

Phase 6: Optimization

Team Name: Deep Diver's

Team Members:

Vaishnavi Hemant Salaskar vsala2@unh.newhaven.edu

Vedant Chidgopkar vchid2@unh.newhaven.edu

Riddhi Joshi rjosh5@unh.newhaven.edu

Research Question and Selected Data Set

Today mental health is becoming a more common problem. However, evaluation of mental well-being is extremely important to understanding and providing therapeutic solutions. Diagnostics are complicated tasks and misdiagnosis can result in serious problems if a mental disorder is not properly detected. Can we recognize mental health issues accurately by using data mining techniques?

The data has been collected from Kaggle by Open Sourcing Mental Illness, LTD. Survey data about mental health attitudes are included in this dataset. Which then has been analyzed and pre-processed. The data contains different labels such as age, gender, country, self-employee, family history, work interference, seek help, etc. For better prediction, we have label encoded the data.

List of Data Mining Techniques Used

- Naïve Bayes
- Decision Tree
- Random Forest

List of Techniques Used for Optimization

- Multiclass Classifier
- CV Parameter selection
- Random Subspace

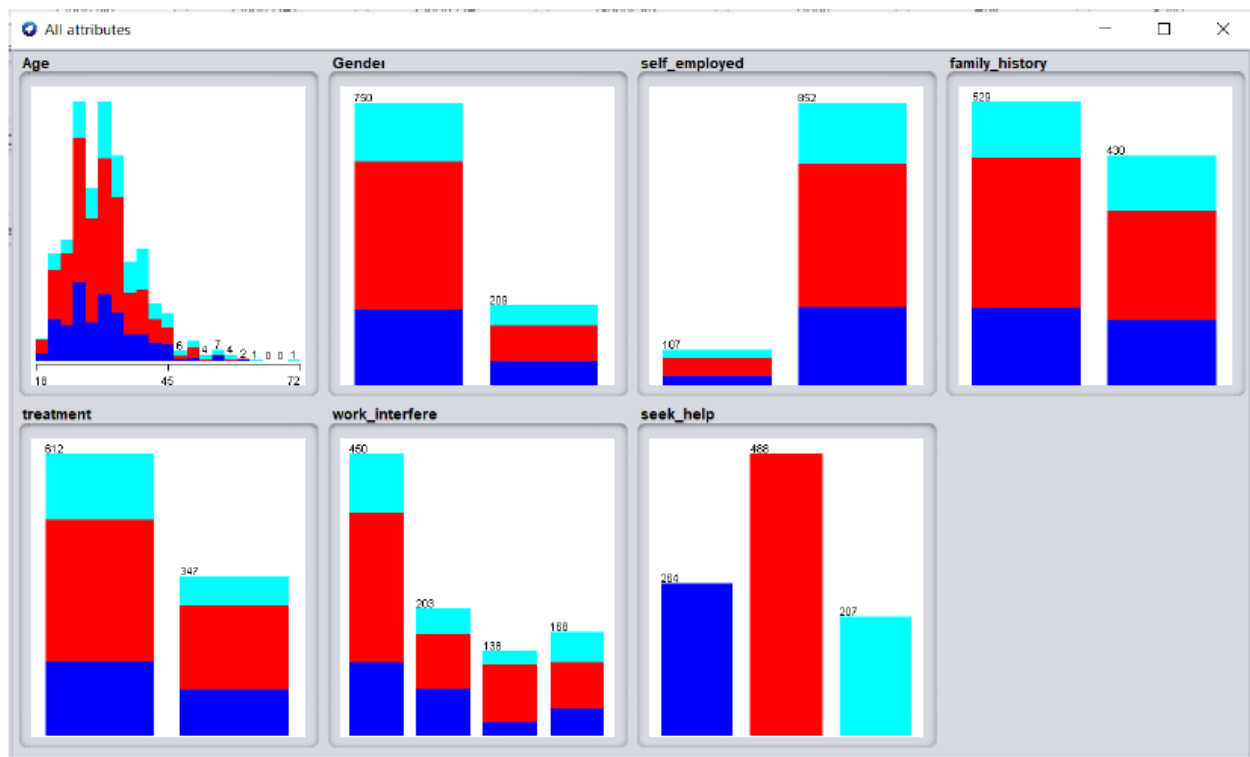
Visualization Techniques

Histogram

In Weka, we can visualize data in multiple ways. In the histogram, the distribution is done mainly from attributes. A single selected attribute is distributed at a time. By default, that will be the class attribute.

