

CS 8001 Big Data

HW #4: Movie Ratings (10 points + 2 bonus points)

Jing Su

js929@mail.missouri.edu

In this exercise, you will write Python code on Spark to find and predict movie ratings. See reference [here](http://grouplens.org/datasets/movielens/). The dataset to be used is MovieLens 1M and 20M datasets at <http://grouplens.org/datasets/movielens/>.

Part I (10 points): Run Spark on MovieLens 1M datasets.

Spark Environment: AWS EMR (4 nodes), Windows 10

1. Basic Recommendations: find the top 20 movies with highest average ratings and more than 500 reviews. Print out the names of the movies, their average ratings, and their numbers of reviews, in descending order of their average ratings.

Code Structure:

- 1) Parsing the two files yields two RDDs and cache them. For each line in the ratings dataset, I create a tuple of (UserID, MovieID, Rating). For each line in the movies dataset, I create a tuple of (MovieID, Title).
- 2) From ratingsRDD create an RDD with tuples of the form (MovieID, Python iterable of Ratings for that MovieID).
- 3) Using movieIDsWithRatingsRDD and getCountsAndAverages() function, compute the number of ratings and average rating for each movie to yield tuples of the form (MovieID, (number of ratings, average rating)).
- 4) To moviesRDD, apply transformations that use movieIDsWithAvgRatingsRDD to get the movie names for movieIDsWithAvgRatingsRDD, yielding tuples of the form (average rating, movie name, number of ratings).
- 5) Apply a single RDD transformation to movieNameWithAvgRatingsRDD to limit the results to movies with ratings from more than 500 people. We then use the sortFunction() helper function to sort by the average rating to get the movies in descending order of their rating.

Execution Process:

- 1) Put the 1M dataset to HDFS.

```
[hadoop@ip-172-31-38-247 ~]$ hadoop fs -put ratings.dat /user/hadoop/ratings.dat  
[hadoop@ip-172-31-38-247 ~]$ hadoop fs -put movies.dat /user/hadoop/movies.dat
```

- 2) Use command: spark-submit RecommendMovie1M.py ratings.dat movies.dat

Result:

```
The top 20 movies with highest average ratings and more than 500 reviews: [(4.560509554140127, u'Seven Samurai (The Magnificent Seven) (Shichinin no samurai) (1954)', 628), (4.554557700942973, u'Shawshank Redemption, The (1994)', 2227), (4.524966261808367, u'Godfather, The (1972)', 2223), (4.52054794520548, u'Close Shave, A (1995)', 657), (4.517106001121705, u'Usual Suspects, The (1995)', 1783), (4.510416666666667, u'Schindler's List (1993)', 2304), (4.507936507936508, u'Wrong Trousers, The (1993)', 882), (4.477724741447892, u'Raiders of the Lost Ark (1981)', 2514), (4.476190476190476, u'Rear Window (1954)', 1050), (4.453694416583082, u'Star Wars: Episode IV - A New Hope (1977)', 2991), (4.4498902706656915, u'Dr. Strangelove or: How I Learned to Stop Worrying and Love the Bomb (1963)', 1367), (4.425646551724138, u'To Kill a Mockingbird (1962)', 928), (4.415607985480944, u'Double Indemnity (1944)', 551), (4.412822049131217, u'Casablanca (1942)', 1669), (4.406262708418057, u'Sixth Sense, The (1999)', 2459), (4.401925391095066, u'Lawrence of Arabia (1962)', 831), (4.395973154362416, u'Maltese Falcon, The (1941)', 1043), (4.390724637681159, u'One Flew Over the Cuckoo's Nest (1975)', 1725), (4.388888888888889, u'Citizen Kane (1941)', 1116), (4.386993603411514, u'Bridge on the River Kwai, The (1957)', 938)]
```

