POLITENICO DI MILANO

MIDDLEWARE TECHNOLOGIES FOR DISTRIBUTED SYSTEMS

PROJECT REPORT

# Analysis of COVID-19 Data

*Author:*
Fabio LOSAVIO
Matteo VELATI
Niccolò ZANGRANDO

*Supervisor:*
Dr. Alessandro MARGARA

April 11, 2021

# 1 Introduction

## 1.1 Project description

We have implemented a program that analyzes open datasets to study the evolution of the COVID-19 situation worldwide. The program starts from the dataset of new reported cases for each country daily and computes the following queries:

1. Seven days moving average of new reported cases, for each country and for each day

2. Percentage increase (with respect to the day before) of the seven days moving average, for each country and for each day

3. Top 10 countries with the highest percentage increase of the seven days moving average, for each day

## 1.2 Technology used

We have decided to use the *Apache Spark* technology. In particular, given the nature of the project, we decided to use the *Spark SQL* engine.

In the following sections, we are going to describe how we have exploited the three different queries, and then we will show and comment the different performance results we have obtained with different configurations.

# 2 Implementation

As mentioned above, given the nature of the problem and the format in .csv of the dataset [1], the solution has been implemented using the *SQL engine of Apache Spark*.

## 2.1 Query 1: Seven days moving average of new reported cases, for each country and for each day

The moving average can be defined as:

$$movingAverage = \frac{1}{n} \sum_{t=-n}^{1} a_t \tag{1}$$

where *a=value at period t*, and *n=number of time periods*.

Since we have a dataset of new reported cases for each country daily, we can split the computation of the seven days moving average for each country simply using the SQL *"partition by" Window function*, order the rows by date and then calculate the average between each row and its six previous. The query will be:

```
final Dataset<Row> dateCountry_ma = covidData
    .select(col("dateRep"), col("countriesAndTerritories"),
        avg("cases").over(Window
            .partitionBy("countriesAndTerritories")
        .orderBy("dateRep").rowsBetween(-6, 0))
            .as("7days_moving_average"))
    .orderBy("countriesAndTerritories", "dateRep");
```
Listing 1: *Query 1 solution*

where *covidData* is the dataframe created from the .csv dataset.

---

[1]https://opendata.ecdc.europa.eu/covid19/casedistribution/csv/data.csv

## 2.2 Query 2: Percentage increase (with respect to the day before) of the seven days moving average, for each country and for each day

The percentage increase (with respect to the day before) of the seven days moving average can be calculated as:

$$increase = 100 * \frac{7ma_x}{7ma_{x-1}} - 100 \tag{2}$$

where

$7ma_x$=is the seven days moving average at day x, and $7ma_{x-1}$=is the seven days moving average at day x-1.

The approach is very similar to the previous one. Using the first query solution [Listing 1], we can split the computation for each country using the the SQL *"partition by" Window function*, order the rows by date and then calculate the ratio between each row and its previous one.
The query will be:

```
final Dataset<Row> dateCountry_increase = date_Country_ma
    .select(col("dateRep"), col("countriesAndTerritories"),
        col("7days_moving_average").divide(lag(
        col("7_days_moving_average"), 1)
            .over(Window
                .partitionBy("countriesAndTerritories")
            .orderBy("dateRep"))).multiply(100)
                .minus(100).as("increase_day_before"))
    .orderBy("countriesAndTerritories", "dateRep")
    .na().fill(0);
```

Listing 2: *Query 2 solution*

## 2.3 Query 3: Top 10 countries with the highest percentage increase of the seven days moving average, for each day

In this case, starting from the query 2 solution [Listing 2] we can use the SQL *"partition by" Window function* for partitioning the dataset by date. Then, for each day, we can order the countries by the percentage increase

(from the biggest value) and keep only the first 10 rows.
The query will be:

```
final Dataset<Row> topCountry_increase =
    date_Country_increase
    .select(col("dateRep"), col("countriesAndTerritories"),
        col("increase_day_before"), row_number().over(Window
            .partitionBy("dateRep").orderBy(
                desc("increase_day_before")))
                .alias("rank"))
    .filter(col("rank").leq(10))
    .orderBy(desc("dateRep"), desc("increase_day_before"));
```
Listing 3: Query 3 solution

# 3   Analysis and Testing

In this section we are going to describe the analysis and tests we have done with different configurations.

