

# Implementing Predictive Analytics with Hadoop in Azure HDInsight

Lab 3 – Building Machine Learning Models

## Overview

In this lab, you will use Spark to create machine learning models.

## What You'll Need

To complete the labs, you will need the following:

- A web browser
- A Microsoft account
- A Microsoft Azure subscription
- A Windows, Linux, or Mac OS X computer
- The lab files for this course

**Note:** To set up the required environment for the lab, follow the instructions in the [Setup](#) document for this course. Specifically, you must have signed up for an Azure subscription and installed the Azure CLI tool.

## Provisioning an HDInsight Spark Cluster

In this exercise, you will create a Spark cluster.

**Note:** The Microsoft Azure portal is continually improved in response to customer feedback. The steps in this exercise reflect the user interface of the Microsoft Azure portal at the time of writing, but may not match the latest design of the portal exactly.

## Provision an HDInsight Cluster

**Note:** If you already have a Spark HDInsight cluster running, you can skip this procedure

1. In a web browser, navigate to <http://portal.azure.com>, and if prompted, sign in using the Microsoft account that is associated with your Azure subscription.
2. In the Microsoft Azure portal, add a new HDInsight cluster with the following settings:
3. In the Microsoft Azure portal, in the Hub Menu, click **New**. Then in the **Data + Analytics** menu, click **HDInsight**.
4. In the **New HDInsight Cluster** blade, enter the following settings, and then click **Create**:

- **Cluster Name:** Enter a unique name (and make a note of it!)
  - **Subscription:** Select your Azure subscription
  - **Select Cluster Type:**
    - **Cluster Type:** Spark
    - **Operating System:** Linux
    - **Version:** Select the latest available version of Spark
    - **Cluster Tier:** Standard
  - **Applications:** None
  - **Credentials:**
    1. **Cluster Login Username:** Enter a user name of your choice (and make a note of it!)
    2. **Cluster Login Password:** Enter and confirm a strong password (and make a note of it!)
    - **SSH Username:** Enter another user name of your choice (and make a note of it!)
    - **SSH Authentication Type:** Password
    - **SSH Password:** Enter and confirm a strong password (and make a note of it!)
  - **Data Source:**
    - **Create a new storage account:** Enter a unique name consisting of lower-case letters and numbers only (and make a note of it!)
    - **Choose Default Container:** Enter the cluster name you specified previously
    - **Location:** Select any available region
  - **Pricing:**
    - **Number of Worker nodes:** 1
    - **Worker Nodes Pricing Tier:** Use the default selection
    - **Head Node Pricing Tier:** Use the default selection
  - **Optional Configuration:** None
  - **Resource Group:** Create a new resource group with a unique name
  - **Pin to dashboard:** Not selected
5. In the Azure portal, view **Notifications** to verify that deployment has started. Then wait for the cluster to be deployed (this can take a long time – often 30 minutes or more. Now might be a good time to take a break!)

**Note:** As soon as an HDInsight cluster is running, the credit in your Azure subscription will start to be charged. The free-trial subscription includes a credit limit of approximately \$200 (or local equivalent) that you can spend over a period of 30 days, which is enough to complete the labs in this course as long as clusters are deleted when not in use. If you decide not to complete this lab, follow the instructions in the *Clean Up* procedure at the end of the lab to delete your cluster in order to avoid using your Azure credit unnecessarily.

## Creating a Linear Regression Model

In this exercise, you will create a linear regression machine learning model that predicts the fuel efficiency of a car based on its features.

**Tip:** You can copy the Scala or Python code for this exercise from **Linear Regression for Autos – Scala.txt** or **Linear Regression for Autos – Python.txt** in the lab files folder for this lab.

### Upload Source Data to Azure Storage

**Note:** You can use any Azure Storage client to upload the files to your Azure storage container. The instructions here assume you are using the Azure command line interface.

1. Open a new command line window.
2. Enter the following command to switch the Azure CLI to resource manager mode.

```
azure config mode arm
```

**Note:** If a *command not found* error is displayed, ensure that you have followed the instructions in the setup guide to install the Azure CLI.

3. Enter the following command to log into your Azure subscription:

```
azure login
```

4. Follow the instructions that are displayed to browse to the Azure device login site and enter the authentication code provided. Then sign into your Azure subscription using your Microsoft account.

