
Implementation and Analysis of Random Forests

Jae Lee

School of Computing Science
Simon Fraser University
Burnaby BC V5A 1S6

Richard Mar

School of Computing Science
Simon Fraser University
Burnaby BC V5A 1S6

Robin White

School of Mechatronic Systems Engineering
Simon Fraser University
Surrey BC V3T 0A3

1 Introduction

In machine learning, there is often a tug-of-war between bias and variance; having high accuracy to observed data but not to lose generalization (or over-fit) to unseen data. This is often referred to as the "bias-variance tradeoff" and its consideration is a significant part of properly engineering machine learning algorithms. Often, regularization is used where there is an addition to the loss function to represent a cost to complexity, or in the way of neural networks, random dropout of neurons to force generalization [1]. Another method is to also use validation datasets so as to understand the loss from unseen data without actually testing on the test set. All of these methods add to the complexity of the algorithm and can also often lead to loss in accuracy of the training data.

In the early 90's, Tin Kam Ho from Bell Labs published a series of papers where he showed surprisingly that by combining independent learners in a unique way increased the accuracy of classifying handwritten digits monotonically; without suffering from over-adaptation to the training data. [2, 3, 4] The application of this method to decision trees in his '93 paper marked the introduction of Random Forests to the community. [2] Decision trees are simple yet effective classifiers, with high execution speed and easily relatable, however they are limited by their complexity for possible loss of generalization to unseen data. Some methods such as pruning have been used previously to try and increase generalization, however methods such as these usually come with a loss in accuracy toward training data. By using principles of stochastic modeling, Ho showed that tree-based classifiers could be arbitrarily expanded for increases in accuracy on unseen testing data without loss in training data accuracy. A characteristic which is still unique among machine learning classifiers. The concept is that multiple learners can compensate for the bias of a single learner and so trees are constructed from randomly selecting subspaces of the feature space. In this way, each tree generalizes in a different way.

Random Forests have been applied to a variety of machine learning tasks including classifications in ecology and geosciences, image segmentation in medical applications, business analytics, sporting analytics, as well as the unmentioned number of general data science applications. [5, 6, 7, 8, 9]. These traits of high accuracy and resistance to overfitting that random forests possess are fascinating and so we would like to further our knowledge in how they are able to achieve these feats. TODO: Figure out a better way to connect this section and the next one.

2 Approach

TODO: This section is supposed to be about what we want to do and why.

In order to develop our knowledge of random forests and investigate their learning capabilities, we implement the random forest algorithm in python and apply it on a specific dataset. Our dataset

of choice is the hockey player dataset from class as it possesses discrete and continuous variables which allow us to use the random forest algorithm as a classifier and regressor on the same dataset. In order to determine whether or not our implementation is correct, we compare it against other freely available implementations. Once we are confident that our implementation is comparable to other implementations, we explore varying specific learning parameters for random forests to learn their impact on random forest's ability to learn. Finally, armed with knowledge on how random forest's learning parameters should be set, we compare random forest's ability to predict a couple attributes from various hockey players to see how it fairs relative to other well known machine learning algorithms.

3 Experiments

TODO: Each subsection for this section should have the following format/flow:

- 1 to 2 sentence description of what we are testing.
- parameter setup and experimental design (probably enough so that people are able to try and reproduce our results) - results (plots and tables)
- brief discussion of our results with supporting citations of any claims we make

3.1 Implementation Comparisons

TODO: This section should be testing the accuracy of our implementation against Weka and scikit-learn to demonstrate are implementation is correct. Possibly come up with a better subsection title.

Table 1: Comparison of Random forest classifier for GP_greater_than_0

Machine Learning Package	Accuracy	Time
Weka	69.4%	<1sec
scikit-learn	70.7%	<1sec
our implementation	68.6%	1min 55sec

The accuracy of the projects implementation against Weka and scikit-learn is compared to demonstrate its correctness. Experiment was designed to first find the optimal parameters for each model by searching the parameter space using the k-fold cross-validation score as a metric. The next step was to compare both the parameter values and the accuracy of each optimal models. All three models were configured as classifiers for this test.

