

# Report Big Data

ANTONIO SESSA

January 2024

Email:a.sessa108@studenti.unisa.it

Matricola: 0622702305

Professore: Giuseppe D'Aniello gidaniello@unisa.it gidaniello@unisa.it

## Indice

<b>1 Hadoop</b>	<b>1</b>
1.1 Descrizione del problema (traccia) . . . . .	1
1.2 Soluzione proposta Hadoop . . . . .	2
1.3 Codice Sorgente Mapper . . . . .	3
1.4 Codice Sorgente Reducer . . . . .	4
1.5 Codice Sorgente Driver . . . . .	5
<b>2 Spark</b>	<b>6</b>
2.1 Descrizione del problema (statistica a scelta) . . . . .	6
2.2 Soluzione proposta Spark . . . . .	6
2.3 Codice Sorgente Spark . . . . .	7
<b>3 Risultati</b>	<b>10</b>
3.1 Risultato Hadoop con input utente ITA . . . . .	10
3.2 Risultati Spark con Input X = 15 e Y = 1 . . . . .	11
3.3 Risultati Spark con Input X = 15 e Y = 4 . . . . .	11
<b>4 Come eseguire il codice</b>	<b>11</b>
4.1 Eseguire Hadoop . . . . .	11
4.2 Eseguire Spark . . . . .	12

## 1 Hadoop

### 1.1 Descrizione del problema (traccia)

Realizzare l'indice degli arbitri, di una certa nazionalità, che consenta di sapere ogni arbitro quali partite ha arbitrato, considerando che un arbitro può ricoprire il ruolo di Referee, Assistant 1 e Assistant 2 in partite diverse. La nazionalità deve essere indicata dall'utente (ad esempio ITA).

## 1.2 Soluzione proposta Hadoop

Per risolvere il problema ho deciso di far utilizzo del pattern di summarization indice invertito e del pattern di filtering filtering. Le parole chiavi per cui costruiremo l'indice invertito sono i campi corrispondenti a Referre, Assistant 1 e Assistant 2, la lista di identificatori che associeremo sono il campo MatchId (che interpreto come identificatore univoco per la partite). La condizione per cui filtriamo le parole chiavi è la presenza della stringa passata da utente (ad esempio ITA). Il dataset presenta alcuni valori con caratteri anomali, ma nessuno all'interno dei campi che prendiamo in considerazione e quindi non c'è necessità di modificarlo. Per utilizzare la stringa passata da utente all'interno del mapper, la configuriamo nel driver, così da poterci accedere nel mapper.

```
1 Configuration conf = new Configuration();
2 conf.set("nationality", args[2]);
```

Non è stato utilizzato un combiner, dato che per il pattern indice invertito e filtro non partano alcun beneficio.

### 1.3 Codice Sorgente Mapper

```
3  class MapperReferee extends Mapper
4      <LongWritable, // Input key type
5          Text, // Input value type
6          Text, // Output key type
7          Text> { // Output value type
8              private String nationality;
9
10             @Override
11             public void setup(Context context){
12                 // Get the searched nationality from the context
13                 nationality = context.getConfiguration().get("nationality");
14             }
15             @Override
16             protected void map(LongWritable key, // Input key type
17                 Text value, // Input value type
18                 Context context) throws IOException,
19                 InterruptedException {
20                 // Split each sentence in words. Use whitespace"," as
21                 // delimiter
22                 // The split method returns an array of strings
23                 String[] words = value.toString().split(",");
24                 // Check field 13,14,15 (respectively Referee,
25                 // assistant 1 and assistant 2)
26                 // if it contains the searched nationality,
27                 // context write it + match Id
28                 // match Id is field 17
29                 if(words[13].contains(nationality)) context.
30                     write(new Text(words[13]), new Text(words
31                         [17]));
32                 if(words[14].contains(nationality)) context.
33                     write(new Text(words[14]), new Text(words
34                         [17]));
35                 if(words[15].contains(nationality)) context.
36                     write(new Text(words[15]), new Text(words
37                         [17]));
38             }
39 }
```

Una possibile aggiunta che si può introdurre nel mapper e anche rimuovere dalla stringa che contiene il nome dell'arbitro/assistente l'indicatore della nazionalità.

