Sayeed Ahmed – SEC01 (NUID 002191535)

Big Data System Engineering with Scala Spring 2023 Spark Assignment No. 2



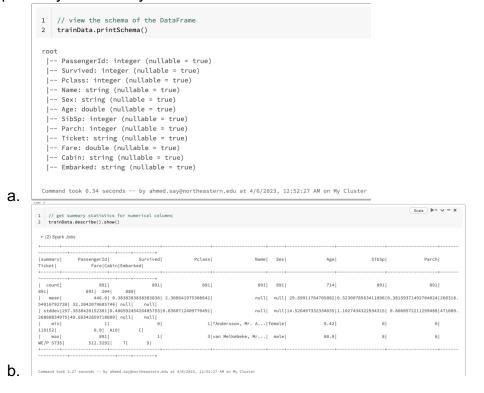
Github: https://github.com/sayeedahmed01/CSYE7200/tree/Spring2022/Spark-Assignments/Spark_Assignment_2

Implementation:

1) Loading datasets and imports:



2) Exploratory Data Analysis:



```
Cmd 6
                1 // check for missing values
                        \begin{tabular}{ll} \be
               3 null_counts.show()
                 ▶ ■ null_counts: org.apache.spark.sql.DataFrame = [PassengerId: long, Survived: long ... 10 more fields]
                |PassengerId|Survived|Pclass|Name|Sex|Age|SibSp|Parch|Ticket|Fare|Cabin|Embarked|\\
                                0| 0| 0| 0| 0|177| 0| 0| 0| 0| 687| 2|
                null_counts: org.apache.spark.sql.DataFrame = [PassengerId: bigint, Survived: bigint ... 10 more fields]
                Command took 2.02 seconds -- by ahmed.say@northeastern.edu at 4/6/2023, 12:52:27 AM on My Cluster
C.
                 1 // check the number of distinct values in each column
                 2 trainData.select(trainData.columns.map(c => countDistinct(col(c)).alias(c)): _*).show()
                   ▶ (3) Spark Jobs
                  |PassengerId\,|\,Survived\,|\,Pclass\,|\,Name\,|\,Sex\,|\,Age\,|\,SibSp\,|\,Parch\,|\,Ticket\,|\,Fare\,|\,Cabin\,|\,Embarked\,|
                                                     2 3 891 2 88 7 7 681 248 147
                  +-----
                 Command took 1.96 seconds -- by ahmed.say@northeastern.edu at 4/6/2023, 12:52:27 AM on My Cluster
d.
                Cmd 8
                           // Count the number of passengers in each category
                           trainData.groupBy("sex").count().show()
                   2
                   3 trainData.groupBy("pclass").count().show()
                    4 trainData.groupBy("embarked").count().show()
                      ▶ (6) Spark Jobs
                     | sex|count|
                     +----+
                     |female| 314|
                     | male| 577|
                     +----+
                     +----+
                     |pclass|count|
                              1 216
                                3 | 491 |
                             2| 184|
                     +----+
                     |embarked|count|
                     +----+
                                  Q| 77|
                               null| 2|
                                   Cl 1681
                   Command took 2.48 seconds -- by ahmed.say@northeastern.edu at 4/6/2023, 12:52:27 AM on I
e.
```

```
1 // Create a pivot table to show the survival rate by sex and class
    val pivotDF = trainData.filter(col("survived").isNotNull).groupBy("sex").pivot("pclass").agg(avg("survived"))
3 pivotDF.show()
▶ (7) Spark Jobs
 ▶ ■ pivotDF: org.apache.spark.sql.DataFrame = [sex: string, 1: double ... 2 more fields]
|female| 0.9680851063829787|0.9210526315789473|
   male|0.36885245901639346|0.1574074074074074|0.13544668587896252|
pivotDF: org.apache.spark.sql.DataFrame = [sex: string, 1: double ... 2 more fields]
Command took 2.53 seconds -- by ahmed.say@northeastern.edu at 4/6/2023, 12:52:27 AM on My Cluster
```

