

# Handout: Generalized Random Forest

**Author:** Malgorzata Olesiewicz, Paper : Generalized Random Forest (Athey et al. 2019)

## Summary

The generalized random forest allows for adaptation of random forest technique to the purpose of heterogeneous parameters prediction. The model has used the generalized method of moments to solve the sample equation, which by the analog principle, also holds true for the population. The generalized random forest has proven to perform better than the other random forest algorithms for the causal heterogeneity estimation. The better results have been attributed to the problem specific split function, which conducts the split of the data according to the specific nature of the parameter one tries to estimate. The model performs well in high dimensional setting and could be used for personalized recommendations, which should be considered as a causal and not prediction problem.

## Random Forest - Survival Kit

Classification and Regression Tree (CART) Decision tree classifies data according to the observations' characteristics. The underling idea here is that observations with similar characteristics will most likely have similar outcomes. The tree is developed though recursive partitioning where observations are classified into subsamples, each of them with different outcome which can be a category or a constant .

Bagging (bootstrap aggregation) focuses on the reduction of variance of the prediction through randomization and averaging the predictions of many unbiased but noisy models.

U-Statistics is a commonly accepted notions of unbiased estimation such as the sample mean and the unbiased sample variance (the "U" stands for "unbiased") in a large sample models.

Honest Tree The ultimate goal of the "honest" CART's algorithms  $\pi(\cdot)$  used on the training sample  $S^{te}$  is to maximize the honest criterion:

$$Q^H(\pi) \equiv -E_{S^{est}, S^{est}, S^{est}} [\text{MSE}(\mathcal{S}^{te}, \mathcal{S}^{est}, \pi(\mathcal{S}^{tr}))] \quad (28)$$

where  $S^{te}$  stand for the test sample and  $S^{est}$  for the estimation sample. The estimated coefficients in each leaf is unbiased as the sample used for the partitioning is different than for its estimation.

## Key equations

Estimation Equation

$$M_{\theta, \nu}(x) = E[\psi_{\theta(x), \nu(x)}(O_i) | X_i = x] = 0 \text{ for all } x \in \mathcal{X} \quad (29)$$

Data:  $(X_i, O_i) \in \mathcal{X} \times \mathcal{O}$ ,  $\psi(\cdot)$ : some scoring function,  $\theta(x)$ : parameter of interest,  $\nu(x)$ : optional nuisance.

The gradient ( $A_P$ ) of the expectation of the moment function is given as:

$$A_P = \frac{1}{|\{i : X_i \in P\}|} \sum_{\{i : X_i \in P\}} \nabla \psi_{\hat{\theta}_P, \hat{\nu}_P}(O_i) \quad (30)$$

The coefficient of the child node ( $\tilde{\theta}_C$ ) is estimated in the Newton Step fashion from the estimate of the parent node ( $\hat{\theta}_P$ ). The relation between the estimators is given by the equation:

$$\tilde{\theta}_C = \hat{\theta}_P - \frac{1}{|\{i : X_i \in C\}|} \sum_{\{i : X_i \in C\}} \xi^\top A_P^{-1} \psi_{\hat{\theta}_P, \hat{\nu}_P}(O_i) \quad (31)$$

With vector  $\xi$  we chose the parameters of interest from the parameters vector. The second component of the right had side of the equation simply gives an average influence function of the  $i$ -th observation which ends up in the children node in estimating the parameter in the parent node.

Pseudo-outcomes in the children nodes:

$$\rho_i = -\xi^\top A_P^{-1} \psi_{\hat{\theta}_P, \hat{\nu}_P}(O_i) \in R \quad (32)$$

The regression step uses the pseudo-outcomes to run the standard CART and obtains split which maximizes the  $\Delta$ -criterion defined as:

$$\tilde{\Delta}(C_1, C_2) = \sum_{j=1}^2 \frac{1}{|\{i : X_i \in C_j\}|} \left( \sum_{\{i : X_i \in C_j\}} \rho_i \right)^2 \quad (33)$$

