. ¿Siguen estando ahí? ¿Dónde están?. Borra la tabla, ¿qué ocurre con los datos almacenados en HDFS?

Inserción de datos...

```
hive> LOAD DATA INPATH '/user/cloudera/test/datos1.txt' INTO TABLE numeros tbl;
Loading data to table numeros.numeros tbl
Table numeros.numeros tbl stats: [numFiles=1, totalSize=18]
Time taken: 0.311 seconds
hive> select * from numeros tbl;
numeros tbl.num1
                          numeros tbl.num2
                                                   numeros tbl.num3
         2
                 3
1
4
         5
                 6
Time taken: 0.285 seconds, Fetched: 3 row(s)
Vemos donde estaban almacenados los datos antes...
hive> DESCRIBE FORMATTED numeros tbl;
                     hdfs://quickstart.cloudera:8020/user/hive/warehouse/numeros.db/numeros tbl
Vemos que después de eliminar la tabla no aparece:
```

[cloudera@quickstart Padron]\$ hdfs dfs -ls /user/cloudera/test/;

- 1. La tabla en Hive no es externa: En Hive, si la tabla no está definida como externa y eliminas la tabla con DROP TABLE, Hive eliminará automáticamente los datos almacenados en HDFS asociados a esa tabla.
- 2. Ubicación predeterminada de la tabla: Si no se especifica una ubicación externa personalizada para la tabla, los datos se almacenan en la ubicación predeterminada de Hive, que se encuentra en /user/hive/warehouse (o similar). Al eliminar la tabla, Hive elimina tanto la definición de la tabla como los datos en la ubicación predeterminada. (es el caso)

• 5.7) Vuelve a mover el fichero de texto datos1 desde el almacenamiento local al directorio anterior en HDFS.

```
[cloudera@quickstart Padron]$ hdfs dfs -put datos1.txt /user/cloudera/test
[cloudera@quickstart Padron]$ hdfs dfs -ls /user/cloudera/test
Found 1 items
-rw-r--r-- 1 cloudera cloudera 18 2024-07-31 14:48 /user/cloudera/test
/datos1.txt
```

• 5.8) Desde Hive, crea una tabla externa sin el argumento location. Y carga datos1 (desde HDFS) en ella. ¿A dónde han ido los datos en HDFS? Borra la tabla ¿Qué ocurre con los datos en hdfs?

```
hive> CREATE EXTERNAL TABLE numeros externos (
           coll INT,
     >
           col2 INT,
     >
           col3 INT
     >
     > )
     > ROW FORMAT DELIMITED
     > FIELDS TERMINATED BY ','
     > LINES TERMINATED BY '\n';
Time taken: 1.52 seconds
hive> LOAD DATA INPATH '/user/cloudera/test/datos1.txt' INTO TABLE numeros exter
nos;
Loading data to table numeros.numeros externos
Table numeros.numeros externos stats: [numFiles=1, numRows=0, totalSize=18, rawD
ataSize=01
0K
Time taken: 0.225 seconds
Una vez hecha la inserción, se puede comprobar que los datos están en hdfs:
[cloudera@quickstart Padron]$ hdfs dfs -ls /user/cloudera/test
Found 1 items
-rwxr-xr-x 1 cloudera cloudera 18 2024-08-01 08:18 /user/cloudera/test
/datos1.txt
Borro la tabla...
hive> drop table numeros externos;
0K
Time taken: 0.359 seconds
Compruebo si siguen los datos...
[cloudera@quickstart Padron]$ hdfs dfs -ls /user/cloudera/test
Found 1 items
-rwxr-xr-x 1 cloudera cloudera
                                    18 2024-08-01 08:18 /user/cloudera/test
/datos1.txt
```

Los datos permanecen en HDFS a pesar de eliminar la tabla externa en Hive porque:

- Tabla externa: Al crear una tabla externa, Hive solo gestiona la metadata y no los datos.
- Eliminar la tabla: Eliminar una tabla externa en Hive solo elimina la metadata de la tabla, pero no los datos en HDFS.

Así, los datos siguen existiendo en HDFS después de la eliminación de la tabla externa en Hive.

• 5.9) Borra el fichero datos1 del directorio en el que estén. Vuelve a insertarlos en el directorio que creamos inicialmente (/test). Vuelve a crear la tabla numeros desde hive pero ahora de manera externa y con un argumento location que haga referencia al directorio donde los hayas situado en HDFS (/test). No cargues los datos de ninguna manera explícita. Haz una consulta sobre la tabla que acabamos de crear que muestre todos los registros. ¿Tiene algún contenido?

```
[cloudera@quickstart Padron] hdfs dfs -rm /user/cloudera/test/datos1.txt
Deleted /user/cloudera/test/datos1.txt
[cloudera@quickstart Padron]$ hdfs dfs -put datos1.txt /user/cloudera/test/
[cloudera@quickstart Padron]$ hdfs dfs -ls /user/cloudera/test
Found 1 items
-rw-r--r--
            1 cloudera cloudera
                                        18 2024-08-01 08:26 /user/cloudera/test
/datos1.txt
hive> CREATE EXTERNAL TABLE numeros externos (
          num1 INT,
    >
          num2 INT,
    >
          num3 INT
    >
    > )
    > ROW FORMAT DELIMITED
    > FIELDS TERMINATED BY ','
    > LINES TERMINATED BY '\n'
    > LOCATION '/user/cloudera/test';
0K
Time taken: 0.066 seconds
hive> SELECT * FROM numeros externos;
0K
numeros externos.num1
                         numeros externos.num2
                                                  numeros externos.num3
        2
1
                3
4
        5
                 6
7
                9
        8
Time taken: 0.37 seconds, Fetched: 3 row(s)
```

La tabla externa se refiere directamente a los datos en HDFS. No es necesario cargar los datos explícitamente en la tabla ya que Hive mapea directamente al archivo o directorio especificado en LOCATION. Por lo tanto, al hacer la consulta, los datos presentes en el archivo datos1 en el directorio /test serán mostrados.

• 5.10) Inserta el fichero de datos creado al principio, "datos2" en el mismo directorio de HDFS que "datos1". Vuelve a hacer la consulta anterior sobre la misma tabla. ¿Qué salida muestra?

```
[cloudera@quickstart Padron]$ hdfs dfs -put datos2.txt /user/cloudera/test/
[cloudera@quickstart Padron]$ hdfs dfs -ls /user/cloudera/test
Found 2 items
- rw - r - - r - -
             1 cloudera cloudera
                                         18 2024-08-01 08:26 /user/cloudera/test
/datos1.txt
             1 cloudera cloudera
                                         27 2024-08-01 08:34 /user/cloudera/test
 -rw-r--r--
/datos2.txt
hive> SELECT * FROM numeros externos;
0K
numeros_externos.num1 numeros_externos.num2 numeros_externos.num3
1
        2
                3
4
        5
                6
        8
                9
10
                32
        21
43
        54
        87
76
                98
Time taken: 0.06 seconds, Fetched: 6 row(s)
```

Como insertamos de nuevo en la carpeta a la que hace referencia la tabla, todos los datos que metamos dentro de esa carpeta, los podremos utilizar ya que ta tabla pilla su referencia, por lo que podemos sacar como conclusión que trabajar con 'LOCATION', es una forma bastante eficiente y fácil de usar, que nos permite trabajar de una manera más cómoda.

• 5.11) Extrae conclusiones de todos estos anteriores apartados.

