DATA WAREHOUSE & MINING
TA COURSE #08

# CLASSIFICATION

TA COURSE #08
CLASSIFICATION
_____

ALGO
EXAMPLES

# LOGISTIC REGRESSION

▸ A popular method to predict a categorical response.

   ▸ It is a special case of Generalized Linear models that predicts the probability of the outcomes.

▸ Using binomial logistic regression to predict a binary outcome

▸ …Or using multinomial logistic regression to predict a multiclass outcome

▸ Use the family parameter to select between these two algorithms, or leave it unset and Spark will infer the correct variant.

# BINOMIAL LOGISTIC REGRESSION

```scala
1  import org.apache.spark.ml.classification.LogisticRegression
2
3  // Load training data
4  val training = spark.read.format("libsvm").load("data/mllib/sample_libsvm_data.txt")
5
6  val lr = new LogisticRegression()
7    .setMaxIter(10)
8    .setRegParam(0.3)
9    .setElasticNetParam(0.8)
10
11 // Fit the model
12 val lrModel = lr.fit(training)
13
14 // Print the coefficients and intercept for logistic regression
15 println(s"Coefficients: ${lrModel.coefficients} Intercept: ${lrModel.intercept}")
16
17 // We can also use the multinomial family for binary classification
18 val mlr = new LogisticRegression()
19    .setMaxIter(10)
20    .setRegParam(0.3)
21    .setElasticNetParam(0.8)
22    .setFamily("multinomial")
23
24 val mlrModel = mlr.fit(training)
25
26 // Print the coefficients and intercepts for logistic regression with multinomial family
27 println(s"Multinomial coefficients: ${mlrModel.coefficientMatrix}")
28 println(s"Multinomial intercepts: ${mlrModel.interceptVector}")
```

# MULTINOMIAL LOGISTIC REGRESSION

```scala
1  import org.apache.spark.ml.classification.LogisticRegression
2
3  // Load training data
4  val training = spark
5    .read
6    .format("libsvm")
7    .load("data/mllib/sample_multiclass_classification_data.txt")
8
9  val lr = new LogisticRegression()
10    .setMaxIter(10)
11    .setRegParam(0.3)
12    .setElasticNetParam(0.8)
13
14  // Fit the model
15  val lrModel = lr.fit(training)
16
17  // Print the coefficients and intercept for multinomial logistic regression
18  println(s"Coefficients: \n${lrModel.coefficientMatrix}")
19  println(s"Intercepts: \n${lrModel.interceptVector}")
20
21  val trainingSummary = lrModel.summary
22
23  // Obtain the objective per iteration
24  val objectiveHistory = trainingSummary.objectiveHistory
25  println("objectiveHistory:")
26  objectiveHistory.foreach(println)
27
28  // for multiclass, we can inspect metrics on a per-label basis
29  println("False positive rate by label:")
30  trainingSummary.falsePositiveRateByLabel.zipWithIndex.foreach { case (rate, label) =>
31    println(s"label $label: $rate")
32  }
33
34  println("True positive rate by label:")
35  trainingSummary.truePositiveRateByLabel.zipWithIndex.foreach { case (rate, label) =>
```

# MULTINOMIAL LOGISTIC REGRESSION

```scala
36      println(s"label $label: $rate")
37  }
38
39  println("Precision by label:")
40  trainingSummary.precisionByLabel.zipWithIndex.foreach { case (prec, label) =>
41      println(s"label $label: $prec")
42  }
43
44  println("Recall by label:")
45  trainingSummary.recallByLabel.zipWithIndex.foreach { case (rec, label) =>
46      println(s"label $label: $rec")
47  }
48
49
50  println("F-measure by label:")
51  trainingSummary.fMeasureByLabel.zipWithIndex.foreach { case (f, label) =>
52      println(s"label $label: $f")
53  }
54
55  val accuracy = trainingSummary.accuracy
56  val falsePositiveRate = trainingSummary.weightedFalsePositiveRate
57  val truePositiveRate = trainingSummary.weightedTruePositiveRate
58  val fMeasure = trainingSummary.weightedFMeasure
59  val precision = trainingSummary.weightedPrecision
60  val recall = trainingSummary.weightedRecall
61  println(s"Accuracy: $accuracy\nFPR: $falsePositiveRate\nTPR: $truePositiveRate\n" +
62      s"F-measure: $fMeasure\nPrecision: $precision\nRecall: $recall")
```

# RANDOM FOREST CLASSIFIER

▸ Ensembles of decision trees.