From Table 1, it can be seen that all three implementations of Random forest yield similar results. It was noted however that when performing on standardized data, such that values are centered by their means and divided by the standard deviation, the accuracy was significantly lower for all three implementations, around 55% accuracy (not shown here). This can be due to the sensitivity of the splitting using Gini impurity or entropy with small values, or perhaps some information is lost on the impact of certain features. The time taken for completion is difficult to compare since coding implementations are different (python vs C), however to improve speed, multi-processing was implemented. Using 4 processes running the same parameters as shown in table 4 for target class GP_greater_than_0, the time of ≈ 2 minutes is slow, however not unreasonable, this addition was a great improvement over a single process taking 6min 58sec.

Table 2: Comparison of Random forest classifier for sum_7yr_GP

Machine Learning Package	Mean square error	Time
Weka	9185.3	<1sec
scikit-learn	?	<1sec
our implementation	?	?

Note: Needed to use Weka Random forest, not bagging with J48 since can't be used for regression, instead REPTree is used. Regression splitting using variance. TODO some explanation and references here.

3.2 Parameters Exploration

TODO: This section should be exploring how the parameters affect a random forest's ability to learn and our best parameters. Possibly come up with a better subsection title.

3.2.1 Number of Trees

Random forests utilize an ensemble method where by a collection of weak learners is combined to produce a more accurate and robust model. We explore this by investigating the classification accuracy of our random forests implementation by varying the number of decision trees which are combined to produce the final prediction of our model. We keep parameters such as maximum tree depth, number of features, and minimum split size constant, with values 5, 5, 2 respectively.

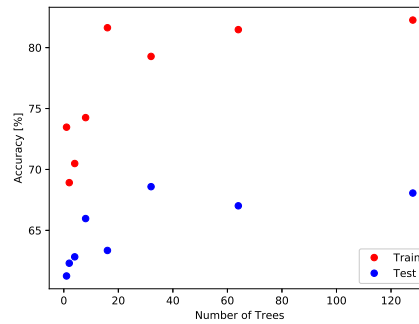


Figure 1: Showing Random forest inhibition to produce over fitting errors even with increased complexity by large number of trees.

The ability of Random forests to resist over fitting is largely due to the number of trees and features, which is demonstrated in figure 1. Randomization increases bias but makes it possible to reduce the variance of the corresponding ensemble model through averaging. [10]

3.2.2 Tree Depth

The depth of each decision tree used to make up the ensemble used in Random forests has a large impact on the accuracy of the model. We investigate this by comparing the accuracy of our model implementation as a function of the maximum depth of each tree. That is, the tree can have a depth that is less than this value, depending on the features used to build the tree, however once it reaches this maximum depth, the leaf is forced to be a terminal and predict the target class. We keep other parameters such as number of trees, number of features, and minimum split size constant, with values 64, 5, 2 respectively.

It has been shown that by increasing the maximum depth of each tree, a level of over fitting can be present, and similarly, having shallow trees can produce low-confidence predictions. [11] It is thus important that an appropriate value of the maximum depth be used, as its optimal value is related to the problem complexity. We notice a small decrease in the accuracy at 20 features, as shown in figure 2 which suggests that Random forests are still susceptible to over fitting, albeit less sensitive than many other machine learning algorithms without active regularization.

3.2.3 Features Per Tree

It is discussed in literature that a optimal value for the number of features to be used is $\frac{\ln(M)}{\ln(2)}$ of the total number of features, M . [12] To investigate this we compare the classification accuracy while

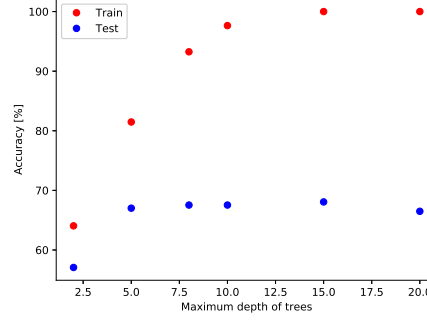


Figure 2: Showing accuracy of Random forest with 64 trees with varying maximum depth of each decision tree. It should be noted that there is a slight decrease in the test accuracy at 20 features which can possibly be related to over fitting on the training data which reaches 100% after 15.

varying the number of features, with all other parameters such as number of trees, maximum depth, and minimum split size constant, with values 64, 5, 2 respectively.

