

Introduction to Machine Learning With Spark

DataSummit 2016

May 9, 2016

Nathan Halko @nhalko

Data Scientist, @SpotRight

All code and slides are on github MahoutApp

The Stack

Apps, UserInterface

API's, webservice

Languages:
Scala, Java, Python

Client Code (libraries, packages):
Mahout, *MLLib*, Matlab?

Framework (execution engine):
Spark, Hadoop

Database:
Cassandra, HDFS, SQL, filesystem

Hardware (servers):
Laptop, Cloud (aws)



spotright

Nathan Halko, Data Scientist

n8halko@gmail.com

Information & Insight

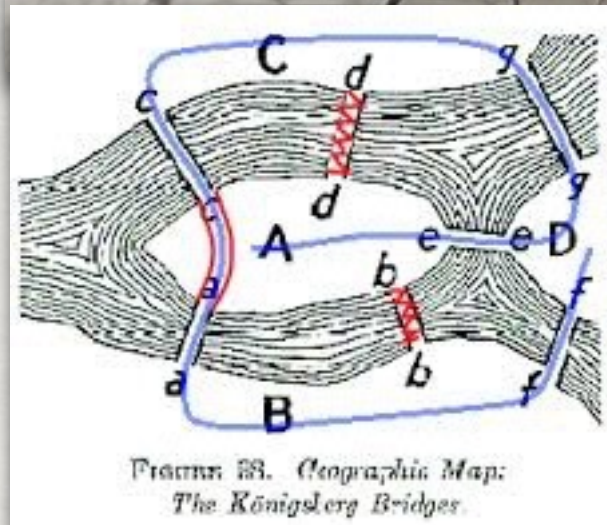


FIGURE 33. Geographic Map:
The Königsberg Bridges.

Raw Unstructured Data

Algorithms

Data Science

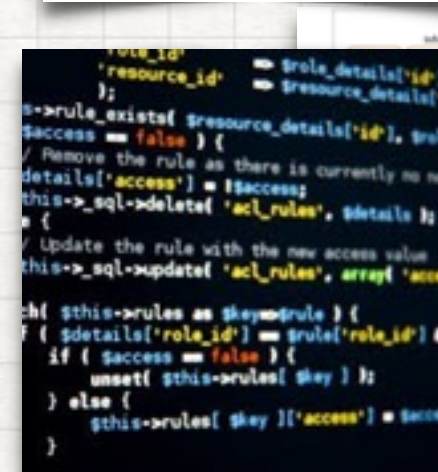
Clients

Sales/Marketing

Product

Engineering

Infrastructure/IT



Machine Learning

Putting some *things* into some *buckets*.

Supervised: you know the names of the buckets.

Unsupervised: you don't.

Tutorials and Credits

1. Decision Tree **classification** for predicting flight delays.

Carol McDonald, MapR

tutorial

2. **Clustering** the News with *unsupervised* kmeans

Dan McClary, Big Data Science Bootcamp

tutorial

3. **Collaborative Filtering** with MovieLens

Guy Ernest, AWS

DataBricks and Grouplens

tutorial

Apache Zeppelin

Download 0.5.6-incubating Binary package [here](#)

Instructions for install

```
> tar -xzvf ~/Downloads/zeppelin-0.5.6-incubating-bin-  
all.tgz  
> ln -s zeppelin-0.5.6-incubating-bin-all.tgz zeppelin  
> cd zeppelin
```

Lets give it more memory:

```
> cp conf/zeppelin-env.sh{.template,}
```

add to zeppelin-env.sh

```
export ZEPPELIN_PORT=8082      *** optional ***  
export ZEPPELIN_MEM=-Xmx4G
```

```
> bin/zeppelin-daemon.sh start
```

Go to <http://localhost:8080>

**** Use a new note per exercise and multiple paragraphs within for easy re-compilation ****

Collaborative Filtering with MovieLens

::credit::

Guy Ernest, Solutions Architect with AWS

tutorial

Code/data/material from
DataBricks/spark-training

Download datasets

```
> curl -O https://raw.githubusercontent.com/databricks/spark-training/master/data/movielens/medium/movies.dat
```

```
> curl -O https://raw.githubusercontent.com/databricks/spark-training/master/data/movielens/medium/ratings.dat
```

Create a sample set for easy debugging

```
> head -10000 ratings.dat > ratings_sample.dat
```

Skip down to the *"Build the Recommender with SparkML"* section of the tutorial, the basic steps are as follows:

Load data in Spark

Partition: Training 60%, Validation 20%, Test 20%

Define *RMSE* to evaluate model performance

Train the model and select the best fit

Evaluate on the Test data set

Recommend!


```

1 import java.io.File
2 import scala.io.Source
3
4 import org.apache.log4j.Logger
5 import org.apache.log4j.Level
6
7 import org.apache.spark.SparkConf
8 import org.apache.spark.SparkContext
9 import org.apache.spark.SparkContext._
10 import org.apache.spark.rdd._
11 import org.apache.spark.mllib.recommendation.{ALS, Rating, MatrixFactorizationModel}
12
13 /**
14  * Load the data into Spark
15  */
16 val movieLensHomeDir = "/Users/nhalko/Documents/NYC_2016/datasets/movielens"
17
18 val movies = sc.textFile(s"${movieLensHomeDir}/movies.dat").map { line =>
19   // MovieID::Title::Genres
20   // 1::Toy Story (1995)::Animation|Children's|Comedy
21   val Array(id, name, _) = line.split("::", 3)
22   id.toInt -> name
23 }.collect.toMap
24
25 val ratings = sc.textFile(s"${movieLensHomeDir}/ratings.dat").map { line =>
26   // format: (timestamp % 10, Rating(userId, movieId, rating))
27   // UserID::MovieID::Rating::Timestamp
28   // 1::1193::5::978300760
29   val Array(userId, movieId, rating, timeStamp) = line.split("::", 4)
30
31   (timeStamp.toLong % 10, Rating(userId.toInt, movieId.toInt, rating.toDouble))
32 }
33
34 /**
35  * Partition into Training 60%, Validation 20%, Test 20%
36  */
37
38 val training = ratings.filter{case (ts, _) => ts < 6}.values.cache()
39 val validation = ratings.filter{case (ts, _) => ts >= 6 && ts < 8}.values.cache()
40 val test = ratings.filter{case (ts, _) => ts >= 8}.values.cache()
41
42 println(s"""%table DataSet\tCount
43 Training\t${training.count()}
44 Validation\t${validation.count()}
45 Test\t${test.count()}""")

```

Load data in Spark

Partition: Training 60%, Validation 20%, Test 20%

RMSE

Train the model

Pick the best

```
1 /**
2  * Define the RMSE for evaluating model performance
3  */
4 def computeRmse(model: MatrixFactorizationModel, data: RDD[Rating]): Double = {
5
6     val predictions: RDD[Rating] = model.predict( data.map {r => r.user -> r.product} )
7
8     val predictionsAndRatings = predictions
9     .map {p => (p.user, p.product) -> p.rating}
10    .join(data.map {r => (r.user, r.product) -> r.rating})
11    .values
12
13    math.sqrt(
14        predictionsAndRatings
15        .map{ case (pred, actual) => math.pow(pred - actual, 2) }
16        .mean
17    )
18 }
19
20 /**
21  * Train the model, find the best one and evaluate performance
22  */
23
24 case class ModelResult(
25     model: Option[MatrixFactorizationModel] = None,
26     rmse: Double = Double.MaxValue,
27     rank: Int = 0,
28     lambda: Double = -1.0,
29     numIter: Int = -1
30 )
31
32
33 val modelResults = for {
34     rank <- List(8, 12)
35     lambda <- List(0.1, 10.0)
36     numIter <- List(10) // List(10, 20)
37 } yield {
38     val model = ALS.train(training, rank, numIter, lambda)
39     val validationRmse = computeRmse(model, validation)
40
41     ModelResult(Some(model), validationRmse, rank, lambda, numIter)
42 }
43
44 // find the model with minimum error
45 val best = modelResults.minBy(_.rmse)
46 val bestModel = best.model.get
47 // display results
48 println(s"""%table RMSE (validation)\trank\tlambda\tnumIters
49 ${modelResults.map(mr => s"${mr.rmse}\t${mr.rank}\t${mr.lambda}\t${mr.numIter}").mkString("\n")}""")
50
```


