

# Traitement des données structurées avec *Spark*

# Données structurées sur Spark: DataSet/DataFrame/Spark SQL

- Ce sont des RDDs structurées dans lesquelles les données sont sous format tabulaire
- Une collection immutable d'objets fortement typés qui sont mappés vers un schéma relationnel
  - · Elles peuvent être assimilées à des tables relationnelles
- En interne:
  - Catalyst: application des optimisations inspirées des technique d'évaluation de requêtes dans les systèmes relationnelles grâce à la connaissance du schéma et des types de données
  - *Tungsten*: représentation en format binaire des données pour minimiser les opérations de dé-sérialisation
- 2 types:
  - DataSet: quand le type des colonnes est connu
  - DataFrame: un DataSet avec des lignes dont le schéma n'est pas connu, i.e. de type générique appelé Row

## DataFrame & DataSet

Le moteur Spark SQL permet d'effectuer des requêtes similaires à SQL sur des données distribuées. Il est dédié aux traitements de données structurées et semi structurées via l'API de haut niveau: DataFrame/DataSet.

#### **DataFrame**

- Collection de donnée distribuée.
- Contrairement aux RDD, les données sont organisées par colonne
- Conceptuellement, un DataFrame est l'équivalent d'une table dans une base de données relationnelle.
- Une DataFrame est non typé: toutes les entrées sont de type générique appelé Row
- Une DataFrame supporte des transformations basées sur des paradigmes fonctionnels ou relationnels.
- Elle peut être créée à partir:
  - 1. d'une RDD.
  - 2. d'une table (hive, SGBD, etc.)
  - 3. d'une ou plusieurs DataFrame via des opérations de transformation.
- Disponible en Java, Scala, Python et R

#### DataSet

- Une DataSet est une DataFrame fortement typé.
- Disponible seulement en Scala et Java.

# API DataFrame (1)

```
public static void main (String[] args) {
                                                                            SparkSession spark = SparkSession
                                                                                        .builder()
                                                                                        .appName("Java Spark SQL basic example")
public class Main {
                                                                            .getOrCreate();
JavaRDD<Person> peopleRDD = spark.read()
    public static void main(String[] args) {
         SparkSession spark = SparkSession.builder()
                                                                                .textFile("examples/src/main/resources/people.txt")
                  appName("Java Spark SQL basic example")
                                                                                .javaRDD()
                   .getOrCreate();
                                                                                .map(new Function<String, Person>() {
   private static final long serialVersionUID = 1L;
         // Read CVS file
         Dataset<Row> dfFromCSV = spark.read()
                                                                                    @Override
                  .format("csv")
                                                                                    public Person call(String line) throws Exception {
                  .option("header", "true")
                                                                                        String[] parts = line.split(",");
Person person = new Person();
                  .load("csvfile.csv");
                                                                                        person.setName(parts[0]);
         // Read Parquet file
                                                                                        person.setAge(Integer.parseInt(parts[1].trim()));
         Dataset<Row> dfFromParquet = spark.read()
                                                                                        return person:
                   .parquet("parquetfile.parquet");
                                                                                   }
         // Read ORC File
                                                                                });
         Dataset<Row> dfFromORC = spark.read()
                                                                            //Apply a schema to an RDD of JavaBeans to get a DataFrame
                 .orc("orcfile.orc");
                                                                            Dataset<Row> peopleDF = spark.createDataFrame(peopleRDD,
                                                                                   Person.class);
    }
                                                                            peopleDF.show();
}
                                                                    }
                                                                                           Construction à partir d'un RDD
        Construction à partir de la lecture d'un fichier
                              Dataset<Row> selectedColRows = dfFromCSV.select(col("column name"),col("age").plus(1));
                              Dataset<Row> filtredRowd = dfFromCSV.filter(col("column_name").geq(13));
         Manipulation
              de
                              Dataset<Row> results = dfFromCSV.groupBy("column_name").avg("column_name") ;
         DataFrames
                              JavaRDD<Row> rowJavaRDD = dfFromCSV.toJavaRDD():
```

