

```

SPARK SESSION AVAILABLE AS spark .
Welcome to
 version 2.4.0-cdh6.3.4

Using Scala version 2.11.12 (Java HotSpot(TM) 64-Bit Server VM, Java 1.8.0_181)
Type in expressions to have them evaluated.
|Type :help for more information.

scala> import org.apache.spark.ml.feature.VectorAssembler
import org.apache.spark.ml.feature.VectorAssembler

scala> import org.apache.spark.sql.types._
import org.apache.spark.sql.types._

scala> var df = spark.read.option("header", "false").csv("Final_Project/all_input/finalcleandata.txt")
df: org.apache.spark.sql.DataFrame = [_c0: string, _c1: string ... 4 more fields]

scala> df.show()
+---+---+---+---+
| _c0| _c1| _c2| _c3|   _c4| _c5|
+---+---+---+---+
| 1| 3| 1| 1| 0.02394888206929671| 1|
| 1| 3| 1| 1| 0.019320723877998994| 1|
| 1| 7| 1| 1| 0.00227734768143619| 1|
| 1| 6| 1| 1| 0.013223399118016587| 1|
| 1| 4| 0| 1| 0.00730106617887...| 0|
| 1| 4| 1| 1| 0.010137870322812716| 1|
| 1| 7| 1| 1| 0.001836570710835...| 1|
| 1| 5| 0| 1| 0.0052893213647286635| 1|
| 1| 3| 0| 1| 0.007933985478889953| 0|
| 1| 3| 0| 1| 0.015427193971019352| 1|
| 1| 5| 1| 1| 0.006978968701175421| 0|
| 1| 4| 1| 1| 0.008154373956110229| 1|
| 1| 3| 1| 1| 0.01072557295128812| 1|
| 1| 3| 1| 1| 0.01386738407180949| 1|
| 1| 3| 1| 1| 0.0181467649534857752| 1|
| 1| 5| 0| 1| 0.0042893213647286635| 1|
| 1| 3| 1| 1| 0.019320723877998994| 1|
| 1| 5| 1| 1| 0.0201788149987558| 1|
| 1| 5| 1| 1| 0.03166247985480638| 0|
| 1| 3| 1| 1| 0.007052471529688846| 1|
+---+---+---+---+
only showing top 20 rows

scala> var df_final_1 = df.withColumn("sex", col("_c3").cast(IntegerType)).withColumn("age", col("_c1").cast(IntegerType)).withColumn("credit", col("_c5").cast(IntegerType)).withColumn("Employment", col("_c3").cast(IntegerType)).withColumn("Housing", col("_c0").cast(IntegerType)).withColumn("income", col("_c4").cast(DoubleType))
df_final_1: org.apache.spark.sql.DataFrame = [_c0: string, _c1: string ... 10 more fields]

scala> df_final_1.show()
+---+---+---+---+-----+
| _c0| _c1| _c2| _c3|   _c4| _c5|sex|age|credit|Employment|Housing| income|
+---+---+---+---+-----+

```

```

[1] | 5| 1| 1| 0| 0.03166247985480638| 8|
[1] | 3| 1| 1| 1| 0| 0.007852431529608846| 1|
+---+---+---+---+---+---+
only showing top 20 rows

[1] scalas> var df_final_1 = df.WithColumn("sex", col("c3").cast(IntegerType)).WithColumn("age", col("c1").cast(IntegerType)).WithColumn("income", col("c4").cast(DoubleType)).WithColumn("Employment", col("c3").cast(IntegerType)).WithColumn("Housing", col("c0").cast(IntegerType))
df_final_1: org.apache.spark.sql.DataFrame = [c0: string, c1: string ... 10 more fields]

scalas> df_final_1.show()
+---+---+---+---+---+---+---+---+---+---+
| _c0| _c1| _c2| _c3| _c4| _c5| sex| age| credit| Employment| Housing| income|
+---+---+---+---+---+---+---+---+---+---+
| 1| 1| 1| 1| 0.02394888206929671| 1| 1| 3| 1| 1| 1| 0.02394888206929671| |
