



## YCBS-257 - Data at Scale

---

### Workshop 9

### In Memory Data Processing

#### General Instructions:

Apache Spark is an open source parallel processing framework for running large-scale data analytics applications across clustered computers. It can handle both batch and real-time analytics and data processing workloads. Spark was designed for executing complex multi-stage applications, like machine learning algorithms, and interactive ad hoc queries.

This series of exercises in Scala and PySpark are designed to familiarize you with Apache Spark ecosystem.

To perform the exercises you will need a Scala and / or PySpark command interpreter (Python for Spark). You can use either the command line or Apache Zeppelin (recommended) to perform the workshop.

Start your Cloudera QuickStart VM to complete the workshop

#### Online resources:

<https://zeppelin.apache.org/>

<https://spark.apache.org/>

#### Exercise 1: WordCount in Scala

In this exercise, we will approach the word count problem by using Scala with Spark. You might use Zeppelin to enter, test and run your code or use the command line interface.

1. Create a new directory `/sparklab` on HDFS
2. Put the file `5000-8.txt` inside this directory
3. Create a new Zeppelin note `wordcount`

Zeppelin Interpreter	<code>%spark</code>
Load the file from HDFS	<code>val textFile = sc.textFile("/sparklab/5000-8.txt")</code>
Count the number of elements	<code>textFile.count()</code>
Show the 1 <sup>st</sup> line	<code>textFile.first()</code>



Find how many 'The' occurrences in the document	<pre>val linesWithSpark = textFile.filter(line =&gt; line.contains("The")).count()</pre>
Get the longest line	<pre>textFile.map(line =&gt; line.split(" ").size).reduce((a, b) =&gt; if (a &gt; b) a else b)</pre>
The same using the Max function Import Java Math library to be able to use the Max function	<pre>import java.lang.Math  textFile.map(line =&gt; line.split(" ").size).reduce((a, b) =&gt; Math.max(a, b))</pre>
Create a flat structure using flatMap() function and run a ReduceByKey to get the final result	<pre>val wordCounts = textFile.flatMap(line =&gt; line.split(" ")).map(word =&gt; (word, 1)).reduceByKey((a, b) =&gt; a + b)  wordCounts.collect()</pre>

## Exercise 2: Babies names analysis using Spark

In this exercise you will analyze the babies' names in the state of New York given between 2007 and 2016.

1. Copy the **Baby\_Names\_\_Beginning\_2007.csv** file into **/sparklab** HDFS directory
2. Create a new Zeppelin note: **babies**

Spark Interpreter	<pre>%spark</pre>
RDD creation	<pre>val babyNames = sc.textFile("/sparklab/Baby_Names__Beginning_2007.csv")  babyNames.count()</pre>
Print the header line	<pre>babyNames.first()</pre>
Split the columns and do a distinct count (First Name column)	<pre>val rows = babyNames.map(line =&gt; line.split(","))  rows.map(row =&gt; row(2)).distinct.count</pre>



Count all occurrences of the first name "DAVID"	<pre>val davidRows = rows.filter(row =&gt; row(1).contains("DAVID"))  davidRows.count()</pre>
Limit results to departments where first name has been assigned more than 10 times	<pre>davidRows.filter(row =&gt; row(4).toInt &gt; 10).count()</pre>
Perform a distinct count of departments where 'David' appear more than 10 times	<pre>davidRows.filter(row =&gt; row(4).toInt &gt; 10).map( r =&gt; r(2) ).distinct.count</pre>
Display the 100 most common first names	<pre>val names = rows.map(name =&gt; (name(1),1)).reduceByKey((a,b) =&gt; a + b).sortBy(-_._2) //.foreach ( println _)  names.take(100).foreach (println _)</pre>
Remove the header line and Count all rows and print the 20 first names	<pre>val filteredRows = babyNames.filter(line =&gt; !line.contains("Count")).map(line =&gt; line.split(","))  filteredRows.map ( n =&gt; (n(1), n(4).toInt)).reduceByKey((a,b) =&gt; a + b).sortBy(- _._2).take(20).foreach (println _)</pre>

### Exercise 3: Explore Spark RDDs

In this exercise you will count negative comments from a public survey. As a data source we will use [Rotten Tomatoes Kaggle Dataset](#) which is a sentiments analysis about movies.

