# Cost-complexity pruning of random forests

B Ravi Kiran[1] and Jean Serra[2]

[1]CRIStAL Lab, UMR 9189, Université Charles de Gaulle, Lille 3
`kiran.ravi@univ-lille3.fr`
[2]Université Paris-Est, A3SI-ESIEE LIGM `jean.serra@esiee.fr`

July 20, 2017

**Abstract**

Random forests perform boostrap-aggregation by sampling the training samples with replacement. This enables the evaluation of out-of-bag error which serves as a internal cross-validation mechanism. Our motivation lies in using the unsampled training samples to improve each decision tree in the ensemble. We study the effect of using the out-of-bag samples to improve the generalization error first of the decision trees and second the random forest by post-pruning. A prelimiary empirical study on four UCI repository datasets show consistent decrease in the size of the forests without considerable loss in accuracy. [1]

**Keywords**:Random Forests, Cost-complexity Pruning, Out-of-bag

# 1 Introduction

Random Forests [5] is an ensemble method which predicts by averaging over multiple instances of classifiers/regressors created by randomized feature selection and bootstrap aggregation (Bagging). The model is one of the most consistently performing predictor in many real world applications [6]. Random forests use CART decision tree classifiers [2] as weak learners. Random forests combine two methods : Bootstrap aggregation [3] (subsampling input samples with replacement) and Random subspace [11] (subsampling the variables without replacement). There has been continued work during the last decade on new randomized ensemble of trees. Extremely randomized trees [9] where instead of choosing the best split among a subset of variables under search for maximum information gain, a random split is chosen. This improves the prediction accuracy. In furthering the understanding of random forests [7] split the training set points into structure points: which decide split points but are not involved in prediction, estimation points: which are used for estimation. The partition into two sets are done randomly to keep consistency of the classifier.

Over-fitting occurs when the statistical model fits noise or misleading points in the input distribution, leading to poor generalization error and performance. In individual decision tree classifiers grown deep, until each input sample can be fit into a leaf, the predictions generalizes poorly on unseen data-points. To handle this decision trees are pruned. There has been a decade of study on the different pruning methods, error functions and measures [14], [16]. The common procedure follow is : 1. Generate a set in "interesting trees", 2. Estimate the true performance of each of these trees, 3. Choose the best tree. This is called post-pruning since we grow complete decision trees

---

[1]Previous version in proceedings ISMM 2017.

and then generate a set of interesting trees. CART uses cost-complexity pruning by associating with each cost-complexity parameter a nested subtree [10].

Though there has been extensive study on the different error functions to perform post-pruning [13], [17], there have been very few studies performed on pruning random forests and tree ensembles. In practice Random forests are quite stable with respect to parameter of number of tree estimators. They are shown to converge asymptotically to the true mean value of the distribution. [10] (page 596) perform an elementary study to show the effect tree size on prediction performance by fixing minimum node size (smaller it is the deeper the tree). This choice of the minimum node size are difficult to justify in practice for a given application. Furthermore [10] discuss that rarity of over-fitting in random forests is a claim, and state that this asymptotic limit can over-fit the dataset; the average of fully grown trees can result in too rich a model, and incur unnecessary variance. [15] demonstrates small gains in performance by controlling the depths of the individual trees grown in random forest.

Finally random forests and tree ensembles are generated by multiple randomization methods. There is no optimization of an explicit loss functions. The core principle in these methods might be interpolation, as shown in this excellent study [18]. Though another important principle is the use of non-parametric density estimation in these recursive procedures [1].

In this paper we are primarily motivated by the internal cross-validation mechanism of random forests. The out-of-the-bag (OOB) samples are the set of data points that were not sampled during the creation of the bootstrap samples to build the individual trees. Our main contribution is the evaluation of the predictive performance cost-complexity pruning on random forest and other tree ensembles under two scenarios : 1. Setting the cost-complexity parameter by minimizing the individual tree prediction error on OOB samples for each tree. 2. Setting the cost-complexity parameter by minimizing average OOB prediction error by the forest on all training samples.

In this paper we do not study ensemble pruning, where the idea is to prune complete instances of decision trees away if they do not improve the accuracy on unseen data.

## 1.1 Notation and Formulation

Let $Z = \{\mathbf{x}^i, y^i\}_N$ be set of $N$ (input, output) pairs to be used in the creation of a predictor. Supervised learning consists of two types of prediction tasks : regression and classification problem, where in the former we predict continuous target variables, while in the latter we predictor categorical variables. The inputs are assumed to belong to space $X := \mathbb{R}^d$ while $Y := \mathbb{R}$ for regression and $Y := \{C_i\}_K$ with $K$ different abstract classes. A supervised learning problem aims to infer the function $f : X \to Y$ using the empirical samples $Z$ that "generalizes" well.

Decision trees fundamentally perform data adaptive non-parametric density estimation to achieve classification and regression tasks. Decision trees evaluate the density function of the joint distribution $P(X, Y)$ by recursively splitting the feature space $X$ greedily, such that after each split or subspace, the $Y$s in the children become "concentrated" or in some sense well partitioned. The best split is chosen by evaluating the information gain(chage in entropy) produced before and after a split. Finally at the leaves of the decision trees one is able to evaluate the class/value by observing the subspace (like the bin for histograms) and predicting the majority class respectively [8].

