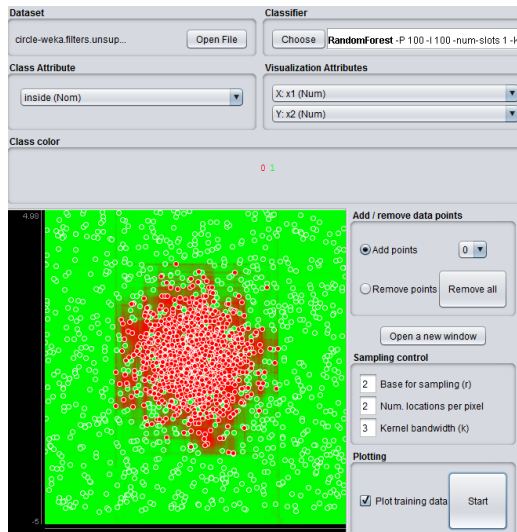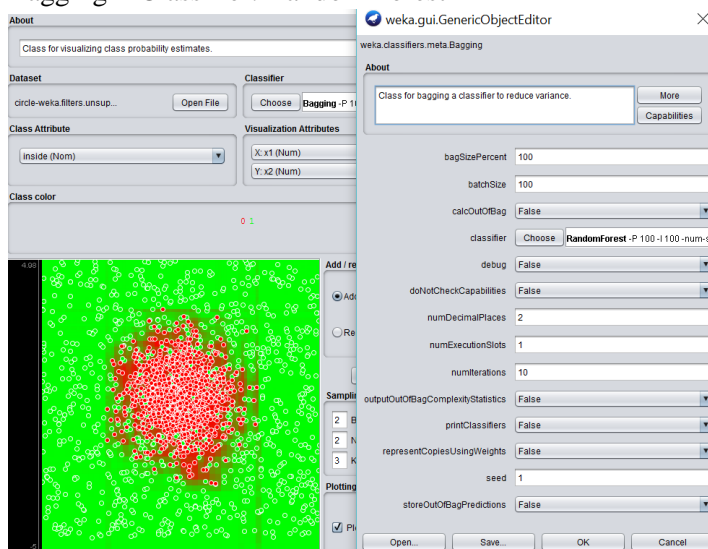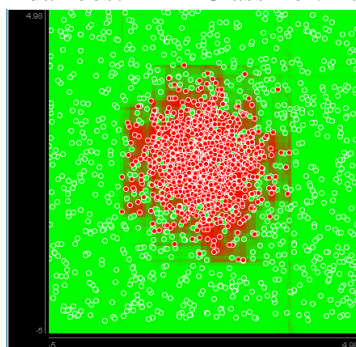# 1. Dataset: Circle

- Classifier: Random Forest



Random Forest performs basically good.
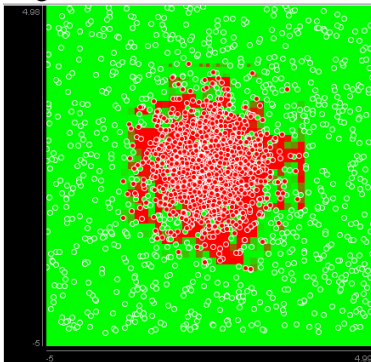
- Bagging + Classifier: Random Forest



As can be seen, for Random Forest Classifier, adding bagging does not perform better, remaining almost the same.
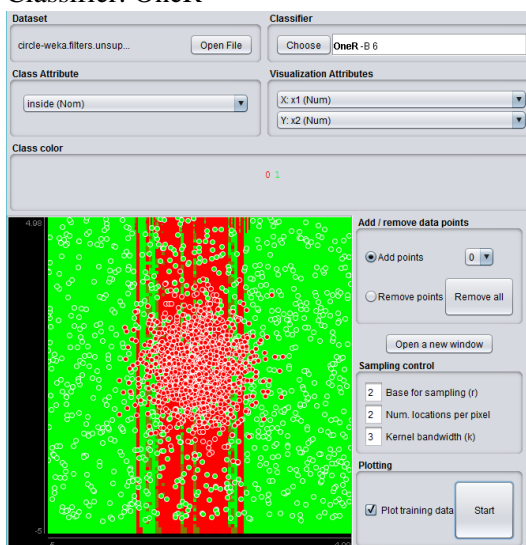
- AdaBoost M1 + Classifier: Random Forest

As can be seen, for Random Forest Classifier, adding AdaBoost seems perform a bit better or remaining the same. While, the running speed increases.

