

Recommender Systems at Scale

```
{"name": "Eoin Hurrell",  
  "twitter": "@eoinhurrell",  
  "github": "eoinhurrell",  
  "email": "eoin.hurrell@gmail.com"}
```

Slides About Recommender Systems Being Everywhere are Everywhere

The value of recommendations

- Netflix: 2/3 of the movies watched recommended
- Google News: recommendations get 38% more clickthrough
- Amazon: 35% sales from recommendations
- Choicestream: 28% of the people who listen to more music if they found what they

- Omnipresent nowadays
- Important for user experience and sales

amazon.com

Recommended for You

Amazon.com has new recommendations for you based on [items](#) you purchased or told us you own.



[Google Apps Deciphered: Compute in the Cloud to Streamline Your Desktop](#)



[Google Apps Administrator Guide: A Private-Label Web Workspace](#)



[Googlepedia: The Ultimate Google Resource \(3rd Edition\)](#)

NETFLIX

Xavier Amatriain – July 2014 – Rec

<http://technocalifornia.blogspot.ie/2014/08/introduction-to-recommender-systems-4.html>

<https://www.youtube.com/watch?v=CmHPxASIj6Y>

The Recommendation Task

We have a catalogue of items. Given what we know about the user what will they like best?

Companies want to get a prediction of the subjective rating a user would have of every item in their collection.

Collaborative Filtering

We can use what we know about what others like to predict what a user will like, but hasn't seen.

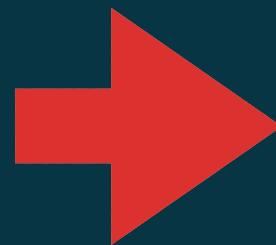
"You're a lot like Alice, and you haven't seen Alice's favourite film, you might like it"

Matrix Factorization

Matrix Factorization and Latent Factors

Matrix (Rating Database)

ITEM	U1	U2	...	UN
I1	3		...	
I2		4	...	
...
IN	2		...	



User Factors

Item Factors

Matrix Factorization

Latent Factors are variables contributing to the rating. They are often opaque, but think of them as:

User Factors

E.g. "I rarely watch films"

Item Factors

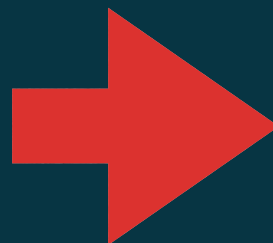
E.g. "This film has lots of explosions."

Matrix Factorization

Latent Factors can rebuild the matrix

User Factors

Item Factors



Matrix with Predictions

ITEM	U1	U2	...	UN
I1	3	1	...	2
I2	3	4	...	4
...
IN	2	5	...	5

Matrix Factorization

Problems:

We only partially know the User and Item Latent Factors.

Estimating Latent Factors can be computationally expensive and difficult to scale.

Alternating Least Squares

A method of matrix factorization that can easily be computed in parallel

<http://bugra.github.io/work/notes/2014-04-19/alternating-least-squares-method-for-collaborative-filtering/>

Alternating Least Squares

- * Estimate Item Factor matrix using current User Factors
- * Estimate User Factor matrix using current Item Factors
- * Iterate to convergence (either User or Item matrix no longer changes)

<http://bugra.github.io/work/notes/2014-04-19/alternating-least-squares-method-for-collaborative-filtering/>

From Theory to Practice

How do we implement this in a way that can deal with
lots of data?

Apache Spark

What is Spark?

"Apache Spark™ is a fast and general engine for large-scale data processing."

Spark **runs on Hadoop**, Mesos, **standalone**, or in the cloud. It can access diverse data sources including HDFS, Cassandra, HBase, and S3.

You can run Spark using its standalone cluster mode, on EC2, on Hadoop YARN, or on Apache Mesos. Access data in HDFS, Cassandra, HBase, Hive, Tachyon, and any Hadoop data source.

Using Spark to Scale

Spark is fast and has support for exactly this type of operation, with libraries such as MLlib

The Approach

Start small, with a small amount of data locally.
Then scale up to run in the cloud.

We can do both with Spark

The Code: Python

Python is sort of a second-class citizen in Spark.

