CISC 5950 - Big Data Programming

Project 2 - Part 2 Report

Team Members:
Revathi Bhuvaneswari
Tianying Luo
Youfei Zhang
Yue Zeng

Spring 2019

P2: Heart Disease Prediction using Logistic Regression

Bash script: Import the data file from HDFS.

```
#!/bin/bash
#../../start.sh
source ../../env.sh
/usr/local/hadoop/bin/hdfs dfs -rm -r /part2/input/
/usr/local/hadoop/bin/hdfs dfs -rm -r /part2/output/
/usr/local/hadoop/bin/hdfs dfs -mkdir -p /part2/input/
/usr/local/hadoop/bin/hdfs dfs -copyFromLocal ../../data/framingham.csv /part2/input/
/usr/local/spark/bin/spark-submit --master=spark://$SPARK_MASTER:7077 ./framingham.py hdfs://$SPARK_MASTER:9000/part2/input/
#../../stop.sh
```

Python Script:

Stage 1: The train / test data are read from the CSV files and loaded into dataframes.

```
reload(sys)
sys.setdefaultencoding('utf8')
if __name__ == "__main__":
    if len(sys.argv) != 2:
        print("Usage: Q2.py <directory_of_files>", file=sys.stderr)
        sys.exit(-1)

spark = SparkSession.builder.appName('lr-predic').getOrCreate()
    df = spark.read.csv(sys.argv[1], header = True, inferSchema = True)
    df.show()
    cols=df.columns
    df.printSchema()
```

Stage 2: The dataset is cleaned, the vectors of 'features' & 'label' are added to the dataframe.

```
#Filter out all the 'NA' values

df = df.filter((df.cigsPerDay != 'NA')& (df.BPMeds != 'NA') & (df.totChol != 'NA') & (df.BMI !
```

The cleaning process is done by first filtering out all the 'NA' values from the following columns: 'cigsPerDay', 'BPMeds', 'totChol', 'BMI', 'heartRate', 'glucose'. 'NA' values represent missing values in the dataset and are identified as string value. Then all the columns of string type are altered by converting them into float type. Additionally, the unwanted column of 'education' is dropped.

The following function is used to one-hot-encode all categorical columns, and get vector representation of all features. String Indexer and One Hot Encoder Estimator are used for this.

```
sparse_to_array(v):
       v = DenseVector(v)
       new_array = list([float(x) for x in v])
return new_array
                   one-hot-encoded stages
def ohe(df):
       # One-hot-encode category columns
categoricalColumns = ['education','currentSmoker', 'BPMeds', 'prevalentStroke', 'prevalentHyp', 'diabetes']
       stages = []
for categoricalCol in categoricalColumns:
               stringIndexer = StringIndexer(inputCol = categoricalCol, outputCol = categoricalCol + 'Index')
encoder = OneHotEncoderEstimator(inputCols=[stringIndexer.getOutputCol()], outputCols=[categoricalCol + "classVec"])
               stages += [stringIndexer, encoder]
       # vectorize numeric columns
      # Vectorize numeric columns
numericCols = ['male', 'age', 'cigsPerDay', 'totChol', 'sysBP', 'diaBP', 'BM.
assemblerInputs = [c + "classVec" for c in categoricalColumns] + numericCols
assembler = VectorAssembler(inputCols=assemblerInputs, outputCol="features")
stages += [assembler]
                                                                                                                                                      'BMI', 'heartRate', 'glucose']
       # perform final pipeline to get vectorized feature column
pipeline = Pipeline().setStages(stages)
       pipelineModel = pipeline.fit(df)
df = pipelineModel.transform(df)
selectedCols = ['TenYearCHD', 'fd')
                                                                     features'] + cols
       df = df.select(selectedCols)
      # convert sparse vector initial_feature column to dense vector
sparse_to_array_udf = udf(sparse_to_array, T.ArrayType(T.FloatType()))
df = df.withColumn('dense_vector_features', sparse_to_array_udf('featured')
ud_f = udf(lambda r : Vectors.dense(r), VectorUDT())
df = df.withColumn('features', ud_f('dense_vector_features'))
return df
        return df
```