Given a classification target variable with $C_k = \{1, 2, 3, ...K\}$ classes, we denote the proportion of of class $k$ in node as :

$$\hat{p}_{tk} = \frac{1}{n_t} \sum_{\mathbf{x}_i \in R_t} I(y_i = k) \tag{1}$$

which represents proportion of classifications in node $t$ in decision region $R_t$ with $n_t$ observations. The prediction in case of classification is performed by taking the majority vote in a leaf, i.e. $\hat{y}_i = \text{argmax}_k \hat{p}_{tk}$. The misclassification error is given by :

$$l(y, \hat{y}) = \frac{1}{N_t} \sum_{i \in R_t} I(y_i \neq \hat{y}_i) = 1 - \hat{p}_{mt} \tag{2}$$

The general idea in classification trees is the recursive partition of $\mathbb{R}^d$ by axis parallel splits while maximizing the gini coefficient :

$$\sum_{k \neq k'} \hat{p}_{mt} \hat{p}_{mt'} = \sum_{k=1}^{K} \hat{p}_{mt}(1 - \hat{p}_{mt}) \tag{3}$$

In a decision split the parameters are the split dimension denoted by $j$ and the split threshold $c$. Given an input decision region $S$ we are looking for the dimension (here in three dimensions) that minimizes entropy. Since we are splitting along $d$ unordered variables, there are $2^{d-1} - 1$ possible partitions of the $d$ values into two groups (splits) and is computationally prohibitive. We greedily decide the best split on a subset of variables. We apply this procedure iteratively till the termination condition.
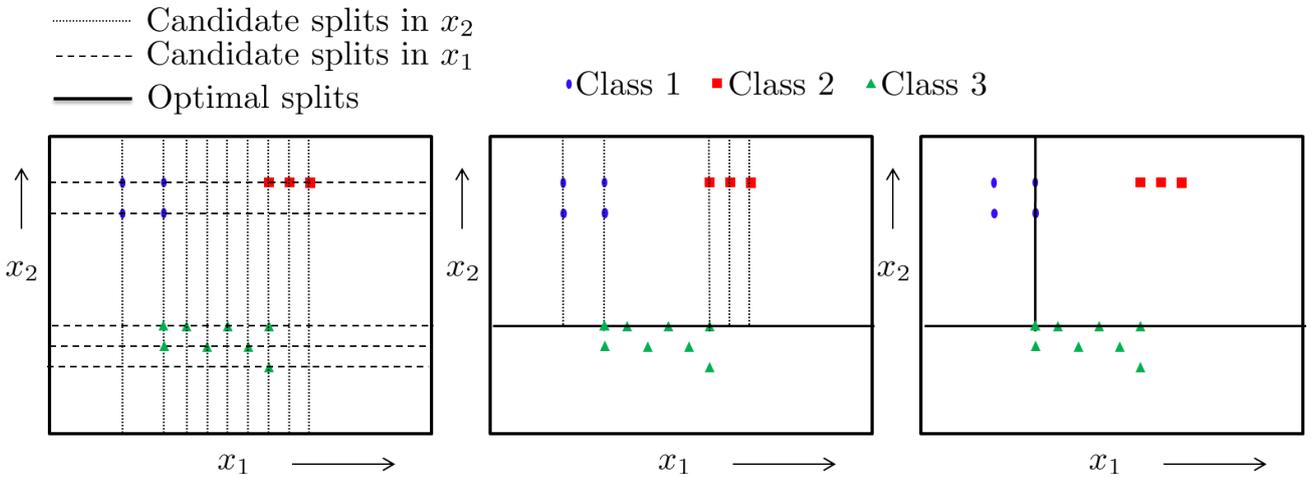


Figure 1: Choosing the axis and splits. Given $\mathbb{R}^2$ feature space, we need to chose the best split given we chose a single axis. This consists in choosing all the coordinates along the axis, at which there are points. The optimal split is one that separates the classes the best, according to the impurity measure, entropy or other splitting measures. Here we show the sequence of two splits. Since there are finite number of points and feature pairs, there are finite number of splits possible.

As shown in figure 1 the set of splits over which the splitting measure is minimized is determined by the coordinates of the training set points. The number of variables or dimension $d$ can be very large (100s-1000 in bio-informatics). Most frequently in CART one considers the sorted coordinates and from them the split points where the class $y$ change and finally one picks the split that minimizes the purity measure best.