Total Instances: 959

Attributes: 7



In our dataset, we have noticed that in the Age attribute the Distinct Values are 44, Unique Values are 7 having the Standard Deviation as 7.403, and Mean is 32.356. Furthermore, in the treatment attribute, the Distinct Values are 2 with no Unique Values. The total count of yes is 612 and the total count of no is 347.

Scatter Plot

In this type of graph, the data is plotted against X and Y axis. We have used 'work_interface' attribute on X axis and 'treatment' on Y axis. In Weka, there is option to select different color to each value of attribute. We have selected blue for 'Male' and red for 'Female'. This helped for better understanding of data.

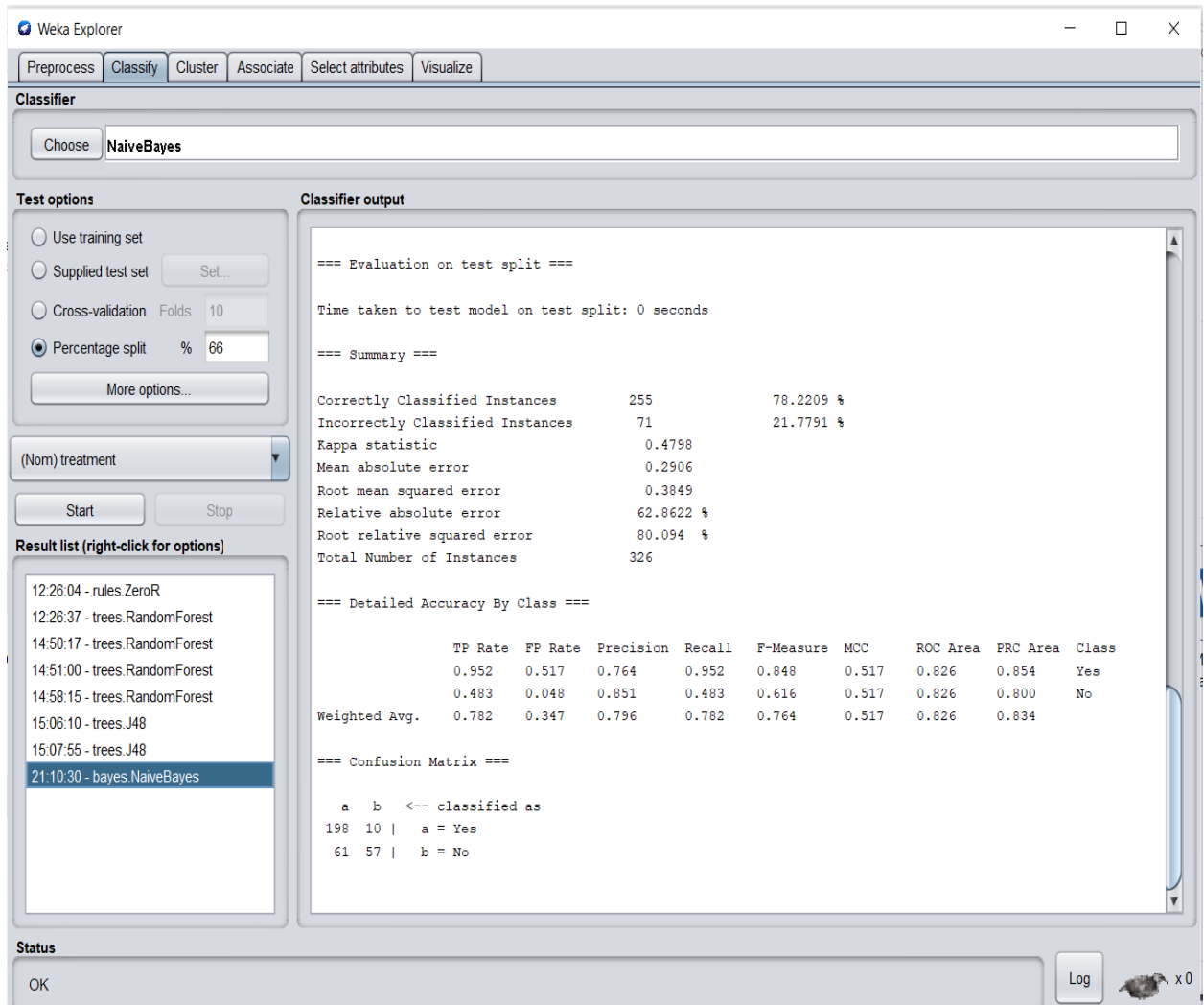


This visualization shows that number of people who has work interfere mostly take treatment. It is also predictable that the number is greater for people who has work interfere but takes treatment often or rarely. These group of people are target for this model.

