

Random forest: an ensemble machine learning algorithm

Atiq Ur Rehman* Insha Ullah

Abstract. Random forest is an ensemble (combines multiple models) machine learning tool popular among data scientists. It is simple to understand and requires less efforts in terms of training and tuning. It is flexible because of having no formal distributional assumptions and can handle skewed, multi-modal, categorical (ordinal or non-ordinal), time series and time to event (survival) data. It does not overfit, which is a major issue for most machine learning algorithms. It can handle missing values along the way.

Here we attempt to introduce the random forest algorithm in detail and try to make the document self-contained as much as it is possible. We present two implementations of the algorithm in R packages, namely `randomForest` and `ranger`, and elaborate on them during the analyses of various publicly available datasets.

Key words and phrases: Random Forest, Ensemble Learning, `ranger`, `randomForest`.

1. INTRODUCTION

The rapid evolution of technology has given us the ability to generate, store and process massive amounts of data. A variety of approaches have been developed in the field of machine learning to help humans understand the important trends and disparities in the data. These approaches can be broadly categorized as supervised and unsupervised learning techniques. The supervised learning techniques are supervised because of their requirement of pre-labelled data for training purposes. In more simple terms, one needs to show some examples so that the technique learns the pattern first; only then it will be able to predict the labels for similar unseen data. On the other hand, the unsupervised learning techniques discover the labels for the training data on their own; therefore, they do not require pre-labelled training data. These algorithms learn by discovering hidden groups in the unlabelled training data and can utilize that learning to label new (previously unseen) observations based on their similarities to the identified groups. We concern ourselves with supervised learning for the rest of this paper and do not discuss unsupervised learning any further.

Supervised learning can be further categorized into classification and regression tasks. In a classification task, the response variable is categorical with a fixed number of categories, $K \geq 2$. The objective is to classify an instance in one of the K classes based on some known features of that instance. A special and commonly encountered case is when the response variable has two possible classes is known as a binary classification problem. If the response variable has more than two categories, the classification problem is a multi-class classification problem. In a

*To whom correspondence should be addressed.

regression task, on the other hand, the response variable is numerical/continuous. It is used to quantify the relationship between a numerical response variable and features that explain the response variable, thereby enabling us to understand how input features affect the response variable of interest. Additionally, it allows us to predict the value of a numerical variable based on some known features.

Various supervised learning methods have been developed. These methods include, among others, linear regression, logistic regression, random forest, support vector machine, naive Bayes, k-nearest neighbours, and neural network. Some may be more suitable in one application than the other, depending on the type of the response variable and its relationship with the available features. For example, linear regression may be used to predict a numerical response variable when its relationship with the features is linear. On the other hand, the random forest can deal with non-linear relationships between a response variable and features.

Random forest is a machine learning tool for supervised learning that can perform regression and classification tasks (Breiman, 2001). It is one of the most frequently used algorithms due to its simplicity. Random forest is an ensemble method that combines a group of models (decision trees) to reduce variance and improve precision. First, the random forest has the advantage and outperforms the linear regression in a way that it can be used in both cases whether the relation between a continuous outcome variable and features is linear or non-linear. Second, it is mainly suitable for the problem of overfitting, that is, it can fit both training and testing data well. Thirdly, it can be used in both cases of large as well as a small datasets. Lastly, random forest is good at handling missing values automatically and unlike most supervised learning methods it is usually robust to outliers.

To our knowledge, there is no article available that comprehensively explains the random forest. Moreover, there two R packages available mostly commonly used, that is, randomForest (Liaw et al., 2002) and ranger (Wright & Ziegler, 2017). In this tutorial, we will explain in more detail random forest algorithm for both classification and regression. We introduce two implementations of the algorithm in R packages, namely randomForest and ranger, and elaborate on them during the analyses of various publicly available datasets.