## 1.2 Cost-Complexity Pruning

The decision splits near the leaves often provide pure nodes with very narrow decision regions that are over-fitting to a small set of points. This over-fitting problem is resolved in decision trees by performing pruning [2]. There are several ways to perform pruning : we study the cost-complexity pruning here. Pruning is usually not performed in decision tree ensembles, for example in random forest since bagging takes care of the variance produced by unstable decision trees. Random subspace produces decorrelated decision tree predictions, which explore different sets of predictor/feature interactions.
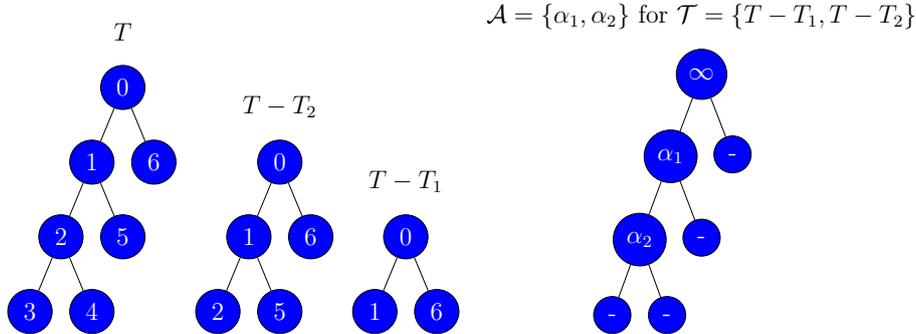
Figure 2: Figure shows shows a sequence of nested subtrees $\mathcal{T}$ and the values of cost-complexity parameters associated with these subtrees $\mathcal{A}$ calculated by equation (6) and algorithm (2). Here $\alpha_2 < \alpha_1 \implies T - T_1 \subset T - T_2$
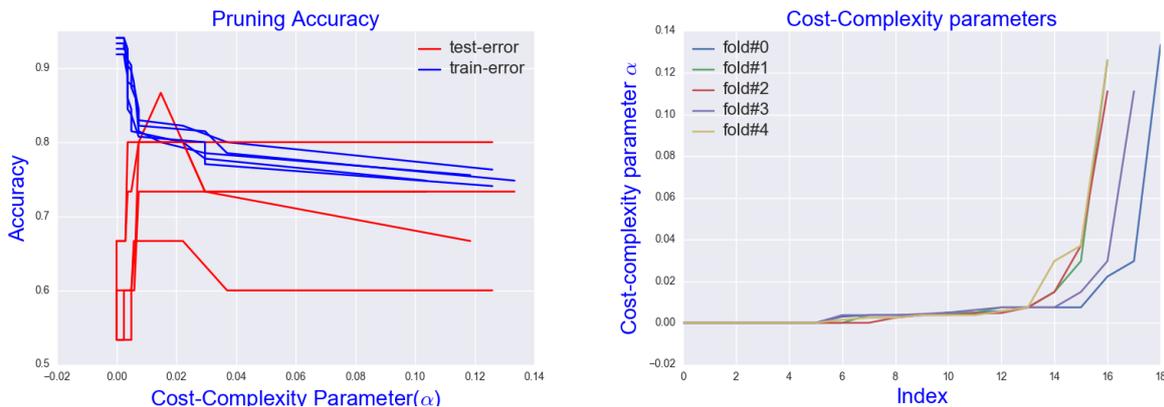


Figure 3: Training error and averaged cross-validation error on 5 folds as cost-complexity parameter varies with its index. The index here refers to the number of subtrees.

The basic idea of cost-complexity pruning is to calculate a cost function for each internal node. An internal node is all nodes that are not the leaves nor the root node in a tree. The cost function is given by [10]:

$$R_\alpha(T) = R(T) + \alpha \cdot |\text{Leaves}(T)| \tag{4}$$

where

$$R(T) = \sum_{t \in \text{Leaves}(T)} r(t) \cdot p(t) = \sum_{t \in \text{Leaves}(T)} R(t) \tag{5}$$

$R(T)$ is the training error, $\text{Leaves}(T)$ gives the leaves of tree $T$, $r(t) = 1 - \max_k p(C_k)$ is the misclassification rate and $p(t) = n_t/N$ is the number of samples in node $n_t$ to total training samples N. Now the variation in cost complexity is given by $R_\alpha(T - T_t) - R_\alpha(T)$, where $T$ is the complete tree, $T_t$ is the subtree with root at node $t$, and a tree pruned at node $t$ would be $T - T_t$. An ordering on the internal nodes for pruning is calculated by equating the cost-complexity function $R_\alpha$ of pruned subtree $T - T_t$ to that of the branch at node $t$:

$$g(t) = \frac{R(t) - R(T_t)}{|\text{Leaves}(T_t)| - 1} \tag{6}$$

The final step is to choose the weakest link to prune by calculating argmin $g(t)$. This calculation of $g(t)$ in equation (6) and then pruning the weakest link is repeated until we are left with the root node. This provides a sequence of nested trees $\mathcal{T}$ and associated cost-complexity parameters $\mathcal{A}$.
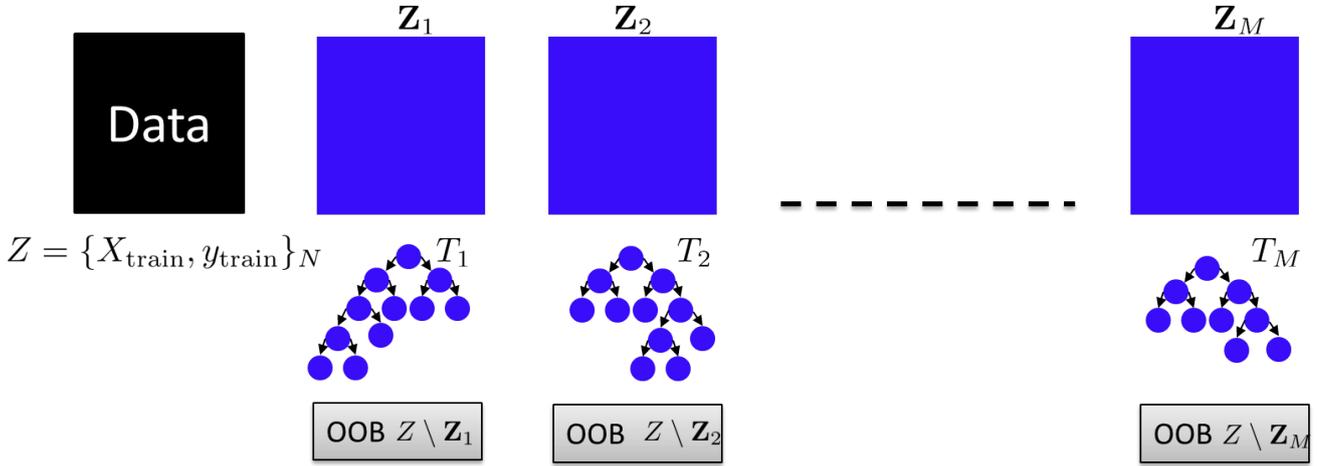
4

