# Project in Spark 2017

Adrianna Janik

Ion Mosnoi

Lei Guo

October 24, 2017

## 1 TASK

Firstly we uncompressed the data stored in ling-spam.zip folder with *Extract all* command. Secondly we open Virtual Box machine with Hortonworks, we signed in with maria_dev username and maria_dev password on Ambari available under 127.0.0.1:8080 ip address. We have selected *Files view*, than navigated to *tmp* folder and created directories *tmp/ling-spam/ham* and *ling-spam/spam*. Following that we logged in with ssh credentials to Hortonworks machine

```
1  $ssh root@127.0.0.1 -p 2222
```

In the meantime upload to the virtual machine ling-spam.zip with:

```
1  $sudo scp -P 2222 ../ling-spam.zip  root@127.0.0.1:/tmp/
```

We unzipped ling-spam.zip with:

```
1  $unzip ling-spam.zip -d /tmp/ling-spam
```

We putted files into /tmp/ling-spam/ folder in hdfs with:

```
1  $hdfs dfs -put ./ling-spam/ham /tmp/ling-spam/ham
2  $hdfs dfs -put ./ling-spam/spam /tmp/ling-spam/spam
```

## 2 TASK

Installation of sbt:

```
1  $wget http://dl.bintray.com/sbt/rpm/sbt-0.13.12.rpm
```

Edit file /etc/yum.repos.d/sandbox.repo:

```
1  ~[sandbox]
2  ~name=Sandbox repository (tutorials)
3  ~gpgcheck=0
4  ~enabled=0
5  ~baseurl=http://dev2.hortonworks.com.s3.amazonaws.com/repo/dev/
   ↪ master/utils/
```

```
1  $yum clean all
2  $yum update
3  $sudo yum localinstall sbt-0.13.12.rpm
4  $sbt -update
5  $sudo scp -P 2222 -r ../spamTopWords/*  root@127.0.0.1:/tmp/
   ↪ spamTopWords/
6  $sbt package
```

## 3 TASK

Firstly we created Spark Context with:

```
1  val conf = new SparkConf(.setAppName(''Spam Filter Application'').
   ↪ setMaster(''local'')
2  val sc = new SparkContext(conf))
```

Than we called function *probaWordDir* with defined spark context as well as folder name for which we want to count words.

```
1  val (probaHW, nbHFiles) = probaWordDir(sc)(args(0)+"ham/*.txt")
2  print("number of files in "+ args(0)+"ham/*.txt" +":")
3  println(nbHFiles)
4
5
6  //process spam files
7  val (probaSW, nbSFiles) = probaWordDir(sc)(args(0)+"spam/*.txt")
8  print("number of files in "+ args(0)+"spam/*.txt" +":")
9  println(nbSFiles)
```

```
1  val rdd = sc.wholeTextFiles(filesDir)
2  // The number of files is counted and stored in a variable nbFiles
3  val nbFiles = rdd.count()
4  // Non informative words must be removed from the set of unique
   ↪ words.
```

```scala
5  val stopWords = Set(".", ":", ",", " ", "/", "\\", "-", ",", "(",
   ↪ ")", "@")
6  // Each text file must be splitted into a set of unique words
7  //(if a word occurs several times, it is saved only one time in
   ↪ the set).
8  val wordBagRdd: RDD[(String, Set[String])] = rdd.map(textTuple =>
9        (textTuple._1, textTuple._2.trim().
10       split("\\s+").toSet.diff(stopWords)))
11 // Get the Number of occurrences amongst all files
12 val wordDirOccurency: RDD[(String, Int)] = wordBagRdd.flatMap(
13 x => x._2.map(y => (y, 1))).reduceByKey(_ + _)
14 val probaWord: RDD[(String, Double)] = wordDirOccurency.map(
15 x => (x._1, x._2.toDouble / nbFiles))
16 return (probaWord, nbFiles)
```

# 4 ,5 Tasks

Function: probaWordDir:

```scala
1  def probaWordDir(sc:SparkContext)(filesDir:String)
2  :(RDD[(String, Double)], Long) = {
3
4
5      val rdd = sc.wholeTextFiles(filesDir)
6      // The number of files is counted and stored in a variable
          ↪ nbFiles
7      val nbFiles = rdd.count()
8      // Non informative words must be removed from the set of
          ↪ unique words.
9      val stopWords = Set(".", ":", ",", " ", "/", "\\", "-", ",",
          ↪ "(", ")", "@")
10     // Each text file must be splitted into a set of unique
          ↪ words (if a word occurs several times, it is saved
          ↪ only one time in the set).
11     val wordBagRdd: RDD[(String, Set[String])] = rdd.map(
          ↪ textTuple =>
12           (textTuple._1, textTuple._2.trim().
13           split("\\s+").toSet.diff(stopWords)))
14     // Get the Number of occurrences amongst all files
15     val wordCountRdd: RDD[(String, Int)] = wordBagRdd.flatMap(x
          ↪ => x._2.map(y => (y, 1))).reduceByKey(_ +_)
16     val probaWord: RDD[(String, Double)] = wordCountRdd.map(x =>
          ↪ (x._1, x._2.toDouble / nbFiles))
17     return (probaWord, nbFiles)
18
19
20 }
```