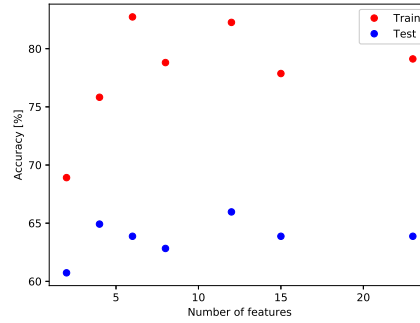


Figure 3: Showing the impact of the number of features used when constructing decision trees in Random forest with 64 trees and max depth 5.

Following our results, shown in figure 3, we observe a maximum train accuracy for 5 features, which agrees with literature for an optimal number of features. We observe, however, that the number of features used has less of an impact on the total accuracy when comparing to the above observations for varying number of trees and maximum tree depth. The maximum depth for optimal test accuracy also does not agree with the training value, and the variance between each run makes a conclusive result difficult on the over-all impact.

3.2.4 Grid Search Parameter Optimization

Following the observations after varying each parameter individually, a grid search method was employed to find the optimal values for number of trees, maximum depth, and number of features for our implementation of Random forest. The range of values used for the grid search were as follows:

Trees 2,8,16,32,64,128

Depth 2,5,8,10,15,20

Features 2,5,8,10,15,23

With optimal values shown in table 3. Figure 4 shows the effect of the different parameters on the accuracy of the classifier for the test set. As can be observed, there are some local regions of high

accuracy. In particular, when looking at the number of features for values $\frac{\ln(M)}{\ln(2)}$ which in this case correspond to a value of approximately 5. It is surprising that the optimal number of features is very low with a value of 2. This may be a result of creating a higher chance of uncorrelated trees which has been shown to improve accuracy. [13]

Table 3: Optimal values found for project’s Random forest by grid search

Target class	Max depth	Num features	Num of trees
GP_greater_than_0	15	2	64
sum_7yr_GP	8	5	64

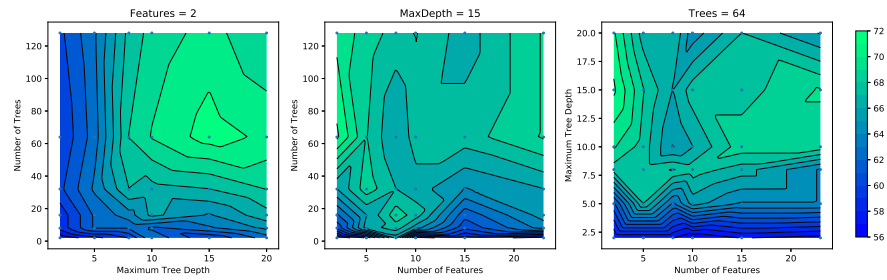


Figure 4: Showing results from grid search with effect on test set accuracy of project’s Random forest implementation

Table 4: Optimal values found for sci-kit learn Random forest by grid search

Target class	Max depth	Num features	Min samples per leaf	Min samples for split	Num of trees
GP_greater_than_0	8	5	2	2	128
sum_7yr_GP	8	5	4	2	64

Comparing the project’s implementation of Random forest to that of scikit-learn implementation, we further compare against the optimal values found by grid search, shown in table 4, and 3. As can be seen, the values are quite different, although result in similar accuracy. A high number of trees is consistent between both showing the over-all effect of the Random forest algorithm on improving performance. The difference between the values can be due to the different algorithms used and shows that finding optimal values for the specific algorithm is an important consideration when running any machine learning algorithm implementation.

3.2.5 Splitting criteria

When building a decision tree, the splitting criteria by which each level is formed to yield the greatest information gain can be calculated by two common methods; Entropy and Gini impurity. It has been discussed in literature that there exists a minimal difference between the effect of accuracy, however Gini impurity function is computationally less expensive than Entropy since Entropy requires additionally calculating the logarithm of percentage of each class labels in the child nodes. The difference between the splitting criteria is investigated through the implementation of both entropy and Gini impurity in our implementation of Random forest, with results shown in table 5. Both models were run on a single process.

As shown in table 5, using the Gini index as the score function is indeed noticeably faster than Entropy.