| 1| 3| 1| 1| 0| 0.019326723877999964| 1| 1| 3| 1| 1| 1| 0.019326723877999964|
| 1| 7| 1| 1| 0| 0.00227734768143619| 1| 1| 7| 1| 1| 1| 0.00227734768143619|
| 1| 6| 1| 1| 0| 0.01322309118016587| 1| 1| 6| 1| 1| 1| 0.01322309118016587|
| 1| 4| 0| 1| 0| 0.0057306106617807...| 1| 1| 4| 0| 1| 1| 0.0057306106617807...|
| 1| 4| 1| 1| 0| 0.013787032817216| 1| 1| 4| 1| 1| 1| 0.013787032817216|
| 1| 7| 1| 1| 0| 0.001836570710853...| 1| 1| 7| 1| 1| 1| 0.001836570710853...|
| 5| 0| 1| 0| 0.0057895323647206635| 1| 1| 5| 1| 1| 1| 0.0057895323647206635|
| 3| 0| 1| 0| 0.007933985470809953| 1| 1| 3| 0| 1| 1| 0.007933985470809953|
| 3| 0| 1| 0| 0.015427193971019352| 1| 1| 3| 1| 1| 1| 0.015427193971019352|
| 5| 1| 1| 1| 0| 0.00697896870175421| 1| 1| 5| 0| 1| 1| 0.00697896870175421|
| 1| 4| 1| 1| 0| 0.008154373956110229| 1| 1| 4| 1| 1| 1| 0.008154373956110229|
| 3| 1| 1| 1| 0| 0.010725572951280812| 1| 1| 3| 1| 1| 1| 0.010725572951280812|
| 3| 1| 1| 1| 0| 0.011386738407189949| 1| 1| 3| 1| 1| 1| 0.011386738407189949|
| 3| 1| 1| 1| 0| 0.019467649534857752| 1| 1| 3| 1| 1| 1| 0.019467649534857752|
| 5| 0| 1| 0| 0.004113918392271827| 1| 1| 5| 1| 1| 1| 0.004113918392271827|
| 3| 1| 1| 1| 0| 0.019326723877999964| 1| 1| 3| 1| 1| 1| 0.019326723877999964|
| 3| 1| 1| 1| 0| 0.02012881499875858| 1| 1| 3| 1| 1| 1| 0.02012881499875858|
| 5| 1| 1| 1| 0| 0.03166247985480638| 1| 1| 5| 0| 1| 1| 0.03166247985480638|
+---+---+---+---+---+---+---+---+---+---+
only showing top 20 rows

[1] scalas> val cols = Array("age", "Housing", "Employment", "income", "sex")
cols: Array[String] = Array(age, Housing, Employment, income, sex)

[1] scalas> val assembler = new VectorAssembler().setInputCols(cols).setOutputCol("features")
assembler: org.apache.spark.ml.feature.VectorAssembler = vecAssembler_fib23953f5d3

[1] scalas> val featureDF = assembler.transform(df_final_1)
featureDF: org.apache.spark.sql.DataFrame = [c0: string, c1: string ... 11 more fields]

scalas> featureDF.show()
+---+---+---+---+---+---+---+---+---+---+---+
| _c0| _c1| _c2| _c3| _c4| _c5| sex| age| credit| Employment| Housing| features|
+---+---+---+---+---+---+---+---+---+---+---+
| 1| 3| 1| 1| 1| 0.02394888206929671| 1| 1| 3| 1| 1| 1| 0.02394888206929671| [3, 0, 1, 0, 1, 0, 0, 02...]
| 1| 3| 1| 1| 0| 0.019326723877999964| 1| 1| 3| 1| 1| 1| 0.019326723877999964| [3, 0, 1, 0, 1, 0, 0, 01...]
| 1| 7| 1| 1| 0| 0.00227734768143619| 1| 1| 7| 1| 1| 1| 0.00227734768143619| [7, 0, 1, 0, 1, 0, 0, 00...]
| 1| 6| 1| 1| 0| 0.01322309118016587| 1| 1| 6| 1| 1| 1| 0.01322309118016587| [6, 0, 1, 0, 1, 0, 0, 01...]
| 1| 4| 0| 1| 0| 0.0057306106617807...| 1| 1| 4| 0| 1| 1| 0.0057306106617807...| [4, 0, 1, 0, 1, 0, 0, 00...]
| 1| 4| 1| 1| 0| 0.013787032817216| 1| 1| 4| 1| 1| 1| 0.013787032817216| [4, 0, 1, 0, 1, 0, 0, 01...]
| 1| 7| 1| 1| 0| 0.001836570710853...| 1| 1| 7| 1| 1| 1| 0.001836570710853...| [7, 0, 1, 0, 1, 0, 0, 00...]

