

CS5542 Big Data Apps and Analytics

LAB ASSIGNMENT #4

The image classification program is run using decision tree algorithm instead of random forest algorithm. The data set used is the data set provided in the source code, since the other data set is not suitable/ giving no proper output. The decision tree algorithm is a classification algorithm, starting from top to bottom, the train examples are always at the root of the tree. In the decision tree algorithm, the partitioning is done using selected attributes. Here, for this lab tutorial, the data set is the default data set provided, the categories are:

bibimap
jiaozi
oden
omelet
rice
sausage
spaghetti
sushi
tempura
toast

these are divided into 70% train and 30% test. The image classification program is run on this data set and the resulting output is the accuracy obtained by running the program: 0.22

The program is also run on other data set(DataLab4) with 5 categories:

airplane
apples
cats
eggs
wolves

the correct/desired output was not obtained, there was a problem with the images/ image format, so the original data set was used to get the output.

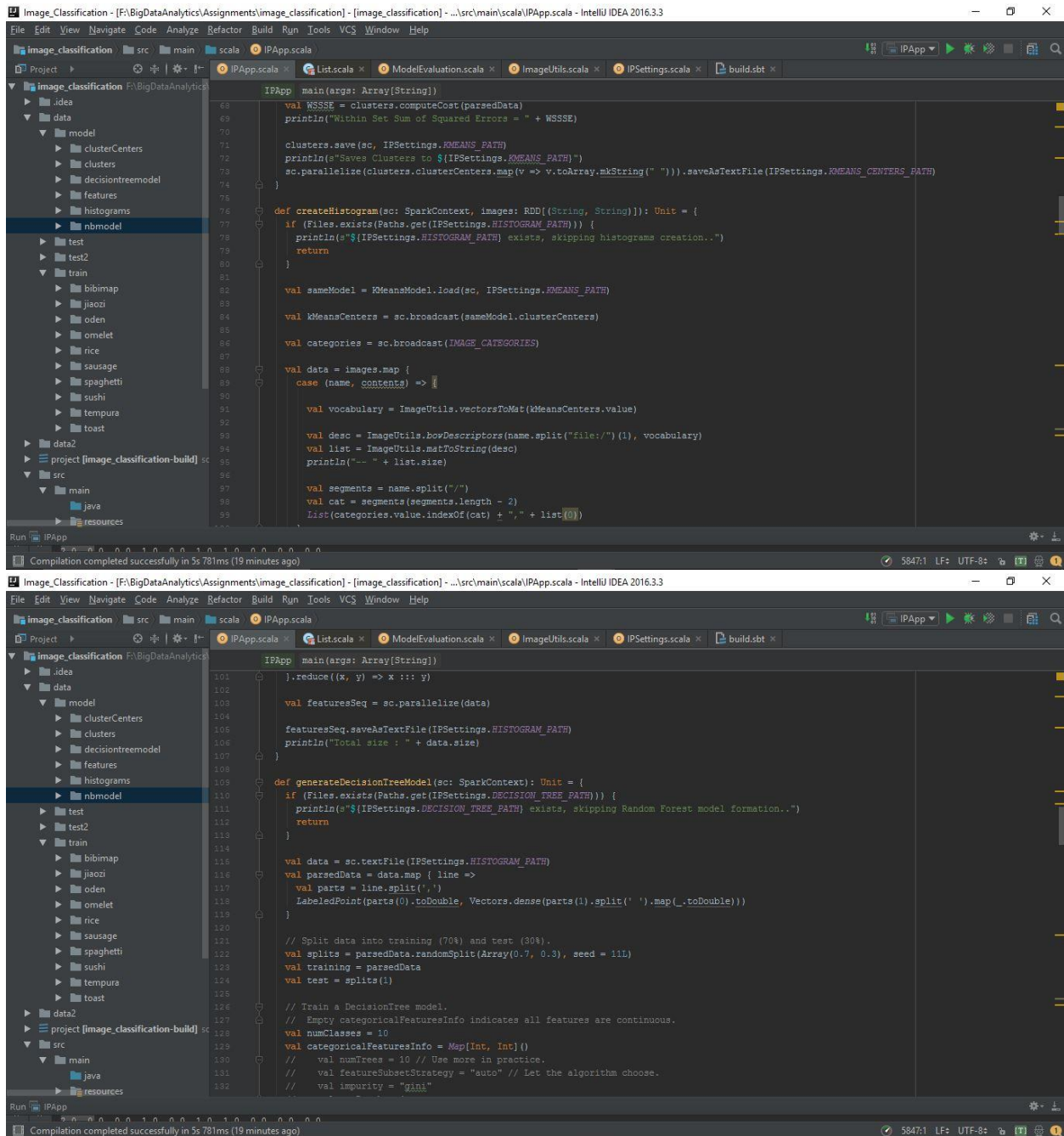
These are few screen shots from the program:

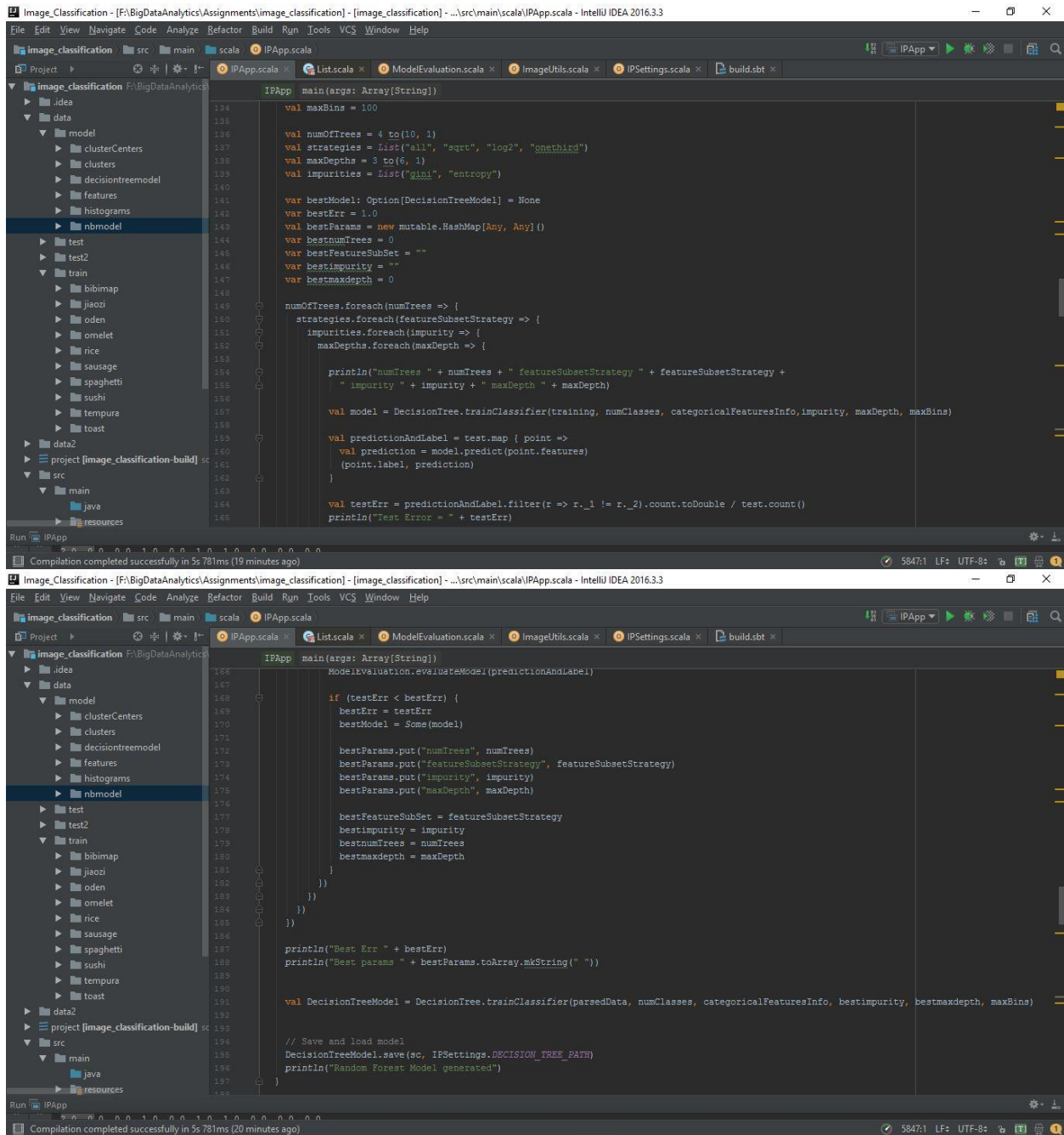
```
Image_Classification - [F:\BigDataAnalytics\Assignments\image_classification] - [image_classification] - \src\main\scala\IPApp.scala - IntelliJ IDEA 2016.3.3
File Edit View Navigate Code Analyze Refactor Build Run Tools VCS Window Help
image_classification src main scala IPApp.scala
Project image_classification
  .idea
  data
    model
      clusterCenters
      clusters
      decisiontreemodel
      features
      histograms
      nbmodel
    test
    test2
    train
      bibimap
      jiaozi
      oden
      omelet
      rice
      sausage
      spaghetti
      sushi
      tempura
      toast
    data2
  project [image_classification-build]
  src
    main
      java
      resources

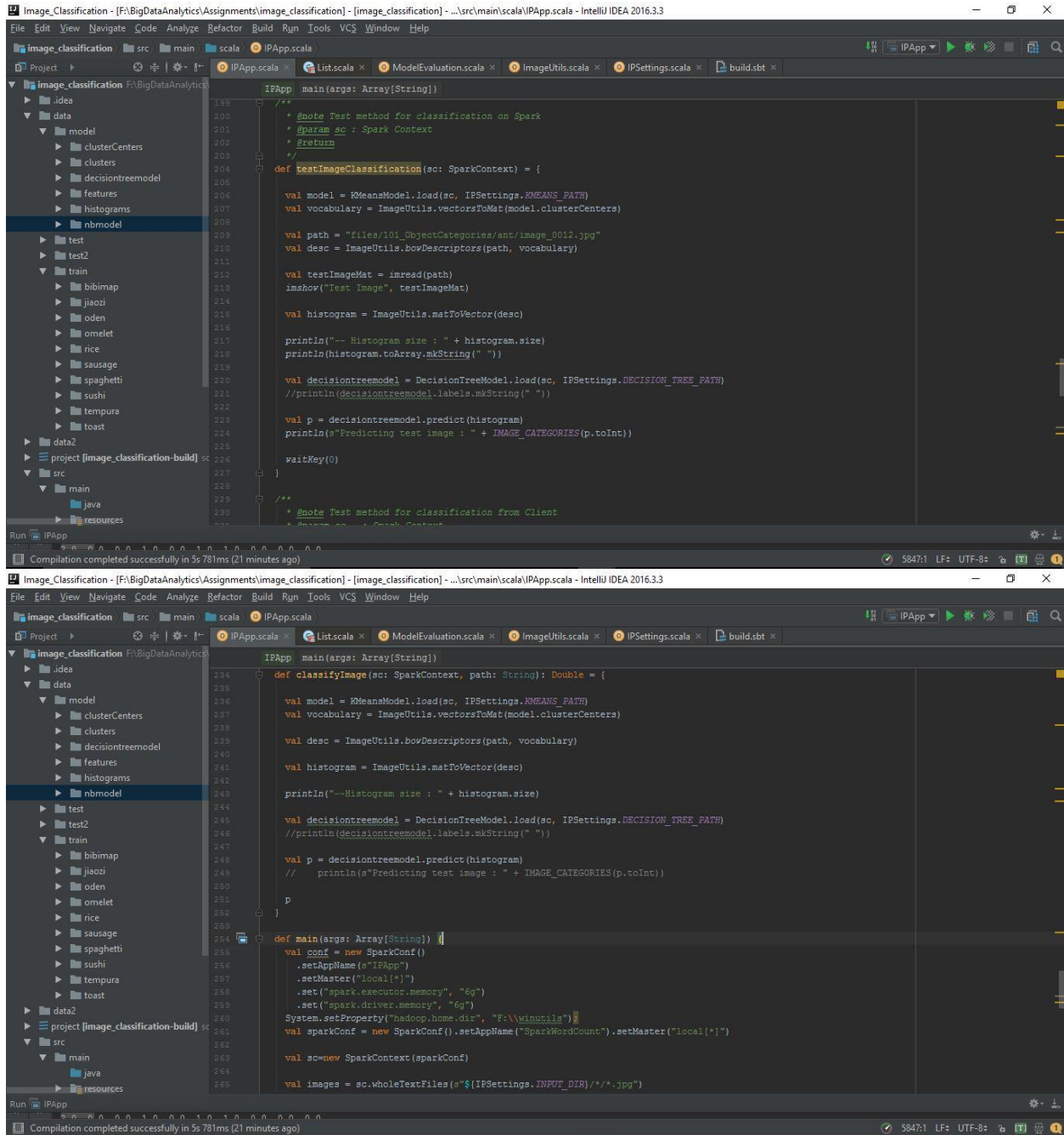
IPApp main(args: Array[String])
1  /**...*/
2
3  import java.nio.file.{Files, Paths}
4
5  import org.apache.spark.mllib.clustering.{KMeans, KMeansModel}
6  import org.apache.spark.mllib.linalg.Vectors
7  import org.apache.spark.mllib.regression.LabeledPoint
8  import org.apache.spark.mllib.tree.DecisionTree
9  import org.apache.spark.mllib.tree.model.DecisionTreeModel
10 import org.apache.spark.rdd.RDD
11 import org.apache.spark.streaming.{Seconds, StreamingContext}
12 import org.apache.spark.{SparkConf, SparkContext}
13 import org.bytedeco.javacpp.opencv_highgui._
14
15 import scala.collection.mutable
16
17 object IPApp {
18   val featureVectorsCluster = new mutable.MutableList[String]
19
20   val IMAGE_CATEGORIES = List("rice", "tempura", "toast", "bibimap", "sushi", "spaghetti", "sausage", "oden", "omelet", "jiaozi")
21   //val IMAGE_CATEGORIES = List("accordion", "airplanes", "anchor", "ant", "barrel", "bass", "beaver", "binoculars", "bonsai")
22
23   /**
24    * @param sc : SparkContext
25    * @param images : Images list from the training set
26    */
27   def extractDescriptors(sc: SparkContext, images: RDD[(String, String)]): Unit = {
28     if (Files.exists(Paths.get(IPSettings.FEATURES_PATH))) {
29       println(s"${IPSettings.FEATURES_PATH} exists, skipping feature extraction..")
30       return
31     }
32
33     Run IPApp
34     Compilation completed successfully in 5s 781ms (17 minutes ago)
```

