



**ESCUELA TÉCNICA SUPERIOR DE INGENIERÍA INFORMÁTICA**

**GRADO EN INGENIERÍA INFORMÁTICA**

**Curso Académico 2022/2023**

**Trabajo Fin de Grado**

**INGENIERÍA DE DATOS CON EL FRAMEWORK DE BIG DATA  
SPARK Y SCALA**

**Autor:** Álvaro Sánchez Pérez

**Directores:** Juan Manuel Serrano Hidalgo



## Índice

<b>1. Introducción .....</b>	<b>3</b>
<b>2. Objetivos .....</b>	<b>3</b>
<b>3. Descripción informática .....</b>	<b>4</b>
1. Fuentes de datos .....	4
1.1. Obtención de los datos.....	8
2. Programación de queries.....	10
3. Visualización de queries .....	17
4. Despliegue en AWS EMR .....	22
<b>2. Experimentos / validación .....</b>	<b>22</b>
2.1. Consultas realizadas .....	22
2.2. Análisis de requisitos no funcionales .....	22
<b>3. Conclusiones .....</b>	<b>22</b>
<b>4. Bibliografía .....</b>	<b>22</b>
<b>5. Apéndices .....</b>	<b>22</b>



# 1. Introducción

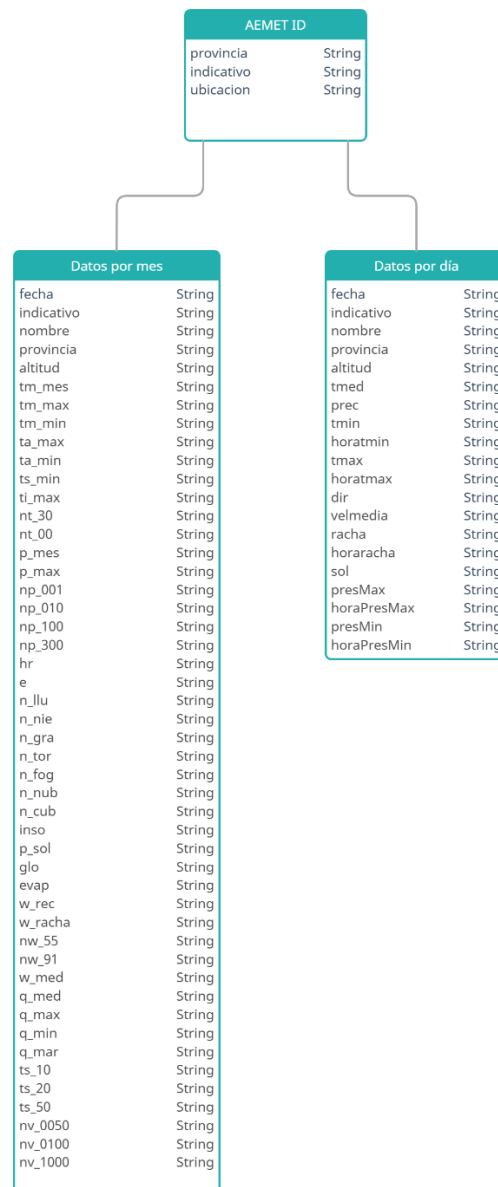
Prueba

## 2. Objetivos

### 3. Descripción informática

#### 1. Fuentes de datos

Los datos meteorológicos, han sido obtenidos a través de la página de *AEMET OpenData*, se tratan tanto de datos que muestran la información de forma mensual como de forma diaria desde el año 2007 hasta el año actual 2022.



Resaltar que el diagrama mostrado, aparecen los id de las columnas y valores que deberían de encontrarse en estas mediante los datos obtenidos, ya que así nos lo indican los metadatos que nos proporcionan. Pero en cambio tanto el nombre, como la provincia no aparecen en la mayoría de los archivos de datos que nos proporcionan, y por ello debemos de hacer uso de la tabla *AEMET ID*, que gracias a esta podremos relacionar mediante el indicativo de que provincia y ubicación se trata.

La estructura en la que vienen estos datos se puede observar en el anterior diagrama, todos los valores vienen en formato String, los cuales transformaremos mediante el casteo correspondiente gracias a Spark, para así poder trabajar posteriormente con ellos de una manera más cómoda. Debemos de realizar lo siguiente para la transformación.

1. Primero creamos el esquema correspondiente a los datos, en este ejemplo lo realizaremos con los datos por día. En este esquema, le indicaremos tanto en nombre de la columna, como el tipo de dato que se trata, en nuestro caso todos String, y por último si la columna puede contener valores nulos.

```
val schema = StructType(
  Array(
    StructField("fecha", StringType, true),
    StructField("indicativo", StringType, true),
    StructField("nombre", StringType, true),
    StructField("provincia", StringType, true),
    StructField("altitud", StringType, true),
    StructField("tmed", StringType, true),
    StructField("prec", StringType, true),
    StructField("tmin", StringType, true),
    StructField("horatmin", StringType, true),
    StructField("tmax", StringType, true),
    StructField("horatmax", StringType, true),
    StructField("dir", StringType, true),
    StructField("velmedia", StringType, true),
    StructField("racha", StringType, true),
    StructField("horaracha", StringType, true),
    StructField("sol", StringType, true),
    StructField("presMax", StringType, true),
    StructField("horaPresMax", StringType, true),
    StructField("presMin", StringType, true),
    StructField("horaPresMin", StringType, true)
  )
)
```