5. Enter the following command to view your Azure resources:

```
azure resource list
```

6. Verify that your HDInsight cluster and the related storage account are both listed. Note that the information provided includes the resource group name as well as the individual resource names. Note the resource group and storage account name

7. Enter a `dir` (Windows) or `ls` (Mac OS X or Linux) command to view the contents of the **Lab03** folder where you extracted the lab files for this course (for example, `dir c:\DAT202.3x\Lab03` or `ls DAT202.3x/Lab03`), and verify that this folder contains a file named **autos.csv**. This file contains 358 observations of different cars, their attributes and the miles per gallon that they achieve.

8. Enter the following command on a single line to determine the connection string for your Azure storage account, specifying the storage account and resource group names you noted earlier:

```
azure storage account connectionstring show account -g resource_group
```

9. Note the connection string, copying it to the clipboard if your command line tool supports it.

10. If you are working on a Windows client computer, enter the following command to set a system variable for the connection string:

```
SET AZURE_STORAGE_CONNECTION_STRING=your_connection_string
```

If you are using a Linux or Mac OS X client computer, enter the following command to set a system variable for the connection string (enter the connection string in quotation marks):

```
export AZURE_STORAGE_CONNECTION_STRING="your_connection_string"
```

11. Enter the following command on a single line to upload the autos.csv file to the container used by your HDInsight cluster. Replace **autos\_path** with the local path to autos.csv (for example `c:\DAT202.3x\Lab03\autos.csv` or `DAT202.3x/Lab03/autos.csv`) and replace **container** with the name of the storage container used by your cluster (which should be the same as the cluster name):

```
azure storage blob upload autos_path container
```

12. Wait for the files to be uploaded.
13. Keep the command window open – you will use the Azure CLI again in a later exercise.

## Create a Notebook

1. In your web browser, in the Azure Portal, view the blade for your cluster. Then, under **Quick Links**, click **Cluster Dashboards**; and in the **Cluster Dashboards** blade, click **Jupyter Notebook**.
2. If you are prompted, enter the HTTP user name and password you specified for your cluster when provisioning it (not the SSH user name).
3. In the Jupyter web page that opens in a new tab, on the **Files** tab, view the existing folders.
4. If you have not previously created a folder named **Labs** in this cluster, in the **New** drop-down menu, click **Folder**. This create a folder named **Untitled Folder**. Then select the checkbox for the **Untitled Folder** folder, and above the list of folders, click **Rename**. Then rename the directory to **Labs**.
5. Open the **Labs** folder, and in the **New** drop-down list, click **Spark** if you prefer to program in Scala, or **PySpark** if you prefer to program in Python. This create a new notebook named **Untitled**.
6. At the top of the new notebook, click **Untitled**, and rename notebook to **Linear Regression for Autos**.

## Load the Autos Data into an RDD




1. In the empty cell at the top of the notebook, enter the following Scala or Python code to import the necessary namespaces:

### Python

```
from pyspark.mllib.regression import LabeledPoint
from pyspark.mllib.regression import LinearRegressionModel
from pyspark.mllib.regression import LinearRegressionWithSGD
from pyspark.mllib.linalg import Vectors
from pyspark.mllib.feature import StandardScaler
```

### Scala

```
import org.apache.spark.mllib.regression.LabeledPoint
import org.apache.spark.mllib.regression.LinearRegressionModel
import org.apache.spark.mllib.regression.LinearRegressionWithSGD
import org.apache.spark.mllib.linalg.Vectors
import org.apache.spark.mllib.feature.StandardScaler
```

2. With cursor still in the cell, on the toolbar click the **run cell, select below** button. It can take a short time to execute while the Spark context is created. As the code runs the  symbol at the top right of the page changes to a  symbol, and then returns to  when the code has finished running.
3. When the code has finished running, in the new cell under the output from the first cell, enter the following code to read the **autos.csv** text file from Azure storage into an RDD named **input** and then turn the comma-separated line into a **LabeledPoint** which comprises of a label and set of features. This will then be persisted across the Spark cluster:

### Python

```
def parsePoint(line):
    values = [float(x) for x in line.split(',')[0:7]]
    return LabeledPoint(values[0], values[1:])

input = sc.textFile("wasb:///autos.csv")

featurevector = input.map(lambda line : parsePoint(line)).persist()
```