Evaluate and compare against naive baseline

```
1 // test it on the test data
2 val testRmse = computeRmse(bestModel, test)
3 println(s"The best model was trained with rank = ${best.rank} and lambda = ${best.lambda} and numIter = ${best.numIter} and its RMSE on the test set is $testRmse")
4
5 // create a naive baseline and compare it with the best model
6 val meanRating = training.union(validation).map(_._rating).mean
7 val baselineRmse = math.sqrt(test.map {r => math.pow(meanRating - r._rating, 2)}.mean)
8 val improvement = (baselineRmse - testRmse) / baselineRmse * 100
9
10 println("The best model improves the baseline by " + "%.2f".format(improvement) + "%.")
11
12 // save for use in making recommendations (output dir must not exist)
13 //bestModel.save(sc, s"$movieLensHomeDir/model")
14
```

```
1 /**
2  * Create a method to make the recommendations
3  */
4 val moviesWithGenres = sc.textFile(s"$movieLensHomeDir/movies.dat").map { line =>
5   // MovieID::Title::Genres
6   // 1::Toy Story (1995)::Animation|Children's|Comedy
7   val Array(id, _, genres) = line.split("::", 3)
8   id.toInt -> genres
9 }.collect.toMap
10
11
12 def recommend(userId: Int, genre: String = "") = {
13   val candidates = sc.parallelize(
14     moviesWithGenres.filter {case (id, genres) => genres.toLowerCase.contains(genre.toLowerCase)}.keys.toSeq
15   )
16
17   val userRatings = ratings
18     .filter {case (_, r) => r.user == userId}
19     .map {case (_, r) => movies(r.product) -> r.rating}
20     .collect()
21
22   val recs = bestModel
23     .predict(candidates.map {movieId => userId -> movieId})
24     .collect()
25     .sortBy(_._rating)
26     .take(5)
27     .map {r =>
28       movies(r.product)
29     }
30
31   (recs, userRatings)
32 }
```

Recommend!


```

1
2 /**
3  * Make some recommendations
4  */
5
6 val genres = moviesWithGenres.values.toList.flatMap(_.split("\\|")).distinct.map(g => (g,g))
7 val genre = z.select("genre", genres)
8 val userId = z.input("userId", "100").asInstanceOf[String]
9
10 val (recs, userRatings) = recommend(userId.toInt, genre.asInstanceOf[String])
11 println(s"\n\nRecommended for user: \n${recs.mkString("\n")}")
12 println(s"\n\nBased on:\n${userRatings.mkString("\n")}")
13
14
15
16

```

genre

Sci-Fi

:

userId

10

Recommended for user:

Star Wars: Episode IV - A New Hope (1977)

Matrix, The (1999)

Star Wars: Episode V - The Empire Strikes Back (1980)

Visitors, The (Les Visiteurs) (1993)

Galaxy Quest (1999)

Based on:

(Midsummer Night's Dream, A (1999),5.0)

(Mission: Impossible (1996),4.0)

(Star Wars: Episode I - The Phantom Menace (1999),3.0)

(Defending Your Life (1991),5.0)

(Breaking Away (1979),3.0)

(Truman Show, The (1998),5.0)

(Prophecy II, The (1998),4.0)

(Alien⁰ (1992),3.0)

Create a little recommender gui

Decision Tree for predicting flight delays

::credit::

Carol McDonald, SOLUTIONS ARCHITECT, MAPR

tutorial

spark mllib-decision-tree

Get dataset

<https://github.com/caroljmcdonald/sparkmldecisiontree/blob/master/data/rita2014jan.csv.zip>

click View Raw to download

```
> unzip rita2014jan.csv.zip
```

Load data into Spark

Transform data into LabeledPoints

Partition data, Training 70% Test 30%

Train model

Evaluate

Imports

```
1 import scala.util.Try
2
3 import org.apache.spark._
4 import org.apache.spark.rdd.RDD
5 // Import classes for MLlib
6 import org.apache.spark.mllib.regression.LabeledPoint
7 import org.apache.spark.mllib.linalg.Vectors
8 import org.apache.spark.mllib.tree.DecisionTree
9 import org.apache.spark.mllib.tree.model.DecisionTreeModel
10 import org.apache.spark.mllib.util.MLUtils
```