2. A continuación, realizaremos la lectura de los datos y casteo de los datos.

```
val allData = spark.read.schema(schema).json("D:/TFGAlvaroSanchez/data/dayJSONLine/*.json")
  .withColumn("fecha", $"fecha".cast(DateType))
  .withColumn("altitud", $"altitud".cast(IntegerType))
  .withColumn("tmed", func.regex_replace($"tmed", ",", ".").cast(DoubleType))
  .withColumn("prec", func.regex_replace($"prec", ",", ".").cast(DoubleType))
  .withColumn("tmin", func.regex_replace($"tmin", ",", ".").cast(DoubleType))
  .withColumn("tmax", func.regex_replace($"tmax", ",", ".").cast(DoubleType))
  .withColumn("dir", $"dir".cast(IntegerType))
  .withColumn("velmedia", func.regex_replace($"velmedia", ",", ".").cast(DoubleType))
  .withColumn("racha", func.regex_replace($"racha", ",", ".").cast(DoubleType))
  .withColumn("sol", func.regex_replace($"sol", ",", ".").cast(DoubleType))
  .withColumn("presMax", func.regex_replace($"presMax", ",", ".").cast(DoubleType))
  .withColumn("horaPresMax", $"horaPresMax".cast(IntegerType))
  .withColumn("presMin", func.regex_replace($"presMin", ",", ".").cast(DoubleType))
  .withColumn("horaPresMin", $"horaPresmin".cast(IntegerType))
```

- 2.1. Para la lectura le proporcionaremos el esquema creado anteriormente en la variable *schema*, y a su vez le indicamos en la ubicación donde se encuentran los datos, además del formato de fichero.

```
spark.read.schema(schema).json("D:/TFGAlvaroSanchez/data/dayJSONLine/*.json")
```

- 2.2. Finalmente, modificaremos el tipo de dato de las columnas que lo requieran. Usaremos *withColumn* donde primero le indicaremos el nuevo nombre de la columna que nos generará, en este caso usaremos el mismo nombre de la comuna que vamos a modificar para que así ocurra una sustitución de esta, y en la segunda parte será donde le indiquemos la columna con la que queremos trabajar y el cambio que deseamos.

```
.withColumn("fecha", $"fecha".cast(DateType))
```

He de mencionar que, los datos que se tuvieron que transformar a tipo Double se debió realizar una modificación previa, sustituyendo las comas por puntos, ya que de otra manera no resultaba posible realizar el casteo debido a que para que se considere un numero decimal debe de venir separada la parte entera de la decimal mediante un punto.

```
.withColumn("tmed", func.regex_replace($"tmed", ",", ".").cast(DoubleType))
```

3. Por último, guardaremos los datos en formato Parquet, particionados mediante el campo *indicativo*.

```
allData.write.format("parquet").partitionBy("indicativo").mode("overwrite").save("D:/TFGAlvaroSanchez/data/dayParquet/")
```

Quedando particionado en memoria de la siguiente manera:

indicativo=0002I	29/10/2022 19:52	Carpeta de archivos
indicativo=0016A	29/10/2022 19:52	Carpeta de archivos
indicativo=0076	29/10/2022 19:52	Carpeta de archivos

Y dentro de cada carpeta algo similar a lo siguiente:

.part-00055-ef893175-2d23-4949-a94b-6cb4b59574dc....	29/10/2022 19:51	Archivo CRC	1 KB
.part-00079-ef893175-2d23-4949-a94b-6cb4b59574dc....	29/10/2022 19:51	Archivo CRC	1 KB
.part-00081-ef893175-2d23-4949-a94b-6cb4b59574dc....	29/10/2022 19:51	Archivo CRC	1 KB
.part-00031-ef893175-2d23-4949-a94b-6cb4b59574dc....	29/10/2022 19:51	Archivo PARQUET	17 KB
.part-00040-ef893175-2d23-4949-a94b-6cb4b59574dc....	29/10/2022 19:51	Archivo PARQUET	25 KB

Respecto a los datos mensuales nos quedaríamos únicamente con las siguientes columnas.

```
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| fecha|indicativo| p_max| glo| hr|nw_55|tm_min| ta_max|ts_min|nt_30|n_des| w_racha|np_100|nw_91|np_001| ta_min|w_rec| e|np_300|p_mes|w_med|nt_00|ti_max|tm_mes|tm_max|np_010|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|2014-01-01|      B954|7.5(16)|26493.0|73.0| 12|  9.2|21.4(03)| 13.2|  0| 1|28/25.8(04)|  0|  1|  6|3.8(31)|
349|115.0|  0| 26.3| 15.0|  0| 14.6| 13.3| 17.4|  4|
|2014-02-01|      B954|4.1(09)|33644.0|72.0| 12|  8.3|20.4(13)| 13.0|  0| 4|26/21.7(10)|  0|  0|  9|3.0(03)|
457|110.0|  0| 10.8| 19.0|  0| 14.4| 13.0| 17.7|  3|
```



De las que podríamos destacar las siguientes: *tm\_mes* la cual nos muestra una temperatura media mensual de la ubicación, *tm\_max* y *tm\_min* siendo las temperaturas medias máximas y mínimas respectivamente, *ta\_max* y *ta\_min* como las temperaturas máximas y mínimas absolutas del mes y *p\_mes* la cual nos indica precipitación total en ese mes. En caso de querer obtener información acerca de alguna otra variable, se puede acceder a esta de las siguientes maneras: a través del fichero *metadataMonth.json* incluido en la carpeta *data* del proyecto o realizando cualquier consulta sobre climatologías mensuales/anuales en la página de *AEMET OpenData* donde obtendrá un enlace acerca de los metadatos donde encontrar este tipo de información.