- 1. Rendimiento de Hive vs. Impala:
 - Impala: Observamos que las consultas en Impala se ejecutan mucho más rápido que en Hive. Esto se debe a que Impala ejecuta las consultas en memoria y utiliza un motor de procesamiento distribuido optimizado. Esto contrasta con Hive, que utiliza MapReduce, un enfoque más lento debido a su operación basada en disco.
 - Tablas particionadas: Las tablas particionadas ofrecen un mejor rendimiento en las consultas porque no necesitan leer toda la tabla, sino solo las particiones relevantes.
 Esto mejora significativamente la eficiencia de las consultas.
- 2. Operaciones con tablas en Hive:
 - Carga de datos en tablas Hive:
 - Al cargar datos en una tabla interna de Hive, los datos se mueven desde su ubicación original en HDFS al directorio gestionado por Hive. Esto hace que los datos desaparezcan de su ubicación original una vez que se cargan en la tabla.
 - Si se elimina una tabla interna, los datos almacenados en HDFS también se eliminan.
 - o Tablas externas en Hive:
 - Al crear y cargar datos en una tabla externa, Hive no mueve los datos, sino que los referencia directamente desde su ubicación original en HDFS. Esto significa que los datos permanecen en su lugar original.
 - Al eliminar una tabla externa, los datos en HDFS no se eliminan porque Hive no es el propietario de los datos, solo los referencia.
- 3. Gestión de archivos en HDFS:
 - Borrar y mover archivos en HDFS:
 - Los comandos HDFS (hdfs dfs -rm y hdfs dfs -put) permiten gestionar los archivos y directorios en el sistema de archivos distribuido.
 - Creación de directorios y archivos:
 - Creamos directorios en HDFS para organizar y almacenar datos, facilitando su referencia desde Hive.
 - o Persistencia de datos:
 - Los datos cargados en tablas externas en Hive permanecen en su ubicación original en HDFS incluso después de eliminar la tabla, lo cual es útil para la gestión y persistencia de grandes volúmenes de datos.
- 4. Resumen general:
 - Impala es altamente eficiente para consultas rápidas debido a su enfoque en memoria.
 - Hive es útil para operaciones ETL y manipulación de datos debido a su integración con Hadoop y su capacidad para manejar grandes volúmenes de datos.
 - Tablas particionadas mejoran el rendimiento al reducir la cantidad de datos leídos durante las consultas.
 - Tablas externas en Hive proporcionan flexibilidad para referenciar datos en HDFS sin moverlos, asegurando que los datos permanezcan accesibles incluso si la tabla se elimina.

Estas conclusiones resaltan las ventajas y consideraciones prácticas al utilizar herramientas de Big Data como Hive e Impala, así como la importancia de una adecuada gestión de datos en HDFS para un rendimiento y eficiencia óptimos.

La siguiente sección de la práctica se abordará si ya se tienen suficientes conocimientos de Spark, en concreto de el manejo de DataFrames, y el manejo de tablas de Hive a través de Spark.sql.

6-SPARK

• 6.1) Comenzamos realizando la misma práctica que hicimos en Hive en Spark, importando el csv. Sería recomendable intentarlo con opciones que quiten las "" de los campos, que ignoren los espacios innecesarios en los campos, que sustituyan los valores vacíos por 0 y que infiera el esquema.

```
from pyspark.sql import SparkSession
from pyspark.sql.functions import col, trim, when
spark = SparkSession.builder.appName("Padron").getOrCreate()
df = spark.read.csv("/FileStore/tables/Formacion_Binaia/Practica_Padron/estadisticas202407.csv",
                   header=True, #primera linea de archivo son columnas
                   inferSchema=True, #permite a spark inferir el tipo de datos de cada columna
                   sep=';', #separacion de campos
                   quote='"',  #indica que las comillas se usan para envolver campos
                   ignoreLeadingWhiteSpace=True, #ignora espacios en blanco al principio
                   ignoreTrailingWhiteSpace=True, #y al final de los campos
                   nullValue='', #celdas vacias como valores nulos
                   nanValue='0') #sustituye valores nan por 0
# Reemplazar valores nulos con 0
df = df.fillna(0)
# Eliminar espacios innecesarios en los campos
df = df.select([trim(col(c)).alias(c) for c in df.columns])
# Mostrar el esquema para confirmar que se infirió correctamente
df.printSchema()
# Mostrar algunas filas para verificar los datos
df.show(10)
```

	+		+-		+		-+		
COD_D	ISTRITO	DESC_DISTRITO CC	D_DIST_BARRIO	DESC_BARRIO CO	D_BARRIO	COD_DIST_SECCION COD	_SECCION COD_E	DAD_INT ESPANO	LESHOMBRES ESPAN
OLESMU	JERES EXT	RANJEROSHOMBRES EXTRA	NJEROSMUJERES	FX_C	ARGA FX_D	ATOS_INI FX_DATOS_FI	N		
+	+-	+	+-			+			
	+		+-		+		-+		
	2	ARGANZUELA	202	ACACIAS	2	2094	94	47	6
16		0	0 2024-07-07	00:00: 20	24-07-01	2024-07-01			
	18	VILLA DE VALLECAS	1802 5	ANTA EUGENIA	2	18029	29	77	7
3		0	2 2024-07-07	00:00: 202	4-07-01	2024-07-01			
	2	ARGANZUELA	205	DELICIAS	5	2062	62	71	6
6		0	2 2024-07-07	00:00: 202	4-07-01	2024-07-01			
	9	MONCLOA-ARAVACA	907	ARAVACA	7	9083	83	55	15
14		0	0 2024-07-07	00:00: 20	24-07-01	2024-07-01			
	3	RETIRO	306	NIÑO JESUS	6	3080	80	60	7
8		1	1 2024-07-07	00:00: 202	4-07-01	2024-07-01			
	4	SALAMANCA	405	LISTA	5	4102	102	78	1
6		0	0 2024-07-07	00:00: 202	4-07-01	2024-07-01			

• 6.2) De manera alternativa también se puede importar el csv con menos tratamiento en la importación y hacer todas las modificaciones para alcanzar el mismo estado de limpieza de los datos con funciones de Spark.

```
%python
from pyspark.sql import SparkSession
from pyspark.sql.functions import col, trim, regexp replace
# Crear una sesión de Spark
spark = SparkSession.builder.appName("Padron").getOrCreate()
# Leer el archivo CSV en un DataFrame de Spark con menos opciones de tratamiento inicial
df = spark.read.csv("/FileStore/tables/Formacion_Binaia/Practica_Padron/estadisticas202407.csv", header=True, sep=';')
# Mostrar el esquema inicial para ver cómo se cargaron los datos
df.printSchema()
# Reemplazar comillas dobles y espacios innecesarios en los campos
df = df.select([trim(regexp_replace(col(c), '"', '')).alias(c) for c in df.columns])
# Reemplazar valores vacíos con 0
df = df.fillna('0')
# Convertir las columnas adecuadas a int
int_columns = ['ESPANOLESHOMBRES', 'ESPANOLESMUJERES', 'EXTRANJEROSHOMBRES', 'EXTRANJEROSMUJERES']
for column in int columns:
  df = df.withColumn(column, col(column).cast('int'))
# Mostrar el esquema para confirmar que se infirió correctamente
df.printSchema()
# Mostrar algunas filas para verificar los datos
```

_	STRITO	DESC_DISTRITO CO		DESC_BARRIO COD_	BARRIO COD_D	IST_SECCION COD_	SECCION COD_E	DAD_INT ESPANOL	ESHOMBRES ESPAN
OLESMUJ +	ERES EXT	RANJEROSHOMBRES EXTRA		_		INI FX_DATOS_FIN	•		
	+				+				
	2	ARGANZUELA	202	ACACIAS	2	2094	94	47	6
16		0	0 2024-07-07	00:00: 2024	-07-01 202	4-07-01			
	18	VILLA DE VALLECAS	1802 SA	NTA EUGENIA	2	18029	29	77	7
3		0	2 2024-07-07 0	0:00: 2024-	07-01 2024	-07-01			
	2	ARGANZUELA	205	DELICIAS	5	2062	62	71	6
6		0	2 2024-07-07 0	0:00: 2024-	07-01 2024	-07-01			
1	9	MONCLOA-ARAVACA	907	ARAVACA	7	9083	83	55	15
14		0	0 2024-07-07	00:00: 2024	-07-01 202	4-07-01			
1	3	RETIRO	306	NIÑO JESUS	6	3080	80	60	7

• 6.3) Enumera todos los barrios diferentes.

```
%python
 # Enumerar todos los barrios diferentes
 barrios_unicos = df.select('DESC_BARRIO').distinct().orderBy('DESC_BARRIO')
 # Mostrar los barrios únicos
 barrios unicos.show(truncate=False)
▶ ■ barrios_unicos: pyspark.sql.dataframe.DataFrame = [DESC_BARRIO: string]
ABRANTES
ACACIAS
ADELFAS
AEROPUERTO
ALAMEDA DE OSUNA
LALMAGRO
LALMENARA
ALMENDRALES
ALUCHE
AMPOSTA
ANGELES
APOSTOL SANTIAGO
| ARAPILES
ARAVACA
IARCOS
ARGUELLES
ATALAYA
ATOCHA
|BELLAS VISTAS
only showing top 20 rows
```

