

Big data: architectures and data analytics

Spark - Exercises

Exercise #30

- Log filtering
 - Input: a simplified log of a web server (i.e., a textual file)
 - Each line of the file is associated with a URL request
 - Output: the lines containing the word "google"
 - Store the output in an HDFS folder

3

Exercise #30 - Example

■ Input file

```
66.249.69.97 - - [24/Sep/2014:22:25:44 +0000] "GET http://www.google.com/bot.html"
66.249.69.97 - - [24/Sep/2014:22:26:44 +0000] "GET http://www.google.com/how.html"
66.249.69.97 - - [24/Sep/2014:22:28:44 +0000] "GET http://dbdmg.polito.it/course.html"
71.19.157.179 - - [24/Sep/2014:22:30:12 +0000] "GET http://www.google.com/faq.html"
66.249.69.97 - - [24/Sep/2014:31:28:44 +0000] "GET http://dbdmg.polito.it/thesis.html"
```

■ Output

```
66.249.69.97 - - [24/Sep/2014:22:25:44 +0000] "GET http://www.google.com/bot.html"
66.249.69.97 - - [24/Sep/2014:22:26:44 +0000] "GET http://www.google.com/how.html"
71.19.157.179 - - [24/Sep/2014:22:30:12 +0000] "GET http://www.google.com/faq.html"
```

4

Exercise #31

- Log analysis
 - Input: log of a web server (i.e., a textual file)
 - Each line of the file is associated with a URL request
 - Output: the list of distinct IP addresses associated with the connections to a google page (i.e., connections to URLs containing the term "www.google.com")
 - Store the output in an HDFS folder

5

Exercise #31 - Example

■ Input file

```
66.249.69.97 - - [24/Sep/2014:22:25:44 +0000] "GET http://www.google.com/bot.html"
66.249.69.97 - - [24/Sep/2014:22:26:44 +0000] "GET http://www.google.com/how.html"
66.249.69.97 - - [24/Sep/2014:22:28:44 +0000] "GET http://dbdmg.polito.it/course.html"
71.19.157.179 - - [24/Sep/2014:22:30:12 +0000] "GET http://www.google.com/faq.html"
66.249.69.95 - - [24/Sep/2014:31:28:44 +0000] "GET http://dbdmg.polito.it/thesis.html"
66.249.69.97 - - [24/Sep/2014:56:26:44 +0000] "GET http://www.google.com/how.html"
56.249.69.97 - - [24/Sep/2014:56:26:44 +0000] "GET http://www.google.com/how.html"
```

■ Output

```
66.249.69.97
71.19.157.179
56.249.69.97
```

6

Exercise #32

- Maximum value
 - Input: a collection of (structured) textual csv files containing the daily value of PM10 for a set of sensors
 - Each line of the files has the following format
sensorId,date,PM10 value ($\mu\text{g}/\text{m}^3$)\n
 - Output: report the maximum value of PM10
 - Print the result on the standard output

7

Exercise #32 - Example

- Input file

```
S1,2016-01-01,20.5
S2,2016-01-01,30.1
S1,2016-01-02,60.2
S2,2016-01-02,20.4
S1,2016-01-03,55.5
S2,2016-01-03,52.5
```

- Output

60.2

8

Exercise #33

- Top-k maximum values
 - Input: a collection of (structured) textual csv files containing the daily value of PM10 for a set of sensors
 - Each line of the files has the following format
sensorId,date,PM10 value ($\mu\text{g}/\text{m}^3$)\n
 - Output: report the top-3 maximum values of PM10
 - Print the result on the standard output

9

Exercise #33 - Example

- Input file

```
S1,2016-01-01,20.5
S2,2016-01-01,30.1
S1,2016-01-02,60.2
S2,2016-01-02,20.4
S1,2016-01-03,55.5
S2,2016-01-03,52.5
```

- Output

```
60.2
55.5
52.5
```

10

Exercise #34

- Readings associated with the maximum value
 - Input: a collection of (structured) textual csv files containing the daily value of PM10 for a set of sensors
 - Each line of the files has the following format
sensorId,date,PM10 value ($\mu\text{g}/\text{m}^3$)\n
 - Output: the line(s) associated with the maximum value of PM10
 - Store the result in an HDFS folder

