

## UD 6 - Apache Spark - Spark API

Ya hemos visto que uno de los mayores atractivos de Apache Spark para los desarrolladores ha sido sus APIs fáciles de usar para operar en grandes conjuntos de datos, en diferentes lenguajes de programación tales como: Scala, Java, Python y R.

### 1. Spark APIs

Vamos a explicar los tres conjuntos de APIs disponibles en **Apache Spark** (**RDDs**, **DataFrames**, and **Datasets**) además de cuándo y por qué debemos usar cada conjunto de API; descubriremos rendimiento y beneficios de optimización y enumeraremos escenarios sobre cuándo usar DataFrames y Datasets en lugar de RDD. Principalmente, nos centraremos en DataFrames y Datasets, porque desde Apache Spark 2.0, estas dos API están unificadas.

#### 1.1 RDD (Resilient Distributed Dataset)

**RDD** fue la principal API orientada al usuario en Spark desde su inicio. En el núcleo, un RDD es una colección distribuida **inmutable** de elementos de sus datos, **particionados a través de nodos en su clúster** que se pueden operar en **paralelo** con una API de bajo nivel que ofrece *transformaciones* y acciones.

##### ¿Cuándo usar RDD?

- Necesitamos una transformación de **bajo nivel** y acciones y control en su conjunto de datos.
- Para datos **no estructurados**, como flujos de medios o flujos de texto.
- desea manipular sus datos con construcciones de programación funcional que expresiones específicas de dominio.
- No necesitamos imponer un esquema, como formato de columna, mientras procesamos o accedemos a los atributos de los datos por nombre o columna.
- Podemos renunciar a algunos beneficios de optimización y rendimiento disponibles con DataFrames y Datasets para datos estructurados y semiestructurados.

#### 1.2 Dataset

Un **Dataset** es una **colección distribuida de datos**. Dataset es una nueva interfaz agregada en Spark 1.6 que proporciona los beneficios de los RDD (fuertemente tipado, capacidad de utilizar potentes funciones lambda) con los beneficios del motor de ejecución optimizado de Spark

SQL. Se puede construir un Dataset a partir de objetos JVM y luego manipularlo mediante transformaciones funcionales (map, flatMap, filter, etc.).

La API del Dataset está disponible en Scala y Java. Python no es compatible con la API del Dataset. Pero debido a la naturaleza dinámica de Python, muchos de los beneficios de la API del Dataset ya están disponibles (es decir, puede acceder al campo de una fila por su nombre `row.columnName`). El caso de R es similar.

## 1.3 Dataframe

Un **DataFrame** es un **Dataset organizado en columnas con nombre**. Es conceptualmente equivalente a una tabla en una base de datos relacional o un marco de datos en R/Python, pero con optimizaciones más ricas bajo el capó. Los DataFrames se pueden construir a partir de una amplia gama de fuentes, como: archivos de datos estructurados, tablas en Hive, bases de datos externas o RDDs existentes.

La API DataFrame está disponible en Scala, Java, Python y R. En Scala y Java, un DataFrame está representado por un Dataset de filas. En la API de Scala, DataFrame es simplemente un alias de tipo de Dataset[Row]. Mientras que, en la API de Java, los usuarios deben usar Dataset para representar un DataFrame.

## 1.4 Spark SQL

**Spark SQL** es un módulo Spark para el procesamiento de datos estructurados. A diferencia de la API básica de Spark RDD, las interfaces proporcionadas por Spark SQL brindan a Spark más información sobre la estructura de los datos y el cálculo que se realiza. Internamente, Spark SQL utiliza esta información adicional para realizar optimizaciones adicionales. Hay varias formas de interactuar con Spark SQL, incluidas SQL y la API del conjunto de datos. Al calcular un resultado, se utiliza el mismo motor de ejecución, independientemente de qué API/lenguaje esté utilizando para expresar el cálculo. Esta unificación significa que los desarrolladores pueden alternar fácilmente entre diferentes API según cuál proporcione la forma más natural de expresar una transformación determinada.

## 1.5 Beneficios Dataframe/Dataset APIs

1. **Tipado estático y tipado seguro en tiempo de ejecución:** Dadas las diferencias de tipado entre las 3 APIs de Spark nos encontramos con las siguientes diferencias

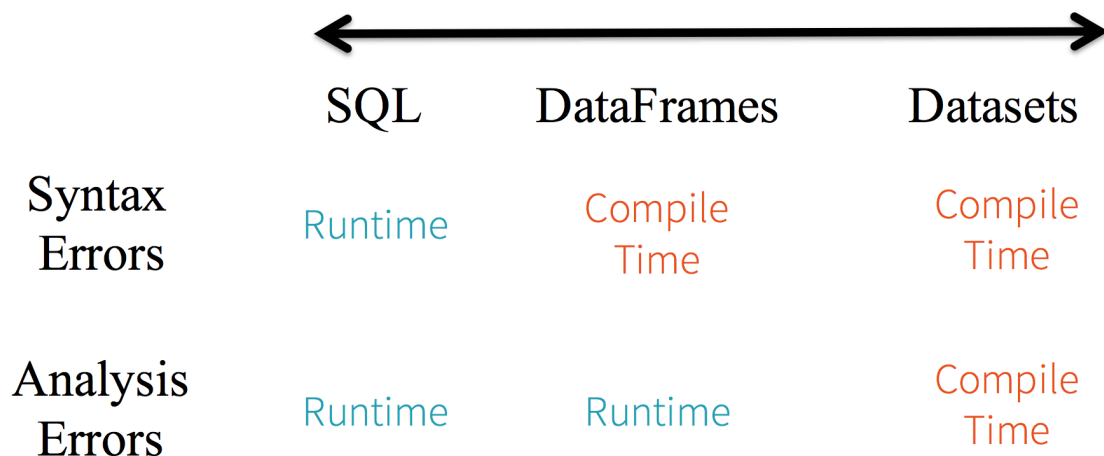


Figura 6.1\_Spark API: Mejoras en tiempo de ejecución. (Fuente: databricks.com)

2. **Alto nivel de Abstracción y vista personalizada de datos estructurados y semiestructurados.**
3. **Facilidad de uso de API con estructura:** Introduce una semántica rica y un conjunto sencillo de operaciones específicas de dominio que pueden expresarse como construcciones de alto nivel. Es mucho más sencillo realizar operaciones `agg`, `select`, `sum`, `avg`, `map`, `filter`, `groupBy`, `filter`, `map`
4. **Rendimiento y optimización:** Mejoras de eficiencia en el uso de almacenamiento y ganancias de rendimiento.

## Space Efficiency

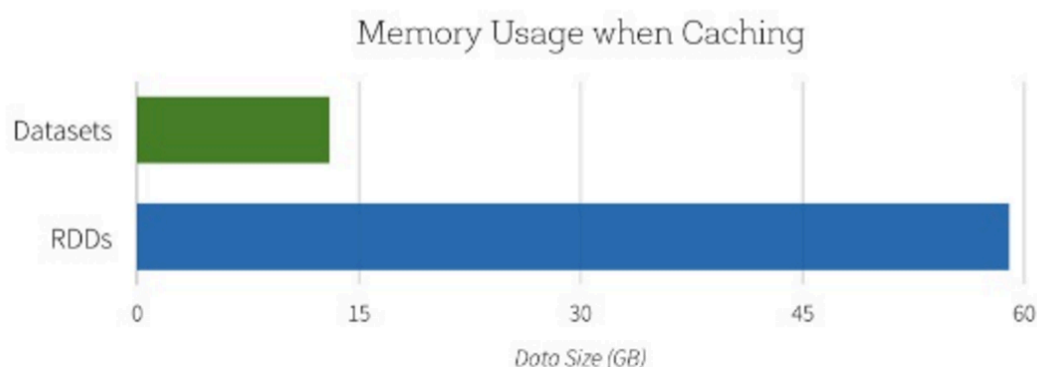


Figura 6.2\_Spark API: Mejoras de eficiencia y rendimiento. (Fuente: databricks.com)

## 2. RDDs. Getting Started

Hay dos formas de crear RDD: paralelizar una colección existente en su driver program o hacer referencia a un conjunto de datos en un sistema de almacenamiento externo, como un sistema de archivos compartido, HDFS, HBase o cualquier fuente de datos que ofrezca un formato de entrada Hadoop.

### 2.1. Parallelized Collections

Parallelized Collections se crean llamando al método de `parallelize` de `SparkContext`. Los elementos de la colección se copian para formar un conjunto de datos distribuido que se puede operar en paralelo en nuestro cluster. Por ejemplo, Creamos una colección paralela que contenga los números del 1 al 8:

#### Python

```
1 data = [1, 2, 3, 4, 5, 6, 7, 8]
2 distData = sc.parallelize(data)
```

#### Scala

```
1 val data = Array(1, 2, 3, 4, 5)
2 val distData = sc.parallelize(data)
```

#### Java

```
1 //JavaSparkContext en lugar de SparkContext
2 List<Integer> data = Arrays.asList(1, 2, 3, 4, 5);
3 JavaRDD<Integer> distData = sc.parallelize(data);
```

Una vez creado, el conjunto de datos distribuido (`distData`) se puede operar en paralelo. Por ejemplo, podemos llamar a `distData.reduce(lambda a, b: a + b)` para sumar los elementos de la lista.

### 2.2 Conjunto de datos externo

Spark puede crear conjuntos de datos distribuidos desde cualquier fuente de almacenamiento compatible con Hadoop, incluido su sistema de archivos local, HDFS, Cassandra, HBase, Amazon S3, etc. Spark admite archivos de texto, SequenceFiles y cualquier otro formato de entrada de Hadoop.

