



# Apache Mahout

Scalable machine learning and data mining

# Origin of Mahout

## Origin of Mahout

- Primary goal is creating scalable and efficient machine learning algorithms

# Origin of Mahout

- Primary goal is creating scalable and efficient machine learning algorithms
- Ng et al.'s paper "**Map-Reduce for Machine Learning on Multicore**" was the driving force

# Origin of Mahout

- Primary goal is creating scalable and efficient machine learning algorithms
- Ng et al.'s paper "**Map-Reduce for Machine Learning on Multicore**" was the driving force
- Developed as a 'driver' for Hadoop

# Origin of Mahout

- Primary goal is creating scalable and efficient machine learning algorithms
- Ng et al.'s paper "**Map-Reduce for Machine Learning on Multicore**" was the driving force
- Developed as a 'driver' for Hadoop
- The latest release (Samsara), has shifted away from MapReduce

# Origin of Mahout

- Primary goal is creating scalable and efficient machine learning algorithms
- Ng et al.'s paper "**Map-Reduce for Machine Learning on Multicore**" was the driving force
- Developed as a 'driver' for Hadoop
- The latest release (Samsara), has shifted away from MapReduce
- Evolved from being a collection of algorithms, to a scala based programming environment

# Hadoop MapReduce vs Mahout on Spark



# Hadoop MapReduce vs Mahout on Spark

MapReduce

Mahout Samsara

# Hadoop MapReduce vs Mahout on Spark

## MapReduce

- Real time streaming operations unsupported

## Mahout Samsara

- In-memory storage makes streaming possible

# Hadoop MapReduce vs Mahout on Spark

## MapReduce

- Real time streaming operations unsupported
- Currently available algorithms :

Item based filtering, Matrix Factorization

*\*Most algorithms have been deprecated*

## Mahout Samsara

- In-memory storage makes streaming possible
- Currently available algorithms :

Distributed BLAS, User & Item based  
filtering, Naive Bayes Classifier, SVD, PCA,  
RowSimilarity Job

# Hadoop MapReduce vs Mahout on Spark

## MapReduce

- Real time streaming operations unsupported
- Currently available algorithms :

Item based filtering, Matrix Factorization

*\*Most algorithms have been deprecated*

- Strong security measures exist

## Mahout Samsara

- In-memory storage makes streaming possible
- Currently available algorithms :

Distributed BLAS, User & Item based filtering, Naive Bayes Classifier, SVD, PCA, RowSimilarity Job

- Security is still in its infancy

# Hadoop MapReduce vs Mahout on Spark

## MapReduce

- Real time streaming operations unsupported
- Currently available algorithms :

Item based filtering, Matrix Factorization

*\*Most algorithms have been deprecated*

- Strong security measures exist
- Java is the primary choice

## Mahout Samsara

- In-memory storage makes streaming possible
- Currently available algorithms :

Distributed BLAS, User & Item based filtering, Naive Bayes Classifier, SVD, PCA, RowSimilarity Job

- Security is still in its infancy
- Scala, Java

# Algorithms in Mahout

# Algorithms in Mahout

- Algorithms are usually classified into 3 categories
  - Collaborative Filtering - Item Recommendation used in ecommerce (Amazon)

# Algorithms in Mahout

- Algorithms are usually classified into 3 categories
  - Collaborative Filtering - Item Recommendation used in ecommerce (Amazon)
  - Clustering - K means clustering



# Algorithms in Mahout

- Algorithms are usually classified into 3 categories
  - Collaborative Filtering - Item Recommendation used in ecommerce (Amazon)
  - Clustering - K means clustering
  - Categorization/Classification - Naive Bayes, Logistic Regression

# Algorithms in Mahout

- Algorithms are usually classified into 3 categories
  - Collaborative Filtering - Item Recommendation used in ecommerce (Amazon)
  - Clustering - K means clustering
  - Categorization/Classification - Naive Bayes, Logistic Regression
- Dimensionality reduction algorithms are also provided by the Mahout Math-Scala Core library (SVD, Distributed RowMatrix)

# Algorithms in Mahout

- Algorithms are usually classified into 3 categories
  - Collaborative Filtering - Item Recommendation used in ecommerce (Amazon)
  - Clustering - K means clustering
  - Categorization/Classification - Naive Bayes, Logistic Regression
- Dimensionality reduction algorithms are also provided by the Mahout Math-Scala Core library (SVD, Distributed RowMatrix)
- We shall try to implement a variant of Collaborative Filtering algorithm - Recommendation Engine

