



# Spark SQL & Machine Learning A Practical Demonstration

Craig Warman



© 2016 MapR Technologies



## What You'll Learn From This Presentation

- A brief Spark background/overview
- Spark SQL architecture
- How to build load and query a dataset with Spark SQL
- Basic machine learning with Spark MLlib
- Leveraging Spark in the real world



- Spark Background/Overview
  - Brief Spark background
  - The Spark+Hadoop team
  - Spark's five main components
- Spark SQL Architecture
  - Features, Languages, How DataFrames work, The SQLContext, Data sources
- Loading And Querying a Dataset with Spark SQL
  - Live demonstration of setting up a SQLContext
  - Loading it with data
  - Running queries against it
- Machine Learning with Spark MLlib
  - Collaborative filtering basics
  - Alternating Least Squares (ALS) algorithm
  - Live demo of a simple recommender model and training-test loop iterations



- How to connect to Spark SQL using ODBC/JDBC
  - Live demonstration of how to leverage Spark SQL ODBC/JDBC connectivity using Tableau
- Next steps
  - Some Real-World Use Cases
  - Basically answer the questions "What's it good for?" and "Who's using this?"
  - How to download a ready-to-use sandbox VM

## Spark Background / Overview



- A powerful open source processing engine
- Promoted to Apache top-level project in 2014
- Key Drivers:
  - Rich and well-documented APIs
  - Supports multiple languages
  - Optimized to run in-memory
- It is not a Hadoop replacement
  - Though it may eventually replace MapReduce



© 2016 MapR Technologies **MAPR** 3

Spark began life in 2009 as a project within the AMPLab at the University of California, Berkeley.

Spark became an incubated project of the Apache Software Foundation in 2013

Was promoted to Foundation top-level project in early 2014

- Is currently one of the most active projects managed by the Foundation

Key drivers (Reasons why you might select Spark for a data project):

- Rich and well-documented API designed specifically for interacting with data at scale.
- Supports Java, Scala, Python, R, and SQL
- Spark was optimized to run in memory from the beginning
  - Helps it process data far more quickly than

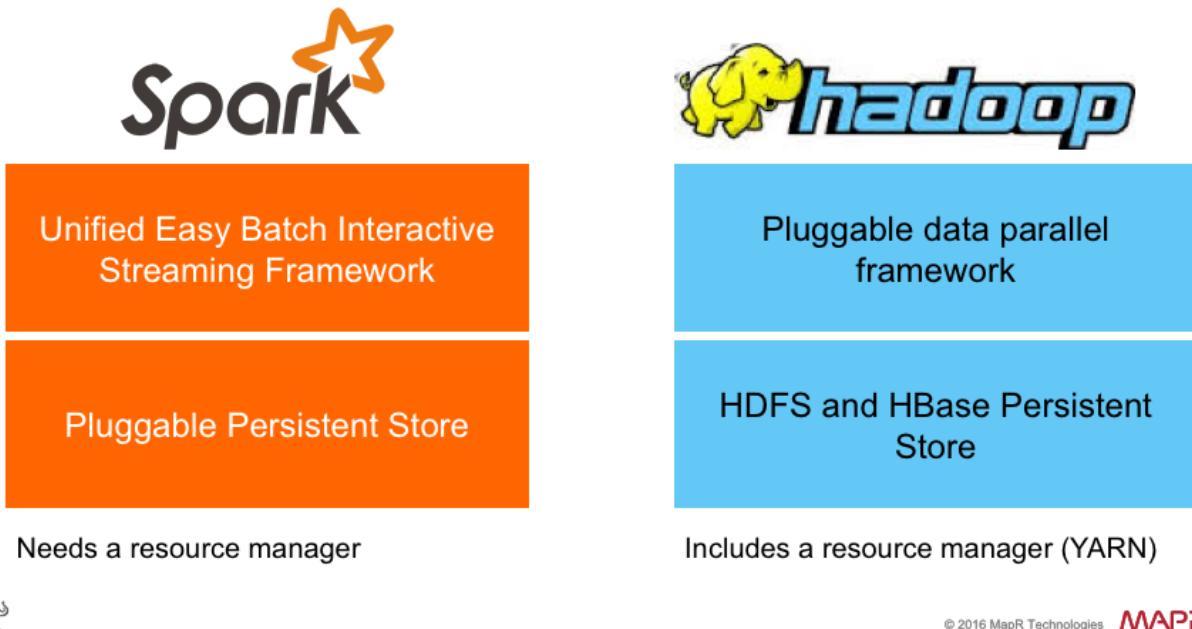


alternative approaches like Hadoop's MapReduce, which tends to write data to and from computer hard drives between each stage of processing.

Spark is not, despite the hype, a replacement for Hadoop. Nor is MapReduce dead.

- Hadoop is a platform that encompasses a wide variety of technologies
- Spark is faster than MapReduce for iterative algorithms that fit data in memory.
  - So, even though you can run Spark in stand-alone mode, doing so means you completely miss out on Hadoop's ability to run multiple types of workloads (incl. advanced analytics with Spark) on the same data at the same time.
  - In other words, Spark without Hadoop is just another silo.

## The Spark+Hadoop Team



Important to remember: Hadoop is more than just MapReduce

Spark and MapReduce are:

- Scalable frameworks for executing custom code on a cluster
- Nodes in the cluster work independently to process fragments of data and also combine those fragments together when appropriate to yield a final result
- Can tolerate loss of a node during a computation
- **Require a distributed storage layer for common data view**

Spark is often deployed in conjunction with a Hadoop cluster

- On its own, Spark isn't well-suited to production workloads
- It needs a resource manager
  - YARN, for instance, which takes responsibility for scheduling tasks across cluster nodes
  - A distributed filesystem



- Data security

## Components



- Spark Core
- Spark SQL
- Spark Streaming
- MLlib
- GraphX



© 2016 MapR Technologies **MAPR** 5

Apache Spark consists of Spark Core and a set of libraries.

### Spark Core

- This is the distributed execution engine
- Implements a programming abstraction known as Resilient Distributed Datasets RDDs
- Provides APIs for Java, Scala, and Python
- Some documentation shows R as part of the core, others show it as a separate module. The R module was added to the 1.4 release.

### Spark SQL

- This is the module for working with structured data using SQL
- Supports the open source Hive project, along with ODBC and JDBC

### Spark Streaming



- Enables scalable and fault-tolerant processing of data streams
- Supports sources like Flume (for data logs) and Kafka (for messaging)

### MLlib

- Scalable machine learning library
- Implements commonly-used machine learning and statistical algorithms
- These algorithms are often iterative
  - Spark's ability to cache the dataset in memory greatly speeds up such iterative data processing
  - So it's an ideal processing engine for implementing such algorithms

### GraphX

- Supports analysis of and computation over graphs of data
- Includes a number of widely understood graph algorithms, including PageRank.
- Graph databases are well-suited for analyzing interconnections
- They're useful for working with data that involve complex relationships (such as social media)

This presentation will focus primarily on the Spark Core API using Scala, along with SQL and the machine learning library

## Resilient Distributed Datasets (RDDs)



- Core programming abstraction
- Supports in-memory data storage
- Fault-tolerant and efficient
- Immutable (they cannot be changed)



© 2016 MapR Technologies **MAPR** 6

RDDs are a core programming abstraction at the heart of Spark.

Allows cross-cluster distribution

**Fault-tolerance** is achieved by tracking the lineage of transformations applied to coarse-grained sets of data.

**Efficiency** is achieved by parallelization of processing across multiple cluster nodes, and by minimization of data replication between those nodes.

RDDs remain in memory to enhance performance

Particularly true in use cases with a requirement for iterative queries or processes.

We will observe this in the demo



**This immutability is important.**

The chain of transformations from RDD1 to RDDn are logged, and can be repeated in the event of data loss or the failure of a cluster node.

## Resilient Distributed Datasets (RDDs)



- Core programming abstraction
- Supports in-memory data storage
- Fault-tolerant and efficient
- Immutable (they cannot be changed)

## RDD Operations

- Transformations
  - Mapping, filtering, sampling, etc.
  - A new RDD is created
  - Are “lazily” evaluated
- Actions
  - Counts, saves to files, sampling, etc.
  - Measure (but don’t change) the original data



© 2016 MapR Technologies 7

There are two basic types of operations

**Transformations** essentially change the data, so a new RDD is created because the original cannot be changed

They’re lazily evaluated – meaning that they aren’t executed until a subsequent action has a need for a result

This improves performance because it can avoid unnecessary processing

**Actions** measure – but don’t change – the original data.

They essentially force processing to take place



## Data Frames



- Distributed collection of data
- Organized into named columns
- Basically similar to a table
- APIs for Python, Java, Scala, and R



© 2016 MapR Technologies **MAPR** 8

The DataFrames API was added to Spark in 2015

If you're an R or Python programmer then you'll be familiar with the concept here

They have the ability to scale from kilobytes of data on a single laptop to petabytes on a large cluster

Basically this is structured data, just like a table in a relational database

Which leads us to...



## Data Frames

- Distributed collection of data
- Organized into named columns
- Basically similar to a table
- APIs for Python, Java, Scala, and R



## Spark SQL

- Enables JDBC/ODBC access from BI tools
- Supports Hive, Avro, Parquet, ORC, JSON, and JDBC sources
- Supported by Python, Java, Scala, and R APIs



© 2016 MapR Technologies **MAPR** 9

### ...Spark SQL

Basically this gives Spark programs the ability to query structured data using SQL

It also enables accessibility from BI tools such as Tableau  
In fact we'll explore this later in the demo

You've got a wide variety of supported data sources  
Joins between dataframes are supported as well

So this can be a very powerful capability



## Demo #1: RDDs, Data Frames, and SQL



© 2016 MapR Technologies The MapR logo is located in the bottom right corner of the slide, just below the footer text.