3) Feature Engineering:

f.

```
// Drop irrelevant columns
     val trainData2 = trainData.drop("PassengerId", "Ticket", "Name", "Cabin")
     val testData2 = testData.drop("PassengerId", "Ticket", "Name", "Cabin")
     // Create a new feature "FamilySize"
     val trainData3 = trainData2.withColumn("FamilySize", col("SibSp") + col("Parch") + 1)
     val testData3 = testData2.withColumn("FamilySize", col("SibSp") + col("Parch") + 1)
     val trainData4 = trainData3.withColumn("IsAlone", when(col("FamilySize") === 1, 1).otherwise(0))
11 val testData4 = testData3.withColumn("IsAlone", when(col("FamilySize") === 1, 1).otherwise(0))
 ▶ ■ trainData2: org.apache.spark.sql.DataFrame = [Survived: integer, Pclass: integer ... 6 more fields]
 ▶ ■ testData2: org.apache.spark.sql.DataFrame = [Pclass: integer, Sex: string ... 5 more fields]
 ▶ ■ trainData3: org.apache.spark.sql.DataFrame = [Survived: integer, Pclass: integer ... 7 more fields]
  ▶ ■ testData3: org.apache.spark.sql.DataFrame = [Pclass: integer, Sex: string ... 6 more fields]
 trainData4: org.apache.spark.sql.DataFrame = [Survived: integer, Pclass: integer ... 8 more fields]
  testData4: org.apache.spark.sgl.DataFrame = [Pclass: integer, Sex: string ... 7 more fields]
 trainData2: org.apache.spark.sql.DataFrame = [Survived: int, Pclass: int ... 6 more fields]
testData2: org.apache.spark.sql.DataFrame = [Pclass: int, Sex: string ... 5 more fields] trainData3: org.apache.spark.sql.DataFrame = [Survived: int, Pclass: int ... 7 more fields]
 testData3: org.apache.spark.sql.DataFrame = [Pclass: int, Sex: string ... 6 more fields]
 trainData4: org.apache.spark.sql.DataFrame = [Survived: int, Pclass: int ... 8 more fields]
 testData4: org.apache.spark.sql.DataFrame = [Pclass: int, Sex: string ... 7 more fields]
Command took 0.58 seconds -- by ahmed.say@northeastern.edu at 4/6/2023, 12:15:35 PM on My Cluster
Cmd 12
```

trainData4.show(5) 2 testData4.show(5)

▶ (2) Spark Jobs

+	+	+	+		+		+
Survived Pc	lass Sex Age Si	bSp Pa	arch	Fare Emb	arked Fami	lySize IsA7	one
+	+	+	+	+	+		+
0	3 male 22.0	1	0	7.25	S	2	0
1	1 female 38.0	1	0 71	1.2833	C	2	0
1	3 female 26.0	0	0	7.925	S	1	1
1	1 female 35.0	1	0	53.1	S	2	0
0	3 male 35.0	0	0	8.05	S	1	1
++							
only showing top 5 rows							

|Pclass| Sex| Age|SibSp|Parch| Fare|Embarked|FamilySize|IsAlone| 3| male|34.5| 0| 7.8292| 3|female|47.0| 7.0 0 | 1| 2| male|62.0| 0 | 0| 9.6875| QI 1 1 0| 8.6625| 3| male|27.0| 0 | 1| 1| SI 1 | 12,2875 | 3|female|22.0| 1| SI 31 0 l onlv showing top 5 rows

Command took 1.11 seconds -- by ahmed.say@northeastern.edu at 4/6/2023, 12:16:21 PM on My Cluster

b.

a.