```

```

assembler: org.apache.spark.ml.feature.VectorAssembler = vecAssembler_f1b23933f5d3
scala> val featureDf = assembler.transform(df_final_1)
featureDf: org.apache.spark.sql.DataFrame = [_c0: string, _c1: string ... 11 more fields]

scala> featureDf.show()
+---+---+---+---+---+---+---+---+
| _c0|_c1|_c2|_c3| _c4|_c5|sex|age|credit|employment|housing| income| features|
+---+---+---+---+---+---+---+---+
| 1| 3| 1| 1| 0.02394888206929671| 1| 1| 3| 1| 1| 1| 1| 0.02394888206929671|[3,0,1,0,1,0,0,02...|
| 1| 3| 1| 1| 0.019320723877999094| 1| 1| 3| 1| 1| 1| 1| 0.019320723877999094|[3,0,1,0,1,0,0,01...|
| 1| 7| 1| 1| 0.00227734768143619| 1| 1| 7| 1| 1| 1| 1| 0.00227734768143619|[7,0,1,0,1,0,0,00...|
| 1| 6| 1| 1| 0.013223389118016587| 1| 1| 6| 1| 1| 1| 1| 0.013223389118016587|[6,0,1,0,1,0,0,01...|
| 1| 4| 0| 1| 0.005730180617807...| 0| 1| 4| 0| 1| 1| 1| 0.005730180617807...|[4,0,1,0,1,0,0,00...|
| 1| 4| 1| 1| 0.01813787879323812716| 1| 1| 4| 1| 1| 1| 1| 0.01813787879323812716|[4,0,1,0,1,0,0,01...|
| 1| 7| 1| 1| 0.00183602376018035...| 1| 1| 7| 1| 1| 1| 1| 0.00183602376018035...|[7,0,1,0,1,0,0,00...|
| 1| 5| 0| 1| 0.0052032376018036635| 1| 1| 5| 0| 1| 1| 1| 0.0052032376018036635|[5,0,1,0,1,0,0,00...|
| 1| 3| 0| 1| 0.007933985476809953| 1| 1| 3| 0| 1| 1| 1| 0.007933985476809953|[5,0,1,0,1,0,0,00...|
| 1| 3| 0| 1| 0.015427193971019352| 1| 1| 3| 0| 1| 1| 1| 0.015427193971019352|[3,0,1,0,1,0,0,01...|
| 1| 5| 1| 1| 0.006578968701175421| 0| 1| 5| 0| 1| 1| 1| 0.006578968701175421|[5,0,1,0,1,0,0,00...|
| 1| 4| 1| 1| 0.008154373956110229| 1| 1| 4| 1| 1| 1| 1| 0.008154373956110229|[5,0,1,0,1,0,0,00...|
| 1| 3| 1| 1| 0.010772557795128012| 1| 1| 3| 1| 1| 1| 1| 0.010772557795128012|[5,0,1,0,1,0,0,01...|
| 1| 3| 1| 1| 0.0113867384807180949| 1| 1| 3| 1| 1| 1| 1| 0.0113867384807180949|[3,0,1,0,1,0,0,01...|
| 1| 3| 1| 1| 0.019467649534857752| 1| 1| 3| 1| 1| 1| 1| 0.019467649534857752|[3,0,1,0,1,0,0,01...|
| 1| 5| 0| 1| 0.004113918392271827| 1| 1| 5| 0| 1| 1| 1| 0.004113918392271827|[5,0,1,0,1,0,0,00...|
| 1| 3| 1| 1| 0.019320723877999094| 1| 1| 3| 1| 1| 1| 1| 0.019320723877999094|[3,0,1,0,1,0,0,01...|
| 1| 3| 1| 1| 0.02012881499875858| 1| 1| 3| 1| 1| 1| 1| 0.02012881499875858|[3,0,1,0,1,0,0,02...|
| 1| 5| 1| 1| 0.03166247905480638| 0| 1| 5| 0| 1| 0| 1| 0.03166247905480638|[5,0,1,0,1,0,0,03...|
| 1| 3| 1| 1| 0.007085241529608846| 1| 1| 3| 1| 1| 1| 1| 0.007085241529608846|[3,0,1,0,1,0,0,00...|
only showing top 28 rows

[| scala> import org.apache.spark.ml.feature.StringIndexer