Figure 4: Figure shows the bagging procedure, OOB samples and its use for cost-complexity pruning on each decision tree in the ensemble. There are two ways to choose the optimal subtree : one uses the OOB samples as a cross-validation(CV) set to evaluate prediction error. The optimal subtree is the smallest subtree that minimizes the prediction error on the CV-set.

In figure 3 we plot the training error and test(cross-validation) error on 5 folds (usually 20 folds are used, this is only for visualization). We observe a deterioration in performance of both training and test errors. The small tree with 1 SE(standard error) of the cross-validation error is chosen as the optimal subtree. In our studies we use the simpler option which simply chooses the smallest tree with the smallest cross validation (CV) error.

## 2   Out-of-Bag(OOB) cost complexity Pruning

In Random forests, for each tree grown, $\frac{1}{e}N$ samples are not selected in bootstrap, and are called out of bag (OOB) samples. The value $\frac{1}{e}$ refers to the probability of choosing an out-of-bag sample when $N \to \infty$. The OOB samples are used to provide an improved estimate of node probabilities and node error rate in decision trees. They are also a good proxy for generalization error in bagging predictors [4]. OOB data is usually used to get an unbiased estimate of the classification error as trees are added to the forest.

The out-of-bag (OOB) error is the average error on the training set $Z$ predicted such that, samples from the OOB-set $Z \setminus Z_j$ that do not belong to the set of trees $\{T_j\}$ are predicted with as an ensemble, using majority voting (using the sum of their class probabilities).

In our study (see figure 4) we use the OOB samples corresponding to a given tree $T_j$ in the random forest ensemble, to calculate the optimal subtree $T_j^*$ by cross-validation. There are two ways we propose to evaluate the optimal cost-complexity parameter, and thus the optimal subtree :

- Independent tree pruning : calculate the optimal subtree by evaluating

$$\mathcal{T}_j^* = \underset{\alpha \in \mathcal{A}_j}{\operatorname{argmin}} \mathbb{E}\left[\|Y_{\text{OOB}} - \mathcal{T}_j^{(\alpha)}(X_{\text{OOB}}^j)\|^2\right] \tag{7}$$

  where $X_{\text{OOB}}^j = X_{\text{train}} \setminus X_j$, and $X_j$ being the samples used in the creation of tree $j$.

- Global threshold pruning : calculate the optimal subtree by evaluating

$$\{\mathcal{T}_j^*\}_{j=1}^M = \underset{\alpha \in \cup_j \mathcal{A}_j}{\operatorname{argmin}} \mathbb{E}\left[\|Y_{\text{train}} - \frac{1}{M}\sum_{j=1}^M \mathcal{T}_j^{(\alpha)}(X_{\text{OOB}}^j)\|^2\right] \tag{8}$$

where the cross-validation uses the out-of-bag prediction error as to evaluate the optimal $\{\alpha_j\}$ values. This basically considers a single threshold of cost-complexity parameters, which chooses a forest of subtrees for each threshold. The optimal threshold is calculated by cross-validating over the training set.