Table 5: Comparison of Entropy and Gini impurity for splitting criteria

Number of trees	Gini impurity	Entropy
4	1min 23secs	1min 47secs
16	3min 58secs	4min 26secs
32	6min 19secs	7min 23secs

3.2.6 Accuracy Comparisons

I think we can maybe remove this and speed comparisons and just have everything in implementation comparison without subdividing - it becomes difficult to discuss otherwise.

3.2.7 Speed Comparisons

TODO: Show that gini is faster than entropy. Possibly come up with a better subsection title or remove this subsection.

3.3 Machine Learning Algorithm Comparisons

TODO: This section should be comparing our accuracy using best parameters against the best accuracy of other machine learning algorithms (e.g. linear regression, single decision trees, naive bayes, etc.) Possibly come up with a better subsection title.

4 Conclusion

TODO: Summarize our findings and mention future work which is the uncorrelated tree algorithm in the paper I found before.

TODO: Figure out if this paragraph should be removed. Random forests have been shown to achieve high classification performance through ensemble with a set of decision trees that are constructed using randomly selected feature subspaces. The performance of an ensemble learner is dependent on the accuracy of each component learner and the diversity of the components, especially when using a small set of trees which may be limited due to computational cost. The randomization can cause occurrence of bad predicting trees as well as correlated trees which can lead to poor ensemble decisions, which can be observed when performing multiple training runs using the same parameters which can lead to different accuracy results. Attempts have been made to improve the performance of this model by building a forest of only uncorrelated high performing trees. [13]

Contributions

All authors contributed equally.

See GitLab repository here for specific commits:

<https://csil-git1.cs.surrey.sfu.ca/rkm3/mlclass-1777-randomforest>

References

- [1] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: A Simple Way to Prevent Neural Networks from Overfitting," *Journal of Machine Learning Research*, vol. 15, pp. 1929–1958, 2014.
- [2] T. Ho, "Recognition of handwritten digits by combining independent learning vector quantizations," *Proceedings of the Second International Conference on Document Analysis and Recognition*, pp. 818–821, 1993.
- [3] T. Ho, "Random decision forests," *Proceedings of the Third International Conference on Document Analysis and Recognition*, pp. 278–282, 1995.

- [4] T. Ho, "The random subspace method for constructing decision forests," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 20, no. 8, pp. 832–844, 1998.
- [5] D. R. Cutler, T. C. Edwards, K. H. Beard, A. Cutler, T. Kyle, J. Gibson, J. J. Lawler, H. Beard, and T. Hess, "Random Forests for Classification in Ecology," *Ecology*, vol. 88, no. 11, pp. 2783–2792, 2007.
- [6] J. R. Harris and E. C. Grunsky, "Predictive lithological mapping of Canada's North using Random Forest classification applied to geophysical and geochemical data," *Computers and Geosciences*, vol. 80, pp. 9–25, 2015.
- [7] G. Luo and K. Wang, "A combined random forest and active contour-model approach to fully automatic segmentation of the left atrium in volumetric MRI," *BioMed Research International*, vol. 2017, 2017.
- [8] N. Ghatasheh, "Business Analytics using Random Forest Trees for Credit Risk Prediction: A Comparison Study," *International Journal of Advanced Science and Technology*, vol. 72, pp. 19–30, 2014.
- [9] D. Lock and D. Nettleton, "Using random forests to estimate win probability before each play of an NFL game," *Journal of Quantitative Analysis in Sports*, vol. 10, no. 2, pp. 197–205, 2014.
- [10] S. Formann-Roe, "Bias and variance," Jun 2012.
- [11] A. Criminisi, J. Shotton, and E. Konukoglu, "Decision Forests for Classification , Regression , Density Estimation , Manifold Learning and Semi-Supervised Learning," *Microsoft Research technical report*, vol. 7, pp. 81–227, 2011.
- [12] L. Brieman, "Random forests," *Machine Learning*, vol. 45, no. 1, pp. 5–32, 2001.
- [13] S. Bharathidasan and J. C. Venkataeswaran, "Improving Classification Accuracy based on Random Forest Model with Uncorrelated High Performing Trees," *International Journal of Computer Applications*, vol. 101, no. 13, pp. 26–30, 2014.