# **Cost-complexity pruning of random forests**

ISMM 2017, 13th International Symposium on Mathematical Morphology, Fontainebleau, France, May 15 - 17, 2017

## Why perform pruning?

Out-of-bag samples are un-used samples from the Bootstrap Aggregation procedure in random forests. 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.

## **Decision Tree Pruning**

$$T T - T_1 T - T_2 Leaves(T) = \{3, 4, 5, 6\}$$

#### **Overview of method**



Internal nodes =  $\{1, 2\}$ Sequence of subtrees :  $T \supseteq T - T_2 \supseteq T - T_1$ Cost complexity values :  $g(t) = 0, \alpha_2, \alpha_1$ Final set of trees and parameters :  $\mathcal{T}, \mathcal{A}$ Cost-Complexity function :  $g(t) = \frac{R(t) - R(T_t)}{|\text{Leaves}(T_t)| - 1}$ 

 $R(T) = \sum_{t \in \text{Leaves}(T)} r(t) \cdot p(t) = \sum_{t \in \text{Leaves}(T)} R(t)$ R(T) is the training error, Leaves(T) is #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.



Independent tree pruning

# **Results and Analysis**

### **Plots of** $\mathcal{A}_i \forall j$ **for RFs, ETs, BTs** :



Figure: RFs and ETs provide a larger subset of CC-parameter values  $\mathcal{A}_i$ and thus subtrees  $\mathcal{T}_i$  for the cross-validation step.

**Performance on datasets from UCI repository :** 

Size rations and accuracies #trees = 100, Dataset : iris

and accuracies #trees = 100, Dataset : digits

#### **Random Forests**



**Decision tree ensembles** : Random Forests (RF), Extremely Randomize Trees (ET), Bagged trees (BT) **Randomization Methods** : Boostrap Aggregation, Random Feature selection, Random Threshold section



- Reduction in forest size for marginal loss in classification accuracy.
- Out-of-Bag samples provide cross-validation mechanism to prune forests.

# Out-Of-Bag Cost Complexity Pruning

#### **Independent tree pruning :**

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

**Global threshold** pruning :

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

## Future work

- Understand non-monotonicity (spikes) of random forest training error.
- Does post-pruning preserve consistency of forests? How to define a global cost-complexity parameter for random forests?

**B Ravi Kiran<sup>\*</sup>**, Jean Serra<sup>+</sup> beedotkiran@gmail.com, https://beedotkiran.github.io/forest.html Université Lille 3, CRIStAL, Université<sup>\*</sup>, Ecole des Mines de Paris, CMM <sup>+</sup>