import org.apache.spark.ml.feature.StringIndexer
|]
scala> val indexer = new StringIndexer().setInputCol("credit").setOutputCol("label")
indexer: org.apache.spark.ml.feature.StringIndexer = strIdx_3a9604d967a0

scala> val labelDf = indexer.fit(featureDf).transform(featureDf)
labelDf: org.apache.spark.sql.DataFrame = [_c0: string, _c1: string ... 12 more fields]

scala> labelDf.show()
+---+---+---+---+---+---+---+---+
| _c0|_c1|_c2|_c3| _c4|_c5|sex|age|credit|employment|housing| income| features|label|
+---+---+---+---+---+---+---+---+
| 1| 3| 1| 1| 0.02394888206929671| 1| 1| 3| 1| 1| 1| 1| 0.02394888206929671|[3,0,1,0,1,0,0,02...| 0.0|
| 1| 3| 1| 1| 0.019320723877999094| 1| 1| 3| 1| 1| 1| 1| 0.019320723877999094|[3,0,1,0,1,0,0,01...| 0.0|
| 1| 7| 1| 1| 0.00227734768143619| 1| 1| 7| 1| 1| 1| 1| 0.00227734768143619|[7,0,1,0,1,0,0,00...| 0.0|
| 1| 6| 1| 1| 0.013223389118016587| 1| 1| 6| 1| 1| 1| 1| 0.013223389118016587|[6,0,1,0,1,0,0,01...| 0.0|
| 1| 4| 0| 1| 0.005730180617807...| 0| 1| 4| 0| 1| 1| 1| 0.005730180617807...|[4,0,1,0,1,0,0,00...| 1.0|
| 1| 4| 1| 1| 0.01813787879323812716| 1| 1| 4| 1| 1| 1| 1| 0.01813787879323812716|[4,0,1,0,1,0,0,00...| 1.0|
| 1| 7| 1| 1| 0.00183602376018035...| 1| 1| 7| 1| 1| 1| 1| 0.00183602376018035...|[7,0,1,0,1,0,0,00...| 0.0|
| 1| 5| 1| 1| 0.005789323647206635| 1| 1| 5| 1| 1| 1| 1| 0.005789323647206635|[5,0,1,0,1,0,0,00...| 0.0|
| 1| 3| 1| 1| 0.007933985476809953| 0| 1| 3| 1| 0| 1| 1| 0.007933985476809953|[3,0,1,0,1,0,0,00...| 1.0|
| 1| 3| 1| 1| 0.015427193971019352| 1| 1| 3| 1| 1| 1| 1| 0.015427193971019352|[3,0,1,0,1,0,0,01...| 0.0|
| 1| 5| 1| 1| 0.006578968701175421| 0| 1| 5| 1| 1| 1| 1| 0.006578968701175421|[5,0,1,0,1,0,0,00...| 0.0|
| 1| 4| 1| 1| 0.008154373956110229| 1| 1| 4| 1| 1| 1| 1| 0.008154373956110229|[5,0,1,0,1,0,0,00...| 0.0|
| 1| 3| 1| 1| 0.010772557795128012| 1| 1| 3| 1| 1| 1| 1| 0.010772557795128012|[5,0,1,0,1,0,0,01...| 0.0|
| 1| 3| 1| 1| 0.0113867384807180949| 1| 1| 3| 1| 1| 1| 1| 0.0113867384807180949|[3,0,1,0,1,0,0,01...| 0.0|
| 1| 3| 1| 1| 0.019467649534857752| 1| 1| 3| 1| 1| 1| 1| 0.019467649534857752|[3,0,1,0,1,0,0,01...| 0.0|
| 1| 5| 0| 1| 0.004113918392271827| 1| 1| 5| 0| 1| 1| 1| 0.004113918392271827|[5,0,1,0,1,0,0,00...| 0.0|

```

```
Q> final
+-----+-----+-----+-----+-----+-----+-----+
| c0 | c1 | c2 | c3 | c4 | c5|sex|age|credit|Employment|Housing| income| features|label|
+-----+-----+-----+-----+-----+-----+-----+
| 1 | 3 | 1 | 1 | 0.02394888286929671 | 1 | 1 | 3 | 1 | 1 | 1 | 1 | 0.02394888206929671|[3.0, 1.0, 1.0, 0.0, 0.2 ... | 0.0|
