Blog Archives

APR 1ST, 2013

Movie Recommendations and More With Spark

This post is inspired by <u>Edwin Chen's post on Scalding</u>. I encourage you to first read that post! The Spark code is adapted from his Scalding code and is available in full here.

As outlined in Ed's post, Scalding is a Scala DSL for Hadoop MapReduce that makes it easier, more natural and more concise to write MapReduce workflows. The Scala code ultimately compiles down to MapReduce jobs via Cascading.

Spark The **Spark Project** is a cluster computing framework that emphasizes low-latency job execution and in-memory caching to provide speed. It can run up to 100x faster than Hadoop MapReduce (when all the data is cached in memory) as a result. It is written in Scala, but also has Java and

Python APIs. It is fully compatible with HDFS and any Hadoop InputFormat/OutputFormat, but is independent of Hadoop MapReduce. The Spark API bears many similarities to Scalding, providing a way to write natural Scala code instead of Mappers and Reducers. Taking Ed's example:

Scalding

```
tweets.map('tweet -> 'length) { tweet : String => tweet.size }.groupBy('length) { _.size }
                                              Spark
     tweets.groupBy(tweet : String => tweet.size).map(pair => (pair._1, pair._2.size))
Movie Similarities
I've recently been experimenting a lot with Spark, and thought it would be interesting to
compare Ed's approach to computing movie similarities in Scalding with Spark. So I've ported
```

his Scalding code over to Spark and we'll compare the two as we go along. For a basic introduction to Spark's API see the **Spark Quickstart**.

val numRaters =

val ratingsWithSize =

ratings

Firstly, we read the ratings from a file. Since I don't have access to a nice Twitter tweet datasource, I used the MovieLens 100k rating dataset. The training set ratings are in a file called ua.base, while the movie item data is in u.item. Scalding

val INPUT FILENAME = "data/ratings.tsv"

* Let's also keep track of the total number of people who rated each movie.

val ratings = Tsv(INPUT_FILENAME, ('user, 'movie, 'rating))

.groupBy('movie) { _.size }.rename('size -> 'numRaters)

```
ratings.joinWithSmaller('movie -> 'movie, numRaters)
                                          Spark
val TRAIN_FILENAME = "ua.base"
val MOVIES_FILENAME = "u.item"
val sc = new SparkContext(master, "MovieSimilarities")
val ratings = sc.textFile(TRAIN_FILENAME)
  .map(line => {
    val fields = line.split("\t")
    (fields(0).toInt, fields(1).toInt, fields(2).toInt)
})
// get num raters per movie, keyed on movie id
val numRatersPerMovie = ratings
  .groupBy(tup => tup._2)
  .map(grouped => (grouped._1, grouped._2.size))
val ratingsWithSize = ratings
  .groupBy(tup => tup._2)
  .join(numRatersPerMovie)
  .flatMap(joined => {
    joined._2._1.map(f => (f._1, f._2, f._3, joined._2._2))
})
```

In order to determine how similar two movies are to each other, we must (as per Ed's post again):

Similarly to Scalding's Tsv method, which reads a TSV file from HDFS, Spark's sc.textFile

Also, Spark's API for joins is a little lower-level than Scalding's, hence we have to groupBy first

method reads a text file from HDFS. However it's up to us to specify how to split the fields.

and transform after the join with a flatMap operation to get the fields we want. Scalding

actually does something similar under the hood of joinWithSmaller.

• For every pair of movies A and B, find all the people who rated both A and B. • Use these ratings to form a Movie A vector and a Movie B vector. • Calculate the correlation between these two vectors. • Whenever someone watches a movie, you can then recommend the movies most

correlated with it.

val ratingPairs =

ratingPairs

ratings : (Double, Double) =>

.groupBy('movie, 'movie2) { group =>

.sum('ratingProd -> 'dotProduct)

ratingsWithSize

.joinWithSmaller('user -> 'user2, ratings2)

Computing similarity

This is item-based collaborative filtering. So let's compute the first two steps above: Scalding

* So first make a dummy copy of the ratings that we can join against. val ratings2 = ratingsWithSize .rename(('user, 'movie, 'rating, 'numRaters) -> ('user2, 'movie2, 'rating2, 'numRaters2))

