



Apache Spark SQL

Diego J. Bodas Sagi

- ¿Por qué?
- Contextos
- Arquitectura
- Datasets
- DataFrames
- Varios
- A programar...





¿Por qué?

- ➊ Porque nos nutrimos de datos que están almacenados persistentemente
- ➋ Porque SQL es un estándar extendido y fácil de usar
- ➌ Porque recurrir a Hive implica depender de Hadoop
- ➍ Porque consideramos que se puede mejorar lo que tenemos...





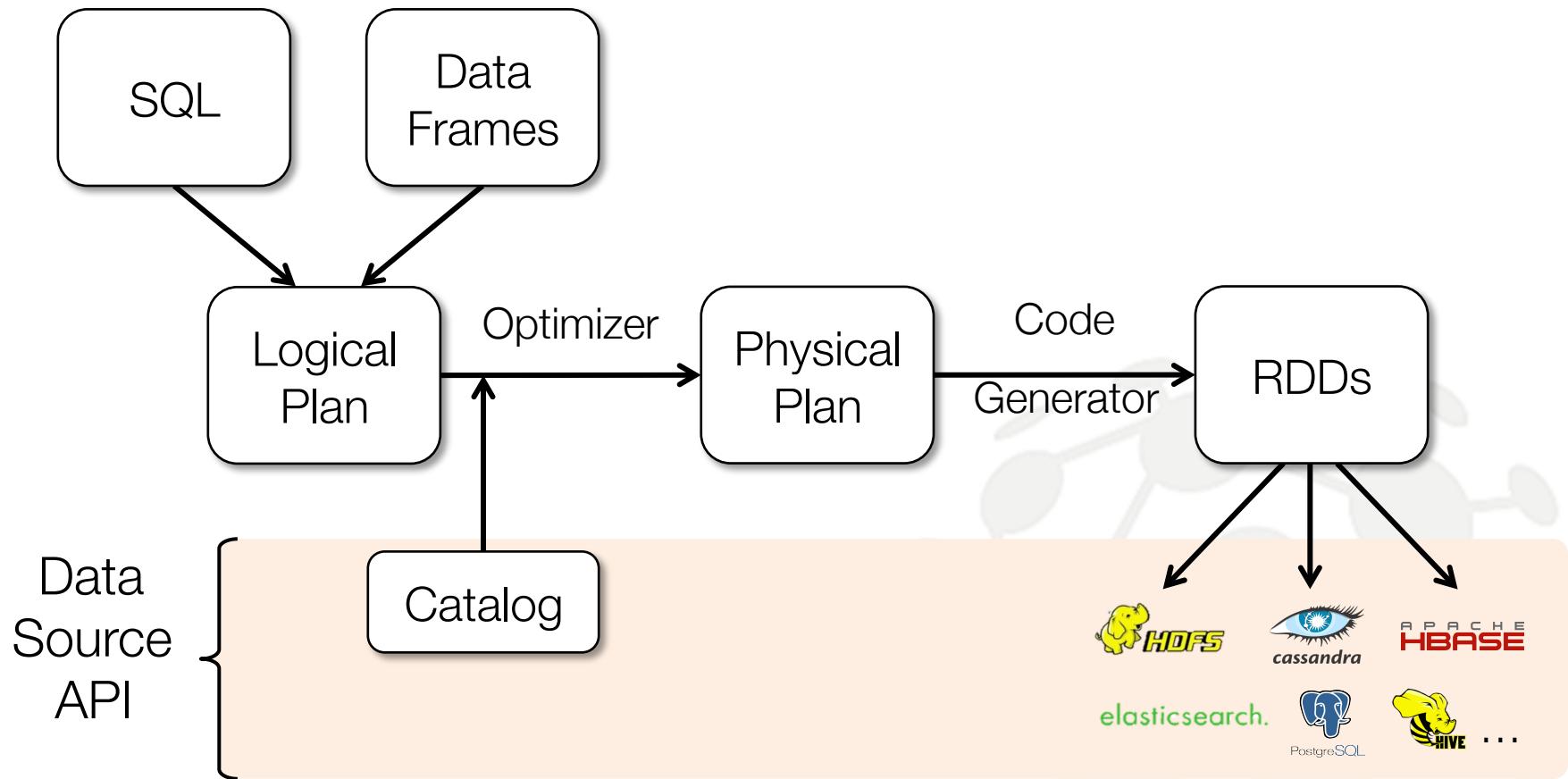
Contextos SQL

- Toda aplicación Spark depende de un contexto
- El contexto recoge toda la información necesaria para acceder al clúster y usar sus recursos
- El *master* es el encargado de gestionar el cluster
- En este caso, necesitamos un `sqlContext`
- Este contexto nos permite crear **DataFrames**