| 1 | 3 | 1 | 1 | 0.019326723877999994 | 1 | 1 | 3 | 1 | 1 | 1 | 1 | 0.01932672387799994|[3.0, 1.0, 1.0, 0.0, 0.01 ... | 0.0|
| 1 | 7 | 1 | 1 | 0.00227734768143619 | 1 | 1 | 7 | 1 | 1 | 1 | 1 | 0.00227734768143619|[7.0, 1.0, 1.0, 0.0, 0.00 ... | 0.0|
| 1 | 6 | 1 | 1 | 0.01322309118016587 | 1 | 1 | 6 | 1 | 1 | 1 | 1 | 0.01322309118016587|[6.0, 1.0, 1.0, 0.0, 0.01 ... | 0.0|
| 1 | 4 | 0 | 1 | 0.005730108617807 ... | 0 | 1 | 4 | 0 | 1 | 1 | 1 | 0.005730108617807 ...|[4.0, 1.0, 1.0, 0.0, 0.0 ... | 1.0|
| 1 | 4 | 1 | 1 | 0.018137870323812716 | 1 | 1 | 4 | 1 | 1 | 1 | 1 | 0.018137870323812716|[4.0, 1.0, 1.0, 0.0, 0.01 ... | 0.0|
| 1 | 7 | 1 | 1 | 0.001836570710835 ... | 1 | 1 | 7 | 1 | 1 | 1 | 1 | 0.001836570710835 ...|[7.0, 1.0, 1.0, 0.0, 0.00 ... | 0.0|
| 1 | 5 | 0 | 0 | 0.005289323647206635 | 1 | 1 | 5 | 1 | 1 | 1 | 1 | 0.005289323647206635|[5.0, 1.0, 1.0, 0.0, 0.00 ... | 0.0|
| 1 | 3 | 0 | 0 | 0.007933985470889953 | 0 | 1 | 3 | 0 | 0 | 1 | 1 | 0.007933985470889953|[3.0, 1.0, 1.0, 0.0, 0.00 ... | 0.0|
| 1 | 3 | 0 | 1 | 0.015427193971019352 | 1 | 1 | 3 | 1 | 1 | 1 | 1 | 0.015427193971019352|[3.0, 1.0, 1.0, 0.0, 0.01 ... | 0.0|
| 1 | 5 | 1 | 1 | 0.0069789681175421 | 0 | 1 | 5 | 1 | 0 | 1 | 1 | 0.0069789681175421|[5.0, 1.0, 1.0, 0.0, 0.00 ... | 1.0|
| 1 | 4 | 1 | 1 | 0.00815473956110229 | 1 | 1 | 4 | 1 | 1 | 1 | 1 | 0.00815473956110229|[4.0, 1.0, 1.0, 0.0, 0.00 ... | 0.0|
| 1 | 3 | 1 | 1 | 0.01386738487180949 | 1 | 1 | 3 | 1 | 1 | 1 | 1 | 0.01386738487180949|[3.0, 1.0, 1.0, 0.0, 0.01 ... | 0.0|
| 1 | 3 | 1 | 1 | 0.01946749534857752 | 1 | 1 | 3 | 1 | 1 | 1 | 1 | 0.01946749534857752|[3.0, 1.0, 1.0, 0.0, 0.01 ... | 0.0|
| 1 | 3 | 1 | 0 | 0.00411391832271827 | 1 | 1 | 3 | 1 | 1 | 1 | 1 | 0.00411391832271827|[5.0, 1.0, 1.0, 0.0, 0.00 ... | 0.0|
| 1 | 3 | 1 | 1 | 0.01932072387799994 | 1 | 1 | 3 | 1 | 1 | 1 | 1 | 0.01932072387799994|[3.0, 1.0, 1.0, 0.0, 0.01 ... | 0.0|
| 1 | 3 | 1 | 1 | 0.02612881499675585 | 1 | 1 | 3 | 1 | 1 | 1 | 1 | 0.02612881499675585|[3.0, 1.0, 1.0, 0.0, 0.02 ... | 0.0|
| 1 | 5 | 1 | 1 | 0.0316624779954860638 | 0 | 1 | 5 | 1 | 0 | 1 | 1 | 0.0316624779954860638|[5.0, 1.0, 1.0, 0.0, 0.03 ... | 1.0|
| 1 | 3 | 1 | 1 | 0.0070752431529608846 | 1 | 1 | 3 | 1 | 1 | 1 | 1 | 0.0070752431529608846|[3.0, 1.0, 1.0, 0.0, 0.00 ... | 0.0|
+-----+-----+-----+-----+-----+-----+-----+
only showing top 20 rows

scala> val indexer = new StringIndexer().setInputCol("c5").setOutputCol("label")
indexer: org.apache.spark.ml.feature.StringIndexer = stridx_9e89edab1515