## 3.1   Hardware

For all the tests we have used 3 machines connected to the same local network, with 1 Master and 3 Workers. The specifications of each machine are as follows:

- Machine A:

    - CPU: Intel Core i5-8259U @ 2.30GHz
    - RAM: 8GB

- Machine B:

    - CPU: Intel Core i5-8259U @ 2.30GHz
    - RAM: 8GB

- Machine C:

    - CPU: Intel Core i7-8550U @ 1.80GHz
    - RAM: 8GB

## 3.2 Tests

### 3.2.1 Test 1

In this test we have submitted the program to the spark cluster with the default parameters using all the 24 available cores (8 for each machine). Here are shown the results obtained:
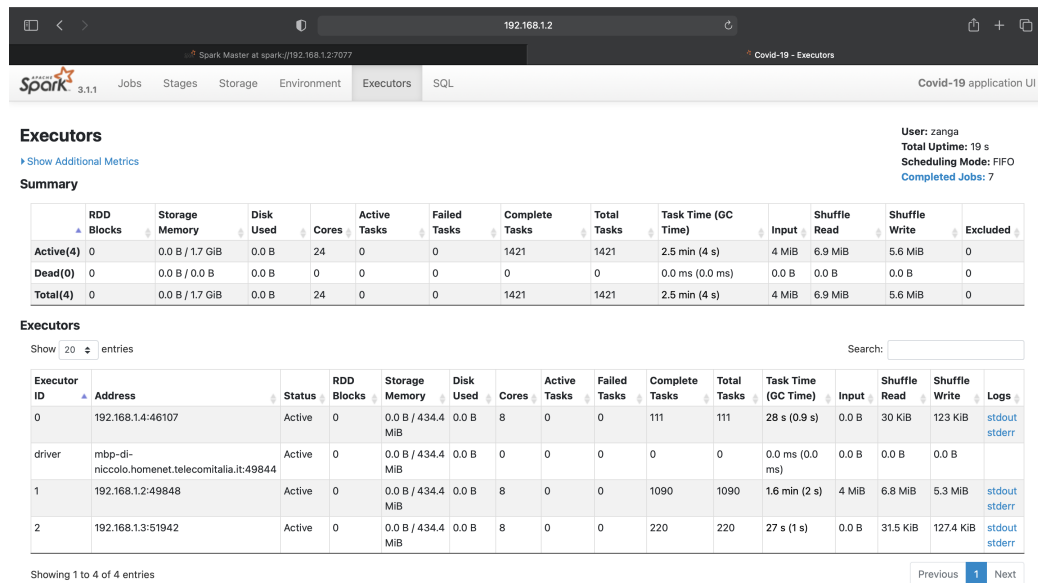


Figure 1: test 1 results

As you can see, 1421 tasks have been executed in 19 seconds and the workload has not been evenly divided among the Executors.