En relación con los datos diarios, tendríamos lo siguiente:

+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+													
fecha indicativo		nombre provincia		altitud	tmed	prec	tmin	horatmin	tmax	horatmax	dir	velmedia	racha
sol presMax horaPresMax		presMin horaPresMin		+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+									
2010-01-01	1387	A CORUÑA	A CORUÑA	58	9.7	3.3	6.2	05:30	13.2	12:40	23	3.9	10.8
3.3	null	null	null	null									
2010-01-01	1387E	A CORUÑA AEROPUERTO	A CORUÑA	98	7.8	0.7	4.0	08:58	11.5	13:17	23	null	10.3
3.9	1003.1	23	987.4	4									

En este caso nos hemos quedado con todas las columnas, de las que resaltaríamos las siguientes: *tmed* la cual muestra la temperatura media diaria, *prec* que nos ofrece información acerca de la precipitación diaria, *tmax* y *tmin* las cuales nos muestran la temperatura máxima y mínima diaria respectivamente, y *horatmax* y *horatmin* que, de forma correspondiente, presentan la hora y minuto de la temperatura máxima y mínima. En caso de querer obtener información acerca de alguna otra variable, se puede acceder a esta de las siguientes maneras: a través del fichero *metadataDay.json* incluido en la carpeta *data* del proyecto o realizando cualquier consulta sobre climatologías diarias en la página de *AEMET OpenData* donde obtendrá un enlace acerca de los metadatos donde encontrar este tipo de información

He de mencionar que, a la hora de descargar los datos se encuentran en un formato JSON multilínea, el cual no llega a ser lo suficientemente óptimo para la lectura con Apache Spark. Por lo tanto, se realizaron cambios en todos los ficheros transformándolos a ficheros JSON que contuvieran toda la información en una única línea, esto se llevó a cabo mediante el siguiente comando utilizando la consola de Windows: **FOR %a IN (../data/\*.json) DO jq . -c "%a" > "../JSONLine/%a"**. También cabe destacar que los datos una vez leídos en este formato JSON, la mayoría fueron transformados al formato Parquet, particionados por el indicativo de cada estación, ya que resulta más óptimo para la realización de consultas con Spark.

## 1.1. Obtención de los datos

Para la obtención de los datos, se crearon unos pequeños programas en el lenguaje Java. Estos programas, se basan en la simulación de los pasos que deberíamos de seguir para la descarga y los podemos encontrar en las carpetas *DescargaDatosPorDias* y *DescargaDatosPorMeses*. Para la simulación de estos pasos, se hizo uso de la librería *Selenium*, a través de la cual se puede manejar un navegador web pudiendo realizar diversas acciones sobre la página mostrada. Por lo tanto, de manera resumida, los pasos que realizarían ambos programas para la descarga serían los siguientes.

```
WebDriver driver = new ChromeDriver();  
driver.get(baseUrl);
```

Abriríamos un nuevo navegador donde accederemos al sitio web de *AEMET OpenData*. En nuestro caso le pasamos la URL del sitio web a través de la variable *baseUrl*.

```
driver.findElement(By.id("apikey")).sendKeys(apiKey);  
desplegable1 = new Select(driver.findElement(By.id("clim1")));  
desplegable1.selectByIndex(provincia);
```

Una vez dentro de esta página, el primer paso sería introducir la clave API correspondiente a nuestro usuario, está la podremos solicitar también a través de la misma página. Posteriormente, dependiendo del tipo de información que quisiéramos obtener, ya fuera información por días o por meses, buscaríamos los desplegables y elementos necesarios mediante su *id* o *xpath*.

Después de seleccionar la información deseada, se nos abriría una nueva página donde nos proporciona diferentes tipos de información, además del enlace de la página donde se encontrarán los datos deseados.



```
String texto = driver.findElement(By.xpath("//pre[contains(@style,'word-wrap')]")).getText();
if(texto.contains("\"descripcion\" : \"exito\"")) {
    String urlDatos = texto.split(regex: "\\n")[9];

    driver.get(urlDatos);
    texto = driver.findElement(By.xpath("//pre[contains(@style,'word-wrap')]")).getText();

    //Imprimimos la informacion en un fichero externo
    PrintWriter printWriter = null;
    String ubicacionGuardar = "D:\\TF6AlvaroSanchez\\data\\day\\";
    String nombreFichero = ubicacionGuardar.concat(estacionMeterologica).concat(" ").concat(ano).concat(" ");

    try {
        printWriter = new PrintWriter(nombreFichero);
    } catch (FileNotFoundException e) {
        System.out.println("Unable to locate the fileName: " + e.getMessage());
    }

    Objects.requireNonNull(printWriter).println(texto);
    printWriter.close();
}
```

He de destacar también, que al comienzo del código existen una serie de variables que pueden ser cambiadas por los usuarios. Desde introducir su correspondiente clave API, cambiar el controlador de Chrome en caso de que se esté usando una versión diferente del navegador, o cambiar las fechas en caso de que se quisieran obtener datos en un rango diferente.