Data Mining Technique and Optimization Technique used

- **Naïve Bayes Classification:**

After performing the Naïve Bayes classification, we observed that Correctly classified Instances are 255 with the accuracy of 78.2209%. And Incorrectly Classified Instances are 71 with inaccuracy of 21.7791%.



The screenshot displays the Weka Explorer interface with the NaiveBayes classifier selected. The 'Test options' section shows 'Percentage split' at 66%. The 'Classifier output' section provides a summary of performance metrics and a detailed accuracy breakdown by class.

Summary Metrics:

Metric	Value	Percentage
Correctly Classified Instances	255	78.2209 %
Incorrectly Classified Instances	71	21.7791 %
Kappa statistic	0.4798	
Mean absolute error	0.2906	
Root mean squared error	0.3849	
Relative absolute error	62.8622 %	
Root relative squared error	80.094 %	
Total Number of Instances	326	

Detailed Accuracy By Class:

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
Weighted Avg.	0.952	0.517	0.764	0.952	0.848	0.517	0.826	0.854	Yes
	0.483	0.048	0.851	0.483	0.616	0.517	0.826	0.800	No

Confusion Matrix:

a \ b	a = Yes	a = No
b = Yes	198	10
b = No	61	57

To make this value optimized, we used Multiclass Classifier optimization technique for this data mining technique.

Multiclass Classifier:

For this technique we split 80% of data as training set. By doing this we got accuracy of 79.1667%. Where correctly classified instances are 152. For 40 instances which are incorrectly classified we found inaccuracy rate is 20.8333%. By using MultiClass Classifier, optimization technique for Naïve Bayes Classification, we were able to improve the accuracy rate from 78.22% to 79.166%.

The screenshot shows the Weka Explorer interface with the 'Classify' tab selected. The classifier chosen is 'MultiClassClassifier'. The test options are set to 'Percentage split' at 80%. The classifier output is displayed, showing the evaluation on the test split.

Test options

- ☐ Use training set
- ☐ Supplied test set
- ☐ Cross-validation Folds 10
- ☒ Percentage split % 80

Classifier output

```
=== Evaluation on test split ===

Time taken to test model on test split: 0 seconds

=== Summary ===

Correctly Classified Instances      152      79.1667 %
Incorrectly Classified Instances    40      20.8333 %
Kappa statistic                    0.5358
Mean absolute error                 0.2951
Root mean squared error             0.3911
Relative absolute error             63.3321 %
Root relative squared error         80.4787 %
Total Number of Instances          192

=== Detailed Accuracy By Class ===

      TP Rate  FP Rate  Precision  Recall   F-Measure  MCC     ROC Area  PRC Area  Cl
      0.908    0.397    0.788     0.908    0.844     0.548    0.833    0.868    Ye
      0.603    0.092    0.800     0.603    0.688     0.548    0.833    0.780    No
Weighted Avg.   0.792    0.281    0.793     0.792    0.784     0.548    0.833    0.834

=== Confusion Matrix ===

  a  b  <-- classified as
108 11 |  a = Yes
 29 44 |  b = No
```

Result list (right-click for options)

- 21:20:35 - bayes.NaiveBayes
- 21:21:03 - bayes.NaiveBayes
- 21:22:29 - meta.MultiClassClassifier

