

Spark Tutorial 2

Let's make some practice!

Download and launch the shell

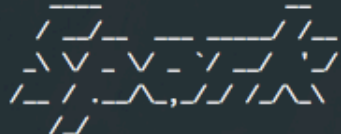
We'll be using spark 2.3.2

see <http://spark.apache.org/downloads.html>

1. Go to your browser and download the latest version
2. Unzip the archive
3. Open your terminal and go to the extracted directory
4. Go in the folder **./bin/**
5. Launch the script: **./pyspark**

PySpark

```
nadai@sibiu:/data/nadai/spark-2.3.0-bin-hadoop2.7$ ./bin/pyspark  
Python 3.6.5 (default, May 26 2018, 18:07:52)  
[GCC 7.3.0] on linux  
Type "help", "copyright", "credits" or "license" for more information.  
18/10/30 17:04:27 WARN Utils: Your hostname, sibiu resolves to a loopback address: 127.0.1.1; using 192.168.159.39 instead (on interface em1)  
18/10/30 17:04:27 WARN Utils: Set SPARK_LOCAL_IP if you need to bind to another address  
18/10/30 17:04:28 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable  
18/10/30 17:04:33 WARN Utils: Service 'SparkUI' could not bind on port 4040. Attempting port 4041.  
  
Welcome to
```



```
version 2.3.0  
  
Using Python version 3.6.5 (default, May 26 2018 18:07:52)  
SparkSession available as 'spark'.  
>>>
```

Favorites	Name	Date Modified	Size
[Folder Icon] marcodona.it	[File Icon] capacity.pdf	Today at 13:55	61 KB
[Folder Icon] marco			
[Folder Icon] Desktop			
[Folder Icon] PhD			
[Folder Icon] Documents			

Initializing Spark

Create a Spark Context (already provided in the notebook)

```
conf = SparkConf().setAppName(appName).setMaster(master)
sc = SparkContext(conf=conf)
```

appName = Name of the application to show on the cluster UI

master = Spark URL ('spark://ip-address:7077') or 'local'

Resilient Distributed Datasets [RDD]

Parallelized Collections

```
data = [1, 2, 3, 4, 5]  
distData = sc.parallelize(data)  
vs.  
sc.parallelize(data, 10)
```

```
distData.reduce(lambda a, b: a + b)
```

Resilient Distributed Datasets [RDD]

External Datasets

- HDFS
- Cassandra
- Hbase
- Text File
- Amazon S3
- Sequence File
- NO Database (SparkSQL does it!)

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

(with optional number of partitions)

`sc.textFile` -> reads a file and returns one record per line

`sc.wholeTextFiles` -> reads a folder as a pair <filename, content>

Resilient Distributed Datasets [RDD]

RDD Operations

- TRANSFORMATIONS create a new dataset from an existing one (e.g. map)
- ACTIONS return a value to the driver program after running a computation on the dataset (e.g. reduce)

```
lines = sc.textFile("data.txt")  
lineLengths = lines.map(lambda s: len(s))  
totalLength = lineLengths.reduce(lambda a, b: a + b)
```

Remember: LAZY COMPUTATION!

Solution: persist or cache RDD in memory (or disk)

If you want to use lineLengths again:

```
lineLengths.cache()
```

Resilient Distributed Datasets [RDD]

Passing Functions to Spark

```
lines = sc.textFile("data.txt")
```

```
def computeLength(s):  
    return len(s)
```

```
lineLengths = lines.map(computeLength)
```

```
def sumLengths(a, b):  
    return a + b
```

```
totalLength = lineLengths.reduce(sumLengths)
```


Resilient Distributed Datasets [RDD]

Understanding Closures

```
counter = 0
rdd = sc.parallelize(data)

# Wrong: Don't do this!!
def increment_counter(x):
    global counter
    counter += x

rdd.foreach(increment_counter)

print("Counter value: ", counter)
```

Remember: do NOT mutate the state in a loop or locally defined method

If you need to change the state -> use an ACCUMULATOR

In the same way, do not try to print an RDD with `rdd.foreach(println)` in a cluster!

Resilient Distributed Datasets [RDD]

Working with Key-Value Pairs

Key-Value pair is built-in in python -> (key, value)

Simply call the desired operation

```
lines = sc.textFile("data.txt")
```

```
pairs = lines.map(lambda s: (s, 1))
```

```
counts = pairs.reduceByKey(lambda a, b: a + b)
```

Resilient Distributed Datasets [RDD]

Transformations

- map
- filter
- flatMap
- mapPartitions
- mapPartitionsWithIndex
- sample
- union
- intersection
- distinct
- groupByKey
- reduceByKey
- aggregateByKey
- sortByKey
- join
- cogroup
- cartesian
- pipe
- coalesce
- repartition
- repartitionAndSortWithinPartitions

Resilient Distributed Datasets [RDD]