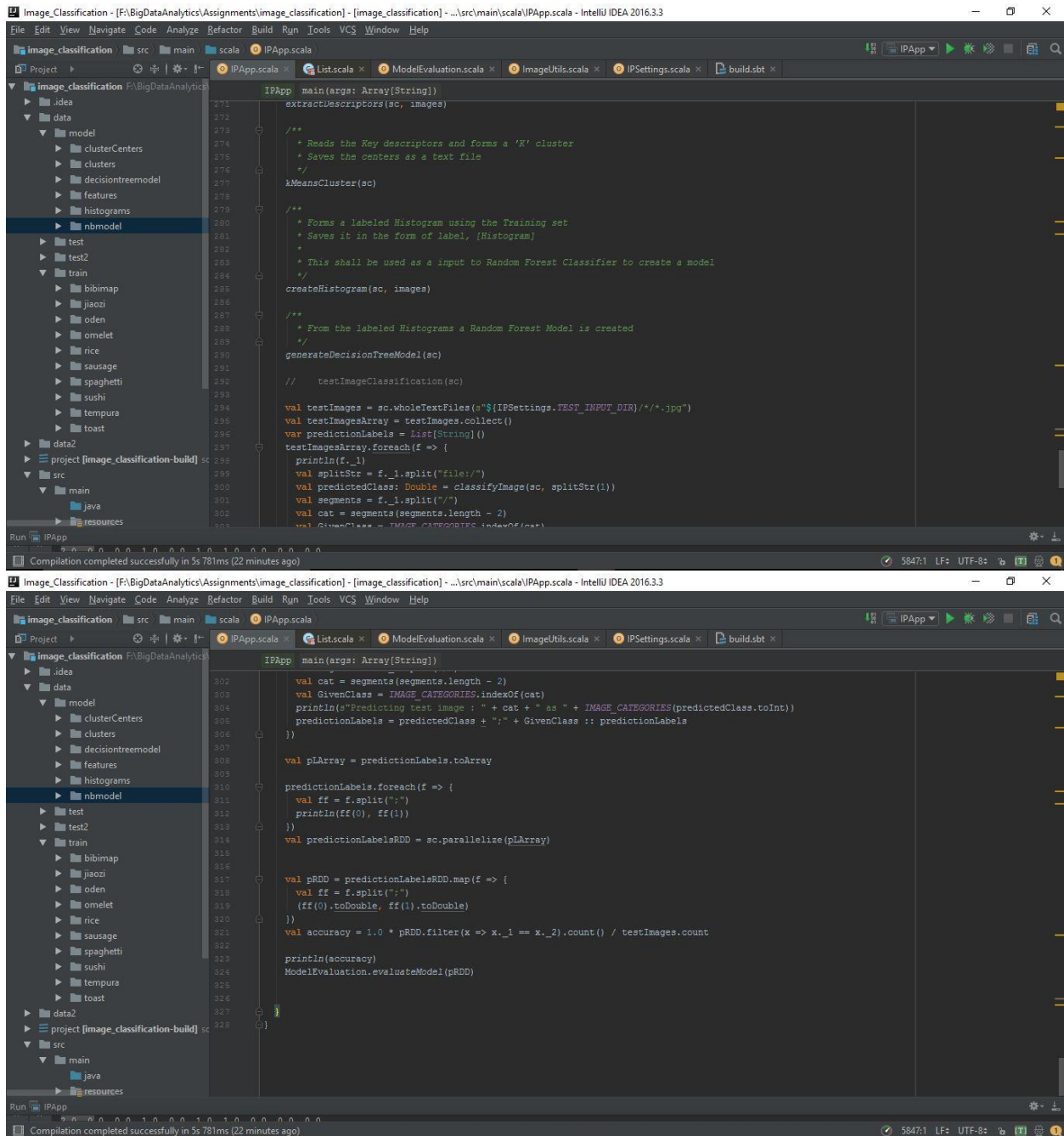
```
Image_Classification - [F:\BigDataAnalytics\Assignments\image_classification] - [image_classification] - \src\main\scala\IPApp.scala - IntelliJ IDEA 2016.3.3
File Edit View Navigate Code Analyze Refactor Build Run Tools VCS Window Help
image_classification src main scala IPApp.scala
Project image_classification
  .idea
  data
    model
      clusterCenters
      clusters
      decisiontreemodel
      features
      histograms
      nbmodel
    test
    test2
    train
      bibimap
      jiaozi
      oden
      omelet
      rice
      sausage
      spaghetti
      sushi
      tempura
      toast
    data2
  project [image_classification-build]
  src
    main
      java
      resources

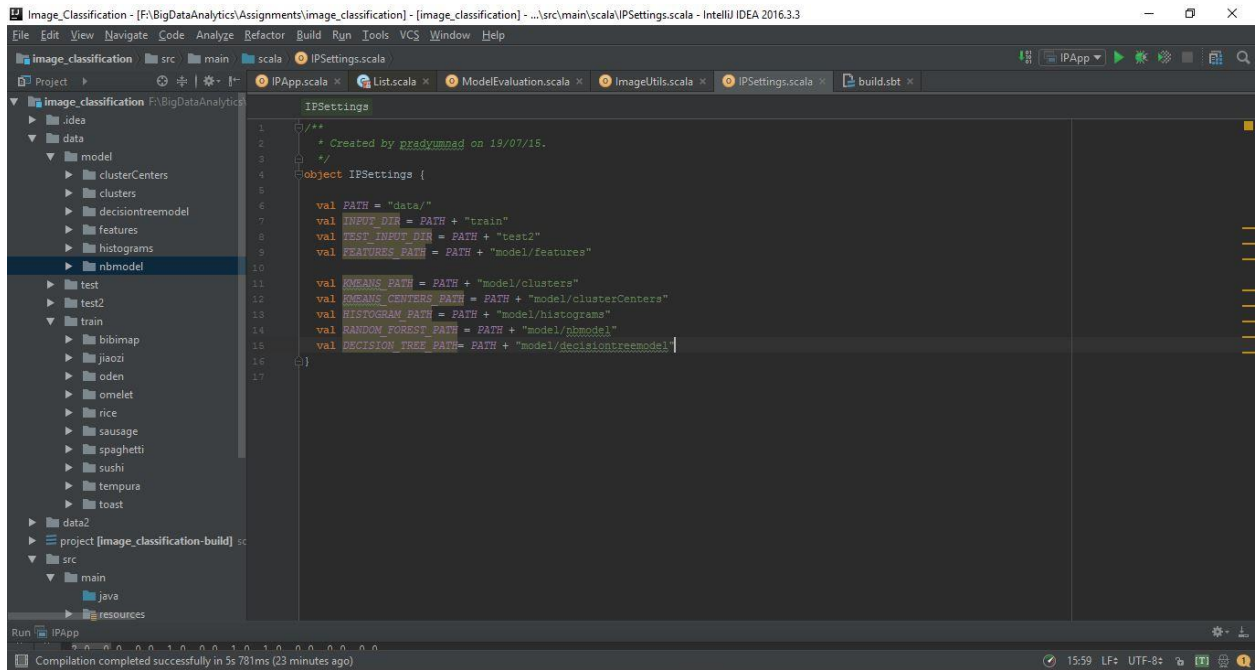
IPApp main(args: Array[String])
36
37   val data = images.map {
38     case (name, contents) => {
39       val desc = ImageUtils.descriptors(name.split("/").(1))
40       val list = ImageUtils.matToString(desc)
41       println("--- " + list.size)
42       list
43     }
44   }.reduce((x, y) => x ::: y)
45
46   val featuresSeq = sc.parallelize(data)
47   featuresSeq.saveAsTextFile(IPSettings.FEATURES_PATH)
48   println("Total size : " + data.size)
49
50
51   def KMeansCluster(sc: SparkContext): Unit = {
52     if (Files.exists(Paths.get(IPSettings.KMEANS_PATH))) {
53       println(s"${IPSettings.KMEANS_PATH} exists, skipping clusters formation..")
54       return
55     }
56
57     // Load and parse the data
58     val data = sc.textFile(IPSettings.FEATURES_PATH)
59     val parsedData = data.map(s => Vectors.dense(s.split(' ').map(_.toDouble)))
60
61     // Cluster the data into (#400) classes using KMeans
62     val numClusters = 400
63     val numIterations = 20
64     val clusters = KMeans.train(parsedData, numClusters, numIterations)
65
66     // Evaluate clustering by computing Within Set Sum of Squared Errors
67
68     Run IPApp
69     Compilation completed successfully in 5s 781ms (18 minutes ago)
```











```

Image_Classification - [FBI\BigDataAnalytics\Assignments\image_classification] - ...src\main\scala\IPApp.scala - IntelliJ IDEA 2016.3.3
File Edit View Navigate Code Analyze Refactor Build Run Tools VCS Window Help

image_classification src main scala IPApp.scala

Run IPApp

"C:\Program Files\Java\jdk1.8.0_77\bin\java" ...
Using Spark's default log4j profile: org/apache/spark/log4j-defaults.properties
17/02/17 01:19:08 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
17/02/17 01:19:09 INFO Slf4jLogger: Slf4jLogger started
17/02/17 01:19:09 INFO Remoting: Starting remoting
17/02/17 01:19:09 INFO Remoting: Remoting started; listening on addresses :[akka.tcp://sparkDriverActorSystem@192.168.1.167:55403]
data/model/features exists, skipping feature extraction..
data/model/clusters exists, skipping clusters formation..
data/model/histograms exists, skipping histograms creation..
numTrees 4 featureSubsetStrategy all impurity gini maxDepth 3
17/02/17 01:19:12 INFO FileInputFormat: Total input paths to process : 4
17/02/17 01:19:12 INFO deprecation: mapred.tip.id is deprecated. Instead, use mapreduce.task.id
17/02/17 01:19:12 INFO deprecation: mapred.task.id is deprecated. Instead, use mapreduce.task.attempt.id
17/02/17 01:19:12 INFO deprecation: mapred.task.is.map is deprecated. Instead, use mapreduce.task.ismap
17/02/17 01:19:12 INFO deprecation: mapred.task.partition is deprecated. Instead, use mapreduce.task.partition
17/02/17 01:19:12 INFO deprecation: mapred.job.id is deprecated. Instead, use mapreduce.job.id
Test Error = 0.6875

|===== Confusion matrix =====
11.0 1.0 3.0 3.0 1.0 1.0 1.0
2.0 8.0 2.0 2.0 5.0 0.0 3.0
0.0 0.0 1.0 0.0 0.0 0.0 0.0
0.0 0.0 0.0 3.0 0.0 0.0 0.0
0.0 0.0 0.0 0.0 4.0 0.0 0.0
0.0 4.0 4.0 2.0 3.0 5.0 3.0
0.0 0.0 0.0 0.0 0.0 0.0 3.0
0.3125
numTrees 4 featureSubsetStrategy all impurity gini maxDepth 4
Test Error = 0.625

|===== Confusion matrix =====
2.0 0.0 0.0 0.0 0.0 0.0 1.0 0.0 0.0
2.0 11.0 1.0 2.0 5.0 0.0 1.0 5.0 3.0
0.0 1.0 4.0 1.0 1.0 0.0 1.0 0.0 0.0
0.0 0.0 0.0 3.0 0.0 0.0 0.0 0.0 0.0
0.0 0.0 0.0 0.0 4.0 0.0 0.0 0.0 0.0
0.0 0.0 2.0 1.0 2.0 5.0 0.0 2.0 3.0
0.0 0.0 0.0 0.0 0.0 0.0 1.0 0.0 0.0
0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
Compilation completed successfully in 5s 781ms (23 minutes ago)
254.34 Lf UTF-8