## My Environment

- 13" Mac Laptop
- 16GB RAM, 256GB SSD
- 2.9 GHz Intel Core i5 (2 cores)
- Spark installed using Homebrew



```
brew install apache-spark
```

Scala Shell: spark-shell

Python Shell: pyspark

SQL Shell: spark-sql



So this is not a fast machine I'm working with

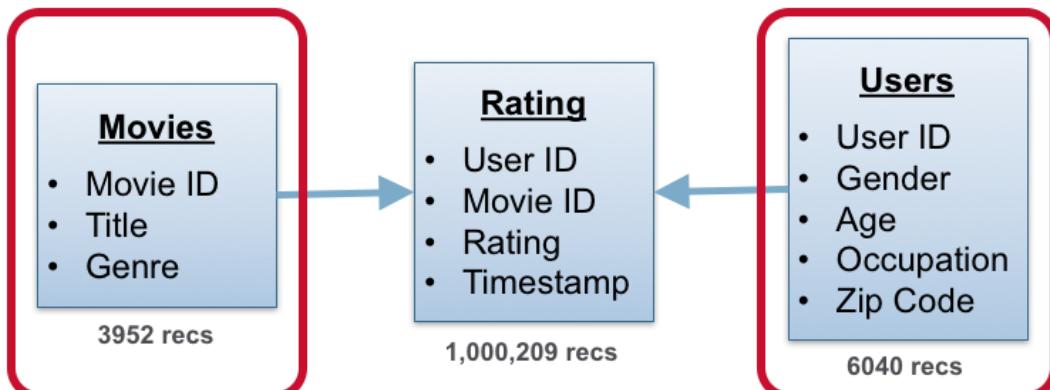
Installation using Homebrew is dead simple.  
Basically Homebrew downloads Apache Spark



## My Data Set

movielens

- GroupLens Research - University of Minnesota
- Sizes available: 100K to over 20MM records



© 2016 MapR Technologies MAPR 12

MovieLens is maintained by members of GroupLens Research at the University of Minnesota

I'm using one of the smaller sets, collected in 2000

We'll focus on the Users and Movies files for this first demo

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



## Define Schemas

```
case class Movie  
(movieId: Int,  
 title: String,  
 genres: Seq[String])
```

```
case class User  
(userId: Int,  
 gender: String,  
 age: Int,  
 occupation: Int,  
 zip: String)
```



These are Scala case classes that define schemas corresponding to their respective files



## Define Parsing Functions

```
def parseMovie(str: String): Movie = {  
    val fields = str.split("::")  
    assert(fields.size == 3)  
    Movie(fields(0).toInt, fields(1).toString, Seq(fields(2)))  
}  
  
def parseUser(str: String): User = {  
    val fields = str.split("::")  
    assert(fields.size == 5)  
    User(fields(0).toInt, fields(1).toString,  
         fields(2).toInt, fields(3).toInt, fields(4).toString)  
}
```



These functions parse lines from each file into the corresponding case classes defined

Note the references to the earlier-defined classes



## Load Files into RDDs

```
val moviesRDD = sc.textFile("movies.dat").map(parseMovie)  
val usersRDD = sc.textFile("users.dat").map(parseUser)
```



© 2016 MapR Technologies **MAPR** 15

These functions parse lines from each file into the corresponding case classes defined

Note the **map** method – this invokes each of the parsing functions defined earlier



## Load Files into RDDs

```
val moviesRDD = sc.textFile("movies.dat").map(parseMovie)  
val usersRDD = sc.textFile("users.dat").map(parseUser)
```

## Convert RDDs into Data Frames

```
val moviesDF = moviesRDD.toDF()  
val usersDF = usersRDD.toDF()
```



The **toDF** method makes the conversion happen

Examine the schemas once loaded



## Load Files into RDDs

```
val moviesRDD = sc.textFile("movies.dat").map(parseMovie)  
val usersRDD = sc.textFile("users.dat").map(parseUser)
```

## Convert RDDs into Data Frames

```
val moviesDF = moviesRDD.toDF()  
val usersDF = usersRDD.toDF()
```

## Register Data Frames

```
moviesDF.registerTempTable("movies")  
usersDF.registerTempTable("users")
```



© 2016 MapR Technologies **MAPR** 17

`registerTempTable()` creates an in-memory table that is scoped to the cluster in which it was created.

The data is stored using Hive's highly-optimized, in-memory columnar format.

This allows it to be queried with SQL statements that are executed using methods provided by `sqlContext`

We'll look at a **more permanent method** for materializing data frames at the end of this presentation.

I referenced this document:

<https://forums.databricks.com/questions/400/what-is-the-difference-between-registertemtable-a.html>



## Query #1 – What are some of the movies?

```
sqlContext.sql("SELECT * FROM movies").show()
```



Just a simple query here



## Query #1 – What are some of the movies?

```
sqlContext.sql("SELECT * FROM movies").show()
```

## Query #2 – What are some user demographics?

```
sqlContext.sql(  
    "SELECT  
        age,  
        sum(if(gender='M',user_count,0)) AS M,  
        sum(if(gender='F',user_count,0)) AS F  
    FROM  
        (SELECT age, gender, count(*) AS user_count  
        FROM users GROUP BY age, gender) AS ud  
    GROUP BY age ORDER BY age").show()
```



A little more involved query that utilizes a nested view



# Machine Learning with Spark MLlib



© 2016 MapR Technologies  MAPR

So we've covered the basics - RDDs, Data Frames, and SQL

Let's move on to the more interesting stuff



## Collaborative Filtering

- Technique used by *recommender systems*
- Seek to predict ratings/preferences
- Based on data from larger population
- Assumption:  
User A and User B agree about X  
So, User B is *more likely* to agree  
with User A about Y than a randomly-chosen opinion



Collaborative filtering (CF) is a technique used by some recommender systems.

**Recommender systems** essentially seek to predict the 'rating' or 'preference' that a user would give to an item.

So, **Collaborative Filtering** tries to make these predictions (we call this "filtering") about the interests of a user by collecting information from many users (hence "collaborating")



## Collaborative Filtering

- Technique used by *recommender systems*
- Seek to predict ratings/preferences
- Based on data from larger population
- Assumption:  
User A and User B agree about X  
So, User B is *more likely* to agree  
with User A about Y than a randomly-chosen opinion



I borrowed this animated GIF from Wikipedia

Source: [https://en.wikipedia.org/wiki/Collaborative\\_filtering](https://en.wikipedia.org/wiki/Collaborative_filtering)

This is an example of collaborative filtering.

At first, people rate different items (like videos, images, games). Then, the system makes predictions about a user's rating for an item not rated yet.

**The new predictions are built upon the existing ratings of other users with similar ratings with the active user.**

In the image, the system predicts that the user will not like the video.



## Alternating Least Squares (ALS) Algorithm

- Overall goal: Complete the rating matrix
- The rating matrix is your central data structure.

		Movies						
		1	2	3	4	5	6	7
Users	1	4	5			3	2	1
	2		4			5		
	3	4			1		4	
	4					5		
	5						5	5
	6	4	2		3	5		
	7			3		4		

Ratings Matrix



© 2016 MapR Technologies MAPR 23

Here's the overall goal:

Accurately predict every user's rating for the movies they haven't watched yet.

I used an explanation from Cambridge Coding Academy's "Predicting User Preferences in Python using Alternating Least Squares" tutorial located here:  
<http://online.cambridgecoding.com/notebooks/mhaller/predicting-user-preferences-in-python-using-alternating-least-squares>



## Alternating Least Squares (ALS) Algorithm

- Overall goal: Complete the rating matrix
- The rating matrix is your central data structure.

	Movies							
Users	4	5	4	2	3	2	3	1
	4	2	4	2	5	3	2	1
	4	3	3	1	2	4	2	2
	3	4	4	3	5	4	3	2
	5	4	3	4	5	4	5	5
	4	2	4	3	5	3	4	3
	4	3	4	3	4	4	3	2

Ratings Matrix



If we do this right, we can suggest the most suitable movies to watch next to each user individually.

This of course can apply to lots of things – in our case it's movies.

I used an explanation from Cambridge Coding Academy's "Predicting User Preferences in Python using Alternating Least Squares" tutorial located here:  
<http://online.cambridgecoding.com/notebooks/mhaller/predicting-user-preferences-in-python-using-alternating-least-squares>



## Alternating Least Squares (ALS) Algorithm

- There are a total of “m” users (rows)
- And a total of “n” movies (columns)

Movies (n)				
Users (m)	$r_{11}$	$r_{12}$	$\dots$	$r_{1n}$
	$r_{21}$	$r_{22}$	$\dots$	$r_{2n}$
	:	:	:	:
	$r_{m1}$	$r_{m2}$	$\dots$	$r_{mn}$

Ratings  
 $m \times n$



We'll shrink this down for the sake of keeping things simple.

What we have are some number of users => “m” rows  
And some number of movies => “n” columns

The overall ratings matrix is  $m \times n$  (rows x columns)

I used an explanation from Cambridge Coding Academy's  
“Predicting User Preferences in Python using Alternating Least  
Squares” tutorial located here:  
<http://online.cambridgecoding.com/notebooks/mhaller/predicting-user-preferences-in-python-using-alternating-least-squares>



## Alternating Least Squares (ALS) Algorithm

- Break out users – We will denote them with the letter “P”

Movies (n)				
Users (m)	r <sub>11</sub>	r <sub>12</sub>	...	r <sub>1n</sub>
	r <sub>21</sub>	r <sub>22</sub>	...	r <sub>2n</sub>
:	:	:		:
	r <sub>m1</sub>	r <sub>m2</sub>	...	r <sub>mn</sub>

Ratings  
m x n

Users (m)	p <sub>11</sub>
	p <sub>21</sub>
:	
	p <sub>m1</sub>



Let's break off the users – We know that there are m number of them, that's a given

But there's something we don't know about these users  
 There are certain influences that might cause these users to rate movies higher or lower

It could be **demographics** - age-related, gender, profession, whatever.

The point is, **we don't know what they are**, we just know they're present in the background somewhere.

I used this reference regarding latent variables:

<http://stats.stackexchange.com/questions/162585/predicting-score-in-the-presence-of-latent-variables>

I used an explanation from Cambridge Coding Academy's



“Predicting User Preferences in Python using Alternating Least Squares” tutorial located here:  
<http://online.cambridgecoding.com/notebooks/mhaller/predicting-user-preferences-in-python-using-alternating-least-squares>

## Alternating Least Squares (ALS) Algorithm

- Break out users – We will denote them with the letter “P”
- They are influenced by “features” – a total quantity of “k”

		Movies (n)						Features (k)			
Users (m)						Users (m)					
		$r_{11}$	$r_{12}$	$\dots$	$r_{1n}$			$p_{11}$	$p_{12}$	$\dots$	$p_{1n}$
$r_{21}$	$r_{22}$	$\dots$	$r_{2n}$			$p_{21}$	$p_{22}$	$\dots$	$p_{2n}$		
:	:	:	:			:	:	:	:		
$r_{m1}$	$r_{m2}$	$\dots$	$r_{mn}$			$p_{m1}$	$p_{m2}$	$\dots$	$p_{mn}$		

Ratings    User Features  
 $m \times n$      $m \times k$



Now we'll introduce a new variable called **“features”**  
 Algorithms call these different things, such as **“latent factors”** or  
**“latent variables”** or **“hidden features”**  
 For the sake of expediency, we'll call them **“Features”**  
 We also don't know how many there are, so we'll just throw them in the matrix for now  
 We'll assign an arbitrary number of to represent their quantity  
 In other words, we're going to assume that there are **K features** for each user that might influence their rating in some way.

