

Spark ML Tutorial

Let's make some practice! 😊

As of Spark 2.0, the RDD-based APIs in the **spark.mllib** package have entered maintenance mode. The primary Machine Learning API for Spark is now the DataFrame-based API in the **spark.ml** package.

- DataFrames provide a more user-friendly API than RDDs
- Many sources, queries, optimizers, ...
- Uniform API across ML algorithms and across multiple languages
- Facilitate practical ML Pipelines
- The RDD-based API is expected to be deprecated in Spark 2.3 and removed in Spark 3.0

Random Forest Classifier

```
from pyspark.ml import Pipeline
from pyspark.ml.classification import RandomForestClassifier
from pyspark.ml.feature import IndexToString, StringIndexer, VectorIndexer
from pyspark.ml.evaluation import MulticlassClassificationEvaluator

# Load and parse the data file, converting it to a DataFrame.
data = spark.read.format("libsvm").load("data/mllib/sample_libsvm_data.txt")

# Index labels, adding metadata to the label column.
# Fit on whole dataset to include all labels in index.
labelIndexer = StringIndexer(inputCol="label", outputCol="indexedLabel").fit(data)
```

Random Forest Classifier

```
# Automatically identify categorical features, and index them.  
# Set maxCategories so features with > 4 distinct values are treated as continuous.  
featureIndexer =\\  
    VectorIndexer(inputCol="features", outputCol="indexedFeatures", maxCategories=4).fit(data)  
  
# Split the data into training and test sets (30% held out for testing)  
(trainingData, testData) = data.randomSplit([0.7, 0.3])
```

Random Forest Classifier

```
# Train a RandomForest model.  
rf = RandomForestClassifier(labelCol="indexedLabel", featuresCol="indexedFeatures", numTrees=10)  
  
# Convert indexed labels back to original labels.  
labelConverter = IndexToString(inputCol="prediction", outputCol="predictedLabel",  
                                labels=labelIndexer.labels)
```

Random Forest Classifier

```
# Chain indexers and forest in a Pipeline
pipeline = Pipeline(stages=[labelIndexer, featureIndexer, rf, labelConverter])

# Train model. This also runs the indexers.
model = pipeline.fit(trainingData)

# Make predictions.
predictions = model.transform(testData)
```

Random Forest Classifier

```
# Select example rows to display.  
predictions.select("predictedLabel", "label", "features").show(5)  
  
# Select (prediction, true label) and compute test error  
evaluator = MulticlassClassificationEvaluator(  
    labelCol="indexedLabel", predictionCol="prediction", metricName="accuracy")  
accuracy = evaluator.evaluate(predictions)  
print("Test Error = %g" % (1.0 - accuracy))
```

Random Forest Classifier

```
rfModel = model.stages[2]
print(rfModel) # summary only
```

Linear Regression

```
from pyspark.ml.regression import LinearRegression

# Load training data
training = spark.read.format("libsvm")\
    .load("data/mllib/sample_linear_regression_data.txt")

lr = LinearRegression(maxIter=10, regParam=0.3, elasticNetParam=0.8)

# Fit the model
lrModel = lr.fit(training)

# Print the coefficients and intercept for linear regression
print("Coefficients: %s" % str(lrModel.coefficients))
print("Intercept: %s" % str(lrModel.intercept))
```

Linear Regression

```
# Summarize the model over the training set and print out some metrics
trainingSummary = lrModel.summary
print("numIterations: %d" % trainingSummary.totalIterations)
print("objectiveHistory: %s" % str(trainingSummary.objectiveHistory))
trainingSummary.residuals.show()
print("RMSE: %f" % trainingSummary.rootMeanSquaredError)
print("r2: %f" % trainingSummary.r2)
```

K-Means Clustering

```
from pyspark.ml.clustering import KMeans
from pyspark.ml.evaluation import ClusteringEvaluator

# Loads data.
dataset = spark.read.format("libsvm").load("data/mllib/sample_kmeans_data.txt")

# Trains a k-means model.
kmeans = KMeans().setK(2).setSeed(1)
model = kmeans.fit(dataset)

# Make predictions
predictions = model.transform(dataset)
```

K-Means Clustering

```
# Evaluate clustering by computing Silhouette score
evaluator = ClusteringEvaluator()

silhouette = evaluator.evaluate(predictions)
print("Silhouette with squared euclidean distance = " + str(silhouette))

# Shows the result.
centers = model.clusterCenters()
print("Cluster Centers: ")
for center in centers:
    print(center)
```

Frequent Pattern Mining

```
from pyspark.ml.fpm import FPGrowth

df = spark.createDataFrame([
    (0, [1, 2, 5]),
    (1, [1, 2, 3, 5]),
    (2, [1, 2])
], ["id", "items"])
```

Frequent Pattern Mining

```
fpGrowth = FPGrowth(itemsCol="items", minSupport=0.5, minConfidence=0.6)
model = fpGrowth.fit(df)

# Display frequent itemsets.
model.freqItemsets.show()

# Display generated association rules.
model.associationRules.show()

# transform examines the input items against all the association rules and summarize the
# consequents as prediction
model.transform(df).show()
```

Recommending Music and the Audioscrobbler Data Set

Exercises

- Frequent Pattern Mining of the words in the tweets

Contacts

For any problem, send a mail to

daniele.foroni@unitn.it