▸ Random forests combine many decision trees in order to reduce the risk of overfitting.

▸ The spark.ml implementation supports random forests for binary and multiclass classification and for regression, using both continuous and categorical features.

# RANDOM FOREST CLASSIFIER

```scala
1  import org.apache.spark.ml.Pipeline
2  import org.apache.spark.ml.classification.{RandomForestClassificationModel, RandomForestClassifier}
3  import org.apache.spark.ml.evaluation.MulticlassClassificationEvaluator
4  import org.apache.spark.ml.feature.{IndexToString, StringIndexer, VectorIndexer}
5
6  // Load and parse the data file, converting it to a DataFrame.
7  val data = spark.read.format("libsvm").load("data/mllib/sample_libsvm_data.txt")
8
9  // Index labels, adding metadata to the label column.
10 // Fit on whole dataset to include all labels in index.
11 val labelIndexer = new StringIndexer()
12   .setInputCol("label")
13   .setOutputCol("indexedLabel")
14   .fit(data)
15 // Automatically identify categorical features, and index them.
16 // Set maxCategories so features with > 4 distinct values are treated as continuous.
17 val featureIndexer = new VectorIndexer()
18   .setInputCol("features")
19   .setOutputCol("indexedFeatures")
20   .setMaxCategories(4)
21   .fit(data)
22
23 // Split the data into training and test sets (30% held out for testing).
24 val Array(trainingData, testData) = data.randomSplit(Array(0.7, 0.3))
25
26 // Train a RandomForest model.
27 val rf = new RandomForestClassifier()
28   .setLabelCol("indexedLabel")
29   .setFeaturesCol("indexedFeatures")
30   .setNumTrees(10)
31
32 // Convert indexed labels back to original labels.
33 val labelConverter = new IndexToString()
34   .setInputCol("prediction")
35   .setOutputCol("predictedLabel")
```

# RANDOM FOREST CLASSIFIER

```scala
36      .setLabels(labelIndexer.labels)
37
38 // Chain indexers and forest in a Pipeline.
39 val pipeline = new Pipeline()
40    .setStages(Array(labelIndexer, featureIndexer, rf, labelConverter))
41
42 // Train model. This also runs the indexers.
43 val model = pipeline.fit(trainingData)
44
45 // Make predictions.
46 val predictions = model.transform(testData)
47
48 // Select example rows to display.
49 predictions.select("predictedLabel", "label", "features").show(5)
50
51 // Select (prediction, true label) and compute test error.
52 val evaluator = new MulticlassClassificationEvaluator()
53    .setLabelCol("indexedLabel")
54    .setPredictionCol("prediction")
55    .setMetricName("accuracy")
56 val accuracy = evaluator.evaluate(predictions)
57 println(s"Test Error = ${(1.0 - accuracy)}")
58
59 val rfModel = model.stages(2).asInstanceOf[RandomForestClassificationModel]
60 println(s"Learned classification forest model:\n ${rfModel.toDebugString}")
```

# GRADIENT-BOOSTED TREE CLASSIFIER

▸ Ensembles of decision trees.

▸ GBTs iteratively train decision trees in order to minimize a loss function.

▸ The spark.ml implementation supports GBTs for binary classification and for regression, using both continuous and categorical features.