Load data into Spark

```

1 // define the Flight Schema
2 case class Flight(dofM: String, dofW: String, carrier: String, tailnum: String, flnum: Int, org_id: String,
3   |           |           | origin: String, dest_id: String, dest: String, crsdeptime: Double, deptime: Double, depdelaymins: Double,
4   |           |           | crsarrrtime: Double, arrtime: Double, arrdelay: Double, crselapsedtime: Double, dist: Int)
5
6 // function to parse input into Flight class
7 def parseFlight(str: String): Flight = {
8   val line = str.split(",")
9   Flight(line(0), line(1), line(2), line(3), line(4).toInt, line(5), line(6), line(7), line(8),
10     line(9).toDouble, line(10).toDouble, line(11).toDouble, line(12).toDouble, line(13).toDouble,
11     line(14).toDouble, line(15).toDouble, line(16).toInt)
12 }
13
14 // parse the RDD of csv lines into an RDD of flight classes
15 val linesCnt = sc.accumulator(0L, "total lines")
16 val flightsRDD = sc.textFile("/Users/nhaliko/Documents/NYC_2016/datasets/flights/rita2014jan.csv")
17   .flatMap { line =>
18     |   linesCnt += 1
19     |   Try { parseFlight(line) }.toOption
20   }
21   .cache()
22
23 // use .count to trigger the computation and tick the counters
24 val validCnt = flightsRDD.count()
25 val badCnt = linesCnt.value - validCnt
26 println(s"$badCnt ${"%".4f".format(badCnt / linesCnt.value.toDouble * 100)}% bad lines out of ${linesCnt.value} total.")

```


Map all data to Doubles

```
29 /**
30  * Transform non-numeric data into numeric values
31  */
32
33 val carrierMap: Map[String, Int] = flightsRDD.map(_._1).distinct.collect()
34   .zipWithIndex
35   .toMap
36
37 val originMap: Map[String, Int] = flightsRDD.map(_._2).distinct.collect
38   .zipWithIndex
39   .toMap
40
41 val destMap: Map[String, Int] = flightsRDD.map(_._3).distinct.collect
42   .zipWithIndex
43   .toMap
44
45 /**
46  * Define the features array as a LabeledPoint. Select only relevant features
47  */
48 val mldata = flightsRDD.map { flight =>
49   LabeledPoint(
50     if (flight.depdelaymins > 40) 1.0 else 0.0,
51     Vectors.dense(
52       flight.dofM.toDouble - 1,           // day of month
53       flight.dofW.toDouble - 1,           // day of week
54       flight.crsdeptime,                   // departure time
55       flight.crsarrtime,                   // arrival time
56       carrierMap(flight.carrier).toDouble, // carrier
57       flight.crselapsedtime,               // elapsed flight time
58       originMap(flight.origin).toDouble,  // departure city
59       destMap(flight.dest).toDouble       // arrival city
60     )
61   )
62 }.cache()
```

Create *LabeledPoints* from
relevant features

Notice *continuous* vs *discrete*
features


```

1 /**
2  * Partition the data into training and test data sets
3  */
4  // mldata0 is %85 not delayed flights
5  val mldata0 = mldata.filter(x => x.label == 0).randomSplit(Array(0.85, 0.15))(1)
6  // mldata1 is %100 delayed flights
7  val mldata1 = mldata.filter(x => x.label != 0)
8  // mldata2 is delayed and not delayed
9  val mldata2 = mldata0 ++ mldata1
10
11 // split mldata2 into training and test data
12 val splits = mldata2.randomSplit(Array(0.7, 0.3))
13 val (trainingData, testData) = (splits(0), splits(1))

```

Partition data
Training 70%
Test 30%

```

1 /**
2  * Define some parameters and train the model.
3  */
4
5 // airity of non-continuous features
6 val categoricalFeaturesInfo = Map(
7   0 -> 31,
8   1 -> 7,
9   4 -> carrierMap.size,
10  6 -> originMap.size,
11  7 -> destMap.size
12 )
13
14 val numClasses = 2 // 0,1
15 val impurity = "gini" // or 'entropy'
16 val maxDepth = 3 //9 // depth of the decision tree
17 val maxBins = 500 //7000 // discretization of continuous features
18
19 // call DecisionTree trainClassifier with the trainingData , which returns the model
20 val model = DecisionTree.trainClassifier(trainingData, numClasses, categoricalFeaturesInfo,
21 impurity, maxDepth, maxBins)
22
23
24 // print out the decision tree
25 println(model.toDebugString)

```

Train the model

Evaluate

```

1 /**
2  * Evaluate model on test instances and compute test error
3  */
4
5
6 val obsCnt = sc.accumulator(0.0, "observation count")
7 val wrongPred = testData.filter { point =>
8   obsCnt += 1
9   point.label != model.predict(point.features)
10 }.count()
11
12 val ratioWrong = wrongPred / obsCnt.value
13 println(s"The model is wrong ${"%0.4f".format(ratioWrong * 100)}% of the time.")