2. RANDOM FOREST

Decision trees, building block of the random forest, can be used for both classification and regression to predict response variable based on number of features, for further detail see Krzywinski & Altman (2017). However, Decision trees have the problem of overfitting and they are unstable with small change in data as it can lead to the estimation of a completely different tree. To circumvent this issue, Breiman (1996) has introduced an ensemble method called bagging. The most important advantage of this method is that it can perform very well in which the training dataset is sampled many times with replacement called Bootstrap samples, and decision trees are fitted to each Bootstrap sample. Following this, the prediction of all decision trees are averaged (in case of regression) or taking a majority vote (in case of classification) to get the overall prediction. This ensemble method has the advantage of reducing the variance and avoid overfitting. But the downside of bagging is that each decision tree uses all predictors for the splits. In this respect, only important variables have a higher chance of being selected

for the splits. As a consequence, we will get similar results from decision trees. In other words, the decision trees are correlated to one another. Thus, there will be no performance gains. One way to overcome this problem is to use another ensemble method that has been proposed by Breiman (2001) known as random forest, a variant of bagging, which selects a subset of predictors for each splits in the decision tree explained below.

The random forest algorithm is an ensemble of decision trees that is used for supervised learning. It can perform both a classification task as well as a regression task; that is, it performs classification when the response variable is categorical/factor/qualitative and regression when the response variable is a continuous/numerical/quantitative variable. As its name hints, it makes a forest by growing many individual decision trees, each of which is built on a random subset of data and a random subset of distinct features. The predictions are then made by aggregation across the entire forest. In the case of classification, the aggregation means majority vote, that is, an instance is classified in a particular class if the majority of the trees suggests doing so. In the case of regression, the aggregation reduces to averaging the prediction of individual regression trees.

A tree in a random forest is built by starting with the entire sample in the root node and recursively partitioning it into nodes in such a way so that instances in a node are homogeneous concerning the response variable. The process continues until the terminal nodes are reached. Note that the best homogeneous partitioning of variable and selection of variable use the same rules as of decision trees.

Each tree is built on a random subset of data and a random subset of features at each node of tree is the fact that makes random forest different from other ensemble algorithms. The major advantage of choosing a subsample of features is that there may be few important variables in the whole set of features; therefore, considering all variables for each tree (as in case of bagging) of the random forest will only select the important variables and leave the remaining less important variables to be selected for each node. Taking that into consideration, the results of trees will be similar to one another and thus correlated trees. Therefore, averaging similar results from highly correlated trees will not reduce variance. Selecting a subset of features without replacement for each split will also give chance to less important variables and because of that different trees in a random forest will give different results in other words the trees will be uncorrelated. The feature randomness the so-called feature bagging ensures the low correlation among trees, and thus more accurate results compared to the single tree model. Moreover, using all predictors for the splits is computationally expensive, particularly, in high-dimensional data.

The number of variables to be selected for each split, which is denoted by *mtry* in most of the R packages should be chosen adaptively from the data. The default value of *mtry* in case of classification is \sqrt{p} (Bernard et al., 2009), where p represents total number of predictors, they also concluded that it can be improved. In case of regression the default optimal *mtry* is $\frac{p}{3}$. It is also been shown that the default values of hyperparameters is reasonable but not optimal (Huang & Boutros, 2016). In practice the best values for *mtry* will depend on the data at hand. For example, there may be more relevant/important predictors in the data in that case the *mtry* should be small. As such, the less important variables will also have chance to be selected for the splits. However, it is also important

to note that too small *mtry* can prevent variance from reduction (Goldstein et al., 2011). In addition to that in small subset of variables there may be situations where none of the features are relevant and due to that reason irrelevant feature selected for splitting will have poor predictive ability. On the other hand, in case of few relevant variables high *mtry* is preferred (Goldstein et al., 2011). Thus, *mtry* is most important tuning parameter of random forest and its optimal value depends on the data problem. To our knowledge, there is no exact optimal value of *mtry*; therefore, it is advisable to select optimal value via greed search.

Another hyperparameter in random forest is the minimum number of observations in the terminal node usually denoted by *nodesize*. This parameter value set the depth of the decision trees in a random forest. In other words, smaller value of *nodesize* can lead to larger trees thus more splits until leaf node and vice versa. In most R packages, the default value of *nodesize* is 1 for classification and 5 for regression. Its value can be selected through greedy search. It is important to state that in machine learning algorithms there are two types of parameters: model parameters that can be estimated from data automatically, and model hyperparameters that are needed to be set manually via grid search.