The independent tree pruning and global threshold pruning are demonstrated in algorithmic form in figure 4 as functions, BestTree_byCrossValidation_Tree and BestTree_byCrossValidation_Forest. The main difference between them lies in the cross-validation samples and predictor (tree vs forest) used.
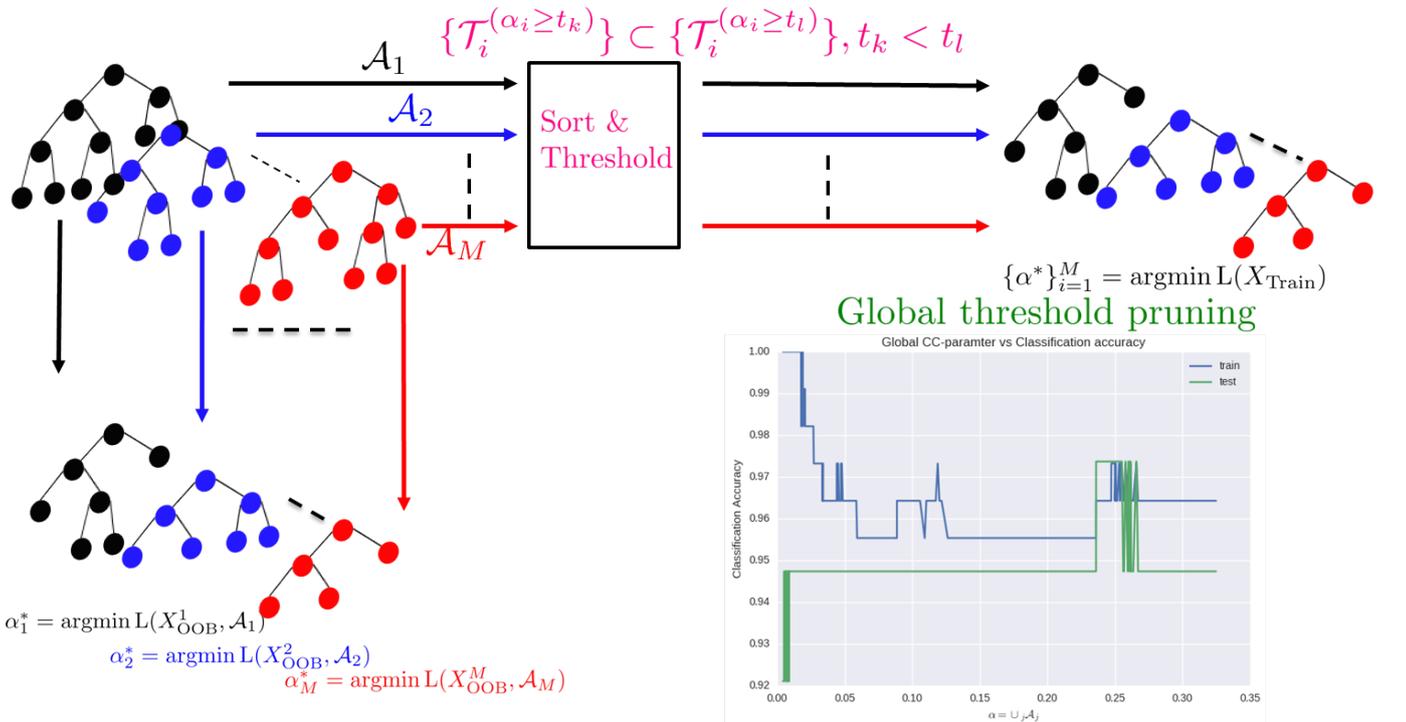
---

**Algorithm 1** Creating random forest

---

**Precondition:** $X_{\text{train}} \in \mathbb{R}^{N \times d}, Y_{\text{train}} \in \{C_k\}_1^K$, M-trees

1: **function** CREATEFOREST($\{T_i\}_{i=1}^M, X_{\text{train}}, y_{\text{train}}$)
2:     **for** $j \in [1, M]$ **do**
3:         $\mathbf{Z}_j \leftarrow \text{BootStrap}(X_{\text{train}}, y_{\text{train}}, \text{N})$
4:         $T_j \leftarrow \text{GrowDecisionTree}(\mathbf{Z}_j)$ ▷ Recursively repeat : 1. select $m = \sqrt{(d)}$ variables 2. pick best split point among $m$ 3. Split node into daughter nodes.
5:         $Y_{\text{pred}} \leftarrow \text{argmax}_{C_k} \frac{1}{M} \sum_{j=1}^M T_j(X_{\text{test}})$
6:     **end for**
7:     $\text{Error} \leftarrow \|Y_{\text{pred}} - Y_{\text{test}}\|^2$
8:     **return** $\{T_j\}_{j=1}^M$
9: **end function**

---



Figure 5: The two pruning methods from equations 7 and 8 are visually demonstrated. For the global method one can plot a training-vs-validation error plot can be generated. This is very similar to figure 3 for a single decision tree.

---

**Algorithm 2** Cost complexity pruning on DTs

---

**Precondition:** $T_j$ a DT, $r$, re-substitution error, $p$ node occupancy proportion

1: **function** COSTCOMPLEXITYPRUNE_TREE($T_j, p, r$)
2:     $\mathcal{T}_j \leftarrow \emptyset$
3:     $\mathcal{A}_j \leftarrow \emptyset$               $\triangleright$ $\mathcal{T}_j$ for set of pruned trees & cost-complexity parameter set $\mathcal{A}_j$.
4:     **while** $T_j \neq \{0\}$ **do**             $\triangleright$ While tree is not pruned to root node.
5:         **for** $t \in T_i \setminus \text{Leaves}(T_i) \setminus \{0\}$ **do**     $\triangleright$ For each internal node i.e. not a leaf/root
6:            $R(t) \leftarrow r(t) * p(t)$
7:            $R(T_t) \leftarrow \sum_{t \in \text{Leaves}(n)} r(t) * p(t)$
8:            $g(t) \leftarrow \frac{R(t) - R(T_t)}{|\text{Leaves}(n)| - 1}$
9:            $t_\alpha, \alpha \leftarrow \arg\min g(t)$            $\triangleright$ $\alpha^*$ is decided by cross-validation
10:         **end for**
11:         $\mathcal{T}_j \leftarrow \mathcal{T}_j \cup \text{Prune}(T_j, t_\alpha)$
12:         $\mathcal{A}_j \leftarrow \mathcal{A}_j \cup \alpha$
13:     **end while**
14:     **return** $\mathcal{T}_j, \mathcal{A}_j$
15: **end function**