```



```
Image_Classification - [F:\BigDataAnalytics\Assignments\image_classification] - [image_classification] - \src\main\scala\IPApp.scala - IntelliJ IDEA 2016.3.3
File Edit View Navigate Code Analyze Refactor Build Run Tools VCS Window Help
image_classification src main scala IPApp.scala
Run IPApp
13.0 0.0 2.0 0.0 2.0 1.0 1.0 0.0 0.0
0.0 11.0 2.0 0.0 0.0 0.0 2.0 0.0 0.0
0.0 0.0 5.0 0.0 0.0 0.0 0.0 0.0 0.0
0.0 0.0 0.0 10.0 0.0 0.0 0.0 2.0 1.0
0.0 0.0 0.0 0.0 11.0 0.0 0.0 0.0 0.0
0.0 0.0 0.0 0.0 0.0 4.0 0.0 0.0 0.0
0.0 0.0 1.0 0.0 0.0 1.0 8.0 0.0 0.0
0.0 0.0 0.0 0.0 0.0 0.0 1.0 9.0 0.0
0.0 2.0 0.0 0.0 0.0 0.0 0.0 5.0 9.0
0.7142857142857143
numTrees 4 featureSubsetStrategy onethird impurity gini maxDepth 3
Test Error = 0.6875
|===== Confusion matrix =====
11.0 1.0 3.0 3.0 1.0 1.0 1.0
2.0 8.0 2.0 2.0 5.0 0.0 3.0
0.0 0.0 1.0 0.0 0.0 0.0 0.0
0.0 0.0 0.0 3.0 0.0 0.0 0.0
0.0 0.0 0.0 0.0 4.0 0.0 0.0
0.0 4.0 4.0 2.0 3.0 5.0 3.0
0.0 0.0 0.0 0.0 0.0 0.0 3.0
0.3125
numTrees 4 featureSubsetStrategy onethird impurity gini maxDepth 4
Test Error = 0.625
|===== Confusion matrix =====
2.0 0.0 0.0 0.0 0.0 0.0 1.0 0.0 0.0
2.0 11.0 1.0 2.0 5.0 0.0 1.0 5.0 3.0
0.0 1.0 4.0 1.0 1.0 0.0 1.0 0.0 0.0
0.0 0.0 0.0 3.0 0.0 0.0 0.0 0.0 0.0
0.0 0.0 0.0 0.0 4.0 0.0 0.0 0.0 0.0
0.0 0.0 2.0 1.0 2.0 5.0 0.0 2.0 3.0
0.0 0.0 0.0 0.0 0.0 0.0 1.0 0.0 0.0
9.0 1.0 3.0 3.0 1.0 1.0 8.0 9.0 1.0
0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 3.0
0.375
numTrees 4 featureSubsetStrategy onethird impurity gini maxDepth 5
Compilation completed successfully in 5s 781ms (24 minutes ago) 25628 LF+ UTF-8+

Image_Classification - [F:\BigDataAnalytics\Assignments\image_classification] - [image_classification] - \src\main\scala\IPApp.scala - IntelliJ IDEA 2016.3.3
File Edit View Navigate Code Analyze Refactor Build Run Tools VCS Window Help
image_classification src main scala IPApp.scala
Run IPApp
2.0 11.0 1.0 2.0 3.0 0.0 1.0 1.0 3.0 2.0
0.0 0.0 4.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
0.0 0.0 0.0 4.0 0.0 0.0 0.0 1.0 0.0 0.0
0.0 1.0 0.0 0.0 7.0 0.0 0.0 0.0 2.0 1.0
0.0 0.0 2.0 1.0 2.0 5.0 1.0 0.0 2.0 0.0
8.0 1.0 3.0 3.0 1.0 1.0 7.0 8.0 3.0 1.0
0.0 0.0 0.0 0.0 0.0 0.0 0.0 2.0 0.0 0.0
1.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 6.0 0.0
0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 6.0
0.48214285714285715
numTrees 5 featureSubsetStrategy all impurity gini maxDepth 6
Test Error = 0.4017857142857143
|===== Confusion matrix =====
3.0 0.0 0.0 0.0 1.0 0.0 0.0 0.0 0.0 0.0
0.0 11.0 0.0 0.0 0.0 0.0 0.0 0.0 1.0 0.0
0.0 0.0 4.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
0.0 0.0 0.0 4.0 0.0 0.0 0.0 0.0 0.0 0.0
0.0 0.0 0.0 0.0 7.0 0.0 0.0 1.0 0.0 1.0
0.0 0.0 1.0 0.0 2.0 5.0 0.0 0.0 1.0 0.0
7.0 1.0 3.0 2.0 1.0 0.0 8.0 1.0 2.0 1.0
1.0 0.0 1.0 2.0 0.0 1.0 0.0 9.0 2.0 0.0
0.0 1.0 0.0 0.0 0.0 0.0 0.0 0.0 8.0 0.0
2.0 0.0 1.0 2.0 2.0 0.0 1.0 1.0 2.0 8.0
0.5982142857142857
numTrees 5 featureSubsetStrategy all impurity entropy maxDepth 3
Test Error = 0.6875
|===== Confusion matrix =====
11.0 1.0 0.0 1.0 1.0 8.0 0.0
0.0 2.0 1.0 1.0 0.0 0.0 0.0
0.0 2.0 8.0 0.0 1.0 1.0 6.0
0.0 4.0 0.0 7.0 0.0 0.0 2.0
0.0 0.0 1.0 2.0 4.0 2.0 0.0
0.0 0.0 0.0 1.0 0.0 1.0 0.0
2.0 4.0 0.0 1.0 0.0 0.0 2.0
0.3125
numTrees 5 featureSubsetStrategy all impurity entropy maxDepth 4
Compilation completed successfully in 5s 781ms (25 minutes ago) 25628 LF+ UTF-8+