I used this reference regarding latent variables:  
<http://stats.stackexchange.com/questions/162585/predicting-score-in-the-presence-of-latent-variables>

I used an explanation from Cambridge Coding Academy's



“Predicting User Preferences in Python using Alternating Least Squares” tutorial located here:  
<http://online.cambridgecoding.com/notebooks/mhaller/predicting-user-preferences-in-python-using-alternating-least-squares>

## Alternating Least Squares (ALS) Algorithm

- Break out movies – We will denote them with the letter “Q”

Movies (n)		Features (k)		Movies (n)			
Users (m)		Users (m)		Users (m)			
$r_{11}$	$r_{12}$	$\dots$	$r_{1n}$	$p_{11}$	$p_{12}$	$\dots$	$p_{1n}$
$r_{21}$	$r_{22}$	$\dots$	$r_{2n}$	$p_{21}$	$p_{22}$	$\dots$	$p_{2n}$
:	:	:	:	:	:	:	:
$r_{m1}$	$r_{m2}$	$\dots$	$r_{mn}$	$p_{m1}$	$p_{m2}$	$\dots$	$p_{mn}$

Ratings    User Features  
 $m \times n$      $m \times k$



Likewise, let's break off the movies – We know that there are n number of them, that's a given

Like users, there's something we don't know about these movies

It's those same **features** that we just talked about...

I used an explanation from Cambridge Coding Academy's "Predicting User Preferences in Python using Alternating Least Squares" tutorial located here:  
<http://online.cambridgecoding.com/notebooks/mhaller/predicting-user-preferences-in-python-using-alternating-least-squares>



## Alternating Least Squares (ALS) Algorithm

- Break out movies – We will denote them with the letter “Q”
- They are also influenced by “features” – a total quantity of “k”

Movies (n)		Features (k)		Movies (n)	
Users (m)		Users (m)		Features (k)	
$r_{11}$	$r_{12}$	$\dots$	$r_{1n}$	$p_{11}$	$p_{12}$
$r_{21}$	$r_{22}$	$\dots$	$r_{2n}$	$p_{21}$	$p_{22}$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$r_{m1}$	$r_{m2}$	$\dots$	$r_{mn}$	$p_{m1}$	$p_{m2}$
Ratings $m \times n$		User Features $m \times k$		Movie Features $k \times n$	



So we'll make the same assumption for the movies:

We'll assume that the same features that influence users to give certain ratings to certain movies would likewise be applicable

For example if we knew 25-29 year old men from California tend to rate certain types of movies a certain way, we could extrapolate the other direction

And likewise say certain types of movies are rated a certain way by 25-29 year old men from California

So these factor matrices represent hidden features which the algorithm tries to discover.

One matrix tries to describe the latent or hidden features of each user, and one tries to describe latent properties of each movie.

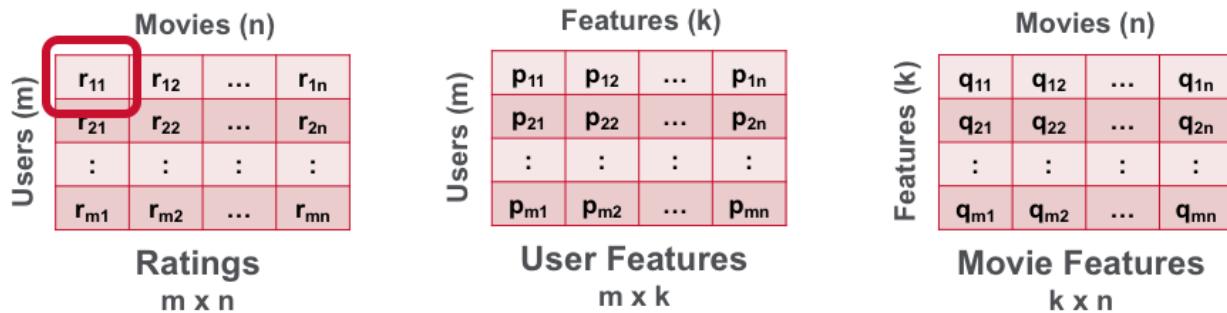


I used an explanation from Cambridge Coding Academy's  
"Predicting User Preferences in Python using Alternating Least  
Squares" tutorial located here:

<http://online.cambridgecoding.com/notebooks/mhaller/predicting-user-preferences-in-python-using-alternating-least-squares>

## Alternating Least Squares (ALS) Algorithm

- Ratings are the **sum of user features** weighted by **movie features**



Each of the **ratings** is actually the sum of latent user features, weighted by the respective latent movie features. You get that by computing the **dot product** of the user and movie features matrices

I used an explanation from Cambridge Coding Academy's "Predicting User Preferences in Python using Alternating Least Squares" tutorial located here:  
<http://online.cambridgecoding.com/notebooks/mhaller/predicting-user-preferences-in-python-using-alternating-least-squares>



## Alternating Least Squares (ALS) Algorithm

- Ratings are the **sum of user features** weighted by **movie features**
- Dot Product  $r_{11} = (p_{11} \times q_{11}) + (p_{12} \times q_{21})$  and so on...

		Movies (n)			
		$r_{11}$	$r_{12}$	$\dots$	$r_{1n}$
Users (m)	$r_{21}$	$r_{22}$	$\dots$	$r_{2n}$	
	:	:	:	:	
	$r_{m1}$	$r_{m2}$	$\dots$	$r_{mn}$	
	Ratings				

$m \times n$

		Features (k)			
		$p_{11}$	$p_{12}$	$\dots$	$p_{1n}$
Users (m)	$p_{21}$	$p_{22}$	$\dots$	$p_{2n}$	
	:	:	:	:	
	$p_{m1}$	$p_{m2}$	$\dots$	$p_{mn}$	
	User Features				

$m \times k$

		Movies (n)			
		$q_{11}$	$q_{12}$	$\dots$	$q_{1n}$
Features (k)	$q_{21}$	$q_{22}$	$\dots$	$q_{2n}$	
	:	:	:	:	
	$q_{m1}$	$q_{m2}$	$\dots$	$q_{mn}$	
	Movie Features				

$k \times n$



Here's how **dot products** work:

A dot product is composed of the corresponding **user features** row and the first **movie features** column

It looks like this:  $R_{11} = (P_{11} \times Q_{11}) + (P_{12} \times Q_{21})$  and so on.

But there's a problem...

The concept of **dot product** is explained here:

<https://www.mathsisfun.com/algebra/matrix-multiplying.html>

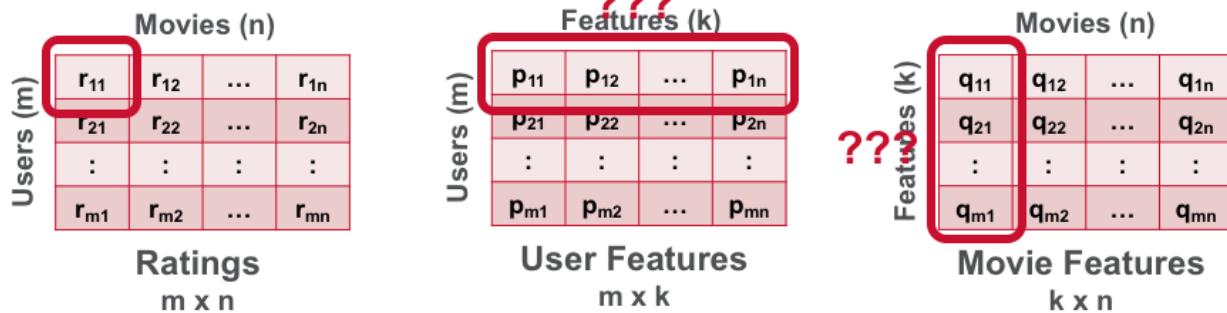
I used an explanation from Cambridge Coding Academy's "Predicting User Preferences in Python using Alternating Least Squares" tutorial located here:



<http://online.cambridgecoding.com/notebooks/mhaller/predicting-user-preferences-in-python-using-alternating-least-squares>

## Alternating Least Squares (ALS) Algorithm

- Rating are the **sum of user features** weighted by **movie features**
- Dot Product  $r_{11} = (p_{11} \times q_{11}) + (p_{12} \times q_{21})$  and so on...
- Problem:



Problem:

We don't know what the features are, or how many we have

For that, we'll use the trick of **alternating least squares**.

The concept of **dot product** is explained here:

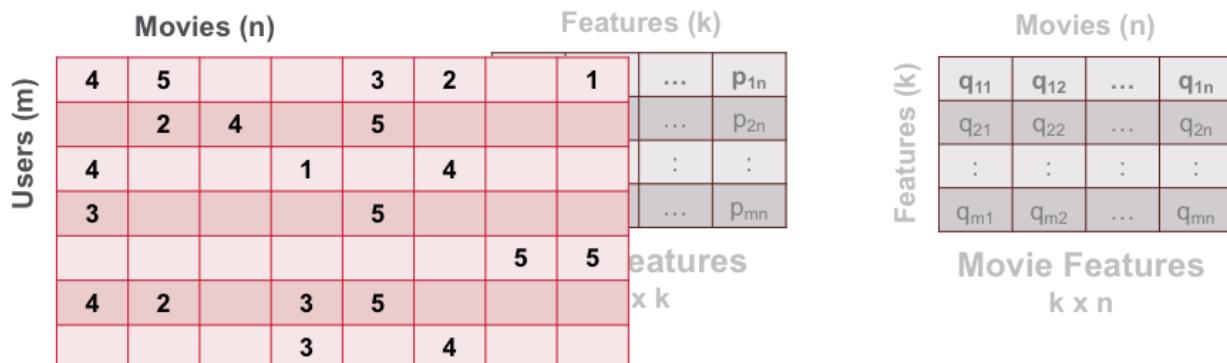
<https://www.mathsisfun.com/algebra/matrix-multiplying.html>

I used an explanation from Cambridge Coding Academy's "Predicting User Preferences in Python using Alternating Least Squares" tutorial located here:  
<http://online.cambridgecoding.com/notebooks/mhaller/predicting-user-preferences-in-python-using-alternating-least-squares>



## Alternating Least Squares (ALS) Algorithm

- We already know *some* of the movie ratings



© 2016 MapR Technologies MAPR 33

Here's the concept:

We already know some of the movie ratings – remember these?