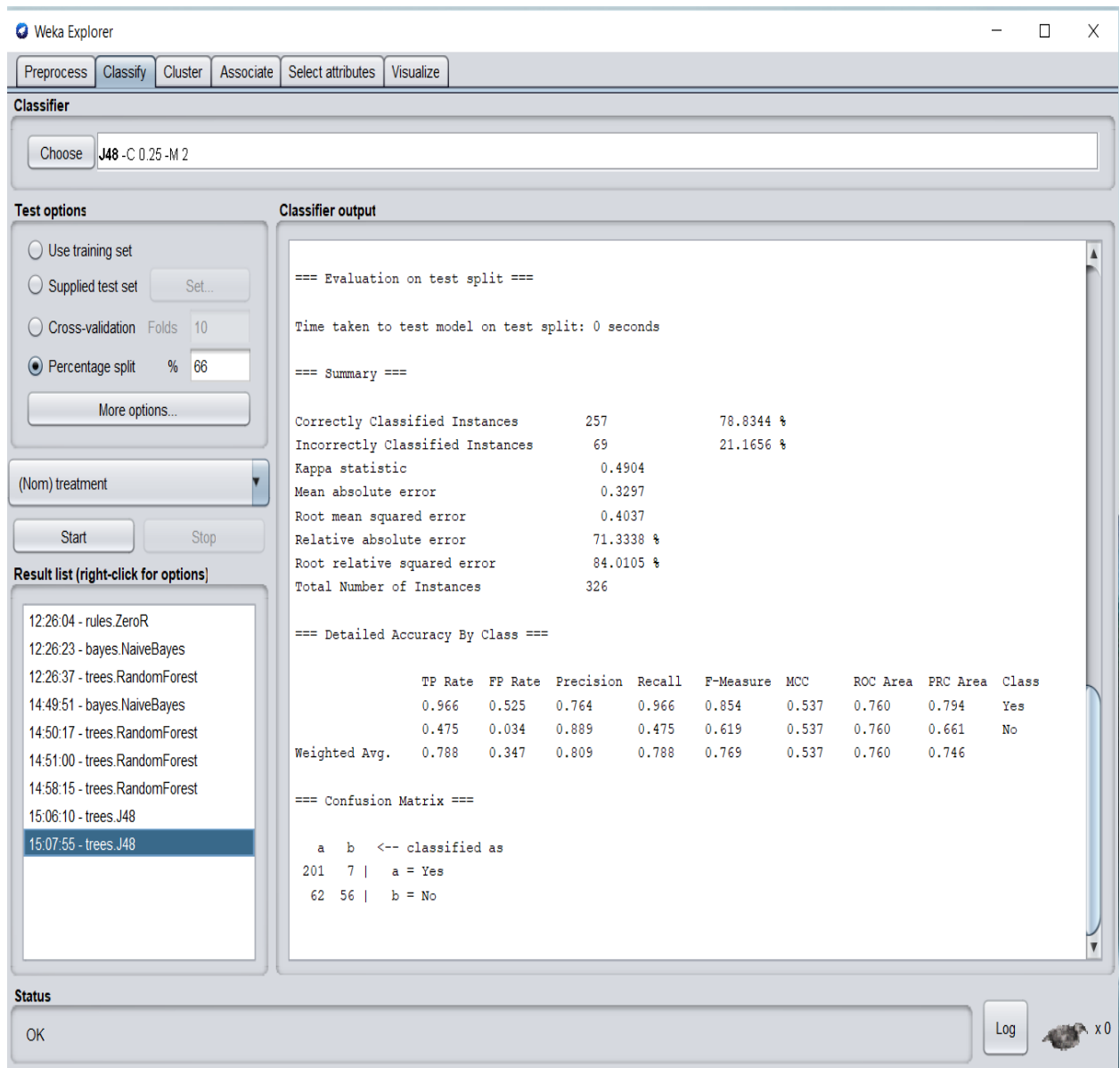
Status

OK

Log

- **Decision Tree:**

Where Decision Tree Classification results were, Total Number of Instances are 326 from which Correctly Classified Instances and Incorrectly Classified Instances are 257 and 69, respectively. We also find the Mean absolute error which is 0.3297. Alongside Root mean squared error 0.4037. Accuracy with uprunning set to false - 77.83% and with uprunning set to true - 78.8344%.



Weka Explorer

Preprocess **Classify** Cluster Associate Select attributes Visualize

Classifier

Choose **J48 -C 0.25-M 2**

Test options

☐ Use training set
☐ Supplied test set **Set...**
☐ Cross-validation Folds **10**
☒ Percentage split % **66**
More options...

(Nom) treatment

Start Stop

Result list (right-click for options)

- 12:26:04 - rules.ZeroR
- 12:26:23 - bayes.NaiveBayes
- 12:26:37 - trees.RandomForest
- 14:49:51 - bayes.NaiveBayes
- 14:50:17 - trees.RandomForest
- 14:51:00 - trees.RandomForest
- 14:58:15 - trees.RandomForest
- 15:06:10 - trees.J48
- 15:07:55 - trees.J48**

Classifier output

```

=== Evaluation on test split ===

Time taken to test model on test split: 0 seconds

=== Summary ===

Correctly Classified Instances      257      78.8344 %
Incorrectly Classified Instances    69      21.1656 %
Kappa statistic                    0.4904
Mean absolute error                 0.3297
Root mean squared error             0.4037
Relative absolute error             71.3338 %
Root relative squared error         84.0105 %
Total Number of Instances          326

=== Detailed Accuracy By Class ===

          TP Rate  FP Rate  Precision  Recall   F-Measure  MCC    ROC Area  PRC Area  Class
          0.966   0.525   0.764     0.966   0.854     0.537   0.760    0.794    Yes
          0.475   0.034   0.889     0.475   0.619     0.537   0.760    0.661    No
Weighted Avg.   0.788   0.347   0.809     0.788   0.769     0.537   0.760    0.746

=== Confusion Matrix ===

  a  b  <-- classified as
201  7  |  a = Yes
 62 56  |  b = No
  
```

Status

OK Log x 0

To improve the accuracy for Decision Tree, we used CV Parameter selection optimization technique.