In case of classification, there are situations where the classes of response variable are imbalanced, that is, the samples are not equally distributed (e.g see section 3 where our data is imbalanced). A class with large amount of data/samples is known as majority class and fewer observations for other class is called minority class. Therefore, it is more challenging for random forest to learn the characteristics of the minority class and thereby the model will have poor predictive performance. These situations are very common in classification and often the minority class is of interest such as anomaly detection, spam detection etc. To solve imbalance class problem, we can use resampling methods such as undersampling the majority class, oversampling minority class or applying both simultaneously. There is also another method known as class weights, in that respect the classes should be weighted by giving more weight to the under-represented class and less weight to over-represented class. These methods are available in R packages such as randomForest (Liaw et al., 2002) and renger (Wright & Ziegler, 2017).

The results from a random forest algorithm may vary across different packages and even across different replications using the same implementation. The reason is that a random forest is algorithmic in nature and does not develop prediction equation as is done in linear regression.

2.1 Variable Importance

One advantage of the random forest is that it can find whether features are strongly or poorly linked to the response variable. Therefore, random forest can also be used for feature selection. To examine how relevant the predictors are that are used to estimate the random forest model. There are two measures of importance given for each variable in the random forest. The first measure is based on how much the accuracy decreases when the variable is excluded from the features set known as permutation importance. The second measure a.k.a Gini importance index is based on the decrease of Gini impurity when a variable is chosen to split a node.

The permutation importance is based on the OOB samples. These OOB samples are most common in machine learning methods based on bagging. Random

forest is also based on the bootstrap aggregation. In random forest, each tree is based on the bootstrap sample from the original training data. With bootstrap sample, another sample is also created that is not part of the bootstrap sample and thus not used in constructing a decision tree called OOB samples. These OOB sample are very important in a way that they are used to measure the prediction error of random forest. Calculating prediction with OOB sample also known as OOB error is computationally less expensive and memory efficient compare to the cross-validation. This is because in case of k-fold cross validation k random forest are needed be constructed while in case of OOB error random forest is estimated once for constructing a model and error estimation.

In permutation importance method ([Breiman, 2001](#)) a variable is considered as important if it has positive effect on the prediction performance. In other words, the more positive effect of the predictor is relevant the variable would be. To compute the permutation importance of a predictor, we first train random forest model on original predictors and obtain OOB error of the trained model. After that we take one predictor variable and shuffle its samples keeping other predictors and outcome variable in place. Now the dataset is changed, we again fit the random forest model and obtain the OOB error. Then we find the difference between the OOB error of the original data and permuted-one-predictor data. The difference is the permutation importance for that one predictor being shuffled. Thus the process is repeated for all the predictors. Feature with high positive permutation importance score is the most relevant variable.

Another important mostly commonly used method for variable importance is Gini importance index. The gini importance index of a single feature is the sum of all impurity reduction of all nodes where the feature is considered the best for splitting normalizing by the number of trees being used. In simple words, we first take sum of the impurity reduction of a feature over all nodes in which the feature is used as a best split in a single tree of a random forest. After that the total impurity score is then averaged over all trees in the random forest.

2.2 The ranger and randomForest Packages in R

The ranger package is implemented in C++ is the fastest and memory-efficient R package compare to the other packages of random forest in R ([Wright & Ziegler, 2017](#)) designed for high-dimensional data. [Wright & Ziegler \(2017\)](#) compared the runtime and memory usage of all possible R packages developed before 2017 and found out that ranger is the fastest and more memory efficient. Moreover, as the trees are independent of one another the ranger package can train multiple trees simultaneously depend on the number of CPUs available (also known as parallel computing), which is not true in randomForest R Package ([Liaw et al., 2002](#)). Thus, that also makes ranger faster compared to the randomForest. The major advantage of ranger package is that we can estimate other types of random forests such as survival and quantile regression random forests, which is not true in case of randomForest package.

In function `ranger()`, if the response variable is a factor it grows classification forest otherwise regression. If the response variable is numeric and we want to fit the classification forest then argument “classification” should be set TRUE as default is FALSE.

In case of imbalanced data where training samples are unevenly distributed

among classes. Also for classification, we will use imbalanced data mentioned in the next section. The problem is that the decision trees are sensitive to imbalanced data and the random forest is the combination of decision trees. Therefore, random forest is also sensitive to imbalanced data. To make balanced data either the majority class can be undersampled or the minority class can be oversampled. This facility is not available in ranger but it can be done manually and a list strata equal to the number of trees can be provided to an argument “inbag” shown in the R codes below. However, this can be directly done in the randomForest package using “strata” argument in randomForest function. Another argument that can be used for class imbalanced data is “class.weight” in ranger and “classwt” in randomForest, but calculation of class weights is a bit tricky and should be carefully chosen.