I used an explanation from Cambridge Coding Academy's "Predicting User Preferences in Python using Alternating Least Squares" tutorial located here:  
<http://online.cambridgecoding.com/notebooks/mhaller/predicting-user-preferences-in-python-using-alternating-least-squares>



## Alternating Least Squares (ALS) Algorithm

- We already know *some* of the movie ratings
- Our goal is to fill in the ones that we don't know

		Movies (n)								Features (k)		Movies (n)			
Users (m)		4	5	4	2	3	2	3	1	...	p <sub>1n</sub>	q <sub>11</sub>	q <sub>12</sub>	...	q <sub>1n</sub>
		4	2	4	2	5	3	2	1	...	p <sub>2n</sub>	q <sub>21</sub>	q <sub>22</sub>	...	q <sub>2n</sub>
		4	3	3	1	2	4	2	2	:	:	⋮	⋮	⋮	⋮
		3	4	4	3	5	4	3	2	...	p <sub>mn</sub>	q <sub>m1</sub>	q <sub>m2</sub>	...	q <sub>mn</sub>
		5	4	3	4	5	4	5	5	Features		Movie Features		k x n	
		4	2	4	3	5	3	4	3						
		4	3	4	3	4	4	3	2						
		4	3	4	3	4	4	3	2						



© 2016 MapR Technologies 34

Our goal here is to fill in the ones that we don't know.

I used an explanation from Cambridge Coding Academy's "Predicting User Preferences in Python using Alternating Least Squares" tutorial located here:  
<http://online.cambridgecoding.com/notebooks/mhaller/predicting-user-preferences-in-python-using-alternating-least-squares>



## Alternating Least Squares (ALS) Algorithm

- We'll use what we do know to solve for what we don't know

Movies (n)			
Users (m)	4	5	
	2	4	
	4		1
	3		

Ratings  
 $m \times n$

Features (k)			
Users (m)	$p_{11}$	$p_{12}$	$\dots$
	$p_{21}$	$p_{22}$	$\dots$
	:	:	:
	$p_{m1}$	$p_{m2}$	$\dots$

User Features  
 $m \times k$

Movies (n)			
Features (k)	$q_{11}$	$q_{12}$	$\dots$
	$q_{21}$	$q_{22}$	$\dots$
	:	:	:
	$q_{m1}$	$q_{m2}$	$\dots$

Movie Features  
 $k \times n$



I used an explanation from Cambridge Coding Academy's "Predicting User Preferences in Python using Alternating Least Squares" tutorial located here:  
<http://online.cambridgecoding.com/notebooks/mhaller/predicting-user-preferences-in-python-using-alternating-least-squares>

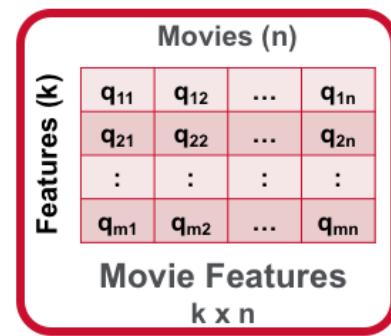
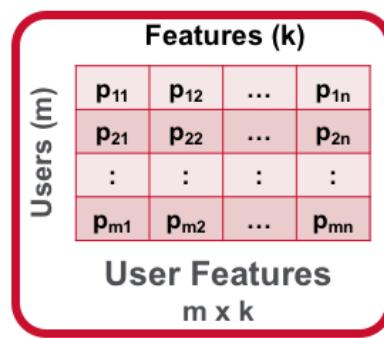


## Alternating Least Squares (ALS) Algorithm

- We'll use what we do know to solve for what we don't know
- Namely, the **User Features** and **Movie Features** matrices

Movies (n)			
Users (m)	4	5	
	2	4	
	4		1
	3		

Ratings  
m x n



I used an explanation from Cambridge Coding Academy's "Predicting User Preferences in Python using Alternating Least Squares" tutorial located here:  
<http://online.cambridgecoding.com/notebooks/mhaller/predicting-user-preferences-in-python-using-alternating-least-squares>



## Alternating Least Squares (ALS) Algorithm

- Let's just assume we know the number of features

Movies (n)			
Users (m)	4	5	
	2	4	
	4		1
	3		

Ratings  
 $m \times n$

Features (k)			
Users (m)	$p_{11}$	$p_{12}$	$\dots$
	$p_{21}$	$p_{22}$	$\dots$
	:	:	:
	$p_{m1}$	$p_{m2}$	$\dots$

User Features  
 $m \times k$

Movies (n)			
Features (k)	$q_{11}$	$q_{12}$	$\dots$
	$q_{21}$	$q_{22}$	$\dots$
	:	:	:
	$q_{m1}$	$q_{m2}$	$\dots$

Movie Features  
 $k \times n$



We'll start by assuming we know we know the number of features

I used an explanation from Cambridge Coding Academy's "Predicting User Preferences in Python using Alternating Least Squares" tutorial located here:  
<http://online.cambridgecoding.com/notebooks/mhaller/predicting-user-preferences-in-python-using-alternating-least-squares>



## Alternating Least Squares (ALS) Algorithm

- Let's just assume we know the number of features
- We'll "fix" the Movie Features matrix Q with random values:

Movies (n)			
Users (m)	4	5	
	2	4	
	4		1
	3		

Ratings  
 $m \times n$

Features (k)			
Users (m)			

User Features  
 $m \times k$

Movies (n)			
Features (k)	.452	.803	.642
	.343	.911	.892
	.546	.710	.353
	.145	.998	.908
			.830

Movie Features  
 $k \times n$



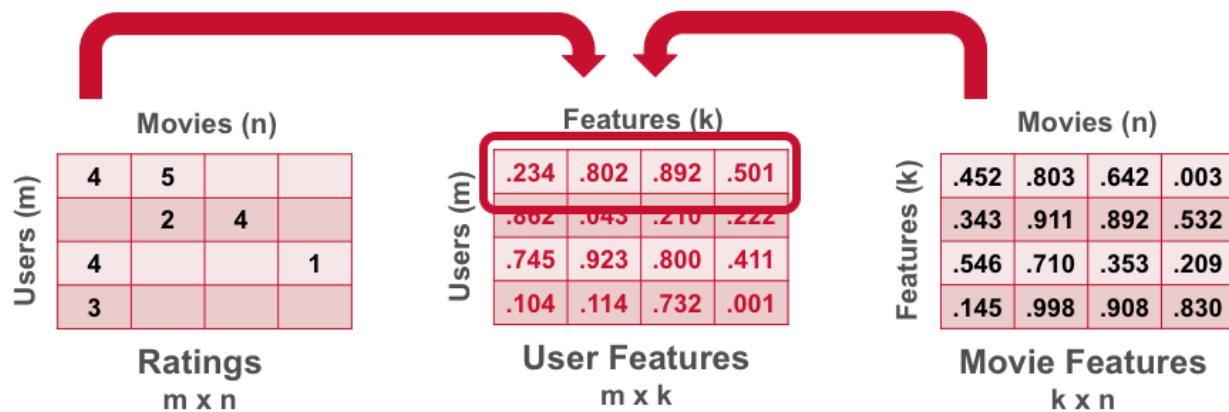
And we'll fill the **Movie Features** matrix with random values

I used an explanation from Cambridge Coding Academy's "Predicting User Preferences in Python using Alternating Least Squares" tutorial located here:  
<http://online.cambridgecoding.com/notebooks/mhaller/predicting-user-preferences-in-python-using-alternating-least-squares>



## Alternating Least Squares (ALS) Algorithm

- Estimate User Features



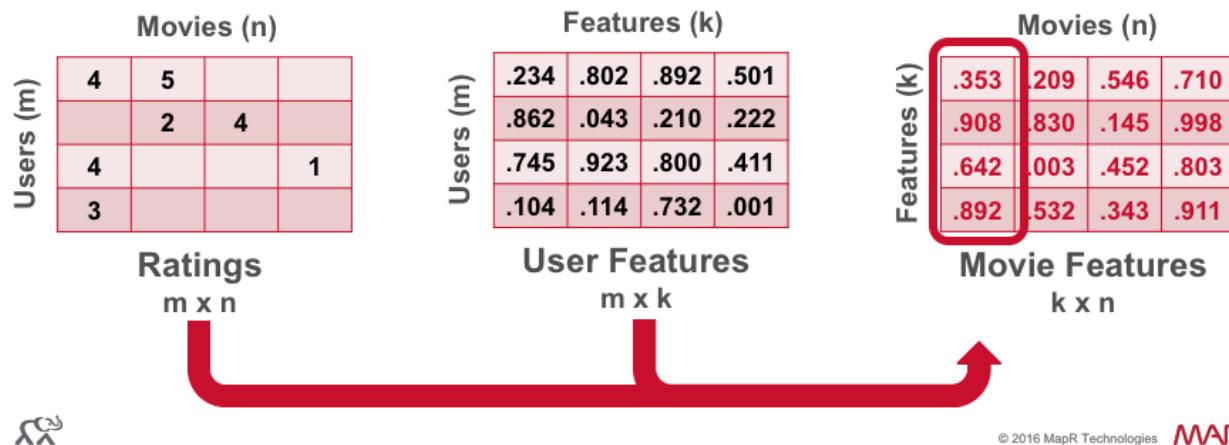
We'll then estimate the **User Features** matrix for every user, row by row, using the non-empty cells from the **Ratings Matrix** and the random values from the **Movie Features** matrix.