CV Parameter Selection:

We divided 80% of dataset as training set. Obtained accuracy of 78.125% for Correctly Classified 150 Instances. Whereas, 21.875% for Incorrectly Classified 42 Instances. The precision and recall values we found by this technique are 0.78 and 0.781 respectively.

The screenshot shows the Weka Explorer interface with the 'Classify' tab selected. The classifier chosen is 'CVPParameterSelection'. The test options are set to 'Percentage split' at 80%. The classifier output window displays the following results:

```
=== Evaluation on test split ===  
Time taken to test model on test split: 0 seconds  
  
=== Summary ===  
Correctly Classified Instances      150      78.125 %  
Incorrectly Classified Instances    42      21.875 %  
Kappa statistic                    0.5153  
Mean absolute error                 0.3165  
Root mean squared error             0.4108  
Relative absolute error             67.9288 %  
Root relative squared error         84.5243 %  
Total Number of Instances          192  
  
=== Detailed Accuracy By Class ===  


|               | TP Rate | FP Rate | Precision | Recall | F-Measure | MCC   | ROC Area | PRC Area | Cl  |
|---------------|---------|---------|-----------|--------|-----------|-------|----------|----------|-----|
| Yes           | 0.891   | 0.397   | 0.785     | 0.891  | 0.835     | 0.524 | 0.778    | 0.795    | Yes |
| No            | 0.603   | 0.109   | 0.772     | 0.603  | 0.677     | 0.524 | 0.778    | 0.660    | No  |
| Weighted Avg. | 0.781   | 0.288   | 0.780     | 0.781  | 0.775     | 0.524 | 0.778    | 0.743    |     |

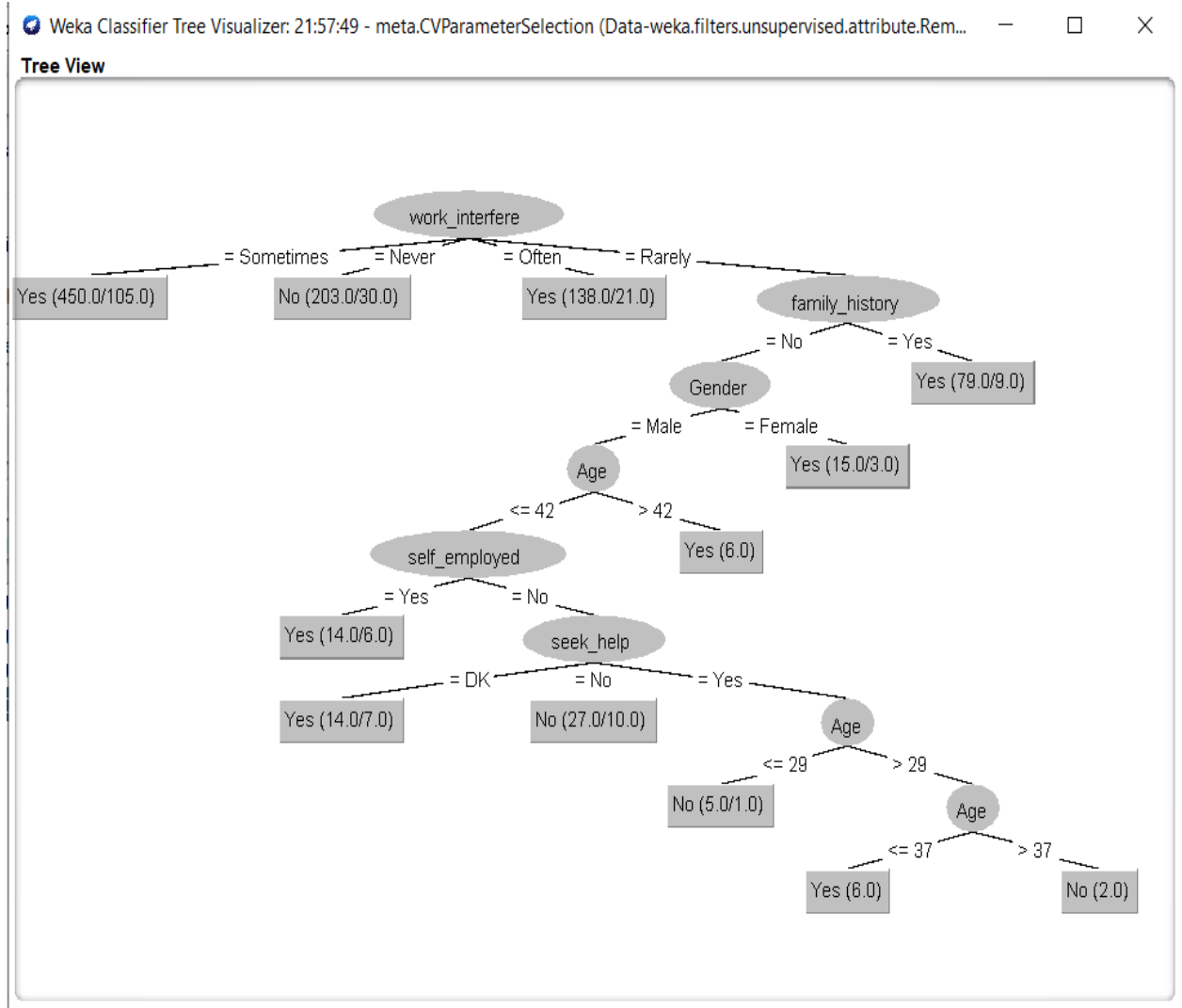
  
=== Confusion Matrix ===  


| a \ b | <-- classified as |        |
|-------|-------------------|--------|
|       | a = Yes           | b = No |
| 106   | 13                |        |
| 29    | 44                |        |


