

1.

```

1 import org.apache.spark.rdd.RDD
2
3 case class Movie(movie_ID: Integer, title: String, genre: String)
4 case class User(user_ID: Integer, gender: String, age: Integer, occupation: String, zip_code: String)
5 case class Rating(user_ID: Integer, movie_ID: Integer, rating: Integer, timestamp: String)
6
7 def parseRatings(row: String): Rating = {
8   val splitted = row.split(":").map(_._trim).toList
9   return Rating(splitted(0).toInt, splitted(1).toInt, splitted(2).toInt, splitted(3))
10 }
11
12 val ratings = sc.textFile("ratings.dat").map(element => parseRatings(element)).cache
13
14 println("Lines that ratings RDD contain:", ratings.count)
15 println("Unique movies that have been rated:", ratings.map(x => x.movie_ID).distinct.count())
16 println("User who rated the most, number of ratings: ", ratings.keyBy(x => x.user_ID).mapValues(x => 1).reduceByKey((x, y) => x + y).sortBy(_._2, false).first)
17
18 println("User who gave the most '5' ratings, number of ratings: ", ratings.filter(x => x.rating == 5).keyBy(x => x.user_ID).aggregateByKey(0)((acc, x) => acc + 1, (x, y) => x + y).sortBy(_._2, false).first)
19

```

2.

### **APPROACH:**

First, the respective classes are defined. Then the RDD file is loaded using the parseRatings, parseMovies, parseUsers function that splits the row by "::".

The number of records in movies and users RDD is calculated by count function scala function on RDD.

Number of comedy movies is found by filtering the movies RDD by its genre with the help of filter and contain function and then the number of comedy movies is found by using the count function.

To find the comedy movie with most ratings, we first filter out the comedy movies as in (b). Then I used the movie id as a key and join this modified table with the ratings RDD. Then I map each row by the movie title and group them which is then mapped to movietitle, number of ratings.

Then using the sort in desc, I get the movie with the most ratings.

To find the number of unique users that rated the movies with movie\_IDs 2858, 356 and 2329, I first filter the ratings RDD to contain only movies having respective movie\_IDs and then group it by user\_ID. Then reduced it by the distinct keys clause which remove duplicate rows having same user id. Then using the count function, I have reported the answer.

Created an inverted index on ratings using field movie\_ID and printed the first entry.

Created respective movie\_ID arrays and joined it with the inverted index. Then mapped using the user ID and found the distinct count.

### **OUTPUT:**

- a. Loaded the movies and user RDD successfully. Number of records in movies RDD = 3883, in users RDD = 6040
- b. Number of comedy movies = 1200
- c. Comedy movie having most ratings = American Beauty (1999),3428
- d. Number of unique users without inverted index = 4213  
Time taken = 0.617s
- e. Inverted index. First entry =
- f. Number of unique users with inverted index = 4213  
Time taken = 9.569s



```
println("Comedy with most ratings (Movie, Number of rating): ", movies.filter(x => x.genre.contains("Comedy")).keyBy(x =>
x.movie_ID).join(ratings.keyBy(x => x.movie_ID).map(x => (x._2._1.title)).groupBy(x => x).map{case (k, v) => (k, v.size)}.sortBy(x =>
x._2, false).first)

val stime1: Long = (System.currentTimeMillis)
println("Answer to e:", ratings.filter(x => (List(2858, 356, 2329).contains(x.movie_ID))).groupBy(x => x.user_ID).keys.distinct.count)
println("Time taken: ", ((System.currentTimeMillis) - stime1)/1000.0)

val ratingsInv = ratings.groupBy(r => r.movie_ID).cache

println("First item of inverted index on ratings: ")
print(ratingsInv.lookup(1).take(1))

val stime1: Long = (System.currentTimeMillis)
val items = sc.parallelize(List(Integer)(2858, 356, 2329)).keyBy(x => x)
println("Answer to g: ", ratingsInv.join(items).flatMap(x => x._2._1).map(x => x.user_ID).distinct.count)

println("Time taken:", ((System.currentTimeMillis) - stime1)/1000.0)
```

3.

