

Introduction to Big Data.

Assignment Report.

Spark Wars: Lada Morozova, Danis Alukaev, Denis Schegletov

Introduction.

The goal of the assignment was to finalize the implementation of a movie recommendation system in Scala. Moreover, it was required to use Spark framework and implement it not only on local Hadoop cluster but also on Private Network Hadoop cluster.

To finalize implementation, team implement `parseTitle` and `rmse` functions; post-processing of recommendations; method that allow to load user data directly instead of surveying user. Moreover, the team will run the application with different model ranks, compare and evaluate the differences between baseline and prediction for the model with a higher rank. In addition, low-frequency items will be excluded from the training and recommendation process.

Implementation of functions

`parseTitle`

The function `parseTitle` is used for splitting film entry and extracting from it the title of the film. Implementation of `parseTitle` from the `MovieLensALS.scala`:

```
def parseTitle(filmEntry: String) = {  
  // Task 0.1. Complete the code  
  "\".*\\\"".r findFirstIn filmEntry match {  
    case Some(s) => s.replace("\\\"", "")  
    case None => filmEntry.split(",").map(_.trim).lift(1)  
  }  
}
```

`rmse`

The function `rmse` used for calculating root-mean-square error between ground truth ratings of the film and the ratings predicted by the system. The root-mean-square error is calculated as $\sqrt{\frac{\sum_{i=1}^N (x_i - \hat{x}_i)^2}{N}}$. Implementation of `RMSE` from the `MovieLensALS.scala`:

```
def rmse(test: RDD[Rating], prediction: scala.collection.Map[Int, Double]) = {  
  // Task 0.2. Complete the code  
  val rating_pairs = test  
    .map(x => (x.rating, prediction.getOrElse(x.product, 0.0)))  
  math.sqrt(rating_pairs  
    .map(x => (x._1 - x._2) * (x._1 - x._2))  
    .reduce(_ + _) / test.count())  
}
```

Post-processing of recommendations

The goal of the recommender system is to present to users 20 movies that can be considered as most suitable for following users based on his/her ratings of other movies. To do so, the prediction based on the user's rating are made, then watched movies are excluded, and 20 movies with the highest ratings are presented to the user. The following code snippet from `MovieLensALS.scala` depicts our solution for the following task:

```

if (doGrading) {
    println("Predictions for user\n")

    // Task 1. Post-processing
    // collect all movie ids in user rating
    val filmIds = sc.textFile(ratingsPath + "/user_rating.tsv")
        .map(line => line.split("\t"))
        .map(x => x(0).toInt)
        .collect()

    // for input data format refer to documentation
    // get movie ids to rank from baseline map
    model.predict(sc.parallelize(baseline.keys.map(x => (0, x)).toSeq))
        // Task 1. Post-processing
        // filter watched movies
        .filter(x => !filmIds.contains(x.product))
        // sort by ratings in descending order
        .sortBy(_._rating, false)
        // take first 20 items
        .take(20)
        // print title and predicted rating
        .foreach(x => println(s"${filmId2Title(x.product)}\t\t${x.rating}"))
}

```

Data loading

To improve loading of user's ratings data the two methods was introduced: `loadPreferences()` and `predictPreferences()`. First method `loadPreferences()` allows to upload user's rating with use of .tsv file instead of surveying the user film by film. However, the ability to evaluate films one by one has been retained and has been moved to a separate function `predictPreferences()`. The implementation of `loadPreferences()` and `predictPreferences()` from `Grader.scala`

```

// Task 2. Load movie preferences
println("Do you want to load your preferences?")
val answer = scala.io.StdIn.readLine()
val answerPattern: Regex = "(y|yes)$".r

def loadPreferences() = {
    println("Loading from user_rating.tsv...")
    graded = sc.textFile(path + "/user_rating.tsv")
        .map(line => line.split("\t"))
        .map(x => (x(0).toInt, x(1).toDouble))
        .collect()
}

def predictPreferences() = {
    while (step < limit) {
        println(s"\nGraded ${step}/${limit}, Viewed ${position}/${films.length}")
        println(s"${grading_message}\n${films(position)._2}")
        var grade = scala.io.StdIn.readLine()

        try {
            var intGrade = grade.toDouble

            if (intGrade < 0.0 || intGrade > 5.0) {
                throw new NumberFormatException()
            }
        }
    }
}