```

```
Image_Classification - [F:\BigDataAnalytics\Assignments\image_classification] - [image_classification] - \src\main\scala\IPApp.scala - IntelliJ IDEA 2016.3.3
File Edit View Navigate Code Analyze Refactor Build Run Tools VCS Window Help
image_classification src main scala IPApp.scala
Run IPApp
0.0 0.0 0.0 0.0 10.0 0.0 0.0 1.0 0.0
0.0 0.0 0.0 0.0 0.0 4.0 0.0 1.0 1.0
0.0 0.0 0.0 0.0 0.0 0.0 1.0 0.0 0.0
0.0 1.0 0.0 0.0 1.0 0.0 2.0 4.0 0.0
2.0 2.0 2.0 0.0 1.0 0.0 0.0 5.0 1.0
0.45535714285714285
numTrees 5 featureSubsetStrategy onethird impurity entropy maxDepth 6
Test Error = 0.2857142857142857
|===== Confusion matrix =====
13.0 0.0 2.0 0.0 2.0 1.0 1.0 0.0 0.0
0.0 11.0 2.0 0.0 0.0 0.0 2.0 0.0 0.0
0.0 0.0 5.0 0.0 0.0 0.0 0.0 0.0 0.0
0.0 0.0 0.0 10.0 0.0 0.0 0.0 2.0 1.0
0.0 0.0 0.0 0.0 11.0 0.0 0.0 0.0 0.0
0.0 0.0 0.0 0.0 0.0 4.0 0.0 0.0 0.0
0.0 0.0 1.0 0.0 0.0 1.0 8.0 0.0 0.0
0.0 0.0 0.0 0.0 0.0 0.0 1.0 9.0 0.0
0.0 2.0 0.0 0.0 0.0 0.0 0.0 5.0 9.0
0.7142857142857143
numTrees 6 featureSubsetStrategy all impurity gini maxDepth 3
Test Error = 0.6075
|===== Confusion matrix =====
11.0 1.0 3.0 3.0 1.0 1.0 1.0
2.0 8.0 2.0 2.0 5.0 0.0 3.0
0.0 0.0 1.0 0.0 0.0 0.0 0.0
0.0 0.0 0.0 3.0 0.0 0.0 0.0
0.0 0.0 0.0 0.0 4.0 0.0 0.0
0.0 4.0 4.0 2.0 3.0 5.0 3.0
0.0 0.0 0.0 0.0 0.0 0.0 3.0
0.3125
numTrees 6 featureSubsetStrategy all impurity gini maxDepth 4
Test Error = 0.625
|===== Confusion matrix =====
2.0 0.0 0.0 0.0 0.0 0.0 1.0 0.0 0.0
2.0 11.0 1.0 2.0 5.0 0.0 1.0 5.0 3.0
Compilation completed successfully in 3s 781ms (25 minutes ago) 256x28 LF+ UTF-8
```

```
Image_Classification - [F:\BigDataAnalytics\Assignments\image_classification] - [image_classification] - \src\main\scala\IPApp.scala - IntelliJ IDEA 2016.3.3
File Edit View Navigate Code Analyze Refactor Build Run Tools VCS Window Help
image_classification src main scala IPApp.scala
Run IPApp
2.0 11.0 1.0 2.0 5.0 0.0 1.0 5.0 3.0
0.0 1.0 4.0 1.0 1.0 0.0 1.0 0.0 0.0
0.0 0.0 0.0 3.0 0.0 0.0 0.0 0.0 0.0
0.0 0.0 0.0 0.0 4.0 0.0 0.0 0.0 0.0
0.0 0.0 2.0 1.0 2.0 5.0 0.0 2.0 3.0
0.0 0.0 0.0 0.0 0.0 0.0 1.0 0.0 0.0
9.0 1.0 3.0 3.0 1.0 1.0 8.0 9.0 1.0
0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 3.0
0.375
numTrees 6 featureSubsetStrategy onethird impurity gini maxDepth 5
Test Error = 0.5178571428571429
|===== Confusion matrix =====
2.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
2.0 11.0 1.0 2.0 3.0 0.0 1.0 1.0 3.0 2.0
0.0 0.0 4.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
0.0 0.0 0.0 4.0 0.0 0.0 0.0 1.0 0.0 0.0
0.0 1.0 0.0 0.0 7.0 0.0 0.0 0.0 2.0 1.0
0.0 0.0 2.0 1.0 2.0 5.0 1.0 0.0 2.0 0.0
8.0 1.0 3.0 3.0 1.0 1.0 7.0 8.0 3.0 1.0
0.0 0.0 0.0 0.0 0.0 0.0 0.0 2.0 0.0 0.0
1.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 6.0 0.0
0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 6.0
0.48214285714285715
numTrees 6 featureSubsetStrategy onethird impurity gini maxDepth 6
Test Error = 0.4017857142857143
|===== Confusion matrix =====
3.0 0.0 0.0 0.0 1.0 0.0 0.0 0.0 0.0 0.0
0.0 11.0 0.0 0.0 0.0 0.0 0.0 0.0 1.0 0.0
0.0 0.0 4.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
0.0 0.0 0.0 4.0 0.0 0.0 0.0 0.0 0.0 0.0
0.0 0.0 0.0 0.0 7.0 0.0 0.0 1.0 0.0 1.0
0.0 0.0 1.0 0.0 2.0 5.0 0.0 0.0 1.0 0.0
7.0 1.0 3.0 2.0 1.0 0.0 8.0 1.0 2.0 1.0
1.0 0.0 1.0 2.0 0.0 1.0 0.0 9.0 2.0 0.0
0.0 1.0 0.0 0.0 0.0 0.0 0.0 0.0 8.0 0.0
2.0 0.0 1.0 2.0 2.0 0.0 1.0 1.0 2.0 8.0
Compilation completed successfully in 3s 781ms (25 minutes ago) 256x28 LF+ UTF-8
```