I used an explanation from Cambridge Coding Academy's "Predicting User Preferences in Python using Alternating Least Squares" tutorial located here:  
<http://online.cambridgecoding.com/notebooks/mhaller/predicting-user-preferences-in-python-using-alternating-least-squares>



## Alternating Least Squares (ALS) Algorithm

- Estimate User Features
- Then use User Features to estimate Movie Features



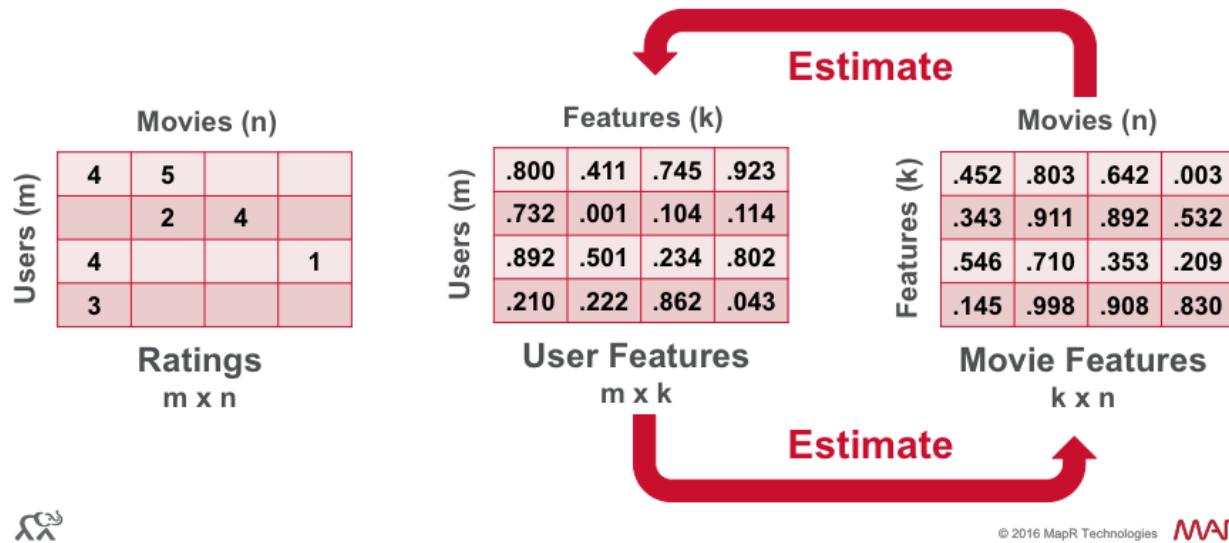
We'll then estimate the **Movie Features** matrix for every movie, column by column, using the non-empty cells from the **Ratings Matrix** and the values from the **User Features** matrix.

I used an explanation from Cambridge Coding Academy's "Predicting User Preferences in Python using Alternating Least Squares" tutorial located here:  
<http://online.cambridgecoding.com/notebooks/mhaller/predicting-user-preferences-in-python-using-alternating-least-squares>



## Alternating Least Squares (ALS) Algorithm

- Rinse and repeat...



We'll repeat this procedure some number of times

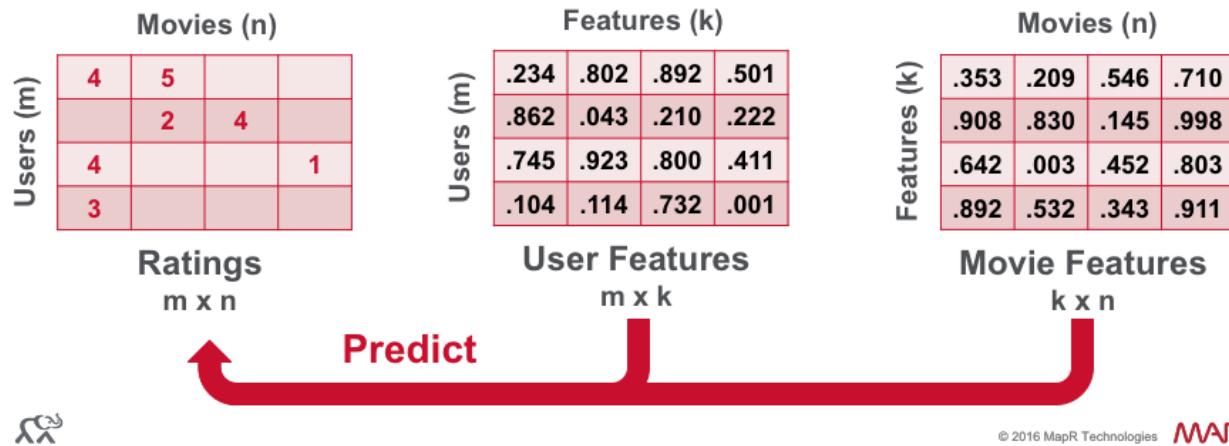
This alternation between which matrix to optimize is where the "alternating" in the name comes from.

I used an explanation from Cambridge Coding Academy's "Predicting User Preferences in Python using Alternating Least Squares" tutorial located here:

<http://online.cambridgecoding.com/notebooks/mhaller/predicting-user-preferences-in-python-using-alternating-least-squares>

## Alternating Least Squares (ALS) Algorithm

- Eventually we'll have a model for predicting the ratings we know



I used an explanation from Cambridge Coding Academy's "Predicting User Preferences in Python using Alternating Least Squares" tutorial located here:  
<http://online.cambridgecoding.com/notebooks/mhaller/predicting-user-preferences-in-python-using-alternating-least-squares>

## Alternating Least Squares (ALS) Algorithm

- Eventually we'll have a model for predicting the ratings we know
- We call this **Convergence**

Movies (n)			
Users (m)	4	5	
	2	4	
	4		1
	3		

Ratings  
 $m \times n$

Features (k)			
Users (m)	.234	.802	.892
	.862	.043	.210
	.745	.923	.800
	.104	.114	.732
			.501
			.222
			.411
			.001

User Features  
 $m \times k$

Movies (n)			
Features (k)	.353	.209	.546
	.908	.830	.145
	.642	.003	.452
	.892	.532	.343
			.710
			.998
			.803
			.911

Movie Features  
 $k \times n$



✓ Yep, it works.



I used an explanation from Cambridge Coding Academy's "Predicting User Preferences in Python using Alternating Least Squares" tutorial located here:  
<http://online.cambridgecoding.com/notebooks/mhaller/predicting-user-preferences-in-python-using-alternating-least-squares>



## Alternating Least Squares (ALS) Algorithm

- Eventually we'll have a model for predicting the ratings we know
- We call this **Convergence**
- The process is called **Model Training** – and it's done *offline*

		Movies (n)				
		Users (m)	4	5		
		2	4			
		4			1	
		3				

Ratings  
 $m \times n$

		Features (k)				
		Users (m)	.234	.802	.892	.501
		.862	.043	.210	.222	
		.745	.923	.800	.411	
		.104	.114	.732	.001	

User Features  
 $m \times k$

		Movies (n)				
		Features (k)	.353	.209	.546	.710
		.908	.830	.145	.998	
		.642	.003	.452	.803	
		.892	.532	.343	.911	

Movie Features  
 $k \times n$

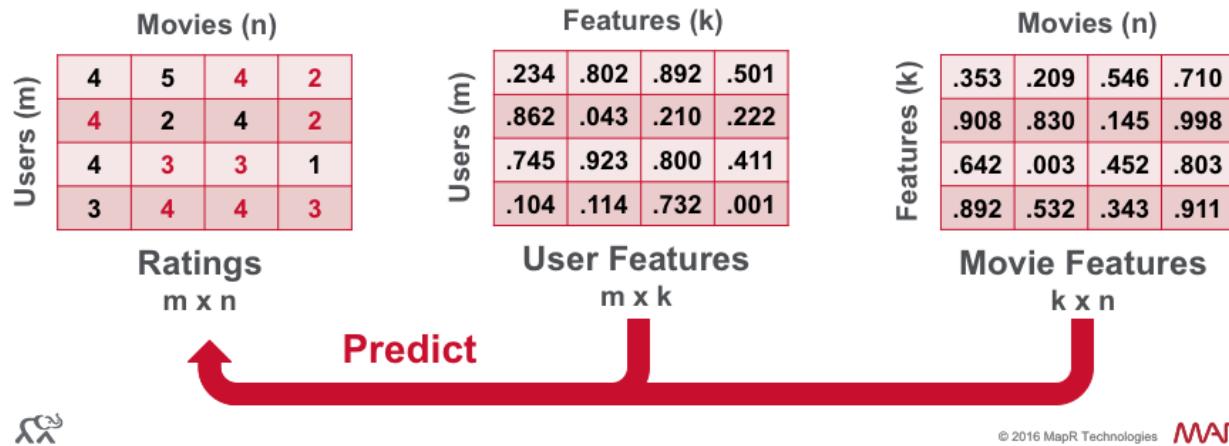


I used an explanation from Cambridge Coding Academy's "Predicting User Preferences in Python using Alternating Least Squares" tutorial located here:  
<http://online.cambridgecoding.com/notebooks/mhaller/predicting-user-preferences-in-python-using-alternating-least-squares>



## Alternating Least Squares (ALS) Algorithm

- Now we can use the model to **predict what we don't know**:



I used an explanation from Cambridge Coding Academy's "Predicting User Preferences in Python using Alternating Least Squares" tutorial located here:  
<http://online.cambridgecoding.com/notebooks/mhaller/predicting-user-preferences-in-python-using-alternating-least-squares>

## Model Quality Influencers

- Training data (quality and quantity)
- Number of features **Rank**
- Number of training iterations to run **Iterations**
- **Regularization** level (control “overfitting”) **Lambda**



So we have things that will influence how good our model is

**Regularization** is a technique used in an attempt to solve the overfitting problem in statistical models.

Here's how it works, in a nutshell:

Let's assume we only have two features that influence movie ratings – age and gender. That's it.

Well probably this model will fail because it's too simple.

So let's throw in more demographics such as profession, what part of the country they're from, education level, etc.

This makes the model more **interesting** and **complex**.

Our model will do better in some ways, but worse in others due to this concept of overfitting

This is because it's sticking too much to the data – in other words, it can't generalize



Which means we're being influenced by "background noise"  
So overfitting happens when the model works well on training data using known noise but doesn't perform well on the actual testing data.  
This isn't acceptable.

So we apply a regularization factor to control the model's complexity, and thus help prevent overfitting  
The higher the regularization factor, the lower the overfitting, but the greater the bias.  
If you're a visual person, imagine smoothing a regression line so that it's less influenced by noisy data points.