9

## 2. Programación de queries

### **Año de la temperatura máxima promedio en cada mes**

En esta consulta, realizaremos una comparación del mismo mes en diferentes años, para saber en qué año se registró la temperatura máxima promedio entre todas las estaciones en este mes. Por último, lo representaremos tanto mediante una tabla, como de manera gráfica para poder observar las diferencias de manera visual entre los distintos años.

Para realizar esta consulta utilizaremos los datos meteorológicos mensuales, ya que nos ofrecen la temperatura máxima registrada en cada mes de una manera bastante sencilla. Por lo tanto, el primer paso que deberíamos de realizar sería la lectura de estos datos.

A continuación, deberíamos de agrupar estos datos por fecha, ya que, al realizar la lectura, tenemos la información de muchas estaciones meteorológicas diferentes en el mismo mes y año, por lo tanto, los agruparemos haciendo que coincida la fecha, y realizando una media de la temperatura máxima en todas las estaciones. Por último, dividiremos los datos por meses y ordenándolos de manera descendente a través de nuestro objetivo, la temperatura máxima, ya que así estaríamos consiguiendo obtener el año, en el cual este mes fue el más caluroso.

## Olas de calor

En esta consulta, desearemos obtener las diferentes olas de calor que han producido cada año en España, con el objetivo de observar si es verdad que se está produciendo un aumento de las temperaturas durante los últimos años. Para ello, calcularemos las olas de calor teniendo en cuenta la temperatura máxima diaria en cada estación meteorológica, y realizaremos los cálculos necesarios para finalmente mostrar los resultados sobre un mapa donde el usuario podrá elegir el año del cual desea obtener la información. Resaltar que, se considerará ola de calor cuando la temperatura máxima supere o iguale los 40 grados durante más de tres días consecutivos.

Debido a la exactitud que necesitamos para esta consulta se utilizarán los datos meteorológicos diarios, ya que requeriremos los datos de las temperaturas máximas con una gran exactitud en el marco temporal. Por lo tanto, el primer paso para comenzar con esta consulta sería realizar la lectura de estos datos.

Una vez leídos los datos, deberemos de eliminar aquellos donde no se llegue a alcanzar la temperatura establecida, en nuestro caso eliminaremos aquellos que tengan una temperatura menor a 40 grados. Posteriormente, deberemos de identificar posibles olas de calor, aun sin tener en cuenta la duración de estas, para ello dividiremos los datos por estación meteorológica y año, y a cada fecha se le asignara un identificador, en caso de que las fechas que sean consecutivas se les asignara el mismo identificador. Por ir resumiendo, ahora mismo lo que tendríamos serían los días que presentan temperaturas mayores o iguales a 40 grados, donde además se les habrá asignado un identificador, con el que posteriormente podremos determinar la duración de la ola de calor.

A continuación, agruparemos los datos por el identificador asignado a cada fecha, donde contaremos el número de apariciones de cada identificador, y por último eliminaremos aquellos no válidos, quedándonos únicamente con los que se encuentran más de 3 veces. Ya que como se comentó anteriormente, gracias al número de veces que se encuentre un identificador podremos determinar la duración de la ola de calor.

```
import org.apache.spark.sql.expressions.Window

val data = spark.read.parquet("D:/TFGAlvaroSanchez/data/dayParquet/")
val stations = spark.read.option("delimiter", ";").csv("D:/TFGAlvaroSanchez/data/aemetID.csv")
    .toDF("provincia", "indicativo", "ubicacion")

val window = Window.partitionBy($"indicativo", $"año").orderBy($"fecha")

val results = data
    .filter(!func.isnull($"tmax") && $"tmax" >= 40)
    .withColumn("año", func.year($"fecha"))
    .withColumn("n_fila", func.row_number().over(window))
    .withColumn("id", func.expr("date_sub(fecha, n_fila)"))
    .groupBy($"indicativo", $"año", $"id")
    .agg(func.count($"id").alias("dias"), func.avg($"tmax"), func.max($"tmax"), func.min($"tmax"))
    .filter($"dias" > 3)
    .join(stations, "indicativo")
    .select($"ubicacion", $"provincia", $"año", $"dias", func.round($"avg(tmax)", 2).alias("avg(tmax)"),
        $"max(tmax)", $"min(tmax)")
```

```

val resultsSave = results
  .groupBy($"provincia", $"año")
  .agg(func.count($"provincia"), func.avg($"dias"), func.avg($"avg(tmax)"), func.max($"max(tmax)"), func.min($"min(tmax)"))
  .select($"provincia", $"año", $"count(provincia)".alias("nº de olas de calor"),
    $"avg(dias)".alias("duracion media"), $"avg(avg(tmax))".alias("avg(tmax)"),
    $"max(max(tmax))".alias("max(tmax)"), $"min(min(tmax))".alias("min(tmax)"))

resultsSave
  .withColumnRenamed("nº de olas de calor", "nOlasCalor")
  .withColumnRenamed("duracion media", "duracionMedia")
  .withColumnRenamed("avg(tmax)", "avgTmax")
  .withColumnRenamed("max(tmax)", "maxTmax")
  .withColumnRenamed("min(tmax)", "minTmax")
  .write.format("parquet").partitionBy("provincia").mode("overwrite").save("D:/TFGAlvaroSanchez/data/resultadoOlasCalor/")

```

La consulta se podría dividir en dos partes, una primera donde calculamos las diferentes olas de calor, y una segunda donde se agruparían estas olas de calor por provincias y obtenemos la información necesaria para representarla posteriormente en forma de mapa.