### Scala

```
val input = sc.textFile("wasb:///autos.csv")

val featurevector = input.map { line =>
  val lineSplit = line.split(',')
  val featureArr = lineSplit.slice(1, 8).map(_.toDouble)
  val values = Vectors.dense(featureArr)
  LabeledPoint(lineSplit(0).toDouble, values)
}.persist()
```

4. With cursor still in the second cell, on the toolbar click the **run cell, select below** button. When the cell contents have finished executing the `[*]` will be replaced by a `[2]`

### Scale the Data

1. In the new cell under the second cell, enter the following code to scale the features so that they use a similar independent scale.

#### Python

```
labels = featurevector.map(lambda x: x.label)
features = featurevector.map(lambda x: x.features)
scaler = StandardScaler().fit(features)
scaledData = scaler.transform(features)
newfeatures = labels.zip(scaledData)
allset = newfeatures.map(lambda line: LabeledPoint(line[0], line[1]))
```

#### Scala

```
val scaler = new StandardScaler().fit(featurevector.map(x =>
x.features))
val scaledData = featurevector.map(x =>
LabeledPoint(x.label, scaler.transform(Vectors.dense(x.features.toArray)
)))
```

2. Click the **run cell, select below** button to run the third cell.

### Split the Data into Training and Test Data

1. In the fourth cell that appears at the bottom of the notebook, enter the following which will split the data into training and test data in the ratio of 70% to 30%:

#### Python

```
training, test = allset.randomSplit(weights=[0.7, 0.3], seed=1)
```

#### Scala

```
val allData = scaledData.randomSplit(Array(0.7, 0.3), seed = 11L)
val (training, test) = (allData(0), allData(1))
```

2. Click the **run cell, select below** button to run the fourth cell.

### Prepare and Train a Linear Regression Model

1. In the fifth cell enter the following code which will prepare the Linear Regression model:

#### Python

```
numIterations = 100
stepSize = 0.001
```

### Scala

```
val numIterations = 100
val stepSize = 0.001
val algorithm = new LinearRegressionWithSGD()
algorithm.setIntercept(true)
algorithm.optimizer.setNumIterations(numIterations)
algorithm.optimizer.setStepSize(stepSize)
```

2. Press **CTRL+ENTER** to run the fifth cell.
3. In the next cell that appear below enter the following code which will train the model with the training data we created earlier:

### Python

```
algorithm = LinearRegressionWithSGD.train(training,
iterations=numIterations, step=stepSize, intercept=True)
```

### Scala

```
val model = algorithm.run(training)
```

4. Click the **run cell, select below** button to run the contents of this cell.

## Test the Model

1. In the new cell enter the following code to create a set of predictions over the miles per gallon of each vehicle. The code creates a tuple of Labels against predicted Labels and then displays the first twenty of them.

### Python

```
valuesAndPreds = test.map(lambda p: (p.label,
algorithm.predict(p.features)))

valuesAndPreds.take(20)
```

### Scala

```
val valuesAndPreds = test.map { point =>
  val prediction = model.predict(point.features)
  (point.label, prediction)
}
valuesAndPreds.take(20)
```

2. Click the **run cell, select below** button to run the contents of this cell.
3. You should be able to see an output below the Cell of twenty Labels and predicted Labels which match fairly closely to each other.
4. In the new cell below add the following code to calculate the Mean Squared Error (MSE) of the model:

### Python

```
MSE = valuesAndPreds.map(lambda (v, p): (v - p)**2).reduce(lambda x, y:
x + y) / valuesAndPreds.count()

print("Mean Squared Error = " + str(MSE))
```

### Scala

```
val MSE = valuesAndPreds.map{case (v, p) => math.pow((v - p), 2)}.mean()
println("training Mean Squared Error = " + MSE)
```

5. Click the **run cell, select below** button to run the contents of this cell and view the MSE.

6. On the toolbar, click the **Save** button.
7. On the **File** menu click **Close and Halt**. If prompted, confirm that you want to close the tab.
8. Verify that the **Linear Regression for Autos** notebook is listed in the **Labs** folder.

## Creating a Pipeline

In this exercise, you will implement a pipeline to create and cross validate a model that performs sentiment analysis on tweets.