# Installation of Mahout and Spark

- Install Java, git, maven, ssh
- Install Scala
- Install Spark
- Install Mahout
- Add JAVA\_HOME, SCALA\_HOME, SPARK\_HOME, MAHOUT\_HOME values to the path

# Installation of Mahout and Spark

- [Download](#) Mahout and extract
- Enter the folder , use **mvn -DskipTests=True clean install**

# Installation of Mahout and Spark

.bashrc file should contain the following paths

#exporting JAVA, MAHOUT, SCALA and SPARK home values

export JAVA\_HOME=/usr/lib/jvm/java-1.8.0-openjdk-amd64

export MAHOUT\_HOME=/home/username/mahout

export SCALA\_HOME=/usr/local/src/scala

export SPARK\_HOME=/home/username/spark-1.6.1

export MASTER=spark://username-VirtualBox:7077

PATH=\$JAVA\_HOME/bin:\$PATH

PATH=\$MAHOUT\_HOME/bin:\$PATH

PATH=\$SCALA\_HOME/bin:\$PATH

PATH=\$SPARK\_HOME/bin:\$PATH

export PATH

# Installation of Mahout and Spark

```
shahbaz@shahbaz-VirtualBox: ~/mahout
[INFO] --- maven-assembly-plugin:2.4.1:single (bin-assembly) @ apache-mahout-distribution ---
[INFO] Assemblies have been skipped per configuration of the skipAssembly parameter.
[INFO] --- maven-assembly-plugin:2.4.1:single (src-assembly) @ apache-mahout-distribution ---
[INFO] Assemblies have been skipped per configuration of the skipAssembly parameter.
[INFO] --- maven-install-plugin:2.5.2:install (default-install) @ apache-mahout-distribution ---
[INFO] Installing /home/shahbaz/mahout/distribution/pom.xml to /home/shahbaz/.m2/repository/org/apache/mahout/apache-mahout-distribution/0.12.1-SNAPSHOT/apache-mahout-distribution-0.12.1-SNAPSHOT.pom
[INFO] -----
[INFO] Reactor Summary:
[INFO]
[INFO] Mahout Build Tools ..... SUCCESS [ 8.936 s]
[INFO] Apache Mahout ..... SUCCESS [ 0.279 s]
[INFO] Mahout Math ..... SUCCESS [ 31.368 s]
[INFO] Mahout HDFS ..... SUCCESS [ 6.351 s]
[INFO] Mahout Map-Reduce ..... SUCCESS [ 40.610 s]
[INFO] Mahout Integration ..... SUCCESS [ 9.608 s]
[INFO] Mahout Examples ..... SUCCESS [ 27.441 s]
[INFO] Mahout Math Scala bindings ..... SUCCESS [01:24 min]
[INFO] Mahout H2O backend ..... SUCCESS [ 38.335 s]
[INFO] Mahout Spark bindings ..... SUCCESS [03:29 min]
[INFO] Mahout Flink bindings ..... SUCCESS [02:04 min]
[INFO] Mahout Spark bindings shell ..... SUCCESS [01:18 min]
[INFO] Mahout Release Package ..... SUCCESS [ 2.971 s]
[INFO] -----
[INFO] BUILD SUCCESS
[INFO] -----
[INFO] Total time: 11:04 min
[INFO] Finished at: 2016-05-02T14:16:32+02:00
[INFO] Final Memory: 79M/375M
[INFO] -----
shahbaz@shahbaz-VirtualBox:~/mahout$
```

# In-Core Algebra

- In-core means In-memory which effectively refers to a shift from MapReduce paradigm
- Two basic interfaces for in-core algebra are **Matrix**, **Vector** (Tensor types)
- Vector Implementations
  - DenseVector, RandomAccessSparseVector, SequentialAccessSparseVector
- Matrix Implementations
  - DenseMatrix, SparseRowMatrix, SparseMatrix, DiagonalMatrix



# Distributed Row Matrix

$$\mathbf{X} = \begin{bmatrix} 2 & 2 & 10.5 & 10 \\ 1 & 2 & 12 & 12 \\ 1 & 1 & 12 & 13 \\ 2 & 1 & 11 & 13 \\ 1 & 2 & 12 & 11 \\ 2 & 1 & 16 & 8 \\ 6 & 2 & 17 & 1 \\ 3 & 2 & 13 & 7 \\ 3 & 3 & 13 & 4 \end{bmatrix}$$