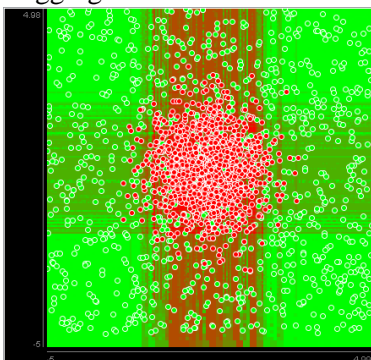- LogitBoost + Classifier: Random Forest



As can be seen, adding LogitBoost performs a bit worse for Random Forest Classifier.
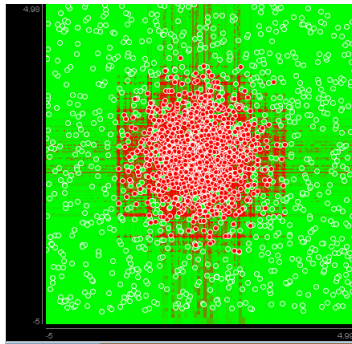
- Classifier: OneR



Basically, one vertical bar is selected for OneR.
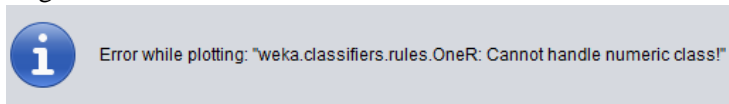
- Bagging + Classifier: OneR



As can be seen, for OneR Classifier, adding bagging looks adding another dimension when classifying. OneR only selects a long band, while Bagging + OneR selects two orthogonal bands with a square in the center, performing better.
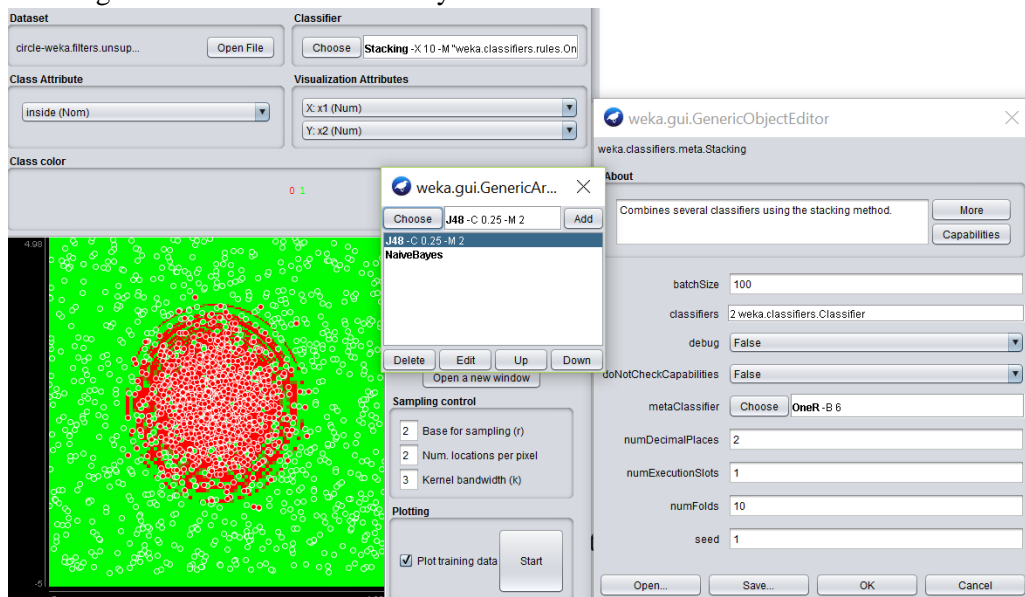
- AdaBoost + Classifier: OneR

As can be seen, adding AdaBoost to OneR performs much better, and runs more quickly. The bias error of the model OneR decreases a lot.

- LogitBoost + Classifier: OneR



Error while plotting: "weka.classifiers.rules.OneR: Cannot handle numeric class!"
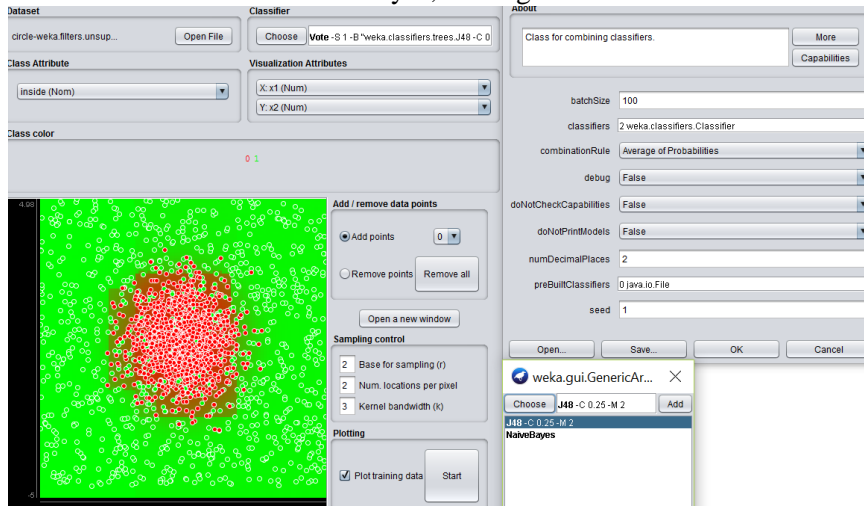
- Stacking: Decision Trees + NaïveBayes + metaclassifier OneR
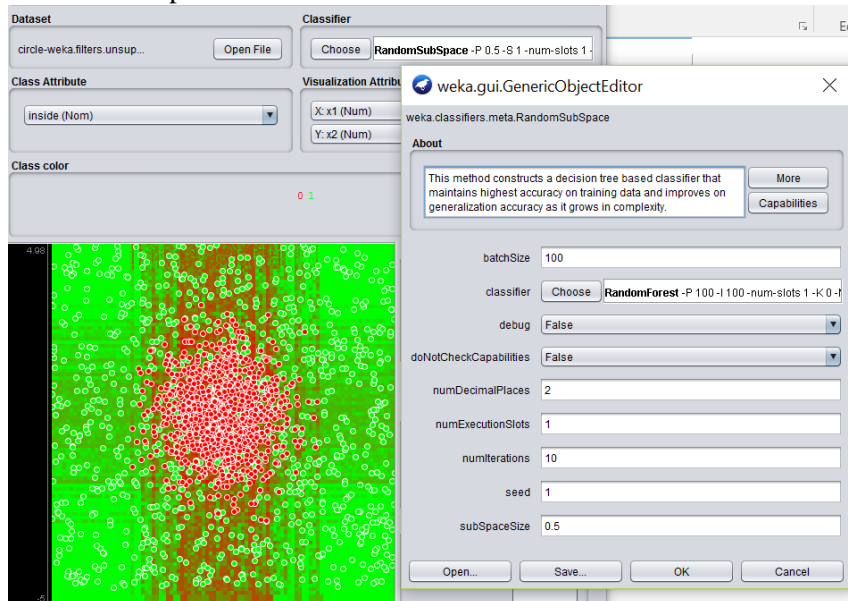


As can be seen, it selects an almost exact circle in the center. This performs much better.

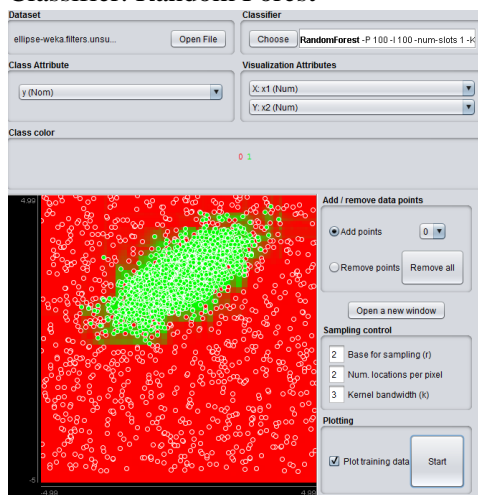- Vote: Decision Trees + NaiveBayes, Average of Probabilities

- As can be seen, the Vote result with combination rule of Average of Probabilities is not quite good. The selected areas consist of several overlapping blocks.

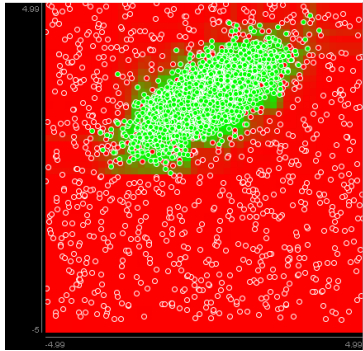- RandomSubSpace: Random Forest



RandomSubSpace is an attribute bagging. As can be seen, it basically selects a rough band, performs worse than Random Forest itself.

## 2. Dataset: Ellipse

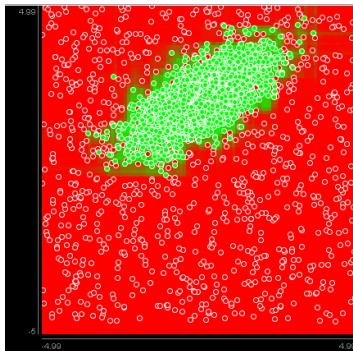- Classifier: Random Forest



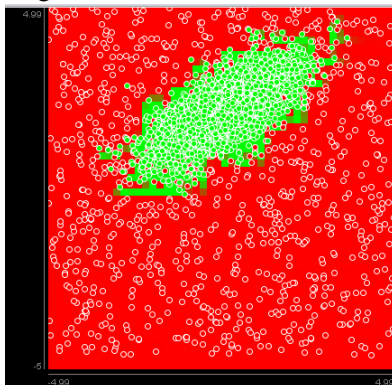- Bagging + Classifier: Random Forest

As can be seen, for Random Forest Classifier, adding bagging does not perform better, remaining almost the same. The running time is much longer.

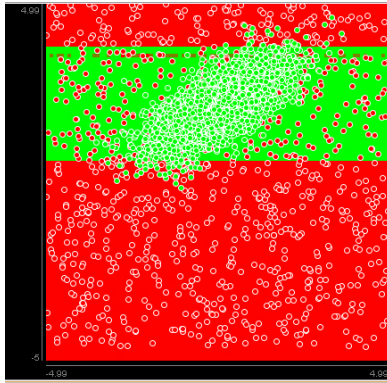- AdaBoost M1 + Classifier: Random Forest



As can be seen, for Random Forest Classifier, adding AdaBoost seems perform a bit better or remaining the same. While, the running speed increases as well.
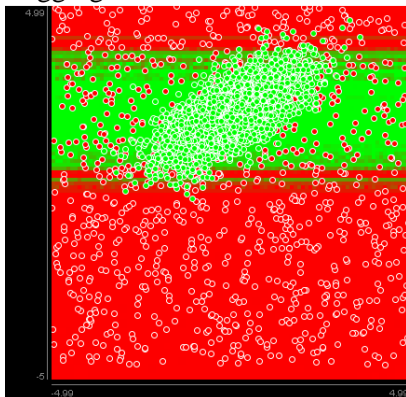
- LogitBoost + Classifier: Random Forest



As can be seen, adding LogitBoost performs a bit worse for Random Forest Classifier.
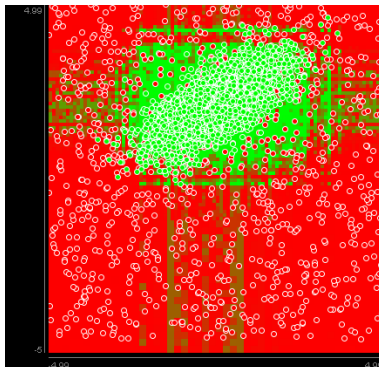
- Classifier: OneR

A horizontal bar is selected.

- Bagging + Classifier: OneR



As can be seen, for OneR Classifier, adding bagging looks more fine-grained. OneR only selects a long band, while Bagging + OneR selects not exactly one horizontal bar.

- AdaBoost + Classifier: OneR



As can be seen, adding AdaBoost to OneR performs better, and runs more quickly. The bias error of the model OneR decreases a lot. The selected area looks like a square.

- LogitBoost + Classifier: OneR



Error while plotting: "weka.classifiers.rules.OneR: Cannot handle numeric class!"

- Stacking: Decision Trees + NaïveBayes + metaclassifier OneR

As can be seen, it selects a ladder-like shape. This performs better.
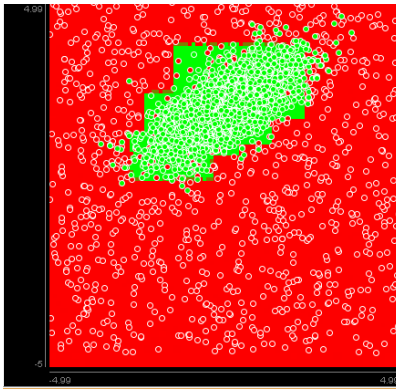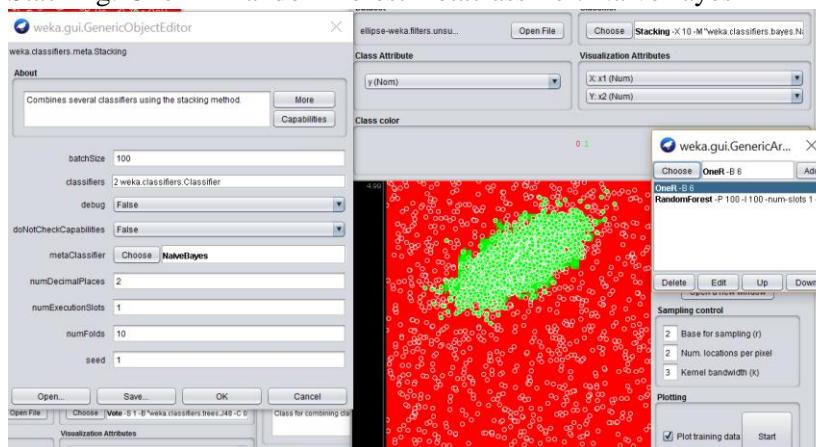
- Stacking: OneR + Random Forest metaclassifier: NaiveBayes



The above stacking of different classifiers performs better than the previous stacking.

- Vote: Decision Trees + NaiveBayes, Average of Probabilities



As can be seen, the Vote result with combination rule of Average of Probabilities is not quite good. The selected areas consist of several overlapping blocks, similar to stacking with majority voting rule.

- RandomSubSpace: Random Forest

RandomSubSpace is an attribute bagging. As can be seen, it basically selects a rough square, performing much worse than Random Forest itself.
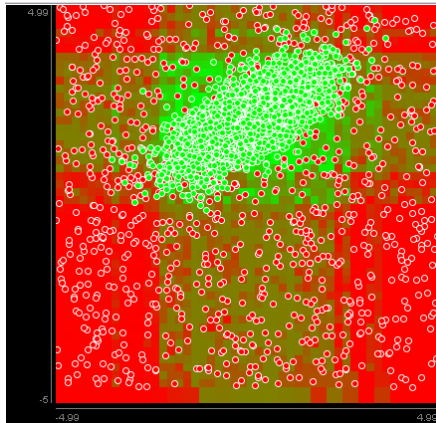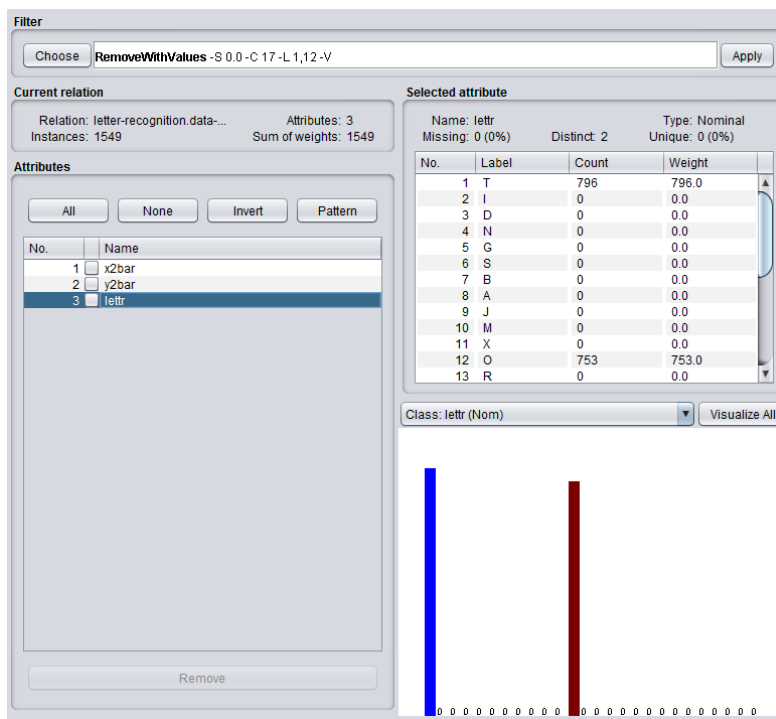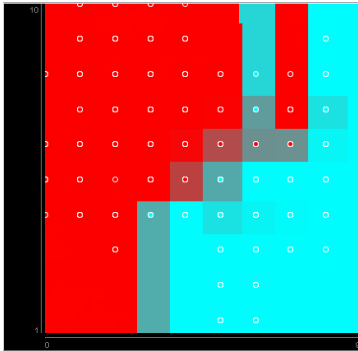
## 3. Dataset: letter-recognition.data (OCR)

Preprocessing:

I Only select letter "T" and "O", and select two attributes "x2bar" and "y2bar". The two attributes look easier to distinguish for the two letters.
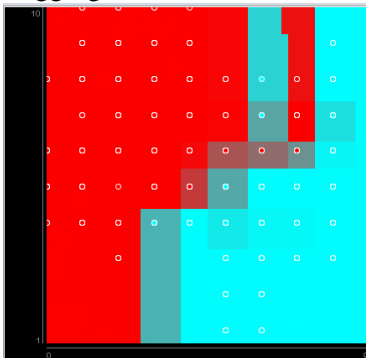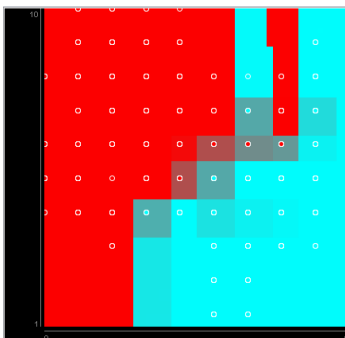


- Classifier: Random Forest

Random Forest performs basically good. Red points are almost in red part, and blue points are almost in blue part.

- Bagging + Classifier: Random Forest



As can be seen, for Random Forest Classifier, adding bagging performs the same. The values are all integers, so it should be easier to classify. The running time is much longer.

- AdaBoost M1 + Classifier: Random Forest
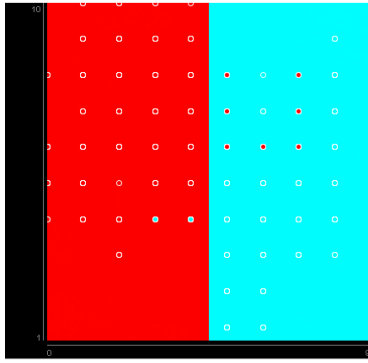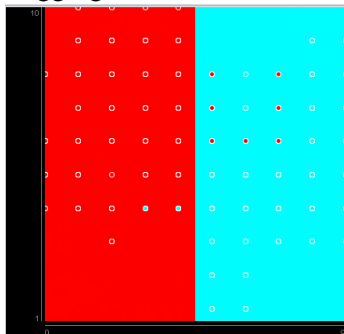


As can be seen, for Random Forest Classifier, adding AdaBoost performs the same. The values are all integers, so it should be easier to classify. While, the running time increases a lot.

- LogitBoost + Classifier: Random Forest
  The running time is much longer, so I stop doing this one. Similarly, the result should remain almost the same as above.

- Classifier: OneR

Basically, one vertical parting line is generated with OneR. There are several error points.

- Bagging + Classifier: OneR



As can be seen, for OneR Classifier, adding bagging performs exactly the same.

- AdaBoost + Classifier: OneR



As can be seen, adding AdaBoost to OneR performs much better. The bias error of the model OneR decreases a lot. There are 4 to 5 points that are incorrectly classified.

- LogitBoost + Classifier: OneR



Error while plotting: "weka.classifiers.rules.OneR: Cannot handle numeric class!"

- Stacking: Decision Trees + NaïveBayes + metaclassifier OneR

As can be seen, there are 4 points that are incorrectly classified with majority voting.
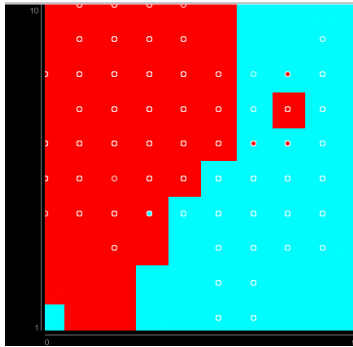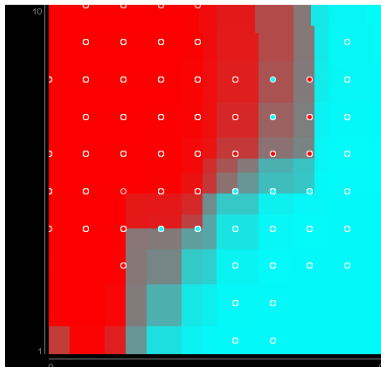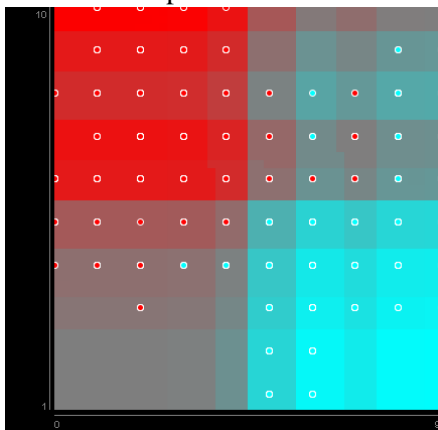
- Vote: Decision Trees + NaiveBayes, Average of Probabilities



As can be seen, the Vote result with combination rule of Average of Probabilities is basically good with some overlapping shadow. The selected areas consist of several overlapping blocks. There are 3 to 4 points that are incorrectly classified.

- RandomSubSpace: Random Forest



RandomSubSpace is an attribute bagging. As can be seen, it basically separates the whole into 4 parts including 2 shadow areas. It performs worse than Random Forest itself.

# 4. Dataset: names_ethnea_genni_country_sample

Preprocessing:

I only select PubCountry attribute as class attribute, and select Portugal and HongKong two countries. I select Last Name as the only attribute for classification.
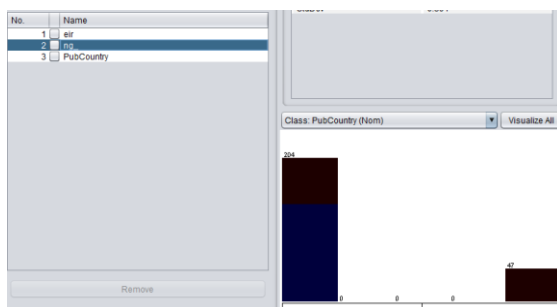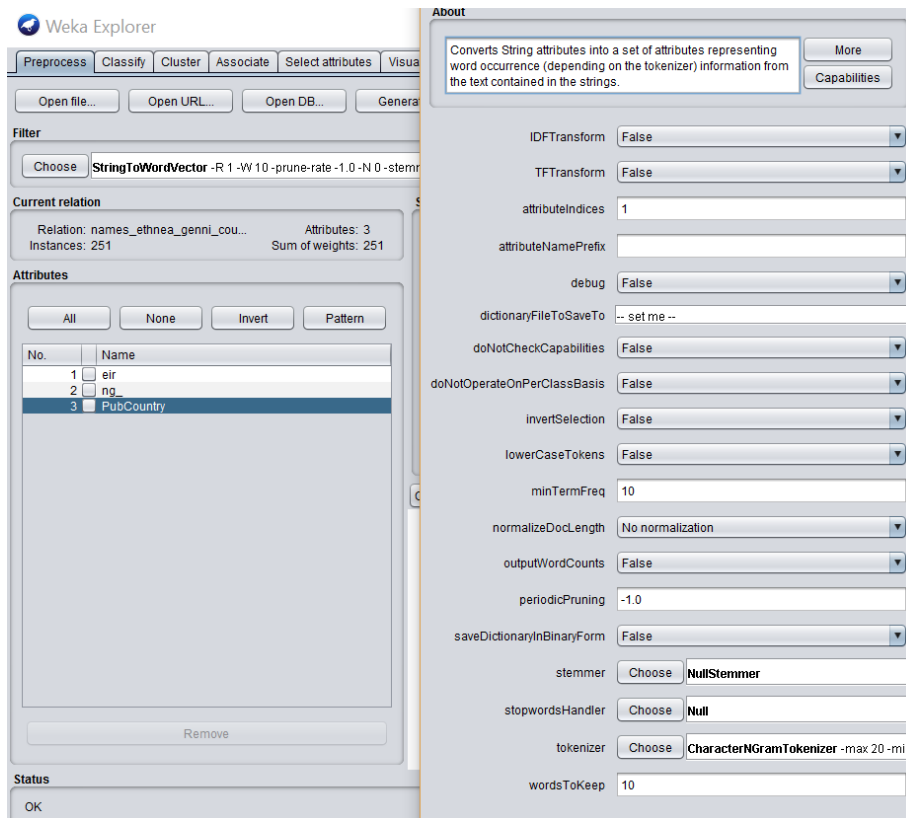
Filter: For First attribute, NominalToString -> StringToWordVector

CharacterNGramTokenizer -> set NGramMinSize and NGramMaxSize
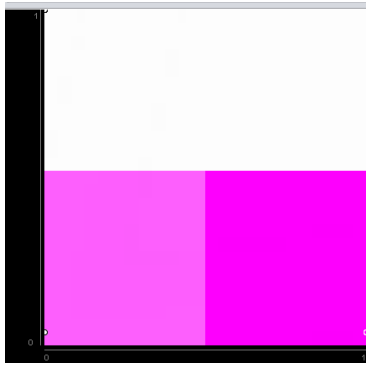
Set minTermFreq to 10.

Set WordsToKeep to 10.

Select "ng_" and "eir".





For "eir" and "ng_", instances only have numeric values "0" and "1".
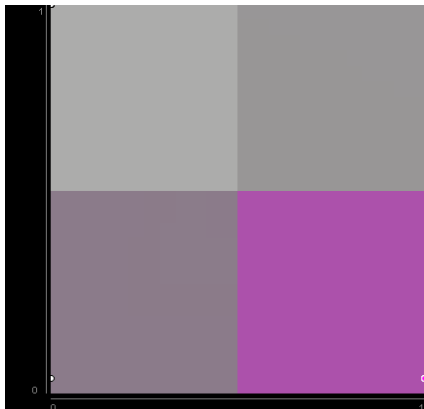
- Classifier: Random Forest

The two points are located at the left and right edges. There is not much difference between the color of the two blocks. It is not a good classification.
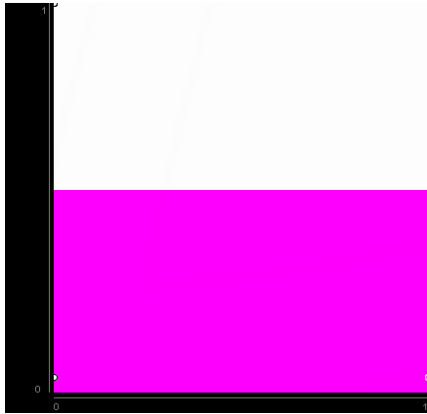
- Bagging + Classifier: Random Forest



As can be seen, for Random Forest Classifier, adding bagging performs the same. The running time is much longer. It is not a good classification either.