| Master Parameter | Description |
|--------------------------------|--|
| <code>local</code> | run Spark locally with one worker thread (no parallelism) |
| <code>local[K]</code> | run Spark locally with K worker threads (ideally set to number of cores) |
| <code>spark://HOST:PORT</code> | connect to a Spark standalone cluster; PORT depends on config (7077 by default) |
| <code>mesos://HOST:PORT</code> | connect to a Mesos cluster; PORT depends on config (5050 by default) |



Arquitectura





Datasets

Datasets

RDDs

- Functional Programming
- Type-safe

Dataframes

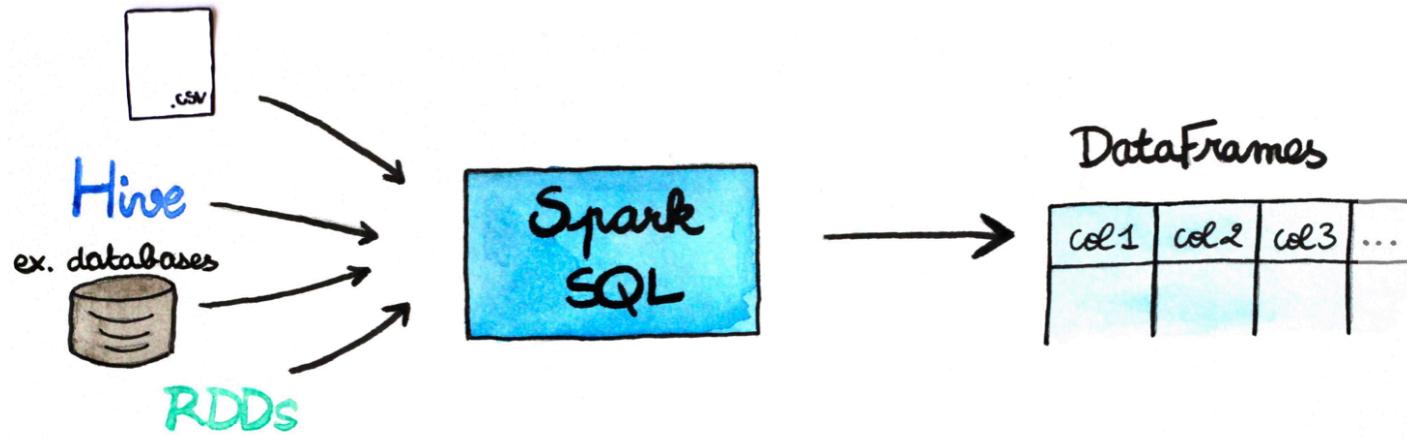
- Relational
- Catalyst query optimization
- Tungsten direct/packed RAM
- JIT code generation
- Sorting/suffling without deserializing

Spark





DataFrames



Spark

Cada fila de un dataset es un objeto **Row**





DataFrames: transformations and actions

| Transformations | Actions |
|-----------------|---------|
| filter | count |
| select | collect |
| drop | show |
| join | take |