# Distributed Row Matrix

$$\mathbf{X} = \begin{bmatrix} 2 & 2 & 10.5 & 10 \\ 1 & 2 & 12 & 12 \\ 1 & 1 & 12 & 13 \\ 2 & 1 & 11 & 13 \\ 1 & 2 & 12 & 11 \\ 2 & 1 & 16 & 8 \\ 6 & 2 & 17 & 1 \\ 3 & 2 & 13 & 7 \\ 3 & 3 & 13 & 4 \end{bmatrix}$$

1	2	2	10.5	10
2	1	2	12	12
3	1	1	12	13
4	2	1	11	13
5	1	2	12	11
6	2	1	16	8
7	6	2	17	1
8	3	2	13	7
9	3	3	13	4

1	->	2	2	10.5	10
2	->	1	2	12	12
3	->	1	1	12	13
4	->	2	1	11	13
5	->	1	2	12	11
6	->	2	1	16	8
7	->	6	2	17	1
8	->	3	2	13	7
9	->	3	3	13	4

Index the matrix by rows

# Distributed Row Matrix

1	2	2	10.5	10
2	1	2	12	12
3	1	1	12	13
4	2	1	11	13
5	1	2	12	11
6	2	1	16	8
7	6	2	17	1
8	3	2	13	7
9	3	3	13	4

1	->	2	2	10.5	10
2	->	1	2	12	12
3	->	1	1	12	13
4	->	2	1	11	13
5	->	1	2	12	11
6	->	2	1	16	8
7	->	6	2	17	1
8	->	3	2	13	7
9	->	3	3	13	4

1	2	2	10.5	10
2	1	2	12	12
3	1	1	12	13
4	2	1	11	13
5	1	2	12	11
6	2	1	16	8
7	6	2	17	1
8	3	2	13	7
9	3	3	13	4

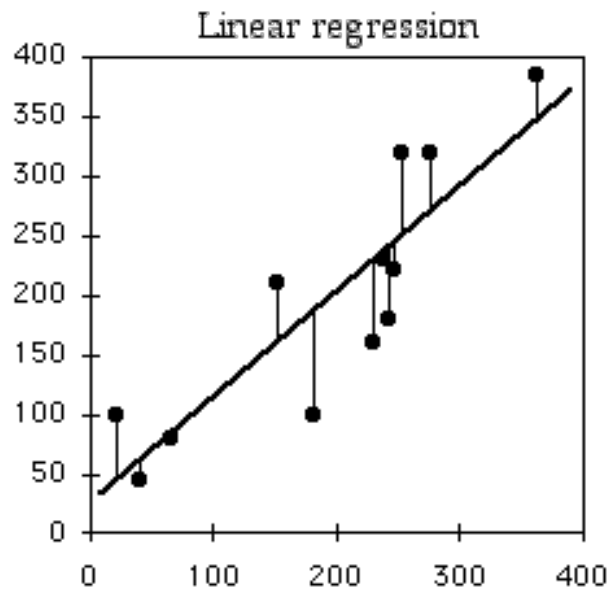
[1,2,3] ->	2	2	10.5	10
	1	2	12	12
	1	1	12	13
[4,5,6] ->	2	1	11	13
	1	2	12	11
	2	1	16	8
[7,8,9] ->	6	2	17	1
	3	2	13	7
	3	3	13	4

Blocks of the matrix can also be used as index to be distributed

## Mahout's interactive Spark shell

- Mahout is a scala DSL, which implicitly is optimized for Spark
- Lets, look at an example for performing Linear regression

# Linear Regression



# Linear Regression

- Our goal is to find an estimate of  $\beta$  in the equation  $y = \beta x + c$ , where  $y$  is the target matrix,  $x$  is the feature matrix,  $\beta$  is the parameter,  $c$  is noise, which explains the data very well
- $\beta$  can be estimated using least squares method which minimizes the sum of residual squares between the true target variable and the prediction of the target variable  $\|X\hat{w} - y\|^2$ .

We need to solve  $X^T X \hat{w} = X^T y \longrightarrow \hat{w} = (X^T X)^{-1} X^T y$ .