public class DFRDD {

# API DataFrame (2)

```
java.util.List;
org.apache.spark.api.java.JavaRDD;
org.apacne.spark.api.java.JavaRDD;
org.apache.spark.api.java.function.Function,
org.apache.spark.sql.*;
org.apache.spark.sql.types.DataTypes;
org.apache.spark.sql.types.StructField;
org.apache.spark.sql.types.StructType;
  elass RDDtoDFSchemaSpec (
lic static void main (String[] args) {
SparkSession spark = SparkSession.bullder().appName("Java Spark SQL basic example").getOrCreate();
  String schemaString = "name age";
  List<StructField> fields = new ArrayList<StructField>();

for (String fieldName : schemaString.split(")) {
    StructField field = DataTypes.createStructField(fieldName, DataTypes.StringType, true);
  // Convert records of the RDD (people) to Rows
JavaRDD<Row> rowRDD = peopleRDD.map(new Function<String, Row>() {
          public Row call(String record) throws Exception {
   String[] attributes = record.split(",");
                 return RowFactory.create(attributes[0], attributes[1].trim());
    // Apply the schema to the RDD
Dataset<Row> peopleDataFrame = spark.createDataFrame(rowRDD, schema).
```

## **API DataSet**

Construction à partir de la lecture d'un fichier

## **API DataSet**

```
import org.apache.spark.api.java.JavaRDD;
import org.apache.spark.sql.*;
import org.apache.spark.api.java.function.Function;
                                                                                               Construction à partir d'un RDD
public class DFRDD {
    public static void main (String[] args) {
       SparkSession spark = SparkSession
                .builder()
                 .appName("Java Spark SQL basic example").getOrCreate();
        JavaRDD<Person> peopleRDD = spark.read()
                .textFile("examples/src/main/resources/people.txt").javaRDD()
                .map(new Function<String, Person>() {
                     public Person call(String line) throws Exception {
                         String[] parts = line.split(",");
Person person = new Person();
                         person.setName(parts[0]);
                         \verb"person.setAge(Integer.parseInt(parts[1].trim()))";
                         return person;
     Encoder<Person> encoder = Encoders.bean(Person.class);
     Dataset<Row> peopleDF = spark.createDataFrame(peopleRDD, Person.class);
     // DataFrame --> DataSet
     Dataset<Person> datasetFromDF = peopleDF.as(encoder);
     // RDD --> DataSet
     Dataset<Person> datasetFromRDD = spark.createDataset(peopleRDD.rdd(), encoder);
```

## **API DataSet**

```
Dataset<Person> selectedColRows = datasetFromDF.select(col("column_name"),col("age")).as(encoder) ;
Dataset<Person> filteredDS = datasetFromDF.filter(new FilterFunction<Person>() {
    @Override
    public boolean call(Person person) throws Exception { return false; } });
Dataset<Integer> results =
    datasetFromDF.groupBy("column_name").avg("column_name").as(Encoders.INT());

// DataSet --> DataFrame
Dataset<Row> dSToDF = selectedColRows.toDF();

// DataSet --> JavaRDD
JavaRDD<Person> dsToJavaRDD = selectedColRows.toJavaRDD();
```