Python Code:

```
1  from __future__ import print_function
2  import sys
3  from pyspark import SparkContext
4  import re
5  import time
6  import math
7  from pyspark.mllib.recommendation import ALS
8
9  def get_ratings_tuple(entry):
10     """ Parse a line in the ratings dataset
11     Args:
12         entry (str): a line in the ratings dataset in the form of
13             UserID::MovieID::Rating::Timestamp
14     Returns:
15         tuple: (UserID, MovieID, Rating)
16     """
17     items = entry.split('::')
18     return int(items[0]), int(items[1]), float(items[2])
19
20 def get_movie_tuple(entry):
21     """ Parse a line in the movies dataset
22     Args:
23         entry (str): a line in the movies dataset in the form of
24             MovieID::Title::Genres
25     Returns:
26         tuple: (MovieID, Title)
27     """
28     items = entry.split('::')
29     return int(items[0]), items[1]
30
31 def sortFunction(tuple):
32     """ Construct the sort string (does not perform actual sorting)
33     Args:
34         tuple: (rating, MovieName)
35     Returns:
```

```

36         sortString: the value to sort with, 'rating MovieName'
37         """
38         key = unicode('%0.3f' % tuple[0])
39         value = tuple[1]
40         return (key + ' ' + value)
41
42     def getCountsAndAverages(IDandRatingsTuple):
43         """ Calculate average rating
44         Args:
45             IDandRatingsTuple: a single tuple of (MovieID, (Rating1, Rating2, Rating3, ...))
46         Returns:
47             tuple: a tuple of (MovieID, (number of ratings, averageRating))
48         """
49         list_tuples = [1.0*rating for rating in IDandRatingsTuple[1]]
50         return (IDandRatingsTuple[0], (len(list_tuples), sum(list_tuples)/len(list_tuples)))
51
52     if __name__ == "__main__":
53         start = time.time()
54         #Print error message when the command is wrong.
55         if len(sys.argv) != 3:
56             print("Usage: RecommendMovie.py ratings.dat movies.dat", file=sys.stderr)
57             exit(-1)
58         sc = SparkContext(appName="RecommendMovies")
59         #Create a SparkContext to store dat files.
60         numPartitions = 2
61         rawRatings = sc.textFile(sys.argv[1]).repartition(numPartitions)
62         rawMovies = sc.textFile(sys.argv[2])
63
64         #Part 1: Basic Recommendations
65         # Parsing the two files yields two RDDs and cache them.
66         ratingsRDD = rawRatings.map(get_ratings_tuple).cache()
67         moviesRDD = rawMovies.map(get_movie_tuple).cache()
68
69         # From ratingsRDD with tuples of (UserID, MovieID, Rating) create an RDD with tuples of
70         # the (MovieID, iterable of Ratings for that MovieID)
71
72         movieIDsWithRatingsRDD = ratingsRDD.map(lambda (userid, movieid, rating):
73             (movieid, rating)).groupByKey()
74
75         # Using 'movieIDsWithRatingsRDD', compute the number of ratings and average rating for
76         # each movie to yield tuples of the form (MovieID, (number of ratings, average rating))
77         movieIDsWithAvgRatingsRDD = movieIDsWithRatingsRDD.map(lambda rec: getCountsAndAverages(rec))
78
79         # To 'movieIDsWithAvgRatingsRDD', apply RDD transformations that use 'moviesRDD' to get the movie
80         # names for 'movieIDsWithAvgRatingsRDD', yielding tuples of the form
81         # (average rating, movie name, number of ratings)
82         movieNameWithAvgRatingsRDD = moviesRDD.join(movieIDsWithAvgRatingsRDD).map(lambda
83             (movieID, (movieName, (numRatings, avgRating))): (avgRating, movieName, numRatings))
84
85         # Apply an RDD transformation to 'movieNameWithAvgRatingsRDD' to limit the results to movies with
86         # ratings from more than 500 people. We then use the 'sortFunction()' helper function to sort by
87         # the average rating to get the movies in order of their rating (highest rating first)
88         movieLimitedAndSortedByRatingRDD = movieNameWithAvgRatingsRDD.filter(lambda (_, __, numRatings):
89             numRatings > 500).sortBy(sortFunction, False)
90         print ('The top 20 movies with highest average ratings and more than 500 reviews: %s\n' %
91             movieLimitedAndSortedByRatingRDD.take(20))

```

2. Collaborative filtering

- a. Break up the ratings dataset into three sets, training, validation, and test set as follows. Report the number of entries in each set.
 - 1) Sort the ratings data in ascending order of timestamp (past to present)
 - 2) Divide the data (80%-10%-10%) into
 - training set: the first 4/5 of the sorted data,
 - validation set: the next 1/10 of the sorted data,
 - test set: the last 1/10 of the sorted data.

Result:

```
Training: 800198, validation: 100030, test: 99981
```

- b. Use MLlib's alternating least squares method, ALS.train(), to train using the training set and validate using the validation set. Train 3 models (different ranks) and select the best model based on their validation performance. Report the ranks of your models, their training error (RMSE), training time in seconds, validation error (RMSE), and validation time in seconds, respectively. The model having the lowest validation error is your best model.

Result:

```
For rank 4 :

the training error is 0.865652272897

the training time is 10.5775930882

the validation error is 0.885625491895

the validation time is 1.11262798309

-----

For rank 8 :

the training error is 0.851508149701

the training time is 8.08406496048

the validation error is 0.884062583335

the validation time is 1.27135682106

-----

For rank 12 :

the training error is 0.843171246507

the training time is 7.05665397644

the validation error is 0.885762632323

the validation time is 1.05017805099

-----

The best training model was trained with rank 8
```

- c. Test your best model on the test set. Report the test error (RMSE) and test time in seconds.

