

Introduction to Spark

Master M2 – Université Grenoble Alpes & Grenoble INP

2018

In this lab, we are going to write our first Spark programs. To this end, we are going to conduct some analysis on a dataset of significant size. The programs are going to be run on a single node. Note however that the same code would be used to run on multiple nodes.

1 About the lab

This lab is not graded. However, you are strongly suggested to work on this lab by yourself to improve your knowledge on Spark. The computations proposed in the lab are some examples of simple things that can be done with Spark. You are strongly encouraged try running other analysis on the provided dataset to continue practicing.

Do not hesitate to ask questions during the lab session. If you have questions after the session, you can always send them by email to Francieli Zanon Boito (francieli.zanon-boito@inria.fr) and Thomas Ropars (thomas.ropars@univ-grenoble-alpes.fr)¹.

2 The dataset

The dataset we will study in this lab comes from the project *Climatological Database for the World's Oceans*. Detailed information about this project can be found at <http://webs.ucm.es/info/cliwoc/>.

In a few words, this dataset has been created starting from the logbook of ships that were navigating the oceans during the 18th and 19th century. These logbooks were maintained by the crew of the ships and contain a lot of information regarding the weather conditions during that period.

Each observation reported in the dataset includes many entries including the date of the observation, the location, the air temperature, the wind speed, etc. A detailed description of all included fields is available here: <http://webs.ucm.es/info/cliwoc/content/CLIWOC15all.htm>

The dataset is to be downloaded at the following address: <https://filesender.renater.fr/?s=download&token=8f4df1cb-678a-5640-97f4-1fa6957a56ad>

¹Sending your question to both addresses increases the probability of getting a rapid answer.

3 Work on the dataset

We are going to use Spark to analyze the data included in this dataset. Instructions on how to install and use Spark are provided in Appendix A.

3.1 Programming language

For this lab, you can either use Scala or Python as programming language. Using Scala has the advantage that in general performance will be better since Spark executes inside a JVM.

Note that to use Scala, you will have to use the machines from the lab room. The installation based on Docker we propose for your own machine does not have support to compile Scala code.

3.2 Provided code

The provided project contains an initial Spark code in Python (file `code/navigation.py`) and in Scala (file `code/Navigation/src/main/scala/navigation.scala`). This code already works and does the following

- It creates a RDD out of the input dataset
- It displays 5 of the nationalities associated with reports in the dataset (column "Nationality")
- Note that the provided code assumes that the dataset has been stored in the directory `code/data`.

Start by running this code and try to understand how it works.

A complete documentation of the Spark API for manipulating RDDs is available online:

- For Scala: <https://spark.apache.org/docs/latest/api/scala/index.html#org.apache.spark.api.java.JavaRDD>
- For Python: <https://spark.apache.org/docs/latest/api/python/pyspark.html#pyspark.RDD>

3.3 A few comments

Before starting doing the exercises, here are a few comments to have in mind while you are programming with Spark:

- While you are debugging a program, it can be good to run with a single worker (`local[1]`) to avoid very large and messy log traces in case of error.
- Once your program works, you can try executing it with more worker to observe the impact on performance. (Elapsed time can be measured in the way described here in Python: <https://stackoverflow.com/a/25823885>, and here in Scala: <https://stackoverflow.com/a/37731494>)

- Some of the questions below can be solved in different ways. Do not hesitate to implement multiple solutions and compare their performance.
- Caching RDDs can have a significant impact on the performance of Spark (see <https://spark.apache.org/docs/latest/rdd-programming-guide.html#rdd-operations>). You can also try to evaluate this impact during your tests.
- Note that in the dataset, an entry with no value is represented by the string "NA" (stands for *Not Available*)

3.4 Exercises

Implement the Spark program that does the following:

1. When running the provided example code, you will observe that there might be several entries that are equivalent. More specifically, you will see two entries for "British" with the only difference that one has an extra white space in the name. Propose a new version of the computation that will consider these two entries as the same one.
2. Count the total number of observations included in the dataset (each line corresponds to one observation)
3. Count the number of years over which observations have been made (Column "Year" should be used)
4. Display the oldest and the newest year of observation
5. Display the years with the minimum and the maximum number of observations (and the corresponding number of observations)
6. Count the distinct departure places (column "VoyageFrom") using two methods (i.e., using the function *distinct()* or *reduceByKey()*) and compare the execution time.
7. Display the 10 most popular departure places
8. Display the 10 roads (defined by a pair "VoyageFrom" and "VoyageTo") the most often taken.
 - Here you can start by implementing a version where a pair "VoyageFrom"- "VoyageTo" A-B and a pair B-A correspond to different roads.
 - Implement then a second version where A-B and B-A are considered as the same road.
9. Compute the hottest month (defined by column "Month") on average over the years considering all temperatures (column "ProbTair") reported in the dataset.