- AdaBoost M1 + Classifier: Random Forest



As can be seen, for Random Forest Classifier, adding AdaBoost produces 4 blocks. But it is not a good classification either. The two points are located at the left and right edges. There is not much difference between the color of the two blocks. The running time is much longer.

- LogitBoost + Classifier: Random Forest
  The running time is much longer, so I stop doing this one. Similarly, the result should remain almost the same as above.
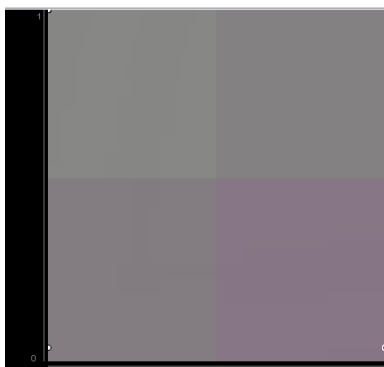
- Classifier: OneR

Basically, one vertical parting line is generated with OneR. The two points are not separated at all.

- Bagging + Classifier: OneR



As can be seen, for OneR Classifier, adding bagging performs exactly the same.

- AdaBoost + Classifier: OneR



As can be seen, adding AdaBoost to OneR performs better. The two points are separated in two different blocks. But there is not much difference between the colors of the two blocks.

- LogitBoost + Classifier: OneR



Error while plotting: "weka.classifiers.rules.OneR: Cannot handle numeric class!"

- Stacking: Decision Trees + NaïveBayes + metaclassifier OneR

The two points are not separated at all.

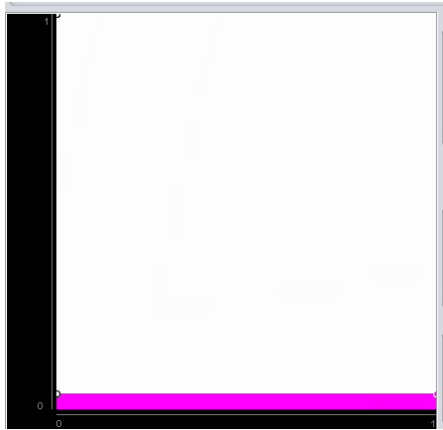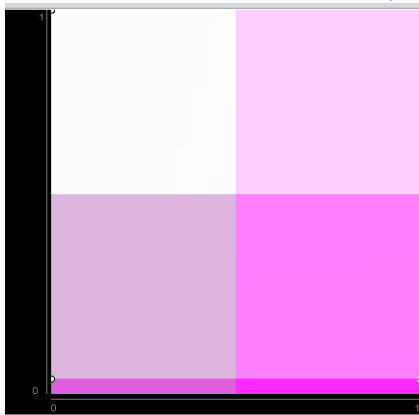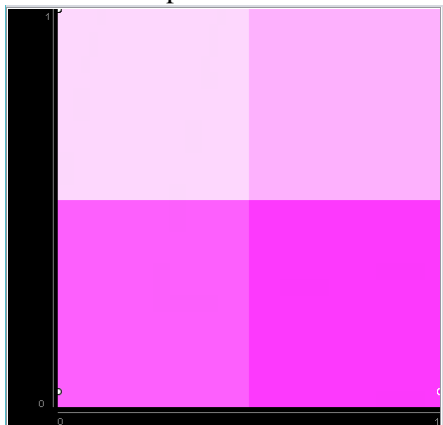- Vote: Decision Trees + NaiveBayes, Average of Probabilities



The two points are separated in two different blocks of different colors. But there is not much difference between the colors. It is not a really good classification.

- RandomSubSpace: Random Forest



RandomSubSpace is an attribute bagging. The two points are separated in two different blocks of different colors. But there is not much difference between the colors. It is not a really good classification.

# 5 General Conclusion

In general, Random Forest performs better than OneR, consuming longer time as well. (Random Forest: Bagging + random subspace for decision trees)

Bias is the model error, and variance represents training set error (whether representative or not).

For bagging, it performs better together with unstable classifiers than the classifiers themselves. Those classifiers include decision trees. Generally, bagging will decrease the variance error in some cases, however, bagging will not change the performance in other cases.

Boosting (not work for regression) probably performs better when adding to some classifiers. It will decrease bias error in most cases.

LogitBoost is additive logistic regression. Hard to explain.

Stacking method is combining models of different types. The default in Weka is majority voting. However, for voting with rule of average probability, some shadow areas will be generated based on the majority voting result.