# Big Data Analysis: Assignment 1 Part Two and Three: Spark Scala and PySpark

Baochen Hu
Ning Zhang

June 20, 2015

**Zusammenfassung**

This part of the assignment aims to practice algorithms in Scala and PySpark with Spark. We used the white vinho verde wine samples data, from the north of Portugal. We separated the dataset into two parts. The first part is our training dataset. The second part is our test dataset for the results.

# 1  Development Environment

Spark version 1.3.1
sbt

# 2  Running Introduction

Use sbt build jar package
Command:

```
sbt package
```

Use spark-submit in /spark/bin run the jar file we built in your scala project Example:

```
./spark-submit ../../../../../Users/ningzhang/Desktop/Assignment1/spark_scala_ml/
    LassoWithSGD/target/scala-2.10/simple-project_2.10-1.0.jar
```

# 3  Tryout

We will try out the algorithms in Scala and PySpark using Spark MLlib with the following models:

Classification:

```
LogisticRegressionWithLBFGS
LogisticRegressionWithSGD
```

Regression:

```
LinearRegressionWithSGD
RidgeRegressionWithSGD
?LassoWithSGD
```

## 3.1 LogisticRegressionWithLBFGS

### 3.1.1 Introduction

Train a classification model for Multinomial/Binary Logistic Regression using Limited-memory BFGS. Standard feature scaling and L2 regularization are used by default. NOTE: Labels used in Logistic Regression should be 0, 1, ..., k - 1 for k classes multi-label classification problem.

### 3.1.2 Source Code in Scala:

```scala
/* SimpleApp.scala */
import org.apache.spark.SparkContext
import org.apache.spark.SparkContext._
import org.apache.spark.SparkConf
import org.apache.spark.mllib.evaluation.MulticlassMetrics
import org.apache.spark.mllib.linalg.Vectors
import org.apache.spark.mllib.regression.LabeledPoint
import org.apache.spark.mllib.classification.{LogisticRegressionWithLBFGS,
      LogisticRegressionModel}
object SimpleApp {
  def main(args: Array[String]) {
    val csvPath = "/Users/baochenhu/Desktop/spark-1.3.1-bin-hadoop2.4/winequality-white.csv"
    // Should be some file on your system
    val conf = new SparkConf().setAppName("Simple Application")
    val sc = new SparkContext(conf)
    val csv = sc.textFile(csvPath)
    val headerAndRows= csv.map(line=>line.split(";").map(_.trim))
    val header = headerAndRows.first
    val data = headerAndRows.filter(_(0) != header(0))
    val parsedData = data.map{line => LabeledPoint(line(11).toDouble, Vectors.dense(line(0).
        toDouble,line(1).toDouble))}
    parsedData.cache()
    val splits = parsedData.randomSplit(Array(0.6, 0.4), seed = 11L)
    val training = splits(0).cache()
    val test = splits(1)
    val numIterations = 100
    val model = new LogisticRegressionWithLBFGS().setNumClasses(10).run(training)
    val predictionAndLabels = test.map { case LabeledPoint(label, features) =>
     val prediction = model.predict(features)
            (prediction, label)}


    // Get evaluation metrics.
val metrics = new MulticlassMetrics(predictionAndLabels)
val precision = metrics.precision
println("Precision = " + precision)
  }
}
\subsubsection{Output:}
\begin{lstlisting}
Precision = 0.4532994923857868
```

### 3.1.3 Source Code in PySpark:

```python
from pyspark.mllib.regression import LabeledPoint, LinearRegressionWithSGD
from numpy import array
from pyspark import SparkContext
from pyspark.mllib.classification import LogisticRegressionWithLBFGS
# Load and parse the data
def parsePoint(line):
    values = [float(x) for x in line.split(';')]
    return LabeledPoint(values[11], values[0:10])

```

```
10
11  sc = SparkContext("local", "Simple App")
12  data = sc.textFile("../winequality.csv")
13  parsedData = data.map(parsePoint)
14
15  # Build the model
16  model = LogisticRegressionWithLBFGS.train(parsedData)
17
18  # Evaluating the model on training data
19  labelsAndPreds = parsedData.map(lambda p: (p.label, model.predict(p.features)))
20  trainErr = labelsAndPreds.filter(lambda (v, p): v != p).count() / float(parsedData.count())
21  print("Training Error = " + str(trainErr))
```

## 3.2 LogisticRegressionWithSGD

### 3.2.1 Introduction

Train a classification model for Binary Logistic Regression using Stochastic Gradient Descent. By default
L2 regularization is used, which can be changed via LogisticRegressionWithSGD.optimizer. NOTE: Labels
used in Logistic Regression should be 0, 1, ..., k - 1 for k classes multi-label classification problem. Using
LogisticRegressionWithLBFGS is recommended over this.