## 1.4 Codice Sorgente Reducer

```
30 class ReducerReferee extends Reducer
31     <Text, // Input key type
32     Text, // Input value type
33     Text, // Output key type
34     Text> { // Output value type
35
36     @Override
37     protected void reduce(Text key, // Input key type (Referee
38                           name)
39                           Iterable<Text> values, // Input value type (Match
40                           //Ids)
41                           Context context) throws IOException,
42                           InterruptedException {
43
44         String invIndex = new String();
45
46         // Iterate over the set of Match Ids and concatenate
47         // them
48         for (Text value : values) {
49             invIndex = invIndex.concat(value + ",,");
50         }
51
52         // avoid writing the last comma by taking the substring
53         context.write(key, new Text(invIndex.substring(0,
54                                     invIndex.length() - 1)));
55     }
56 }
```

Il codice per il reducer è praticamente identico a quello visto a lezione, l'unica aggiunta è la rimozione della virgola finale.

## 1.5 Codice Sorgente Driver

```
53 public class FinalProject {  
54  
55     /**  
56      * @param args the command line arguments  
57      */  
58     public static void main(String[] args) throws IOException,  
59             InterruptedException, ClassNotFoundException {  
60         Path inputPath;  
61         Path outputDir;  
62  
63         if(args.length < 3) {  
64             System.out.println("Usage : inputPath (hdfs) ,  
65                             outputPath (hdfs), nationality e.g ITA");  
66             return ;  
67         }  
68  
69         // Parse the parameters  
70         inputPath = new Path(args[0]);  
71         outputDir = new Path(args[1]);  
72         Configuration conf = new Configuration();  
73         conf.set("nationality", args[2]);  
74  
75         // Define a new job  
76         Job job = Job.getInstance(conf);  
77  
78         // Assign a name to the job  
79         job.setJobName("InvertedIndex");  
80  
81         // Set path of the input file/folder (if it is a folder  
82             // , the job reads all the files in the specified  
83             // folder) for this job  
84         FileInputFormat.addInputPath(job, inputPath);  
85  
86         // Set path of the output folder for this job  
87         FileOutputFormat.setOutputPath(job, outputDir);  
88  
89         // Specify the class of the Driver for this job  
90         job.setJarByClass(FinalProject.class);  
91  
92         // Set job input format  
93         job.setInputFormatClass(TextInputFormat.class);  
94  
95         // Set job output format  
96         job.setOutputFormatClass(TextOutputFormat.class);  
97  
98         // Set map class  
99         job.setMapperClass(MapperReferee.class);  
100    }
```

```

97     // Set map input key and value classes
98     // Set map output key and value classes
99     job.setMapOutputKeyClass(Text.class);
100    job.setMapOutputValueClass(Text.class);
101
102    // Set reduce class
103    job.setReducerClass(ReducerReferee.class);
104
105    // Set reduce output key and value classes
106    job.setOutputKeyClass(Text.class);
107    job.setOutputValueClass(Text.class);
108
109    // Set number of reducers
110    //job.setNumReduceTasks(numberOfReducers);
111
112    // Execute the job and wait for completion
113    System.exit(job.waitForCompletion(true) ? 0 : 1);
114}
115
116}

```

## 2 Spark

### 2.1 Descrizione del problema (statistica a scelta)

Realizzare la classifica delle top X squadre che hanno vinto più partite con Y o più goal di differenza, ad esempio se Y = 1, ottengo la classifica delle squadre che hanno vinto più partite. Considerare la possibilità che una partita finisca ai rigori. X e Y devono essere indicati dall'utente e sono entrambi interi positivi. La classifica deve essere stampata a video.