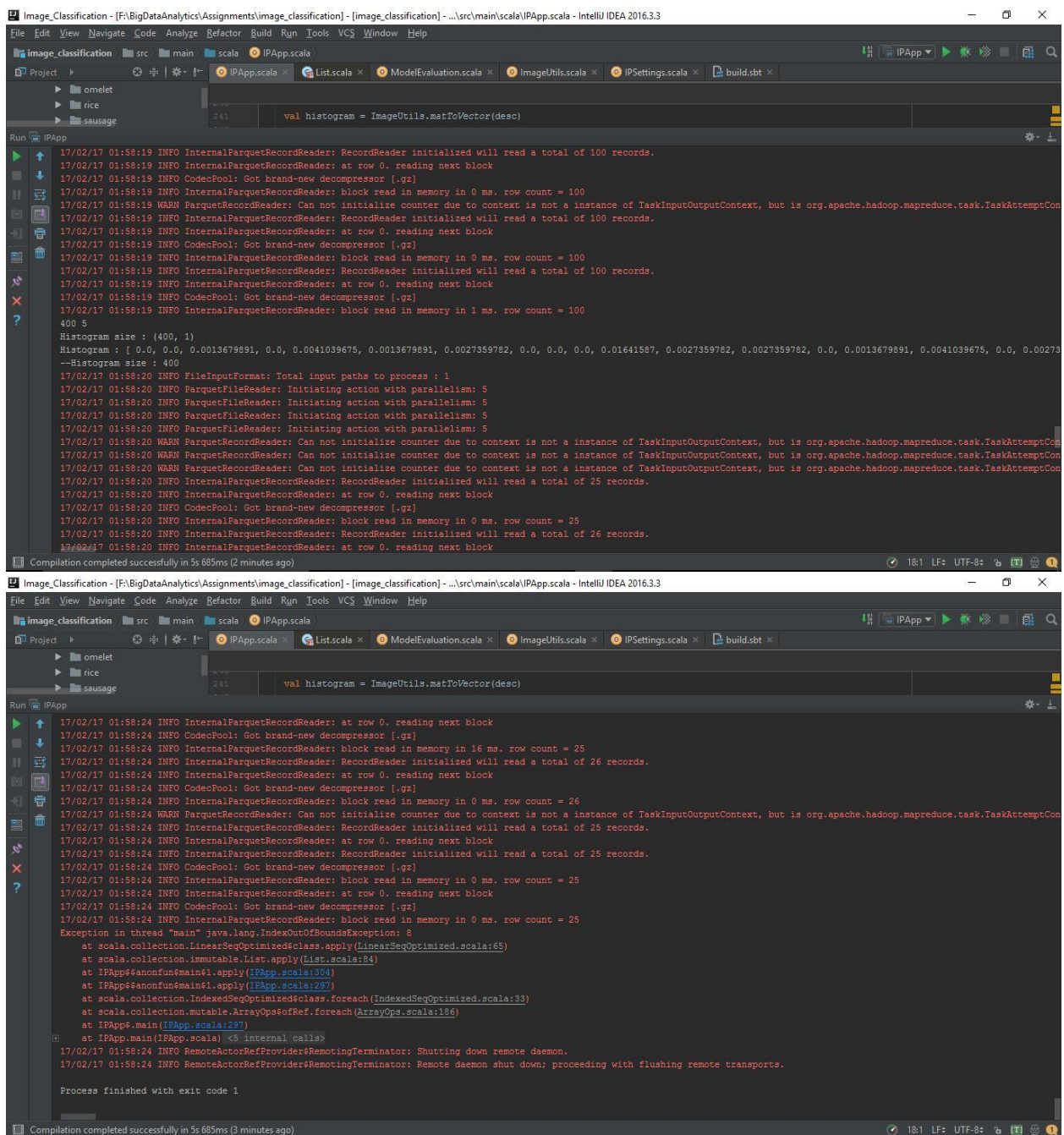
```
Image_Classification - [F:\BigDataAnalytics\Assignments\image_classification] - [image_classification] - \src\main\scala\IPApp.scala - IntelliJ IDEA 2016.3.3
File Edit View Navigate Code Analyze Refactor Build Run Tools VCS Window Help
image_classification src main scala IPApp.scala
Run IPApp
0.0 2.0 0.0 0.0 0.0 0.0 0.0 0.0 5.0 9.0
0.7142857142857143
numTrees 9 featureSubsetStrategy sqrt impurity gini maxDepth 3
Test Error = 0.6875
|===== Confusion matrix =====
11.0 1.0 3.0 3.0 1.0 1.0 1.0
2.0 8.0 2.0 2.0 5.0 0.0 3.0
0.0 0.0 1.0 0.0 0.0 0.0 0.0
0.0 0.0 0.0 3.0 0.0 0.0 0.0
0.0 0.0 0.0 0.0 4.0 0.0 0.0
0.0 4.0 4.0 2.0 3.0 5.0 3.0
0.0 0.0 0.0 0.0 0.0 0.0 3.0
0.3125
numTrees 9 featureSubsetStrategy sqrt impurity gini maxDepth 4
Test Error = 0.625
|===== Confusion matrix =====
2.0 0.0 0.0 0.0 0.0 0.0 1.0 0.0 0.0
2.0 11.0 1.0 2.0 5.0 0.0 1.0 5.0 3.0
0.0 1.0 4.0 1.0 1.0 0.0 1.0 0.0 0.0
0.0 0.0 0.0 3.0 0.0 0.0 0.0 0.0 0.0
0.0 0.0 0.0 0.0 4.0 0.0 0.0 0.0 0.0
0.0 0.0 2.0 1.0 2.0 5.0 0.0 2.0 3.0
0.0 0.0 0.0 0.0 0.0 0.0 1.0 0.0 0.0
9.0 1.0 3.0 3.0 1.0 1.0 8.0 9.0 1.0
0.0 0.0 0.0 0.0 0.0 0.0 0.0 3.0
0.375
numTrees 9 featureSubsetStrategy sqrt impurity gini maxDepth 5
Test Error = 0.5178571428571429
|===== Confusion matrix =====
2.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
2.0 11.0 1.0 2.0 3.0 0.0 1.0 1.0 3.0 2.0
0.0 0.0 4.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
0.0 0.0 0.0 4.0 0.0 0.0 0.0 1.0 0.0 0.0
0.0 1.0 0.0 0.0 7.0 0.0 0.0 0.0 2.0 1.0
0.0 0.0 2.0 1.0 2.0 5.0 1.0 0.0 2.0 0.0
Compilation completed successfully in 5s 781ms (25 minutes ago) 25628 LF: UTF-8
```

```
Image_Classification - [F:\BigDataAnalytics\Assignments\image_classification] - [image_classification] - \src\main\scala\IPApp.scala - IntelliJ IDEA 2016.3.3
File Edit View Navigate Code Analyze Refactor Build Run Tools VCS Window Help
image_classification src main scala IPApp.scala
Run IPApp
17/02/17 01:34:50 INFO InternalParquetRecordReader: block read in memory in 0 ms. row count = 100
17/02/17 01:34:50 WARN ParquetRecordReader: Can not initialize counter due to context is not a instance of TaskInputOutputContext, but is org.apache.hadoop.mapreduce.task.TaskAttemptContextImpl
17/02/17 01:34:50 INFO CodecPool: Got brand-new decompressor [.gz]
17/02/17 01:34:50 INFO InternalParquetRecordReader: block read in memory in 0 ms. row count = 100
17/02/17 01:34:50 INFO CodecPool: Got brand-new decompressor [.gz]
17/02/17 01:34:50 INFO InternalParquetRecordReader: block read in memory in 0 ms. row count = 100
17/02/17 01:34:50 INFO InternalParquetRecordReader: RecordReader initialized will read a total of 100 records.
17/02/17 01:34:50 INFO InternalParquetRecordReader: at row 0. reading next block
17/02/17 01:34:50 INFO CodecPool: Got brand-new decompressor [.gz]
17/02/17 01:34:50 INFO InternalParquetRecordReader: block read in memory in 0 ms. row count = 100
400 5
Histogram size : (400, 1)
--Histogram size : 400
17/02/17 01:34:50 INFO FileInputFormat: Total input paths to process : 1
17/02/17 01:34:50 INFO ParquetFileReader: Initiating action with parallelism: 5
17/02/17 01:34:50 INFO ParquetFileReader: Initiating action with parallelism: 5
17/02/17 01:34:50 INFO ParquetFileReader: Initiating action with parallelism: 5
17/02/17 01:34:50 INFO ParquetFileReader: Initiating action with parallelism: 5
17/02/17 01:34:50 WARN ParquetRecordReader: Can not initialize counter due to context is not a instance of TaskInputOutputContext, but is org.apache.hadoop.mapreduce.task.TaskAttemptContextImpl
17/02/17 01:34:50 WARN ParquetRecordReader: Can not initialize counter due to context is not a instance of TaskInputOutputContext, but is org.apache.hadoop.mapreduce.task.TaskAttemptContextImpl