```

The result list on the left shows the execution history, with the most recent entry being '21:57:49 - meta.CVPParameterSelection'. The status bar at the bottom indicates 'OK'.

Below is the optimized decision tree we found using CV Parameter Selection. For decision tree, observations that we found is work interfere is a parameter where data is being split and the weightage of attributes are as shown in figure.



We observed that, Decision Tree has more accuracy 78.8344% than the optimization technique CV Parameter Selection which is 78.125%. We came to the conclusion that the Decision Tree is already optimized.

- **Random Forest:**

After performing the classification, we observed that Correctly classified Instances are 142 with the accuracy of 73.9583%. And Incorrectly Classified Instances are 50 with inaccuracy of 26.0417%.

Total Number of Instances are 192.

Classifier

Choose **RandomForest-P 100 -I 100 -num-slots 1 -K 0 -M 1.0 -V 0.001 -S 1**

Test options

☐ Use training set
☐ Supplied test set (Set...)
☐ Cross-validation Folds 10
☒ Percentage split % 80

More options...

(Nom) treatment

Start Stop

Result list (right-click for options)

- 21:49:53 - trees.J48
- 21:50:05 - trees.J48
- 21:50:36 - trees.J48
- 21:57:49 - meta.CVPParameterSelection
- 22:10:38 - trees.RandomForest**
- 22:18:35 - meta.RandomSubSpace

Classifier output

```

=== Evaluation on test split ===

Time taken to test model on test split: 0.08 seconds

=== Summary ===

Correctly Classified Instances      142      73.9583 %
Incorrectly Classified Instances    50      26.0417 %
Kappa statistic                    0.4355
Mean absolute error                 0.2927
Root mean squared error             0.444
Relative absolute error             62.8224 %
Root relative squared error        91.3736 %
Total Number of Instances          192

=== Detailed Accuracy By Class ===

          TP Rate  FP Rate  Precision  Recall  F-Measure  MCC      ROC Area  PRC Area  Cl
          0.824    0.397    0.772     0.824    0.797     0.437    0.772    0.806    Ye
          0.603    0.176    0.677     0.603    0.638     0.437    0.771    0.695    No
Weighted Avg.   0.740    0.313    0.736     0.740    0.736     0.437    0.771    0.764

=== Confusion Matrix ===

  a  b  <-- classified as
98 21 |  a = Yes
29 44 |  b = No
  
```

Status

OK Log x 0

Here, the Precision Value is 0.736. Also, 98 True(yes) are instances, 29 are False(No) instances and 44 are true negative instances.

To improve the accuracy for Random Forest, we used Random Subspace Optimization technique.

Random Subspace:

We divided 80% of dataset as training set. Obtained accuracy of 75.5208% for Correctly Classified 145 Instances. Whereas, 24.4792% for Incorrectly Classified 47 Instances. The precision and recall values we found by this technique are 0.754 and 0.755 respectively.