### APPROACH:

Loaded the file and split it by “, | --”.

invertIndex function takes a field and uses groupBy to compile all records with the key into an iterable list mapped to the key and returns this rdd.

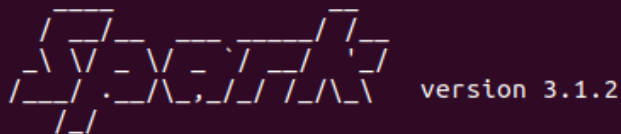
To get the answer without inverted key, I filter out the Rdds with download\_id “gthorren-22” map them to the repo name of the log if it has any. Then using distinct count, I report the answer.

To get the answer with inverted index, I used filter to search key, “gthorren-22” on the inverted index we have created over download\_id. Since all this data is on a single executor, I collect it and parallelize it to distribute it on many executors. Then I map it to the repo name and use distinct(), count as in (b).

### OUTPUT:

- Created function invertIndex
- 4577
- 4577

```
Spark session available as 'spark'.
(Unique repos with download id gthorren-22 without using inverted index:,4577)
21/08/31 23:03:57 WARN TaskSetManager: Stage 4 contains a task of very large size (94910 KiB). The maximum recommended task size is 1000 KiB.
(Unique repos with download id gthorren-22 using inverted index:,4577)
Welcome to
```



## CODE:

```
import java.util.TimeZone
import java.text.SimpleDateFormat

case class schema (debug_level:String,timestamp:java.util.Date,download_id:String,retrieval_stage:String,rest:String)//
TimeZone.setDefault(TimeZone.getTimeZone("GMT"))
var timeFormat=new SimpleDateFormat("yyyy-MM-dd'T'HH:mm:ssX")

def repo_names(x:String):String={
  try{
    if(x.contains("github.com/repos/")){
      var words=x.split("github.com/repos/")(1).split("\\?")(0)
      var path=words.split("/")
      if(path.size==1){
        return path(0)
      }
      else if(path.size>1){
        return path(0)+"/"+path(1)
      }
    }
    else{
      var words=x.split(" ")
      var index=words.indexOf("Repo")
      return words(index+1)
    }
  }
  return null
}
catch{
  case e:Exception=>return null
}
}

def invertIndex(rdd:org.apache.spark.rdd.RDD[schema],key: String) = {rdd.groupBy(r=>key match{
  case "debug_level" => r.debug_level
  case "timestamp" => r.timestamp
  case "download_id" => r.download_id
  case "retrieval_stage" => r.retrieval_stage
}))}

var file=sc.textFile("file:///home/sahil/Desktop/sdm/a1/ghorrent-logs.txt")
var parsedRDD=file.filter(x=>x.size>0).map(line=>line.split(" | - ",4))
parsedRDD=parsedRDD.map(line=> if(line.size==4)Array(line(0),line(1),line(2))+line(3).split(": ",2) else line)
var myRDD=parsedRDD.map(x=>try{schema(x(0),timeFormat.parse(x(1)),x(2),x(3),x(4))}catch{case
e:Exception=>schema(null,null,null,null,null)})

val solution1=myRDD.filter(x=>x.rest!=null && x.download_id=="ghorrent-22" && (x.rest.contains("Repo") ||
x.rest.contains("repos"))).map(x=>repo_names(x.rest)).filter(x=>x!=null).distinct().count()
println("Unique repos with download id ghortrent-22 without using inverted index:",solution1)

val invert_index=invertIndex(myRDD,"download_id")
val list_log=invert_index.filter(x=>x._1=="ghorrent-22").flatMap(x=>x._2).collect()
val solution2=sc.parallelize(list_log).filter(x=>(x.rest.contains("Repo") ||
x.rest.contains("repos"))).map(x=>repo_names(x.rest)).filter(x=>x!=null).distinct().count()
println("Unique repos with download id ghortrent-22 using inverted index:",solution2)
```