**Incidentally**, three of these are actually parameters into the algorithm

**Rank** - Starting at 5 and increasing by 5 until the recommendation improvement rate slows down, memory and CPU permitting, is a good approach

**Iterations** – Likewise, 5, 10, and 20 iterations for various combinations of rank and lambda are good starting points

**Lambda** - Values of 0.01, 1 and 10 are good values to test.

I used an explanation from Quora here:

<https://www.quora.com/What-is-regularization-in-machine-learning>

And here:

<https://cloud.google.com/solutions/recommendations-using-machine-learning-on-compute-engine>

## The Process

- Split the dataset into three parts:
  - Training **60%**
  - Validation **20%**
  - Testing **20%**



Split your dataset into three pieces:

- One for training
- One for validation
- One for testing



## The Process

- Split the dataset into three parts:
  - Training **60%**
  - Validation **20%**
  - Testing **20%**
- Train the model with the training dataset



Train the model with the training data – This is where the model gets created



## The Process

- Split the dataset into three parts:
  - Training 60%
  - Validation 20% (highlighted)
  - Testing 20%
- Train the model with the training dataset
- Run the model against the validation dataset and rate its accuracy



Run the model against the validation data

(Why not run it against the training data? Because we already know it works pretty well with that data set)



## The Process

- Split the dataset into three parts:
  - Training **60%**
  - Validation **20%**
  - Testing **20%**
- Train the model with the training dataset
- Run the model against the validation dataset and rate its accuracy
- Repeat with adjusted **Rank**, **Iteration**, and **Lambda** parameters



Change your **rank**, **iteration**, and **lambda** values and re-run  
Did your accuracy improve?  
Do that over and over again until you're happy with the result



## The Process

- Split the dataset into three parts:
  - Training **60%**
  - Validation **20%**
  - Testing **20%**
- Train the model with the training dataset
- Run the model against the validation dataset and rate its accuracy
- Repeat with adjusted **Rank, Iteration, and Lambda** parameters
- Finally, score the model using the test dataset



Finally, run it against the test dataset and check its accuracy.  
That's your final score.  
If score isn't good enough then you might add additional  
training data or refine parameter adjustments.



## The Process

- Split the dataset into three parts:
  - Training **60%**
  - Validation **20%**
  - Testing **20%**
- Train the model with the training dataset
- Run the model against the validation dataset and rate its accuracy
- Repeat with adjusted **Rank**, **Iteration**, and **Lambda** parameters
- Finally, score the model using the test dataset
- What's missing?



There's one more thing we need here – a standard of measurement for rating the model's accuracy



## Root Mean Square Error (RMSE)

- A frequently-used performance measurement
- Evaluates difference between **predicted** and **actual** values



For our example we'll use something called **Root Mean Square Error**

This is a frequently used measure of the differences between the values predicted by a model and the values actually observed.

I used an explanation from Wikipedia here:

[https://en.wikipedia.org/wiki/Root-mean-square\\_deviation](https://en.wikipedia.org/wiki/Root-mean-square_deviation)



## Root Mean Square Error (RMSE)

- A frequently-used performance measurement
- Evaluates difference between **predicted** and **actual** values
- In mathematical terms:

*“RMSE is the **sample standard deviation** of the **differences** between **predicted** and **observed** values”*



In mathematical terms, it represents the **sample standard deviation** of the **differences** between **predicted** and **observed** values.

I used an explanation from Wikipedia here:

[https://en.wikipedia.org/wiki/Root-mean-square\\_deviation](https://en.wikipedia.org/wiki/Root-mean-square_deviation)



## Root Mean Square Error (RMSE)

- A frequently-used performance measurement
- Evaluates difference between **predicted** and **actual** values
- In mathematical terms:

*“RMSE is the **sample standard deviation** of the **differences** between **predicted** and **observed** values”*

- We will use it as a **cost function** to evaluate our model
- Our goal will be to keep the cost as low as possible.



I used an explanation from Wikipedia here:

[https://en.wikipedia.org/wiki/Root-mean-square\\_deviation](https://en.wikipedia.org/wiki/Root-mean-square_deviation)



## Demo #2: Under The Hood With MLlib ALS



© 2016 MapR Technologies The MapR logo is located in the bottom right corner of the slide. It consists of the word "MAPR" in a bold, red, sans-serif font, with a registered trademark symbol (®) at the top right corner of the letter "R".



## Import library for MLlib recommendation data types

```
import org.apache.spark.mllib.recommendation.{ALS, MatrixFactorizationModel, Rating}
```



## Import library for MLlib recommendation data types

```
import org.apache.spark.mllib.recommendation.{ALS, MatrixFactorizationModel, Rating}
```

## Create Ratings RDD, Load Data Frame, Register

```
def parseRating(str: String): Rating = {  
    val fields = str.split("::")  
    Rating(fields(0).toInt, fields(1).toInt, fields(2).toDouble)  
}  
  
val ratingsRDD = sc.textFile("ratings.dat").map(parseRating).cache()  
  
val ratingsDF = ratingsRDD.toDF  
  
ratingsDF.registerTempTable("ratings")
```



© 2016 MapR Technologies **MAPR** 58

Note the use of the **cache** method here – I want to keep this in memory



## Define the RMSE function

```
import org.apache.spark.rdd._  
def computeRmse(model: MatrixFactorizationModel, data: RDD[Rating], n: Long): Double = {  
    val predictions: RDD[Rating] = model.predict(data.map(x => (x.user, x.product)))  
    val predictionsAndRatings = predictions.map(x => ((x.user, x.product), x.rating))  
    .join(data.map(x => ((x.user, x.product), x.rating))).values  
    math.sqrt(predictionsAndRatings.map(x => (x._1 - x._2) * (x._1 - x._2)).reduce(_ + _) / n)  
}
```



I borrowed this one from an AWS blog here:

<https://blogs.aws.amazon.com/bigdata/post/Tx6J5RM20WPG5V/Building-a-Recommendation-Engine-with-Spark-ML-on-Amazon-EMR-using-Zeppelin>



## Split Ratings RDD into Training/Validation/Test RDDs

```
val splitsRDD = ratingsRDD.randomSplit(Array(0.6, 0.2, 0.2), 0L)
val trainingRDD = splitsRDD(0).cache()
val validationRDD = splitsRDD(1).cache()
val testRDD = splitsRDD(2).cache()
val numTraining = trainingRDD.count()
val numValidation = validationRDD.count()
val numTest = testRDD.count()
```

**60%, 20%, 20%**



Note that this is a 60% / 20% / 20% split



## Validation Runs

```
val ranks = List(5, 10, 15)
val iters = List(10, 20)
val lambdas = List(0.1, 1.0)
var bestModel: Option[MatrixFactorizationModel] = None
var bestRmse = Double.MaxValue
var bestRank = 0
var bestIter = -1
var bestLambda = -1.0
for (rank <- ranks; iter <- iters; lambda <- lambdas) {
    val model = ALS.train(trainingRDD, rank, iter, lambda)
    val validationRmse = computeRmse(model, validationRDD, numValidation)
    println(s"RMSE = $validationRmse for the model trained with rank = $rank, lambda = $lambda, and
iter = $iter.")
    if (validationRmse < bestRmse) {
        bestModel = Some(model)
        bestRmse = validationRmse
        bestRank = rank
        bestLambda = lambda
        bestIter = iter
    }
}
```



I've got three parameters, with a variety of values for each, for a total of 12 runs

The model will train

At the end of each run it'll tell us the RMSE

The lowest RMSE will be our final test run parameters.



## The Netflix Prize



- \$1,000,000 contest announced October 2, 2006
- An open competition for best collaborative filtering algorithm
- Predict user ratings for films
  - Based on previous ratings w/o info about the users or films
  - Training data set of 100,480,507 ratings that 480,189 users gave to 17,770 movies
  - Qualifying set had 2,817,131 ratings split into test and quiz sets
- Winner must beat Netflix's *Cinematch* algorithm RMSE by 10%
  - Cinematch RMSE was 0.9525 on the test set
  - Winner had to achieve 0.8572 RMSE (or less) on the test set



© 2016 MapR Technologies MAPR 62

Thought it might be fun to take a look at another movie ratings project called “The Netflix Prize”

There were fairly sophisticated rules around how yearly progress prizes would be awarded

To win a progress or grand prize a participant had to provide source code and a description of the algorithm to the jury. Following verification the winner also had to provide a non-exclusive license to Netflix.

Netflix would publish only the description, not the source code, of the system.

A team could choose to not claim a prize, in order to keep their algorithm and source code secret.

Once one of the teams succeeded to improve the RMSE by 10%



or more, the jury would issue a last call, giving all teams 30 days to send their submissions.

Only then, the team with best submission was asked for the algorithm description, source code, and non-exclusive license, and, after successful verification; declared a grand prize winner.

The contest would last until the grand prize winner was declared.

Had no one received the grand prize, it would have lasted for at least five years (until October 2, 2011).

After that date, the contest could have been terminated at any time at Netflix's sole discretion.

Two progress prizes were awarded in 2007 and 2008

Finally, on September 18, 2009, Netflix announced team "BellKor's Pragmatic Chaos" as the prize winner

#### References:

[https://en.wikipedia.org/wiki/Netflix\\_Prize](https://en.wikipedia.org/wiki/Netflix_Prize)

<https://www.techdirt.com/blog/innovation/articles/20120409/03412518422/why-netflix-never-implemented-algorithm-that-won-netflix-1-million-challenge.shtml>

## The Netflix Prize



Here's a picture of the BellKor's Pragmatic Chaos team

Their algorithm had achieved a Test RMSE of 0.8567

The team consisted of two Austrian researchers from a consulting firm,

Two researchers from AT&T Labs

A Yahoo! researcher

And two researchers from Pragmatic Theory

On December 17, 2009, four Netflix users filed a class action lawsuit against Netflix, alleging that Netflix had violated U.S. fair trade laws and the Video Privacy Protection Act by releasing the datasets.

There was public debate about privacy for research participants.

On March 19, 2010, Netflix reached a settlement with the plaintiffs, after which they voluntarily dismissed the lawsuit.



Postnote:

Netflix never implemented the winning solution!

They did, however, make use of a blend of two algorithms that members of the grand prize winning BellKor team had submitted to win the first \$50K progress prize in 2007 (they were the AT&T Labs and Yahoo! Researchers)

This submission had shown an 8.43% improvement.

The grand prize winning algorithm? Well, it wasn't worth the work it would take to get the additional marginal improvement. Plus, by then Netflix's business had shifted to online streaming. And recommendations for streaming videos is different than for rental viewing a few days later.