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose **RandomSubSpace** -P 0.5 -S 1 -num-slots 1 -I 10 -W weka.classifiers.trees.REPTree -- -M 2 -V 0.001 -N 3 -S 1 -L -1 -I 0.0

Test options

☐ Use training set
☐ Supplied test set Set...
☐ Cross-validation Folds 10
☒ Percentage split % 80
 More options...

(Nom) treatment

Start Stop

Result list (right-click for options)

- 21:49:53 - trees.J48
- 21:50:05 - trees.J48
- 21:50:36 - trees.J48
- 21:57:49 - meta.CVPParameterSelection
- 22:10:38 - trees.RandomForest
- 22:18:35 - meta.RandomSubSpace

Classifier output

```

=== Evaluation on test split ===

Time taken to test model on test split: 0.35 seconds

=== Summary ===

Correctly Classified Instances      145           75.5208 %
Incorrectly Classified Instances    47           24.4792 %
Kappa statistic                    0.45
Mean absolute error                 0.3833
Root mean squared error             0.4216
Relative absolute error             82.2695 %
Root relative squared error         86.7468 %
Total Number of Instances          192

=== Detailed Accuracy By Class ===

               TP Rate  FP Rate  Precision  Recall   F-Measure  MCC      ROC Area  PRC Area  Cl
Weighted Avg.   0.755   0.330   0.754     0.755   0.745     0.464   0.792    0.787   Ye