The OOB error for each class and overall of individual trees cannot be obtained directly from the ranger but can be obtained in the randomForest package. However, in ranger it is possible to calculate it manually from OOB samples misclassified in each individual tree.

3. DATA AND RESULTS

For classification, we demonstrate random forest using a dataset from the Kaggle competition named “Don’t Overfit! II”¹. This competition is completed two years ago and the new version of the data for competition is uploaded, but we are using the same old version dataset. Kaggle, a subsidiary of Google, is a data science community of data scientists all over the world. A community where different resources and tools are available and it allows users to upload and use those datasets to build machine learning models. It allows companies and individuals to host a competition of data science problems and data scientists are allowed to compete solving those problems. In this full dataset, three csv files are given, that is, train.csv, test.csv and sample_submission.csv. The train.csv file contains the training data with the number of variables exceeds the total number of observations; that is, there are 250 observations and 300 continuous variables plus one binary response variable. The test.csv file contains test data contain 19,750 observations with the same number of predictors. The sample_submission.csv contains column id, which is an id of each observation and column named target with predictions of a test data. The sample_submission.csv file is given for a reason that final predictions of test data on a trained model (in our case random forest) will be submitted in that fashion. The dataset is equation-based simulated data. The problem with this dataset is that most models overfit it due to the small training samples relative to the number of variables and noise in the data. It is a binary classification problem, that is, the response variable is binary with two categories as namely “0” and “1” shown in figure 1 with 90 and 160 samples associated with each class, respectively. We train a random forest algorithm and predict classes for the 20,000 samples in the test data.

We are using R version 4.0.3 with two common random forest packages namely ranger version 0.12.1 (Wright & Ziegler, 2017) and randomForest 4.6.14 (Liaw et al., 2002) to estimate random forest model. The algorithm of random forest in these two packages are different as randomForest package uses Breiman (2001) and in ranger package a lot of methodological choices have been made with speed

¹<https://www.kaggle.com/c/dont-overfit-ii/data>

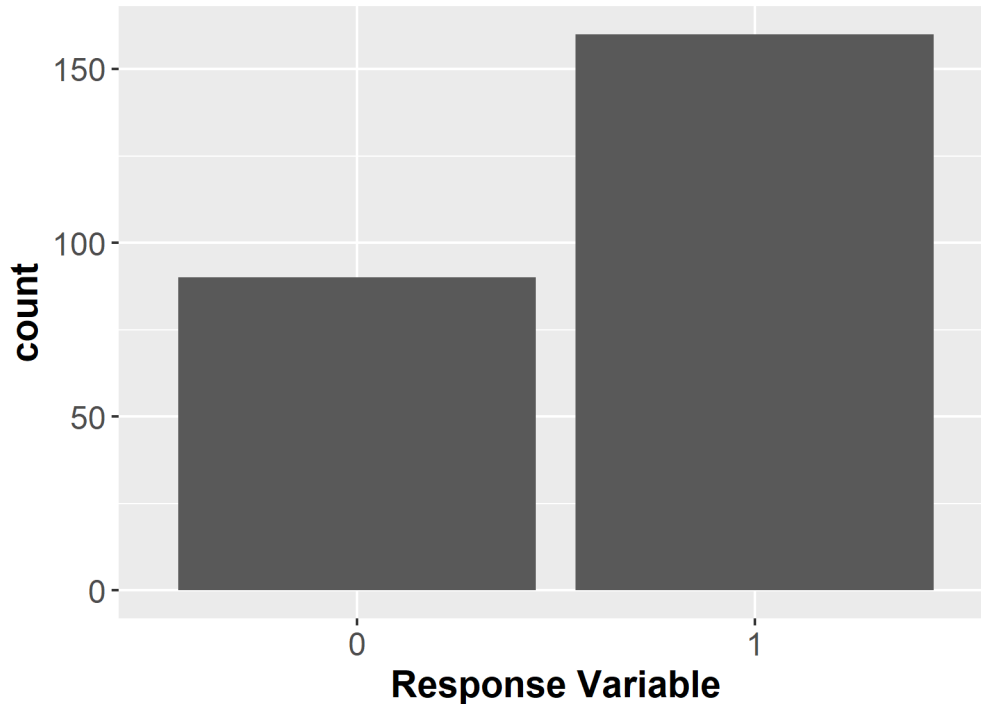


Figure 1: Binary response variable with 90 sample belongs to minority class and 160 samples belong to majority class.

in mind ([Wright & Ziegler, 2017](#)). We want to evaluate random forest model on both of these packages.