### 2.2 Soluzione proposta Spark

Per risolvere il problema ho deciso di fare utilizzo dell'equivalente per Spark dei pattern di filtering, di Record count (per group) e top K. Per gestire le ties aggiungo anche un pattern di indice invertito, così da poter inserire più squadre nella stessa posizioni della classifica. Il dataset presenta alcuni valori anomali nei campi che indicano il nome di due squadre. In particolare, prima di REPUBLIC OF IRELAND e UNITED ARAB EMIRATES è presente il seguente tipo RN">. Questo tipo è consistente, ed ogni volta che una delle due squadre è presente nella riga è sempre preceduta dal tipo, che quindi non è un problema per le statistiche (non dobbiamo gestire il caso in cui una volta la Repubblica di Irlanda è chiamata REPUBLIC OF IRELAND e un'altra volta RN">REPUBLIC OF IRELAND). Si potrebbe gestire questo tipo rimuovendolo prima di stampare a video la classifica.

## 2.3 Codice Sorgente Spark

La classe Match è realizzata per rendere il codice più leggibile, e per filtrare le partite che non terminano con Y goal di scarto tramite il metodo checkGoalDifference().

```
1 public static class Match implements Serializable {
2     String homeTeam;
3     String awayTeam;
4     Integer goalDifference;
5     Integer treshold;
6
7     public Match(String homeTeam, String awayTeam, Integer
8         goalHomeTeam, Integer goalAwayTeam, Integer treshold) {
9         this.homeTeam = homeTeam;
10        this.awayTeam = awayTeam;
11        this.goalDifference = goalHomeTeam - goalAwayTeam;
12        this.treshold = treshold;
13    }
14    /**
15     * @return The team that won the match with 3 or more goal
16     * lead, null otherwise
17     */
18    public Tuple2<String, Integer> checkGoalDifference() {
19        if (goalDifference >= treshold) {
20            return new Tuple2(homeTeam, 1);
21        }
22        if (goalDifference <= -treshold) {
23            return new Tuple2(awayTeam, 1);
24        }
25    }
}
```

Creo l'oggetto Spark Context per leggere il file, e scarto l'header.

```
26 // Create a configuration object and set the name of the
27 // application
28 SparkConf conf = new SparkConf().setAppName("Spark World Cup
29 // Stats");
30
31 // Create a Spark Context object
32 JavaSparkContext sc = new JavaSparkContext(conf);
33
34 // Build an RDD of Strings from the input textual file
35 // Each element of the RDD is a line of the input file
36 JavaRDD<String> lines = sc.textFile(inputFile);
37 // filter out the header
38 JavaRDD<String> linesWithoutHeader = lines.filter(s -> !s.
39     contains("Away Team Name"));
```

La lambda expression passata come parametro della funzione map, crea un RDD di matches, trasformando l'RDD di stringhe. All'interno della lambda expression viene anche considerata la possibilità che una partita finisca ai rigori.

```

37 // Create an RDD of matches
38 JavaRDD<Match> matches = linesWithoutHeader.map(s -> {
39     String[] splitted = s.split(",");
40     // check for penalties
41     if(splitted[9].contains("penalties"))
42     {
43         // Assuming that the penalites are the only integers in
44         // the string, this leaves us with a string containing
45         // the penalties shootout result separated by "-" (e.g
46         // 4-1)
47         String [] penalties = splitted[9].replaceAll("[^0-9-]",
48             "").split("-");
49         // adding the result of the penalties to the scores
50         return new Match(splitted[5], splitted[8], Integer.
51             valueOf(splitted[6]) + Integer.valueOf(penalties[0])
52             , Integer.valueOf(splitted[7])+ Integer.valueOf(
53                 penalties[1]),threshold);
54     }
55     // the value we are intrest in are homeTeam (5), awayTeam
56     // (6) and the result of the match (7 and 8)
57     return new Match(splitted[5], splitted[8], Integer.valueOf(
58         splitted[6]), Integer.valueOf(splitted[7]),threshold);
59 });