# Recommendation Engine

# Recommendation Engine

- An application of Collaborative filtering aspect of Machine learning



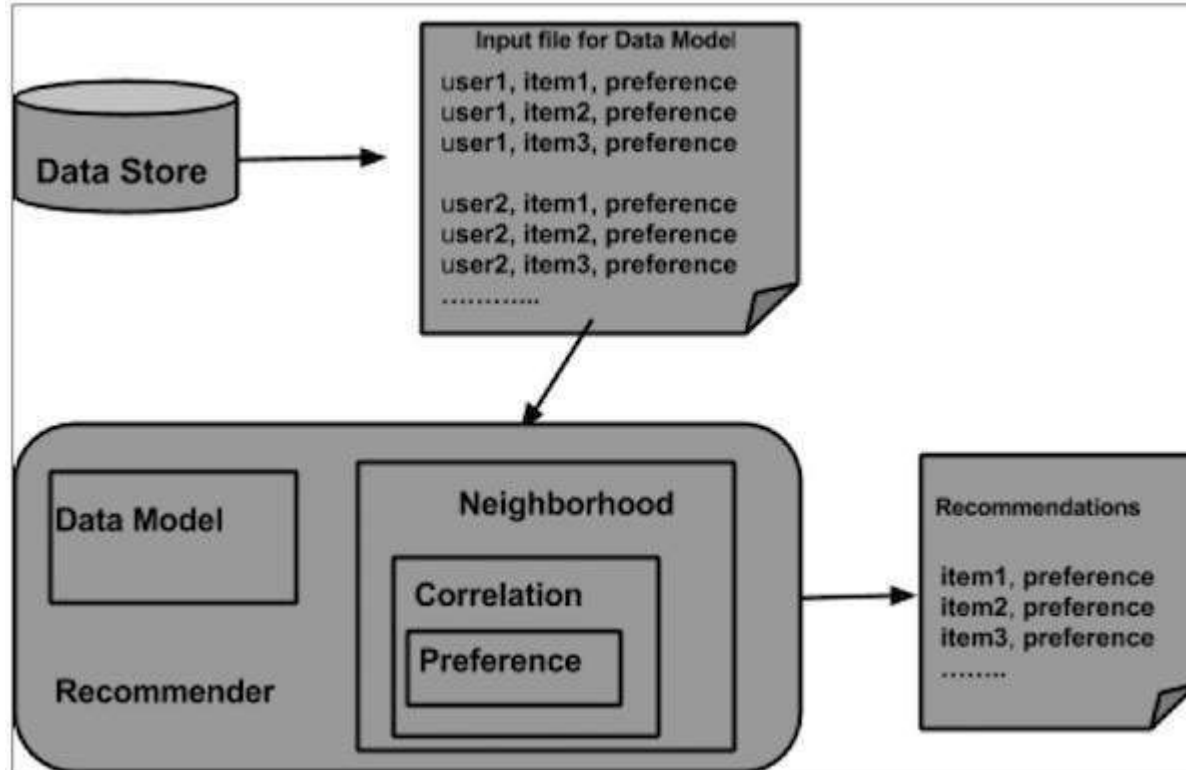
# Recommendation Engine

- An application of Collaborative filtering aspect of Machine learning
- Can be based on 3 models
  - UserSimilarity
  - ItemSimilarity
  - Model based