As already mentioned, do not hesitate to run other analysis on the dataset to practice with Spark.

3.5 To go further: Dataframes

Dataframes have been introduced in Spark to make it simpler to manipulate structured data organized into named columns. This is the case for the data that are manipulated during this lab.

Dataframes are not going to be studied during this course. However, if you are curious about it, you can try solving the questions of this lab using Dataframes.

Here is a starting point for a description of Dataframes: <https://spark.apache.org/docs/latest/sql-programming-guide.html#datasets-and-dataframes>

A Installing and using Apache Spark

Instructions about installing and using Spark. Note that the installation based on Docker only allows running Spark applications coded in Python. If you want to use Scala, you have to follow the instructions for using the machines of the lab rooms.

A.1 Using your own machine

A.1.1 Installation

You already received the instructions to install Spark using Docker on your machine. If it is not done yet, the instructions are available at <https://tropars.github.io/downloads/lectures/LSDM/LSDM-Spark-on-your-Laptop.pdf>.

A.1.2 Using Spark in Python

We describe here how to use Spark after you have installed it following the instructions described above:

1. Start the docker container by running in a terminal:

```
docker run -v absolute_path_to_folder:/home/jovyan/work -it \
--rm -p 8888:8888 -p 4040:4040 jupyter/pyspark-notebook
```

- In the previous command, `absolute_path_to_folder` should be replaced by the actual path to the directory where you are going to store your Python scripts
 - When you launch this command, logs are written into the terminal. The last displayed line is an url that you should copy into a web browser. This url allows you to connect to a Jupyter Notebook that gives access to Spark.
2. In Jupyter Notebook, open a new terminal (New > Terminal)
 3. In the new terminal, move into the directory storing your Python scripts (for this simply run: `cd work`)
 4. Assuming that you have a script called `navigation.py`, you can run it by simply executing in the Notebook terminal:

```
python3 ./navigation.py
```

A.2 Using the machines from the lab rooms

A.2.1 Installation

Here are the instructions to install and configure Spark in the lab rooms:

1. Download the latest already compiled version of Spark here: <https://www.apache.org/dyn/closer.lua/spark/spark-2.3.2/spark-2.3.2-bin-hadoop2.7.tgz>
2. Extract the downloaded archive:

```
tar zxvf spark-2.3.2-bin-hadoop2.7.tgz
```

3. Configure the require environment file by updating the file `$HOME/.bashrc` by adding the following lines at the beginning of the file²:

```
export SPARK_HOME=PATH_TO_DIR/spark-2.3.2-bin-hadoop2.7
```

```
export PYTHONPATH="${SPARK_HOME}/python/:$PYTHONPATH"
```

```
export PYTHONPATH="${SPARK_HOME}/python/lib/py4j-0.10.7-src.zip:$PYTHONPATH"
```

```
export PATH=${SPARK_HOME}/bin:$PATH
```

- where `PATH_TO_DIR` correspond to the directory where your stored Spark.
4. Start a new terminal to make your changes active
 5. In the new terminal, launch `pyspark` to check that everything works correctly

A.2.2 Using Spark in Python

We describe here how to use Spark after you have installed it following the instructions described above:

1. Assuming that you have a script called `navigation.py`, you can run it by simply executing in a terminal:

```
python3 ./navigation.py
```

A.2.3 Using Spark in Scala

To use Spark in Scala

1. Insert your code in the file `Navigation/src/main/scala/navigation.scala` in the provided project tree.
2. Compile your project by running the command `sbt package` in the directory `Navigation`
3. To run your program, use the following command:

```
spark-submit target/scala-2.11/navigation_2.11-1.0.jar
```

²To open this file, simply run `nano ~/.bashrc` in a terminal