11

Exercise #34 - Example

- Input file

```
S1,2016-01-01,20.5
S2,2016-01-01,30.1
S1,2016-01-02,60.2
S2,2016-01-02,20.4
S1,2016-01-03,60.2
S2,2016-01-03,52.5
```

- Output

```
S1,2016-01-02,60.2
S1,2016-01-03,60.2
```

12

Exercise #35

- Dates associated with the maximum value
 - Input: a collection of (structured) textual csv files containing the daily value of PM10 for a set of sensors
 - Each line of the files has the following format
sensorId,date,PM10 value ($\mu\text{g}/\text{m}^3$)\n
 - Output: the date(s) associated with the maximum value of PM10
 - Store the result in an HDFS folder

33

Exercise #35 - Example

Input file

```
S1,2016-01-01,20.5
S2,2016-01-01,30.1
S1,2016-01-02,60.2
S2,2016-01-02,20.4
S1,2016-01-03,60.2
S2,2016-01-03,52.5
```

Output

```
2016-01-02
2016-01-03
```

34

Exercise #36

- Average value
 - Input: a collection of (structured) textual csv files containing the daily value of PM10 for a set of sensors
 - Each line of the files has the following format
sensorId,date,PM10 value ($\mu\text{g}/\text{m}^3$)\n
 - Output: compute the average PM10 value
 - Print the result on the standard output

35

Exercise #36 - Example

Input file

```
S1,2016-01-01,20.5
S2,2016-01-01,30.1
S1,2016-01-02,60.2
S2,2016-01-02,20.4
S1,2016-01-03,55.5
S2,2016-01-03,52.5
```

Output

39.86

36

Big data: architectures and data analytics

Exercise #31

- Log analysis
 - Input: log of a web server (i.e., a textual file)
 - Each line of the file is associated with a URL request
 - Output: the list of distinct IP addresses associated with the connections to a google page (i.e., connections to URLs containing the term "www.google.com")
 - Store the output in an HDFS folder

2

Exercise #31 - Example

Input file

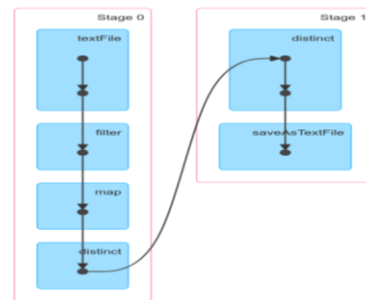
```
66.249.69.97 - [24/Sep/2014:22:25:44 +0000] "GET http://www.google.com/bot.html"
66.249.69.97 - [24/Sep/2014:22:26:44 +0000] "GET http://www.google.com/how.html"
66.249.69.97 - [24/Sep/2014:22:28:44 +0000] "GET http://dbdmg.polito.it/course.html"
71.19.157.179 - [24/Sep/2014:22:30:12 +0000] "GET http://www.google.com/faq.html"
66.249.69.95 - [24/Sep/2014:31:28:44 +0000] "GET http://dbdmg.polito.it/thesis.html"
66.249.69.97 - [24/Sep/2014:56:26:44 +0000] "GET http://www.google.com/how.html"
56.249.69.97 - [24/Sep/2014:56:26:44 +0000] "GET http://www.google.com/how.html"
```

Output

```
66.249.69.97
71.19.157.179
56.249.69.97
```

3

Exercise #31 - DAG



4

Exercise #31 - Simulation

- Suppose that Sparks splits the RDD associated with the input file in two partitions

Part. #1

```
66.249.69.97 - [24/Sep/2014:22:25:44 +0000] "GET http://www.google.com/bot.html"
66.249.69.97 - [24/Sep/2014:22:26:44 +0000] "GET http://www.google.com/how.html"
66.249.69.97 - [24/Sep/2014:22:28:44 +0000] "GET http://dbdmg.polito.it/course.html"
```