```

Creo ora il PairRDD contenente le tuple<NOMEQUADRA,1>. Devo usare il la funzione flatMapToPair perché la funzione può restituire tuple vuote, quando la partita finisce senza Y goal di differenza. Questo equivale al mapper del pattern record count (Per group)

```

54 // Create the RDD <teamName,1>
55 // flat map because not every match ends with a three goal
56 // difference (or larger)
57 JavaPairRDD<String, Integer> partialWinners = matches.
58     flatMapToPair(f -> {
59         List<Tuple2<String, Integer>> pairs = new ArrayList();
60         if(f.checkGoalDifference() != null){
61             pairs.add(f.checkGoalDifference());
62         }
63         return pairs.iterator();
64     });

```

Ora non resta che eseguire un reduceByKey che equivale al reduce del pattern record count (Per group) per ottenere il PairRDD contenente le tuple tuple<NOMESQUADRA,NUMERODIVITTORECONYGOALDIASCARTO>

```

63 // Create the RDD <teamName, numberOfWorksWithThreeGoalLead>
64 JavaPairRDD<String, Integer> reducedWinners = partialWinners.
65     reduceByKey((a, b) -> {
66         return a + b;
67     });

```

A questo punto si può eseguire la seguente azione<sup>1</sup>, che equivale ad un pattern di top K, gestire le ties con il comparatore e stampare a video il risultato.

```

1 List<Tuple2<String, Integer>> Results = reducedWinners.top(K,
    new Tuple2Comparator());

```

In alternativa si è scelto di poter ottenere posizioni in classifica a pari merito, per fare ciò si possono invertire i nomi delle squadre e il numero di partite, eseguire una istruzione di reduce per concatenare i nomi delle squadre (simil reducer del pattern indice invertito), ordinarli per chiave in ordine discendente e prelevare i primi X elementi.

```

74 // Create the RDD <numberOfWinsWithThreeGoalLead, teamName>
75 JavaPairRDD<Integer, String> invertedReducedWinners =
76     reducedWinners.mapToPair(f -> {
77         return new Tuple2(f._2(), f._1());
78     });
79
80 // Create the RDD <numberOfWinsWithThreeGoalLead, list of
81 // teamNames>
82 JavaPairRDD<Integer, String> invertedIndex =
83     invertedReducedWinners.reduceByKey((a,b)->{
84         return a + " " + b;
85     });
86 // Sort the RDD by key invertedIndex
87 JavaPairRDD<Integer, String> sortedInvertedIndex =
88     invertedIndex.sortByKey(false);
89 // Take the top K of the sorted invertedIndex
90 List<Tuple2<Integer, String>> TopK = sortedInvertedIndex.take(K
91 );

```

---

<sup>1</sup>Ho provato ad usare una lambda expression per il comparatore, ma causa un'eccezione 'task non serializzabile'. Sembra che l'espressione lambda che funge da comparatore non sia serializzabile e causi questo problema, che può essere risolto creando una classe che implementa COMPARATOR<TUPLE2<STRING, INTEGER>>, SERIALIZABLE

Finiamo stampando a video i risultati e chiudendo la spark Context

```
88 System.out.println("Le Top " + K + " Squadre dei mondiali di  
89     calcio tenutasi tra il 1930 e il 1998 che hanno vinto il  
90     maggior numero di partite con "+ threshold +" o maggior goal  
91     di scarto (rigori inclusi)");  
92 for (Tuple2<Integer, String> t : TopK) {  
93     System.out.println(t);  
94 }  
95 sc.close();
```