### 3.2.2 Test 2

In this test we have tried to increase the performace. We know that the number of tasks depends also on the number partitions. By setting the configuration parameter *"spark.sql.shuffle.partitions=3"* we have obtained the results shown in Figure 2.
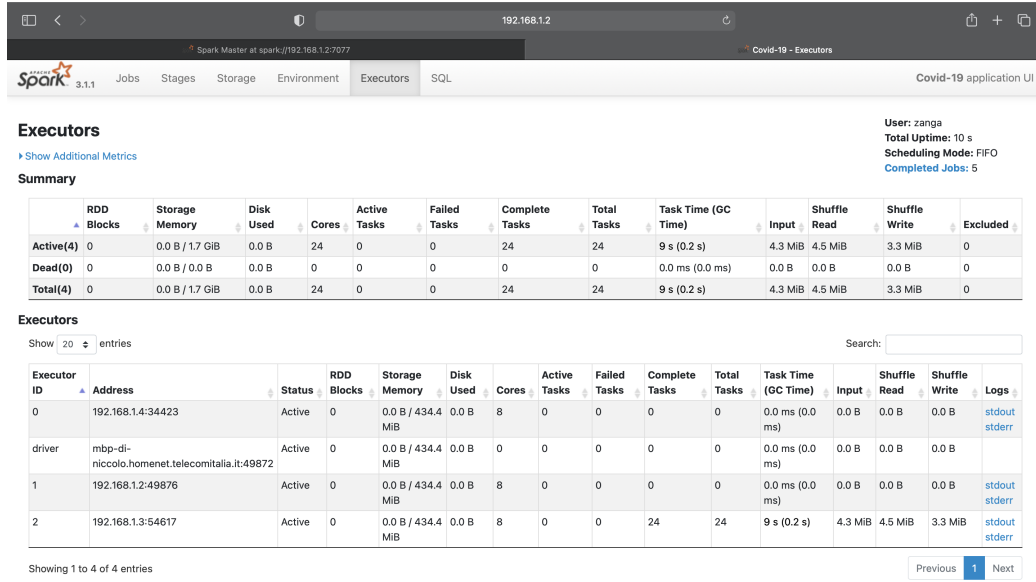


Figure 2: test 2 results

In this test 24 tasks have been executed in 10 seconds. We observed that all the tasks have been executed by one Executor and the execution time has been halved. This behavior is justified by the Apache Spark architecture, in fact since the computation workload is very small (the dataset is only 4MB), Apache Spark try not to split the work among different nodes.

## 3.3 Test 3

Since the original dataset is only 4MB, we have increased the dataset size by simply duplicating many times the original entries. Now the dataset size is about 130MB. By setting *"spark.sql.shuffle.partitions=3"* we have obtained the results shown in Figure 3.
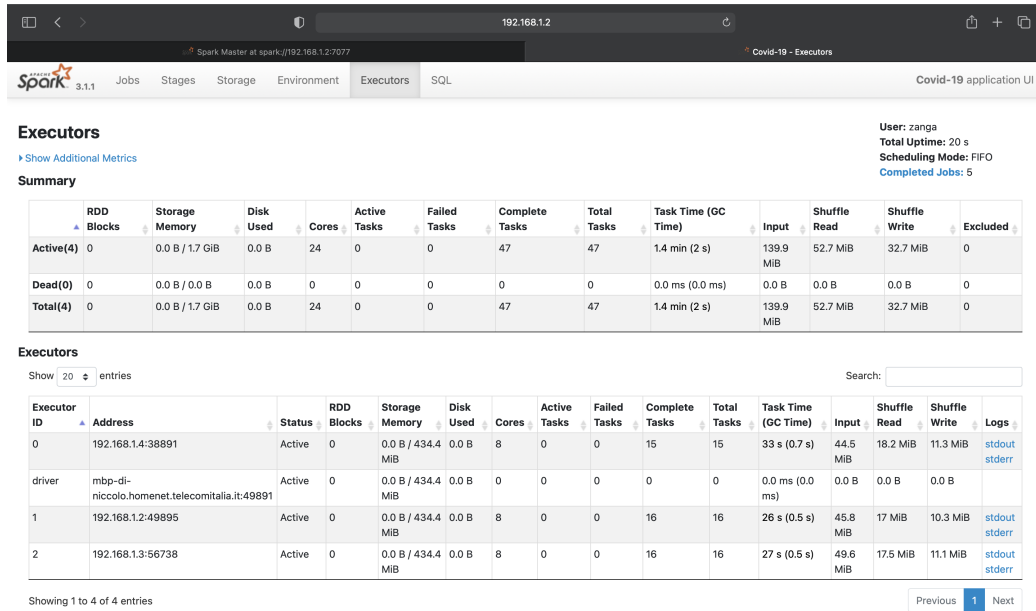


Figure 3: test 3 results

In this case the 47 tasks have been equally divided among all the Executors and the computation finished in 20 seconds.