### 3.2.2 Source Code:

```
1   /* SimpleApp.scala */
2   import org.apache.spark.SparkContext
3   import org.apache.spark.SparkContext._
4   import org.apache.spark.SparkConf
5   import org.apache.spark.mllib.linalg.Vectors
6   import org.apache.spark.mllib.regression.LabeledPoint
7   import org.apache.spark.mllib.classification.LogisticRegressionWithSGD
8   object SimpleApp {
9     def main(args: Array[String]) {
10       val csvPath = "/Users/ningzhang/Desktop/Assignment1/winequality-white.csv" // Should be
              some file on your system
11       val conf = new SparkConf().setAppName("Simple Application")
12       val sc = new SparkContext(conf)
13       val csv = sc.textFile(csvPath)
14       val headerAndRows= csv.map(line=>line.split(";").map(_.trim))
15       val header = headerAndRows.first
16       val data = headerAndRows.filter(_(0) != header(0))
17       val parsedData = data.map{line => LabeledPoint(line(11).toDouble, Vectors.dense(line(0).
              toDouble,line(1).toDouble))}
18     parsedData.cache()
19     // Run training algorithm to build the model
20         val numIterations = 20
21         val model = LogisticRegressionWithSGD.train(parsedData, numIterations)
22
23         // Evaluate model on training examples and compute training error
24         val labelAndPreds = parsedData.map { point =>
25           val prediction = model.predict(point.features)
26           (point.label, prediction)
27         }
28         val trainErr = labelAndPreds.filter(r => r._1 != r._2).count.toDouble / parsedData.
              count
29         val n1 = labelAndPreds.filter(r => (r._1==1)&&(r._2==1)).count.toDouble
30         val n2 = labelAndPreds.filter(r => (r._1==0)&&(r._2==0)).count.toDouble
31         val d1 = labelAndPreds.filter(r => (r._1==1)).count.toDouble
32         val d2 = labelAndPreds.filter(r => (r._1==0)).count.toDouble
33         val sensitivity = n1/d1
34         val specificity = n2/d2
35
36         println("\nTraining Error = " + trainErr)
```

```
37          println("Sensitivity = " + sensitivity)
38          println("Specificity = " + specificity)
39          println();
40      }
41 }
```

### 3.2.3 Output:

```
1 Training Error = 0.461624026696
2 Sensitivity = 0.48275862069
3 Specificity = 0.47620692758
```

## 3.3 LinearRegressionWithSGD

### 3.3.1 Introduction

Train a linear regression model with no regularization using Stochastic Gradient Descent. This solves the least squares regression formulation f(weights) = 1/n ||A weights-y||2 (which is the mean squared error). Here the data matrix has n rows, and the input RDD holds the set of rows of A, each with its corresponding right hand side label y. See also the documentation for the precise formulation.

### 3.3.2 Source Code in Scala:

```
1 /* SimpleApp.scala */
2 import org.apache.spark.SparkContext
3 import org.apache.spark.SparkContext._
4 import org.apache.spark.SparkConf
5 import org.apache.spark.mllib.linalg.Vectors
6 import org.apache.spark.mllib.regression.LabeledPoint
7 import org.apache.spark.mllib.regression.LinearRegressionWithSGD
8 import org.apache.spark.mllib.regression.RidgeRegressionWithSGD
9 object SimpleApp {
10   def main(args: Array[String]) {
11     val csvPath = "/Users/ningzhang/Assignment1/winequality-white.csv" // Should be some
           file on your system
12     val conf = new SparkConf().setAppName("Simple Application")
13     val sc = new SparkContext(conf)
14     val csv = sc.textFile(csvPath)
15     val headerAndRows= csv.map(line=>line.split(";").map(_.trim))
16     val header = headerAndRows.first
17     val data = headerAndRows.filter(_(0) != header(0))
18     val parsedData = data.map{line => LabeledPoint(line(11).toDouble, Vectors.dense(line(0).
           toDouble, line(1).toDouble))}
19    parsedData.cache()
20     val numIterations = 20
21     val model = LinearRegressionWithSGD.train(parsedData, numIterations)
22     val valuesAndPreds = parsedData.map { point =>
23          val prediction = model.predict(point.features)
24           (point.label, prediction)}
25      val MSE = valuesAndPreds.map{ case(v, p) => math.pow((v - p), 2)}.reduce(_ + _)/
           valuesAndPreds.count
26     println("LinearRegressionWithSGD training Mean Squared Error = " + MSE)
27      println("----------------------------------------")
28     val model1 = RidgeRegressionWithSGD.train(parsedData, numIterations)
29     val valuesAndPreds1 = parsedData.map { point =>
30          val prediction = model1.predict(point.features)
31           (point.label, prediction)}
32      val MSE1 = valuesAndPreds1.map{ case(v, p) => math.pow((v - p), 2)}.reduce(_ + _)/
           valuesAndPreds1.count
```

```
33        println("LinearRegressionWithSGD training Mean Squared Error = " + MSE1)
34        println("————————————————————————————")
35    }
36 }
```