Main function:

```scala
def main(args: Array[String]) {

    if(args.size > 0){
            val conf = new SparkConf().setAppName("Spam Filter
                ↪ Application").setMaster("local")
            val sc = new SparkContext(conf)
            println("Got the path:"+args(0))
            // args(0) should be something like "hdfs:///project
                ↪ /, see readme

            //process ham files
            val (probaHW, nbHFiles) = probaWordDir(sc)(args(0)+"
                ↪ ham/*.txt")

            //process spam files
            val (probaSW, nbSFiles) = probaWordDir(sc)(args(0)+"
                ↪ spam/*.txt")
            print("number of files in "+ args(0)+"ham/*.txt" +":
                ↪ ")
            println(nbHFiles)
            print("number of files in "+ args(0)+"spam/*.txt" +"
                ↪ :")
            println(nbSFiles)

            val nbFiles = nbSFiles + nbHFiles
            val probaW = probaSW.union(probaHW).reduceByKey((x,y
                ↪ ) => (x*nbSFiles.toDouble+y*nbSFiles.toDouble)
                ↪ /(nbFiles.toDouble)) //not sure

            //Compute the probability P(occurs, class) for each
                ↪ word.

            val probaH = nbHFiles.toDouble / nbFiles.toDouble //
                ↪  the probability that an email belongs to the
                ↪ given class.
            val probaS = nbSFiles.toDouble / nbFiles.toDouble
            // Compute mutual information for each class and
                ↪ occurs
            val MITrueHam = computeMutualInformationFactor(
                ↪ probaHW, probaW, probaH, 0.2 / nbFiles) // the
                ↪  last is a default value
            val MITrueSpam = computeMutualInformationFactor(
                ↪ probaSW, probaW, probaS, 0.2 / nbFiles)
            val MIFalseHam = computeMutualInformationFactor(
                ↪ probaHW.map(x => (x._1, 1 - x._2)), probaW,
                ↪ probaH, 0.2 / nbFiles)
            val MIFalseSpam = computeMutualInformationFactor(
                ↪ probaSW.map(x => (x._1, 1 - x._2)), probaW,
```

```scala
                            ↪ probaS , 0.2 / nbFiles)
31
32                  //compute the mutual information of each word as a
                        ↪ RDD with the map structure: word => MI(word)
33                  //sum the prob for all words
34                  val MI :RDD[(String , Double)] = MITrueHam.union(
                        ↪ MITrueSpam).union(MIFalseHam).union(
                        ↪ MIFalseSpam).reduceByKey( (x, y) => x + y)
35
36                  // print on screen the 10 top words (maximizing the
                        ↪ mutual information value)
37                  //These words must be also stored on HDFS in the
                        ↪ file âĂIJ/tmp/topWords.txtâĂİ.
38                  val path: String = "/tmp/topWords.txt"
39                  val topTenWords: Array[(String , Double)] = MI.top
                        ↪ (10)(Ordering[Double].on(x => x._2))
40                  //coalesce to put the results in a single file
41                  sc.parallelize(topTenWords).keys.coalesce(1, true).
                        ↪ saveAsTextFile(path)
42          }
43          else
44                  println("Please write te directory where the ham and
                        ↪  spam")
45  }
```

---

Function: computeMutualInformationFactor

---

```scala
1  def computeMutualInformationFactor(
2    probaWC: RDD[(String , Double)],//prob of just a class , some word
          ↪  could not be
3    probaW: RDD[(String , Double)],//all words prob, all word
4    probaC: Double , //prb of a class : class mails / all mails
5    probaDefault: Double // default value when a probability is
          ↪ missing
6  ): RDD[(String , Double)] = {
7          //p(occurs) =
8          val probWJoin: RDD[(String , (Double , Option[Double]))] =
                ↪  probaW.leftOuterJoin(probaWC)// got all class
                ↪ probs, if not -> default
9                          //p(accurs)  p(accurs,class)
10         val valueClassAndOcu: RDD[(String , (Double , Double))] =
                ↪ probWJoin.map(x => (x._1, (x._2._1, x._2._2.
                ↪ getOrElse(probaDefault))))
11         //We have to change ln to log2 (by using ln(x)/ln(2)=
                ↪ log2(x)
12         valueClassAndOcu.map(x => (x._1, x._2._2 * (math.log(x.
                ↪ _2._2 / (x._2._1 * probaC)) / math.log(2.0))))
13
14  }
```

---