**Tip:** You can copy the Scala or Python code for this exercise from **Sentiment Analysis of Tweets – Scala.txt** or **Sentiment Analysis of Tweets – Python.txt** in the lab files folder for this lab.

## Upload the Source Data

1. Return to the command prompt window in which you previously used the Azure CLI to upload data to Azure storage.
2. Enter a `dir` (Windows) or `ls` (Mac OS X or Linux) command to view the contents of the **Lab03** folder where you extracted the lab files for this course, and verify that this folder contains files named **training-tweets.csv** and **test-tweets.csv**. These files contain tweets on which you will perform sentiment analysis.
3. Enter the following command on a single line to upload the training-tweets.csv file to the blob store container for your cluster. Replace **tweets\_path** with the local path to training-tweets.csv and replace **container** with the name of the storage container used by your cluster (which should be the same as the cluster name):

```
azure storage blob upload tweets_path container
```

4. Enter a second `azure storage blob upload` command to upload test-tweets.csv to the blob store container for your cluster.

```
azure storage blob upload tweets_path container
```

14. Keep the command window open – you will use the Azure CLI again in a later exercise.

## Create a New Notebook

1. Switch to the Jupyter tab in your web browser, and in the **Labs** folder, in the **New** drop-down list, click **Spark** if you prefer to program in Scala, or **PySpark** if you prefer to program in Python. This creates a new notebook named **Untitled**.
2. At the top of the new notebook, click **Untitled**, and rename notebook to **Sentiment Analysis of Tweets**.

## Implement a Pipeline that Creates and Cross-Validates a Model

1. In the empty cell at the top of the notebook, enter the following code. This code imports the necessary namespaces:

### Python

```
import org.apache.spark.ml.Pipeline
import org.apache.spark.ml.classification.LogisticRegression
import org.apache.spark.ml.feature.{HashingTF, Tokenizer}
import org.apache.spark.mllib.linalg.Vector
import org.apache.spark.sql.Row
case class Tweets(id: Int, label: Double, source: String, text: String)
import sqlContext.implicits._
```

## Scala

```
import org.apache.spark.ml.Pipeline
import org.apache.spark.ml.classification.LogisticRegression
import org.apache.spark.ml.evaluation.BinaryClassificationEvaluator
import org.apache.spark.ml.feature.{HashingTF, Tokenizer}
import org.apache.spark.ml.tuning.{ParamGridBuilder, CrossValidator}
import org.apache.spark.mllib.linalg.Vector
import org.apache.spark.sql.Row
import org.apache.spark.sql.types._;
import org.apache.spark.sql.Row;
```

2. With cursor still in the cell, on the toolbar click the **run cell, select below** button.
3. When the code has finished running, and the Spark context has been created, in the new cell under the output from the first cell, enter the following code to read the **training-tweets.csv** and **test-tweets.csv** text file from Azure storage into RDD's labeled training and test. These are then converted into Spark DataFrames so that we can use them with ML Pipelines:

## Python

```
trainingrdd = sc.textFile("wasb:///training-
tweets.csv").zipWithIndex().filter(lambda line: line[1] > 0).map(lambda
line: line[0].split(","))
testrdd = sc.textFile("wasb:///test-
tweets.csv").zipWithIndex().filter(lambda line: line[1] > 0).map(lambda
line: line[0].split(","))

fields = [StructField("id", StringType(), True), StructField("label",
StringType(), True), StructField("source", StringType(), True),
StructField("text", StringType(), True)]
schema = StructType(fields)

training = sqlContext.createDataFrame(trainingrdd, schema)
test = sqlContext.createDataFrame(testrdd, schema)
```

## Scala

```
val training = sc.textFile("/training-
tweets.csv").zipWithIndex().filter(_._2 > 0).map(line =>
line._1.split(",")).map(tw => Tweets(tw(0).toInt, tw(1).toDouble,
tw(2), tw(3))).toDF()

val test = sc.textFile("/test-tweets.csv").zipWithIndex().filter(_._2 >
0).map(line => line._1.split(",")).map(tw => Tweets(tw(0).toInt,
tw(1).toDouble, tw(2), tw(3))).toDF()
```