### 3.3.3  Output:

```
1 LinearRegressionWithSGD training Mean Squared Error = 1.4220919486200828E49
```

### 3.3.4  Source Code in PySpark:

```python
1 from pyspark.mllib.regression import LabeledPoint, LinearRegressionWithSGD
2 from numpy import array
3 from pyspark import SparkContext
4 # Load and parse the data
5 def parsePoint(line):
6     values = [float(x) for x in line.split(';')]
7     return LabeledPoint(values[11], values[0:10])
8
9
10 sc = SparkContext("local", "Simple App")
11 data = sc.textFile("../winequality.csv")
12 parsedData = data.map(parsePoint)
13
14 # Build the model
15 model = LinearRegressionWithSGD.train(parsedData)
16
17 # Evaluate the model on training data
18 valuesAndPreds = parsedData.map(lambda p: (p.label, model.predict(p.features)))
19 MSE = valuesAndPreds.map(lambda (v, p): (v - p)**2).reduce(lambda x, y: x + y) /
      valuesAndPreds.count()
20 print("Mean Squared Error = " + str(MSE))
```

## 3.4  RidgeRegressionWithSGD

### 3.4.1  Introduction

Train a regression model with L2-regularization using Stochastic Gradient Descent. This solves the l1-regularized least squares regression formulation $f(weights) = 1/2n \ ||A \ weights-y||2 + regParam/2 \ ||weights||2$
Here the data matrix has n rows, and the input RDD holds the set of rows of A, each with its corresponding right hand side label y. See also the documentation for the precise formulation.

### 3.4.2  Source Code in Scala:

```scala
1 /* SimpleApp.scala */
2 import org.apache.spark.SparkContext
3 import org.apache.spark.SparkContext._
4 import org.apache.spark.SparkConf
5 import org.apache.spark.mllib.linalg.Vectors
6 import org.apache.spark.mllib.regression.LabeledPoint
7 import org.apache.spark.mllib.regression.LinearRegressionWithSGD
8 import org.apache.spark.mllib.regression.RidgeRegressionWithSGD
9 object SimpleApp {
10    def main(args: Array[String]) {
```

```scala
11    val csvPath = "/Users/baochenhu/Desktop/spark-1.3.1-bin-hadoop2.4/winequality-white.csv"
         // Should be some file on your system
12    val conf = new SparkConf().setAppName("Simple Application")
13    val sc = new SparkContext(conf)
14    val csv = sc.textFile(csvPath)
15    val headerAndRows= csv.map(line=>line.split(";").map(_.trim))
16    val header = headerAndRows.first
17    val data = headerAndRows.filter(_(0) != header(0))
18    val parsedData = data.map{line => LabeledPoint(line(11).toDouble, Vectors.dense(line(0).
         toDouble,line(1).toDouble))}
19   parsedData.cache()
20    val numIterations = 20
21    val model1 = RidgeRegressionWithSGD.train(parsedData, numIterations)
22    val valuesAndPreds1 = parsedData.map { point =>
23        val prediction = model1.predict(point.features)
24         (point.label, prediction)}
25     val MSE1 = valuesAndPreds1.map{ case(v, p) => math.pow((v - p), 2)}.reduce(_ + _)/
         valuesAndPreds1.count
26     println("RidgeRegressionWithSGD training Mean Squared Error = " + MSE1)
27     println("————————————————————————————")
28   }
29 }
```