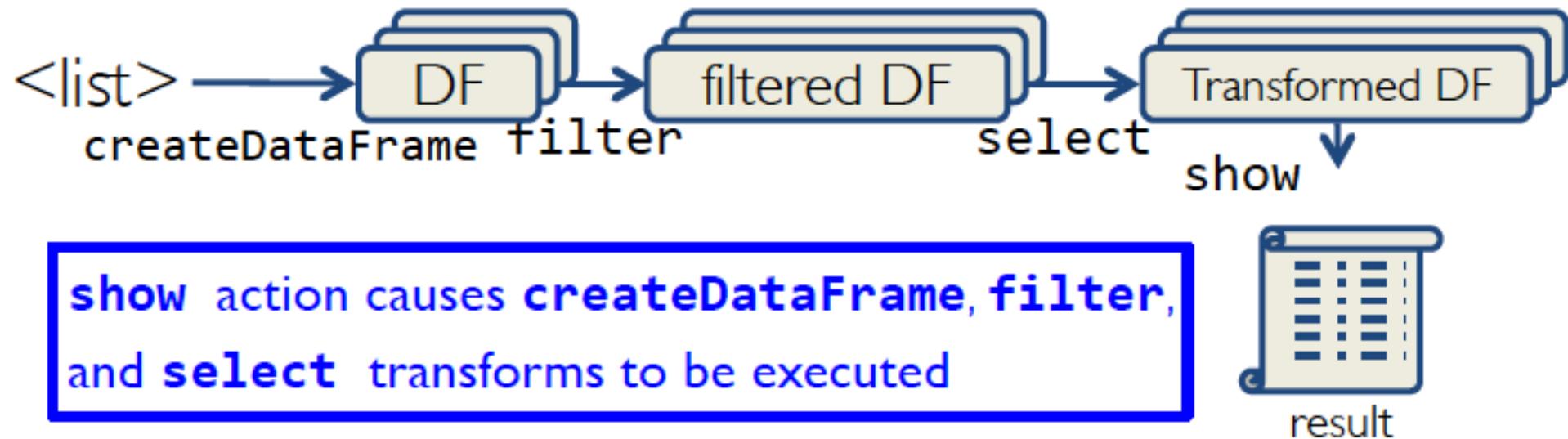
Lazy evaluation

Transformed DF is executed when action runs on it





DataFrames: transformations and actions





Ejemplo

```
>>> data = [('Alice', 1), ('Bob', 2)]  
>>> data  
[('Alice', 1), ('Bob', 2)]  
  
>>> df = sqlContext.createDataFrame(data)
```

```
[Row(_1=u'alice', _2=1), Row(_1=u'Bob', _2=2)]
```

```
>>> sqlContext.createDataFrame(data, ['name', 'age'])
```

```
[Row(name=u'alice', age=1), Row(name=u'Bob', age=2)]
```

No computation occurs with
`sqlContext.createDataFrame()`

- Spark only records how to create the DataFrame





Crear un DataFrame desde Pandas o R

- ➊ No confundir DataFrames de Spark con DataFrames de Python o R

```
# Create a Spark DataFrame from Pandas  
>>> spark_df = sqlContext.createDataFrame(pandas_df)
```





Leer y cargar datos en el DataFrame

```
>>> df = sqlContext.read.text("README.txt")  
  
>>> df.collect()  
[Row(value=u'hello'), Row(value=u'this')]
```





Más datos

- Persistencia de un DataFrame: **cache()**
- Nunca usar collect() (*todos*) en producción, en su lugar usar take(n)
- Lanzar una consulta es muy sencillo

```
# SQL can be run over SchemaRDDs that have been registered as a table.  
results = sqlContext.sql("SELECT name FROM people")
```

- Pero primero se debe haber registrado como tabla

```
// Register the DataFrames as a table.  
peopleDataFrame.registerTempTable("people")
```





Salvar y persistir

DataFrames can be saved as Parquet files, maintaining the schema information.
schemaPeople.write.parquet("people.parquet")





Simplifica el código



Using RDDs

```
data = sc.textFile(...).split("\t")
data.map(lambda x: (x[0], [int(x[1]), 1])) \
    .reduceByKey(lambda x, y: [x[0] + y[0], x[1] + y[1]]) \
    .map(lambda x: [x[0], x[1][0] / x[1][1]]) \
    .collect()
```

Using SQL

```
SELECT name, avg(age)
FROM people
GROUP BY name
```

Using Pig

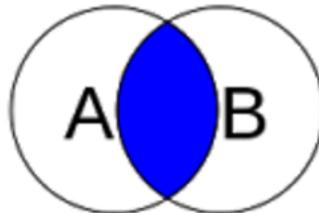
```
P = load '/people' as (name, name);
G = group P by name;
R = foreach G generate ... AVG(G.age);
```