4) Prediction:

```
// Convert the string columns to numerical indices
                            val sexIndexer = new StringIndexer().setInputCol("Sex").setOutputCol("SexIndex").setHandleInvalid("skip")
val embarkedIndexer = new StringIndexer().setInputCol("Embarked").setOutputCol("EmbarkedIndex").setHandleInvalid("skip")
                           // Convert the numerical indices to binary vectors
val sexEncoder = new OneHotEncoder().setInputCol("SexIndex").setOutputCol("SexVec")
                             val embarkedEncoder = new OneHotEncoder().setInputCol("EmbarkedIndex").setOutputCol("EmbarkedVec")
                  10 val assembler = new VectorAssembler()
                 11 .setInputCols(Array("Pclass", "SexVec", "Age", "SibSp", "Parch", "Fare", "EmbarkedVec", "FamilySize", "IsAlone"))
12 .setOutputCol("features")
13 .setHandleInvalid("skip")
                    sexIndexer: org.apache.spark.ml.feature.StringIndexer = strIdx_df497a020a70
                    embarkedIndexer: org.apache.spark.ml.feature.StringIndexer = strIdx_373ff8a64bde
                    sexEncoder: org.apache.spark.ml.feature.OneHotEncoder = oneHotEncoder_da9434b204b1
                   embarkedEncoder: org.apache.spark.ml.feature.OneHotEncoder = oneHotEncoder_e04dcbbb970a
                    assembler: org.apache.spark.ml.feature.VectorAssembler = VectorAssembler: uid=vecAssembler_2536348786b8, handleInvalid=skip, numInputCols=9
                   Command took 0.36 seconds -- by ahmed.say@northeastern.edu at 4/6/2023, 12:16:21 PM on My Cluste
a.
                               // Create a pipeline
                                 val pipeline = new Pipeline().setStages(Array(sexIndexer, embarkedIndexer, sexEncoder, embarkedEncoder, assembler))
                      4 // Fit the pipeline to the training data
                    5 val pipelineModel = pipeline.fit(trainData4)
                      pipeline: org.apache.spark.ml.Pipeline = pipeline_3d0ab8fe066b
                      pipelineModel: org.apache.spark.ml.PipelineModel = pipeline_3d0ab8fe066b
                      Command took 1.26 seconds -- by ahmed.say@northeastern.edu at 4/6/2023, 12:16:21 PM on My Cluster
b.
                                                                                                                                                                                                                                                                                                                                                 Scala >
                       // Transform the training and test data
val trainDataFinal = pipelineModel.transform(trainData4).select("features"
val testDataFinal = pipelineModel.transform(testData4).select("features")
                        // Split the data into training and testing sets
val Array(trainingData, testingData) = trainDataFinal.randomSplit(Array(0.7, 0.3))
                        // Create the Random Forest Classifier model
val rf = new RandomForestClassifier().setLabelCol("Survived").setFeaturesCol("features")
                        // Make predictions on the testing set
val predictions = model.transform(testingData)
                      // Evaluate the model using MulticlassClassificationEvaluator
val evaluator = new MulticlassClassificationEvaluator().setLabelCol("Survived").setPredictionCol("prediction").setMetricName("accuracy")
val accuracy = evaluator.evaluate(predictions)

■ 1 rest/bash-final: org apache.spark.sql.Dash-fame = [features: udf] Over |
■ 1 rest/bash-final: org apache.spark.sql.Dash-fame = [features: udf] Over |
■ 1 rest/mgData: org apache.spark.sql.Dash-fame = [features: udf] Over |
■ 1 rest/mgData: org apache.spark.sql.Dash-fame = [features: udf, Survived: integer] |
■ 1 predictions: org apache.spark.sql.Dash-fame = [features: udf, Survived: integer]. 3 more fields]
                 ** Impredictions: orgapache.spark.sql.Dataframe = [features: vector, Survived: int]
testBatafrant: org.apache.spark.sql.Dataframe = [features: vector, Survived: int]
testBatafrant: org.apache.spark.sql.Dataframe = [features: vector, Survived: int]
testBatafrant: org.apache.spark.sql.Dataster[org.apache.spark.sql.Row] = [features: vector, Survived: int]
testEngData: org.apache.spark.sql.Dataster[org.apache.spark.sql.Row] = [features: vector, Survived: int]
for org.apache.spark.al.classtfication.RandomforestClassifier = (features: vector, Survived: int]
for org.apache.spark.al.classtfication.RandomforestClassifier = (features: vector, Survived: int)
predictions: org.apache.spark.al.classtfication.RandomforestClassificationModel: vid=fc_ed3679acf7e0, numTrees=20, numFeatures=10
predictions: org.apache.spark.al.classtficationModel: vid=fc_ed3679acf7e0, numTrees=20, numFeatures=10
predictions: org.apache.spark.al.classtficationModel: vid=fc_ed3679acf7e0, numTrees=20, numFeatures=10
predictions: org.apache.spark.al.classtficationModel: vid=fc_ed3679acf7e0, numTrees=20, numFeatures=10
predictions: org.apache.spark.al.classtfi
C.
```

5) Result:

```
Cmd 17

1     println(s"Accuracy: ${accuracy*100}%")

Accuracy: 82.85714285714286%

Command took 0.24 seconds -- by ahmed.say@northeastern.edu at 4/6/2023, 12:16:21 PM on My Cluster

Cmd 18
```