

Universidad
Rey Juan Carlos

Escuela Técnica Superior
de Ingeniería Informática

Grado en Ingeniería de Computadores

Curso 2021-2022

Trabajo de Fin de Grado

Comparativa entre las API de Spark de Scala y Python

Autor: Oscar Nydza Nicpoñ

Tutor: Juan Manuel Serrano Hidalgo

Indice

- Introducción
- Objetivos
- Descripción técnica
 - Programación de queries
 - Visualización de resultados
 - Migración de código
- Comparativa
- Conclusión

Introducción

¿Qué es Spark?

De Wikipedia:

“Apache Spark se puede considerar un sistema de computación en clúster de propósito general y orientado a la velocidad. Proporciona APIs en Java, Scala, Python y R.”



Introducción

¿Qué API usamos?

Dos categorías:

- Nativas: Scala y Java
- No nativas: Python y R

Compararemos las API
de Scala y Python

Objetivos

Principal: Comparar el rendimiento entre APIs

Derivados:

- Programar queries en ambas APIs
- Visualizar los resultados obtenidos de las consultas
- Migrar queries de una API a otra

Objetivos

Orden a seguir:

1. Programar queries en ambas APIs
2. Visualizar los resultados obtenidos de las consultas
3. Migrar queries de una API a otra
4. Comparar el rendimiento entre APIs

Descripción técnica: Conjunto de datos

Se va a utilizar un conjunto de datos de la Fórmula 1

Las cuestiones sobre las que buscamos respuestas son:

- Piloto más consistente en la temporada 2012
- Mejor piloto de la historia
- Fabricante más dominante en un periodo de tiempo
- Análisis de temporada por piloto
- Temporada más interesante para el espectador



Descripción técnica: Programación

Ejemplo de código en Scala:

```
1  val lastRace = Window.partitionBy("year")
2
3  val lastRaces = spark.read.format("csv")
4    .option("header", "true")
5    .option("sep", ",")
6    .load("../data/races.csv")
7    .where(col("year") >= 1990 && col("year") <= 1999)
8    .withColumn("round", col("round").cast(IntegerType))
9    .withColumn("max", max(col("round")).over(lastRace))
10   .where(col("round") === col("max"))
11   .select("raceId", "year")
```


Descripción técnica: Programación

Ejemplo de código en Scala:

```
1 val constructors = Window.partitionBy("constructorId")
2
3 val constructorWinners = spark.read.format("csv")
4   .option("header", "true")
5   .option("sep", ",")
6   .load("../data/constructor_standings.csv")
7   .join(lastRaces, Seq("raceId"), "right")
8   .where(col("position") === 1)
9   .select("constructorId", "wins", "year")
10  .withColumn("totalChampWins", count(col("constructorId")).over(constructors))
11  .withColumn("totalRaceWins", sum(col("wins")).over(constructors).cast(IntegerType))
12  .drop("wins")
```

Descripción técnica: Programación

Ejemplo de código en Python:

```
1 lastRace = Window.partitionBy("year")
2
3 lastRaces = spark.read.format("csv")\
4     .option("header", "true")\
5     .option("sep", ",")\
6     .load("../data/races.csv")\
7     .where((F.col("year") >= 1990) & (F.col("year") <= 1999))\
8     .withColumn("round", F.col("round").cast(T.IntegerType()))\
9     .withColumn("max", F.max(F.col("round")).over(lastRace))\
10    .where(F.col("round") == F.col("max"))\
11    .select("raceId", "year")
```

Descripción técnica: Visualización

Utilizaremos Plotly como herramienta de visualización

- Uso sencillo
- Disponible para Python y Scala
- Podemos hacerlo funcionar con datos obtenidos de Spark