Stage 3: In this stage, a new column of 'Scaled_features' is created by using StandardScaler to standardize the newly feature values according to their different calibration. Also, since the classification values are imbalanced. Balancing ratio is calculated, based on which class weights will be calculated and applied to different classes. In addition, ChiSqSelector method is used to select the feature with p-value below 0.05, which proves a strong predictive power. At last, use the class column as labelCol, 'Aspect' (generated by chi-square feature selection) as featuresCol, and classweights as weightCol, and the iteration time is set to be 10. After that, split the dataset into train and test by 80/20. Train and fit the model by using the training dataset.

```
standardscaler=StandardScaler().setInputCol("features").setOutputCol("Scaled_features")
df=standardscaler.fit(df).transform(df)
df.select("features","Scaled_features").show(5)
train, test = df.randomSplit([0.8, 0.2], seed=123)
total=float(train.select("TenYearCHD").count())
                                                    ').where('TenYearCHD == 1').count()
numPositives=train.select("TenYearC
numPositives=(float(numPositives)/float(total))*100
numNegatives=float(total-numPositives)
print('\n\nThe number of Class 1 are {}'.format(numPositives))
print('\n\nPercentage of Class 1 are {}'.format(per_ones))
BalancingRatio= numNegatives/total
   rint('\n\nBalancingRatio = {}'.format(BalancingRatio))
#Creat a new column named "classWeights" in the "train" dataset
train=train.withColumn("classWeights", when(train.TenYearCHD== , BalancingRatio).otherwise(1-BalancingRatio))
train.select("classWeights", "TenYearCHD").show(10)
#Feature selection
css = ChiSqSelector(featuresCol='features',outputCol='Aspect',labelCol='TenYearCHD',fpr=8.8
train=css.fit(train).transform(train)
test=css.fit(test).transform(test)
test.select("Aspect").show(",truncate=False)
lr = LogisticRegression(labelCol="TenYearCHD", featuresCol="Aspect",weightCol="classWeights", maxIter=10)
model=lr.fit(train)
```

Stage 4: In this last stage, the trained Logistic Regression model is used to make predictions by applying the transform() function on both the train as well as the test data frames. Then extract the counts of TP, TN, FP, and FN by filtering the corresponding pairs and counting them. The Area Under the ROC is calculated using the Binary Classification Evaluator by passing in the actual predictions data frame, which has both the labels and predictions.

```
predict_train=model.transform(train)
predict_test=model.transform(test)
predict_train.select("TenYearCHD", "prediction").show(10)
predict_test.select("TenYearCHD", "prediction").show(10)

#Evaluation of the train model
TP_train = predict_train.filter((predict_train.TenYearCHD == 1) & (predict_train.prediction == 1.0)).count()
TN_train = predict_train.filter((predict_train.TenYearCHD == 0) & (predict_train.prediction == 0.0)).count()
FP_train = predict_train.filter((predict_train.TenYearCHD == 0) & (predict_train.prediction == 1.0)).count()
FN_train = predict_train.filter((predict_train.TenYearCHD == 0) & (predict_train.prediction == 1.0)).count()
evaluation(TP_train, TN_train, FP_train, FN_train)

#Evaluation of the test model
TP_test = predict_test.filter((predict_test.TenYearCHD == 0) & (predict_test.prediction == 1.0)).count()
TN_test = predict_test.filter((predict_test.TenYearCHD == 0) & (predict_test.prediction == 1.0)).count()
FP_test = predict_test.filter((predict_test.TenYearCHD == 0) & (predict_test.prediction == 1.0)).count()
FN_test = predict_test.filter((predict_test.TenYearCHD == 0) & (predict_test.prediction == 1.0)).count()
FN_test = predict_test.filter((predict_test.TenYearCHD == 0) & (predict_test.prediction == 1.0)).count()
FN_test = predict_test.filter((predict_test.TenYearCHD == 0) & (predict_test.prediction == 0.0)).count()
FN_test = predict_test.filter((predict_test.TenYearCHD == 0) & (predict_test.prediction == 0.0)).count()
FN_test = predict_test.filter((predict_test.TenYearCHD == 0) & (predict_test.prediction == 0.0)).count()
FN_test = predict_test.filter((predict_test.TenYearCHD == 0) & (predict_test.prediction == 0.0)).count()
FN_test = predict_test.filter((predict_test.TenYearCHD == 0) & (predict_test.prediction == 0.0)).count()
FN_test = predict_test.filter((predict_test.TenYearCHD == 0) & (predict_test.prediction == 0.0)).count()
FN_test = predict_test.filter((predict_test.TenYearCHD == 0) & (predict_test.prediction == 0.0)).count()
FN_test = predict_test.filter((predict_test.TenYe
```