## 3.4   Test 4

For this test we have used the "big" dataset with a different configuration. In fact we have used only 2 cores per Worker, but we have set *"spark.sql.shuffle.partitions=3"* as before. In [Fig. 4] we can see that has been processed 29 tasks in 19 seconds. The execution time is more or less the same despite lower computing power (6 cores). On the other side the execution has been divided in less tasks than *Test 3* [Fig.3].
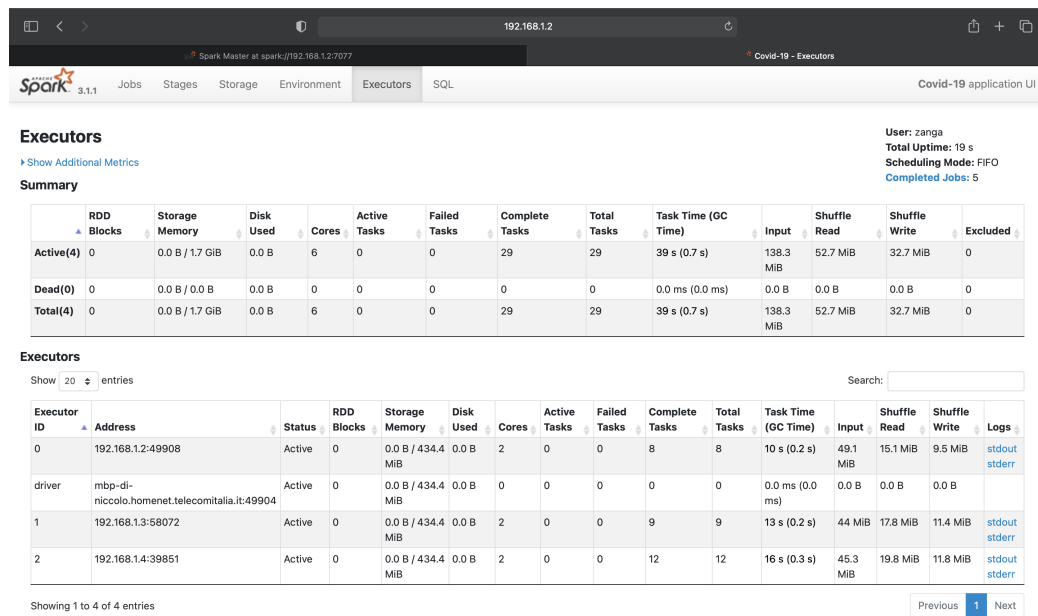


Figure 4: test 4 results

## 3.5   Test 5

Since in *Test 4* (Fig.4) we have reached better performances than *Test 3* (Fig.3), we have changed again the configuration. We set *"spark.sql.shuffle.partitions=24"* and we used all 24 available cores.
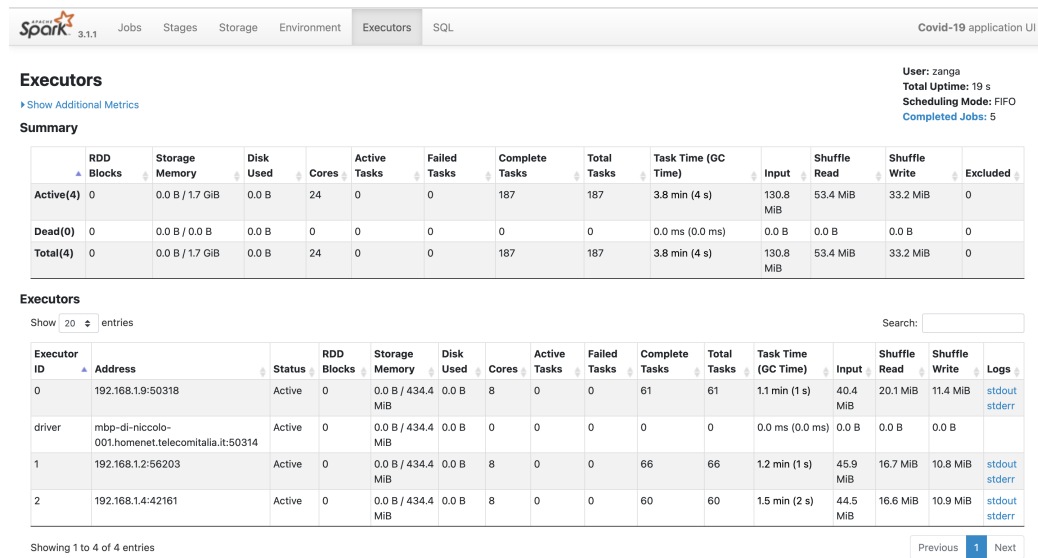