There are different tuning parameters described in the previous section that needs to be optimized to best the data at hand. For example, the number of trees to be used in the random forest. A random forest uses a subset of features to estimate on a decision tree. Thus if the number of features is large and the number of trees used in the random forest is small meaning that some features will be missed in a subset of features used. As a consequence, the predictive performance of the random forest will also be reduced. However, increasing the number of trees in the random forest does not overfit the model (Friedman et al., 2001) and in this respect, all the features will be used to estimate the model. But caution must be exercised as a large number of trees sometimes can increase the computational time instead of increasing the performance of the model. As we are allowed in number of trees to use in the random forest. Thus, from that we can conclude that it is unnecessary to tune number of trees for the random forest. Other hyperparameters of the random forest are *mtry* and *min.node.size*, that need to be tuned and get the best optimal values in order to achieve minimum prediction error. The hyperparameters optimal values cannot be estimated from the data and they must be set before training a model. Hyperparameters tuning depends more on experimental results than theory, and thus the best method to determine the optimal settings is to try many different combinations to evaluate the performance of each model. Therefore, a grid search, that is, a range of value of the hyperparameters can be used and on each combination random forest

model should be fitted. After that a combination that gives minimum predictive error would be our optimal values of hyperparameters. In our case, we used *mtry* from 5 to 30 with increment of 3 and *min.node.size* (1-5) and estimated random forest models with number of trees equal 3500. The results are shown in figure 2 (a) with Out-of-Bag error plotted on y-axis and *mtry* on x-axis for different *min.node.size* values. It can be seen that *mtry* = 29 and *min.node.size* = 2 gives us the minimum predictive error compare to the other combinations.

Thus, using the optimal hyperparameters values with *mtry* = 29 and *min.node.size* = 2, we have fitted the final random forest model. The total number of trees used in random forest are 3500. The bootstrap sampling is done without replacement. Moreover, for variable importance, we have used permutation feature importance, and the splitting rule to select a variable for each node in a tree is the Gini index. After fitting the model, we have use the trained model to predict the test data and achieved 0.768 score. The score can still be improved by using various strategies, we have tried two approaches in order to improve the score as in the following paragraphs.

As the training data is imbalanced, that is, 160 samples belong to class “1”, and 90 samples belong to class “0” shown in figure 1. In that case, the majority class get pretty high accuracy in prediction compare to the minority class. Therefore, the traditional classification algorithms with the correct classification rates may bias towards the majority classes (Lu et al., 2019). Various methods are available to make the imbalanced data balanced such as undersampling and oversampling etc. We carry out undersampling with replacement in both over and under represented classes by eliminating samples. In other words, we take 60 out of 160 samples from the majority class 1 and 60 out of 90 samples from minority class 0 with replacement. Making the training data balanced, now we find the optimal hyperparameters values via grid search. The *mtry* ranges from 5 to 30 with difference of 3 and *min.node.size* from 1 to 5. From figure 2 (b), it can be seen that *mtry* = 29 and *min.node.size* = 3 is the optimal values, which gives us the least OOB error. Training random forest model with those hyperparameters values, We predicted the test data on the trained model and achieved score of 0.778.

