

Interacting with Spark: simple examples

The following examples use Scala: for the Python examples, please refer to the website or the book indicated at the beginning of the document.

The interactive shell contains a built-in Spark Context (“sc”): we will use the Spark Context to load the file into a RDD, then we will do some RDD transformations or we will apply some actions.

Let’s start by loading the file “README.md” and count the number of lines and display the first 3 lines.

```
scala> val lines = sc.textFile("README.md")
lines: org.apache.spark.rdd.RDD[String] = README.md MapPartitionsRDD[1] at textFile
at <console>:24

scala> lines.count()
res0: Long = 104

scala> lines.take(3)
res2: Array[String] = Array(# Apache Spark, "", Spark is a fast and general cluster
computing system for Big Data. It provides)
```

Next, let’s count the words:

```
scala> val words = lines.flatMap(lines => lines.split(" "))
words: org.apache.spark.rdd.RDD[String] = MapPartitionsRDD[2] at flatMap at
<console>:26

scala> val wordskv = words.map(word => (word,1))
wordskv: org.apache.spark.rdd.RDD[(String, Int)] = MapPartitionsRDD[3] at map at
<console>:28

scala> val counts = wordskv.reduceByKey((a,b) => a + b)
counts: org.apache.spark.rdd.RDD[(String, Int)] = ShuffledRDD[4] at reduceByKey at
<console>:30

scala> counts.collect()
res3: Array[(String, Int)] = Array((package,1), (this,1),
(Version)](http://spark.apache.org/docs/latest/building-spark.html#specifying-the-
hadoop-version),1), (Because,1), (Python,2),
(page)](http://spark.apache.org/documentation.html),1), (cluster.,1), (its,1),
([run,1), (general,3), (have,1), (pre-built,1), (YARN,,1),
([http://spark.apache.org/developer-tools.html](the,1), (changed,1), (locally,2),
(sc.parallelize(1,1), (only,1), (locally.,1), (several,1), (This,2), (basic,1),
(Configuration,1), (learning,,1), (documentation,3), (first,1), (graph,1), (Hive,2),
(info,1), (["Specifying,1), ("yarn",1), ([params]`,1), ([project,1), (prefer,1),
(SparkPi,2), (<http://spark.apache.org/>,1), (engine,1), (version,1), (file,1),
(documentation,,1), (MASTER,1), (example,3), (["Parallel,1), (are...
scala>
scala> counts.saveAsTextFile("../data/wordcount")
```

The action `collect()` on a RDD shows the initial content of the RDD. We can save the result in a file with the following command

```
scala> counts.saveAsTextFile("../wordcount")
```

Exercise – Analysis of the tstat data

Consider the data and the exercises about tstat done with Pig. Re-do all the exercises with Spark.