.filter('movie, 'movie2) { movies : (String, String) => movies._1 < movies._2 }</pre>

Spark

// By grouping on ('movie, 'movie2), we can now get all the people who rated any pair of movie

.project('movie, 'rating, 'numRaters, 'movie2, 'rating2, 'numRaters2)

```
// dummy copy of ratings for self join
     val ratings2 = ratingsWithSize.keyBy(tup => tup._1)
     val ratingPairs =
       ratingsWithSize
       .keyBy(tup => tup._1)
       .join(ratings2)
       .filter(f => f._2._1._2 < f._2._2._2)
Notice how similar the APIs are with respect to the functional operations like filter - they
each simply take a Scala closure. We then compute the various vector metrics for each ratings
vector (size, dot-product, norm etc). We'll use these to compute the various similarity metrics
between pairs of movies.
                                              Scalding
     val vectorCalcs =
```

.sum('rating -> 'ratingSum) .sum('rating2 -> 'rating2Sum) .sum('ratingSq -> 'ratingNormSq) .sum('rating2Sq -> 'rating2NormSq) .max('numRaters) // Just an easy way to make sure the numRaters field stays. .max('numRaters2) // All of these operations chain together like in a builder object.

Spark

(ratings._1 * ratings._2, math.pow(ratings._1, 2), math.pow(ratings._2, 2))

.map(('rating, 'rating2) -> ('ratingProd, 'ratingSq, 'rating2Sq)) {

```
ratingPairs
     .map(data => {
      val key = (data._2._1._2, data._2._2._2)
       val stats =
        (data._2._1._3 * data._2._2._3, // rating 1 * rating 2
          data._2._1._3, // rating movie 1
          data._2._2._3,
          math.pow(data._2._1._3, 2), // square of rating movie 1
          math.pow(data._2._2._3, 2), // square of rating movie 2
          data._2._1._4, // number of raters movie 1
          data._2._2._4)
                                    // number of raters movie 2
      (key, stats)
     .groupByKey()
     .map(data => {
      val key = data._1
      val vals = data._2
      val dotProduct = vals.map(f => f._1).sum
      val ratingSum = vals.map(f => f._2).sum
      val rating2Sum = vals.map(f => f._3).sum
      val ratingSq = vals.map(f => f._4).sum
      val rating2Sq = vals.map(f => f._5).sum
      val numRaters = vals.map(f => f._6).max
      val numRaters2 = vals.map(f => f._7).max
      (key, (size, dotProduct, ratingSum, rating2Sum, ratingSq, rating2Sq, numRaters, numRaters2))
     })
Similarity metrics
```

For each movie pair we compute correlation, regularized correlation, cosine similarity and Jaccard

Scalding

.map(('size, 'dotProduct, 'ratingSum, 'rating2Sum, 'ratingNormSq, 'rating2NormSq, 'numRate

val (size, dotProduct, ratingSum, rating2Sum, ratingNormSq, rating2NormSq, numRaters,

val corr = correlation(size, dotProduct, ratingSum, rating2Sum, ratingNormSq, rating2Nor val regCorr = regularizedCorrelation(size, dotProduct, ratingSum, rating2Sum, ratingNorm val cosSim = cosineSimilarity(dotProduct, math.sqrt(ratingNormSq), math.sqrt(rating2Norm

('correlation, 'regularizedCorrelation, 'cosineSimilarity, 'jaccardSimilarity)) {

fields: (Double, Double, Double, Double, Double, Double, Double) =>

Spark

val jaccard = jaccardSimilarity(size, numRaters, numRaters2)

similarity (see Ed's post and the code for full details).