Actions

- reduce
- collect
- count
- first
- takeSample
- takeOrdered
- saveAsTextFile
- saveAsSequenceFile
- saveAsObjectFile
- countByKey
- foreach

Shared Variables

Broadcast Variables

- Keep a read-only variable cached on each machine rather than shipping a copy of it with tasks (e.g. give to every node a copy of a large dataset in an efficient manner)

```
broadcastVar = sc.broadcast([1, 2, 3])
```

```
broadcastVar.value()
```

Shared Variables

Accumulators

- Variables that are only “added” to through an associative and commutative operation (e.g. use for implementing counters)

```
accum = sc.accumulator(0)
```

```
sc.parallelize([1, 2, 3, 4]).foreach(lambda x: accum.add(x))
```

```
accum.value
```

Shared Variables

Custom Accumulator

```
class VectorAccumulatorParam(AccumulatorParam):  
    def zero(self, initialValue):  
        return Vector.zeros(initialValue.size)  
  
    def addInPlace(self, v1, v2):  
        v1 += v2  
        return v1
```

Then, create an Accumulator of this type:

```
vecAccum = sc.accumulator(Vector(...), VectorAccumulatorParam())
```

Deploying to a cluster

Use the CommandLine Tool

```
$SPARK_HOME/bin/spark-shell [options] \  
--master spark://ip-address:7077 \  
/path/to/script.py \  
[extra-options for the script]
```

Using Pyspark

```
$SPARK_HOME/bin/pyspark  
and then write your own code!
```

Using a Notebook

Code Repository

<https://github.com/forons/BigDataExamples/>

Today's exercise:

create an account on

<https://databricks.com/try-databricks>

My first application

Python

1. `Data = [1, 2, 3, 4, 5]`
2. `dataRDD = sc.parallelize(data)`
3. `def filterOperation(number):`
 `return number < 4`
4. `dataRDD.filter(filterOperation).collect()`

Simple WordCount application

Python

```
1. from operator import add
2. f = sc.textFile("README.md")
3. wc = f.flatMap(lambda x: x.split(' ')).map(lambda x: (x,1))
   .reduceByKey(add)
4. wc.saveAsTextFile("wc_out")
```

Simple WordCount application in MP (java)

```
1 public class WordCount {
2     public static class TokenizerMapper
3         extends Mapper<Object, Text, Text, IntWritable>{
4
5         private final static IntWritable one = new IntWritable(1);
6         private Text word = new Text();
7
8         public void map(Object key, Text value, Context context
9             ) throws IOException, InterruptedException {
10             StringTokenizer itr = new StringTokenizer(value.toString());
11             while (itr.hasMoreTokens()) {
12                 word.set(itr.nextToken());
13                 context.write(word, one);
14             }
15         }
16     }
17
18     public static class IntSumReducer
19         extends Reducer<Text, IntWritable, Text, IntWritable> {
20         private IntWritable result = new IntWritable();
21
22         public void reduce(Text key, Iterable<IntWritable> values,
23             Context context
24             ) throws IOException, InterruptedException {
25             int sum = 0;
26             for (IntWritable val : values) {
27                 sum += val.get();
28             }
29             result.set(sum);
30             context.write(key, result);
31         }
32     }
33
34     public static void main(String[] args) throws Exception {
35         Configuration conf = new Configuration();
36         String[] otherArgs = new GenericOptionsParser(conf, args).getRemainingArgs();
37         if (otherArgs.length < 2) {
38             System.err.println("Usage: wordcount <in> [<in>...] <out>");
39             System.exit(2);
40         }
41         Job job = new Job(conf, "word count");
42         job.setJarByClass(WordCount.class);
43         job.setMapperClass(TokenizerMapper.class);
44         job.setCombinerClass(IntSumReducer.class);
45         job.setReducerClass(IntSumReducer.class);
46         job.setOutputKeyClass(Text.class);
47         job.setOutputValueClass(IntWritable.class);
48         for (int i = 0; i < otherArgs.length - 1; ++i) {
49             FileInputFormat.addInputPath(job, new Path(otherArgs[i]));
50         }
51         FileOutputFormat.setOutputPath(job,
52             new Path(otherArgs[otherArgs.length - 1]));
53         System.exit(job.waitForCompletion(true) ? 0 : 1);
54     }
55 }
```

```
1 val f = sc.textFile(inputPath)
2 val w = f.flatMap(l => l.split(" ")).map(word => (word, 1)).cache()
3 w.reduceByKey(_ + _).saveAsText(outputPath)
```

WordCount in 3 lines of Spark

We can start

- Find all words that are palindrome

{ rotator } { gig } { pop } { level } { civic }

- Find the words that occur exactly 5 times
- Group word by their number of occurrences
- Group anagrams

{ mean, name } { flow, wolf }

Contacts

For any problem, send a mail to:

daniele.foroni@unitn.it