The following function is used to calculate and output the accuracy, misclassification rate, sensitivity, specificity, positive/negative predictive value and positive/negative likelihood ratio.

```
def evaluation(TP, TN, FP, FN):
    print('\n\n', 'The confusion matrix is: \n\n', '[[',TN,FP,']\n\n','[',FN,TP,']]', '\n\n')
    #Model Evaluation - Statistics
    sensitivity=TP/float(TP+TN)
    specificity=TN/float(TN+FP)

    print('\n\nThe acuuracy of the model = (TP+TN)/(TP+TN+FP+FN) = ',(TP+TN)/float(TP+TN+FP+FN),'\n\n',
    'The Missclassification = 1-Accuracy = ',1-((TP+TN)/float(TP+TN+FP+FN)),'\n\n',
    'Sensitivity or True Positive Rate = TP/(TP+FN) = ',TP/float(TP+FN),'\n\n',
    'Specificity or True Negative Rate = TN/(TN+FP) = ',TN/float(TN+FP),'\n\n',
    'Positive Predictive value = TP/(TP+FP) = ',TP/float(TN+FP),'\n\n',
    'Negative predictive Value = TN/(TN+FN) = ',TN/float(TN+FN),'\n\n',
    'Positive Likelihood Ratio = Sensitivity/(1-Specificity) = ',sensitivity/(1-specificity),'\n\n',
    'Negative likelihood Ratio = (1-Sensitivity)/Specificity = ',(1-sensitivity)/specificity,'\n\n')
```

Output:

| + | + | + | | | + | t | | | + | · | · | · | + | + | + | |
|------|----|---------|----------|---------------|------------|--------|-----------------|--------------|----------|---------|-------|-------|-------|-----------|---------|------------|
| male | ag | ge ec | lucation | currentSmoker | cigsPerDay | BPMeds | prevalentStroke | prevalentHyp | diabetes | totChol | sysBP | diaBP | BMI | heartRate | glucose | TenYearCHD |
| 1 | 3 | 39 | 4 | 0 | 0 | . 0 | 0 | 0 | | 195 | 106.0 | 70.0 | 26.97 | 80 | 77 | . 01 |
| 0 | | 46 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | | | | 28.73 | | | 0 |
| 1 | | 48 | 1 | 1 | 20 | | 0 | 0 | 0 | | | | 25.34 | | | 0 |
| 0 | | 51 | 3 | 1 | 30 | | 0 | 1 | 0 | | | | 28.58 | | | 1 |
| 0 | | 46 | 3 | 1 | 23 | | 0 | 0 | 0 | | | | 23.1 | | | 0 |
| 0 | | 43 | 2 | 0 | 0 | | | 1 | 0 | | | | 30.3 | | | 0 |
| 0 | | 53 | 1 | 0 | 0 | | | 0 | | | | | 33.11 | | | 1 |
| 0 | | 45 | 2 | 1 | 20 | | | 0 | | | | | 21.68 | | | |
| 1 | | 52 | 1 | 0 | 0 | | | 1 | 0 | | | | 26.36 | | | |
| 1 | | 43 | 1 | 1 |] 30 | | | 1 | 0 | | | | 23.61 | | | |
| 0 | | 50 | 1 | 0 | 0 | | | 0 | | | | | 22.91 | | | 0 |
| 0 | | 43 | 2 | 0 | 0 | | 0 | 0 | 0 | | | | 27.64 | | | 0 |
| ! 1 | | 46 | 1 | 1 | 15 | | 0 | 1 | 0 | | | | 26.31 | | | 0 |
| 0 | | 41 | 3 | 0 | 0 | | 0 | 1 | 0 | | | | 31.31 | | | 0 0 |
| 0 | | 39 | 2 | 1 | 9 | | 0 | 0 | 0 | | | | 22.35 | | | 9 |
| 0 | | 38 | 2 | 1 | 20 | | 0 | 1 | 0 | | | | 21.35 | | | 1 1 |
| ! 1 | | 48 | 3 | 1 | 10 | | 0 | 1 | | | | | 22.37 | | | 0 |
| . 0 | | 46 | 2 | 1 | 20 | . 0 | 0 | 0 | . 0 | | | | 23.38 | | | 1 |
| . 0 | | 38 | 2 | 1 | 5 | 0 | 0 | 0 | . 0 | | | | 23.24 | | 78 | 0 |
| T 1 | Τ4 | 41 | 2 | 0 | . 0 | . 0 | 0 | 0 | . 0 | 195 | 139.0 | 88.0 | 26.88 | 85 | 65 | ا ا |
| · | + | _ | | | | · | | | | | | | | | | |