```

```

if (intGrade != 0) {
  graded = graded :+ (films(position)._1, intGrade)
  step += 1
  position += 1
} else {
  position += 1
}

if (position == films.length) {
  step = limit + 1
  println("\nNo more movies to grade\n")
}

if (step == limit) {
  println("Grade 20 more? y/[n]")
  var ans = scala.io.StdIn.readLine()

  if (ans == "y" || ans == "yes" || ans == "Yes" || ans == "YES") {
    limit += 20
  } else {
    println("\nFinished Grading\n")
  }
}
} catch {
  case e: NumberFormatException => println("Try again")
  case e: Exception => println("Unknown Error")
}
dumpRatings()
}
}

```

Comparison of models

There is no significant difference in baseline error. The least training error was reached for rank of 2.

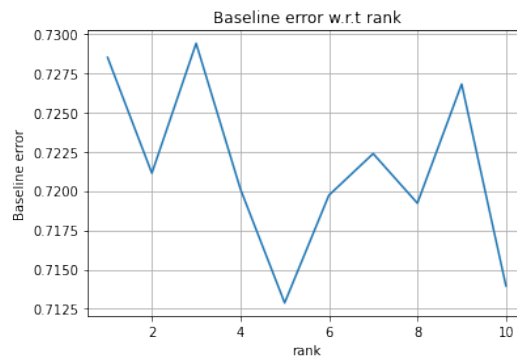


Figure 1

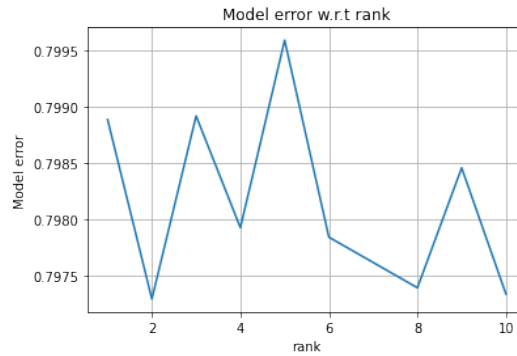


Figure 2

Elimination of low-frequency items

To improve the recommender system the extra filtering was introduced. Movies with a total score of less than 50 were excluded from the learning process, as well as from the recommendation process. The code snippet with elimination process from **MovieLensALS.scala**:

```
// Task 4. Extra Filtering
val films = ratingsData.map(x => x.product).countByValue().filter(x => x._2 > 50).toArray.map(x=>x._1)
ratingsData = ratingsData.filter(x => films.contains(x.product))
// End of Task 4. Extra Filtering
```

Private Network Hadoop cluster

For the test purposes there was built cluster of two machines with Ubuntu 18.04 on the board. Often distributed computations involve several hosts with approximately equal resources, so it was decided to meet this requirement using Virtual Machines. As a provider of virtualization we have used the Virtualbox. Configuration files and source code is available online.

In Vagrantfile we configured the bridge to public network, so that the DHCP server in our private network assigns the IP address to Virtual Machine independently from the host. The IP address could be determined using hostname utility. However, we faces the difficulty with *InnopolisUniversity* network that does not assign IP addresses to Virtual Machines. Consequently, all the further actions were performed on Wi-Fi of our room.

On master node with IP address 192.168.0.194 we installed Hadoop 3.3.1 and Spark of version 3.2.0. The machine was configured and sample jobs runned. Then, all the settings were passed on the slave node with IP address 192.168.0.118. The former plays the role of Namenode, Secondary Namenode, Resource Manager, Datanode, and Worker, whereas the second - as Datanode and Worker. Figure 3 demonstrates the overall HDFS filesystem in the configured cluster.