Part. #2

```
71.19.157.179 - [24/Sep/2014:22:30:12 +0000] "GET http://www.google.com/faq.html"
66.249.69.95 - [24/Sep/2014:31:28:44 +0000] "GET http://dbdmg.polito.it/thesis.html"
66.249.69.97 - [24/Sep/2014:56:26:44 +0000] "GET http://www.google.com/how.html"
56.249.69.97 - [24/Sep/2014:56:26:44 +0000] "GET http://www.google.com/how.html"
```

5

Exercise #31 - Simulation

Part. #1

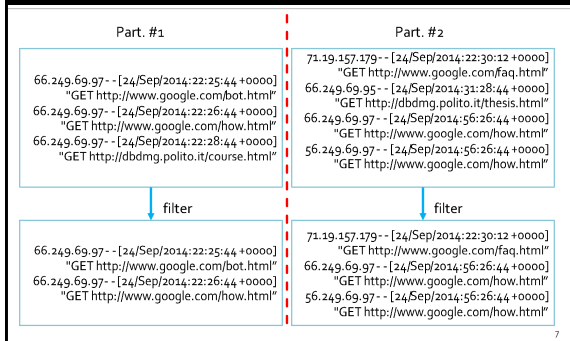
```
66.249.69.97 - [24/Sep/2014:22:25:44 +0000] "GET http://www.google.com/bot.html"
66.249.69.97 - [24/Sep/2014:22:26:44 +0000] "GET http://www.google.com/how.html"
66.249.69.97 - [24/Sep/2014:22:28:44 +0000] "GET http://dbdmg.polito.it/course.html"
66.249.69.97 - [24/Sep/2014:22:28:44 +0000] "GET http://dbdmg.polito.it/course.html"
```

Part. #2

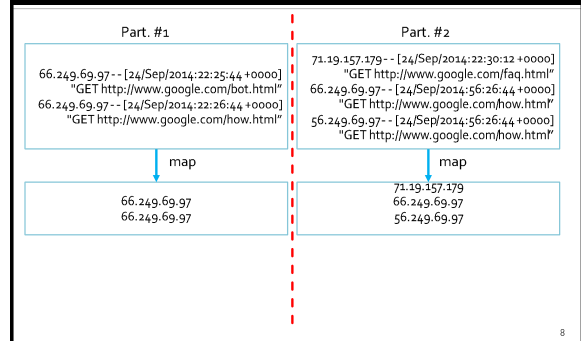
```
71.19.157.179 - [24/Sep/2014:22:30:12 +0000] "GET http://www.google.com/faq.html"
66.249.69.95 - [24/Sep/2014:31:28:44 +0000] "GET http://dbdmg.polito.it/thesis.html"
66.249.69.97 - [24/Sep/2014:56:26:44 +0000] "GET http://www.google.com/how.html"
56.249.69.97 - [24/Sep/2014:56:26:44 +0000] "GET http://www.google.com/how.html"
```