---

---

**Algorithm 3** Pruning DT ensembles

---

**Precondition:** $\{T_j\}_{j=1}^M$ are $M$ DTs from **RF**, **BT** or **ET**

1: **function** COSTCOMPLEXITYPRUNE_FOREST($\{T_j\}_{i=1}^M, X_{\text{train}}, y_{\text{train}}$)
2:     **for** $i \in [1, \text{n\_iter}]$ **do**
3:         **for** $t \in \{T_j\}_{j=1}^M$ **do**
4:            $\mathcal{T}_j, \mathcal{A}_j \leftarrow \text{CostComplexityPrune\_Tree}(T_j, X_{\text{OOB}}^j)$
5:            $T_j^* \leftarrow \text{SmallestTree\_byCrossValidation}(\mathcal{T}_j, \mathcal{A}_j, Z \setminus Z_j)$    $\triangleright$ $Z \setminus Z_j$ is replaced by $Z$ (the complete training set) to have a larger CV set(2nd algorithm).
6:         **end for**
7:         $T_j^{**} \leftarrow \text{SmallestTree\_byCrossValidation\_Forest}(\mathcal{T}, \mathcal{A}, \{Z \setminus Z_j\})$
8:         $Y_{\text{pred}} \leftarrow \arg\max_{C_k} \frac{1}{M} \sum_{j=1}^M T_j^*(X_{\text{OOB}}^j)$
9:         $\text{Error(i)} \leftarrow \|Y_{\text{pred}} - Y_j^{\text{CV}}\|^2$
10:     **end for**
11:     **return** $T_j^*$
12: **end function**

---

The decision function of the decision tree (also denoted by the same symbol) $T_j$ would ideally map an input vector $\mathbf{x} \in \mathbb{R}^d$ to any of the $C_k$ classes to be predicted. To perform the prediction for a given sample, we find nodes hit by the sample until it reaches each leaf in the DT, and predict its majority vote. The class-probability weighted majority vote across trees is frequently used since it provides a confidence score on the majority vote in each node across the different trees.

In algorithm (3) we evaluate the cost complexity pruning across the $M$ different trees $\{T_j\}$ in the ensemble, and obtain the optimal subtrees $\{T_j^*\}$ which minimize the prediction error on the OOB sample set $Z \setminus Z_j$.

One of the dangers of using the OOB set to evelate optimal subtrees individually, is that in small datasets the OOB samples might no more be representative of the original training samples distribution, and might produce large cross-validation errors. Though it remains to be studied whether using the OOB samples as a cross-validation set would effectively reduce the generalization error for the forest, even if we observe reasonable performance.

---

**Algorithm 4** Best subtree minimizing CV error on OOB-set

---

**Precondition:** $\mathcal{T}_j$ set of nested subtrees of DT $T_j$, indexed by their cost-complexity parameter $\alpha \in \mathcal{A}_j$

1: **function** **BestTree_byCrossValidation_Tree**$(\mathcal{T}_j, X_{\text{OOB}}^j)$
2:      **for** $\alpha \in \mathcal{A}_j$ **do**
3:          $Y_{\text{pred}}(\alpha) \leftarrow \text{argmax}_{C_k} \mathcal{T}_j^\alpha(X_{\text{OOB}}^j)$
4:          CV-Error$(\alpha) \leftarrow \|Y_{\text{pred}}(\alpha) - Y_{\text{OOB}}^j\|^2$
5:      **end for**
6:      $\alpha_j^* \leftarrow \text{argmin}_{\alpha \in \mathcal{A}_j}$ CV-Error
7:      **return** $\mathcal{T}_j(\alpha_j^*)$                          ▷ Returns the best subtree of $T_j$
8: **end function**

**Precondition:** $\{\mathcal{T}_j\}_{j=1}^M$ are $M$ DTs $T_j$ and their nested subtrees indexed by their cost-complexity parameter $\alpha \in \mathcal{A}_j$

1: **function** **BestTree_byCrossValidation_Forest**$(\{\mathcal{T}_j\}^M, \{\mathcal{A}_j\}^M, \{X_{\text{OOB}}^j\}^M)$
2:      Unique_alpha $\leftarrow$ Unique$(\bigcup_{j=1}^M \mathcal{A}_j)$
3:      **for** $a \in$ Unique_alpha **do**
4:          $\{\alpha_j\} \leftarrow \{\alpha \in \mathcal{A}_j | \alpha \leq a\}$
5:          $Y_{\text{pred}}(\{\alpha_j\}) \leftarrow \text{argmax}_{C_k} \frac{1}{M} \sum_{j=1}^M T_j^\alpha(X_{\text{OOB}}^j)$
6:          CV-Error$(\{\alpha_j\}) \leftarrow \|Y_{\text{pred}}(\{\alpha_j\}) - Y_{\text{OOB}}^j\|^2$
7:      **end for**
8:      $\{T_j^*\} \leftarrow \text{argmin}_{\{\alpha_j\}}$ CV-Error$(\{\alpha_j\})$
9:      **return** $\{T_j^*\}$                          ▷ Returns the $M$-best subtrees of $\{T_j\}_1^M$
10: **end function**