scala> val labelDf = indexer.fit(featureDf).transform(featureDf)
labelDf: org.apache.spark.sql.DataFrame = [c0: string, c1: string ... 12 more fields]

scala> val labelDf = indexer.fit(featureDf).transform(featureDf)
labelDf: org.apache.spark.sql.DataFrame = [c0: string, c1: string ... 12 more fields]

scala> labelDf.show()
+-----+-----+-----+-----+-----+-----+-----+
| c0 | c1 | c2 | c3 | c4 | c5|sex|age|credit|Employment|Housing| income| features|label|
+-----+-----+-----+-----+-----+-----+-----+
| 1 | 3 | 1 | 1 | 0.02394888286929671 | 1 | 1 | 3 | 1 | 1 | 1 | 1 | 0.02394888206929671|[3.0, 1.0, 1.0, 0.0, 0.02 ... | 0.0|
| 1 | 3 | 1 | 1 | 0.01932672387799994 | 1 | 1 | 3 | 1 | 1 | 1 | 1 | 0.01932672387799994|[3.0, 1.0, 1.0, 0.0, 0.01 ... | 0.0|
| 1 | 7 | 1 | 1 | 0.00227734768143619 | 1 | 1 | 7 | 1 | 1 | 1 | 1 | 0.00227734768143619|[7.0, 1.0, 1.0, 0.0, 0.00 ... | 0.0|
| 1 | 6 | 1 | 1 | 0.01322309118016587 | 1 | 1 | 6 | 1 | 1 | 1 | 1 | 0.01322309118016587|[6.0, 1.0, 1.0, 0.0, 0.01 ... | 0.0|
| 1 | 4 | 0 | 1 | 0.005730108617807 ... | 0 | 1 | 4 | 1 | 0 | 1 | 1 | 0.005730108617807 ...|[4.0, 1.0, 1.0, 0.0, 0.0 ... | 1.0|
| 1 | 4 | 1 | 1 | 0.018137870323812716 | 1 | 1 | 4 | 1 | 1 | 1 | 1 | 0.018137870323812716|[4.0, 1.0, 1.0, 0.0, 0.01 ... | 0.0|
| 1 | 7 | 1 | 1 | 0.001836570710835 ... | 1 | 1 | 7 | 1 | 1 | 1 | 1 | 0.001836570710835 ...|[7.0, 1.0, 1.0, 0.0, 0.00 ... | 0.0|
| 1 | 5 | 0 | 0 | 0.005289323647206635 | 1 | 1 | 5 | 1 | 1 | 1 | 1 | 0.005289323647206635|[5.0, 1.0, 1.0, 0.0, 0.00 ... | 0.0|
| 1 | 3 | 0 | 0 | 0.007933985470889953 | 0 | 1 | 3 | 0 | 0 | 1 | 1 | 0.007933985470889953|[3.0, 1.0, 1.0, 0.0, 0.00 ... | 0.0|
| 1 | 3 | 0 | 1 | 0.015427193971019352 | 1 | 1 | 3 | 1 | 1 | 1 | 1 | 0.015427193971019352|[3.0, 1.0, 1.0, 0.0, 0.01 ... | 0.0|
| 1 | 5 | 1 | 1 | 0.0069789681175421 | 0 | 1 | 5 | 1 | 0 | 1 | 1 | 0.0069789681175421|[5.0, 1.0, 1.0, 0.0, 0.00 ... | 1.0|
| 1 | 4 | 1 | 1 | 0.00815473956110229 | 0 | 1 | 4 | 1 | 1 | 1 | 1 | 0.00815473956110229|[4.0, 1.0, 1.0, 0.0, 0.00 ... | 0.0|
| 1 | 3 | 1 | 1 | 0.01386738487180949 | 1 | 1 | 3 | 1 | 1 | 1 | 1 | 0.01386738487180949|[3.0, 1.0, 1.0, 0.0, 0.01 ... | 0.0|
| 1 | 3 | 1 | 0 | 0.01946749534857752 | 1 | 1 | 3 | 1 | 1 | 1 | 1 | 0.01946749534857752|[3.0, 1.0, 1.0, 0.0, 0.01 ... | 0.0|
| 1 | 5 | 0 | 0 | 0.00411391832271827 | 1 | 1 | 5 | 1 | 1 | 1 | 1 | 0.00411391832271827|[5.0, 1.0, 1.0, 0.0, 0.01 ... | 0.0|
| 1 | 3 | 1 | 1 | 0.01932072387799994 | 1 | 1 | 3 | 1 | 1 | 1 | 1 | 0.01932072387799994|[3.0, 1.0, 1.0, 0.0, 0.01 ... | 0.0|