Using DataFrames

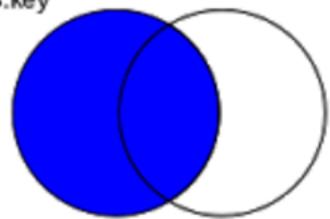
```
sqlCtx.table("people") \
    .groupBy("name") \
    .agg("name", avg("age")) \
    .collect()
```



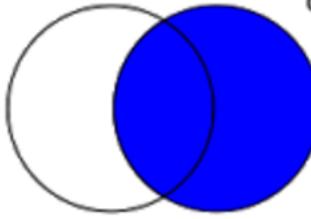
SELECT <fields>
FROM TableA A
INNER JOIN TableB B
ON A.key = B.key



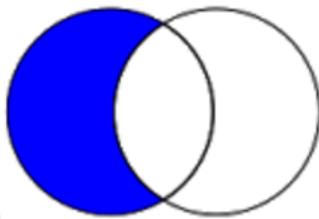
SELECT <fields>
FROM TableA A
LEFT JOIN TableB B
ON A.key = B.key



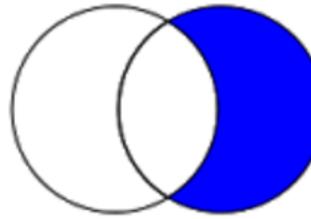
SELECT <fields>
FROM TableA A
RIGHT JOIN TableB B
ON A.key = B.key



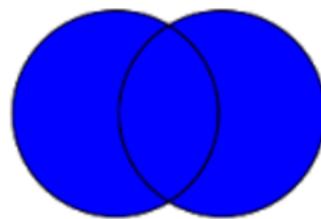
SELECT <fields>
FROM TableA A
LEFT JOIN TableB B
ON A.key = B.key
WHERE B.key IS NULL



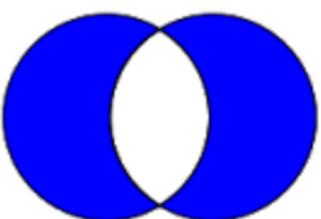
SELECT <fields>
FROM TableA A
RIGHT JOIN TableB B
ON A.key = B.key
WHERE a.key IS NULL



SELECT <fields>
FROM TableA A
FULL OUTER JOIN TableB B
ON A.key = B.key



SELECT <fields>
FROM TableA A
FULL OUTER JOIN TableB B
ON A.key = B.key
WHERE A.key IS NULL
OR B.key IS NULL





Your turn

The screenshot shows a developer's workspace with a code editor and a file browser. The code editor displays a PHP controller class named `UpdateController.php`. The file browser on the left shows the project structure of a Zend Framework application, including modules like `Application`, `Comments`, `Inbox`, `Index`, `Posts`, and `Update`, and various configuration files and storage components.

```
<?php  
  
namespace Application\Controller;  
  
use Application\Module;  
use Zend\Mail\Storage\Map;  
use Zend\Mail\Storage;  
use Zend\Mvc\Controller\ActionController;  
use Zend\View\Model\JsonModel;  
use Application\InboxController;  
use Application\PostController;  
  
class UpdateController extends AbstractActionController  
{  
    const COMMAND_PUBLISH = 'publish';  
    const COMMAND_EDIT = 'edit';  
    const COMMAND_DRAFT = 'draft';  
    const COMMAND_TRASH = 'trash';  
    const COMMAND_REPLY = 'reply';  
    const COMMAND_SUBSCRIBE = 'subscribe';  
  
    const REGEX_PATTERN_COMMAND = '[a-zA-Z]+\\s';  
  
    public function indexAction()  
    {  
        echo "banana":exit;  
        $inbox = new InboxController();  
        $queue = $inbox->indexAction();  
        foreach ($queue as $message) {  
            print_r($message);  
            /*preg_match(self::REGEX_PATTERN_COMMAND, $message[InboxController::KEY_SUBJECT], $commandsFound);*/  
            switch (trim($commandsFound[0])){  
                case self::COMMAND_PUBLISH:  
                    break;  
                case self::COMMAND_
```