### 3.4.3  Output:

```
1 RidgeRegressionWithSGD training Mean Squared Error = 1.4220919486200828E49
```

### 3.4.4  Source Code in PySpark:

```python
1 from pyspark.mllib.regression import LabeledPoint, RidgeRegressionWithSGD
2 from numpy import array
3 from pyspark import SparkContext
4 from pyspark.mllib.classification import LogisticRegressionWithLBFGS
5 # Load and parse the data
6 def parsePoint(line):
7     values = [float(x) for x in line.split(';')]
8     return LabeledPoint(values[11], values[0:10])
9
10
11 sc = SparkContext("local", "Simple App")
12 data = sc.textFile("../winequality.csv")
13 parsedData = data.map(parsePoint)
14
15 # Build the model
16 model = RidgeRegressionWithSGD.train(parsedData)
17
18 # Evaluating the model on training data
19 labelsAndPreds = parsedData.map(lambda p: (p.label, model.predict(p.features)))
20 trainErr = labelsAndPreds.filter(lambda (v, p): v != p).count() / float(parsedData.count())
21 print("Training Error = " + str(trainErr))
```

## 3.5  LassonWithSGD

### 3.5.1  Introduction

Train a regression model with L1-regularization using Stochastic Gradient Descent. This solves the 1-regularized least squares regression formulation f(weights) = 1/2n A weights-y 2 + regParam weights - 1

Here the data matrix has n rows, and the input RDD holds the set of rows of A, each with its corresponding right hand side label y. See also the documentation for the precise formulation.

### 3.5.2   Source Code in Scala:

```scala
/* SimpleApp.scala */
import org.apache.spark.SparkContext
import org.apache.spark.SparkContext._
import org.apache.spark.SparkConf
import org.apache.spark.mllib.linalg.Vectors
import org.apache.spark.mllib.regression.LabeledPoint
import org.apache.spark.mllib.regression.LinearRegressionWithSGD
import org.apache.spark.mllib.regression.LassoWithSGD
object SimpleApp {
  def main(args: Array[String]) {
    val csvPath = "/Users/baochenhu/Desktop/spark-1.3.1-bin-hadoop2.4/winequality-white.csv"
       // Should be some file on your system
    val conf = new SparkConf().setAppName("Simple Application")
    val sc = new SparkContext(conf)
    val csv = sc.textFile(csvPath)
    val headerAndRows= csv.map(line=>line.split(";").map(_.trim))
    val header = headerAndRows.first
    val data = headerAndRows.filter(_(0) != header(0))
    val parsedData = data.map{line => LabeledPoint(line(11).toDouble, Vectors.dense(line(0).
        toDouble,line(1).toDouble))}
    parsedData.cache()
    val numIterations = 20
    val model = LassoWithSGD.train(parsedData, numIterations)
    val valuesAndPreds = parsedData.map { point =>
         val prediction = model.predict(point.features)
          (point.label, prediction)}
     val MSE = valuesAndPreds.map{ case(v, p) => math.pow((v - p), 2)}.reduce(_ + _)/
        valuesAndPreds.count
     println("LassoWithSGD training Mean Squared Error = " + MSE)
     println("————————————————————————————")


  }
}
```

### 3.5.3   Output:

```
LassoWithSGD training Mean Squared Error = 1.409246080888596E49
```

### 3.5.4   Source Code in PySpark:

```python
from pyspark.mllib.regression import LabeledPoint, LassoWithSGD
from numpy import array
from pyspark import SparkContext
from pyspark.mllib.classification import LogisticRegressionWithLBFGS
# Load and parse the data
def parsePoint(line):
    values = [float(x) for x in line.split(';')]
    return LabeledPoint(values[11], values[0:10])


sc = SparkContext("local", "Simple App")
data = sc.textFile("../winequality.csv")
parsedData = data.map(parsePoint)

```

```
15  # Build the model
16  model = LassoWithSGD.train(parsedData)
17
18  # Evaluating the model on training data
19  labelsAndPreds = parsedData.map(lambda p: (p.label, model.predict(p.features)))
20  trainErr = labelsAndPreds.filter(lambda (v, p): v != p).count() / float(parsedData.count())
21  print("Training Error = " + str(trainErr))
```