I chose Python because it's Data Science Esperanto

The Data: MovieLens

<http://grouplens.org/datasets/movielens/>

Dataset	Size	Users	Movies	Tags
10K	5 MB	1,000	1,700	–
1M	6 MB	6,000	4,000	–
10M	63 MB	72,000	10,000	100,000
20M	132 MB	138,000	27,000	465,000
Latest (21M)	144 MB	230,000	30,000	510,000


```

1 """ Example PySpark ALS application
2 """
3 from pyspark import SparkContext # pylint: disable=import-error
4 from pyspark.mllib.recommendation import ALS, Rating # pylint: disable=import-error
5
6
7 def parse_rating(line):
8     """ Parse Movielens Rating line to Rating object.
9         UserID::MovieID::Rating::Timestamp
10     """
11     line = line.split(':')
12     return Rating(int(line[0]), int(line[1]), float(line[2]))
13
14
15 def parse_movie(line):
16     """ Parse Movielens Movie line to Movie tuple.
17         MovieID::Title::Genres
18     """
19     line = line.split(':')
20     return (line[0], line[1])
21
22
23 def main():
24     """ Train and evaluate an ALS recommender.
25     """
26     # Set up environment
27     sc = SparkContext("local[*]", "RecSys")
28
29     # Load and parse the data
30     data = sc.textFile("./data/ratings.dat")
31     ratings = data.map(parse_rating)
32
33     # Build the recommendation model using Alternating Least Squares
34     rank = 10
35     iterations = 50
36     model = ALS.train(ratings, rank, iterations)
37
38     movies = sc.textFile("./data/movies.dat")\
39         .map(parse_movie)
40     # Evaluate the model on training data
41     testdata = ratings.map(lambda p: (p[0], p[1]))
42     predictions = model.predictAll(testdata)\
43         .map(lambda r: ((r[0], r[1]), r[2]))
44     rates_and_preds = ratings.map(lambda r: ((r[0], r[1]), r[2]))\
45         .join(predictions)
46     MSE = rates_and_preds.map(lambda r: (r[1][0] - r[1][1])**2).mean()
47     print("Mean Squared Error = " + str(MSE))
48

```

```

1 """ Example PySpark ALS application
2 """
3 from pyspark import SparkContext # pylint: disable=import-error
4 from pyspark.mllib.recommendation import ALS, Rating # pylint: disable=import-error
5
6
7 def parse_rating(line):
8     """ Parse Movielens Rating line to Rating object.
9         UserID::MovieID::Rating::Timestamp
10    """
11    line = line.split(':')
12    return Rating(int(line[0]), int(line[1]), float(line[2]))
13
14
15 def parse_movie(line):
16     """ Parse Movielens Movie line to Movie tuple.
17         MovieID::Title::Genres
18    """
19    line = line.split(':')
20    return (line[0], line[1])
21
22
23 def main():
24     """ Train and evaluate an ALS recommender.
25    """
26    # Set up environment
27    sc = SparkContext("local[*]", "RecSys")
28
29    # Load and parse the data
30    data = sc.textFile("./data/ratings.dat")
31    ratings = data.map(parse_rating)
32
33    # Build the recommendation model using Alternating Least Squares
34    rank = 10
35    iterations = 50
36    model = ALS.train(ratings, rank, iterations)
37
38    movies = sc.textFile("./data/movies.dat")\
39        .map(parse_movie)
40
41    # Evaluate the model on training data
42    testdata = ratings.map(lambda p: (p[0], p[1]))
43    predictions = model.predictAll(testdata)\
44        .map(lambda r: ((r[0], r[1]), r[2]))
45    rates_and_preds = ratings.map(lambda r: ((r[0], r[1]), r[2]))\
46        .join(predictions)
47    MSE = rates_and_preds.map(lambda r: (r[1][0] - r[1][1])**2).mean()
48    print("Mean Squared Error = " + str(MSE))

```

Performing 50
iterations of ALS

```

1 """ Example PySpark ALS application
2 """
3 from pyspark import SparkContext # pylint: disable=import-error
4 from pyspark.mllib.recommendation import ALS, Rating # pylint: disable=import-error
5
6
7 def parse_rating(line):
8     """ Parse Movielens Rating line to Rating object.
9         UserID::MovieID::Rating::Timestamp
10    """
11    line = line.split(':')
12    return Rating(int(line[0]), int(line[1]), float(line[2]))
13
14
15 def parse_movie(line):
16     """ Parse Movielens Movie line to Movie tuple.
17         MovieID::Title::Genres
18    """
19    line = line.split(':')
20    return (line[0], line[1])
21
22
23 def main():
24     """ Train and evaluate an ALS recommender.
25    """
26    # Set up environment
27    sc = SparkContext("local[*]", "RecSys")
28
29    # Load and parse the data
30    data = sc.textFile("./data/ratings.dat")
31    ratings = data.map(parse_rating)
32
33    # Build the recommendation model using Alternating Least Squares
34    rank = 10
35    iterations = 50
36    model = ALS.train(ratings, rank, iterations)
37
38    movies = sc.textFile("./data/movies.dat")\
39        .map(parse_movie)
40    # Evaluate the model on training data
41    testdata = ratings.map(lambda p: (p[0], p[1]))
42    predictions = model.predictAll(testdata)\
43        .map(lambda r: ((r[0], r[1]), r[2]))
44    rates_and_preds = ratings.map(lambda r: ((r[0], r[1]), r[2]))\
45        .join(predictions)
46    MSE = rates_and_preds.map(lambda r: (r[1][0] - r[1][1])**2).mean()
47    print("Mean Squared Error = " + str(MSE))
48