17/02/17 01:34:50 WARN ParquetRecordReader: Can not initialize counter due to context is not a instance of TaskInputOutputContext, but is org.apache.hadoop.mapreduce.task.TaskAttemptContextImpl
17/02/17 01:34:50 INFO InternalParquetRecordReader: RecordReader initialized will read a total of 25 records.
17/02/17 01:34:50 INFO InternalParquetRecordReader: at row 0. reading next block
17/02/17 01:34:50 INFO InternalParquetRecordReader: RecordReader initialized will read a total of 25 records.
17/02/17 01:34:50 INFO InternalParquetRecordReader: at row 0. reading next block
17/02/17 01:34:50 INFO InternalParquetRecordReader: block read in memory in 1 ms. row count = 25
17/02/17 01:34:50 INFO CodecPool: Got brand-new decompressor [.gz]
17/02/17 01:34:50 INFO InternalParquetRecordReader: block read in memory in 0 ms. row count = 25
17/02/17 01:34:50 INFO InternalParquetRecordReader: RecordReader initialized will read a total of 25 records.
17/02/17 01:34:50 INFO InternalParquetRecordReader: at row 0. reading next block
17/02/17 01:34:50 INFO CodecPool: Got brand-new decompressor [.gz]
17/02/17 01:34:50 INFO InternalParquetRecordReader: block read in memory in 2 ms. row count = 25
17/02/17 01:34:50 WARN ParquetRecordReader: Can not initialize counter due to context is not a instance of TaskInputOutputContext, but is org.apache.hadoop.mapreduce.task.TaskAttemptContextImpl
Compilation completed successfully in 5s 781ms (26 minutes ago) 25628 LF: UTF-8
```

```
Image_Classification - [F:\BigDataAnalytics\Assignments\image_classification] - [image_classification] - ...src\main\scala\IPApp.scala - IntelliJ IDEA 2016.3.3
File Edit View Navigate Code Analyze Refactor Build Run Tools VCS Window Help
image_classification src main scala IPApp.scala
Run IPApp
17/02/17 01:35:48 INFO CodecPool: Got brand-new decompressor [.gz]
17/02/17 01:35:48 INFO InternalParquetRecordReader: block read in memory in 16 ms. row count = 25
17/02/17 01:35:48 INFO InternalParquetRecordReader: block read in memory in 16 ms. row count = 25
17/02/17 01:35:48 INFO InternalParquetRecordReader: block read in memory in 16 ms. row count = 26
17/02/17 01:35:48 INFO CodecPool: Got brand-new decompressor [.gz]
17/02/17 01:35:48 INFO InternalParquetRecordReader: block read in memory in 16 ms. row count = 25
Predicting test image : toast as oden
(7.0,2)
(0.0,2)
(2.0,2)
(0.0,2)
(3.0,2)
(1.0,1)
(0.0,1)
(5.0,1)
(5.0,1)
(8.0,1)
(5.0,4)
(2.0,4)
(0.0,4)
(7.0,4)
(4.0,4)
(5.0,5)
(2.0,5)
(5.0,5)
(8.0,5)
(4.0,5)
(4.0,6)
(6.0,6)
(4.0,6)
(9.0,6)
(0.0,6)
(9.0,0)
(0.0,0)
(9.0,0)
Compilation completed successfully in 5s 781ms (26 minutes ago) 25628 LF: UTF-8

Image_Classification - [F:\BigDataAnalytics\Assignments\image_classification] - [image_classification] - ...src\main\scala\IPApp.scala - IntelliJ IDEA 2016.3.3
File Edit View Navigate Code Analyze Refactor Build Run Tools VCS Window Help
image_classification src main scala IPApp.scala
Run IPApp
(5.0,8)
(0.0,8)
(2.0,8)
(3.0,7)
(9.0,7)
(5.0,7)
(3.0,7)
(0.0,7)
(8.0,9)
(9.0,9)
(3.0,9)
(9.0,9)
(0.0,9)
(6.0,3)
(5.0,3)
(0.0,3)
(0.0,3)
(3.0,3)
0.22
===== Confusion matrix =====
2.0 0.0 1.0 0.0 0.0 0.0 0.0 0.0 0.0 2.0
1.0 1.0 0.0 0.0 0.0 2.0 0.0 0.0 1.0 0.0
2.0 0.0 1.0 1.0 0.0 0.0 0.0 1.0 0.0 0.0
2.0 0.0 0.0 1.0 0.0 1.0 1.0 0.0 0.0 0.0
1.0 0.0 1.0 0.0 1.0 1.0 0.0 1.0 0.0 0.0
0.0 0.0 1.0 0.0 1.0 2.0 0.0 0.0 1.0 0.0
1.0 0.0 0.0 0.0 2.0 0.0 1.0 0.0 0.0 1.0
1.0 0.0 1.0 0.0 1.0 1.0 0.0 0.0 0.0 1.0
1.0 0.0 0.0 1.0 0.0 0.0 0.0 0.0 1.0 2.0
0.22
17/02/17 01:35:50 INFO RemoteActorRefProviders$RemotingTerminator: Shutting down remote daemon.
Process finished with exit code 0
Compilation completed successfully in 5s 781ms (27 minutes ago) 25628 LF: UTF-8
```