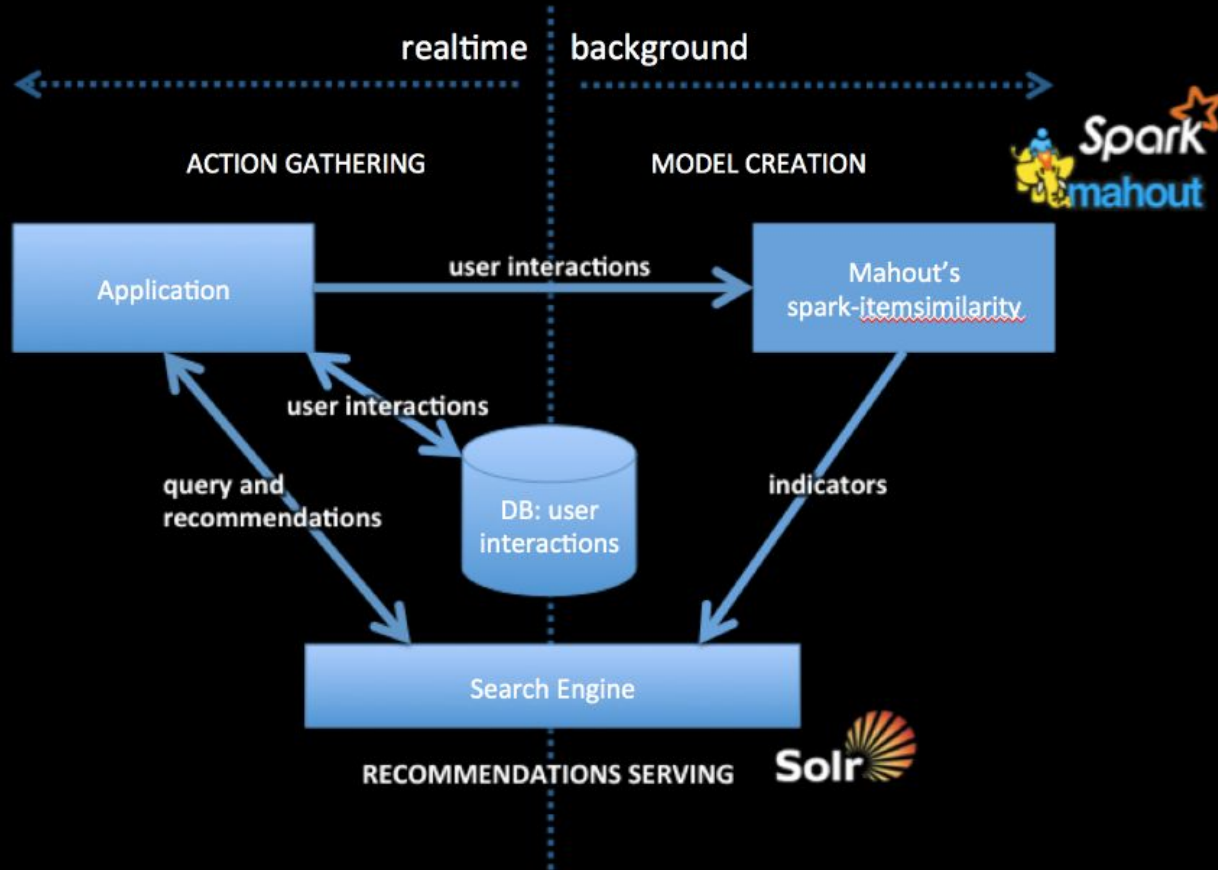
# Recommendation Engine

- An application of Collaborative filtering aspect of Machine learning
- Can be based on 3 models
  - UserSimilarity
  - ItemSimilarity
  - Model based
- Key components provided by Mahout for building recommendations are :
  - DataModel, UserSimilarity, ItemSimilarity, Recommender

# High level view of a Recommendation Engine



# Recommender Architecture



## Under the hood - Item Recommendation

- Given a set of items someone is known to like, recommend a new set of items which are similar to them

## Under the hood - Item Recommendation

- Given a set of items someone is known to like, recommend a new set of items which are similar to them
- Two items can be regarded as similar if two different people like both of them

# Under the hood - Item Recommendation

- Given a set of items someone is known to like, recommend a new set of items which are similar to them
- Two items can be regarded as similar if **two different people like both of them**

## Algorithm

```
for every item  $i$  that  $u$  has no preference for yet
  for every item  $j$  that  $u$  has a preference for
    compute a similarity  $s$  between  $i$  and  $j$ 
    add  $u$ 's preference for  $j$ , weighted by  $s$ , to a running average
return top items, ranked by weighted average
```

## Under the hood - Item Recommendation

- Item similarity can be measured by Matrix multiplication operations
- Multiply item-item similarity matrix, with item-user vector

$$\begin{matrix} \begin{bmatrix} a & b \\ c & d \end{bmatrix} & \times & \begin{bmatrix} x \\ y \end{bmatrix} & = & \begin{bmatrix} xa + yb \\ xc + yd \end{bmatrix} \\ \text{ItemSim} & & \text{UserPref} & & \end{matrix}$$



spark-itemsimilarity

## spark-itemsimilarity

- Spark counterpart of Mapreduce **itemsimilarity**

## spark-itemsimilarity

- Spark counterpart of Mapreduce **itemsimilarity**
- Takes in elements of user interactions(userid, itemid, value), and outputs indicator matrices by comparing user-user interactions

## spark-itemsimilarity

- Spark counterpart of Mapreduce **itemsimilarity**
- Takes in elements of user interactions(userid, itemid, value), and outputs indicator matrices by comparing user-user interactions
- Indicator matrix is an item-item matrix where values are **log-likelihood** ratio strengths

## spark-itemsimilarity

- Spark counterpart of Mapreduce **itemsimilarity**
- Takes in elements of user interactions(userid, itemid, value), and outputs indicator matrices by comparing user-user interactions
- Indicator matrix is an item-item matrix where values are **log-likelihood** ratio strengths
- Extends co-occurrence to cross-co-occurrence by allowing multi-model interactions

## Conclusion

- Apache Mahout fully utilizes the Spark framework to move away from the fallbacks of MapReduce
- It also supports other backend frameworks such as H2o, Flink

Thank you