The implemented earlier code was compiled using sbt package, and the submitted as a job on private network cluster. The sample output of the program is presented in Figure 4.

```

vagrant@ubuntu1804:~$ hdfs dfsadmin -report
Configured Capacity: 62331412480 (58.05 GB)
Present Capacity: 52220649472 (48.63 GB)
DFS Remaining: 52220600320 (48.63 GB)
DFS Used: 49152 (48 KB)
DFS Used%: 0.00%
Replicated Blocks:
    Under replicated blocks: 0
    Blocks with corrupt replicas: 0
    Missing blocks: 0
    Missing blocks (with replication factor 1): 0
    Pending deletion blocks: 0
Erasure Coded Block Groups:
    Low redundancy block groups: 0
    Block groups with corrupt internal blocks: 0
    Missing block groups: 0
    Pending deletion blocks: 0

-----
Live datanodes (2):

Name: 192.168.0.118:9866 (slave2)
Hostname: ubuntu1804.localdomain
Decommission Status : Normal
Configured Capacity: 31165706240 (29.03 GB)
DFS Used: 24576 (24 KB)
Non DFS Used: 3450871808 (3.21 GB)
DFS Remaining: 26108080128 (24.32 GB)
DFS Used%: 0.00%
DFS Remaining%: 83.77%
Configured Cache Capacity: 0 (0 B)
Cache Used: 0 (0 B)
Cache Remaining: 0 (0 B)
Cache Used%: 100.00%
Cache Remaining%: 0.00%
Xceivers: 1
Last contact: Sun Nov 28 08:47:08 PST 2021
Last Block Report: Sun Nov 28 08:46:59 PST 2021
Num of Blocks: 0

Name: 192.168.0.194:9866 (slave1)
Hostname: ubuntu1804.localdomain
Decommission Status : Normal
Configured Capacity: 31165706240 (29.03 GB)
DFS Used: 24576 (24 KB)
Non DFS Used: 3446431744 (3.21 GB)
DFS Remaining: 26112520192 (24.32 GB)
DFS Used%: 0.00%
DFS Remaining%: 83.79%
Configured Cache Capacity: 0 (0 B)
Cache Used: 0 (0 B)
Cache Remaining: 0 (0 B)
Cache Used%: 100.00%
Cache Remaining%: 0.00%
Xceivers: 1
Last contact: Sun Nov 28 08:47:08 PST 2021
Last Block Report: Sun Nov 28 08:46:53 PST 2021
Num of Blocks: 0