4. With cursor still in the second cell, on the toolbar click the **run cell, select below** button and wait for the code to finish running.
5. In the new empty cell, enter the following code to create an ML Pipeline. The code takes the tweet text from the input column "text" and outputs as a token to a new column called "words" using this as input to a "Hashing Term Frequency" (**HashingTF**) evaluator which outputs the hashed terms as a set of features. Logistic Regression will then be applied to calculate the sentiment of the tweet(s). All of these estimators are then set as stages in a single **Pipeline**:

## Python

```
training = training.withColumn("label",
training.label.cast(DoubleType()))
test = training.withColumn("label", training.label.cast(DoubleType()))
```



```
tokenizer = Tokenizer(inputCol="text", outputCol="words")
hashingTF = HashingTF(inputCol="words", outputCol="features")
lr = LogisticRegression(maxIter=10)
pipeline = Pipeline(stages=[tokenizer, hashingTF, lr])
```

#### Scala

```
val tokenizer = new
Tokenizer().setInputCol("text").setOutputCol("words")
val hashingTF = new
HashingTF().setInputCol(tokenizer.getOutputCol).setOutputCol("features"
)
val lr = new LogisticRegression().setMaxIter(10)
val pipeline = new Pipeline().setStages(Array(tokenizer, hashingTF,
lr))
```

6. Click the **run cell, select below** button to run the third cell.
7. In the next cell that appears below, enter the following code which will train the model with the training data we created earlier:

#### Python

```
model = pipeline.fit(training)
```

#### Scala

```
val model = pipeline.fit(training)
```

8. Click the **run cell, select below** button to run the contents of this cell. This can take a while to run.
9. In the new cell enter the following code to transform the test data and retrieve the probability and predictions of each label:

#### Python

```
prediction = model.transform(test)
selected = prediction.select("id", "label", "text", "probability",
"prediction")
selected.show(500)
```

#### Scala

```
val prediction = model.transform(test)
val selected = prediction.select("id", "label", "text", "probability",
"prediction")
selected.show(500)
```

10. Click the **run cell, select below** button to run the contents of this cell, and then review the first 500 predictions made by the model.
11. On the toolbar, click the **Save** button.
12. On the **File** menu click **Close and Halt**. If prompted, confirm that you want to close the tab.

## Creating a Streaming Linear Regression Model

In this exercise, you will create a streaming linear regression machine learning model that predicts whether a person earns less than or greater than \$50K.

**Tip:** You can copy and paste the Scala or Python code for this exercise from **Income Predictor – Scala.txt** or **Income Predictor – Python.txt** in the lab files folder for this lab.

## Start a Streaming Program

1. Use an SSH client (Putty on Windows or the native client on a Mac or Linux) to connect to the cluster. The endpoint to SSH into your cluster will be:  
***cluster\_name-ssh.azurehdinsight.net.***
2. Log in using the SSH User credentials you specified when creating the cluster (**not** the HTTP user credentials)
3. If you intend to use Scala at the bash prompt enter **spark-shell**. Alternatively, if you intend to use python enter **pyspark** at the prompt.
4. After the shell has been initialized, enter the following code to import the required namespaces:

### Python

```
from pyspark.mllib.linalg import Vectors
from pyspark.mllib.regression import LabeledPoint
from pyspark.mllib.classification import
StreamingLogisticRegressionWithSGD
from pyspark.streaming import *
```

### Scala

```
import org.apache.spark.mllib.linalg.Vectors
import org.apache.spark.mllib.regression.LabeledPoint
import
org.apache.spark.mllib.classification.StreamingLogisticRegressionWithSGD
import org.apache.spark.streaming.{Seconds, StreamingContext}
```

5. When the code has finished running, enter the following code to process files in the training and test directory. If you are using Python, it's important to make sure that the indentation is exactly as shown below:

### Scala

```
val ssc = new StreamingContext(sc, Seconds(30))
val training =
ssc.textFileStream("/training/lr").map(LabeledPoint.parse)
val testing = ssc.textFileStream("/testing/lr").map(LabeledPoint.parse)
```

### Python

```
def parse(lp):
    label = float(lp[lp.find('(') + 1: lp.find(',')])
    vec = Vectors.dense(lp[lp.find '[' + 1: lp.find(')')].split(','))
    return LabeledPoint(label, vec)

ssc = StreamingContext(sc, 30)
training = ssc.textFileStream("/training/lr").map(parse)
testing = ssc.textFileStream("/testing/lr").map(parse)
```