Simplifying feature set can be useful in some situations, especially, when the sample size is small as our dataset. The reason is that irrelevant variables in a data can have minimal effect on model prediction performance. Moreover, By using those redundant features can also degrade the model predictions by increasing variance. It should be also noted that random forests are particularly good at feature selection as they internally repeat the simple decision tree feature selection algorithm such as “gini” and “permutation” over bootstrapped copies of the data set. But due to randomness, the random forest model still use those irrelevant variables for splitting and thereby increase variance of the model as well as computational time. In our dataset, There are some non-informative variables and selection of those variables for decision trees in the random forest makes weak learner and reduce model performance. Therefore, they need to be removed before training random forest. For that reason, we perform feature engineering by applying a t-test and take 16 significant features with p-values less than 0.02. We also estimate a random forest model with 300 features and take eight important features apart from the first 16 predictors. Also we have taken 50 samples from

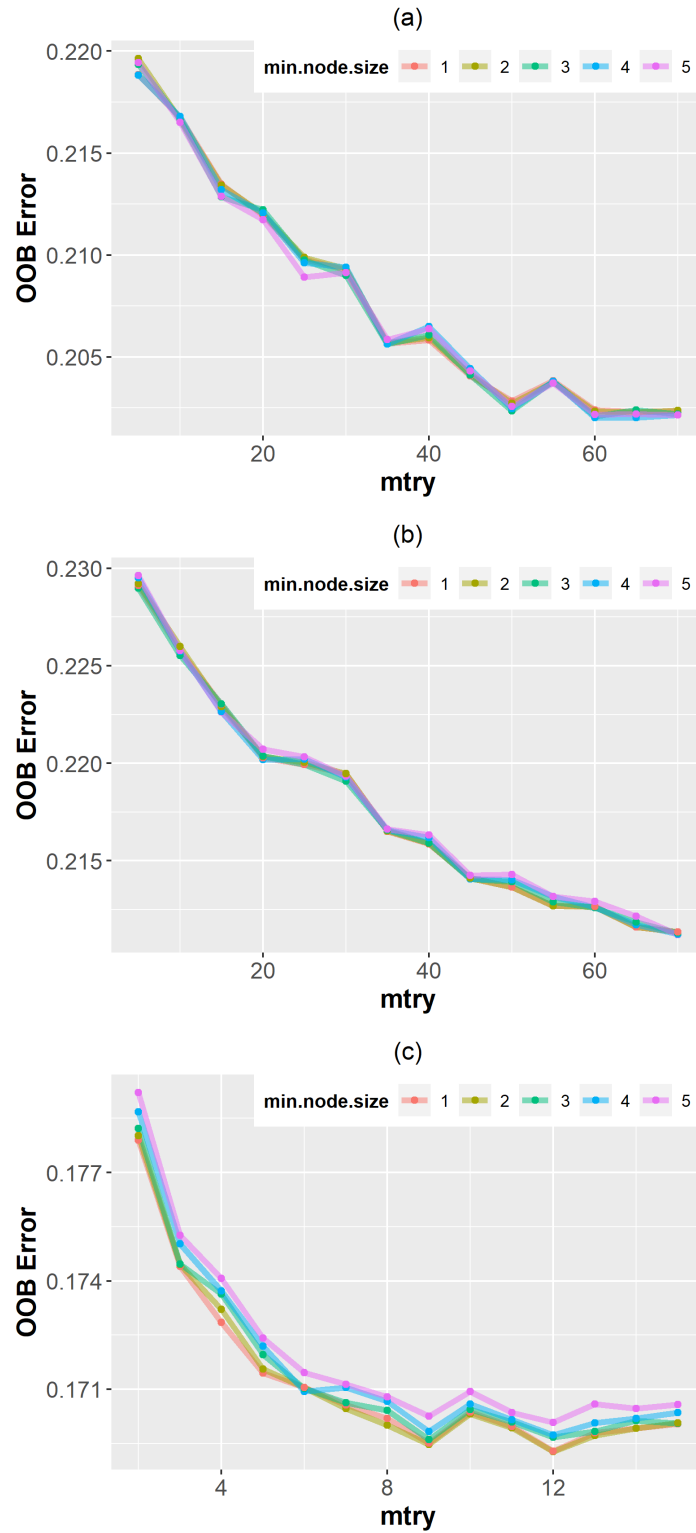


Figure 2: Each subfigure shows the grid search of Out-of-Bag error against $mtry$ and $min.node.size$ of the three random forest model: (a) Grid search for unbalanced training dataset, (b) Grid search for balanced data with 60 samples from each class, and (c) Grid search for data with 24 important variables and balanced with 50 samples from each class.