```


Clustering the News with unsupervised k-means

Dan McClary, Big Data Science Bootcamp
tutorial

Get the data from datahub.io

```
> curl -O 'https://ckanet-storage.commondatastorage.googleapis.com/2014-05-14T08:57:08.036Z/enwiki-20140116-xml.zip'
> unzip enwiki-20140116-xml.zip
```

Create just file of titles (poor man's xml parser)

```
> cat enwiki-20140116.xml | egrep "<title>" | sed -e 's/.*<title> //' -e 's/<\/title> //' >
wikinews_titles.txt
```

We'll be skipping the json stuff in the tutorial in favor of this brute force method.
The the basic steps are:

Setup with some helper functions.

Load the data and create an 'entry' for each word.

Create a Word2Vec model

Create title vectors

Train the kmeans clustering model

See what topics the titles belong to.


```

1 import scala.util.Try
2 import org.apache.spark.rdd.RDD
3 import org.apache.spark.mllib.clustering.KMeans
4 import org.apache.spark.mllib.feature.Word2Vec
5 import org.apache.spark.mllib.feature.Word2VecModel
6 import org.apache.spark.mllib.linalg._

```

Setup with some helper functions.
Load the data and create an 'entry'
for each word.

```

1 /**
2  * Define some helper functions for later use
3  */
4
5 // vector addition
6 def sumArray(m: Array[Double], n: Array[Double]): Array[Double] = m.zip(n).map {case (i,j) => i + j}
7
8 // scalar multiplication 1/a
9 def divArray(m: Array[Double], divisor: Double) : Array[Double] = m.map(_ / divisor)
10
11 def wordToVector(w: String, m: Word2VecModel): Vector = {
12   Try {
13     m.transform(w)
14   } getOrElse {
15     println(s"$w failed to transform")
16     Vectors.zeros(100)
17   }
18 }

```

```

1 /**
2  * ... Seq(President, of, China, lunches, with, Brazilian, President)
3  *     Seq(Palestinians, to, elect, new, president, on, January, 9) ...
4  */
5 val news_titles: RDD[Seq[String]] = sc.textFile("/Users/nhalko/Documents/NYC_2016/datasets/wikinews/wikinews_titles.txt")
6   .map(_ .split(" ").toSeq)
7   .cache()
8
9 /**
10  * ... Seq(President)
11  *     Seq(of)
12  *     Seq(China) ...
13  */
14 val news_titles_words: RDD[Seq[String]] = news_titles.flatMap {
15   words =>
16     words.map(word => Seq(word))
17 }.cache()
18

```


Create a Word2Vec model

```
1 // fit a model that can compute synonyms
2 val word2vec = new Word2Vec()
3 val model = word2vec.fit(news_titles_words)
```

Notice we are skipping adding extra words to the titles

```
1 /**
2  * Create a title vector from each word vector by taking
3  * an average. So the title vector is the average vector of all
4  * its words.
5  */
6 val title_vectors = news_titles.map { title =>
7     val tVec = title
8         .map(word => wordToVector(word, model).toArray)
9         .reduceLeft(sumArray)
10        .map(_ / title.length)
11    new DenseVector(tVec).asInstanceOf[Vector]
12 }.cache()
13
14
15 // Make a tuple of (actualTitle, title_vector)
16 val title_pairs = news_titles.zip(title_vectors)
```

Create title vectors by averaging the vectors for each word in the title

Train the kmeans clustering model

```
1 val numClusters = 100
2 val numIterations = 25
3
4 val clusters = KMeans.train(title_vectors, numClusters, numIterations)
5
6 // Evaluate clustering by computing Within Set Sum of Squared Errors
7 val wssse = clusters.computeCost(title_vectors)
```



```

1 // predict the topic for each title
2 val article_membership: RDD[(Int, String)] = title_pairs.map { case (titleSeq, titleVector) =>
3     // topicId -> title String
4     clusters.predict(titleVector) -> titleSeq.mkString(" ")
5 }.cache()
6
7 // number each cluster center
8 val cluster_centers: RDD[(Int, Vector)] = sc.parallelize(
9     clusters.clusterCenters.zipWithIndex.map {case (center, idx) => idx -> center}
10 )
11
12 // synonyms are an Array of length 5 with (word, distance) tuples
13 // Map away the 'distance' and concatenate to form the cluster topic
14 val cluster_topics: RDD[(Int, String)] = cluster_centers.map {case (idx, center) =>
15     idx -> model.findSynonyms(center, 5)
16     .sortBy {case (word, dist) => dist} // sort by the distance to put the best fits first
17     .map {case (word, dist) => word}
18     .mkString(" ")
19 }
20 println(s"${cluster_topics.take(3).mkString("\n")}")