In other words, the shift from delayed gratification to instant gratification makes a difference in the kinds of recommendations that work

References:

[https://en.wikipedia.org/wiki/Netflix\\_Prize](https://en.wikipedia.org/wiki/Netflix_Prize)

<https://www.techdirt.com/blog/innovation/articles/20120409/03412518422/why-netflix-never-implemented-algorithm-that-won-netflix-1-million-challenge.shtml>



© 2016 MapR Technologies MAPR 64

Speaking of movies, here's one you might consider watching at some point

This one came out in 2014

It's a loose biography of **Alan Turing**

The reason I bring it up is because **Turing** pioneered the concept of **solving mathematical algorithms** using general-purpose computers

He's widely considered the father of artificial intelligence, from which machine learning has evolved

The movie takes a very liberal theatrical license with the story it tells of Turing's life, but it does demonstrate one of the first real-world use cases  
(namely, saving people's lives)



## Next Steps

- Save the model
- Run against additional test data until RMSE is acceptable
  - Training and parameter optimization can be automated
  - Can also automate periodic data enrichment
- Engine Deployment options:
  - REST Service
  - Package as a Java, Scala, or Python application
  - Implement as a real-time streaming application



Here are some potential next steps

First, you'll want to save the model

Spark makes this pretty straightforward

Next, you may do some additional test runs

A lot of this can be automated

You've also got several deployment options at your disposal

Really, it depends on your use case for the prediction model



## Demo #3: Putting It All Together

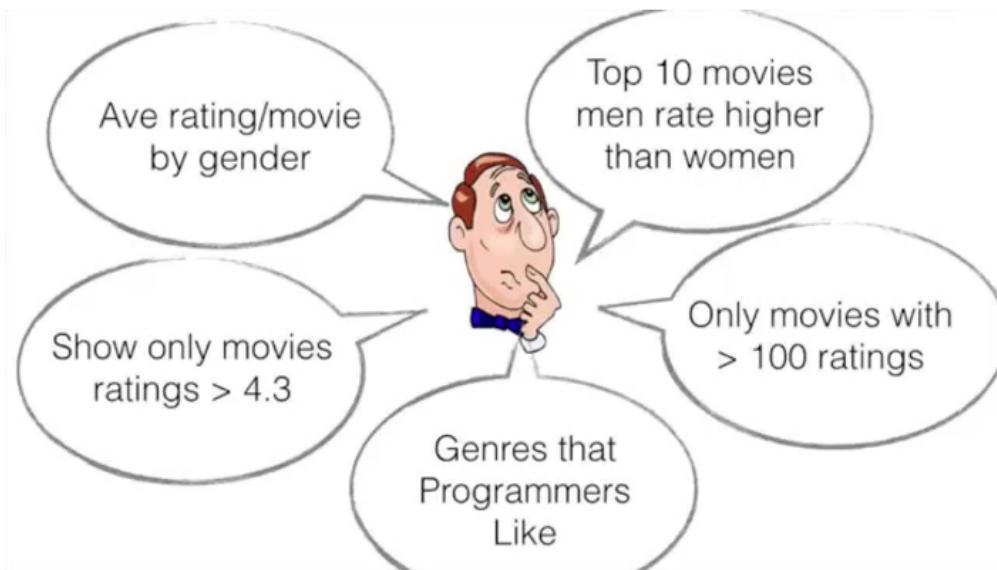


This will be a quick demonstration of how to leverage Spark SQL ODBC/JDBC connectivity to query recommendation data using familiar tools

I've staged a demo recommendations table in a Parquet file  
This table contains predicted ratings of all unwatched movies for a relatively small subset of the users (200 to be exact).



## Leveraging Spark In The Real World



© 2016 MapR Technologies **MAPR** 67

Spark is a great tool for data scientist and engineers  
We could play with this dataset all day long

But the truth is, **Spark isn't an end-user tool**  
We need to expose the datasets to the outside world  
There are some options here

I used this image from this video:  
<https://www.youtube.com/watch?v=BN4-GrZ1mGU>



## Leveraging Spark In The Real World

- Data Frames can be saved/exported
- Option #1: JSON Document Files
- Option #2: Parquet Files
  - Well-supported by many other data processing systems
  - Columnar format, can be queried using SQL
  - Compress nicely
- Option #3: Hive Metastore
  - Use `saveAsTable` command
  - Creates a pointer in the Hive metastore
  - Writes using Parquet format



First, we'll assume your data is present within the context of a Data Frame

Remember this is structured data, just like a table in a relational database

We can save Data Frame as simple JSON document files in the filesystem.

This makes them reasonably portable.

A better option is to materialize dataframes as **Parquet files**. Parquet is a **columnar format** that is supported by many other data processing systems.

The nice thing about a columnar-formatted datastore is that they are well-suited for OLAP-type workloads against very large datasets (e.g., data warehouses)

This is because the data is stored as **columns** rather than rows



Another side benefit is that they tend to **compress very nicely as well**

You can still run SQL against these files

Still, a third option is to save them to a Hive metastore using the **saveAsTable** command

An existing Hive deployment is not necessary to use this feature.

Spark will create a default local Hive metastore using **Apache Derby** for you.

The **saveAsTable** command materializes the contents of a Data Frame and **creates a pointer to the data** in the Hive metastore. If you liked the sound of the columnar **Parquet** file format then you'll be happy to know that **saveAsTable** writes using Parquet as well

I'll demonstrate this third option with Tableau as my BI tool

I used this image from this video:

<https://www.youtube.com/watch?v=BN4-GrZ1mGU>

I referred to this Wikipedia entry:

[https://en.wikipedia.org/wiki/Column-oriented\\_DBMS](https://en.wikipedia.org/wiki/Column-oriented_DBMS) and this Spark documentation: <https://spark.apache.org/docs/latest/sql-programming-guide.html>

## Saving Data Frames as Hive Metastore Tables

```
val options = Map("path" -> "/Users/Shared/Hive/movies")
moviesDF.write.options(options).mode("overwrite").saveAsTable("movies")

val options = Map("path" -> "/Users/Shared/Hive/users")
usersDF.write.options(options).mode("overwrite").saveAsTable("users")

val options = Map("path" -> "/Users/Shared/Hive/ratings")
ratingsDF.write.options(options).mode("overwrite").saveAsTable("ratings")
```



As a side note, we looked at **registerTempTable** during the first demo

Remember that this creates an in-memory table that is scoped to the cluster in which it was created.

The data is stored using Hive's highly-optimized, in-memory columnar format.

The difference here is that **saveAsTable** creates a permanent, physical table in the filesystem

Notice that I'm writing to a shared area on my local laptop  
I don't have a local Hive metastore, so as I mentioned earlier,  
Spark will create one for me using **Apache Derby**.



I referred to this document:

<https://forums.databricks.com/questions/400/what-is-the-difference-between-registertemptable-a.html>

And this Spark documentation:

<https://spark.apache.org/docs/latest/sql-programming-guide.html>

## Saving Data Frames as Hive Metastore Tables

```
val options = Map("path" -> "/Users/Shared/Hive/movies")
moviesDF.write.options(options).mode("overwrite").saveAsTable("movies")

val options = Map("path" -> "/Users/Shared/Hive/users")
usersDF.write.options(options).mode("overwrite").saveAsTable("users")

val options = Map("path" -> "/Users/Shared/Hive/ratings")
ratingsDF.write.options(options).mode("overwrite").saveAsTable("ratings")
```

## Start the Thrift JDBC/ODBC Server

```
/usr/local/Cellar/apache-spark/1.6.2/libexec/sbin/start-thriftserver.sh
```



© 2016 MapR Technologies **MAPR** 70

Spark SQL can also act as a distributed query engine  
We'll utilize its Thrift JDBC/ODBC server for this.

By default, it listens on port 10,000 and requires a local  
username (without password) for login  
You can lock it down more thoroughly if needed

I referred to this Spark documentation:  
<https://spark.apache.org/docs/latest/sql-programming-guide.html>



## Connecting to Thrift Server from Tableau

The screenshot shows the Tableau Connect interface. On the left, there's a sidebar with options like 'Connect', 'To a file', 'Text File', 'Statistical File', 'Other files', 'To a server', 'Tableau Server', 'Spark SQL', 'Microsoft SQL Server', 'MySQL', 'Oracle', 'More Servers...', 'Saved data sources', 'Sample - Superstore', and 'World Indicators'. The 'Spark SQL' option is highlighted. On the right, a 'Spark SQL' configuration dialog is open. It has fields for 'Server: localhost' and 'Port: 10000'. Below these are sections for 'Enter information to sign in to the server': 'Type: SparkThriftServer (Spark 1.1 and later)', 'Authentication: User Name', 'Username: cwarman', and 'Password:'. There are also fields for 'Realm:', 'Host FQDN:', 'Service Name:', and 'HTTP Path:'. At the bottom of the dialog are 'Initial SQL...', 'Cancel', and 'OK' buttons.

© 2016 MapR Technologies 71

Connect from Tableau:

Main "Connect" page > More Servers... > Spark SQL

Server: localhost

Port: 10000

Type: SparkThriftServer (Spark 1.1 and later)

Authentication: User Name

Username: cwarman (notice it's not crwarman)

Select default for the Schema (use magnifying glass "search" functionality)

I referred to this Spark documentation:

<https://spark.apache.org/docs/latest/sql-programming-guide.html>

## Movie Ratings Data Source

movies (default.movies)+ (default)

Connected to Spark SQL

Server  
localhost

Schema  
default

Table  
Enter table name

Exact  Contains  Starts with

- movies (default.movies)
- ratings (default.ratings)
- recommendations (d...t.recommendations)
- users (default.users)
- New Custom SQL