• 6.4) Crea una vista temporal de nombre "padron" y a través de ella cuenta el número de barrios diferentes que hay.

```
%python
#Una vista temporal permite ejecutar consultas SQL sobre datos que ya están en un DataFrame de Spark, sin tener que realizar el procesamiento de datos nuevamente.
# Crear una vista temporal
df.createOrReplaceTempView("padron")

# Contar el número de barrios diferentes usando la vista temporal
count_barrios = spark.sql@"""

| SELECT COUNT(DISTINCT DESC_BARRIO) AS num_barrios|
| FROM padron
| """ |
| # Mostrar el resultado
count_barrios.show()

* (3) Spark Jobs

* (2) Count_barrios: pyspark.sql.dataframe.DataFrame = [num_barrios: long]
| hum_barrios|
| hum_barrios|
| 131|
```

• 6.5) Crea una nueva columna que muestre la longitud de los campos de la columna DESC_DISTRITO y que se llame "longitud".

withcolumn, se usa para añadir una nueva columna o para reemplazar una columna existente con una nueva versión.

```
from pyspark.sql.functions import length
 # Crear una nueva columna 'longitud' que muestra la longitud de los valores en 'DESC_DISTRITO'
 df = df.withColumn("longitud", length(col("DESC_DISTRITO")))
 # Mostrar el esquema actualizado y algunas filas para verificar el resultado
 df.printSchema()
 df.show(10)
▶ (1) Spark Jobs
▶ 🔳 df: pyspark.sql.dataframe.DataFrame = [COD_DISTRITO: string, DESC_DISTRITO: string ... 14 more fields]
--+-----
       2094
                                                                                 94
                                                                                            47
                                                                                                          6
                                                                                                                       16
0
                                                           10
                                                                   18029
        18 | VILLA DE VALLECAS
                                   1802|SANTA EUGENIA|
                                                          2
                                                                                 29
                                                                                            77
                                                                                                         7
                                                                                                                        3
       | 2|2024-07-07 00:00:...| 2024-07-01| 2024-07-01| 17|
| 2| ARGANZUELA| 205| DELICIAS| 5|
| 2|2024-07-07 00:00:...| 2024-07-01| 2024-07-01| 10|
| 9| MONCLOA-ARAVACA| 907| ARAVACA| 7|
0
                                                                    2062
                                                                                 62
                                                                                            71
                                                                                                          6
                                                                                                                        6
0
              907| ARAVACA| 7|
0|2024-07-07 00:00:...| 2024-07-01| 2024-07-01| 15|
RETIRO| 306| NTÃO 35515|
       9 | MONCLOA-ARAVACA |
                                                                     9083
                                                                                 83
                                                                                            55
                                                                                                         15
                                                                                                                       14
       3 |
                                                                     3080
                                                                                 80
                                                                                            60
                                                                                                          7
                                                                                                                        8
```

• 6.6) Crea una nueva columna que muestre el valor 5 para cada uno de los registros de la tabla.

from pyspark.sql.functions import lit # Crear una nueva columna 'longitud' que muestra la longitud de los valores en 'DESC_DISTRITO' df = df.withColumn("valor fijo", lit(5)) # Mostrar el esquema actualizado y algunas filas para verificar el resultado df.show(10) ▶ (1) Spark Jobs ▶ 📾 df: pyspark.sql.dataframe.DataFrame = [COD_DISTRITO: string, DESC_DISTRITO: string ... 15 more fields] ARGANZUELA| 202| ACACIAS| 2| 0|2024-07-07 00:00:...| 2024-07-01| 2024-07-01| 1 | 18 | VILLA DE VALLECAS| 1802|SANTA EUGENIA| | 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-0 7 2062 | 62 | 71 ARGANZUELA| 205| DELICIAS| 2|2024-07-07 00:00:...| 2024-07-01| 2024-07-01| 2|2024-07-07 00:00:...| 2024-07-01| 2024-07-01|
9| MONCLOA-ARAVACA| 997| ARAVACA|
0|2024-07-07 00:00:...| 2024-07-01| 2024-07-01|
3| RETIRO| 306| NIÑO 155US|
1|2024-07-07 00:00:..| 2024-07-01| 2024-07-01|
4| SALAMANCA| 405| LISTA| 9083| 83| ر الموادي | 15 | 5 | | 3080 | 80 | 55 15 60 5| 4102| 1 | 2024-07-07 00:100:....|
4 | SALMANIKA| 405 | LISTA|
0 | 2024-07-07 00:00:...| 2024-07-01 | 2024-07-01 |
20 | SAN BLAS-CANILLEJAS| 2001 | STIMMKAS| 78 102 1 6 9 100 0 +

• 6.7) Borra esta columna.

```
df = df.drop("valor_fijo")
df.printSchema()
▶ ■ df: pyspark.sql.dataframe.DataFrame = [COD_DISTRITC
|-- COD_DISTRITO: string (nullable = true)
-- DESC_DISTRITO: string (nullable = true)
|-- COD_DIST_BARRIO: string (nullable = true)
|-- DESC_BARRIO: string (nullable = true)
|-- COD_BARRIO: string (nullable = true)
|-- COD_DIST_SECCION: string (nullable = true)
|-- COD_SECCION: string (nullable = true)
|-- COD_EDAD_INT: string (nullable = true)
|-- ESPANOLESHOMBRES: string (nullable = true)
|-- ESPANOLESMUJERES: string (nullable = true)
|-- EXTRANJEROSHOMBRES: string (nullable = true)
|-- EXTRANJEROSMUJERES: string (nullable = true)
|-- FX_CARGA: string (nullable = true)
|-- FX_DATOS_INI: string (nullable = true)
|-- FX_DATOS_FIN: string (nullable = true)
|-- longitud: integer (nullable = true)
```

• 6.8) Particiona el DataFrame por las variables DESC_DISTRITO y DESC_BARRIO.

```
%python
# Guardar el DataFrame particionado en Parquet
output_path = "/FileStore/tables/Formacion_Binaia/Practica_Padron/"

df.write.partitionBy("DESC_DISTRITO", "DESC_BARRIO") \
    .format("parquet") \
    .mode("overwrite") \
    .save(output_path)
```

DESC_DISTRITO=__HIVE_DEFAU... ☐ DESC_DISTRITO=ARGANZUELA▼ □ DESC_DISTRITO=BARAJAS ▼ DESC_DISTRITO=CARABANCHET □ DESC_DISTRITO=CENTRO ▼ □ DESC DISTRITO=CHAMARTIN ▼ DESC_DISTRITO=CHAMBERI • DESC_DISTRITO=CIUDAD LINEAT DESC_DISTRITO=FUENCARRAL-T.. □ DESC DISTRITO=HORTALEZA ▼ DESC_DISTRITO=LATINA □ DESC_DISTRITO=MONCLOA-A..▼ DESC DISTRITO=MORATALAZ ☐ DESC_DISTRITO=PUENTE DE V... DESC_DISTRITO=RETIRO □ DESC_DISTRITO=SALAMANCA ▼ DESC DISTRITO=SAN BLAS-CA. DESC_DISTRITO=TETUAN DESC_DISTRITO=USERA □ DESC_DISTRITO=VICALVARO ▼

© DESC_BARRIO=_HIVE_DEFAULT_P.▼

© DESC_BARRIO=ACACIAS

© DESC_BARRIO=ATOCHA

© DESC_BARRIO=CHOPERA

© DESC_BARRIO=DELICIAS

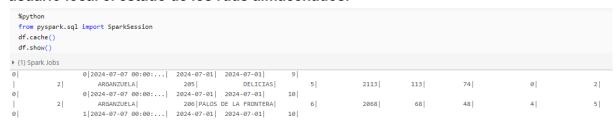
© DESC_BARRIO=IMPERIAL

© DESC_BARRIO=LEGAZPI

© DESC_BARRIO=PALOS DE LA FRO...▼

Se puede ver que ha ido bien, ya que tenemos, los distritos, y dentro de cada uno de ellos, tenemos los barrios dentro de cada distrito.

• 6.9) Almacénalo en caché. Consulta en el puerto 4040 (Ul de Spark) de tu usuario local el estado de los rdds almacenados.



Como estoy trabajando en databricks, no puedo porque el puerto 4040 es local al entorno del cluster de Spark y no está disponible para acceso remoto o desde el navegador local. En Databricks, la administración de Spark y la visualización de trabajos se realiza a través de la interfaz proporcionada por Databricks, no mediante acceso local a puertos.

Podemos visualizar el job de la cache:



• 6.10) Lanza una consulta contra el DF resultante en la que muestre el número total de "espanoleshombres", "espanolesmujeres", extranjeroshombres" y "extranjerosmujeres" para cada barrio de cada distrito. Las columnas distrito y barrio deben ser las primeras en aparecer en el show. Los resultados deben estar ordenados en orden de más a menos según la columna "extranjeroshombres".