Manipulation de DataFrames

## Spark SQL

- Spark SQL est un module de Spark permettant de requêter et d'appliquer des transformations sur une collection de données en utilisant un langage dit SQL-Like
- Il peut être également vu comme un moteur SQL distribué
- Avec les trois APIs DataFrame, DataSet et SQL, les développeurs ont la capacité de choisir l'API qui répondrait le mieux à leurs besoins et switcher d'une API à une autre dans la même application
  - Une DataFrame ou un DataSet peut être exposé comme une table relationnelle qui peut être manipulée en utilisant un langage SQL-Like

```
Encoder<Person> encoder = Encoders.bean(Person.class);
DatasettRow> dfFromRDD = spark.createDataFrame (peopleRDD,
Person.class);
// DataFrame --> DataSet
Dataset<Person> datasetFromDF = dfFromRDD.as(encoder);
// RDD --> DataFrame
Dataset<Person> datasetFromRDD = spark.createDataset(peopleRDD.rdd(),
encoder);
// Creer les tables relationnelles
try {
    dfFromRDD.createGlobalTempView("people_ds");
    datasetFromDF.createGlobalTempView("people_ds");
    datasetFromDF.createGreplaceTempView("people_local_temp");
} catch (AnalysisException e) {
        e.printStackTrace();
}

// Requeter les tables crees
Dataset<Row> result_ds = spark.sql("select * from people_ds");
Dataset<Row> result_ds_tmp =
```

## Comparaison RDD, DataFrame/DataSet

	RDD	Dataframe	Dataset
Représentation des données	Ensemble d'objets JAVA	<ul> <li>Une collection organisée en colonnes nommées</li> </ul>	Une extension des DataFrame mais fortement typées
Interopérabilité	Peut créer un DF à partir d'une RDD et vice-versa	<ul> <li>Peut être créée à partir d'un RDD</li> <li>Pas de conversion DF vers RDD</li> </ul>	<ul> <li>Peut être créé à partir d'un RDD ou DF</li> <li>Peut être converti vers un RDD</li> </ul>
Erreur de Syntaxe	À la compilation	À la compilation	À la compilation
Erreur d'analyse	À la compilation	À l'exécution	À la compilation
Langages	Scala, Java, Python	Scala, Java, python, R	Scala, Java
Paradigme d'interrogation	Fonctionnelle	Fonctionnelle, DSL, SQL	Fonctionnelle, DSL, SQL
Sérialisation	Utilise la sérialisation de JAVA pour stocker et distribuer les objets (couteux)	Les données sont transformées en format binaire et application des transformation en utilisant Tungsten     Génération dynamique de bytecode	Utilise les encodeurs et stocke les données sous format tabulaire en utilisant Tungsten
Optimisation	Pas d'optimisation	Utilise Tungsten et Catalyst	Utilise Tungsten et Catalyst
Garbage collection	• Oui	Pas besoin	Pas besoin de GC pour détruire les objets car la sérialisation est en offheap et géré par Tungsten
Utilisation de la mémoire	Pas très efficace à cause de la sérialisation	Efficace car utilise une sérialisation en off-heap	Très efficace car les opérations sont directement faites sur les données sérialisées

## RDD, DataFrame/DataSet: quoi utiliser quand?

#### **RDD**

#### • Type de données à traiter:

- Les données <u>ne sont pas structurées</u> (media, textes).
- Pas besoin d'imposer un schéma ou une structure lors du traitement de données.
- On veut utiliser la programmation fonctionnelle plutôt qu'un DSL

#### · Optimisation ou contrôle bas niveau:

- Besoin d'avoir un contrôle niveau bas sur les données
- On peut renoncer à certains avantages en termes d'optimisation de performance disponibles sur l'API DataFrames/Dataset.

#### DataFrame/DataSet

#### Type de données à traiter:

- Les traitements exige à la fois des expressions du haut niveau (sql, aggregation, sum, etc) et des fonctions lambda sur des <u>données semi-structurées</u> <u>ou structurées</u>.
- Besoin d'une sémantique riche, des abstractions de haut niveau et des API DSL.
- Besoin d'utiliser des objets typés (DataSet).

### Optimisation ou contrôle bas niveau:

- Représentation binaire des données
- Tirez parti de l'optimisation Catalyst et de la génération de code efficace de Tungsten.

# Spark SQL

https://towards datascience.com/sql-at-scale-with-apache-spark-sql-and-data frames-concepts-architecture-and-examples-c567853a702 for the state of the state of