(corr, regCorr, cosSim, jaccard)

val PRIOR_COUNT = 10

val similarities =

vectorCalcs

val PRIOR_CORRELATION = 0

val PRIOR_COUNT = 10 val PRIOR_CORRELATION = 0

val similarities = vectorCalcs .map(fields => { val key = fields._1 val (size, dotProduct, ratingSum, rating2Sum, ratingNormSq, rating2NormSq, numRaters, numR

val regCorr = regularizedCorrelation(size, dotProduct, ratingSum, rating2Sum,

ratingNormSq, rating2NormSq, PRIOR_COUNT, PRIOR_CORRELATION)

val jaccard = jaccardSimilarity(size, numRaters, numRaters2)

val corr = correlation(size, dotProduct, ratingSum, rating2Sum, ratingNormSq, rating2NormS

val cosSim = cosineSimilarity(dotProduct, scala.math.sqrt(ratingNormSq), scala.math.sqrt(r

(key, (corr, regCorr, cosSim, jaccard)) The nice thing here is that, once the raw input metrics themselves are computed, we can use the exact same functions from the Scalding example to compute the similarity metrics as can be seen above - I simply copy-and-pasted Ed's correlation, regularizedCorrelation, cosineSimilarity and jaccardSimilarity functions! Some results So, what do the results look like after putting all of this together? Since I used a different input data source we won't get the same results, but we'd hope that most of them would make sense. Similarly to Ed's results, I found that using raw correlation resulted in sub-optimal similarities (at least from eye-balling and "sense checking"), since some movie pairs have very few common raters (many had just 1 rater in common). I also found that cosine similarity didn't do so well on a "sense check" basis either, which was

somewhat surprising since this is usually the standard similarity metric for collaborative

have messed up the calculation somewhere (if you spot an error please let me know).

In any case, here are the top 10 movies most similar to Die Hard (1998), ranked by

Movie 2

Die Hard: With a Vengeance

(1995)

Die Hard 2 (1990)

Bananas (1971)

Good, The Bad and The Ugly, The

(1966)

Hunt for Red October, The (1990)

City Slickers II: The Legend of

Curly's Gold (1994)

Looks fairly reasonable! And here are the 10 most similar to Il Postino:

Movie 2

Bottle Rocket (1996)

filtering. This seems to be due to a lot of movies having cosine similarity of 1.0, so perhaps I

Correlation

0.5413

0.4868

0.5516

0.4608

0.4260

0.5349

Reg

Correlation

0.4967

Cosine

0.9692

0.9687

0.9745

0.9743

0.9721

0.9506

Cosine

Similarity

0.9855

Cosine

0.9888

0.9851

0.9816

0.9840

0.9912

0.9995

Reg

0.7168

0.6539

0.4917

0.4118

-0.3392

-0.3310

-0.3175

-0.9045

-0.8039

-0.6062

0.8781

0.8670

0.8998

0.0121

0.0141

0.0220

0.7419

0.6714

0.5074

1.0000

Correlation Correlation Similarity Similarity

Jaccard

0.5306

0.6708

0.5607

0.0442

0.0181

0.0141

Correlation Similarity Similarity

Reg

0.4946

0.4469

0.4390

0.4032

0.3944

0.3903

Jaccard

0.4015

0.4088

0.1618

0.2518

0.4098

0.1116

Jaccard

Similarity

0.0699

Die Hard Grease 2 (1982) 0.6502 0.3901 0.9449 0.0647 (1988)Die Hard Star Trek: The Wrath of Khan 0.4160 0.3881 0.9675 0.4441 (1988)(1982)Die Hard Sphere (1998) 0.7722 0.3861 0.9893 0.0403 (1988)Die Hard Field of Dreams (1989) 0.4126 0.3774 0.9630 0.3375 (1988)

Correlation

0.8789

(1994)Postino, Il Looking for Richard 0.7112 0.4818 0.9820 0.1123 (1994)(1996)Postino, Il Ridicule (1996) 0.4780 0.9759 0.1561 0.6550 (1994)Postino, Il When We Were Kings 0.7581 0.4773 0.0929 0.9888 (1994)(1996)Postino, Il Mother Night (1996) 0.8802 0.4611 0.9848 0.0643 (1994)Postino, Il Kiss Me, Guido (1997) 0.9759 0.9974 0.0452 0.4337 (1994)Postino, Il Blue in the Face (1995) 0.6372 0.9585 0.4317 0.1148 (1994)Postino, Il Othello (1995) 0.5875 0.4287 0.9774 0.1330 (1994)Postino, Il English Patient, The 0.9603 0.4586 0.4210 0.2494 (1994)(1996)Postino, Il Mediterraneo (1991) 0.9879 0.6200 0.4200 0.1235 (1994)