Como lo tenia como String las columnas, he hecho un cast para que pueda usarlas como es adecuado.

Ejecuto la consulta y vemos los resultados

```
from pyspark.sql.functions import sum, col
   # Agrupar por DESC DISTRITO y DESC BARRIO, sumar las columnas relevantes
   result_df = df.groupBy("DESC_DISTRITO", "DESC_BARRIO").agg(
          sum("ESPANOLESHOMBRES").alias("Total_EspanolesHombres"),
         sum("ESPANOLESMUJERES").alias("Total_EspanolesMujeres"),
          sum("EXTRANJEROSHOMBRES").alias("Total_ExtranjerosHombres"),
          sum("EXTRANJEROSMUJERES").alias("Total_ExtranjerosMujeres")
   # Ordenar los resultados por Total_ExtranjerosMujeres y Total_ExtranjerosHombres
   result_df = result_df.orderBy(col("Total_ExtranjerosMujeres").desc(), col("Total_ExtranjerosHombres").desc())
   # Mostrar los resultados
   result_df.show()
▶ (2) Spark Jobs
 ▶ 🔳 result_df: pyspark.sql.dataframe.DataFrame = [DESC_DISTRITO: string, DESC_BARRIO: string ... 4 more fields]
           DESC DISTRITO
                                                      DESC BARRIO|Total EspanolesHombres|Total EspanolesMujeres|Total ExtranjerosHombres|Total ExtranjerosMujeres|
PUENTE DE VALLECAS | SAN DIEUW |
PUEBLO NUEVO |
PUE
SAN DIEGOL
                                                                                                                14886
                                                                                                                                                              16515
                                                                                                                                                                                                                 9795
                                                                                                                                                                                                                                                                  9759
                                                                                                                 23698
                                                                                                                                                              27603
                                                                                                                                                                                                                 7666
                                                                                                                                                                                                                                                                   8301
                        LATINA| ALUCHE|
ABANCHEL| VISTA ALEGRE|
NUMANCIA|
                                                                                                             24896
                                                           ALUCHE
                                                                                                                                                             29873
                                                                                                                                                                                                                 7315
                                                                                                                                                                                                                                                                   8123
               CARABANCHEL
                                                                                                                16076
                                                                                                                                                             19626
                                                                                                                                                                                                                  7580
                                                                                                                                                                                                                                                                   7988
                                                                                                             17256
 PUENTE DE VALLECASI
                                                      NUMANCIA
                                                                                                                                                             19428
                                                                                                                                                                                                                 7772
                                                                                                                                                                                                                                                                  7715
               VILLAVERDE | VILLAVERDE ALTO C.H. |
                                                                                                                17489
                                                                                                                                                              19502
                                                                                                                                                                                                                  7354
                                                                                                                                                                                                                                                                   7379
                        CENTRO EMBAJADORES
                                                                                                                16514
                                                                                                                                                             16102
                                                                                                                                                                                                                 9510
                                                                                                                                                                                                                                                                   7121
           CIUDAD LINEAL
                                                             VENTAS
                                                                                                                 18949
                                                                                                                                                              22299
                                                                                                                                                                                                                  5552
                                                                                                                                                                                                                                                                   6381
                                                                                                                                                              14589
              CARABANCHEL
                                                  PUERTA BONITA
                                                                                                                12745
                                                                                                                                                                                                                  6129
                                                                                                                                                                                                                                                                   6375
                         LATINA PUERTA DEL ANGEL
                                                                                                                 15941
                                                                                                                                                              18308
                                                                                                                                                                                                                  5666
                                                                                                                                                                                                                                                                   6188
| FUENCARRAL-EL PARDO| VALVERDE| | CARABANCHEL| SAN ISIDRO|
                                                                                                                27266
                                                                                                                                                             29556
                                                                                                                                                                                                                  5277
                                                                                                                                                                                                                                                                   5903
                                                                                                                 14750
                                                                                                                                                              16721
                                                                                                                                                                                                                  5709
                                                                                                                                                                                                                                                                   5848
              HORTALEZA VALDEFUENTES 
CARABANCHEL OPAÑEL 
VILLAVERDE LOS ROSALES
                                                                                                                32214
                                                                                                                                                             32976
                                                                                                                                                                                                                  4491
                                                                                                                                                                                                                                                                   5329
                                                                                                                 11876
                                                                                                                                                              14224
                                                                                                                                                                                                                  4915
                                                                                                                                                                                                                                                                   5257
                                                                                                                 14267
                                                                                                                                                            15765
     VILLA DE VALLECAS|ENSANCHE DE VALLECAS|
                                                                                                                  23888
                                                                                                                                                              24623
                                                                                                                  12256
                                                                                                                                                           12210
```

6.11) Elimina el registro en caché.

```
%python
df.unpersist()

Out[36]: DataFrame[COD_DISTRITO: string, DESC_DISTRITO: string, COD_DIST_BARRIO: string, DESC_BARRIO: string, COD_BARRIO: string, COD_DIST_SECCION: string, COD_EDAD_INT: string, ESPANOLESHOMBRES: int, ESPANOLESHOMBRES: int, EXTRANJEROSHOMBRES: int, EXTRANJEROSHOMBRES: int, FX_CARGA: string, FX_DATOS_INI: string, FX_DATOS_FIN: string, Jongitud: int]
```

• 6.12) Crea un nuevo DataFrame a partir del original que muestre únicamente una columna con DESC_BARRIO, otra con DESC_DISTRITO y otra con el número total de "espanoleshombres" residentes en cada distrito de cada barrio. Únelo (con un join) con el DataFrame original a través de las columnas en común.

```
from pyspark.sql.functions import col, sum
 # Suponiendo que 'df' es tu DataFrame original
  # Agrupar por DESC_BARRIO y DESC_DISTRITO, sumar ESPANOLESHOMBRES
  df_sum = df.groupBy("DESC_BARRIO", "DESC_DISTRITO").agg(
     sum(col("ESPANOLESHOMBRES")).alias("Total_EspanolesHombres")
 # Mostrar el esquema y las primeras filas para confirmar
 df_sum.printSchema()
 df_sum.show(10)
▶ ■ df_sum: pyspark.sql.dataframe.DataFrame = [DESC_BARRIO: string, DESC_D
 |-- DESC_BARRIO: string (nullable = true)
 |-- DESC_DISTRITO: string (nullable = true)
 |-- Total_EspanolesHombres: long (nullable = true)
             ARRIO| DESC_DISTRITO|Total_EspanolesHombres|
          VALVERDE | FUENCARRAL-EL PARDO |
   FUENTELARREINA | FUENCARRAL-EL PARDO |
          PAVONES
                             MORATALAZI
                                                         3572
        EL GOLOSO|FUENCARRAL-EL PARDO|
|SAN JUAN BAUTISTA| CIUDAD LINEAL|
| CANILLAS| HORTALEZA|
                                                          5452
                        HORTALEZA|
CHAMARTIN|
                                                         16510
      PROSPERIDAD
                                                         14458
                           ARGANZUELA
          LEGAZPI
      COSTILLARES
                      CIUDAD LINEAL|
                                                          9871
only showing top 10 rows
```