6. When the code has finished executing, enter the following code into the shell to train the Logistic Regression model on six features and continually look in the training and testing directories to ensure that new files are picked up:

### Python

```
numFeatures = 6
model = StreamingLogisticRegressionWithSGD()
```

```
model.setInitialWeights([0.0, 0.0, 0.0, 0.0, 0.0, 0.0])
```

### Scala

```
val numFeatures = 6
val model = new
StreamingLogisticRegressionWithSGD().setInitialWeights(Vectors.zeros(numFeatures))
```

7. When the code has finished executing, enter the following code to configure the trained model to be updated by the new model files in the training directory and the test files to be used and picked up from the test directory, and then store the results in an output directory:

### Python

```
model.trainOn(training)
model.predictOnValues(testing.map(lambda lp : (lp.label,
lp.features))).saveAsTextFiles("/testing/out")
```

### Scala

```
model.trainOn(training)
model.predictOnValues(testing.map(lp => (lp.label,
lp.features))).saveAsTextFiles("/testing/out")
```

8. Lastly enter the following code at the prompt to start the streaming receiver and block to update the console for messages:

### Python

```
ssc.start()
ssc.awaitTermination()
```

### Scala

```
ssc.start()
ssc.awaitTermination()
```

9. In the output you should see an error suggesting that **/training/lr** doesn't exist. You are now going to create the training and test folder.

## Submit New Training and Testing Data

**Note:** You can use any Azure Storage client to upload the files to your Azure storage container. The instructions here assume you are using the Azure command line interface.

1. Return to the command line window in which you used the Azure CLI to upload files to Azure previously.
2. Enter a `dir` (Windows) or `ls` (Mac OS X or Linux) command to view the contents of the **Lab03** folder where you extracted the lab files for this course (for example, `dir c:\DAT202.3x\Lab03` or `ls DAT202.3x/Lab03`), and verify that this folder contains a file named **income-training1.txt**. This file contains LabeledPoints of US income data, including features. Also verify that the folder contains a file named **income-test.txt**.
3. Enter the following commands on separate lines to upload the requisite training and test files to blob storage in the container used by your HDInsight cluster. Replace **file** with the local path to `income-training1.txt` and `income-test.txt` (for example `c:\DAT202.3x\Lab03\income-test.txt` or `DAT202.3x/Lab03/income-test.txt`) and replace **container** with the name of the storage container used by your cluster (which should be the same as the cluster name):

```
azure storage blob upload file container training/lr/income-  
training1.txt
```

```
azure storage blob upload file container testing/lr/income-test.txt
```

4. Wait for the files to be uploaded, and then observe the SSH console to note any increased activity in the Spark program.
5. In the command prompt window where you have been using the Azure CLI, enter the following to check for output in the **testing/out** folder:

```
azure storage blob list container -p testing/out
```

6. Identify any of the blobs that have the file name **part-00000** and a length greater than 0. If there are none, wait a few minutes and resubmit the command.
7. Enter the following command to download any of the files named **part-00000** with a length greater than 0 to a local folder, replacing **container** with your container name, **testing-out-123/part-00000** with the blob name, and **file** with the path to the local file you want to create):

```
azure storage blob download container testing/out-123/part-00000 file
```

8. Open the downloaded file in a text editor and verify that it contains output in the form of **(label, prediction)**.
9. Close the SSH window to end the session.

## Clean Up

Now that you have finished the lab, you should delete your cluster and the associated storage account. This ensures that you avoid being charged for cluster resources when you are not using them. If you are using a trial Azure subscription that includes a limited free credit value, deleting the cluster maximizes your credit and helps to prevent using it all before the free trial period has ended.

### Delete the Resource Group for your HDInsight Cluster

1. If it is not already open in a tab in your web browser, browse to the new Azure portal at <https://portal.azure.com>.
2. In the Azure portal, view your **Resource groups** and select the resource group you created for your cluster. This resource group contains your cluster and the associated storage account.
3. In the blade for your resource group, click **Delete**. When prompted to confirm the deletion, enter the resource group name and click **Delete**.
4. Wait for a notification that your resource group has been deleted.