6

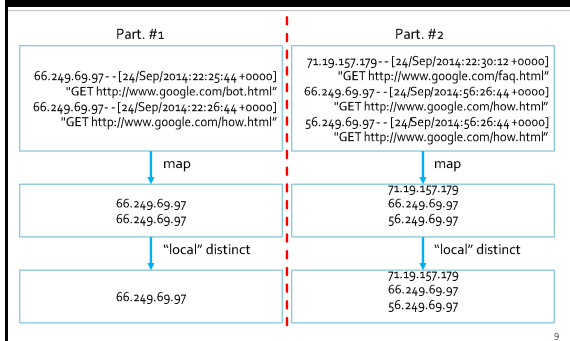
Exercise #31 - Simulation



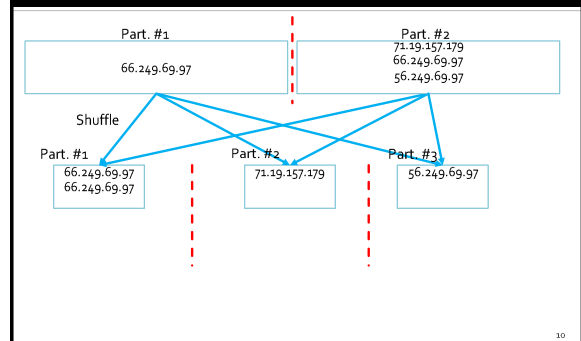
Exercise #31 - Simulation



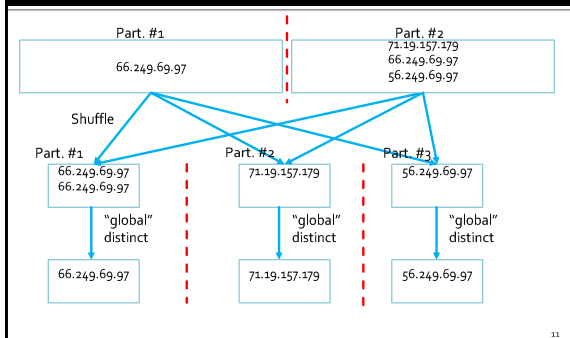
Exercise #31 - Simulation



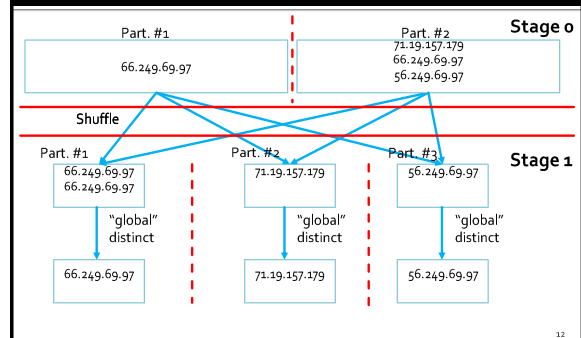
Exercise #31 - Simulation



Exercise #31 - Simulation



Exercise #31 - Simulation



Big data: architectures and data analytics

Spark - Exercises

Exercise #37

- Maximum values
 - Input: a textual csv file containing the daily value of PM₁₀ for a set of sensors
 - Each line of the files has the following format
sensorId,date,PM₁₀ value (µg/m³)\n
 - Output: the maximum value of PM₁₀ for each sensor
 - Store the result in an HDFS file

3

Exercise #37 - Example

■ Input file

```
S1,2016-01-01,20.5
S2,2016-01-01,30.1
S1,2016-01-02,60.2
S2,2016-01-02,20.4
S1,2016-01-03,55.5
S2,2016-01-03,52.5
```

■ Output

```
(S1,60.2)
(S2,52.5)
```

4

Exercise #38

- Pollution analysis
 - Input: a textual csv file containing the daily value of PM₁₀ for a set of sensors
 - Each line of the files has the following format
sensorId,date,PM₁₀ value (µg/m³)\n
 - Output: the sensors with at least 2 readings with a PM₁₀ value greater than the critical threshold 50
 - Store in an HDFS file the sensorIds of the selected sensors and also the number of times each of those sensors is associated with a PM₁₀ value greater than 50

5

Exercise #38 - Example

■ Input file

```
S1,2016-01-01,20.5
S2,2016-01-01,30.1
S1,2016-01-02,60.2
S2,2016-01-02,20.4
S1,2016-01-03,55.5
S2,2016-01-03,52.5
```

■ Output

```
(S1,2)
```

6

Exercise #39

- Critical dates analysis
 - Input: a textual csv file containing the daily value of PM₁₀ for a set of sensors
 - Each line of the files has the following format
sensorId,date,PM₁₀ value (µg/m³)\n
 - Output: an HDFS file containing one line for each sensor
 - Each line contains a sensorId and the list of dates with a PM₁₀ values greater than 50 for that sensor

7

Exercise #39 - Example

- Input file

```
s1,2016-01-01,20.5
s2,2016-01-01,30.1
s1,2016-01-02,60.2
s2,2016-01-02,20.4
s1,2016-01-03,55.5
s2,2016-01-03,52.5
```

- Output

```
(s1,[2016-01-02,2016-01-03])
(s2,[2016-01-03])
```

8