# Unit el DataFrame original con el DataFrame sumado df joined = df.join(df_sum, on=["DESC_BARRIO", "DESC_DISTRITO"], how="inner") # Mostrar el esquema y las primeras filas para verificar la unión df joined.show(10) • (4) Spark Jobs • @ df.joined: pyspark.sql.dataframe.DataFrame = [DESC_BARRIO: string, DESC_DISTRITO: string15 more fields] el 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANIZUELA 2 202 2 2019 19 81 0 7 2024-07-09 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANIZUELA 2 202 2 2018 18 34 7 4 1/2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANIZUELA 2 202 2 2088 88 59 2 2 2 el 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANIZUELA 2 202 2 2088 88 59 2 2 2 el 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANIZUELA 2 202 2 2088 88 32 5 10 2/2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANIZUELA 2 202 2 2088 88 32 5 10 ACACIAS ARGANIZUELA 2 2024-07-01 10 15297 ACACIAS ARGANIZUELA 2 2024-07-01 2024-07-01 10 15297 ACA	%python									
# Mostrar el esquema y las primeras filas para verificar la unión df-joined.printSchema() df-joined.show(10) (4) Spark Jobs	•									
df_joined.printSchema() df_joined.show(10) (A) Spark Jobs	df_joined = df.join(df_sum, on=["DE	SC_BARRIO", "	DESC_DISTRI	TO"], how=	"inner")					
df_joined.printSchema() df_joined.show(10) (A) Spark Jobs 2024-07-08 90:90: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 2021 2 2018 18 34 7 4 2024-07-07 90:90: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 2021 2 2018 18 34 7 4 2024-07-07 90:90: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 2021 2 2018 18 34 7 4 2024-07-07 90:90: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 2021 2 2088 88 59 2 2 2 2024-07-07 90:90: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 2021 2 2088 88 32 5 10 2024-07-07 90:90: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 2021 2 2088 88 32 5 10 2024-07-07 90:90: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 2021 2 2088 88 32 5 10 2024-07-07 90:90: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 2021 2 2016 16 53 10 14 2024-07-07 90:90: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 2021 2 2016 16 53 10 14 2024-07-07 90:90: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 2024-07-01 10 15297 ACACIAS ARGANZUELA 3 2024-07-01 10 15297 ACACIAS ARGANZUELA 3 2024-07-01 10 15297 ACACIAS ARGANZUELA 3 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 3 2024-07-01 2024-07-01 10 15297 ACACIAS ARGAN	# M	£:1	1.							
df_joined.show(10) Mark Jobs		illas para v	enilican la	union						
(4) Spark Jobs □ df joined: pyspark.sql.dataframe.DataFrame = [DESC_BARRIC: string, DESC_DISTRITO: string 15 more fields] 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 2021 2 2019 19 81 0 7 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 2021 2 2018 18 34 7 4 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 2021 2 2088 88 59 2 2 2 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 2021 2 2088 88 32 5 10 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 2021 2 2088 88 32 5 10 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 2021 2 2016 16 53 10 14 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 2021 2 2016 16 53 10 14 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 2021 2 2092 92 69 11 13 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 2024-07-01 2024-07-01 50 5 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 2024-07-01 2024-07-01 50 5 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 2024-07-01 2024-07-01 50 5 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 2024-07-01										
Description Description	dT_Joined.snow(10)									
2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 2021 2 2019 18 81 0 7 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 2021 2 2018 18 34 7 4 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 202 2 2088 88 59 2 2 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 202 2 2088 88 59 2 2 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 202 2 2088 88 32 5 10 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 202 2 2016 16 53 10 14 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 202 2 2016 16 53 10 14 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 202 2 2092 92 69 11 13 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 2021 2 2016 16 28 14 11 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 2024-07-01 2024-0	(4) Spark Jobs									
2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANIZUELA 2 202 2 2019 19 81 0 7 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANIZUELA 2 202 2 2018 18 34 7 4 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANIZUELA 2 202 2 2088 88 59 2 2 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANIZUELA 2 202 2 2088 88 59 2 2 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANIZUELA 2 202 2 2088 88 32 5 10 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANIZUELA 2 202 2 2016 16 53 10 14 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANIZUELA 2 202 2 2016 16 53 10 14 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANIZUELA 2 202 2 2016 16 28 11 13 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANIZUELA 2 202 2 2016 16 28 14 11 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANIZUELA 2 202 2 2016 16 28 14 11 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANIZUELA 2 202 2 2016 16 28 14 11 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANIZUELA 2 202 2 2025 25 75 5 5 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297	 df ioined: pvspark.sql.dataframe.Data 	Frame = IDESC I	BARRIO: strino	a. DESC DIST	RITO: string 15 mor	e fields1				
2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANIZUELA 2 2021 2 2018 18 34 7 4 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANIZUELA 2 2021 2 2088 88 59 2 2 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANIZUELA 2 2021 2 2088 88 32 5 10 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANIZUELA 2 2021 2 2016 16 53 10 14 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANIZUELA 2 202 2 2016 16 53 10 14 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANIZUELA 2 202 2 2016 16 28 11 13 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANIZUELA 2 202 2 2016 16 28 14 11 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANIZUELA 2 2024 2 2016 16 28 14 11 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANIZUELA 2 2024 2025 25 75 5 5 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANIZUELA 2 2024-07-01 10 15297 ACACIAS ARGANIZUELA	2024-07-07 00:00: 2024-07-01	2024-07-01	10		15297					
ACACIAS ARGANZUELA 2 202 2 2018 18 34 7 4 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 202 2 2088 88 59 2 2 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 202 2 2088 88 32 5 10 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 2021 2 2016 16 53 10 14 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 2021 2 2016 16 53 10 14 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 2021 2 2092 92 69 11 13 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 202 2 2016 16 28 14 11 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 202 2 2016 16 28 14 11 ACACIAS ARGANZUELA 2 2024-07-01 10 15297 ACACIAS ARG	ACACIAS ARGANZUELA	2	202	2	2019	19	81	0	7	0
2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297	2024-07-07 00:00: 2024-07-01	2024-07-01	10		15297					
ACACIAS ARGANZUELA 2 202 2 2088 88 59 2 2 2 2024-07-07 00:00: 2024-07-01 10 15297	ACACIAS ARGANZUELA	2	202	2	2018	18	34	7	4	2
2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACITAS ARGANZUELA 2 202 2 2088 88 32 5 10 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACITAS ARGANZUELA 2 202 2 2016 16 53 10 14 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACITAS ARGANZUELA 2 202 2 2002 92 69 11 13 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACITAS ARGANZUELA 2 202 2 2016 16 28 14 11 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACITAS ARGANZUELA 2 202 2 2016 16 28 14 11 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACITAS ARGANZUELA 2 202 2 202 2 2016 16 28 14 11 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACITAS ARGANZUELA 2 202 2 2025 25 75 5 5 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297	2024-07-07 00:00: 2024-07-01	2024-07-01	10		15297					
ACACIAS ARGANZUELA 2 202 2 2088 88 32 5 10 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 202 2 2016 16 53 10 14 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 202 2 2092 92 69 11 13 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 202 2 2016 16 28 14 11 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 202 2 2016 16 28 14 11 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 2024 2025 25 75 5 5 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 202 2 2025 25 75 5 5 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297	ACACIAS ARGANZUELA	2	202	2	2088	88	59	2	2	1
2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACITAS ARGANIZUELA 2 2022 2 2016 16 53 10 14 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACITAS ARGANIZUELA 2 2022 2 2092 92 69 11 13 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACITAS ARGANIZUELA 2 2022 2 2016 16 28 14 11 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACITAS ARGANIZUELA 2 2024 2 2016 16 28 14 11 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACITAS ARGANIZUELA 2 2024 2 2025 25 75 5 5 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297	2024-07-07 00:00: 2024-07-01	2024-07-01	10		15297					
ACACIAS ARGANZUELA 2 202 2 2016 16 53 10 14 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 202 2 2092 92 69 11 13 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 202 2 2016 16 28 14 11 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 2024 - 07-01 10 15297	ACACIAS ARGANZUELA	2	202	2	2088	88	32	5	10	1
2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 202 2 2092 92 69 11 13 2024-07-07 00:00: 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 202 2 2016 16 28 14 11 2024-07-07 00:00: 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 202 2 2016 16 28 14 11 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 202 2 2025 25 75 5 5 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297	2024-07-07 00:00: 2024-07-01	2024-07-01	10		15297					
ACACIAS ARGANZUELA 2 202 2 2092 92 69 11 13 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 202 2 2016 16 28 14 11 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 2024-07-01 2024-07-01 5 5 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297	· · · · · · · · · · · · · · · · · · ·		202	2	2016	16	53	10	14	0
2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACITAS ARGANZUELA 2 2022 2 2016 16 28 14 11 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACITAS ARGANZUELA 2 202 2 2025 25 75 5 5 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297	the state of the s									
ACACIAS ARGANZUELA 2 202 2 2016 16 28 14 11 2024-07-07 00:00: 2024-07-07 10 15297	· ·			2		92	69	11	13	0
2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297 ACACIAS ARGANZUELA 2 202 2 2025 25 75 5 5 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297	the state of the s									-1
ACACIAS ARGANZUELA 2 202 2 2025 25 75 5 5 2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297				2	the second second	16	28	14	11	3
2024-07-07 00:00: 2024-07-01 2024-07-01 10 15297	the state of the s							-1	-1	-1
				2		25	/5	5	5	0
	2024-07-07 00:00: 2024-07-01									

only showing top to rows

• 6.13) Repite la función anterior utilizando funciones de ventana. (over(Window.partitionBy.....)).