```

Figure 3: HDFS Filesystem

```

root@bubuntu190:~# spark-submit --driver-memory 2g /home/vagrant/spark-recommendation.jar --true 8
2021-11-28 10:05:46,757 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
2021-11-28 10:05:47,964 INFO util.log: Logging initialized @2427ms to org.sparkproject.jetty.util.log.Slf4jLog
2021-11-28 10:05:48,098 INFO server.Server: jetty-9.4.43.v20210629; built: 2021-06-30T11:07:22.254Z; git: 52600ecfa3af7f1a27ef3a288e2bef7ea9dd7e8; jvm 1.8.0_191-8u191-b12-bubuntu0.18.04.1-b12
2021-11-28 10:05:48,145 INFO server.Server: Started @2610ms
2021-11-28 10:05:48,197 INFO server.AbstractConnector: Started ServerConnector@5afa86a1[HTTP/1.1, (http/1.1)](master:4040)
2021-11-28 10:05:48,234 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@668591fb[/jobs,null,AVAILABLE,@Spark]
2021-11-28 10:05:48,238 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@51ec2df1[/jobs/json,null,AVAILABLE,@Spark]
2021-11-28 10:05:48,239 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@15fa55a6[/jobs/job,null,AVAILABLE,@Spark]
2021-11-28 10:05:48,243 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@94d74c3a[/jobs/job/json,null,AVAILABLE,@Spark]
2021-11-28 10:05:48,245 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@664138b0c[/stages,null,AVAILABLE,@Spark]
2021-11-28 10:05:48,246 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@22d9c961[/stages/json,null,AVAILABLE,@Spark]
2021-11-28 10:05:48,248 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@bdc8014f[/stages/stage,null,AVAILABLE,@Spark]
2021-11-28 10:05:48,252 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@1698d7c0[/stages/stage/json,null,AVAILABLE,@Spark]
2021-11-28 10:05:48,253 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@87abca8f[/stages/pool,null,AVAILABLE,@Spark]
2021-11-28 10:05:48,257 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@782168b7[/stages/pool/json,null,AVAILABLE,@Spark]
2021-11-28 10:05:48,259 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@7435a578[/storage,null,AVAILABLE,@Spark]
2021-11-28 10:05:48,261 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@13847d7d[/storage/json,null,AVAILABLE,@Spark]
2021-11-28 10:05:48,264 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@65bb9029[/storage/rdd,null,AVAILABLE,@Spark]
2021-11-28 10:05:48,265 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@2b214b94[/storage/rdd/json,null,AVAILABLE,@Spark]
2021-11-28 10:05:48,269 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@49601f82[/environment,null,AVAILABLE,@Spark]
2021-11-28 10:05:48,272 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@2b8d084f[/environment/json,null,AVAILABLE,@Spark]
2021-11-28 10:05:48,273 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@24fabd0f[/executors,null,AVAILABLE,@Spark]
2021-11-28 10:05:48,275 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@661f3fb8[/executors/json,null,AVAILABLE,@Spark]
2021-11-28 10:05:48,277 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@43203a4f[/executors/threadDump,null,AVAILABLE,@Spark]
2021-11-28 10:05:48,280 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@66e5272f[/executors/threadDump/json,null,AVAILABLE,@Spark]
2021-11-28 10:05:48,294 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@69c93ca4[/static,null,AVAILABLE,@Spark]
2021-11-28 10:05:48,297 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@19569ebd[/,null,AVAILABLE,@Spark]
2021-11-28 10:05:48,300 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@6e1f9469[/api,null,AVAILABLE,@Spark]
2021-11-28 10:05:48,301 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@652e6a6c[/jobs/job/kill,null,AVAILABLE,@Spark]
2021-11-28 10:05:48,303 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@31ee96f4[/stages/stage/kill,null,AVAILABLE,@Spark]
2021-11-28 10:05:48,696 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@f1a45f8f[/metrics/json,null,AVAILABLE,@Spark]
2021-11-28 10:05:49,929 INFO mapred.FileInputFormat: Total input files to process : 1
2021-11-28 10:05:51,324 INFO mapred.FileInputFormat: Total input files to process : 1
2021-11-28 10:08:12,726 WARN netlib.InstanceBuilder$NativeBLAS: Failed to load implementation from:dev.ludovic.netlib.blas.JNIBLAS
2021-11-28 10:08:12,788 WARN netlib.InstanceBuilder$NativeBLAS: Failed to load implementation from:dev.ludovic.netlib.blas.ForeignLinkerBLAS
2021-11-28 10:08:14,370 WARN netlib.InstanceBuilder$NativeLAPACK: Failed to load implementation from:dev.ludovic.netlib.lapack.JNILAPACK
Error after training: 0.799361829919522

```

Figure 4: Recommender System run

Contribution of each member

The code was implemented collaboratively using IntelliJ Idea "Code with me" feature.

Danis Alukaev	Implementation of functions Implementation of data loading
Denis Schegletov	Implementation of post-processing of recommendation Elimination of low-frequency items
Lada Morozova	Composition of report Comparison of models

Conclusion

There was finished provided in specification template. For instance, we successfully added parse title and rmse functionality, post-processing of recommendations, data loading, and elimination of low-frequency items. In addition, the final version of the recommendation system was executed in private network Hadoop cluster. All the work was implemented collectively and tasks distributed among team members.