---

# 3   Experiments and evaluation

Here we evaluate the Random Forest(RF), Extremely randomized tree(ExtraTrees, ET) and Bagged Trees (Bagger, BTs) models from scikit-learn on datasets from the UCI machine learning repository [12]. The data sets of different sizes are chosen. Datasets chosen were : Fisher's Iris (150 samples, 4 features, 3 classes), red wine (1599 samples, 11 features, 6 classes), white wine (4898 samples, 11 features, 6 classes), digits dataset (1797 samples, 64 features, 10 classes). Code for the pruning

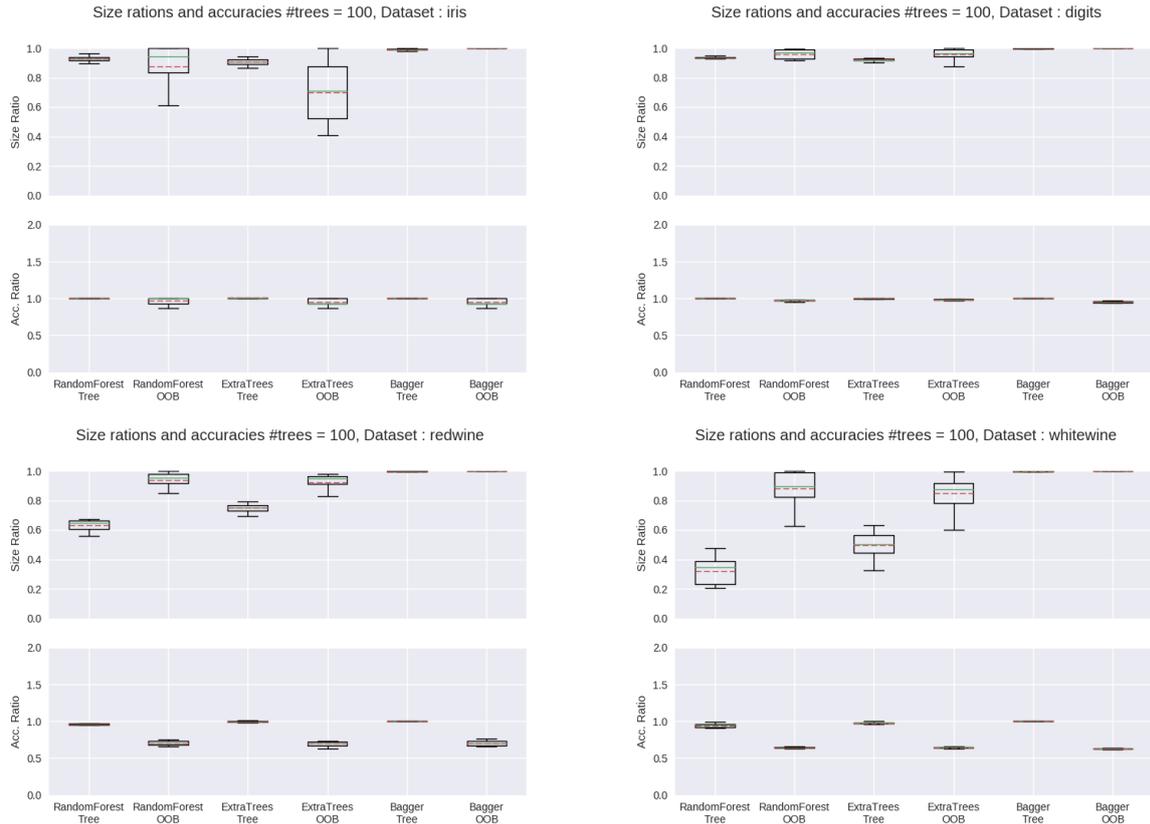experiments are available on github. [2]



Figure 6: Performance of pruning on the IRIS dataset and digits dataset for the three models : RandomForest, Bagged trees and Extra Trees. For each model we show the performance measures for minimizing the tree prediction error and the forest OOB prediction error.

In figure 6 we demonstrate the effect of pruning RFs, BTs and ETs on the different datasets. We observe that random forests an extra trees are often compressed by factors of 0.6 the original size, while maintaining test accuracies, while this is not the case with BTs. To understand the effect of pruning we plot in figure 7 the values $\mathcal{A}_j$ for the different trees in each of the ensembles. We observe that more randomization in RFs and ETs provide a larger set of potential subtrees to cross-validate over.

Another important observation is seen in figure 5, as we prune the forest globally, the forest's accuracy on training set does not monotonically descend (as in the case of a decision tree). As we prune the forest, we could have a set of trees that improve their prediction while the others degrade.

# 4   Conclusions

In this preliminary study of pruning of forests, we studied cost-complexity pruning of decision trees in bagged trees, random forest and extremely randomized trees. In our experiments we observe a reduction in the size of the forest which is dependent on the distribution of points in the dataset. ETs and RFs were shown to perform better than BTs, and were observed to provide a larger set of subtrees to cross-validate. This is the main observation and contribution of the paper.