# GRADIENT-BOOSTED TREE CLASSIFIER

```scala
1  import org.apache.spark.ml.Pipeline
2  import org.apache.spark.ml.classification.{GBTClassificationModel, GBTClassifier}
3  import org.apache.spark.ml.evaluation.MulticlassClassificationEvaluator
4  import org.apache.spark.ml.feature.{IndexToString, StringIndexer, VectorIndexer}
5
6  // Load and parse the data file, converting it to a DataFrame.
7  val data = spark.read.format("libsvm").load("data/mllib/sample_libsvm_data.txt")
8
9  // Index labels, adding metadata to the label column.
10 // Fit on whole dataset to include all labels in index.
11 val labelIndexer = new StringIndexer()
12   .setInputCol("label")
13   .setOutputCol("indexedLabel")
14   .fit(data)
15 // Automatically identify categorical features, and index them.
16 // Set maxCategories so features with > 4 distinct values are treated as continuous.
17 val featureIndexer = new VectorIndexer()
18   .setInputCol("features")
19   .setOutputCol("indexedFeatures")
20   .setMaxCategories(4)
21   .fit(data)
22
23 // Split the data into training and test sets (30% held out for testing).
24 val Array(trainingData, testData) = data.randomSplit(Array(0.7, 0.3))
25
26 // Train a GBT model.
27 val gbt = new GBTClassifier()
28   .setLabelCol("indexedLabel")
29   .setFeaturesCol("indexedFeatures")
30   .setMaxIter(10)
31   .setFeatureSubsetStrategy("auto")
32
33 // Convert indexed labels back to original labels.
34 val labelConverter = new IndexToString()
35   .setInputCol("prediction")
```

# GRADIENT-BOOSTED TREE CLASSIFIER

```scala
36      .setOutputCol("predictedLabel")
37      .setLabels(labelIndexer.labels)
38
39  // Chain indexers and GBT in a Pipeline.
40  val pipeline = new Pipeline()
41      .setStages(Array(labelIndexer, featureIndexer, gbt, labelConverter))
42
43  // Train model. This also runs the indexers.
44  val model = pipeline.fit(trainingData)
45
46  // Make predictions.
47  val predictions = model.transform(testData)
48
49  // Select example rows to display.
50  predictions.select("predictedLabel", "label", "features").show(5)
51
52  // Select (prediction, true label) and compute test error.
53  val evaluator = new MulticlassClassificationEvaluator()
54      .setLabelCol("indexedLabel")
55      .setPredictionCol("prediction")
56      .setMetricName("accuracy")
57  val accuracy = evaluator.evaluate(predictions)
58  println(s"Test Error = ${1.0 - accuracy}")
59
60  val gbtModel = model.stages(2).asInstanceOf[GBTClassificationModel]
61  println(s"Learned classification GBT model:\n ${gbtModel.toDebugString}")
```

# MULTILAYER PERCEPTRON CLASSIFIER

▸ A classifier based on the feedforward artificial neural network.

▸ MLPC consists of multiple layers of nodes. Each layer is fully connected to the next layer in the network. Nodes in the input layer represent the input data. All other nodes map inputs to outputs by a linear combination of the inputs with the node's weights w and bias b and applying an activation function.

▸ MLPC employs backpropagation for learning the model. We use the logistic loss function for optimization and L-BFGS as an optimization routine.

# MULTILAYER PERCEPTRON CLASSIFIER

```scala
1  import org.apache.spark.ml.classification.MultilayerPerceptronClassifier
2  import org.apache.spark.ml.evaluation.MulticlassClassificationEvaluator
3
4  // Load the data stored in LIBSVM format as a DataFrame.
5  val data = spark.read.format("libsvm")
6    .load("data/mllib/sample_multiclass_classification_data.txt")
7
8  // Split the data into train and test
9  val splits = data.randomSplit(Array(0.6, 0.4), seed = 1234L)
10 val train = splits(0)
11 val test = splits(1)
12
13 // specify layers for the neural network:
14 // input layer of size 4 (features), two intermediate of size 5 and 4
15 // and output of size 3 (classes)
16 val layers = Array[Int](4, 5, 4, 3)
17
18 // create the trainer and set its parameters
19 val trainer = new MultilayerPerceptronClassifier()
20   .setLayers(layers)
21   .setBlockSize(128)
22   .setSeed(1234L)
23   .setMaxIter(100)
24
25 // train the model
26 val model = trainer.fit(train)
27
28 // compute accuracy on the test set
29 val result = model.transform(test)
30 val predictionAndLabels = result.select("prediction", "label")
31 val evaluator = new MulticlassClassificationEvaluator()
32   .setMetricName("accuracy")
33
34 println(s"Test set accuracy = ${evaluator.evaluate(predictionAndLabels)}")
```

# LINEAR SUPPORT VECTOR MACHINE

▸ A support vector machine constructs a hyperplane or set of hyperplanes in a high- or infinite-dimensional space, which can be used for classification, regression, or other tasks.

▸ Intuitively, a good separation is achieved by the hyperplane that has the largest distance to the nearest training-data points of any class (so-called functional margin), since in general the larger the margin the lower the generalization error of the classifier.

