# CENG 790: Big Data Analytics, Fall 2020

# Report of Assignment 3: Random Forest Classifier

**Extract Features**

To build a classifier model, you first extract the features that most contribute to the classification. In the credit data set the data is labeled with two classes – 1 (**creditable**) and 0 (**not creditable**).

The features for each item consist of the fields shown below:

- Label → **creditable**: 0 or 1
- Features → {"**balance**", "**duration**", "**history**", "**purpose**", "**amount**", "**savings**", "**employment**", "**instPercent**", "**sexMarried**", "**guarantors**", "**residenceDuration**", "**assets**", "**age**", "**concCredit**", "**apartment**", "**credits**", "**occupation**", "**dependents**", "**hasPhone**", "**foreign**"}

In order for the features to be used by a machine learning algorithm, the features are transformed and put into Feature Vectors, which are vectors of numbers representing the value for each feature.

1. Use a **VectorAssembler** to transform and return a new dataframe with all of the feature columns in a vector column.

```scala
56    // EXTRACT FEATURES
57    // 1. Use a VectorAssembler to transform and return a new dataframe with all of the feature columns in a vector
58    // column.
59    val featureColumnNames = creditDF.columns.filter(colName => !colName.equals("creditability"))
60    val featureColumnName = "features"
61
62    val featuresAssembler = new VectorAssembler()
63      .setInputCols(featureColumnNames)
64      .setOutputCol(featureColumnName)
```

2. Use a **StringIndexer** to return a Dataframe with the creditability column added as a label.

```scala
66    // 2. Use a StringIndexer to return a Dataframe with the creditability column added as a label
67    val labelColumnName = "label"
68
69    val labelIndexer = new StringIndexer()
70      .setInputCol("creditability")
71      .setOutputCol(labelColumnName)
72      .fit(creditDF)
```

3. Use **randomSplit** function to split the data into two sets: 75% of the data is used to train (and tune) the model, 25% will be used for testing.

```scala
75    // 3. Use randomSplit function to split the data into two sets: 75% of the data is used to train (and tune) the
76    // model, 25% will be used for testing.
77    val Array(trainCreditDF, testCreditDF) = creditDF
78      .randomSplit(Array(0.75, 0.25), seed = 4321)
```

**Train the model and optimize hyperparameters**

The model is trained by making associations between the input features and the labeled output associated with those features. In order to find the best model, we search for the optimal combinations of the classifier parameters.

You will optimize the model using a pipeline. A pipeline provides a simple way to try out different combinations of parameters, using a process called grid search, where you set up the parameters to test, and MLLib will test all the combinations. **Pipelines** make it easy to tune an entire model building workflow at once, rather than tuning each element in the Pipeline separately.

4. Next, we train a RandomForestClassifier with the parameters:
   a. **maxDepth**: Maximum depth of a tree. Increasing the depth makes the model more powerful, but deep trees take longer to train.
   b. **maxBins**: Maximum number of bins used for discretizing continuous features and for choosing how to split on features at each node.
   c. impurity: Criterion used for information gain calculation
   d. **auto**: Automatically select the number of features to consider for splits at each tree node
   e. **seed**: Use a random seed number, allowing to repeat the results. Use the random seed 4321 in this assignment.

```
80      // TRAIN THE MODEL AND OPTIMIZE HYPERPARAMETERS
81      // 4. Next, we train a RandomForest Classifier with the parameters:
82      //    a. maxDepth: Maximum depth of a tree. Increasing the depth makes the model more powerful, but deep trees take
83      //    longer to train.
84      //    b. maxBins: Maximum number of bins used for discretizing continuous features and for choosing how to split on
85      //    features at each node.
86      //    c. impurity: Criterion used for information gain calculation
87      //    d. auto: Automatically select the number of features to consider for splits at each tree node
88      //    e. seed: Use a random seed number, allowing to repeat the results. Use the random seed 4321 in this
89      //    assignment.
90      // The model is trained by making associations between the input features and the labeled output associated with
91      // those features. In order to find the best model, we search for the optimal combinations of the classifier
92      // parameters.
93      // You will optimize the model using a pipeline. A pipeline provides a simple way to try out different combinations
94      // of parameters, using a process called grid search, where you set up the parameters to test, and MLLib will test
95      // all the combinations. Pipelines make it easy to tune an entire model building workflow at once, rather than
96      // tuning each element in the Pipeline separately.
97      val randomForestClassifier = new RandomForestClassifier()
98        .setFeaturesCol(featureColumnName)
99        .setLabelCol(labelColumnName)
100       .setSeed(4321)
```