1. Copy the `train.tsv` file into `/sparklab` HDFS directory
2. Create a new Zeppelin note: `rddSpark`

```
%spark
```

```
SC
```

1. Load data into a new RDD

```
val data = sc.textFile("/sparklab/train.tsv")
```



2. Run Action method to 'force' Spark to evaluate the previous operation

```
data.getClass
```

3. Split the loaded data according to the separator character

```
val reviews = data.map(_.split("\t"))
```

4. Count the loaded data in the RDDs

```
reviews.count()
```

```
// Long: 156061
```

5. Count the words of the loaded document

```
val wordCounts = reviews.flatMap(x => x(2).split(" ")).map((_, 1)).reduceByKey((a, b) => a + b).sortBy(_._2, false)
```

```
wordCounts.take(5) // get top 5
```

```
reviews.flatMap(x => x(2).split(" ")).take(5)
```

6. Let's look at the third column

```
reviews.flatMap(x => x(2).split(" ")).map((_, 1)).take(5)
```

```
// MapReduce Style
reviews.flatMap(x => x(2).split(" ")).map((_, 1)).reduceByKey((a, b) => a + b).take(5)
```

7. Remove headers and count negative comments

```
val trueReviews = reviews.filter(x => x(0) != "PhraseId")
```

```
val negativeReviews = trueReviews.filter(x => x(3).toInt == 0)
```

```
val negWordCounts = negativeReviews.flatMap(x => x(2).split(" ")).map((_, 1)).reduceByKey((a, b) => a + b).sortBy(_._2, false)
```

```
negWordCounts.take(50)
```

## Exercise 4: Querying data using Spark SQL

In this exercise you will write a few lines of code in Scala to query a Spark RDD with SQL language. As a data source you will use the [Stations\\_2017.csv](#) file

1. Copy the [customers.txt](#) file into [/sparklab](#) HDFS directory
2. Create a new Zeppelin note: [sparksql](#)

Zeppelin  
Interpreter

```
%spark
```



Create a DataFrame of Station objects from the dataset text file	<pre>// Create a DataFrame of Station objects from the dataset text file. val StationsDF = sqlContext.read.format("csv").option("header","true").option("inferSchema","true").load("/user/cloudera/spark/Stations_2017.csv")</pre>
Register DataFrame as a table	<pre>// Register DataFrame as a table. dfStations.registerTempTable("Stations")</pre>
Display the content of DataFrame	<pre>// Display the content of DataFrame dfStations.show()</pre>
Print the DF schema	<pre>// Print the DF schema dfStations.printSchema()</pre>
Select Station name column	<pre>// Select Station name column dfStations.select("name").show()</pre>
Select Station code and name columns	<pre>// Select Station code and name columns dfStations.select("code", "name").show()</pre>
Select a Station by id	<pre>// Select a Station by id dfStations.filter(dfStations("code").equalTo(6080)).show()</pre>
Switch to SQL interpreter	<pre>%sql</pre>
Print all rows in the table	<pre>select * from Stations or spark.sql("Select * from Stations Limit 10").show</pre>
Print Stations count	<pre>select count(*) from Stations or spark.sql("Select count(*) from Stations").show</pre>

## Exercise 5: Rating movies

In this exercise we will analyze the ratings attributed by the audience to movies. The dataset has 4 columns:

Column 1: User ID  
Column 2: Movie ID  
Column 3: Rating  
Column 4: Timestamp

1. Copy the **u.data** file into **/sparklab** HDFS directory
2. Create a new Zeppelin note: **ratings**



Spark Interpreter	%pyspark	
RDD creation	<pre>import collections  lines = sc.textFile('/user/cloudera/spark/u.data')</pre>	
Extract the 3 <sup>rd</sup> column and count ratings (from 1 à 5)	<pre>ratings = lines.map(lambda x : x.split()[2]) res = ratings.countByValue()</pre>	
Sort the result and print it	<pre>sortedres = collections.OrderedDict(sorted(res.items()))</pre>	
	<pre>for key,value in sortedres.items():     print ("%s %i" %(key,value))</pre>	<pre>1 6110 2 11370 3 27145 4 34174 5 21201</pre>
Extracting the movie ID and creating an Pair-RDD that contains the movie ID and value = 1	<pre>movies = lines.map(lambda x : (int(x.split()[1]), 1)) movies.take(10)</pre>	
	<pre>[(242, 1), (302, 1), (377, 1), (51, 1), (346, 1), (474, 1), (265, 1), (465, 1), (451, 1), (86, 1)]</pre>	
Count movies	<pre>moviesCount = movies.reduceByKey(lambda x , y : x + y )  moviesCount.take(10)</pre>	
	<pre>[(2, 131), (4, 209), (6, 26), (8, 219), (10, 89), (12, 267), (14, 183), (16, 39), (18, 10), (20, 72)]</pre>	
To sort the RDD on 'value' we flip the RDD	<pre>flippedMovies = moviesCount.map(lambda x : (x[1], x[0])).sortByKey()  flippedMovies.take(10)</pre>	
Print the length of the final result	<pre>finalResult = sortedMovies.collect()  len(finalResult)</pre>	
	1682	



Transform the RDD into a data frame and register the dataframe as an in-memory table	<pre>resultDF = sc.parallelize(finalResult).toDF()  sqlContext.registerDataFrameAsTable(resultDF, "MoviesRating")</pre>
Switch to SQL interpreter and query the table	<pre>%sql  select * from MoviesRating order by _1 desc limit 15</pre>