Result:

```
The model had a RMSE on the test set of 0.885953442082
```

Code Structure:

- 1) Break up the ratingsRDD dataset into three pieces:
A training set (RDD), which we will use to train models;
A validation set (RDD), which we will use to choose the best model;
A test set (RDD), which we will use for our experiments.
- 2) Pick a set of model parameters. The most important parameter to ALS.train() is the rank, which is the number of rows in the Users matrix or the number of columns in the Movies matrix. We will train models with ranks of 4, 8, and 12 using the trainingRDD dataset.
- 3) Create a model using ALS.train(trainingRDD, rank, seed=seed, iterations=iterations, lambda_=regularizationParameter) with three parameters: an RDD consisting of tuples of the form (UserID, MovieID, rating) used to train the model, an integer rank (4, 8, or 12), a number of iterations to execute (we will use 5 for the iterations parameter), and a regularization coefficient (we will use 0.1 for the regularizationParameter).
- 4) For the prediction step, create an input RDD, validationForPredictRDD, consisting of (UserID, MovieID) pairs that we extract from validationRDD.
- 5) Using the model and validationForPredictRDD, we can predict rating values by calling model.predictAll() with the validationForPredictRDD dataset, where model is the model we generated with ALS.train(). predictAll accepts an RDD with each entry in the format (userID, movieID) and outputs an RDD with each entry in the format (userID, movieID, rating).
- 6) Evaluate the quality of the model by using the computeError() function to compute the error between the predicted ratings and the actual ratings in validationRDD.
- 7) Use the bestRank=8 to create a model for predicting the ratings for the test dataset and then we will compute the RMSE.

Execution Process:

- 1) Put the 20M dataset to HDFS.
- 2) Use command: spark-submit RecommendMovie20M.py ratings.csv

Python Code:

```

1  from __future__ import print_function
2  import sys
3  from pyspark import SparkContext
4  import re
5  import time
6  import math
7  from pyspark.mllib.recommendation import ALS
8
9  def get_ratings_tuple(entry):
10     """ Parse a line in the ratings dataset
11     Args:
12         entry (str): a line in the ratings dataset in the f
13             orm of UserID::MovieID::Rating::Timestamp
14     Returns:
15         tuple: (UserID, MovieID, Rating, Timestamp)
16     """
17     items = entry.split(',')
18     return int(items[0]), int(items[1]), float(items[2]), long(items[3])
19
20
21 def computeError(predictedRDD, actualRDD):
22     """ Compute the root mean squared error between predicted and actual
23     Args:
24         predictedRDD: predicted ratings for each movie and each user
25             where each entry is in the form (UserID, MovieID, Rating)
26         actualRDD: actual ratings where each entry is in the form
27             (UserID, MovieID, Rating)
28     Returns:
29         RSME (float): computed RSME value
30     """
31     # Transform predictedRDD into the tuples of the form ((UserID, MovieID), Rating)
32     predictedReformattedRDD = predictedRDD.map(lambda (UserID, MovieID, Rating):
33         ((UserID, MovieID), Rating))
34
35     # Transform actualRDD into the tuples of the form ((UserID, MovieID), Rating)
36     actualReformattedRDD = actualRDD.map(lambda (UserID, MovieID, Rating):
37         ((UserID, MovieID), Rating))
38
39     # Compute the squared error for each matching entry (i.e., the same (User ID,
40     # Movie ID) in each RDD) in the reformatted RDDs using RDD transformations
41     # do not use collect()
42     squaredErrorsRDD = (predictedReformattedRDD
43         .join(actualReformattedRDD)
44         .map(lambda ((UserID, MovieID), (PredRating, ActRating)):
45             (PredRating - ActRating)**2))
46
47     # Compute the total squared error - do not use collect()
48     totalError = squaredErrorsRDD.reduce(lambda a,b: a+b)
49
50     # Count the number of entries for which you computed the total squared error
51     numRatings = squaredErrorsRDD.count()
52
53     # Using the total squared error and the number of entries, compute the RSME
54     return math.sqrt(1.0*totalError/numRatings)
55
56
57 if __name__ == "__main__":
58     #Print error message when the command is wrong.
59     if len(sys.argv) != 2:
60         print("Usage: RecommendMovie20M.py ratings.csv", file=sys.stderr)
61         exit(-1)
62     sc = SparkContext(appName="RecommendMovies")
63     #Create a SparkContext to store dat files.
64     numPartitions = 2
65     rawRatings = sc.textFile(sys.argv[1])
66     #extract header from rating.csv
67     header = rawRatings.first()
68     print ("-----")
69     print ('\nHeader: %s\n' % header)
70     #filter out header