| 1 | 3 | 1 | 1 | 0.02612881499675585 | 1 | 1 | 3 | 1 | 1 | 1 | 1 | 0.02612881499675585|[3.0, 1.0, 1.0, 0.0, 0.01 ... | 0.0|
| 1 | 5 | 1 | 1 | 0.0316624779954860638 | 0 | 1 | 5 | 1 | 0 | 1 | 1 | 0.0316624779954860638|[5.0, 1.0, 1.0, 0.0, 0.03 ... | 1.0|
| 1 | 3 | 1 | 1 | 0.0070752431529608846 | 1 | 1 | 3 | 1 | 1 | 1 | 1 | 0.0070752431529608846|[3.0, 1.0, 1.0, 0.0, 0.00 ... | 0.0|
+-----+-----+-----+-----+-----+-----+-----+
only showing top 20 rows
```











```

| zeruiji — zj637@hlog-2:~ — ssh zj637@peel.hpc.nyu.edu — 236x62
| Q final

scala> val randomForestClassifier = new RandomForestClassifier()
<console>:37: error: not found: type RandomForestClassifier
      val randomForestClassifier = new RandomForestClassifier()
                           ^
scala> import org.apache.spark.ml.classification.{RandomForestClassificationModel, RandomForestClassifier}
import org.apache.spark.ml.classification.{RandomForestClassificationModel, RandomForestClassifier}

scala> val randomForestClassifier = new RandomForestClassifier()
randomForestClassifier: org.apache.spark.ml.classification.RandomForestClassifier = rfc_1627125507a1

scala> val randomForestClassifier = new RandomForestClassifier().setImpurity("gini").setMaxDepth(10).setNumTrees(20)
<console>:38: error: value setNumTrees is not a member of Int
      val randomForestClassifier = new RandomForestClassifier().setImpurity("gini").setMaxDepth(10).setNumTrees(20)
                                         ^
scala> val randomForestClassifier = new RandomForestClassifier().setImpurity("gini").setMaxDepth(10).setNumTrees(20).setFeatureSubsetStrategy("auto")
<console>:39: error: ';' expected but ')' found.
      val randomForestClassifier = new RandomForestClassifier().setImpurity("gini").setMaxDepth(10).setNumTrees(20).setFeatureSubsetStrategy("auto")
                                         ^
scala> val randomForestClassifier = new RandomForestClassifier().setImpurity("gini").setMaxDepth(10).setNumTrees(20)...setFeatureSubsetStrategy("auto")
<console>:40: error: ';' expected but ')' found.
      val randomForestClassifier = new RandomForestClassifier().setImpurity("gini").setMaxDepth(10).setNumTrees(20)...setFeatureSubsetStrategy("auto")
                                         ^
scala> val randomForestClassifier = new RandomForestClassifier().setImpurity("gini").setMaxDepth(10).setNumTrees(20).setFeatureSubsetStrategy("auto").setSeed(seed)
<console>:41: error: ';' expected but ')' found.