Join

movies —— ratings

Inner      Left      Right      Full Outer

Data Source      ratings

Movie Id = Product

Add new join clause

Sort fields      Data source order

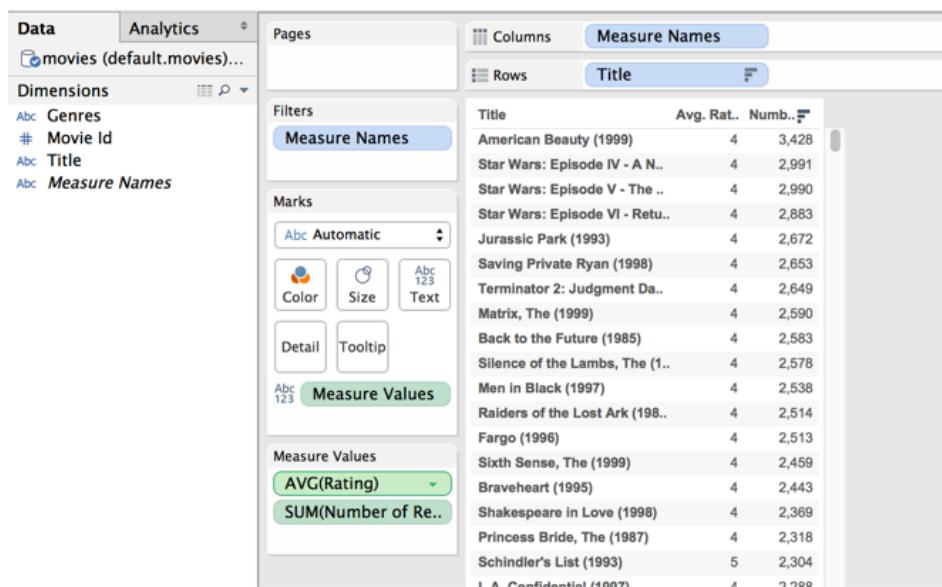
#	movies	Abc	movies	Abc	movies	#	#	#	#
	Movie Id	Title	Genres	User	Product	ratings	ratings	Rating	Rating
1	Toy Story (19...	["Animation Child...	1	1	1	5.00000			
1	Toy Story (19...	["Animation Child...	6	1	1	4.00000			
1	Toy Story (19...	["Animation Child...	8	1	1	4.00000			

© 2016 MapR Technologies  72

Open Ratings workbook if needed



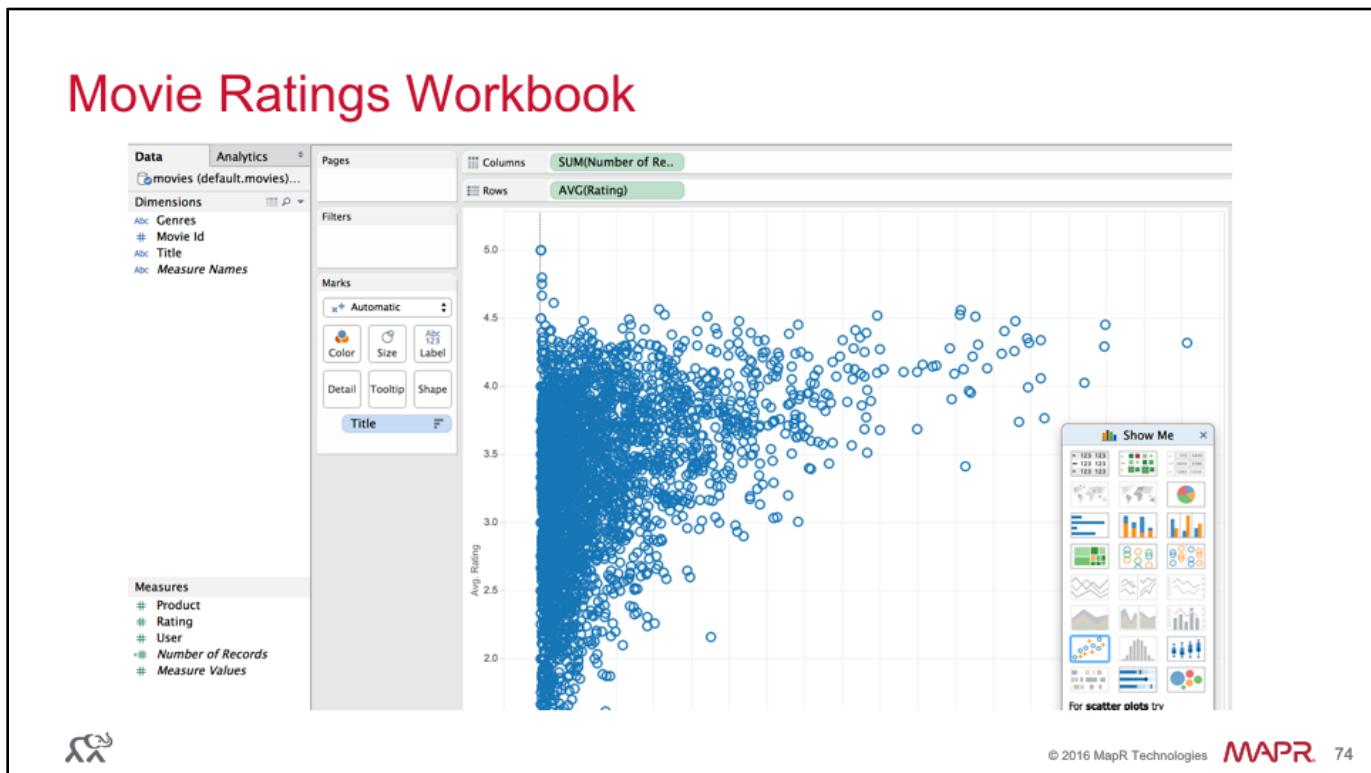
## Movie Ratings Workbook



© 2016 MapR Technologies MAPR 73

Open Ratings workbook if needed





Our own version of a “Magic Quadrant”!



## Movie Recommendations Data Source

movies (default.movies)+ (default)

Connected to Spark SQL.

**Server**  
localhost

**Schema**  
default

**Table**  
Enter table name

Exact    Contains    Starts with

- movies (default.movies)
- ratings (default.ratings)
- recommendations (d...t.recommendations)
- users (default.users)
- New Custom SQL

Join

	movies	recommendations
Movie Id	=	moviedId (recomm...

Add new join clause

**Sort fields** **Data source order**

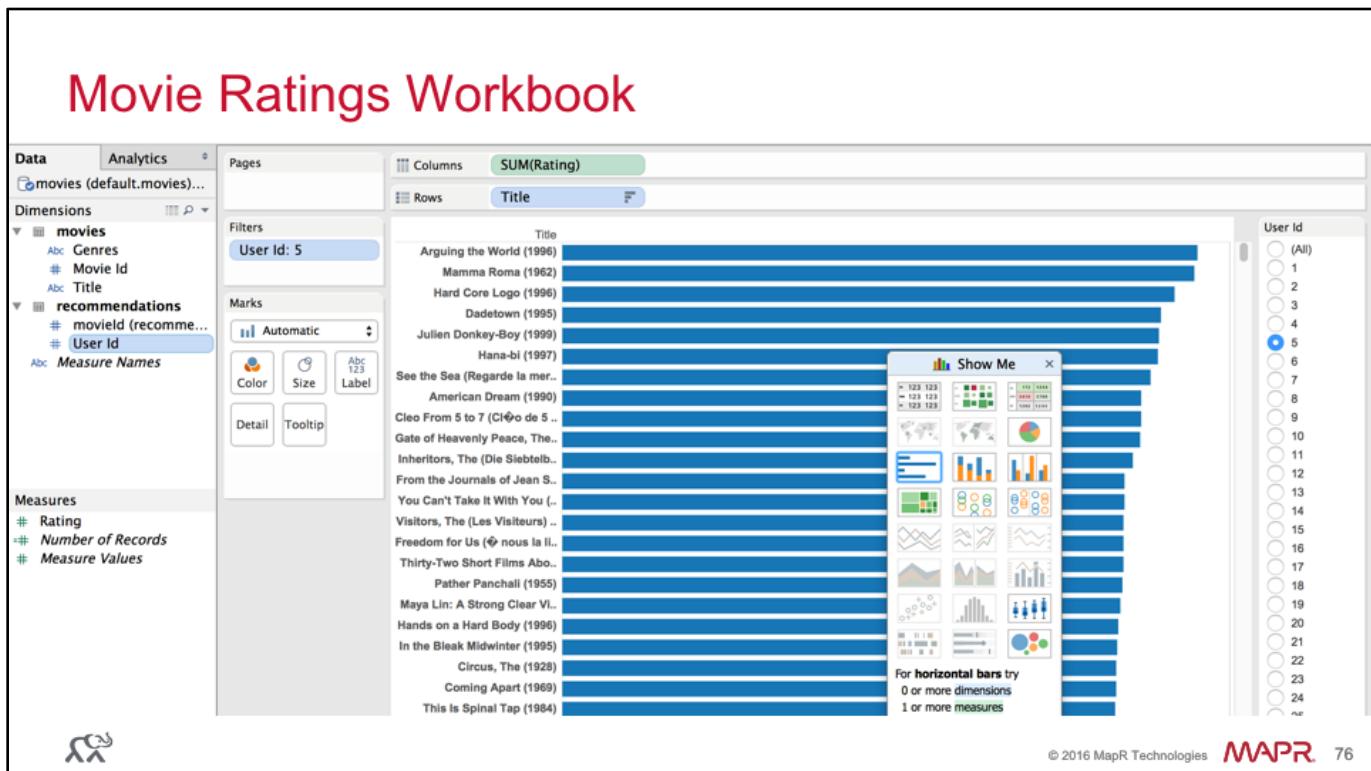
#	movies Movie Id	movies Title	movies Genres	# recommendatio... User Id	# recommendatio... moviedId (recomm...	# recommendatio... Rating
1	1	"Animation Child...		91	1	5.06898
1	1	"Animation Child...		128	1	4.83941
1	1	"Animation Child...		115	1	4.78361
1	1	"Animation Child...		43	1	4.77807



© 2016 MapR Technologies **MAPR** 75

Open Recommendations workbook if needed





Open Recommendations workbook if needed



## What We've Covered During This Presentation

- A brief Spark background/overview
- Spark SQL architecture
- How to build load and query a dataset with Spark SQL
- Basic machine learning with Spark MLlib
- Leveraging Spark in the real world



- Spark Background/Overview
  - Brief Spark background
  - The Spark+Hadoop team
  - Spark's five main components
- Spark SQL Architecture
  - Features, Languages, How DataFrames work, The SQLContext, Data sources
- Loading And Querying a Dataset with Spark SQL
  - Live demonstration of setting up a SQLContext
  - Loading it with data
  - Running queries against it
- Machine Learning with Spark MLlib
  - Collaborative filtering basics
  - Alternating Least Squares (ALS) algorithm
  - Live demo of a simple recommender model and training-test loop iterations



- How to connect to Spark SQL using ODBC/JDBC
  - Live demonstration of how to leverage Spark SQL ODBC/JDBC connectivity using Tableau
- Next steps
  - Some Real-World Use Cases
  - Basically answer the questions "What's it good for?" and "Who's using this?"
  - How to download a ready-to-use sandbox VM



# Spark SQL & Machine Learning A Practical Demonstration

Craig Warman



© 2016 MapR Technologies