El primer paso como se indicó anteriormente sería la lectura de los diferentes datos. Primero de ello, leeremos los datos meteorológicos, los cuales guardaremos en la variable *data*. Podemos observar únicamente 3 columnas ya que, al tratarse de una lectura de tipo Parquet, Spark únicamente lee las columnas que va a utilizar.

indicativo	fecha	tmax
C447A	2013-01-01	15.5
C447A	2013-01-02	15.8
C447A	2013-01-03	15.6
C447A	2013-01-04	16.3
C447A	2013-01-05	17.8

También leeremos el fichero a través del cual relacionaremos el indicativo de las estaciones con su provincia y ubicación. Estos datos los guardaremos en la variable *stations* y contendrá una información similar a la siguiente.

provincia	indicativo	ubicacion
A Coruna	1351	Estaca de Bares
A Coruna	1363X	As Pontes
A Coruna	1387	A Coruna
A Coruna	1387E	A Coruna Aeropuerto

```

val window = Window.partitionBy($"indicativo", $"año").orderBy($"fecha")

```

Además, al comienzo de la consulta nos crearemos nuestra función ventana en la variable *window*. La utilizaremos posteriormente puesto que, a través de esta dividiremos los datos mediante el año y su estación meteorológica, y a su vez ordenaremos estas particiones mediante la fecha de manera ascendente.

```
val results = data
  .filter(!func.isNull($"tmax") && $"tmax" >= 40)
```

Seguidamente comenzamos con la primera parte de la consulta donde obtendremos las distintas olas de calor en cada estación, guardaremos esta información en la variable *results*. El primero de los pasos será eliminar las filas que no utilizaremos, para ello realizamos un filtro en la columna *tmax*, la cual nos proporciona la información de la temperatura máxima, por ello eliminaremos aquellas filas que no contengan ningún tipo de información acerca de esta, y aquellas que no superen la temperatura mínima establecida.

```
.withColumn("año", func.year($"fecha"))
```

A continuación, haremos uso de la función de Spark *year* con la cual obtendremos el año a través de la fecha. Guardaremos esta nueva información en una nueva columna a la que llamaremos *año*.

```
.withColumn("n_fila", func.row_number().over(window))
```

A su vez, también crearemos una nueva columna a la cual llamaremos *n\_fila*, en esta le asignaremos un número a cada fila utilizando la función *row\_number*, la cual a su vez hará uso de la función ventana que hemos creado anteriormente. Mediante la función *row\_number*, asignaremos valores de manera secuencial comenzando desde el 1. Y gracias a la función ventana, estamos consiguiendo asignar un numero de fila que va irá en aumento dependiendo de la fecha ya que era el orden que le habíamos establecido. A su vez cada vez que se trate de un nuevo año o de una nueva estación comenzará a contar de nuevo desde el 1, ya que eran las divisiones que le habíamos indicado cuando nos creamos la función. Quedándonos la consulta por el momento de la siguiente manera.

indicativo	fecha	tmax	año	n_fila
0149X	2015-07-05	40.3	2015	1
0149X	2015-07-07	40.2	2015	2
0149X	2019-06-28	42.4	2019	1
0149X	2021-08-12	41.0	2021	1
0149X	2021-08-13	40.4	2021	2
1002Y	2020-07-30	40.0	2020	1
1002Y	2022-06-18	42.0	2022	1
1002Y	2022-07-17	40.7	2022	2
1002Y	2022-07-18	43.4	2022	3
1002Y	2022-07-24	41.1	2022	4

```
.withColumn("id", func.expr("date_sub(fecha, n_fila)"))
```

Continuaremos creándonos una nueva columna llamada *id*, en esta realizaremos una resta de la columna *fecha* con el *n\_fila*. De esta manera estaríamos consiguiendo identificar las distintas

olas de calor, ya que en caso de tratarse de días consecutivos en la nueva columna *id* contendrán el mismo valor. Por ejemplo, en la anterior imagen tenemos estas filas:

indicativo	fecha	tmax	año	n_fila
1002Y	2022-07-17	40.7	2022	2
1002Y	2022-07-18	43.4	2022	3

Como podemos observar se tratan de fechas consecutivas, por lo tanto, si realizamos la resta de la columna *fecha* – *n\_fila* nos quedaría lo siguiente en ambas:

indicativo	fecha	tmax	año	n_fila	id
1002Y	2022-07-17	40.7	2022	2	2022-07-15
1002Y	2022-07-18	43.4	2022	3	2022-07-15

Como podemos, observar tendrían el mismo valor en la columna *id* y gracias a este estaríamos identificando días consecutivos con temperaturas lo suficientemente altas. Por lo tanto, este proceso lo realizamos con todos nuestros datos, y tendríamos identificadas las distintas posibles olas de calor, ya que aun deberemos tener él cuenta el número de días que duran para que se consideren validas.

```
.groupBy($"indicativo", $"año", $"id")
.agg(func.count($"id").alias("dias"), func.avg($"tmax"), func.max($"tmax"), func.min($"tmax"))
```