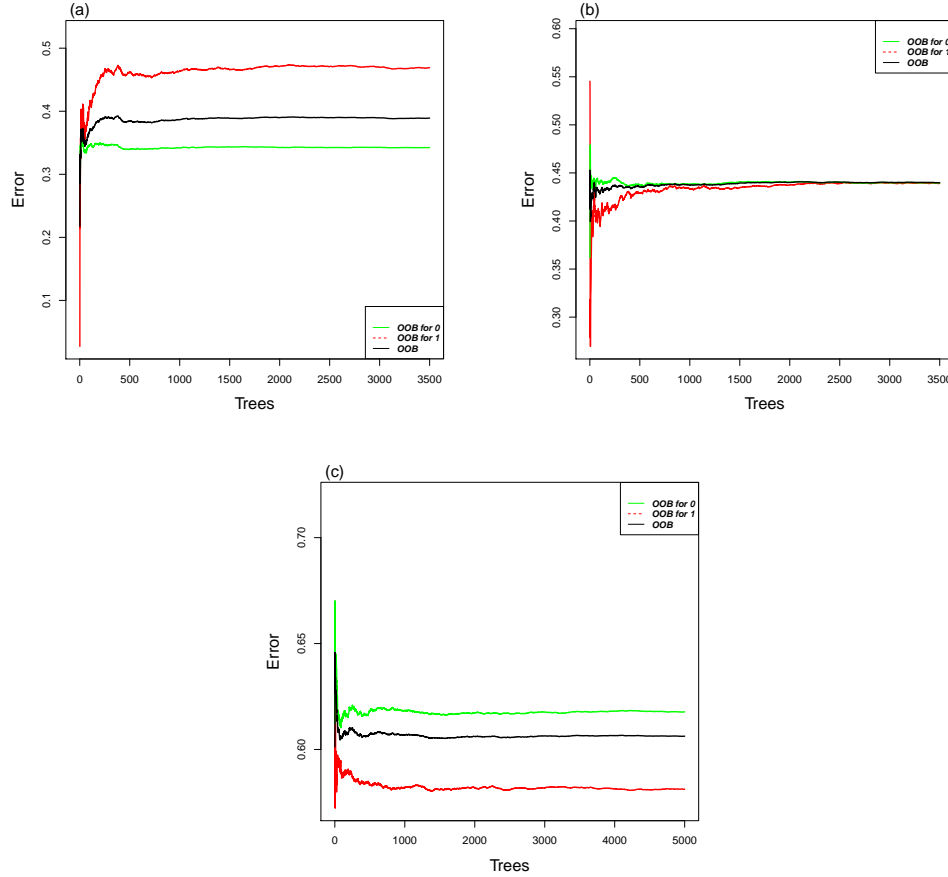


Figure 3: Each subfigure shows the Out-of-Bag error of the three random forest model used to improve the score: (a) OOB error plot of original data, (b) OOB error plot of balanced data with 60 samples from each class, and (c) OOB error plot of data with 24 important variables and balanced with 50 samples from each class

each class with replacement to make it balanced dataset. Thus, we are left with 24 important features explaining the response variable. For selection of optimal combination of hyperparameters, a range of $mtry$ and $min.node.size$ has been used to estimate various random forest models. The predictive error is shown in figure 2 (c), it can be seen that with $mtry = 5$ and $min.node.size = 1$ gives us minimal OOB error. Using the optimal combination of hyperparameters, we have estimated the random forest model. After training random forest model with 5000 number of trees, predicting test data, we achieved 0.805.

4. DISCUSSION AND RECOMMENDATIONS

We have focused on both explaining the technique and demonstrating how to fit and evaluate the models, because this information is presented elsewhere in rather technical language, and sometimes portrays boosting as a ‘black-box’ procedure.

Random forest is a variant of bagging as $mtry < p$, where p denotes the num-

ber of features. When $mtry = p$ random forest becomes bagging; in contrast, when $n_{tree} = 1$ and $mtry = p$, random forest becomes a single decision tree. In random forest using more samples reduces the variance of the fit, but because many bootstrap samples are similar, at some point more bootstraps will merely increase computation time without improving the estimates. Because the increase in number of trees is linearly related to the computational time. There is no fixed value/number for random subset selection of variables for each split in decision trees of random forest. Both small and large number of $mtry$ can be problematic in different situations. However, the default value of $mtry$ is 1 for classification and 5 for regression in R packages ranger (Wright & Ziegler, 2017) and randomForest (Liaw et al., 2002) and reported often to be an appropriate choice (Liaw et al., 2002) but there is no solid theory behind it. Therefore, it is advisable to choose optimal value through a greed search. The same method of greed search should also be done for minimum node size in each tree. Random forest can be applied to binary or N-class ($N > 2$) prediction classification problems, to model interactions among variables, can take on a mixture of categorical and continuous features, and can be used to measure the variable importance.