Ventanas (Window)

- Propósito: Permiten realizar cálculos avanzados como agregados y rankings dentro de subconjuntos de datos.
- Uso: Definidas con Window y aplicadas usando funciones como SUM() o ROW_NUMBER().
- Ejemplo: Calcular una suma acumulativa dentro de cada grupo.

Particiones (partitionBy)

- Propósito: Organizan datos en disco para mejorar la eficiencia de lectura y escritura.
- Uso: Definidas al escribir datos usando partitionBy, dividiendo los datos en subdirectorios.
- Ejemplo: Almacenar datos en diferentes carpetas basadas en valores de columnas.

Diferencia Clave:

- Ventanas: Para cálculos analíticos dentro de datos.
- Particiones: Para organizar datos en almacenamiento distribuido.

```
from pyspark.sql.window import Window
 # Definir la ventana de partición
 windowSpec = Window.partitionBy("DESC_BARRIO", "DESC_DISTRITO")
 # Crear una nueva columna con la suma total de ESPANOLESHOMBRES por partición
 df_with_sum = df.withColumn("Total_EspanolesHombres", sum(col("ESPANOLESHOMBRES")).over(windowSpec))
 # Mostrar el esquema y las primeras filas para verificar
 df_with_sum.printSchema()
df_with_sum.show(10)
(3) Spark Jobs
 ▶ ■ df_with_sum: pyspark.sql.dataframe.DataFrame = [COD_DISTRITO: string, DESC_DISTRITO: string ... 15 more fields]
                                                                0 0
                                               6|
0|2024-07-07 00:00:...| 2024-07-01| 2024-07-01| 6
| G| TETUAN| 0| null|
0|2024-07-07 00:00:...| 2024-07-01| 2024-07-01| 6
1
                                                                   2094
                                                                              94
                                                                                           47
                                                                                                          6
                                                                                                                        16
                                                                                                                                         0
                                                                    2108
                                                                              108
                                                                                           5
                                                                    2019 19
                                                                                                                         7
                                                                                           81
                                                                                                          0|
                                                                   2018 18
                                                                                          34
                                                                                                                         4
                                                               15297|
2088|
                                                                                                          7
                                                                                                                                         2
                                                                              88
                                                                                           59
                                                                                                          2
                                                                                                                         2
                                                                                                                                         1|
                                                               15297|
2088|
                                                                              88|
                                                                                           321
                                                                                                          5
                                                                                                                        10
                                                                                                                                         1
                                                                   2016
1|2024-07-07 00:00:...| 2024-07-01| 2024-07-01|
```

```
,
Unir el DataFrame original con el DataFrame con la columna de suma
              df joined = df.join(df with sum.select("DESC BARRIO", "DESC DISTRITO", "Total EspanolesHombres"),
                                                                                                                                      on=["DESC_BARRIO", "DESC_DISTRITO"],
how="inner")
                # Mostrar el esquema y las primeras filas para verificar la unión
              df ioined.printSchema()
              df joined.show(10)
        ▶ ■ df_joined: pyspark.sql.dataframe.DataFrame = [DESC_BARRIO: string, DESC_DISTRITO: string ... 15 more fields]
  | 2024-07-07 00:00:...| 2024-07-01| 2024-07-01| 16
| ACACIAS | ARGANZUELA| 2 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2
                                                                                                                                                                                                                                                                                                                                                                                                                                 2 |
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   2094
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     6
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                    16
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              0|
                                                                                                                                                                                                                                                                                                                                                                                   | 2 | 2094| | 15297| | 02| | 2 | 2094| | 15297| | 02| | 2 | 2094| | 10| | 15297| | 02| | 2 | 2094| | 10| | 15297| | 02| | 2 | 2094| | 10| | 15297| | 02| | 2 | 2094| | 10| | 15297| | 02| | 2 | 2094| | 10| | 15297| | 02| | 2 | 2094| | 10| | 15297| | 02| | 2 | 2094| | 10| | 15297| | 02| | 2 | 2094| | 10| | 15297| | 02| | 2 | 2094| | 10| | 15297| | 02| | 2 | 2094| | 10| | 15297| | 02| | 2 | 2094| | 10| | 15297| | 02| | 2 | 2094| | 10| | 15297| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| | 02| 
                                                                                                                                                                                                                                                                                                                                                                                           101
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               15297
  | ACACIAS| ARGANZUELA| 2| 202| 0|2024-07-07 00:00:...| 2024-07-01| 2024-07-01| 10
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                          47
0|2024-07-09 00:00:...| 2024-07-01| 2024-07-01| 10
| ACACIAS| ARGANIZUELA| 2| 202|
0|2024-07-07 00:00:...| 2024-07-01| 2024-07-01| 10
| ACACIAS| ARGANIZUELA| 2| 202|
0|2024-07-07 00:00:...| 2024-07-01| 2024-07-01| 10
| ACACIAS| ARGANIZUELA| 2| 202|
0|2024-07-07 00:00:...| 2024-07-01| 2024-07-01| 10
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               2094 94
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               47
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   61
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                  16
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                  94
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   47
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   6
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                  16
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              0
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                 2094| 94|
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                 47
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                 6
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            2094 94
  | ACACIAS| ARGANZUELA| 2| 2021| 0|2024-07-01 0| 10| | ACACIAS| ARGANZUELA| 2| 2024-07-01| 10| | ACACIAS| ARGANZUELA| 2| 2021| 0|2024-07-07 00:00:...| 2024-07-01| 2024-07-01| 10|
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               47
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   6
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                  16
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              0
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                    94
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   47
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   6
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                  16
  | 0|2024-07-07 00:00:...| 2024-07-01| 2024-07-01| 10| 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 15297| | 152
  only showing top 10 rows
```

• 6.14) Mediante una función Pivot muestra una tabla (que va a ser una tabla de contingencia) que contenga los valores totales ()la suma de valores) de espanolesmujeres para cada distrito y en cada rango de edad (COD_EDAD_INT). Los distritos incluidos deben ser únicamente CENTRO, BARAJAS y RETIRO y deben figurar como columnas . El aspecto debe ser similar a este:

La función pivot en PySpark se usa para reorganizar los datos de una tabla, convirtiendo valores únicos de una columna en nuevas columnas. Permite realizar cálculos agregados, como sumas o promedios, para cada combinación de valores en las columnas de origen.

- pivot: Transforma los valores únicos de una columna en encabezados de nuevas columnas.
- Uso común: Crear tablas de contingencia que muestran resultados agregados (por ejemplo, sumas o promedios) organizados en una forma tabular más clara.

Por ejemplo, al pivotar por DESC_DISTRITO, cada distrito se convierte en una columna, y se puede calcular la suma de ESPANOLESMUJERES para cada COD_EDAD_INT en esas nuevas columnas.

```
from pyspark.sql.functions import col. sum
 # Filtrar los datos para los distritos específicos
 # isin() se utiliza para verificar si los valores en una columna de un DataFrame coinciden con cualquiera de los valores proporcionados en una lista dada. Si un valor
 en la columna del DataFrame se encuentra en la lista, devuelve verdadero; de lo contrario, devuelve falso. Esta función es útil para filtrar datos basados en valores
 específicos en los que estás interesado.
 df filtered = df.filter(col("DESC DISTRITO").isin("CENTRO", "BARAJAS", "RETIRO"))
 # Crear la tabla de contingencia (pivot table)
pivot_df = df_filtered.groupBy("COD_EDAD_INT").pivot("DESC_DISTRITO").agg(
     sum("ESPANOLESMUJERES").alias("Total_EspanolesMujeres"
 pivot_df.show()
▶ (7) Spark Jobs
▶ 🔳 df_filtered: pyspark.sql.dataframe.DataFrame = [COD_DISTRITO: string, DESC_DISTRITO: string ... 13 more fields]
▶ ■ pivot_df: pyspark.sql.dataframe.DataFrame = [COD_EDAD_INT: string, BARAJAS: long ... 2 more fields]
|COD EDAD INT|BARAJAS|CENTRO|RETIRO|
           51
                  271
                          255
                                  445
                                  435
                  219
                                  854
                  323
                          768
                                  731
                  260
                          492
                                  767
                  151
                          181
                                  314
                                  601
                          837
```