Retomando la consulta, ahora agruparemos las distintas posibles olas de calor, para ello realizaremos *groupBy* con las columnas *indicativo*, *año* e *id*. Posteriormente, mediante la función *agg*, llamaremos a diferentes funciones Spark a través de las cuales obtendremos información. Por ejemplo, mediante la función *count* con la columna *id*, obtendremos la duración en días de las altas temperaturas, ya que estaremos contando el número de veces que podemos encontrar ese mismo *id* en la misma estación meteorológica y año, a su vez, haciendo uso de la función *alias* le establecemos el nombre a la nueva columna que se nos generará. Mediante las funciones *avg*, *max* y *min* obtendremos la temperatura media, máxima y mínima respectivamente acerca de esa posible ola de calor.

```
.filter($"dias" > 3)
```

Para ir finalizando con la primera parte de la consulta, realizamos un nuevo filtrado de los datos, en este caso quedándonos con lo que nosotros hemos considerado olas de calor, aquellas cuya duración es mayor a 3 días.

```
.join(stations, "indicativo")
```

Realizamos una unión mediante la función *join* con nuestros datos guardados en la variable *stations*, la cual contenía información extra acerca de las estaciones. Estos datos tenían tres



columnas, *provincia*, *indicativo* y *ubicación*, por lo tanto, los unimos mediante la columna común a ambas tablas, *indicativo*.

```
.select($"ubicacion", $"provincia", $"año", $"dias", func.round($"avg(tmax)", 2).alias("avg(tmax)"),  
      $"max(tmax)", $"min(tmax)"))
```

Por último, elegimos las columnas relevantes para nuestro resultado y mediante la función *round* redondeamos el valor perteneciente a la columna *avg(tmax)*, estableciéndole que contenga únicamente dos decimales.

Ahora mismo, tendríamos en la variable *results* las distintas olas de calor en cada estación meteorológica, viéndose de la siguiente manera:

ubicacion	provincia	año	dias	avg(tmax)	max(tmax)	min(tmax)
Loja	Granada	2021	4	44.28	45.6	42.3
Navalmoral de la ...	Caceres	2015	4	40.7	41.2	40.3
Navalmoral de la ...	Caceres	2015	4	40.95	42.5	40.3
Baza	Granada	2007	4	40.85	41.3	40.1
Antequera	Malaga	2017	5	42.68	44.4	40.1
Merida	Badajoz	2017	5	42.46	44.0	40.5
La Roda de Andalucia	Sevilla	2015	5	41.14	43.2	40.0
La Roda de Andalucia	Sevilla	2015	4	40.6	41.1	40.2
ecija	Sevilla	2007	6	42.08	43.2	41.0
Madrid Aeropuerto	Madrid	2022	4	41.0	42.2	40.5

Donde la columna *dias* representaría la duración de esa ola de calor, *avg(tmax)* la temperatura media, *max(tmax)* y *min(tmax)* las temperaturas máximas y mínimas respectivamente.

En la segunda parte de la consulta, preparamos los datos para su representación en el mapa provincial de España, además de guardarlos en un formato Parquet para su posterior lectura en Python, ya que actualmente en Scala no es posible realizar una representación sobre un mapa.

```
val resultsSave = results  
  .groupBy($"provincia", $"año")  
  .agg(func.count($"provincia"), func.avg($"dias"), func.avg($"avg(tmax)"), func.max($"max(tmax)"), func.min($"min(tmax)"))
```

Como primer paso, nos creamos una variable que hemos llamado *resultsSave*, en ella se encontrarán los datos agrupados por provincias preparados para su posterior representación. Por lo tanto, el siguiente paso a realizar sería agrupar los datos por provincia y año mediante la función *groupBy*. Mediante la función *agg* ejecutaremos diferentes funciones, por medio de la función *count* obtendremos el número de olas de calor en un mismo año en una provincia, ya que gracias a contar número de veces que aparece una provincia en los anteriores resultados podremos determinar este valor y mediante las funciones *avg*, *max* y *min* obtendremos de manera respectiva la temperatura media, máxima y mínima para esa provincia, con respecto a las olas de calor pertenecientes a un año.

```
.select($"provincia", $"año", $"count(provincia)".alias("nº de olas de calor"),
      $"avg(días)".alias("duracion media"), $"avg(avg(tmax))".alias("avg(tmax)"),
      $"max(max(tmax))".alias("max(tmax)"), $"min(min(tmax))".alias("min(tmax)"))
```

Por último, seleccionamos las columnas deseadas y cambiamos el nombre de algunas de ellas mediante la función *alias*.

```
resultsSave
  .withColumnRenamed("nº de olas de calor", "nOlasCalor")
  .withColumnRenamed("duracion media", "duracionMedia")
  .withColumnRenamed("avg(tmax)", "avgTmax")
  .withColumnRenamed("max(tmax)", "maxTmax")
  .withColumnRenamed("min(tmax)", "minTmax")
  .write.format("parquet").partitionBy("provincia").mode("overwrite").save("D:/TFGAlvaroSanchez/data/resultadoOlasCalor/")
```

Para guardar los datos en formato Parquet el nombre de las columnas no puede contener ningún espacio en blanco, por lo tanto, deberemos de darles un nuevo nombre a estas columnas. Lo realizamos a través de la función *withColumnRenamed* donde primero le indicamos el nombre actual de la columna, y a continuación el nuevo nombre que le queremos establecer.