### 3 Risultati

#### 3.1 Risultato Hadoop con input utente ITA

AGNOLIN Luigi (ITA) 551,610,120,564,398,627  
ANGONESE Aurelio (ITA) 2066,1987,2004,1984  
BALDAS Fabio (ITA) 3063  
BARLASSINA Rinaldo (ITA) 1129,1130,1175,1119  
BONIVENTO Ferruccio (ITA) 1133  
CAIRONI Camillo (ITA) 1104,1105  
CARMINATI Ettore (ITA) 1111  
CARRARO Albino (ITA) 1105,1102  
CASARIN Paolo (ITA) 741,764,959,900,922  
COLLINA Pierluigi (ITA) 8762,8733  
D ELIA Pietro (ITA) 111,129  
DATTILO Generoso (ITA) 1199,1230,1202,1207,1119  
GALEATTI Giovanni (ITA) 1199,1231,1230,1202  
GONELLA Sergio (ITA) 2201,2198,2433,2246  
JONNI Cesare (ITA) 1564,1559,1480,1562  
LANESE Tullio (ITA) 290,102,114,75  
LO BELLO Concetto (ITA) 1659,1689,1634  
LO BELLO Rosario (ITA) 152  
LONGHI Carlo (ITA) 29  
MAGNI Pierluigi (ITA) 175  
MATTEA Francesco (ITA) 1108,1106  
MAZZEI Gennaro (ITA) 8733,8759,8752  
MELANDRI Ermenegildo (ITA) 1108  
ORLANDINI Vincenzo (ITA) 1313,1233,1264,1304,1434,1423,1387,1278  
PAIRETTO Pierluigi (ITA) 3088,25,198  
RAMICONE Domenico (ITA) 3057,3063,3088  
SASSI Otello (ITA) 1119  
SBARDELLA Antonio (ITA) 1840,1893,1780,1843  
SCARPI Giuseppe (ITA) 1176,1152,1141,1156  
SCORZONI Raffaele (ITA) 1141  
TURBIANI Giuseppe (ITA) 1102

### 3.2 Risultati Spark con Input X = 15 e Y = 1

(55, Brazil)  
(39, Germany FR)  
(38, Italy)  
(32, Argentina)  
(23, France)  
(20, England)  
(16, Yugoslavia Spain)  
(15, Soviet Union Hungary Uruguay Sweden)  
(14, Netherlands)  
(13, Poland)  
(12, Austria)  
(11, Czechoslovakia)  
(10, Belgium)  
(9, Germany)  
(8, Mexico Romania)

### 3.3 Risultati Spark con Input X = 15 e Y = 4

(7, Hungary Brazil)  
(6, Germany FR)  
(4, Argentina Uruguay)  
(3, Netherlands Yugoslavia France)  
(2, Soviet Union Spain Sweden Czechoslovakia Poland)  
(1, Denmark Turkey Mexico Italy Bulgaria Russia Austria)

## 4 Come eseguire il codice

### 4.1 Eseguire Hadoop

```
1 # Inserire il dataset nel cartella data (dove e' montato il  
2 # volume)  
3 # Inserire il jar nella cartella data (dove e' montato il  
# volume)  
3 # Eseguire i seguenti comandi all'interno della bash del nodo  
# master  
4  
5 # avviare HADOOP e YARN  
6 $HADOOP_HOME/sbin/start-dfs.sh  
7 $HADOOP_HOME/sbin/start-yarn.sh  
8  
9 # Inserire il dataset nel dfs  
10 hdfs dfs -put worldcup1.csv /dataset  
11  
12 # Lancio Hadoop jar  
13 # sostituire $0 con le iniziali di una nazione e.g ITA
```

```
14  hadoop jar FinalProject.jar /dataset /refereeIndex $0
15
16  # Prendo dal dfs il risultato di hadoop
17  hdfs dfs -get /refereeIndex
18
19  # Pulizia del dfs (cancello il dataset e il risultato di Hadoop
    )
20  hdfs dfs -rm -r /refereeIndex
21  hdfs dfs -rm /dataset
```

## 4.2 Eseguire Spark

1. Inserire il jar SparkWorldCupStats.jar, la cartella input, i due eseguibili shell launch\_cluster.sh e launch\_single.sh nella stessa cartella.
2. Se si vogliono cambiare i parametri di ingresso, modificando gli eseguibili shell, l'usage del programma è il seguente: Cartella che contiene il csv da processare, X, Y.
3. Eseguire uno dei due eseguibili shell.