Output for other data set:



Image_Classification - [F:\BigDataAnalytics\Assignments\image_classification] - [image_classification] - \src\main\scala\IPSettings.scala - IntelliJ IDEA 2016.3.3

File Edit View Navigate Code Analyze Refactor Build Run Tools VCS Window Help

image_classification data model nbmodel

Project

- omelet
- rice
- sausage
- spaghetti
- sushi
- tempura
- toast
- data2
- DataLab4
 - rmodel
 - test
 - airplane
 - h-AIRPLANES-MOBILE-628
 - HVxk2T5c9051Pgq1cmU
 - images.jpeg
 - MjgwMTc4OA.jpeg
 - noaa-hurricane-hunter-jet
 - apple
 - images (1).jpeg
 - images (2).jpeg
 - images.jpeg
 - Table-of-RedYellow-Apple
 - Yellow-Green-And-Red-Apple
 - cats
 - eggs
 - noodles
 - wolves
 - train
 - airplane
 - apple

```
1  /**
2   * Created by pradyumnad on 19/07/16.
3   */
4   object IPSettings {
5
6       val PATH = "DataLab4/"
7       val INPUT_DIR = PATH + "train"
8       val TEST_INPUT_DIR = PATH + "test"
9       val FEATURES_PATH = PATH + "model/features"
10
11       val KMEANS_PATH = PATH + "model/clusters"
12       val KMEANS_CENTERS_PATH = PATH + "model/clusterCenters"
13       val HISTOGRAM_PATH = PATH + "model/histograms"
14       val RANDOM_FOREST_PATH = PATH + "model/rmodel"
15       val DECISION_TREE_PATH = PATH + "model/decisiontreemodel"
16   }
17
```

Compilation completed successfully in 5s 685ms (2 minutes ago)

836 LF+ UTF-8

Image_Classification - [F:\BigDataAnalytics\Assignments\image_classification] - [image_classification] - \src\main\scala\IPApp.scala - IntelliJ IDEA 2016.3.3

File Edit View Navigate Code Analyze Refactor Build Run Tools VCS Window Help

image_classification src main scala IPApp.scala

Project

- omelet
- rice
- sausage
- spaghetti
- sushi
- tempura
- toast
- data2
- DataLab4
 - rmodel
 - test
 - airplane
 - h-AIRPLANES-MOBILE-628
 - HVxk2T5c9051Pgq1cmU
 - images.jpeg
 - MjgwMTc4OA.jpeg
 - noaa-hurricane-hunter-jet
 - apple
 - images (1).jpeg
 - images (2).jpeg
 - images.jpeg
 - Table-of-RedYellow-Apple
 - Yellow-Green-And-Red-Apple
 - cats
 - eggs
 - noodles
 - wolves
 - train
 - airplane
 - apple

```
1  /**...*/
2
3   import java.nio.file.{Files, Paths}
4
5   import org.apache.spark.mllib.clustering.{KMeans, KMeansModel}
6   import org.apache.spark.mllib.linalg.Vectors
7   import org.apache.spark.mllib.regression.LabeledPoint
8   import org.apache.spark.mllib.tree.DecisionTree
9   import org.apache.spark.mllib.tree.model.DecisionTreeModel
10  import org.apache.spark.rdd.RDD
11  import org.apache.spark.streaming.{Seconds, StreamingContext}
12  import org.apache.spark.{SparkConf, SparkContext}
13  import org.bytedeco.javacpp.opencv_highgui._
14
15  import scala.collection.mutable
16
17  object IPApp {
18      val featureVectorsCluster = new mutable.MutableList[String]
19
20      val IMAGE_CATEGORIES = List("airplane", "apple", "cats", "eggs", "noodles", "wolves")
21      //val IMAGE_CATEGORIES = List("accordion", "airplanes", "anchor", "ant", "barrel", "bass", "beaver", "binoculars", "bonsai")
22
23      /**
24       * @param sc : SparkContext
25       * @param images : Images list from the training set
26       */
27      def extractDescriptors(sc: SparkContext, images: RDD[String]): Unit = {
28          if (Files.exists(Paths.get(IPSettings.FEATURES_PATH))) {
29              println(s"${IPSettings.FEATURES_PATH} exists, skipping feature extraction..")
30              return
31          }
32      }
33  }
```

Compilation completed successfully in 5s 685ms (2 minutes ago)

181 LF+ UTF-8