### 3. Visualización de queries

#### Año de la temperatura máxima promedio en cada mes

```
import org.apache.spark.sql.expressions.Window

val data = spark.read.parquet("D:/TFGAlvaroSanchez/data/monthParquet/*").na.drop()

val window = Window.partitionBy("mes").orderBy($"ta_max".desc)

val dataWindow = data
  .withColumn("ta_max", func.split($"ta_max", "\\(")(0).cast(DoubleType))
  .groupBy($"fecha")
  .agg(func.avg($"ta_max").alias("ta_max"))
  .select(func.month($"fecha").alias("mes"), func.year($"fecha").alias("año"), $"ta_max")
  .withColumn("dense_rank", func.dense_rank().over(window))
  .filter($"dense_rank" === 1)
```

Como se puede observar, tal y como se describió anteriormente, el primer paso que realizamos es la lectura de los datos, los cuales los estamos leyendo en un formato Parquet para que resulte más óptimo y además estamos eliminando las filas que contienen valores nulos mediante la función `na.drop()`, la lectura de estos datos la almacenamos en la variable `data`. Al realizar la lectura de estos, se encontrarían de la siguiente manera.

fecha	ta_max
2020-01-01	19.6(31)
2020-02-01	24.3(03)
2020-03-01	21.0(19)
2020-04-01	23.7(09)

```
val window = Window.partitionBy("mes").orderBy($"ta_max".desc)
```

Por otro lado, también nos creamos la función ventana en la variable `window` que utilizaremos posteriormente. Le indicaremos la manera de la cual deseamos que se dividan los datos, en nuestro caso será a través de los meses y a su vez le señalamos la manera en la que se deberán de ordenar, teniendo en cuenta la temperatura máxima.

```
val dataWindow = data
  .withColumn("ta_max", func.split($"ta_max", "\\(")(0).cast(DoubleType))
```

A continuación, procedemos a realizar la consulta que almacenaremos en la variable `dataWindow`. El primer paso, que realizamos con los datos es la modificación de la columna de

temperatura máxima, ya que esta nos viene con el siguiente formato temperaturaMáximaAlcanzada(díaDelMesEnElQueSeAlcanzó) por lo que Spark estaría tratando esta columna como un String, cuando nuestro objetivo es manejarlo como un tipo Double. Para esto realizamos lo siguiente, separamos mediante la función *Split* la temperatura máxima del día en que se alcanzó, y nos quedamos únicamente con la temperatura máxima, por último, realizamos un casteo al tipo de datos que deseamos.

```
.groupBy($"fecha")
.agg(func.avg($"ta_max").alias("ta_max"))
```

El segundo paso, de la consulta sería agrupar los datos de las diferentes estaciones por fecha, para ello utilizamos la función *groupBy* indicándole la columna a través de la cual deseamos realizar la agrupación. Mediante la función *agg* y *avg*, obtenemos y calculamos la media de la temperatura máxima de las filas que hemos agrupado. Actualmente tendríamos el DataFrame que se muestra a continuación, que contendría únicamente las fechas y sus respectivas temperaturas.

fecha	ta_max
2013-01-01	20.84285714285715
2010-06-01	32.37096774193548
2021-03-01	25.4304347826087
2016-10-01	28.869696969696975
2019-01-01	18.59714285714286

```
.select(func.month($"fecha").alias("mes"), func.year($"fecha").alias("año"), $"ta_max")
```

En el siguiente paso, obtendremos tanto el mes y el año, a partir de la fecha. Para esto utilizamos las funciones *month* y *year* con las cuales obtenemos tanto el mes y el año respectivamente. Mediante la función *select* seleccionaremos únicamente las columnas que vamos a necesitar.

```
.withColumn("dense_rank", func.dense_rank().over(window))
```

Para ir finalizando, añadimos una columna con el valor que nos devuelve la función *dense\_rank*, que hace uso de la ventana *window* anteriormente definida. Mediante esta función, obtenemos la posición en la que se encuentra un valor teniendo en cuenta la expresión *orderBy* establecida en la función ventana (*window*). Por lo tanto, en nuestro nos devolverá una clasificación de los años más calurosos dividido por meses, ya que era la partición que se estableció en la función *window*. Quedando la consulta de la siguiente manera después de aplicar la función ventana.

mes	año	ta_max	dense_rank
7	2022	39.46785714285714	1
7	2015	37.08	2
7	2016	37.00499999999999	3
7	2020	36.78799999999999	4
7	2017	36.67575757575758	5
7	2013	35.952777777777776	6
7	2021	35.90769230769231	7
7	2019	35.85454545454545	8
7	2012	35.5258064516129	9
7	2014	35.071875000000006	10
7	2010	34.84375	11
7	2018	34.155882352941184	12
7	2011	33.51379310344828	13
11	2015	25.378787878787882	1
11	2020	25.307407407407414	2

```
.filter($"dense_rank" === 1)
```