```

Evaluating accuracy

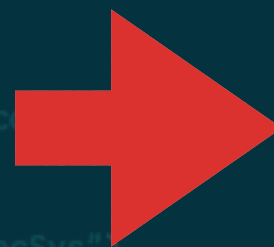
```

1 """ Example PySpark ALS application
2 """
3 from pyspark import SparkContext # pylint: disable=import-error
4 from pyspark.mllib.recommendation import ALS, Rating # pylint: disable=import-error
5
6
7 def parse_rating(line):
8     """ Parse Movielens Rating line to Rating object.
9         UserID::MovieID::Rating::Timestamp
10    """
11    line = line.split(':')
12    return Rating(int(line[0]), int(line[1]), float(line[2]))
13
14
15 def parse_movie(line):
16     """ Parse Movielens Movie line to Movie tuple.
17         MovieID::Title::Genres
18    """
19    line = line.split(':')
20    return (line[0], line[1])
21
22 def main():
23     """ Train and evaluate an ALS recommendation model.
24    """
25    # Set up environment
26    sc = SparkContext("local[*]", "RecSys")
27
28    # Load and parse the data
29    data = sc.textFile("./data/ratings.dat")
30    ratings = data.map(parse_rating)
31
32    # Build the recommendation model using Alternating Least Squares
33    rank = 10
34    iterations = 50
35    model = ALS.train(ratings, rank, iterations)
36
37    # Evaluate the model on training data
38    testdata = ratings.map(lambda p: (p[0], p[1]))
39    predictions = model.predictAll(testdata)\
40        .map(lambda r: ((r[0], r[1]), r[2]))
41    rates_and_preds = ratings.map(lambda r: ((r[0], r[1]), r[2]))\
42        .join(predictions)
43    MSE = rates_and_preds.map(lambda r: (r[1][0] - r[1][1])**2).mean()
44    print("Mean Squared Error = " + str(MSE))
45
46 if __name__ == '__main__':
47     main()
48

```

User 1 likes:

Sympathy for Lady Vengeance (2005),
Boys Don't Cry (2000),
District B13 (2004),
BloodRayne (2005),
Last Holiday (2006),
Underworld: Evolution (2006)



User 1 would like:

Volver (2006), predicted rating: 4.5

github.com/eoinhurrell/python-spark-recsys

Running this at scale

Spark cluster computing runs from 1 to N nodes.

Spark on AWS

AWS EC2 - `spark-ec2` commandline tool allows you to do this trivially.

AWS EMR - This can be done through the AWS Console. Even though the UI mentions 'application JARs' you can use Python code.

Spark on AWS

Cost: >\$10 to predict items for 1,000 users.
Amazon make 35% of their sales from recommendations [1].
If 5% of those users buy something
18 of those sales will be from recommendations.

If you make more than \$0.56 per sale on average recommenders
are worth investigating.

[1] <http://technocalifornia.blogspot.ie/2014/08/introduction-to-recommender-systems-4.html>

Conclusion

- * ALS Matrix Factorization is efficient for recommendation
- * Spark has built-in support for such high-level operations
- * Working with Spark is easy at any scale, with no code changes

github.com/eoinhurrell/python-spark-recsys

```
{"name": "Eoin Hurrell",  
  "twitter": "@eoinhurrell",  
  "github": "eoinhurrell",  
  "email": "eoin.hurrell@gmail.com"}
```


Questions?

eoin.hurrell@gmail.com

github.com/eoinhurrell/python-spark-recsys

```
{"name": "Eoin Hurrell",  
  "twitter": "@eoinhurrell",  
  "github": "eoinhurrell",  
  "email": "eoin.hurrell@gmail.com"}
```