Figure 5: test 5 results

# 4    Conclusion

Is very important to find the correct configuration, based on the submitted application, of the Apache Spark cluster. For example, for small computations like in *Test 1* (Fig.1) and *Test 2* (Fig.2), is useless and lowers the performance creating a huge quantity of tasks. On the other side also creating too few tasks can be useless because, in this case, we would not be making full use of parallelization (Fig. 3). Of course this technology improves performance for big computations, but if we run the application in local mode we can see the benefits. In fact, if we use 1 core with *"spark.sql.shuffle.partitions=1"* we have the solutions in 7 seconds; with 4 cores, we can reach 5,8 seconds.

In the figure below are shown other tests in cluster mode with different drivers (For the applications under the red line we have used the "big" dataset, for the ones above the red line we have used the original dataset).
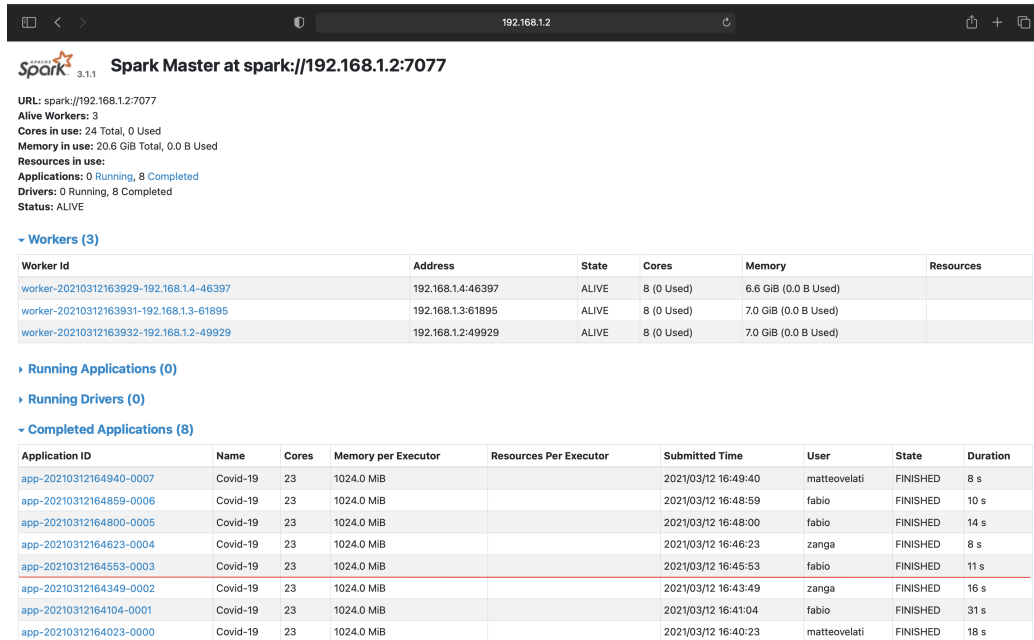


Figure 6: tests results