Por último, realizamos un filtrado quedándonos únicamente con las filas que poseen en la columna *dense\_rank* un valor igual a 1, obteniendo así el año en el que se obtuvo la mayor temperatura máxima en un determinado mes.

```
def countByYear(year : Int) : Long = {
  dataWindow
    .filter($"año" === year)
    .count()
}

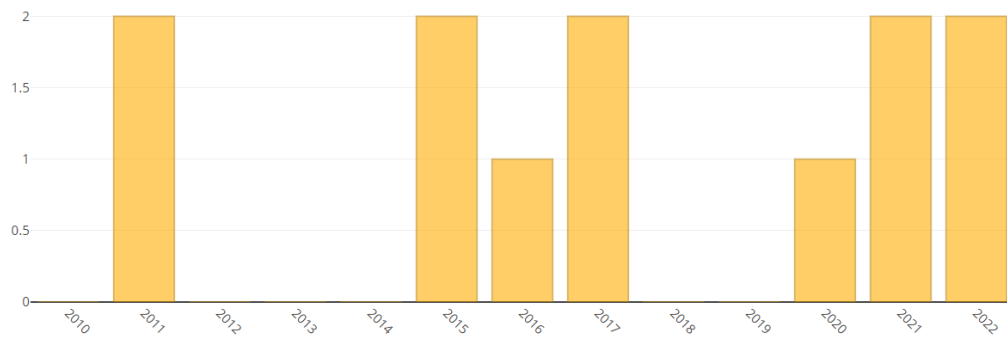
val year2010 = countByYear(2010)
```

Se implementa también la función *countByYear*, la cual nos facilitará los datos que utilizaremos para la representación gráfica, ya que nos devuelve la cantidad de meses por cada año que se encuentran en la consulta realizada anteriormente.

Para finalizar ofrecemos los resultados tanto en forma de tabla como de una forma gráfica que sería para la cual habríamos hecho uso de la función de la *countByYear*.

mes	año	temperatura maxima
1	2021	21.84
2	2020	24.42
3	2017	27.82
4	2011	29.48
5	2015	33.73
6	2017	36.86
7	2022	39.47
8	2022	37.74
9	2016	36.39
10	2011	31.52
11	2015	25.38
12	2021	22.29

Número de meses con temperaturas máximas en cada año



## Olas de calor

```
results
.withColumnRenamed("dias", "duracion (dias)")
.withColumnRenamed("avg(tmax)", "temperatura media")
.withColumnRenamed("max(tmax)", "temperatura maxima")
.withColumnRenamed("min(tmax)", "temperatura minima")
.orderBy($"ubicacion", $"año")
.show()
```

Podremos observar los resultados de diferentes maneras, una de ellas sería mostrar las diferentes olas de calor a través de una tabla, donde se nos muestre toda la información de manera detalla, ordenada por ubicación y año:

ubicacion	provincia	año	duracion (dias)	temperatura media	temperatura maxima	temperatura minima
Alajar	Huelva	2018	4	41.65	42.5	40.7
Alcaniz	Teruel	2019	4	40.95	41.2	40.8
Andujar	Jaen	2008	5	41.32	42.7	40.6
Andujar	Jaen	2009	4	41.5	42.7	40.4
Andujar	Jaen	2009	6	41.82	42.6	40.5
Andujar	Jaen	2010	4	42.35	43.1	40.9
Andujar	Jaen	2010	4	42.35	43.8	41.1
Andujar	Jaen	2010	7	41.39	42.7	40.4
Andujar	Jaen	2012	4	41.78	42.3	40.6
Andujar	Jaen	2015	4	42.38	43.6	40.8
Andujar	Jaen	2015	6	41.5	42.9	40.0
Andujar	Jaen	2015	6	41.92	44.2	40.0
Andujar	Jaen	2016	4	40.88	41.2	40.6

```
results
  .groupBy($"año")
  .agg(func.count($"año").alias("nº olas de calor"))
  .agg(func.sum("nº olas de calor"))
  .orderBy($"año")
  .show()
```

Otra manera sería la siguiente, con la que podremos observar el numero de olas de calor en los diferentes años. Como se puede ver en los ultimos años estas han aumentado de una manera considerable.

año	nº olas de calor
2007	6
2008	1
2009	9
2010	6
2012	6
2013	4
2014	1
2015	28
2016	11
2017	36
2018	23
2019	10
2020	11
2021	24
2022	68

Y por ultimo, una representación mediante un mapa, para la cual fue necesario guardar los datos de los resultados, ya que en Scala aun no existen librerías que nos permitan este tipo de representaciones, por lo tanto haremos uso de Python.

Para esta representacion, haremos uso de las siguientes librerías en Python:

```
import json
import pandas as pd
import plotly.express as px
import ipywidgets as widgets
```

Utilizaremos la librería *json* para leer el fichero geoJSON, el cual es un fichero que nos ofrece datos geograficos con el cual representaremos las diferentes provincias españolas. La librería *pandas* la usaremos para la lectura y manipulacion de los datos, mediante *plotly* podremos representar nuestro gráfico y por ultimo utilizaremos *ipywidgets* para poder utilizar elementos interactivos de HTML en Jupyter notebook.

#### **4. Despliegue en AWS EMR**

## **2. Experimentos / validación**

### **2.1. Consultas realizadas**

### **2.2. Análisis de requisitos no funcionales**

## **3. Conclusiones**

## **4. Bibliografía**

## **5. Apéndices**