	<b>Random Forest</b>	<b>Generalized Ransom Forest</b>
Purpose	Prediction	Estimation of highly heterogeneous $\theta(x)$ parameters
Estimating Equation	$\hat{f}_{\text{rf}}^B(x) = \frac{1}{B} \sum_{b=1}^B T_b(x)$	$E[\psi_{\theta(x), \nu(x)}(O_i)   X_i = x] = 0$ for all $x \in \mathcal{X}$
Assumptions	N/A	1.Lipschitz x-signal 2.Smooth Identification 3.Lipschitz( $\theta, \nu$ )-variogram 4. <i>Regularity of <math>\psi</math></i> 5.Existence of solution 6.Concavity
Algorithm		
Splitting the Data	Training, Testing Subset	Training, Estimating and Testing Subset
Building the Tree	<p>1.Grow a tree by re-cursively repeating the following steps for each terminal node of the tree, until the minimum node size <math>n_{\min}</math> is reached.</p> <p>i. Select m variables at random from the p variables.</p> <p>ii. Pick the best variable/split-point among the <math>m</math>.</p> <p>iii. Split the node into two daughter nodes</p>	<p>1.GradientTree(set of examples <math>J</math>, domain <math>X</math>)</p> <p>2. node <math>P_0 \leftarrow \text{CreateNode}(J, X)</math></p> <p>3. queue <math>Q \leftarrow \text{InitializeQueue}(P_0)</math></p> <p>4: while NotNull(node <math>P \leftarrow \text{Pop}(Q)</math>)do</p> <p>5: <math>(\hat{\theta}_P, \hat{\nu}_P, A_P) \leftarrow \text{SolveEstimatingEquation}(P)</math></p> <p>6: <math>R_P \leftarrow \text{GetPseudoOutcomes}(\hat{\theta}_P, \hat{\nu}_P, A_P)</math></p> <p>7: split <math>\Sigma \leftarrow \text{MakeCartSplit}(P, R_P)</math></p> <p>8: if SplitSucceeded(<math>\Sigma</math>)then</p> <p>9: SetChildren (<math>P</math>, GetLeftChild (<math>\Sigma</math>), GetRightChild(<math>\Sigma</math>))</p> <p>10: AddToQueue(<math>Q</math>, GetLeftChild (<math>\Sigma</math>))</p> <p>11: AddToQueue(<math>Q</math>, GetRightChild (<math>\Sigma</math>))</p> <p>12: output tree with root node <math>P_0</math></p>
Building the Forest	<p>1: For <math>b = 1</math> to <math>B</math> : Draw a bootstrap sample <math>Z</math> of size <math>N</math> from the training data</p> <p>2: Grow a random-forest tree <math>T_b</math> to the bootstrapped data.</p> <p>3: Output the ensemble of trees <math>\{T_b\}_1^B</math></p> <p>To make a prediction at a new point <math>x</math>:</p> <p>Regression: <math>\hat{f}_{\text{rf}}^B(x) = \frac{1}{B} \sum_{b=1}^B T_b(x)</math></p> <p>Classification: Let <math>\hat{C}_b(x)</math> be the class prediction of the <math>b</math>th random-forest tree.</p> <p>Then <math>C_{\text{rf}}^B(x) = \text{majority vote } \{C_b(x)\}_1^B</math></p>	<p>1. GeneralizedRandomForest(set of examples <math>S</math>, test point <math>x</math>)</p> <p>2: weight vector <math>\leftarrow \text{Zeros}( S )</math></p> <p>3: for <math>b = 1</math> to total number of trees <math>B</math> do</p> <p>4: set of examples <math>I \leftarrow \text{Subsample}(S, s)</math></p> <p>5: sets of examples <math>J_1, J_2 \leftarrow \text{SplitSample}(I)</math></p> <p>6: tree <math>T \leftarrow \text{GradientTree}(J_1, X)</math></p> <p>7: <math>N \leftarrow \text{Neighbors}(x, T, J_2)</math> . Returns those elements of <math>J_2</math> that fall into the same leaf as <math>x</math> in the tree <math>T</math> .</p> <p>8: for all example <math>e \in N</math>do</p> <p>9: <math>\alpha[e] += 1/ N </math></p> <p>10: output <math>\theta(x)</math>, the solution to (2) with weights <math>\alpha/B</math></p>
Split Equation	$\min_{j,s} \left[ \min_{c_1} \sum_{x_i \in R_1(j,s)} (y_i - c_1)^2 + \min_{c_2} \sum_{x_i \in R_2(j,s)} (y_i - c_2)^2 \right]$	$\tilde{\Delta}(C_1, C_2) = \sum_{j=1}^2 \frac{1}{ \{i: X_i \in C_j\} } \left( \sum_{\{i: X_i \in C_j\}} \rho_i \right)^2$