Use the **ParamGridBuilder** utility to construct the parameter grid with the following values: maxBins [24, 28, 32], maxDepth, [3, 5, 7], impurity ["entropy", "gini"]

```
102      // Use the ParamGridBuilder utility to construct the parameter grid with the following values: maxBins [24, 28, 32],
103      // maxDepth, [3, 5, 7], impurity ["entropy", "gini"]
104      val paramGridBuilder = new ParamGridBuilder()
105        .addGrid(randomForestClassifier.featureSubsetStrategy, Array("auto"))
106        .addGrid(randomForestClassifier.impurity, Array("entropy", "gini"))
107        .addGrid(randomForestClassifier.maxBins, Array(24, 28, 32))
108        .addGrid(randomForestClassifier.maxDepth, Array(3, 5, 7))
109        .build()
110
111      // Additional ParamGridBuilders
112      val extraParamGridBuilder = new ParamGridBuilder()
113        .addGrid(randomForestClassifier.featureSubsetStrategy, Array("auto") ++ Array("all", "onethird", "sqrt", "log2"))
114        .addGrid(randomForestClassifier.impurity, Array("entropy", "gini"))
115        .addGrid(randomForestClassifier.maxBins, Array(24, 28, 32) ++ Array(12, 16, 20) ++ Array(36, 40))
116        .addGrid(randomForestClassifier.maxDepth, Array(3, 5, 7) ++ Array(2, 4, 6, 8, 9, 10))
117        .addGrid(randomForestClassifier.numTrees, Array(20, 24, 28, 32))
118        .addGrid(randomForestClassifier.subsamplingRate, Array(0.1, 0.25, 0.5, 0.75, 1.0))
119        .build()
```

Also, I tried other parameters such as **numTrees**, **subsamplingRate** and different values for **featureSubsetStrategy**, **maxBins**, **maxDepth** in the tuning phase. But none of them made a significant improvement to accuracy of the model over train and test data.

5. Next, you will create and set up a pipeline to make things easier. A Pipeline consists of a sequence of stages, each of which is either an Estimator or a Transformer.
Use **TrainValidationSplit** that creates a (training, test) dataset pair. It splits the dataset into these two parts using the trainRatio parameter. For example, with trainRatio=0.75, TrainValidationSplit will generate a training and test dataset pair where 75% of the data is used for training and 25% for validation. Use these values in your code. **Note that, this is different from the original random data split.** Here, we further divide the training dataset into training and validation set, for tuning purposes. The final model that has been tuned will be used to evaluate the result on the test set.
The TrainValidationSplit uses an Estimator, a set of ParamMaps, and an Evaluator. Estimatior should be your random forest model, the ParamMaps is the parameter grid that you built in the previous step. The Evaluator should be **new BinaryClassificationEvaluator().**

```scala
122    // 5. Next, you will create and set up a pipeline to make things easier. A Pipeline consists of a sequence of
123    // stages, each of which is either an Estimator or a Transformer.
124    // Use TrainValidationSplit that creates a (training, test) dataset pair. It splits the dataset into these two parts
125    // using the trainRatio parameter. For example, with trainRatio=0.75, TrainValidationSplit will generate a training
126    // and test dataset pair where 75% of the data is used for training and 25% for validation. Use these values in your
127    // code. Note that, this is different from the original random data split. Here, we further divide the training
128    // dataset into training and validation set, for tuning purposes. The final model that has been tuned will be used
129    // to evaluate the result on the test set.
130    val binaryClassificationEvaluator = new BinaryClassificationEvaluator()
131      .setLabelCol(labelColumnName)

133    // The TrainValidationSplit uses an Estimator, a set of ParamMaps, and an Evaluator. Estimator should be your random
134    // forest model, the ParamMaps is the parameter grid that you built in the previous step. The Evaluator should be
135    // new BinaryClassificationEvaluator().
136    val trainValidationSplit = new TrainValidationSplit()
137      .setEstimator(randomForestClassifier)
138      .setEstimatorParamMaps(paramGridBuilder)
139      .setEvaluator(binaryClassificationEvaluator)
140      .setTrainRatio(0.75)
141      .setSeed(4321)

143    val pipeline = new Pipeline()
144      .setStages(Array(featuresAssembler, labelIndexer, trainValidationSplit))

146    val model = pipeline.fit(trainCreditDF)
147    model.write.overwrite().save(MODEL_PATH)
```