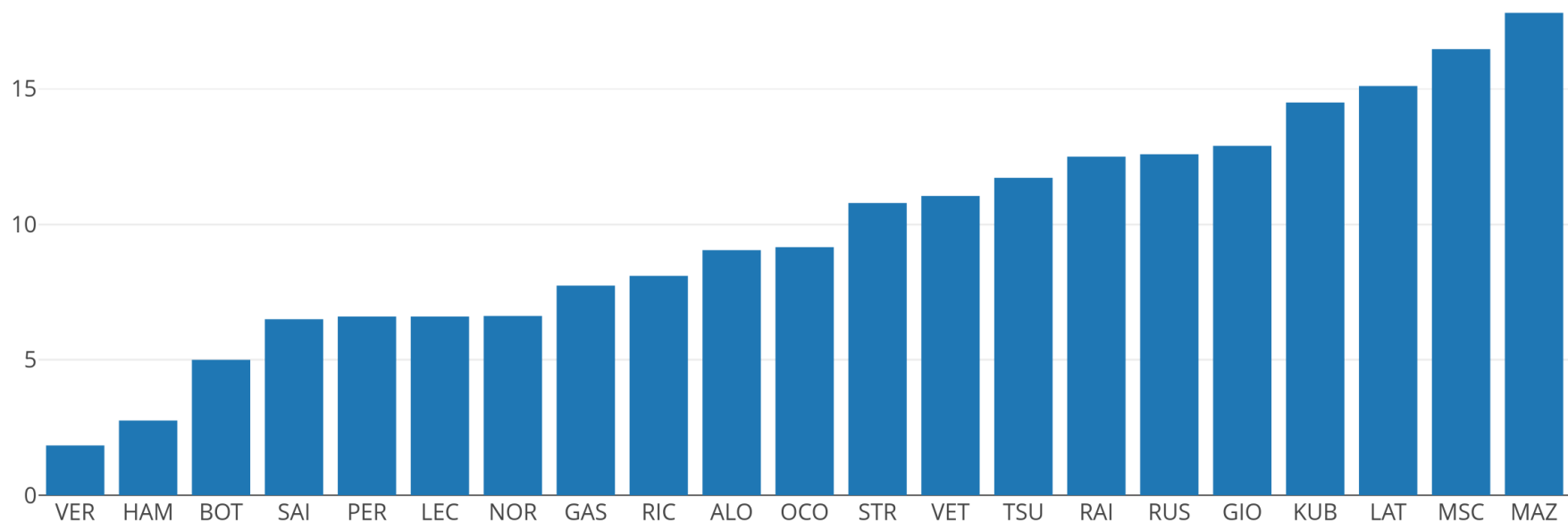
Descripción técnica: Visualización

Procedimiento:

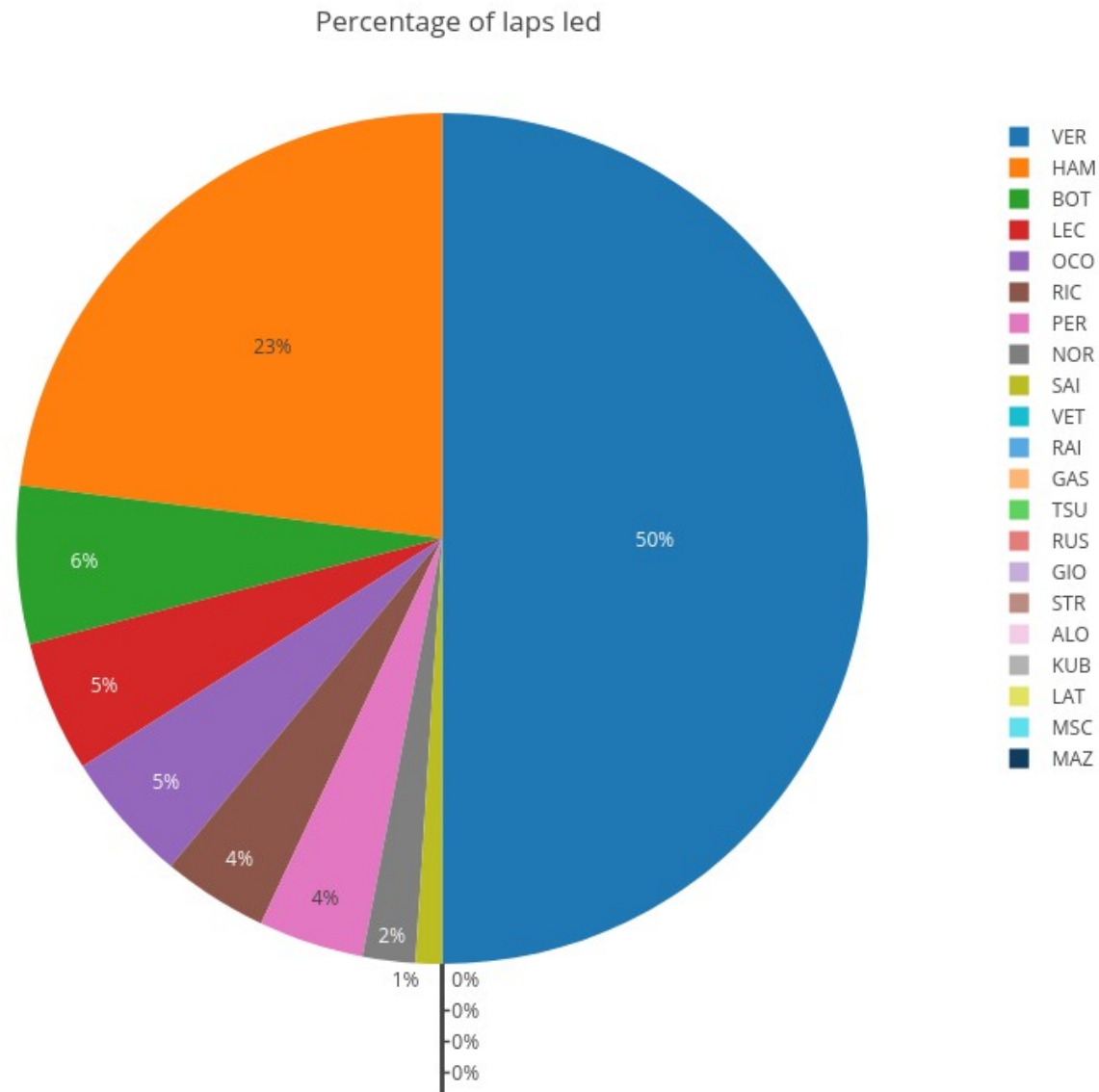
- Ordenar el DataFrame
- Conseguir ambos ejes como lista
- Crear objetos Bar y Layout
- Imprimir por pantalla

```
1 val driverCodes = results
2   .sort("avgPosition")
3   .select("code")
4   .as[String]
5   .collect()
6   .toList
7
8 val avgPosition = results
9   .sort("avgPosition")
10  .select("avgPosition")
11  .as[Double]
12  .collect()
13  .toList
14
15 val data = Seq(
16   Bar(
17     driverCodes,
18     positionDelta
19   )
20 )
21
22 val layout = Layout(
23   barmode = BarMode.Group
24 )
25
26 plot(data, layout)
```

Descripción técnica: Visualización



Descripción técnica: Visualización



Descripción técnica: Migración

Objetivo principal: ver diferencias a nivel de programación

```
1 val lastRace = Window.partitionBy("year")
2
3 val lastRaces = spark.read.format("csv")
4   .option("header", "true")
5   .option("sep", ",")
6   .load("../data/races.csv")
7   .where(col("year") >= 1990 && col("year") <= 1999)
8   .withColumn("round", col("round").cast(IntegerType))
9   .withColumn("max", max(col("round")).over(lastRace))
10  .where(col("round") === col("max"))
11  .select("raceId", "year")
```

← Scala

Python →

```
1 lastRace = Window.partitionBy("year")
2
3 lastRaces = spark.read.format("csv")\
4   .option("header", "true")\
5   .option("sep", ",")\
6   .load("../data/races.csv")\
7   .where((F.col("year") >= 1990 & (F.col("year") <= 1999))\
8   .withColumn("round", F.col("round").cast(T.IntegerType()))\
9   .withColumn("max", F.max(F.col("round")).over(lastRace))\
10  .where(F.col("round") == F.col("max"))\
11  .select("raceId", "year")
```

Descripción técnica: Migración

Automatización del proceso: primer paso

Implementación de expresiones regulares para sustitución

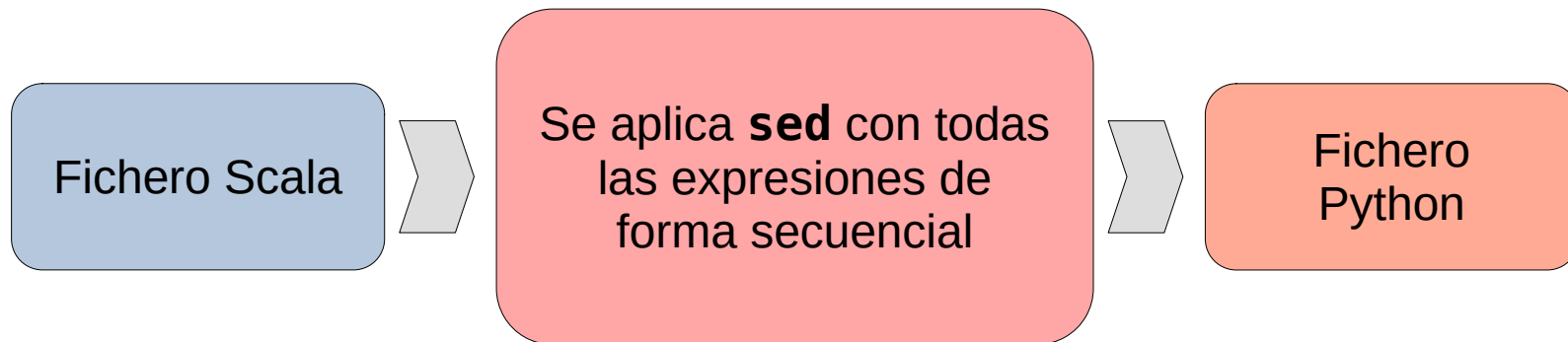
Ejemplos:

- Eliminación de var y val:
`(val|var)[]+?(?=[a-zA-Z]) → <cadena vacía>`
- Añadir carácter \ a los finales de línea:
`(\))(\n) → $1\\$2`
- Intercambiar parámetros de `na.replace()`:
`na.replace\[([|\n]*(.*)[|\n]*,[|\n]*(.*)[|\n]*\)
na.replace($2, $1)`

Descripción técnica: Migración

Automatización del proceso: segundo paso

Implementación de un script de Bash que aplique las expresiones regulares automáticamente.



Otra opción:

Exportar el plan de Spark de una API y usarla en otra
Requiere investigación

Pendiente como trabajo futuro

Comparativa

Recolección de métricas

Necesitamos obtener métricas sobre:

- Gestión del proceso
- Uso de memoria
- Tiempos de ejecución

La Spark UI nos dará lo que necesitamos

Usaremos:

- Trabajos completados
- Tareas completadas
- Tareas saltadas
- Etapas completadas
- Tiempo de tarea
- Total leído
- Total leído en operaciones de shuffle
- Total escrito en operaciones de shuffle