      val randomForestClassifier = new RandomForestClassifier().setImpurity("gini").setMaxDepth(10).setNumTrees(20).setFeatureSubsetStrategy("auto").setSeed(seed)
                                         ^
scala> val randomForestClassifier = new RandomForestClassifier().setImpurity("gini").setMaxDepth(10).setNumTrees(20).setFeatureSubsetStrategy("auto").setSeed(seed)
randomForestClassifier: org.apache.spark.ml.classification.RandomForestClassifier = rfc_53970db6fe7b

scala> val randomForestModel = RandomForestClassifier.fit(trainingData)
randomForestModel: org.apache.spark.ml.classification.RandomForestClassificationModel = RandomForestClassificationModel (uid=rfc_53970db6fe7b) with 20 trees

scala> val predictionDF = randomForestModel.transform(testData)
predictionDF: org.apache.spark.sql.DataFrame = [c0: string, _c1: string ... 15 more fields]

scala> predictionDF.show(10)
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| c0 | c1 | c2 | c3 | c4 | c5 | sex | age | credit | Employment | Housing | income | features|label| rawPrediction| probability| prediction|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| 0 | 1 | 0 | 0 | 8.0 | 0 | 0 | 0 | 1 | 0 | 0 | 8.0 | (5, [0], [1, 0]) | 1.0 | [0, 4161521922103, ... | [0, 52088760961051, ... | 0.8 | |
| 0 | 1 | 0 | 0 | 8.0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 8.0 | (5, [0], [1, 0]) | 0.0 | [0, 4161521922103, ... | [0, 52088760961051, ... | 0.8 |
| 0 | 1 | 0 | 0 | 11.0 | 0.01175405254934... | 0 | 1 | 1 | 0 | 1 | 0 | 0.01175405254934... | [1, 0, 0, 0, 1, 0, 0, 0, 0] | 1.0 | [4, 52985955171071, ... | [0, 22645297758553, ... | 1.0 |
| 0 | 1 | 0 | 0 | 11.0 | 0.01322339091801, ... | 0 | 1 | 1 | 0 | 1 | 0 | 0.01322339091801, ... | [1, 0, 0, 0, 1, 0, 0, 0, 0] | 1.0 | [4, 52985955171071, ... | [0, 22645297758553, ... | 1.0 |
| 0 | 1 | 0 | 0 | 11.0 | 0.01689645053968786 | 0 | 1 | 1 | 0 | 1 | 0 | 0.01689645053968786 | [1, 0, 0, 0, 1, 0, 0, 0, 0] | 1.0 | [4, 52985955171071, ... | [0, 22645297758553, ... | 1.0 |
| 0 | 1 | 0 | 0 | 11.0 | 0.001836570718835... | 1 | 1 | 1 | 1 | 1 | 1 | 0.001836570718835... | [1, 0, 0, 0, 1, 0, 0, 0, 0] | 0.0 | [0, 4, 52985955171071, ... | [0, 22645297758553, ... | 1.0 |
| 0 | 1 | 0 | 0 | 11.0 | 0.002203884853802... | 1 | 1 | 1 | 1 | 1 | 1 | 0.002203884853802... | [1, 0, 0, 0, 1, 0, 0, 0, 0] | 0.0 | [4, 32311684336508, ... | [0, 21615584216825, ... | 1.0 |
| 0 | 1 | 0 | 0 | 11.0 | 0.00227734768143619 | 0 | 1 | 1 | 0 | 1 | 1 | 0.00227734768143619 | [1, 0, 0, 0, 1, 0, 0, 0, 0] | 0.0 | [4, 32311684336508, ... | [0, 21615584216825, ... | 1.0 |
| 0 | 1 | 0 | 0 | 11.0 | 0.00227734768143619 | 1 | 1 | 1 | 1 | 1 | 1 | 0.00227734768143619 | [1, 0, 0, 0, 1, 0, 0, 0, 0] | 0.0 | [4, 32311684336508, ... | [0, 21615584216825, ... | 1.0 |
| 0 | 1 | 0 | 0 | 11.0 | 0.00227734768143619 | 1 | 1 | 1 | 1 | 1 | 1 | 0.00227734768143619 | [1, 0, 0, 0, 1, 0, 0, 0, 0] | 0.0 | [4, 32311684336508, ... | [0, 21615584216825, ... | 1.0 |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
only showing top 10 rows

scala> val randomForestClassifier = new RandomForestClassifier().setImpurity("gini").setMaxDepth(3).setNumTrees(20).setFeatureSubsetStrategy("auto").setSeed(seed)
randomForestClassifier: org.apache.spark.ml.classification.RandomForestClassifier = rfc_72431f480e45

scala> val randomForestModel = RandomForestClassifier.fit(trainingData)
randomForestModel: org.apache.spark.ml.classification.RandomForestClassificationModel = RandomForestClassificationModel (uid=rfc_72431f480e45) with 20 trees

scala> val predictionDF = randomForestModel.transform(testData)
predictionDF: org.apache.spark.sql.DataFrame = [c0: string, _c1: string ... 15 more fields]

```