What are the best combinations for the hyper parameters you optimized?

```scala
149    val bestModel = model.stages(2).asInstanceOf[TrainValidationSplitModel]
150      .bestModel.asInstanceOf[RandomForestClassificationModel]
151    val impurity = bestModel.getImpurity
152    val maxBins = bestModel.getMaxBins
153    val maxDepth = bestModel.getMaxDepth
154    println(s"""Model's Parameters => Impurity:\"$impurity\", MaxBins:$maxBins, MaxDepth:$maxDepth""")
```

```
Model's Parameters => Impurity:"entropy", MaxBins:32, MaxDepth:7
```

The combination for the best model:
- Impurity is **entropy**,
- MaxBins is **32**,
- MaxDepth is **7**.

6. Finally, evaluate the pipeline best-fitted model by comparing test predictions with test labels. You can use **transform** function to get the predictions for test dataset. You can use evaluator's **evaluate** function to get the metrics.

What are your accuracy values on train and test sets? Feel free to provide more stats and comment on your performance.

```
155        // 6. Finally, evaluate the pipeline best-fitted model by comparing test predictions with test labels. You can use
156        // transform function to get the predictions for test dataset. You can use evaluator's evaluate function to get the
157        // metrics.
158        val trainPredictions = model.transform(trainCreditDF)
159        val testPredictions = model.transform(testCreditDF)
160
161        val trainResult = binaryClassificationEvaluator.evaluate(trainPredictions)
162        val testResult = binaryClassificationEvaluator.evaluate(testPredictions)
163        println(s"Model's Accuracies on => Train Data:$trainResult, Test Data:$testResult")
```

```
Model's Accuracies on => Train Data:0.9632126156601339, Test Data:0.7933723196881084
```

The model's accuracy on train data is **0.9632126156601339**, and the accuracy on test data is **0.7933723196881084**. However, the accuracy on the test data didn't satisfy me as I had expected higher accuracy value. The reason for that might be the train and validation data size, which is 75% of all data (750 sample), isn't enough or the tuning parameters were not good enough. Therefore, I tried extending the parameter tuning phase in order to achieve better result:

- I extended the ParamGridBuilder with the following parameters:
  - **FeatureSubsetStrategy** -> auto, all, onethird, sqrt, log2
  - **Impurity** -> entropy, gini
  - **MaxBins** -> 24, 28, 32, 34, 36, 38, 40
  - **MaxDepth** -> 2, 3, 4, 5, 6, 7, 8, 9, 10
  - **NumTrees** -> 4, 8, 12, 16, 20, 24, 28, 32
  - **SubsamplingRate** -> 0.1, 0.25, 0.5, 0.75, 1.0
- Normally, there was 18 different combinations. With the extended parameters, there was more than 2500 different combinations and I was sure that the model would be better after the learning process. Yet, it didn't get any better! The best combination was:
  - FeatureSubsetStrategy is **log2**,
  - Impurity is **gini**,
  - MaxBins is **24**,
  - MaxDepth is **7**,
  - NumTrees is **32**,
  - SubsamplingRate is **0.75**.
- The new best model achieved approximately **0.964** accuracy on train data and approximately **0.787** accuracy, which is lower than the previous best model's accuracy, on test data.

I also tried some random parameters intuitively and achieved a better model in terms of accuracy on both train and test data. The model had the following parameters: FeatureSubsetStrategy **auto**, impurity **gini**, maxBins **12**, maxDepth **7**, numTrees **24**, subsamplingRate **1.0**. It achieved **0.9772576526382118** on train data and **0.7949155295646524** on test data.

In addition, I also tried changing randomSplit and trainValidationSplit ratios, but it didn't improve the model's accuracy neither. I even tried removing the seed with the thought of **seed 4321** was the unlucky one and running the training process multiple times in order to achieve a better result but that was not the case. I accepted the fact that it couldn't improve any further with the Random Forrest Classification algorithm.