• 6.15) Utilizando este nuevo DF, crea 3 columnas nuevas que hagan referencia a qué porcentaje de la suma de "espanolesmujeres" en los tres distritos para cada rango de edad representa cada uno de los tres distritos. Debe estar redondeada a 2 decimales. Puedes imponerte la condición extra de no apoyarte en ninguna columna auxiliar creada para el caso.

```
from pyspark.sql.functions import col, round
 # Calcular la suma total de ESPANOLESMUJERES para cada rango de edad
 pivot_df = pivot_df.withColumn("Total_Sum", col("CENTRO") + col("BARAJAS") + col("RETIRO"))
 # Calcular los porcentajes de cada distrito respecto al total y redondear a 2 decimales
 pivot_df = pivot_df.withColumn("Centro_Percentage", round((col("CENTRO") / col("Total_Sum")) * 100, 2))
 pivot_df = pivot_df.withColumn("Barajas_Percentage", round((col("BARAJAS") / col("Total_Sum")) * 100, 2))
 pivot_df = pivot_df.withColumn("Retiro_Percentage", round((col("RETIRO") / col("Total_Sum")) * 100, 2))
 # Mostrar el resultado
 pivot df.show()
(3) Spark Jobs
▶ ■ pivot_df: pyspark.sql.dataframe.DataFrame = [COD_EDAD_INT: string, BARAJAS: long ... 6 more fields]
+------
|COD_EDAD_INT|BARAJAS|CENTRO|RETIRO|Total_Sum|Centro_Percentage|Barajas_Percentage|Retiro_Percentage|
+-----+
        51 449 803 805 2057
                                        39.04
                                                           21.83
                                                                          48.94
                                            25.06
        7 221 213 416 850
                                          25.06|

26.26|

40.04|

26.89|

35.05|

52.22|

42.15|

32.12|

32.39|

40.21|

28.02|

51.64|

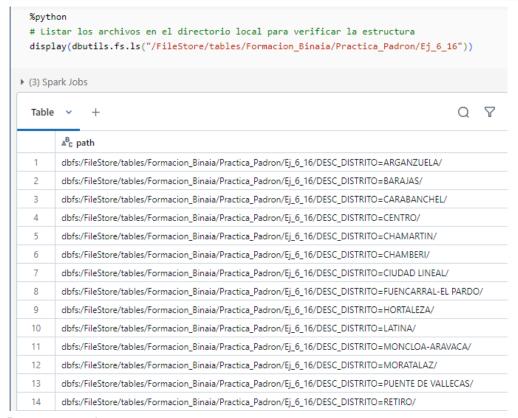
50.91|

40.07|
                                                             26.0
                                971
                                                                          45.83
        15 | 271 | 255 | 445 |
                                                           27.91
                                                                          40.51
        54 372 766 775 1913
                                                           19.45
        11 | 242 | 249 | 435 |
69 | 219 | 579 | 854 |
                                                                           46.98
                                                           26.13
                                926
                                1652
                                                           13.26
                                                                           51.69
                                                            10.47|
17.73|
12.73|
17.12|
                                                                           37.32
        29
             170 848 606 1624
                        731 1822
                                                                           40.12|
55.15|
             323 768
        42
                        455
             105
                   265
                                825
        87
                        767|
                               1519
        73
             260
                   492
                                                            13.95
             255 735 838
                                                                          45.84
                               1828
        64
                                                            23.37
             151 181
                        314
                                                                           48.61
        3
                                646
                                                                           36.56
             194 849 601
                               1644
                                                             11.8
        30
                                                           12.47
             205 | 837 | 602 |
                               1644
                                                                           36.62
        34
            362 791 821 1974
                                             40.07
                                                            18.34
                                                                          41.59
        59
                                                        18.34<sub>|</sub>
28.49|
        8 239 219 381
                                             26.1
                                839
                                                                          45.41
```

• 6.16) Guarda el archivo csv original particionado por distrito y por barrio (en ese orden) en un directorio local. Consulta el directorio para ver la estructura de los ficheros y comprueba que es la esperada.

Cada distrito, tendrá sus barrios.

%python
Guardar el DataFrame particionado por DESC_DISTRITO y DESC_BARRIO en un directorio local
df.write.partitionBy("DESC_DISTRITO", "DESC_BARRIO").csv("/FileStore/tables/Formacion_Binaia/
Practica_Padron/Ej_6_16", header=True)



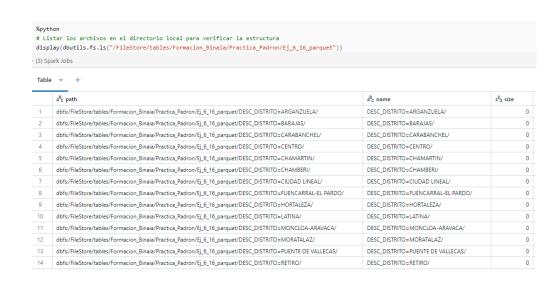
Da como esperaba.

• 6.17) Haz el mismo guardado pero en formato parquet. Compara el peso del archivo con el resultado anterior.

%python

Guardar el DataFrame particionado por DESC_DISTRITO y DESC_BARRIO en un directorio local en formato Parquet

df.write.partitionBy("DESC_DISTRITO", "DESC_BARRIO").parquet("/fileStore/tables/Formacion_Binaia/Practica_Padron/Ej_6_16_parquet", mode='overwrite')



Ahora compararemos los pesos de los archivos.

```
%python
 # Listar los archivos y subdirectorios en el directorio CSV
 csv dir = "/FileStore/tables/Formacion Binaia/Practica Padron/Ej 6 16/"
 csv_files = dbutils.fs.ls(csv_dir)
 # Función para calcular el tamaño total de los archivos en un directorio, incluyendo subdirectorios
 def calculate_size(dir_path):
     files = dbutils.fs.ls(dir_path)
     total_size = 0
     for file in files:
         if file.isDir():
            total_size += calculate_size(file.path)
         else:
            total_size += file.size
     return total_size
 # Calcular el tamaño total de los archivos CSV
 csv total size = calculate size(csv dir)
 print(f"Tamaño total de los archivos CSV: {csv_total_size} bytes")
Tamaño total de los archivos CSV: 19014515 bytes
  %python
  parquet_dir = "/FileStore/tables/Formacion_Binaia/Practica_Padron/Ej_6_16_parquet/
  parquet_files = dbutils.fs.ls(parquet_dir)
  # Calcular el tamaño total de los archivos Parquet
  parquet_total_size = calculate_size(parquet_dir)
  print(f"Tamaño total de los archivos Parquet: {parquet total size} bytes")
```

Se puede ver que en parquet, nos ahorramos casi 13 millones de bytes, por lo que es mucho más eficiente, gasta mucha menos memoria y es mucho más rápido para trabajar.

Tamaño total de los archivos Parquet: 6463347 bytes

Por otra parte, el código para CSV es más largo porque la estructura de archivos CSV generada por Spark incluye múltiples subdirectorios y archivos para cada partición, requiriendo un recorrido recursivo para calcular su tamaño. En cambio, los archivos Parquet tienen una estructura más compacta y eficiente, lo que simplifica tanto el almacenamiento como el cálculo de su tamaño.

7- ¿Y si juntamos Spark y Hive?

• 7.1) Por último, prueba a hacer los ejercicios sugeridos en la parte de Hive con el csv "Datos Padrón" (incluyendo la importación con Regex) utilizando desde Spark EXCLUSIVAMENTE sentencias spark.sql, es decir, importar los archivos desde local directamente como tablas de Hive y haciendo todas las consultas sobre estas tablas sin transformarlas en ningún momento en DataFrames ni DataSets

Creo sesión spark:

```
% 11:40 AM (2s)
%python
from pyspark.sql import SparkSession
from pyspark.sql.functions import col, trim, when

spark = SparkSession.builder.appName("Spark y Hive").getOrCreate()
```

Creo la tabla sin corregirla:

```
%python
# Crear la base de datos si no existe
spark.sql("CREATE DATABASE IF NOT EXISTS datos_padron")
# Usar la base de datos
spark.sql("USE datos_padron")
# Crear la tabla padron_txt con formato delimitado
spark.sql("""
   CREATE TABLE IF NOT EXISTS padron_txt (
      COD_DISTRITO INT,
       DESC_DISTRITO STRING,
       COD_DIST_BARRIO INT,
       DESC_BARRIO STRING,
       COD_BARRIO INT,
       COD_DIST_SECCION INT,
       COD_SECCION INT,
       COD_EDAD_INT STRING,
       ESPANOLESHOMBRES INT,
       ESPANOLESMUJERES INT,
       EXTRANJEROSHOMBRES INT,
       EXTRANJEROSMUJERES INT,
       FX_CARGA STRING,
       FX_DATOS_INI STRING,
       FX_DATOS_FIN STRING
   ROW FORMAT DELIMITED
   FIELDS TERMINATED BY ';'
   LINES TERMINATED BY '\n'
   STORED AS TEXTFILE
   TBLPROPERTIES ("skip.header.line.count"="1");
```