Star Wars Meet John Doe (1941) 0.6396 0.4397 (1977)Star Wars Love in the Afternoon (1957) 0.9234 0.4374 (1977)

Movie 2

Empire Strikes Back, The (1980)

Return of the Jedi (1983)

Raiders of the Lost Ark (1981)

Man of the Year (1995)

Star Wars (1977)	When We Were Kings (1996)	0.5278	0.4021	0.9737	0.0637
Star Wars (1977)	Cry, the Beloved Country (1995)	0.7001	0.3957	0.9763	0.0257
Star Wars (1977)	To Be or Not to Be (1942)	0.6999	0.3956	0.9847	0.0261
Star Wars (1977)	Haunted World of Edward D. Wood Jr., The (1995)	0.6891	0.3895	0.9758	0.0262
Movie 1	Movie 2	Correlation	Reg Correlation	Cosine Similarity	Jaccard Similarit
(1977)	Movie 2 Fathers' Day (1997)	Correlation -0.6625	[일시] : 나는 아이 (이 100 100 100 100 100 100 100 100 100 10		
Star Wars (1977)			Correlation	Similarity	Similarit
Star Wars (1977) Star Wars (1977)	Fathers' Day (1997)	-0.6625	Correlation -0.4417	Similarity 0.9074	Similarit 0.0397
Star Wars (1977) Star Wars (1977) Star Wars (1977)	Fathers' Day (1997) Jason's Lyric (1994)	-0.6625 -0.9661	-0.4417 -0.3978	0.9074 0.8110	0.0397 0.0141
Star Wars (1977) Star Wars (1977) Star Wars (1977) Star Wars	Fathers' Day (1997) Jason's Lyric (1994) Lightning Jack (1994)	-0.6625 -0.9661 -0.7906	-0.4417 -0.3978 -0.3953	0.9074 0.8110 0.9361	0.0397 0.0141 0.0202

I'll leave it to you to decide on the accuracy.

What Happened Was... (1994)

Female Perversions (1996)

Celtic Pride (1996)

Conclusion and Next Steps Hopefully this gives a taste for Spark and how it can be used in a very similar manner to Scalding

fields as Scala Symbols, e.g.

Copyright © 2013 - Nick Pentreath - Powered by Octopress

Tweet

Star wars Poison Ivy II (1995) -0.7443 -0.37220.7169 0.0201 (1977)In the Realm of the Senses (Ai no -0.8090 -0.3596 0.8108 0.0162 (1977)corrida) (1976) Star Wars

variables and accumulators); not to mention interactive analysis via the Scala/Spark console, and a Java and Python API! Check out the documentation, tutorials and examples here. One issue that is apparent from the above code snippets is that Scalding's API is somewhat

cleaner when doing complex field manipulations and joins, due to the ability to have named

capabilities, low-latency execution and other distributed-memory primitives (such as broadcast

and MapReduce - with all the advantages of HDFS compatability, in-memory caching

The lack of named fields in Spark's API does lead to some messy tuple-unpacking and makes keeping track of which fields are which more complex. This could be an interesting potential addition to Spark.

tweets.map('tweet -> 'length) { tweet : String => tweet.size }

Finally, please do let me know if you find any issues or errors. And thanks to the Spark team for a fantastic project! Posted by Nick Pentreath • Apr 1st, 2013 • Scala, Spark, collaborative filtering, recommendations

Movie recommendations and more with Spark

Recent Posts

GitHub Repos elasticsearch-vector-scoring

Glint

mlnick.github.com

Search

Score documents with pure dot product / cosine similarity with ES beijing-meetup-2016 Jupyter Notebooks for Apache Spark ML Meetup - Beijing Nov 2016 glint-fm

Factorization Machines on Spark and

hive-udf Approximate cardinality estimation with HyperLogLog, as a Hive function @MLnick on GitHub

regularized correlation:

Movie 1

Die Hard

(1988)

Movie 1

Postino, Il

How about Star Wars?

Movie 1

Star Wars

(1977)

Star Wars

(1977)

Star Wars

(1977)

Star Wars

(1977)

(1977)

Star Wars

(1977)

Star Wars

(1977)