Our study shows that the out-of-bag samples can be a possible candidate to set the cost-complexity parameter and thus an determine the best subtree for all DTs within ensemble. This

---

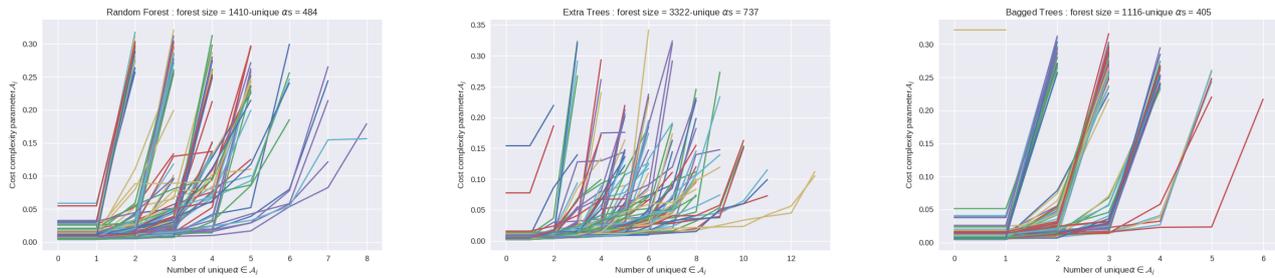[2]https://github.com/beedotkiran/randomforestpruning-ismm-2017

Figure 7: Plot of the cost-complexity parameters for 100-tree ensemble for RFs, ETs and BTs. The distribution of $\mathcal{A}_j$ across trees $j$ shows BTs constitute of similar trees and thus contain similar cost complexity parameter values, while RFs and futhermore ETs have a larger range of parameter values, reflecting the known fact that these ensembles are further randomized. The consequence for pruning is that RFs and more so ETs on average produce subtrees of different depths, and achieve better prediction accuracy and size ratios as compared to BTs.

combines the two ideas originally introduced by Breiman OOB estimates [4] and bagging predictors [5], while using the internal cross-validation OOB score of random forests to set the optimal cost-complexity parameters for each tree.

The speed of calculation of the forest of subtrees is an issue. In the calculation of the forest of subtrees $\{\mathcal{T}\}_{j=1}^{M}$ we evaluate the predicitions at $\text{Unique}(\cup_j \{\mathcal{A}_j\})$ different values of the cost-complexity parameter, which represents the number of subtrees in the forest. In future work we propose to calculate the cost complexity parameter for the forest instead of individual trees.

Though these performance results are marginal, the future scope and goal of this study is to identify the sources of over-fitting in random forests and reduce this by post-pruning. This idea might not be incompatible with smooth-spiked averaged decision function provided by random forests [18].

# References

[1] Biau, G., Devroye, L., Lugosi, G.: Consistency of random forests and other averaging classifiers. JMLR 9, 2015–2033 (2008)

[2] Breiman, L., H. Friedman, J., A. Olshen, R., J. Stone, C.: Classification and Regression Trees. Chapman and Hall, New York (1984)

[3] Breiman, L.: Bagging predictors. Machine learning 24(2), 123–140 (1996)

[4] Breiman, L.: Out-of-bag estimation. Tech. rep., Statistics Department, University of California Berkeley (1996)

[5] Breiman, L.: Random forests. Machine learning 45(1), 5–32 (2001)

[6] Criminisi, A., Shotton, J., Konukoglu, E.: Decision forests: A unified framework for classification, regression, density estimation, manifold learning and semi-supervised learning. Found. Trends. Comput. Graph. Vis. 7, 81–227 (feb 2012)

[7] Denil, M., Matheson, D., De Freitas, N.: Narrowing the gap: Random forests in theory and in practice. In: ICML. pp. 665–673 (2014)

[8] Devroye, L., Gyrfi, L., Lugosi, G.: A probabilistic theory of pattern recognition. Applications of mathematics, Springer, New York, Berlin, Heidelberg (1996)

[9] Geurts, P., Ernst, D., Wehenkel, L.: Extremely randomized trees. Mach. Learn. 63(1), 3–42 (2006)

[10] Hastie, T., Tibshirani, R., Friedman, J.: The elements of statistical learning: data mining, inference and prediction. Springer, 2 edn. (2009)

[11] Ho, T.K.: The random subspace method for constructing decision forests. IEEE transactions on pattern analysis and machine intelligence 20(8), 832–844 (1998)

[12] Lichman, M.: UCI machine learning repository (2013), http://archive.ics.uci.edu/ml

[13] Mingers, J.: An empirical comparison of pruning methods for decision tree induction. Machine learning 4(2), 227–243 (1989)

[14] Rakotomalala, R.: Graphes d'induction. Ph.D. thesis, L'Universit Claude Bernard - Lyon I (1997)

[15] Segal, M.R.: Machine learning benchmarks and random forest regression. Center for Bioinformatics and Molecular Biostatistics (2004)

[16] Torgo, L.F.R.A.: Inductive learning of tree-based regression models. Ph.D. thesis, Universidade do Porto (1999)

[17] Weiss, S.M., Indurkhya, N.: Decision tree pruning: biased or optimal? In: AAAI. pp. 626–632 (1994)

[18] Wyner, A.J., Olson, M., Bleich, J., Mease, D.: Explaining the success of adaboost and random forests as interpolating classifiers. arXiv preprint arXiv:1504.07676 (2015)