Los RDD de archivos de texto se pueden crear utilizando el método `textFile` de `SparkContext`. Este método toma un URI para el archivo (ya sea una ruta local en la máquina

o un URI `hdfs://`, `s3a://`, etc.) y lo lee como una colección de líneas. Aquí hay un ejemplo de invocación:

#### Python

```
1 >>> distFile = sc.textFile("data.txt")
```

#### Scala

```
1 scala> val distFile = sc.textFile("data.txt")
2 distFile: org.apache.spark.rdd.RDD[String] = data.txt MapPartitionsRDD[10]
  at textFile at <console>:26
```

#### Java

```
1 JavaRDD<String> distFile = sc.textFile("data.txt");
```

## 3. Spark SQL, DataFrames and Datasets. Getting started

Puedes consultar toda la información detallada en la [documentación oficial](#)

### 3.1 Starting Point: SparkSession

#### Python

El punto de entrada a todas las funciones de Spark es la clase `SparkSession`. Para crear una `SparkSession` básica, simplemente use `SparkSession.builder`:

```
1 from pyspark.sql import SparkSession
2
3 spark = SparkSession \
4     .builder \
5     .appName("Python Spark SQL basic example") \
6     .config("spark.some.config.option", "some-value") \
7     .getOrCreate()
```

Encuentra los códigos completos en `examples/src/main/python/sql/basic.py` del repositorio `Spark`

#### Scala

El punto de entrada a todas las funciones de Spark es la clase `SparkSession`. Para crear una `SparkSession` básica, simplemente use `SparkSession.builder()`:

```
1 import org.apache.spark.sql.SparkSession
```

```

2
3  val spark = SparkSession
4    .builder()
5    .appName("Spark SQL basic example")
6    .config("spark.some.config.option", "some-value")
7    .getOrCreate()

```

Encuentra los códigos completos en

`examples/src/main/scala/org/apache/spark/examples/sql/SparkSQLExample.scala` del repositorio Spark

## Java

El punto de entrada a todas las funciones de Spark es la clase `SparkSession`. Para crear una `SparkSession` básica, simplemente use `SparkSession.builder`:

```

1  import org.apache.spark.sql.SparkSession;
2
3  SparkSession spark = SparkSession
4    .builder()
5    .appName("Java Spark SQL basic example")
6    .config("spark.some.config.option", "some-value")
7    .getOrCreate();

```

Encuentra los códigos completos en

`examples/src/main/java/org/apache/spark/examples/sql/JavaSparkSQLExample.java` del repositorio Spark

## R

El punto de entrada a todas las funciones de Spark es la clase `SparkSession`. Para crear una `SparkSession` básica, simplemente use `sparkR.session()`:

```

1  sparkR.session(appName = "R Spark SQL basic example", sparkConfig =
  list(spark.some.config.option = "some-value"))

```

Encuentra los códigos completos en `examples/src/main/r/RSparkSQLExample.R` del repositorio Spark

## 3.2 Creando Dataframes

Con `SparkSession`, las aplicaciones pueden crear **DataFrames** a partir de un RDD existente, de una tabla de Hive o de fuentes de datos de Spark.

Como ejemplo, lo siguiente crea un DataFrame basado en el contenido de un archivo JSON

**Python**

```

1 # spark is an existing SparkSession
2 df = spark.read.json("examples/src/main/resources/people.json")
3 # Displays the content of the DataFrame to stdout
4 df.show()
5 # +----+-----+
6 # | age|   name|
7 # +----+-----+
8 # |null|Michael|
9 # |  30|   Andy|
10 # |  19|  Justin|
11 # +----+-----+

```

**Scala**

```

1 val df = spark.read.json("examples/src/main/resources/people.json")
2
3 // Displays the content of the DataFrame to stdout
4 df.show()
5 // +----+-----+
6 // | age|   name|
7 // +----+-----+
8 // |null|Michael|
9 // |  30|   Andy|
10 // |  19|  Justin|
11 // +----+-----+

```

**Java**

```

1 import org.apache.spark.sql.Dataset;
2 import org.apache.spark.sql.Row;
3
4 Dataset<Row> df =
spark.read().json("examples/src/main/resources/people.json");
5
6 // Displays the content of the DataFrame to stdout
7 df.show();
8 // +----+-----+
9 // | age|   name|
10 // +----+-----+
11 // |null|Michael|
12 // |  30|   Andy|
13 // |  19|  Justin|
14 // +----+-----+

```

**R**

```

1 df <- read.json("examples/src/main/resources/people.json")
2
3 # Displays the content of the DataFrame
4 head(df)
5 ##   age   name
6 ## 1  NA Michael

```

```

7  ## 2  30    Andy
8  ## 3  19   Justin
9
10 # Another method to print the first few rows and optionally truncate the
    printing of long values
11 showDF(df)
12 ## +-----+-----+
13 ## | age|    name|
14 ## +-----+-----+
15 ## |null|Michael|
16 ## | 30|    Andy|
17 ## | 19|   Justin|
18 ## +-----+-----+

```

### 3.3 DataFrame Operations

Los DataFrames proporcionan un lenguaje específico de dominio para la manipulación de datos estructurados en Scala, Java, Python y R.

Estas operaciones también se denominan "transformaciones no tipadas" en contraste con las "transformaciones fuertemente tipadas" que vienen con conjuntos de datos Scala/Java.

Aquí incluimos algunos ejemplos básicos de procesamiento de datos estructurados utilizando Datasets:

#### Python

```

1  # spark, df are from the previous example
2  # Print the schema in a tree format
3  df.printSchema()
4  # root
5  # |-- age: long (nullable = true)
6  # |-- name: string (nullable = true)
7
8  # Select only the "name" column
9  df.select("name").show()
10 # +-----+
11 # |    name|
12 # +-----+
13 # |Michael|
14 # |   Andy|
15 # |  Justin|
16 # +-----+
17
18 # Select everybody, but increment the age by 1
19 df.select(df['name'], df['age'] + 1).show()
20 # +-----+-----+
21 # |    name|(age + 1)|
22 # +-----+-----+
23 # |Michael|         null|
24 # |   Andy|          31|
25 # |  Justin|          20|
26 # +-----+-----+

```



```

27
28 # Select people older than 21
29 df.filter(df['age'] > 21).show()
30 # +---+----+
31 # |age|name|
32 # +---+----+
33 # | 30|Andy|
34 # +---+----+
35
36 # Count people by age
37 df.groupBy("age").count().show()
38 # +---+----+
39 # | age|count|
40 # +---+----+
41 # | 19|    1|
42 # |null|   1|
43 # | 30|    1|
44 # +---+----+

```

## Scala

```

1 // This import is needed to use the $-notation
2 import spark.implicits._
3 // Print the schema in a tree format
4 df.printSchema()
5 // root
6 // |-- age: long (nullable = true)
7 // |-- name: string (nullable = true)
8
9 // Select only the "name" column
10 df.select("name").show()
11 // +-----+
12 // |   name|
13 // +-----+
14 // |Michael|
15 // |   Andy|
16 // |  Justin|
17 // +-----+
18
19 // Select everybody, but increment the age by 1
20 df.select($"name", $"age" + 1).show()
21 // +-----+-----+
22 // |   name|(age + 1)|
23 // +-----+-----+
24 // |Michael|        null|
25 // |   Andy|         31|
26 // |  Justin|        20|
27 // +-----+-----+
28
29 // Select people older than 21
30 df.filter($"age" > 21).show()
31 // +---+----+
32 // |age|name|
33 // +---+----+
34 // | 30|Andy|
35 // +---+----+
36

```

```

37 // Count people by age
38 df.groupBy("age").count().show()
39 // +----+-----+
40 // | age|count|
41 // +----+-----+
42 // | 19|    1|
43 // |null|    1|
44 // | 30|    1|
45 // +----+-----+

```

## Java

```

1 // col("...") is preferable to df.col("...")
2 import static org.apache.spark.sql.functions.col;
3
4 // Print the schema in a tree format
5 df.printSchema();
6 // root
7 // |-- age: long (nullable = true)
8 // |-- name: string (nullable = true)
9
10 // Select only the "name" column
11 df.select("name").show();
12 // +-----+
13 // |   name|
14 // +-----+
15 // |Michael|
16 // |   Andy|
17 // |  Justin|
18 // +-----+
19
20 // Select everybody, but increment the age by 1
21 df.select(col("name"), col("age").plus(1)).show();
22 // +-----+-----+
23 // |   name|(age + 1)|
24 // +-----+-----+
25 // |Michael|      null|
26 // |   Andy|       31|
27 // |  Justin|       20|
28 // +-----+-----+
29
30 // Select people older than 21
31 df.filter(col("age").gt(21)).show();
32 // +---+-----+
33 // |age|name|
34 // +---+-----+
35 // | 30|Andy|
36 // +---+-----+
37
38 // Count people by age
39 df.groupBy("age").count().show();
40 // +----+-----+
41 // | age|count|
42 // +----+-----+
43 // | 19|    1|
44 // |null|    1|
45 // | 30|    1|

```

```
46 // +-----+-----+
```

**R**

```
1 # Create the DataFrame
2 df <- read.json("examples/src/main/resources/people.json")
3
4 # Show the content of the DataFrame
5 head(df)
6 ##   age   name
7 ## 1  NA Michael
8 ## 2  30   Andy
9 ## 3  19   Justin
10
11
12 # Print the schema in a tree format
13 printSchema(df)
14 ## root
15 ## |-- age: long (nullable = true)
16 ## |-- name: string (nullable = true)
17
18 # Select only the "name" column
19 head(select(df, "name"))
20 ##      name
21 ## 1 Michael
22 ## 2   Andy
23 ## 3  Justin
24
25 # Select everybody, but increment the age by 1
26 head(select(df, df$name, df$age + 1))
27 ##      name (age + 1.0)
28 ## 1 Michael          NA
29 ## 2   Andy           31
30 ## 3  Justin          20
31
32 # Select people older than 21
33 head(where(df, df$age > 21))
34 ##   age name
35 ## 1  30 Andy
36
37 # Count people by age
38 head(count(groupBy(df, "age"))))
39 ##   age count
40 ## 1  19     1
41 ## 2  NA     1
42 ## 3  30     1
```

### 3.4 Ejecutando sentencias SQL

La función SQL en `SparkSession` permite que las aplicaciones ejecuten consultas SQL mediante programación y devuelve el resultado como un `DataFrame`.

**Python**

```

1  # Register the DataFrame as a SQL temporary view
2  df.createOrReplaceTempView("people")
3
4  sqlDF = spark.sql("SELECT * FROM people")
5  sqlDF.show()
6  # +----+-----+
7  # | age|   name|
8  # +----+-----+
9  # |null|Michael|
10 # |  30|   Andy|
11 # |  19|  Justin|
12 # +----+-----+

```

### Scala

```

1  // Register the DataFrame as a SQL temporary view
2  df.createOrReplaceTempView("people")
3
4  val sqlDF = spark.sql("SELECT * FROM people")
5  sqlDF.show()
6  // +----+-----+
7  // | age|   name|
8  // +----+-----+
9  // |null|Michael|
10 // |  30|   Andy|
11 // |  19|  Justin|
12 // +----+-----+

```

### Java

```

1  import org.apache.spark.sql.Dataset;
2  import org.apache.spark.sql.Row;
3
4  // Register the DataFrame as a SQL temporary view
5  df.createOrReplaceTempView("people");
6
7  Dataset<Row> sqlDF = spark.sql("SELECT * FROM people");
8  sqlDF.show();
9  // +----+-----+
10 // | age|   name|
11 // +----+-----+
12 // |null|Michael|
13 // |  30|   Andy|
14 // |  19|  Justin|
15 // +----+-----+

```

### R

```

1  df <- sql("SELECT * FROM table")

```

## 3.5 Vista temporal global

Las vistas temporales en Spark SQL tienen un ámbito de sesión y desaparecerán si finaliza la sesión que las crea. Si deseamos tener una vista temporal que se comparta entre todas las sesiones y mantenerla activa hasta que finalice la aplicación Spark, podemos crear una **vista temporal global**. La vista temporal global está vinculada a una base de datos `global_temp` preservada por el sistema, y debemos usar el nombre calificado para referirla, por ejemplo:

```
SELECT * FROM global_temp.view1.
```

### Python

```
1 # Register the DataFrame as a global temporary view
2 df.createGlobalTempView("people")
3
4 # Global temporary view is tied to a system preserved database
`global_temp`
5 spark.sql("SELECT * FROM global_temp.people").show()
6 # +----+-----+
7 # | age|   name|
8 # +----+-----+
9 # |null|Michael|
10 # | 30|   Andy|
11 # | 19|  Justin|
12 # +----+-----+
13
14 # Global temporary view is cross-session
15 spark.newSession().sql("SELECT * FROM global_temp.people").show()
16 # +----+-----+
17 # | age|   name|
18 # +----+-----+
19 # |null|Michael|
20 # | 30|   Andy|
21 # | 19|  Justin|
22 # +----+-----+
```

### Scala

```
1 // Register the DataFrame as a global temporary view
2 df.createGlobalTempView("people")
3
4 // Global temporary view is tied to a system preserved database
`global_temp`
5 spark.sql("SELECT * FROM global_temp.people").show()
6 // +----+-----+
7 // | age|   name|
8 // +----+-----+
9 // |null|Michael|
10 // | 30|   Andy|
11 // | 19|  Justin|
12 // +----+-----+
13
14 // Global temporary view is cross-session
15 spark.newSession().sql("SELECT * FROM global_temp.people").show()
16 // +----+-----+
17 // | age|   name|
18 // +----+-----+
```

```

19 // |null|Michael|
20 // | 30|   Andy|
21 // | 19|  Justin|
22 // +----+-----+

```

### Java

```

1 // Register the DataFrame as a global temporary view
2 df.createGlobalTempView("people");
3
4 // Global temporary view is tied to a system preserved database
`global_temp`
5 spark.sql("SELECT * FROM global_temp.people").show();
6 // +----+-----+
7 // | age|   name|
8 // +----+-----+
9 // |null|Michael|
10 // | 30|   Andy|
11 // | 19|  Justin|
12 // +----+-----+
13
14 // Global temporary view is cross-session
15 spark.newSession().sql("SELECT * FROM global_temp.people").show();
16 // +----+-----+
17 // | age|   name|
18 // +----+-----+
19 // |null|Michael|
20 // | 30|   Andy|
21 // | 19|  Justin|
22 // +----+-----+

```

## 3.6 Creando Dataset

Los Dataset son similares a los RDD, sin embargo, en lugar de utilizar la serialización Java o Kryo, utilizan un codificador [Encoder](#) para serializar los objetos para procesarlos o transmitirlos a través de la red. Si bien tanto los codificadores como la serialización estándar son responsables de convertir un objeto en bytes, los codificadores son código generado dinámicamente y utilizan un formato que permite a Spark realizar muchas operaciones como filtrado, clasificación y hash sin deserializar los bytes nuevamente en un objeto.

### Scala

```

1 case class Person(name: String, age: Long)
2
3 // Encoders are created for case classes
4 val caseClassDS = Seq(Person("Andy", 32)).toDS()
5 caseClassDS.show()
6 // +----+----+
7 // |name|age|
8 // +----+----+
9 // |Andy| 32|
10 // +----+----+

```

```

11
12 // Encoders for most common types are automatically provided by importing
spark.implicit._
13 val primitiveDS = Seq(1, 2, 3).toDS()
14 primitiveDS.map(_ + 1).collect() // Returns: Array(2, 3, 4)
15
16 // DataFrames can be converted to a Dataset by providing a class. Mapping
will be done by name
17 val path = "examples/src/main/resources/people.json"
18 val peopleDS = spark.read.json(path).as[Person]
19 peopleDS.show()
20 // +----+-----+
21 // | age|   name|
22 // +----+-----+
23 // |null|Michael|
24 // |  30|   Andy|
25 // |  19|  Justin|
26 // +----+-----+

```

### Java

```

1  import java.util.Arrays;
2  import java.util.Collections;
3  import java.io.Serializable;
4
5  import org.apache.spark.api.java.function.MapFunction;
6  import org.apache.spark.sql.Dataset;
7  import org.apache.spark.sql.Row;
8  import org.apache.spark.sql.Encoder;
9  import org.apache.spark.sql.Encoders;
10
11  public static class Person implements Serializable {
12      private String name;
13      private long age;
14
15      public String getName() {
16          return name;
17      }
18
19      public void setName(String name) {
20          this.name = name;
21      }
22
23      public long getAge() {
24          return age;
25      }
26
27      public void setAge(long age) {
28          this.age = age;
29      }
30  }
31
32  // Create an instance of a Bean class
33  Person person = new Person();
34  person.setName("Andy");
35  person.setAge(32);
36

```

```

37 // Encoders are created for Java beans
38 Encoder<Person> personEncoder = Encoders.bean(Person.class);
39 Dataset<Person> javaBeanDS = spark.createDataset(
40 Collections.singletonList(person),
41 personEncoder
42 );
43 javaBeanDS.show();
44 // +---+-----+
45 // |age|name|
46 // +---+-----+
47 // | 32|Andy|
48 // +---+-----+
49
50 // Encoders for most common types are provided in class Encoders
51 Encoder<Long> longEncoder = Encoders.LONG();
52 Dataset<Long> primitiveDS = spark.createDataset(Arrays.asList(1L, 2L,
53 3L), longEncoder);
54 Dataset<Long> transformedDS = primitiveDS.map(
55     (MapFunction<Long, Long>) value -> value + 1L,
56     longEncoder);
57 transformedDS.collect(); // Returns [2, 3, 4]
58
59 // DataFrames can be converted to a Dataset by providing a class. Mapping
60 based on name
61 String path = "examples/src/main/resources/people.json";
62 Dataset<Person> peopleDS = spark.read().json(path).as(personEncoder);
63 peopleDS.show();
64 // +-----+-----+
65 // | age|    name|
66 // +-----+-----+
67 // |null|Michael|
68 // | 30|    Andy|
69 // | 19|   Justin|
70 // +-----+-----+

```

### 3.7 Interoperar con RDD

Spark SQL admite dos métodos diferentes para convertir RDD existentes en Dataset. El primer método utiliza la **reflexión** para inferir el esquema de un RDD que contiene tipos específicos de objetos. Este enfoque basado en la reflexión genera un código más conciso y funciona bien cuando ya conoce el esquema mientras escribe su aplicación Spark.

El segundo método para crear Dataset es a través de una **interfaz programática** que le permite construir un esquema y luego aplicarlo a un RDD existente. Si bien este método es más detallado, le permite construir Dataset cuando las columnas y sus tipos no se conocen hasta el tiempo de ejecución.

#### 1. Usando Reflexión

##### Python

```

1 from pyspark.sql import Row

```



```

2
3  sc = spark.sparkContext
4
5  # Load a text file and convert each line to a Row.
6  lines = sc.textFile("examples/src/main/resources/people.txt")
7  parts = lines.map(lambda l: l.split(","))
8  people = parts.map(lambda p: Row(name=p[0], age=int(p[1])))
9
10 # Infer the schema, and register the DataFrame as a table.
11 schemaPeople = spark.createDataFrame(people)
12 schemaPeople.createOrReplaceTempView("people")
13
14 # SQL can be run over DataFrames that have been registered as a table.
15 teenagers = spark.sql("SELECT name FROM people WHERE age >= 13 AND age <=
16 19")
17
18 # The results of SQL queries are Dataframe objects.
19 # rdd returns the content as an :class:`pyspark.RDD` of :class:`Row`.
20 teenNames = teenagers.rdd.map(lambda p: "Name: " + p.name).collect()
21 for name in teenNames:
22     print(name)
23 # Name: Justin

```

## Scala

```

1  // For implicit conversions from RDDs to DataFrames
2  import spark.implicits._
3
4  // Create an RDD of Person objects from a text file, convert it to a
5  // DataFrame
6  val peopleDF = spark.sparkContext
7  .textFile("examples/src/main/resources/people.txt")
8  .map(_.split(","))
9  .map(attributes => Person(attributes(0), attributes(1).trim.toInt))
10 .toDF()
11
12 // Register the DataFrame as a temporary view
13 peopleDF.createOrReplaceTempView("people")
14
15 // SQL statements can be run by using the sql methods provided by Spark
16 val teenagersDF = spark.sql("SELECT name, age FROM people WHERE age
17 BETWEEN 13 AND 19")
18
19 // The columns of a row in the result can be accessed by field index
20 teenagersDF.map(teenager => "Name: " + teenager(0)).show()
21 // +-----+
22 // |      value|
23 // +-----+
24 // |Name: Justin|
25 // +-----+
26
27 // or by field name
28 teenagersDF.map(teenager => "Name: " + teenager.getAs[String]
29 ("name")).show()
30 // +-----+
31 // |      value|
32 // +-----+
33 // |Name: Justin|

```

```

30 // +-----+
31
32 // No pre-defined encoders for Dataset[Map[K,V]], define explicitly
33 implicit val mapEncoder = org.apache.spark.sql.Encoders.kryo[Map[String,
Any]]
34 // Primitive types and case classes can be also defined as
35 // implicit val stringIntMapEncoder: Encoder[Map[String, Any]] =
ExpressionEncoder()
36
37 // row.getValuesMap[T] retrieves multiple columns at once into a
Map[String, T]
38 teenagersDF.map(teenager => teenager.getValuesMap[Any])(List("name",
"age"))).collect()
39 // Array(Map("name" -> "Justin", "age" -> 19))

```

### Java

```

1  import org.apache.spark.api.java.JavaRDD;
2  import org.apache.spark.api.java.function.Function;
3  import org.apache.spark.api.java.function.MapFunction;
4  import org.apache.spark.sql.Dataset;
5  import org.apache.spark.sql.Row;
6  import org.apache.spark.sql.Encoder;
7  import org.apache.spark.sql.Encoders;
8
9  // Create an RDD of Person objects from a text file
10 JavaRDD<Person> peopleRDD = spark.read()
11   .textFile("examples/src/main/resources/people.txt")
12   .javaRDD()
13   .map(line -> {
14       String[] parts = line.split(",");
15       Person person = new Person();
16       person.setName(parts[0]);
17       person.setAge(Integer.parseInt(parts[1].trim()));
18       return person;
19   });
20
21 // Apply a schema to an RDD of JavaBeans to get a DataFrame
22 Dataset<Row> peopleDF = spark.createDataFrame(peopleRDD, Person.class);
23 // Register the DataFrame as a temporary view
24 peopleDF.createOrReplaceTempView("people");
25
26 // SQL statements can be run by using the sql methods provided by spark
27 Dataset<Row> teenagersDF = spark.sql("SELECT name FROM people WHERE age
BETWEEN 13 AND 19");
28
29 // The columns of a row in the result can be accessed by field index
30 Encoder<String> stringEncoder = Encoders.STRING();
31 Dataset<String> teenagerNamesByIndexDF = teenagersDF.map(
32     (MapFunction<Row, String>) row -> "Name: " + row.getString(0),
33     stringEncoder);
34 teenagerNamesByIndexDF.show();
35 // +-----+
36 // |      value|
37 // +-----+
38 // |Name: Justin|
39 // +-----+

```

```

40
41 // or by field name
42 Dataset<String> teenagerNamesByFieldDF = teenagersDF.map(
43     (MapFunction<Row, String>) row -> "Name: " + row.
44     <String>getAs("name"),
45     stringEncoder);
46 teenagerNamesByFieldDF.show();
47 // +-----+
48 // |      value|
49 // +-----+
50 // |Name: Justin|
51 // +-----+

```

## 2. Interfaz Programática

Se puede crear un DataFrame mediante programación con tres pasos:

- Creamos un RDD de tuplas/filas o listas a partir del RDD original;
- Creamos el esquema representado por un `StructType` que coincida con la estructura de tuplas/filas o listas en el RDD creado en el paso 1.
- Aplicamos el esquema al RDD mediante el método `createDataFrame` proporcionado por `SparkSession`

### Python

```

1  # Import data types
2  from pyspark.sql.types import StringType, StructType, StructField
3
4  sc = spark.sparkContext
5
6  # Load a text file and convert each line to a Row.
7  lines = sc.textFile("examples/src/main/resources/people.txt")
8  parts = lines.map(lambda l: l.split(","))
9  # Each line is converted to a tuple.
10 people = parts.map(lambda p: (p[0], p[1].strip()))
11
12 # The schema is encoded in a string.
13 schemaString = "name age"
14
15 fields = [StructField(field_name, StringType(), True) for field_name in
16 schemaString.split()]
17 schema = StructType(fields)
18
19 # Apply the schema to the RDD.
20 schemaPeople = spark.createDataFrame(people, schema)
21
22 # Creates a temporary view using the DataFrame
23 schemaPeople.createOrReplaceTempView("people")
24
25 # SQL can be run over DataFrames that have been registered as a table.
26 results = spark.sql("SELECT name FROM people")
27
28 results.show()

```

```

28 # +-----+
29 # |   name|
30 # +-----+
31 # |Michael|
32 # |   Andy|
33 # |  Justin|
34 # +-----+

```

## Scala

```

1  import org.apache.spark.sql.Row
2
3  import org.apache.spark.sql.types._
4
5  // Create an RDD
6  val peopleRDD =
spark.sparkContext.textFile("examples/src/main/resources/people.txt")
7
8  // The schema is encoded in a string
9  val schemaString = "name age"
10
11 // Generate the schema based on the string of schema
12 val fields = schemaString.split(" ")
13 .map(fieldName => StructField(fieldName, StringType, nullable = true))
14 val schema = StructType(fields)
15
16 // Convert records of the RDD (people) to Rows
17 val rowRDD = peopleRDD
18 .map(_._split(","))
19 .map(attributes => Row(attributes(0), attributes(1).trim))
20
21 // Apply the schema to the RDD
22 val peopleDF = spark.createDataFrame(rowRDD, schema)
23
24 // Creates a temporary view using the DataFrame
25 peopleDF.createOrReplaceTempView("people")
26
27 // SQL can be run over a temporary view created using DataFrames
28 val results = spark.sql("SELECT name FROM people")
29
30 // The results of SQL queries are DataFrames and support all the normal
RDD operations
31 // The columns of a row in the result can be accessed by field index or
by field name
32 results.map(attributes => "Name: " + attributes(0)).show()
33 // +-----+
34 // |          value|
35 // +-----+
36 // |Name: Michael|
37 // |   Name: Andy|
38 // | Name: Justin|
39 // +-----+

```

## Java

```

1  import java.util.ArrayList;

```

```
2 import java.util.List;
3
4 import org.apache.spark.api.java.JavaRDD;
5 import org.apache.spark.api.java.function.Function;
6
7 import org.apache.spark.sql.Dataset;
8 import org.apache.spark.sql.Row;
9
10 import org.apache.spark.sql.types.DataTypes;
11 import org.apache.spark.sql.types.StructField;
12 import org.apache.spark.sql.types.StructType;
13
14 // Create an RDD
15 JavaRDD<String> peopleRDD = spark.sparkContext()
16     .textFile("examples/src/main/resources/people.txt", 1)
17     .toJavaRDD();
18
19 // The schema is encoded in a string
20 String schemaString = "name age";
21
22 // Generate the schema based on the string of schema
23 List<StructField> fields = new ArrayList<>();
24 for (String fieldName : schemaString.split(" ")) {
25     StructField field = DataTypes.createStructField(fieldName,
26     DataTypes.StringType, true);
27     fields.add(field);
28 }
29 StructType schema = DataTypes.createStructType(fields);
30
31 // Convert records of the RDD (people) to Rows
32 JavaRDD<Row> rowRDD = peopleRDD.map((Function<String, Row>) record -> {
33     String[] attributes = record.split(",");
34     return RowFactory.create(attributes[0], attributes[1].trim());
35 });
36
37 // Apply the schema to the RDD
38 Dataset<Row> peopleDataFrame = spark.createDataFrame(rowRDD, schema);
39
40 // Creates a temporary view using the DataFrame
41 peopleDataFrame.createOrReplaceTempView("people");
42
43 // SQL can be run over a temporary view created using DataFrames
44 Dataset<Row> results = spark.sql("SELECT name FROM people");
45
46 // The results of SQL queries are DataFrames and support all the normal
47 // RDD operations
48 // The columns of a row in the result can be accessed by field index or
49 // by field name
50 Dataset<String> namesDS = results.map(
51     (MapFunction<Row, String>) row -> "Name: " + row.getString(0),
52     Encoders.STRING());
53 namesDS.show();
54 // +-----+
55 // |          value|
56 // +-----+
57 // |Name: Michael|
58 // |      Name: Andy|
```

```
56 // | Name: Justin|
57 // +-----+
```

Esta es sólo una pequeña muestra de la Spark API. Puedes consultar la [documentación completa en la página oficial](#)

### Ejemplo

```
1 import sys
2 from pyspark.sql import SparkSession
3
4 spark = SparkSession.builder.getOrCreate()
5
6 df = spark.read.options(header='True', inferSchema='True').csv(
sys.argv[1] )
7
8 def myFunc(s):
9     if s["brand"]=="riche" and s["event_type"]=="cart":
10         return [ ( s["product_id"], 1) ]
11     return []
12
13 lines=df.rdd.flatMap(myFunc).reduceByKey(lambda a, b: a + b)
14
15 lines.saveAsTextFile( sys.argv[2] )
```

## 4. Ejemplo

Siguiendo con los ejemplos anteriores, vamos a lanzar la aplicación `wordcount` de python que nos facilitan los propios ejemplos de Spark , para contar las palabras de `El_quijote.txt` ya usado en otras ocasiones (Así, también nos podría servir para establecer una comparativa entre la ejecución en Yarn on MR y en Spark)

### 1. Iniciamos HDFS y Spark

```
1 start-dfs.sh
2 ./sbin/start-master.sh
3 ./sbin/start-workers.sh
```

### 2. Copiamos el fichero `El_quijote.txt` en HDFS

```
1 wget
https://gist.githubusercontent.com/jaimerabasco/cb528c32b4c4092e6a0763d8b6bc25c0/
2 hdfs dfs -mkdir -p /bda/spark/ejemplos
3 hdfs dfs -copyFromLocal El_Quijote.txt /bda/spark/ejemplos
```

### 3. Ejecutamos `wordcount` en Spark

```
1 spark-submit --master spark://192.168.11.10:7077
examples/src/main/python/wordcount.py /bda/spark/ejemplos/El_Quijote.txt
```

**Spark Master at spark://192.168.11.10:7077**

URL: spark://192.168.11.10:7077  
 Alive Workers: 3  
 Cores in use: 3 Total, 3 Used  
 Memory in use: 8.5 GiB Total, 3.0 GiB Used  
 Resources in use:  
 Applications: 1 Running, 0 Completed  
 Drivers: 0 Running, 0 Completed  
 Status: ALIVE

Workers (3)

Worker ID	Address	State	Cores	Memory
worker-20250131124900-192.168.11.11-35953	192.168.11.11:35953	ALIVE	1 (1 Used)	2.8 GiB (1024.0 MiB Used)
worker-20250131124900-192.168.11.12-35757	192.168.11.12:35757	ALIVE	1 (1 Used)	2.8 GiB (1024.0 MiB Used)
worker-20250131124901-192.168.11.13-46677	192.168.11.13:46677	ALIVE	1 (1 Used)	2.8 GiB (1024.0 MiB Used)

Running Applications (1)

Application ID	Name	Cores	Memory per Executor	Resources Per Executor	Submitted Time
app-20250131125043-0000	(kill) PythonWordCount	3	1024.0 MiB		2025/01/31 12:50:43

Completed Applications (0)

Application ID	Name	Cores	Memory per Executor	Resources Per Executor	Submitted Time
----------------	------	-------	---------------------	------------------------	----------------

Figura 6.3\_Spark API: Ejecución app wordcount

## 5. Pandas on PySpark



Figura 6.1\_Spark API Pandas: Logo Pandas (Fuente: [pandas.pydata.org](https://pandas.pydata.org))

### 5.1 ¿Qué es Pandas?

Pandas es una biblioteca de software de código abierto diseñada específicamente para la manipulación y el análisis de datos en el lenguaje Python. Es potente, flexible y fácil de usar.

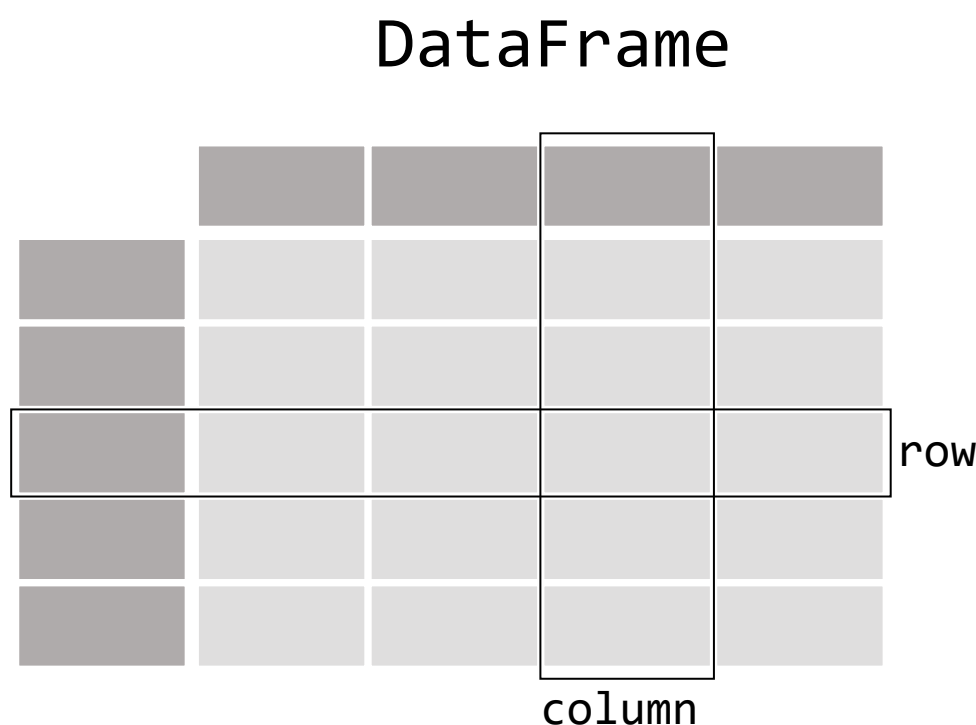
Gracias a Pandas, podemos utilizar el lenguaje Python para cargar, alinear, manipular o incluso fusionar datos. Tiene un alto rendimiento.

El nombre "Pandas" es en realidad una contracción del término «Panel Data» para series de datos que incluyen observaciones a lo largo de varios periodos de tiempo. La biblioteca se creó como herramienta de alto nivel para el análisis en Python. Los creadores de Pandas pretenden que esta biblioteca evolucione hasta convertirse en la herramienta de análisis y manipulación de datos de código abierto más potente y flexible en cualquier lenguaje de programación.

Además del análisis de datos, Pandas se utiliza mucho para la «Data Wrangling». Este término engloba los métodos de transformación de datos no estructurados para hacerlos procesables.

Puedes consultar toda la documentación en la página oficial de [Pandas](#)

## 5.2 Funcionamiento



*Figura 6.2\_Spark API Pandas: Dataframe (Fuente: [pandas.pydata.org](#))*

Pandas esta construida con el paquete **Numpy** que es una librería de python para la manipulación de matrices y vectores n-dimensional. Entonces un **Dataframe** en pandas es la estructura de datos clave que va a permitir la manipulación de datos tabulados en filas y columnas.



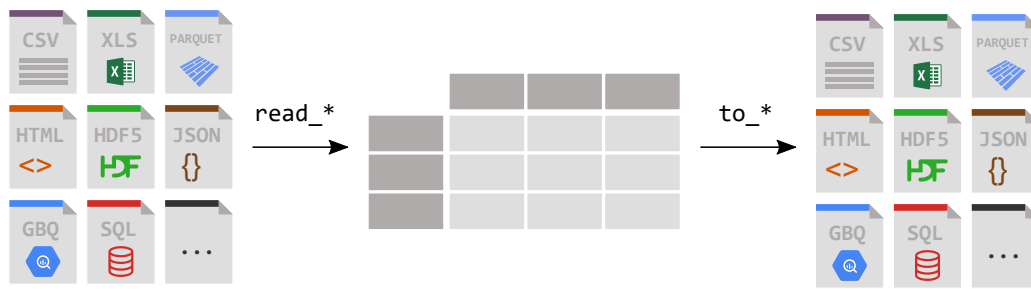


Figura 6.3\_Spark API Pandas: Formatos (Fuente: pandas.pydata.org)

Admite la integración con muchos formatos de archivos o fuentes de datos listas para usar (csv, excel, sql, json, parquet,...). La importación de datos de cada una de estas fuentes de datos se realiza mediante una función con el prefijo `read_*`. De manera similar, los métodos `to_*` se utilizan para almacenar datos.

Además, te permite realizar numerosas operaciones sobre los **DataFrame muy eficientemente**, tales como:

- **Selección de un subconjunto**

Permite seleccionar y filtrar datos según cualquier condición y extraer los datos que queramos.



Figura 6.4\_Spark API Pandas: Subconjunto (Fuente: pandas.pydata.org)

- **Crear tramas (plot)**

Pandas permite "dibujar" sus datos de forma inmediata gracias a *Matplotlib*. Podemos elegir el tipo de gráfico (dispersión, barras, diagrama de caja,...) correspondiente a los datos.

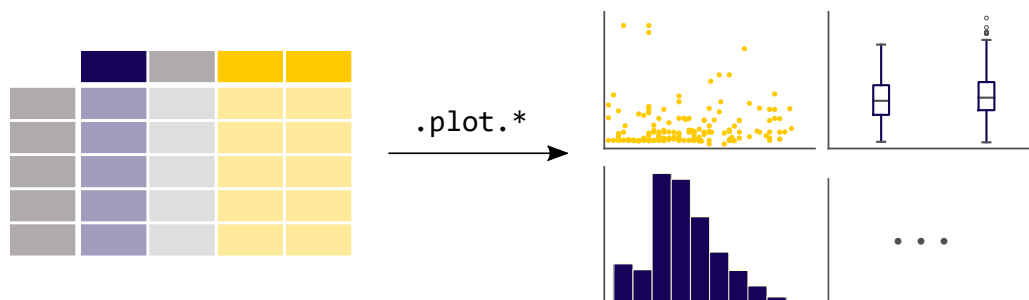


Figura 6.5\_Spark API Pandas: Tramas (Fuente: [pandas.pydata.org](https://pandas.pydata.org))

- **Crear nuevas columnas derivadas de columnas existentes**

No es necesario recorrer todas las filas de su tabla de datos para realizar cálculos. Las manipulaciones de datos en una columna funcionan por elementos. Así que podemos agregar una columna a un DataFrame en función de los datos existentes de forma sencilla y eficiente.

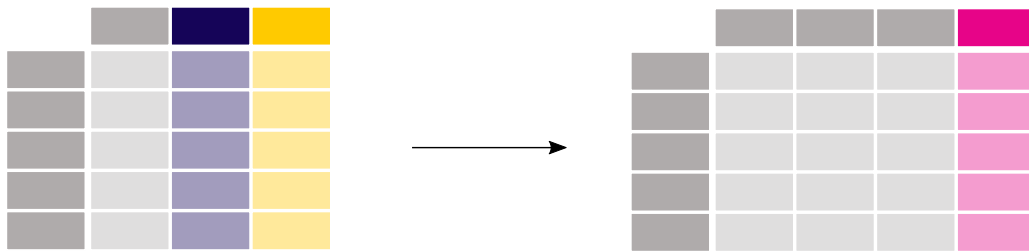


Figura 6.3\_Spark API Pandas: Columnas nuevas (Fuente: [pandas.pydata.org](https://pandas.pydata.org))

- **Calcular estadísticas**

Podemos realizar estadísticas básicas de forma sencilla (media, mediana, mínimo, máximo, recuentos...). Estas agregaciones personalizadas pueden aplicarse a todo el conjunto de datos, a un subconjunto de datos o agruparse por categorías. Este último también se conoce como enfoque de división, aplicación y combinación.

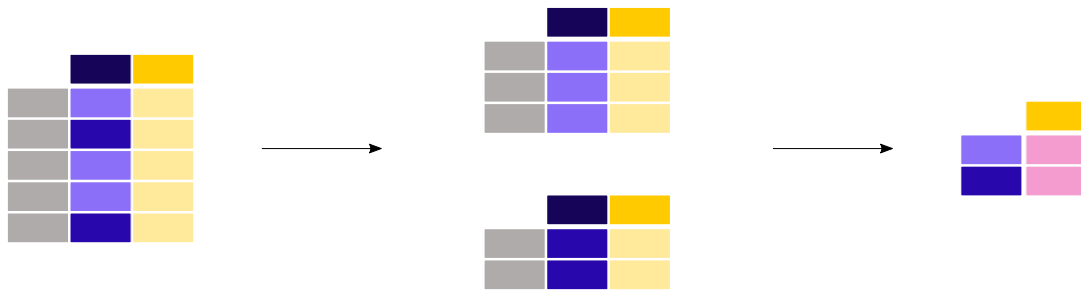


Figura 6.7\_Spark API Pandas: Formatos (Fuente: [pandas.pydata.org](https://pandas.pydata.org))

- **Remodelar la estructura de los datos**

Podemos cambiar la estructura de la tabla de datos de varias maneras. Podemos fundir (*melt()*) o pivotar (*pivot()*) el formato de la tabla de datos. Con agregaciones integradas, podemos crear una tabla dinámica con un solo comando.

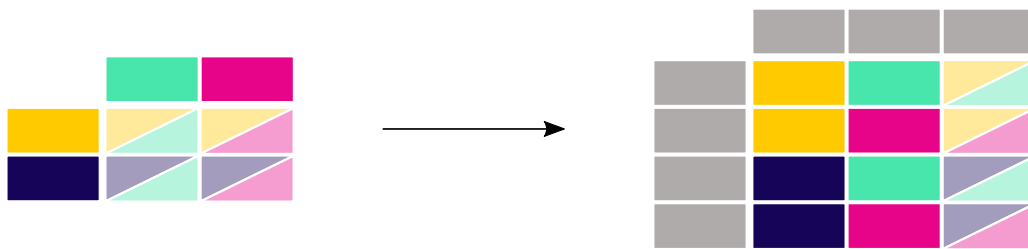


Figura 6.8\_Spark API Pandas: Melt (Fuente: pandas.pydata.org)

- **Combinar datos de varias tablas**

Podemos concatenar varias tablas tanto por columnas como por filas, ya que pandas proporciona operaciones de unión/fusión similares a las de una base de datos para combinar varias tablas.

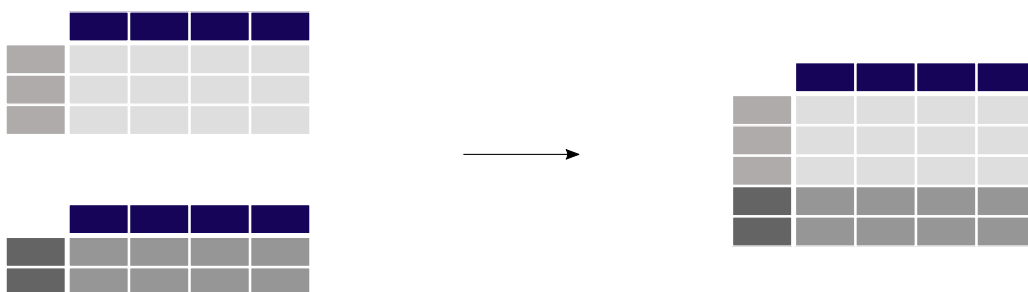


Figura 6.9\_Spark API Pandas: Concatenar (Fuente: pandas.pydata.org)

- **Manejar datos de series de tiempo**

Pandas tiene un gran soporte para series temporales y un amplio conjunto de herramientas para trabajar con fechas, horas y datos indexados en el tiempo.

- **Manipular datos de texto**

Pandas también proporciona una amplia gama de funciones para limpiar datos textuales y extraer información útil de ellos.

✓ **Success**

En temas de rendimiento **la máxima cantidad de registros que soporta pandas es de alrededor de 2 millones**. Si tienes un archivo con numero de registros mayor a esta cantidad una buena opción sería trabajarlo con **Spark**.

## 5.3 Quick Started Pandas on Pyspark

Vamos a ver una introducción a la API de pandas en Spark, dirigida principalmente a nuevos usuarios. Mostraremos algunas diferencias clave entre pandas y la API de pandas en Spark.

1. Importamos la API de pandas en Spark de la siguiente manera:

```
1 import pandas as pd
2 import numpy as np
3 import pyspark.pandas as ps
4 from pyspark.sql import SparkSession
```

### Creación de objetos

2. Creamos una serie de pandas-on-Spark pasando una lista de valores, permitiendo que la API de pandas en Spark cree un índice entero predeterminado:

```
1 s = ps.Series([1, 3, 5, np.nan, 6, 8])
2 s
3 .....
4 0    1.0
5 1    3.0
6 2    5.0
7 3    NaN
8 4    6.0
9 5    8.0
10 dtype: float64
```

3. Creamos un DataFrame de pandas-on-Spark pasando un dictado de objetos que se pueden convertir en series.

```
1 psdf = ps.DataFrame(
2     {'a': [1, 2, 3, 4, 5, 6],
3      'b': [100, 200, 300, 400, 500, 600],
4      'c': ["one", "two", "three", "four", "five", "six"]},
5     index=[10, 20, 30, 40, 50, 60])
6 psdf
7 .....
8      a    b    c
9  10  1  100 one
10  20  2  200 two
11  30  3  300 three
12  40  4  400 four
13  50  5  500 five
14  60  6  600 six
```

4. Creamos un DataFrame de pandas pasando una matriz numpy, con un índice de fecha y hora y columnas etiquetadas:

```
1 dates = pd.date_range('20130101', periods=6)
2 dates
3 .....
4 DatetimeIndex(['2013-01-01', '2013-01-02', '2013-01-03', '2013-01-04',
```

```

5         '2013-01-05', '2013-01-06'],
6         dtype='datetime64[ns]', freq='D')

```

```

1 pdf = pd.DataFrame(np.random.randn(6, 4), index=dates,
2                   columns=list('ABCD'))
3 pdf
4 A    B    C    D
5 2013-01-01  0.912558 -0.795645 -0.289115  0.187606
6 2013-01-02 -0.059703 -1.233897  0.316625 -1.226828
7 2013-01-03  0.332871 -1.262010 -0.434844 -0.579920
8 2013-01-04  0.924016 -1.022019 -0.405249 -1.036021
9 2013-01-05 -0.772209 -1.228099  0.068901  0.896679
10 2013-01-06  1.485582 -0.709306 -0.202637 -0.248766

```

5. Ahora, este DataFrame de pandas puede convertirse en un DataFrame de pandas on Spark

```

1 psdf = ps.from_pandas(pdf)
2 type(psdf)
3 ....
4 pyspark.pandas.frame.DataFrame

```

6. Se ve y se comporta igual que un DataFrame de pandas.

```

1 psdf
2 A    B    C    D
3 2013-01-01  0.912558 -0.795645 -0.289115  0.187606
4 2013-01-02 -0.059703 -1.233897  0.316625 -1.226828
5 2013-01-03  0.332871 -1.262010 -0.434844 -0.579920
6 2013-01-04  0.924016 -1.022019 -0.405249 -1.036021
7 2013-01-05 -0.772209 -1.228099  0.068901  0.896679
8 2013-01-06  1.485582 -0.709306 -0.202637 -0.248766

```

7. Además, es posible crear fácilmente un DataFrame de pandas en Spark desde Spark DataFrame.

8. Creando un Spark DataFrame a partir de pandas DataFrame

```

1 spark = SparkSession.builder.getOrCreate()
2 sdf = spark.createDataFrame(pdf)
3 sdf.show()
4 .....
5 +-----+-----+-----+
6 |          A|          B|          C|
7 +-----+-----+-----+
8 |  0.91255803205208|-0.7956452608556638|-0.28911463069772175|
9 | 0.18760566615081622|
10 |-0.05970271470242...|-1.233896949308984|  0.3166246451758431|
11 |-1.2268284000402265|

```

```

10 | 0.33287106947536615 | -1.2620100816441786 | -0.4348444277082644 |
    | -0.5799199651437185 |
11 | 0.9240158461589916 | -1.0220190956326003 | -0.405248880650239 |
    | -1.0360212104348547 |
12 | -0.7722090016558953 | -1.2280986385313222 | 0.0689011451939635 |
    | 0.8966790729426755 |
13 | 1.4855822995785612 | -0.7093056426018517 |
    | -0.2026366848847041 | -0.24876619876451092 |
14 | +-----+-----+-----+-----+
    | -----+

```

## 9. Creando pandas-on-Spark DataFrame desde Spark DataFrame.

```

1  psdf = sdf.pandas_api()
2  psdf
3  .....
4  A    B    C    D
5  0    0.912558 -0.795645 -0.289115  0.187606
6  1   -0.059703 -1.233897  0.316625 -1.226828
7  2    0.332871 -1.262010 -0.434844 -0.579920
8  3    0.924016 -1.022019 -0.405249 -1.036021
9  4   -0.772209 -1.228099  0.068901  0.896679
10 5    1.485582 -0.709306 -0.202637 -0.248766

```

## 10. Tipos específicos `dtypes`. Actualmente se admiten tipos que son comunes tanto para Spark como para pandas.

```

1  psdf.dtypes
2  .....
3  A    float64
4  B    float64
5  C    float64
6  D    float64
7  dtype: object

```

## 11. A continuación explicamos cómo mostrar las filas superiores del cuadro siguiente. Hay que tener en cuenta que los DataFrame de Spark no conservan el orden natural de forma predeterminada. El orden natural se puede preservar configurando la opción `compute.ordered_head`, pero provoca una sobrecarga de rendimiento con la clasificación interna.

```

1  psdf.head()
2  .....
3  A    B    C    D
4  0    0.912558 -0.795645 -0.289115  0.187606
5  1   -0.059703 -1.233897  0.316625 -1.226828
6  2    0.332871 -1.262010 -0.434844 -0.579920
7  3    0.924016 -1.022019 -0.405249 -1.036021
8  4   -0.772209 -1.228099  0.068901  0.896679

```

## 12. Mostrando el índice, las columnas y los numerosos datos subyacentes.

```

1 psdf.index
2 .....
3 Int64Index([0, 1, 2, 3, 4, 5], dtype='int64')

```

```

1 psdf.columns
2 .....
3 Index(['A', 'B', 'C', 'D'], dtype='object')

```

```

1 psdf.to_numpy()
2 .....
3 array([[ 0.91255803, -0.79564526, -0.28911463,  0.18760567],
4        [-0.05970271, -1.23389695,  0.31662465, -1.2268284 ],
5        [ 0.33287107, -1.26201008, -0.43484443, -0.57991997],
6        [ 0.92401585, -1.0220191 , -0.40524889, -1.03602121],
7        [-0.772209 , -1.22809864,  0.06890115,  0.89667907],
8        [ 1.4855823 , -0.70930564, -0.20263668, -0.2487662 ]])

```

## 12. Mostrando un resumen estadístico rápido de sus datos

```

1 psdf.describe()
2 .....
3 A    B    C    D
4 count  6.000000  6.000000  6.000000  6.000000
5 mean    0.470519  -1.041829  -0.157720  -0.334542
6 std    0.809428  0.241511  0.294520  0.793014
7 min   -0.772209  -1.262010  -0.434844  -1.226828
8 25%   -0.059703  -1.233897  -0.405249  -1.036021
9 50%    0.332871  -1.228099  -0.289115  -0.579920
10 75%    0.924016  -0.795645  0.068901  0.187606
11 max    1.485582  -0.709306  0.316625  0.896679

```

## 13. Transponiendo los datos

```

1 psdf.T
2 .....
3 0    1    2    3    4    5
4 A    0.912558  -0.059703  0.332871  0.924016  -0.772209  1.485582
5 B   -0.795645  -1.233897  -1.262010  -1.022019  -1.228099  -0.709306
6 C   -0.289115  0.316625  -0.434844  -0.405249  0.068901  -0.202637
7 D    0.187606  -1.226828  -0.579920  -1.036021  0.896679  -0.248766

```

## 14. Ordenando por índice

```

1 psdf.sort_index(ascending=False)
2 .....
3 A    B    C    D
4 5    1.485582  -0.709306  -0.202637  -0.248766
5 4   -0.772209  -1.228099  0.068901  0.896679
6 3    0.924016  -1.022019  -0.405249  -1.036021
7 2    0.332871  -1.262010  -0.434844  -0.579920
8 1   -0.059703  -1.233897  0.316625  -1.226828

```

```
9  0  0.912558  -0.795645  -0.289115  0.187606
```

## 15. Ordenando por valor

```
1  psdf.sort_values(by='B')
2  .....
3  A    B    C    D
4  2  0.332871  -1.262010  -0.434844  -0.579920
5  1  -0.059703  -1.233897  0.316625  -1.226828
6  4  -0.772209  -1.228099  0.068901  0.896679
7  3  0.924016  -1.022019  -0.405249  -1.036021
8  0  0.912558  -0.795645  -0.289115  0.187606
9  5  1.485582  -0.709306  -0.202637  -0.248766
```

## Datos perdidos (Missing Data)

16. La API de Pandas en Spark utiliza principalmente el valor `np.nan` para representar los datos perdidos. Por defecto no se incluye en los cálculos.

```
1  pdf1 = pdf.reindex(index=dates[0:4], columns=list(pdf.columns) + ['E'])
2  pdf1.loc[dates[0]:dates[1], 'E'] = 1
3  psdf1 = ps.from_pandas(pdf1)
4  psdf1
5  .....
6  A    B    C    D    E
7  2013-01-01  0.912558  -0.795645  -0.289115  0.187606  1.0
8  2013-01-02  -0.059703  -1.233897  0.316625  -1.226828  1.0
9  2013-01-03  0.332871  -1.262010  -0.434844  -0.579920  NaN
10 2013-01-04  0.924016  -1.022019  -0.405249  -1.036021  NaN
```

17. Eliminar cualquier fila a la que le falten datos.

```
1  psdf1.dropna(how='any')
2  .....
3  A    B    C    D    E
4  2013-01-01  0.912558  -0.795645  -0.289115  0.187606  1.0
5  2013-01-02  -0.059703  -1.233897  0.316625  -1.226828  1.0
```

18. Llenando datos perdidos.

```
1  psdf1.fillna(value=5)
2  .....
3  A    B    C    D    E
4  2013-01-01  0.912558  -0.795645  -0.289115  0.187606  1.0
5  2013-01-02  -0.059703  -1.233897  0.316625  -1.226828  1.0
6  2013-01-03  0.332871  -1.262010  -0.434844  -0.579920  5.0
7  2013-01-04  0.924016  -1.022019  -0.405249  -1.036021  5.0
```

## Operaciones

19. **Estadísticas.** Realizando una estadística descriptiva:



```

1 | psdf.mean()
2 | .....
3 | A      0.470519
4 | B     -1.041829
5 | C     -0.157720
6 | D     -0.334542
7 | dtype: float64

```

**20. Configuraciones Spark.** Se podrían aplicar varias configuraciones en PySpark internamente en la API de pandas en Spark. Por ejemplo, podemos habilitar la optimización de Arrow para acelerar enormemente la conversión interna de pandas. Consulte también la Guía de uso de PySpark para Pandas con [Apache Arrow](#) en la documentación de PySpark.

```

1 | rev = spark.conf.get("spark.sql.execution.arrow.pyspark.enabled") # Keep
  | its default value.
2 | ps.set_option("compute.default_index_type", "distributed") # Use default
  | index prevent overhead.
3 | import warnings
4 | warnings.filterwarnings("ignore") # Ignore warnings coming from Arrow
  | optimizations.

```

```

1 | spark.conf.set("spark.sql.execution.arrow.pyspark.enabled", True)
2 | %timeit ps.range(300000).to_pandas()
3 | .....
4 | 900 ms ± 186 ms per loop (mean ± std. dev. of 7 runs, 1 loop each)

```

```

1 | spark.conf.set("spark.sql.execution.arrow.pyspark.enabled", False)
2 | %timeit ps.range(300000).to_pandas()
3 | .....
4 | 3.08 s ± 227 ms per loop (mean ± std. dev. of 7 runs, 1 loop each)

```

```

1 | ps.reset_option("compute.default_index_type")
2 | spark.conf.set("spark.sql.execution.arrow.pyspark.enabled", prev) # Set
  | its default value back.

```

## Agrupamiento (Grouping)

21. Por “agrupar por” nos referimos a un proceso que involucra uno o más de los siguientes pasos:

- Dividir los datos en grupos según algunos criterios.
- Aplicar una función a cada grupo de forma independiente
- Combinar los resultados en una estructura de datos.

```

1 | psdf = ps.DataFrame({'A': ['foo', 'bar', 'foo', 'bar',
2 |                           'foo', 'bar', 'foo', 'foo'],
3 |                     'B': ['one', 'one', 'two', 'three',

```

```

4         'two', 'two', 'one', 'three'],
5         'C': np.random.randn(8),
6         'D': np.random.randn(8)})
7 psdf
8 .....
9  A   B   C   D
10 0  foo one 1.039632 -0.571950
11 1  bar one 0.972089 1.085353
12 2  foo two -1.931621 -2.579164
13 3  bar three -0.654371 -0.340704
14 4  foo two -0.157080 0.893736
15 5  bar two 0.882795 0.024978
16 6  foo one -0.149384 0.201667
17 7  foo three -1.355136 0.693883

```

22. Agrupar y luego aplicar la función `sum()` a los grupos resultantes.

```

1 psdf.groupby('A').sum()
2 .....
3      C      D
4  A
5 bar 1.200513 0.769627
6 foo -2.553589 -1.361828

```

23. La agrupación por varias columnas forma un índice jerárquico y nuevamente podemos aplicar la función de suma.

```

1 psdf.groupby(['A', 'B']).sum()
2      C      D
3  A   B
4 foo one 0.890248 -0.370283
5 two -2.088701 -1.685428
6 bar three -0.654371 -0.340704
7 foo three -1.355136 0.693883
8 bar two 0.882795 0.024978
9 one 0.972089 1.085353

```

## Plotting

Estos métodos deben ejecutarse en un notebook tipo Jupyter, por ejemplo: databricks o jupyter en pyspark(lo vemos en el siguiente punto). También en google colab, entre otros.

```

1 pser = pd.Series(np.random.randn(1000),
2                  index=pd.date_range('1/1/2000', periods=1000))
3 psrer = ps.Series(pser)
4 psrer = psrer.cummax()
5 psrer.plot()

```

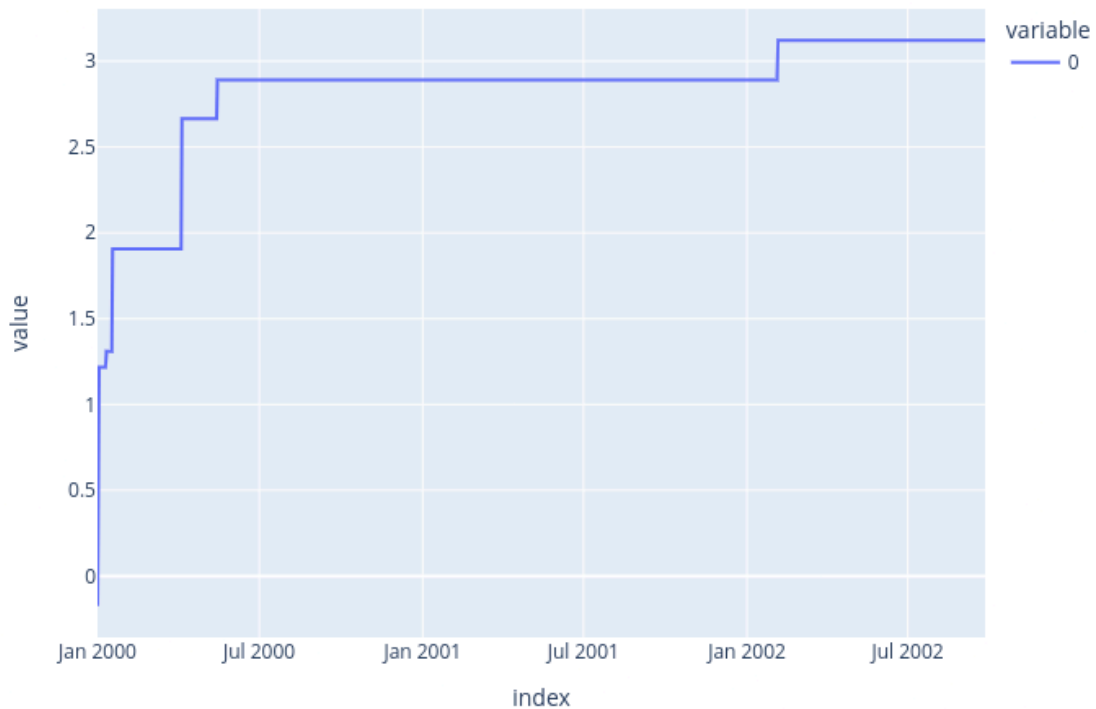


Figura 6.9\_Spark API Pandas: Plot 1

1. En un DataFrame, es conveniente usar el método `plot()` para trazar todas las columnas con etiquetas:

```
1 pdf = pd.DataFrame(np.random.randn(1000, 4), index=pser.index,  
2                      columns=['A', 'B', 'C', 'D'])  
3 psdf = ps.from_pandas(pdf)  
4 psdf = psdf.cummax()  
5 psdf.plot()
```

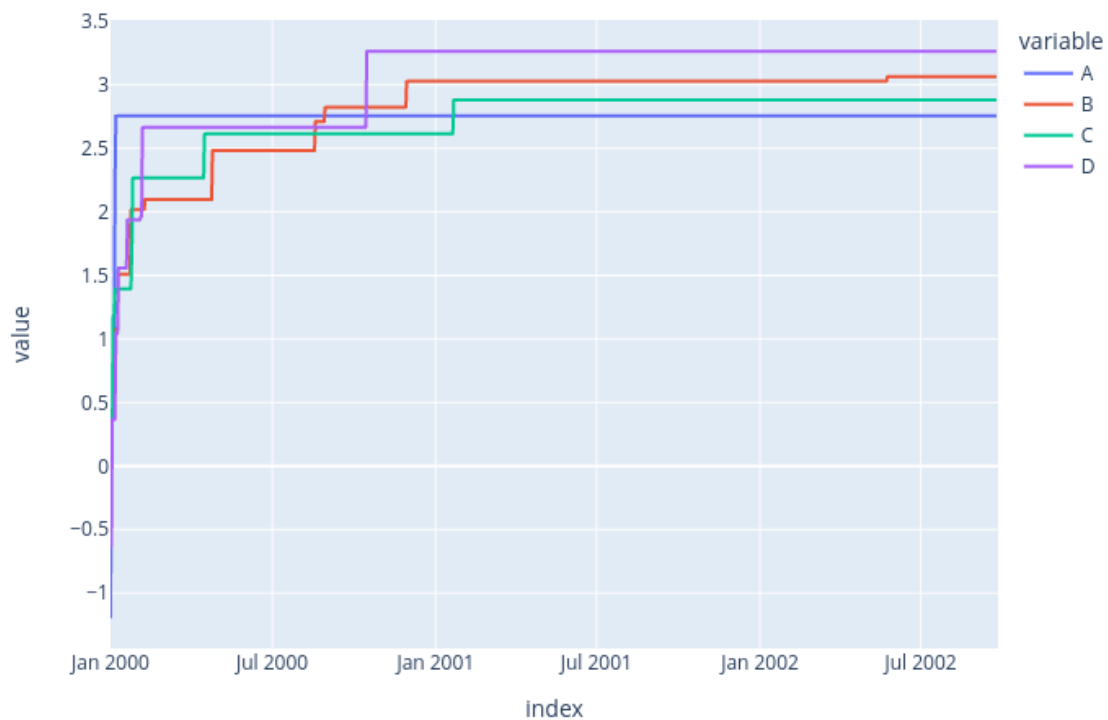


Figura 6.10\_Spark API Pandas: Plot 2

## Entrada/Salida de datos

### 25. CSV

```
1 psdf.to_csv('foo.csv')
2 ps.read_csv('foo.csv').head(10)
3 .....
```

	A	B	C	D
0	-1.187097	-0.134645	0.377094	-0.627217
1	0.331741	0.166218	0.377094	-0.627217
2	0.331741	0.439450	0.377094	0.365970
3	0.621620	0.439450	1.190180	0.365970
4	0.621620	0.439450	1.190180	0.365970
5	2.169198	1.069183	1.395642	0.365970
6	2.755738	1.069183	1.395642	1.045868
7	2.755738	1.069183	1.395642	1.045868
8	2.755738	1.069183	1.395642	1.045868
9	2.755738	1.508732	1.395642	1.556933

**26. Parquet.** Parquet es un formato de archivo eficiente y compacto para leer y escribir más rápido.

```
1 df.to_parquet('bar.parquet')
2 ps.read_parquet('bar.parquet').head(10)
3 .....
```

4	A	B	C	D
5	0	-1.187097	-0.134645	0.377094
6	1	0.331741	0.166218	0.377094
7	2	0.331741	0.439450	0.377094
8	3	0.621620	0.439450	1.190180
9	4	0.621620	0.439450	1.190180
10	5	2.169198	1.069183	1.395642
11	6	2.755738	1.069183	1.395642
12	7	2.755738	1.069183	1.395642
13	8	2.755738	1.069183	1.395642
14	9	2.755738	1.508732	1.395642

27. **Spark IO.** Además, la API de pandas en Spark es totalmente compatible con las diversas fuentes de datos de Spark, como ORC y una fuente de datos externa.

1	psdf.to_spark_io('zoo.orc', format="orc")			
2	ps.read_spark_io('zoo.orc', format="orc").head(10)			
3	.....			
4	A	B	C	D
5	0	-1.187097	-0.134645	0.377094
6	1	0.331741	0.166218	0.377094
7	2	0.331741	0.439450	0.377094
8	3	0.621620	0.439450	1.190180
9	4	0.621620	0.439450	1.190180
10	5	2.169198	1.069183	1.395642
11	6	2.755738	1.069183	1.395642
12	7	2.755738	1.069183	1.395642
13	8	2.755738	1.069183	1.395642
14	9	2.755738	1.508732	1.395642