0.891   0.466   0.757     0.891   0.819     0.464   0.792    0.821   Ye
0.534   0.109   0.750     0.534   0.624     0.464   0.792    0.731   No

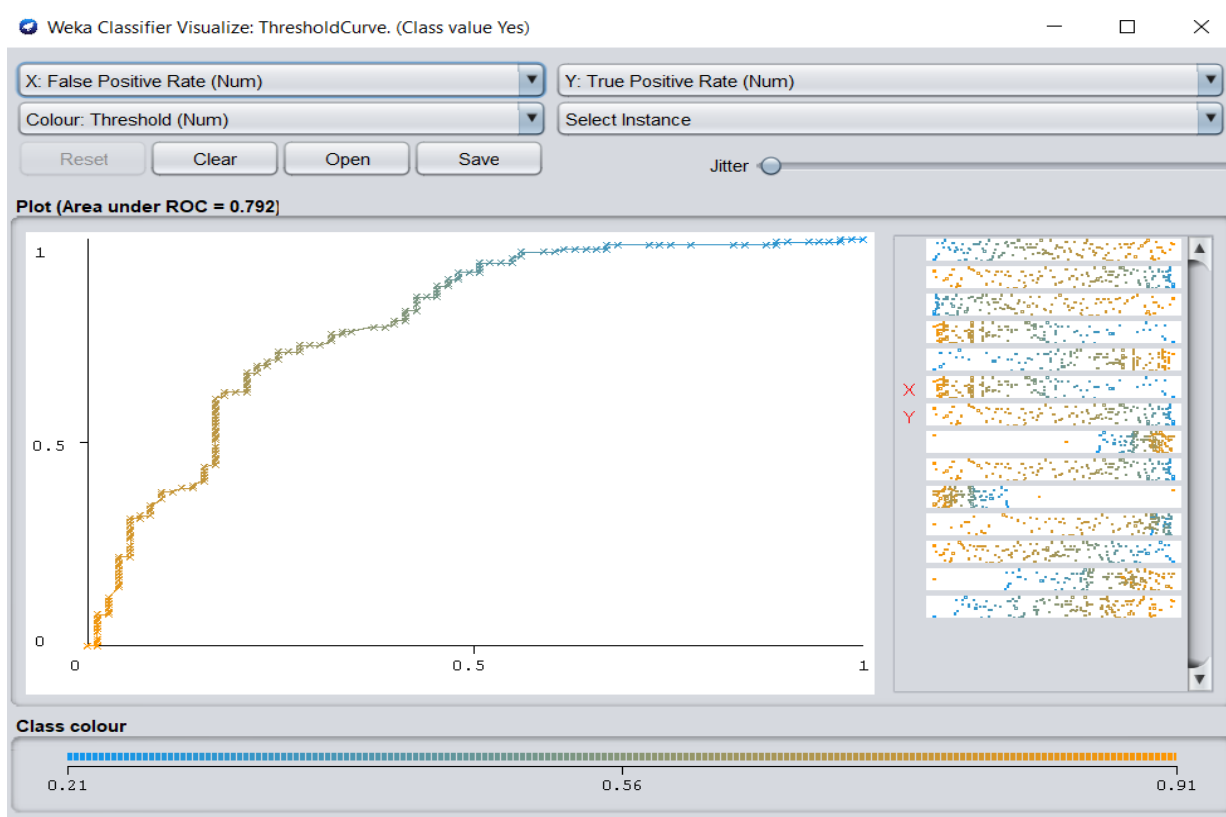
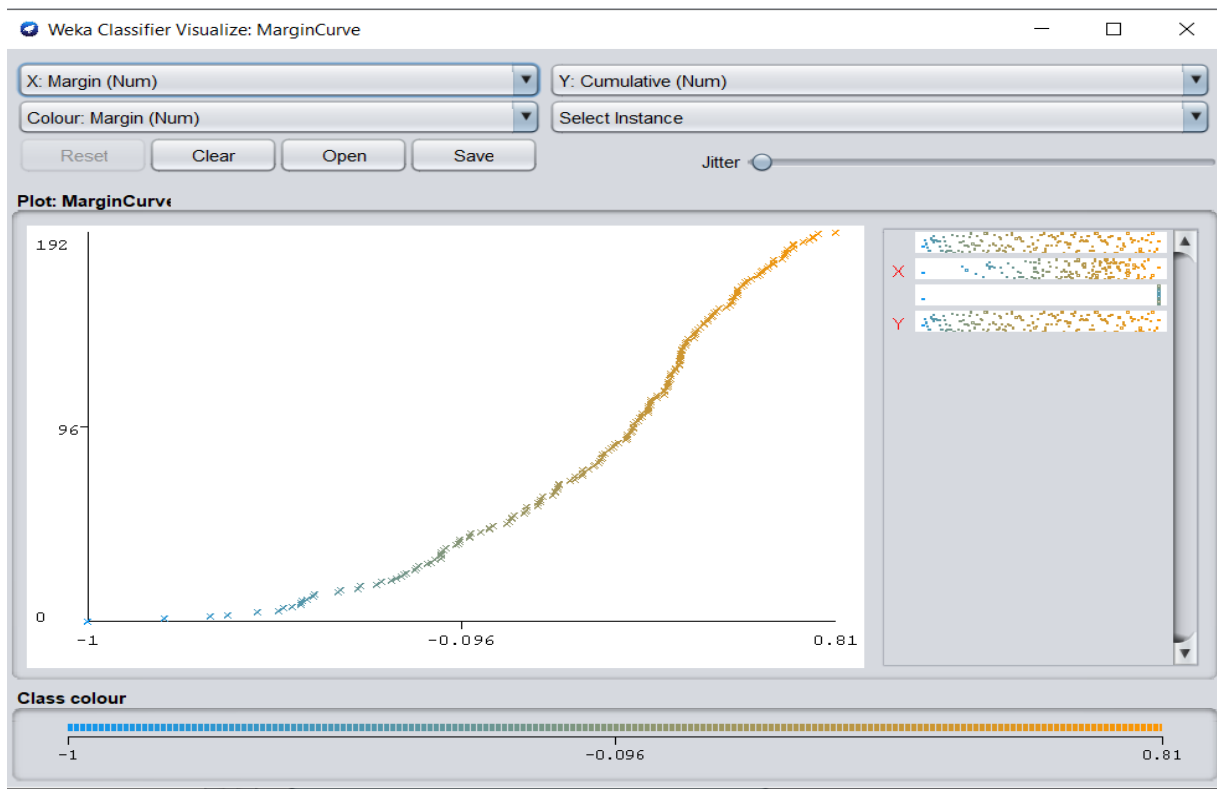
=== Confusion Matrix ===

  a  b  <-- classified as
106 13 | a = Yes
 34 39 | b = No
  
```

Status

OK Log x 0

In Random SubSpace, we are also calculating margin curve and threshold curve for better understanding. The Area under ROC for threshold is 0.792 model.



By using Random Subspace, optimization technique, we were able to improve the accuracy rate from 73.9583% to 75.5208%.

Conclusion:

By performing the optimization phase, we were able to improve the accuracy of the data model for different data mining techniques such as Naïve Bayes, Decision Tree, and Random Forest. We observed the increase in accuracy by using optimization techniques Multiclass Classifier (79.16%), and Random Subspace (75.52%). However, for CV Parameter Selection the accuracy (78.125%) was less than the original Decision Tree Classification which has an accuracy of 78.8344%. Therefore, we conclude that accuracy values are not always final, you can improve accuracy by using different optimization techniques.

Repository:

<https://github.com/vsala2/DataMining>