Figure 1: Initial DataFrame (first 20 rows)

```
root
 |-- male: integer (nullable = true)
  -- age: integer (nullable = true)
  -- education: string (nullable = true)
   - currentSmoker: integer (nullable = true)
  -- cigsPerDay: string (nullable = true)
   BPMeds: string (nullable = true)
    prevalentStroke: integer (nullable = true)
   - prevalentHyp: integer (nullable = true)
   - diabetes: integer (nullable = true)
  -- totChol: string (nullable = true)
  -- sysBP: double (nullable = true)
  — diaBP: double (nullable = true)
   - BMI: string (nullable = true)
   - heartRate: string (nullable = true)
    glucose: string (nullable = true)
    TenYearCHD: integer (nullable = true)
```

Figure 2: Schema of the Initial DataFrame

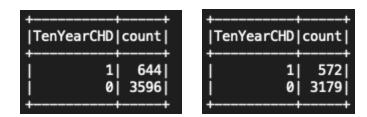


Figure 3: Classification Counts Before & After Data Cleaning

| Summary male | age | currentSmoker | cigsPerDay | BPMeds | prevalentStroke | prevalentHyp |
|---|-----|---------------|------------|----------------------------|-----------------|--------------|
| count 3751 mean 0.4452146094374833 stddev 0.49705575870704916 min 0 max 1 | | | | 0.17168601101128225 0.0 | | |

| diabetes | totChol | sysBP | diaBP | BMI | heartRate | glucose |
|----------|--------------------|--|---|--|---|---|
| | 236.92801919488136 | 132.3684350839776 22.04652246327564 83.5 | 82.93854972007465 11.932779343310061 48.0 | 25.80828845268947 4.06559867696019 15.54 | 75.70407891229006 11.956382146147071 44.0 | 81.88003199146894 23.882233365111478 40.0 |

Figure 4: Generative Statistics Analysis of Each Attribute

```
root
 — TenYearCHD: integer (nullable = true)
   features: vector (nullable = true)
   - male: integer (nullable = true)
   - age: integer (nullable = true)
    currentSmoker: integer (nullable = true)
  -- cigsPerDay: float (nullable = true)
   - BPMeds: float (nullable = true)
   - prevalentStroke: integer (nullable = true)
    prevalentHyp: integer (nullable = true)
    diabetes: integer (nullable = true)
  -- totChol: float (nullable = true)
    sysBP: double (nullable = true)
   - diaBP: double (nullable = true)
   - BMI: float (nullable = true)
   - heartRate: float (nullable = true)
    glucose: float (nullable = true)
    dense_vector_features: array (nullable = true)
      |-- element: float (containsNull = true)
```