Inserto los datos en la tabla, copiando el csv en otra ubicación, porque si ejerzo el comando de "LOAD DATA", se elimina de la ubicación local:

```
%python
dbutils.fs.cp('/FileStore/tables/Formacion_Binaia/Practica_Padron/padron_estadistica.csv', '/FileStore/tables/Formacion_Binaia/Practica_Padron/padron.csv')
spark.sql("""

LOAD DATA INPATH '/FileStore/tables/Formacion_Binaia/Practica_Padron/padron.csv'

INTO TABLE padron_txt
""")
```

Hago la comprobación de que todo ha salido bien:

%python spark.sql("SELECT * FROM padron_txt LINC	ET 5"), show()										
(1) Spark Jobs											
(1) Spank Jobs											
COD_DISTRITO DESC_DISTRITO COD_DISS_FIN	ST_BARRIO DESC_B	ARRIO COD_	BARRIO COD_DI	ST_SECCION COD_	SECCION COD_EDA	ND_INT ESP	ANOLESHOMBRES ESPANOLE	SMUJERES EXTRAN	JEROSHOMBRES EXTRANJER	OSMUJERES	FX_CARGA FX_DATOS_INI FX
null DESC_DISTRITO S_FIN	null DESC_B	ARRIO	null	null	null COD_EDA	ND_INT	null	null	null	null	FX_CARGA FX_DATOS_INI FX
2 "ARGANZUELA 7-01"	202 "ACACIAS		2	2894	94 "	47"	6	16	9	0 "2024-0	-87 88:88 "2824-87-81" "2
18 "VILLA DE VALLECA 7-01"	1802 "SANTA EUGENIA	1	2	18029	29 "	77"	7	3	0	2 "2024-0"	-07 00:00 "2024-07-01" "2
2 "ARGANZUELA 7-01"	205 "DELICIAS	1	5	2862	62 "	71"	6	6	8	2 "2824-8	-87 88:88 "2824-87-81" "2
9 "MONCLOA-ARAVACA 7-01"	987 "ARAVACA	1	7	9883	83 "	55"	15	14	0	0 "2024-0	-07 00:00 "2024-07-01" "2

Ahora creo la tabla padron_txt_2 mediante CTAS, para crear una tabla sin espacios en blanco y sin comillas dobles:

```
%python
from pyspark.sql.functions import trim, regexp_replace, col
spark.sql("""
   CREATE TABLE IF NOT EXISTS padron_txt_2 AS
    SELECT
        CAST(trim(regexp_replace(COD_DISTRITO, '"', '')) AS INT) AS COD_DISTRITO,
        trim(regexp_replace(DESC_DISTRITO, '"', '')) AS DESC_DISTRITO,
        CAST(trim(regexp_replace(COD_DIST_BARRIO, '"', '')) AS INT) AS COD_DIST_BARRIO,
        trim(regexp_replace(DESC_BARRIO, '"', '')) AS DESC_BARRIO,
        CAST(trim(regexp_replace(COD_BARRIO, '"', '')) AS INT) AS COD_BARRIO,
        CAST(trim(regexp_replace(COD_DIST_SECCION, '"', '')) AS INT) AS COD_DIST_SECCION,
        CAST(trim(regexp_replace(COD_SECCION, '"', ''')) AS INT) AS COD_SECCION,
        trim(regexp_replace(COD_EDAD_INT, '"', '')) AS COD_EDAD_INT,
        CAST(trim(regexp_replace(ESPANOLESHOMBRES, '"', '')) AS INT) AS ESPANOLESHOMBRES,
        CAST(trim(regexp_replace(ESPANOLESMUJERES, '"', '')) AS INT) AS ESPANOLESMUJERES,
       CAST(trim(regexp_replace(EXTRANJEROSHOMBRES, '"', '')) AS INT) AS EXTRANJEROSHOMBRES,
        CAST(trim(regexp_replace(EXTRANJEROSMUJERES, '"', '')) AS INT) AS EXTRANJEROSMUJERES,
        trim(regexp_replace(FX_CARGA, '"', '')) AS FX_CARGA,
        trim(regexp_replace(FX_DATOS_INI, '"', '')) AS FX_DATOS_INI,
        trim(regexp_replace(FX_DATOS_FIN, '"', '')) AS FX_DATOS_FIN
    FROM padron_txt;
```

Hago de nuevo comprobación y veo que ha ido bien:

```
spark.sql("SELECT * FROM padron_txt_2 LIMIT 5").show()
(2) Spark Jobs
    [COD_DISTRITO] DESC_DISTRITO[COD_DIST_BARRIO] DESC_BARRIO[COD_BARRIO] COD_DIST_SECCION[COD_SECCION[COD_EDAD_INT|ESPANOLESHOMBRES|ESPANOLESMUJERES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|EXTRANJEROSHOMBRES|
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               FX_CARGA|FX_DATOS_INI|FX_DATOS_FIN|
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     null|COD_EDAD_INT|
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                           mull| (2024-07-07 08:00:...| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01| 2024-07-01|
                                                                                                                                                          ARGANZUELA
                                                                                                                                                                                                                                                                                                                                  202
                                                                                                                                                                                                                                                                                                                                                                                                           ACACIAS
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               2094
                                                                               18 VILLA DE VALLECAS
                                                                                                                                                                                                                                                                                                                       1802 SANTA EUGENTAL
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      18029
                                                                                      2| ARGANZUELA|
9| MONCLOA-ARAVACA|
                                                                                                                                                                                                                                                                                                                                                                                                         ARAVACA
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            9083
```

Ahora lo hago con expresiones regulares:

```
%python
#expresiones regulares
spark.sql(""'
CREATE TABLE padron_hive_regex (
 COD DISTRITO STRING.
 DESC_DISTRITO STRING,
 COD_DIST_BARRIO STRING,
 DESC_BARRIO STRING,
 COD BARRIO STRING,
 COD_DIST_SECCION STRING,
 COD_SECCION STRING,
 COD_EDAD_INT STRING,
 ESPANOLESHOMBRES STRING,
 ESPANOLESMUJERES STRING,
 EXTRANJEROSHOMBRES STRING,
 EXTRANJEROSMUJERES STRING,
 FX_CARGA STRING,
 FX_DATOS_INI STRING,
 FX_DATOS_FIN STRING
ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.RegexSerDe'
WITH SERDEPROPERTIES (
 STORED AS TEXTFILE;
```

Meto los datos:

```
%python
dbutils.fs.cp('/FileStore/tables/Formacion_Binaia/Practica_Padron/padron_estadistica.csv', '/FileStore/tables/Formacion_Binaia/Practica_Padron/padron.csv')
spark.sql("""
    LOAD DATA INPATH '/FileStore/tables/Formacion_Binaia/Practica_Padron/padron.csv'
    INTO TABLE padron_hive_regex
""")
```

Y compruebo que ha salido bien:

%python spark.so	q1("SELECT * FROM padro	_hive_regex	LIMIT 5").show()												
(1) Spark	Jobs														
COD_DISTF	RITO DESC_DISTRI	ro cob_bist_B	ARRIO DESC_	BARRIO COD_	BARRIO COD_DI	ST_SECCION COD_	SECCION COD_	EDAD_INT ESPAN	DLESHOMBRES ESPANOL	ESMUJERES EXTRANJ	EROSHOMBRES EXTRANJE	ROSMUJERES	FX_CARGA F	x_DATOS_INI	FX_DAT
r r null	null nu. 2 ARGANZUELA		null 202 ACACIAS	null	null	null 2094	nu11 94	null	null	null	null 0	null 0 2024-07-07	null	null	
-07-01 -07-01	18 VILLA DE VALLECAS	1	1802 SANTA EUGENIA	1	2	18029	29	77	7	3	0	2 2024-07-07			
 -07-01 -07-01	2 ARGANZUELA 9 MONCLOA-ARAVACA	I	205 DELICIAS 907 ARAVACA	1	7	9083	62 83	71 55	15	6	0	2 2024-07-07 0 2024-07-07			
*		-+	+	+	+									+	