```



```

71 rawRatings = rawRatings.filter(lambda x:x !=header).repartition(numPartitions)
72
73 ratingsRDD = rawRatings.map(get_ratings_tuple).cache()
74 print ("-----")
75 print ('\nRatings: %s\n' % ratingsRDD.take(3))
76
77 #Part 2: Collaborative Filtering
78 #Sort the ratings data in ascending order of timestamp (past to present)
79 sortedRatingsRDD = ratingsRDD.sortBy(lambda r: r[3])
80 print ("-----")
81 print ('\nSorted Ratings: %s\n' % sortedRatingsRDD.take(3))
82 noTimeSortedRatingsRDD = sortedRatingsRDD.map(lambda (UserID, MovieID,
83 Rating, Timestamp): (UserID, MovieID, Rating))
84
85 #The training set contains the first 80% of the original dataset.
86 #The validation set contains the next 10% of the original dataset.
87 #The test set contains the last 10% of the original dataset.
88 #To randomly split the dataset into the multiple groups, we can use the pySpark
89 #randomSplit() transformation.
90 #randomSplit() takes a set of splits and seed and returns multiple RDDs.
91 trainingRDD, validationRDD, testRDD = noTimeSortedRatingsRDD
92                                     .randomSplit([8, 1, 1], seed=0L)
93 print ('\nTraining: %s, validation: %s, test: %s\n' %
94       (trainingRDD.count(), validationRDD.count(), testRDD.count()))
95
96 # For the prediction step, create an input RDD, validationForPredictRDD,
97 # consisting of (UserID, MovieID) pairs that we extract from validationRDD.
98 validationForPredictRDD = validationRDD.map(lambda (UserID, MovieID, Rating)
99 : (UserID, MovieID))
100 trainingForPredictRDD = trainingRDD.map(lambda (UserID, MovieID, Rating)
101 : (UserID, MovieID))
102
103 seed = 5L
104 iterations = 5
105 regularizationParameter = 0.1

```

```

106 ranks = [4, 8, 12]
107 training_errors = [0, 0, 0] #training error
108 validation_errors = [0, 0, 0]
109 err = 0
110 tolerance = 0.03
111 training_time = [0,0,0];
112 validation_time = [0,0,0];
113
114 minError = float('inf')
115 bestRank = -1
116 bestIteration = -1
117
118 for rank in ranks:
119     train_start = time.time();
120     #Training
121     model = ALS.train(trainingRDD, rank, seed=seed, iterations=iterations,
122                       lambda=regularizationParameter)
123     training_time[err] = time.time()-train_start
124     #Training RMSE
125     tra_predictedRatingsRDD = model.predictAll(trainingForPredictRDD)
126     #Use the Root Mean Square Error (RMSE) or Root Mean Square Deviation
127     #(RMSD) to compute the error of each model.
128     tra_error = computeError(tra_predictedRatingsRDD, trainingRDD)
129     training_errors[err] = tra_error
130
131     validation_start = time.time()
132     #Validation
133     val_predictedRatingsRDD = model.predictAll(validationForPredictRDD)
134     validation_time[err] = time.time()-validation_start
135     #Validation RMSE
136     val_error = computeError(val_predictedRatingsRDD, validationRDD)
137     validation_errors[err] = val_error
138     if val_error < minError:
139         minError = val_error
140         bestRank = rank

```

```

141         err += 1
142
143     # Use the bestRank=8 to create a model for predicting the ratings
144     # for the test dataset and then we will compute the RMSE.
145     myModel = ALS.train(trainingRDD, 8, seed=seed, iterations=iterations,
146         lambda_=regularizationParameter)
147     testForPredictingRDD = testRDD.map(lambda (UserID, MovieID, Rating)
148         : (UserID, MovieID))
149     test_start = time.time()
150     predictedTestRDD = myModel.predictAll(testForPredictingRDD)
151     # Test error
152     testRMSE = computeError(testRDD, predictedTestRDD)
153     # Test time
154     test_time = time.time() - test_start
155
156     i=0;
157     for rank in ranks:
158         print ("-----")
159         print ('\nFor rank %s : \n' % rank)
160         print ('the training error is %s\n' % training_errors[i])
161         print ('the training time is %s\n' % training_time[i])
162         print ('the validation error is %s\n' % validation_errors[i])
163         print ('the validation time is %s\n' % validation_time[i])
164         i += 1
165
166     print ("-----")
167     print ('\nThe best training model was trained with rank %s\n' % bestRank)
168     print ('\nThe model had a RMSE on the test set of %s\n' % testRMSE)
169     print ('\nThe validation time is %s\n' % test_time)

```

Part II (2 bonus points): Run Spark on MovieLens 20M Dataset.

Perform the same tasks as in Part I. Compare the results with those in Part I. Fill in the following table.

Best model on 1M dataset			
	training	validation	test
RMSE	0.8515	0.8841	0.8860
Time (sec)	8.0841	1.2714	1.2673
Best model on 20 M dataset			
	training	validation	test
RMSE	0.7821	0.8130	0.8208
Time (sec)	76.4633	25.3092	26.5362