▸ LinearSVC in Spark ML supports binary classification with linear SVM. Internally, it optimizes the Hinge Loss using OWLQN optimizer.

# LINEAR SUPPORT VECTOR MACHINE

```scala
 1  import org.apache.spark.ml.classification.LinearSVC
 2
 3  // Load training data
 4  val training = spark.read.format("libsvm").load("data/mllib/sample_libsvm_data.txt")
 5
 6  val lsvc = new LinearSVC()
 7    .setMaxIter(10)
 8    .setRegParam(0.1)
 9
10  // Fit the model
11  val lsvcModel = lsvc.fit(training)
12
13  // Print the coefficients and intercept for linear svc
14  println(s"Coefficients: ${lsvcModel.coefficients} Intercept: ${lsvcModel.intercept}")
```

# ONE-VS-REST CLASSIFIER

▸ OneVsRest is an example of a machine learning reduction for performing multiclass classification given a base classifier that can perform binary classification efficiently.

  ▸ It is also known as "One-vs-All."

▸ OneVsRest is implemented as an Estimator. For the base classifier, it takes instances of Classifier and creates a binary classification problem for each of the k classes. The classifier for class i is trained to predict whether the label is i or not, distinguishing class i from all other classes.

▸ Predictions are done by evaluating each binary classifier and the index of the most confident classifier is output as label.

# ONE-VS-REST CLASSIFIER

```scala
1  import org.apache.spark.ml.classification.{LogisticRegression, OneVsRest}
2  import org.apache.spark.ml.evaluation.MulticlassClassificationEvaluator
3
4  // load data file.
5  val inputData = spark.read.format("libsvm")
6    .load("data/mllib/sample_multiclass_classification_data.txt")
7
8  // generate the train/test split.
9  val Array(train, test) = inputData.randomSplit(Array(0.8, 0.2))
10
11 // instantiate the base classifier
12 val classifier = new LogisticRegression()
13   .setMaxIter(10)
14   .setTol(1E-6)
15   .setFitIntercept(true)
16
17 // instantiate the One Vs Rest Classifier.
18 val ovr = new OneVsRest().setClassifier(classifier)
19
20 // train the multiclass model.
21 val ovrModel = ovr.fit(train)
22
23 // score the model on test data.
24 val predictions = ovrModel.transform(test)
25
26 // obtain evaluator.
27 val evaluator = new MulticlassClassificationEvaluator()
28   .setMetricName("accuracy")
29
30 // compute the classification error on test data.
31 val accuracy = evaluator.evaluate(predictions)
32 println(s"Test Error = ${1 - accuracy}")
```

# NAIVE BAYES

▸ Naive Bayes classifiers are a family of simple probabilistic, multiclass classifiers based on applying Bayes' theorem with strong (naive) independence assumptions between every pair of features.

▸ Naive Bayes can be trained very efficiently. With a single pass over the training data, it computes the conditional probability distribution of each feature given each label. For prediction, it applies Bayes' theorem to compute the conditional probability distribution of each label given an observation.

▸ MLlib supports both multinomial naive Bayes and Bernoulli naive Bayes.

# NAIVE BAYES

```scala
1  import org.apache.spark.ml.classification.NaiveBayes
2  import org.apache.spark.ml.evaluation.MulticlassClassificationEvaluator
3
4  // Load the data stored in LIBSVM format as a DataFrame.
5  val data = spark.read.format("libsvm").load("data/mllib/sample_libsvm_data.txt")
6
7  // Split the data into training and test sets (30% held out for testing)
8  val Array(trainingData, testData) = data.randomSplit(Array(0.7, 0.3), seed = 1234L)
9
10 // Train a NaiveBayes model.
11 val model = new NaiveBayes()
12   .fit(trainingData)
13
14 // Select example rows to display.
15 val predictions = model.transform(testData)
16 predictions.show()
17
18 // Select (prediction, true label) and compute test error
19 val evaluator = new MulticlassClassificationEvaluator()
20   .setLabelCol("label")
21   .setPredictionCol("prediction")
22   .setMetricName("accuracy")
23 val accuracy = evaluator.evaluate(predictions)
24 println(s"Test set accuracy = $accuracy")
```

# THX!