```

Assign each title to a topic

```

1 val sample_topics = cluster_topics.collect
2 def sample_members(topicId: Int) = {
3     article_membership
4     .filter {case (id, title) => id == topicId}
5     .map {case (_, title) => title}
6     .take(25)
7 }
8
9 val topic = z.select("topic", sample_topics.map {case (id, title) => id.toString -> title})
10
11 println(s"%table Related Articles:\n${sample_members(topic.toString.toInt).mkString("\n")}")
12

```

topic

actress 76 aged wrestler poet

Make a gui to show some titles for a selected topic

```

1 val topicMap = sample_topics.toMap
2
3 def topicMe(title: String) = {
4     val tVec = {
5         val tSeq = title.split(" ")
6         val tArr = tSeq
7             .map(word => wordToVector(word, model).toArray)
8             .reduceLeft(sumArray)
9             .map(_ / tSeq.length)
10
11         new DenseVector(tArr).asInstanceOf[Vector]
12     }
13
14     val topicId = clusters.predict(tVec)
15
16     topicMap(topicId)
17 }
18
19 println(s"\n\nBelongs to topic: ${topicMe(z.input("Type some words:", "Here comes the sun").toString)}")

```

Type some words:

spark machine learning is fun

See where some random titles might fall!

Next steps:

- integrate into an sbt project with editor to explore more source code (*done*)
- Code dive into algorithms for optimizations:
 - .cache() vs .persist(DISK) *speed*
 - local vs distributed *memory*
- Do I *need* to break up the algorithms or replace/optimize anything?

Mahout

<https://mahout.apache.org/> , <https://github.com/apache/mahout>, [mahout-samsara-book](#) , <https://github.com/andrewpalumbo/mahout-samsara-book>

Spark:: <http://spark.apache.org/> , <http://spark.apache.org/docs/latest/quick-start.html>

Other:: <https://git-scm.com/> , <https://maven.apache.org/>

Git Mahout and Spark (from download):

```
project> git clone https://github.com/apache/mahout.git; cd mahout
project/mahout> echo "this is your project directory"; cd ..
project> tar -xzf ~/Downloads/spark-1.5.2-bin-hadoop2.6.tgz
project> ln -s spark-mahout spark-1.5.2-bin-hadoop2.6; cd spark-mahout
project/spark-mahout> bin/spark-shell --master local[4]
```

Try it out!

```
spark> val rdd = sc.parallelize(List.fill(1000000)(scala.util.Random.nextDouble))
spark> rdd.count()
```

SparkUI -> localhost:4041

Setup Mahout (from project/mahout directory):

```
- create a file named setup with:
export MAHOUT_HOME=/project/mahout
export MAHOUT_LOCAL=true
export MAHOUT_OPTS="-Xmx4g"
export SPARK_HOME=/project/spark-mahout
> source setup # activate the environment variables
> mvn clean install -DskipTests
> MASTER=local[4] bin/mahout spark-shell
```

Try it out!

```
mahout> val mtxA = Matrices.symmetricUniformView(5000, 5000, 1234)
mahout> mtxA.rowSums
```


Creating a project

<https://www.jetbrains.com/idea/> , <http://www.scala-sbt.org/>

Create the bones:

```
> mkdir MahoutApp; cd MahoutApp
MahoutApp> mkdir -p src/{main,test}/scala/com/nhalko/mahoutapp
```

Create a *build.sbt* file with contents:

```
lazy val root = (project in file(".")).
settings(
  name := "mahoutapp",
  version := "0.1",
  scalaVersion := "2.10.4"
))
```

Open the project in IDEA and follow the wizard.

Add dependencies and some code:

Update build.sbt to look like [this](#)

Create file src/test/scala/com/nhalko/mahoutapp/[RidgeRegression.scala](#)

Run it or play with it in the repl!

```
> sbt "testOnly *RidgeRegression"
> sbt test:console
```

Create and run (from the shell) *.mscala scripts*

```
mahout> :load my_script.mscala
```

Use your .jar in Mahout shell

```
MahoutApp> sbt package # create a .jar in target/scala-2.10 dir of your project
```

```
mahout> :cp path/to/your/MahoutApp.jar // access your project's classes and methods from mahout shell
```