Figure 5: Schema of the DataFrame After Data Cleaning

Figure 6: Scaled features of Standardization

```
The number of Class 1 are 460

Percentage of Class 1 are 15.343562374916612

BalancingRatio = 0.8465643762508339
```

Figure 7: Percentage of Class 1 & Balancing Ratio

| + classWeights | TenYearCHD |
|---------------------|-----------------|
| 0.15343562374916608 | i 0i |
| 0.15343562374916608 | |
| 0.15343562374916608 | i øj |
| 0.15343562374916608 | j 0j |
| 0.15343562374916608 | j 0j |
| 0.15343562374916608 | j 0j |
| 0.15343562374916608 | [0] |
| 0.15343562374916608 | 0 |
| 0.15343562374916608 | |
| 0.15343562374916608 | 0 |
| only showing top 10 | ++ rows |

Figure 8: Classweights & Aspect

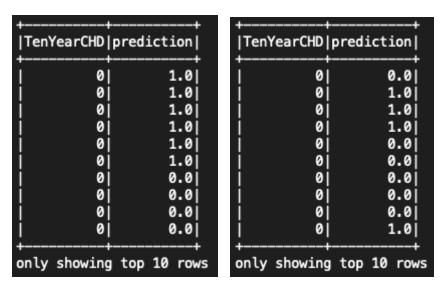


Figure 9: Train & Test Prediction Outputs

| TenYearCHD | rawPrediction | prediction | probability | | | | | |
|-------------------------|---------------|---------------------|---|--|--|--|--|--|
| 0 0 0 | | 1.0 1.0 0.0 | [0.46089184842530045,0.5391081515746997] [0.14262674731316835,0.8573732526868316] [0.16575169569904596,0.8342483043009541] [0.7269861164242096,0.2730138835757904] [0.6527325313989591,0.34726746860104085] | | | | | |
| only showing top 5 rows | | | | | | | | |

Figure 10: Probability of Predictions

The area under ROC for train set is 0.7269039999999907

```
Evaluation Results for Training Dataset

The confusion matrix is:

[[ 1677 823 ]

[ 154 296 ]]

The acuuracy of the model = (TP+TN)/(TP+TN+FP+FN) = 0.6688135593220339

The Missclassification = 1-Accuracy = 0.3311864406779661

Sensitivity or True Positive Rate = TP/(TP+FN) = 0.657777777777778

Specificity or True Negative Rate = TN/(TN+FP) = 0.6708

Positive Predictive value = TP/(TP+FP) = 0.2645218945487042

Negative predictive Value = TN/(TN+FN) = 0.9158929546695794

Positive Likelihood Ratio = Sensitivity/(1-Specificity) = 1.9981098960442822

Negative likelihood Ratio = (1-Sensitivity)/Specificity = 0.5101702776121381
```

Figure 11: Evaluation Results for Training Dataset

The area under ROC for test set is 0.774487554021107

```
Evaluation Results for Test Dataset

The confusion matrix is:

[[ 482 197 ]
  [ 33 89 ]]

The accuracy of the model = (TP+TN)/(TP+TN+FP+FN) = 0.7128589263420724

The Missclassification = 1-Accuracy = 0.2871410736579276

Sensitivity or True Positive Rate = TP/(TP+FN) = 0.7295081967213115

Specificity or True Negative Rate = TN/(TN+FP) = 0.7098674521354934

Positive Predictive value = TP/(TP+FP) = 0.3111888111888112

Negative predictive Value = TN/(TN+FN) = 0.9359223300970874

Positive Likelihood Ratio = Sensitivity/(1-Specificity) = 2.514396271948074

Negative likelihood Ratio = (1-Sensitivity)/Specificity = 0.38104550710835994
```

Figure 12: Evaluation Results for Testing Dataset