From the above data used, we have also learnt that random forest can be used and works quite well when the number of variable is large compare to the number of observations. But caution must be exercised when sample size is very small because taking more bootstrap samples (equal to number of trees in random forest) would not change the ability of random forest to learn something new. Secondly, the the depth of the tree would be very small.

The interesting fact about random forest is that there no need to use k-fold cross-validation for model training and validation. It uses that OOB samples to measure the prediction error of the model known as OOB error. Thus it makes random forest computationally less expensive and memory efficient compare to the k-fold cross-validation. Besides, random forest is nonparametric method thus immune to statistical assumptions. It can automatically handle categorical features, that is, no dummy coding is required before training a random forest model. One disadvantage of the random forest algorithm is lack interpretability. In other words, random forest cannot be used in finding relationship and strength of features with response variable.

Putting flash on random forest throughout the article, we have learnt that random forest does overfit the data. But this does not mean that we do not agree with Breiman who introduced random forest claims that it does not overfit². This is true as long as the tuning parameters is chosen correctly. For example, one hyperparameter is maximum depth of the trees in random forest that avoids overfitting and need to be tuned.

5. CONCLUSION

In this article, we walked through the concept of random forest based on the decision trees and bagging. We have found that random forest is a useful algorithm for both classification and regression and it combines group of weak learners (decision trees) to make strong learner. It is a non-parametric method with no distributional assumptions and can handle any type of data that makes random forest to be applied to a variety of prediction problems. It has various tuning

²https://www.stat.berkeley.edu/~breiman/RandomForests/cc_home.htmremarks

parameters and needs to be selected via grid search. Due to the OOB samples as validation samples, random forest is fast and memory efficient. Another advantage of the random forest is that it does not overfit as long as tuning parameters are chosen appropriately. Furthermore, a common feature of the random forest is the ability of finding the variable importance internally; therefore it can be used in finding most important features.

We have used two common R packages, that is, ranger and randomForest to fit random forest on a dataset. We concluded that both have their own advantages on different occasions. For example, the ranger is faster, memory efficient, and specially made for high-dimensional data. Tuning hyperparameters in case of big data is computationally expensive; therefore, using ranger can make things faster. On the other hand, As random forest is sensitive to imbalanced data. Therefore, imbalanced data should be made balanced and it could be done directly in randomForest function, but in ranger, it can be done manually. Also, the randomForest function provides us OOB error of each class in each tree but is not true in ranger.

ACKNOWLEDGEMENTS

APPENDIX

REFERENCES

- Bernard, S., Heutte, L., & Adam, S. (2009). Influence of hyperparameters on random forest accuracy. In *International workshop on multiple classifier systems* (pp. 171–180).
- Breiman, L. (1996). Bagging predictors. *Machine learning*, 24(2), 123–140.
- Breiman, L. (2001). Random forests. *Machine learning*, 45(1), 5–32.
- Goldstein, B. A., Polley, E. C., & Briggs, F. B. (2011). Random forests for genetic association studies. *Statistical applications in genetics and molecular biology*, 10(1).
- Huang, B. F., & Boutros, P. C. (2016). The parameter sensitivity of random forests. *BMC bioinformatics*, 17(1), 1–13.
- Krzywinski, M., & Altman, N. S. (2017). Classification and regression trees. *Nature Methods*, 14(8), 757–758.
- Liaw, A., Wiener, M., et al. (2002). Classification and regression by randomforest. *R news*, 2(3), 18–22.
- Lu, H., Xu, Y., Ye, M., Yan, K., Gao, Z., & Jin, Q. (2019). Learning misclassification costs for imbalanced classification on gene expression data. *BMC bioinformatics*, 20(25), 1–10.
- Wright, M. N., & Ziegler, A. (2017). Ranger: A fast implementation of random forests for high dimensional data in c++ and R. *J. Stat. Softw.*, 77(1). Retrieved from <http://dx.doi.org/10.18637/jss.v077.i01> doi: