

# Tree-based models: From CART to random forests

Dongjae Son

Yonsei University

Department of Applied Statistics

February 25, 2020

# Before we go on

- These slides rely heavily on [Hastie et al., 2009], [Bishop, 2006] and [Murphy, 2012]

# Introduction: Adaptive basis function models

- Linear regression models:
    - Models conditional mean given covariates
    - Easy to handle and interpret
    - Various regularization methods exist: Ridge, Lasso, Elastic net, etc.
  - Disadvantages of the linear models:
    - Relationship between target variable and predictors are not always linear
    - Not easy to incorporate interaction effect
- ∴ Bad prediction performance on nonlinear patterns

# Introduction: Generalized Additive Models(GAM)

- Generalized Additive Models(GAM) by [Hastie and Tibshirani, 1987]
  - $g(E[Y|X_1, X_2, \dots, X_p]) = \alpha + \sum_{j=1}^p f_j(X_j)$
  - Choice of  $g$  can be:
    - ①  $g(\mu) = \mu$  is the identity link for continuous response
    - ②  $g(\mu) = \log\left(\frac{\mu}{1-\mu}\right)$  is the logit link for binary response
    - ③  $g(\mu) = \log(\mu)$  is for log-linear models for Poisson count data
  - Same with the choice of link functions in GLM

# Introduction: Generalized Additive Models(GAM, cont'd)

- Fitting additive models: use scatterplot smoother, usually smoothing splines
  - $y = \alpha + \sum_{j=1}^p f_j(X_j) + \epsilon$  where  $E[\epsilon] = 0$
  - Minimize the penalized residual sum of squares(PRSS)

$$J = \sum_{i=1}^N \left( y_i - \alpha - \sum_{j=1}^p f_j(X_j) \right)^2 + \sum_{j=1}^p \lambda_j \int f_j''(x_j) dx_j$$

where  $\lambda_j \geq 0$  are tuning parameters

- $\alpha$  is not uniquely determined; add the constraint  $\sum_{i=1}^N f_j(x_{ij}) = 0$  for all  $j \rightarrow \hat{\alpha} = \bar{y}$  and it never changes

# Introduction: Generalized Additive Models(GAM, cont'd)

- Solving the PRSS: the backfitting algorithm

- ① Initialize:  $\hat{\alpha} = \frac{1}{N} \sum_{i=1}^N y_i$  and  $\hat{f}_j = 0$  for all  $j$
- ② Cycle:  $j = 1, 2, \dots, p, \dots, 1, 2, \dots, p, \dots$

$$\hat{f}_j \leftarrow S_j \left[ \{y_i - \hat{\alpha} - \sum_{k \neq j} \hat{f}_k(x_{ik})\}_1^N \right]$$
$$\hat{f}_j \leftarrow \hat{f}_j - \frac{1}{N} \sum_{i=1}^N \hat{f}_j(x_{ij})$$

until the functions  $\hat{f}_j$  change less than a prespecified threshold, say,  $10^{-3}$ . The second line of the cycle process ensures the constraint  $\sum_{i=1}^N \hat{f}_j(x_{ij}) = 0$  for all  $j$ .

# Introduction: Generalized Additive Models(GAM, cont'd)

Although the GAM looks fascinating, it has certain disadvantages:

- 1 Too much computational costs if we introduce interaction effects

$$g(\mu(X)) = \alpha + \sum_{j=1}^p f_j(X_j) + \sum_{j,k} f_{jk}(X_j, X_k) + \sum_{j,k,l} f_{jkl}(X_j, X_k, X_l) + \dots$$

- 2 Not suitable for large  $p$  problems: COSCO procedures([Lin et al., 2006]) or SpAM(Sparse Additive Model, [Ravikumar et al., 2009]) approach

# Classification and Regression Trees

- Tree-based models predict  $y$  with features  $\mathbf{X} \in \mathbb{R}^p$  by dividing the feature space into disjoint rectangles  $R_m$ , or leaves of a tree.
  - $\hat{f}(\mathbf{X}) = \sum_{m=1}^M c_m \mathbb{I}\{\mathbf{X} \in R_m\}$  where  $\mathbb{I}$  is an indicator function
  - $\mathcal{F} = \cup_{m=1}^M R_m$  such that  $R_m \cap R_l = \emptyset$  for  $m \neq l$
  - $R_m$  can be constructed by binary splitting of predictors, i.e.  $\{X_j \leq s\}$
- Suggested by [Breiman et al., 1984]



# Classification and Regression Trees(cont'd)

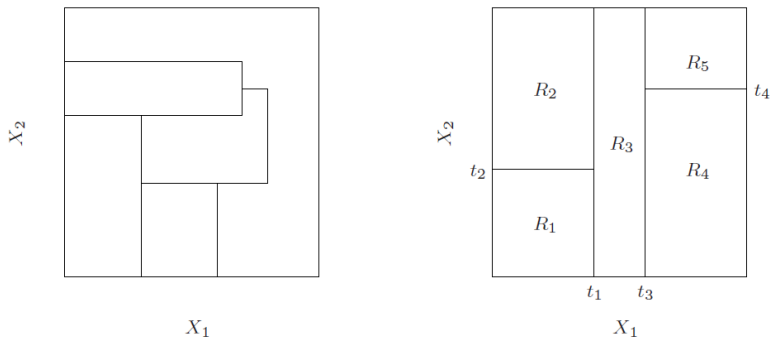


Figure: Partitioning and CART from [Hastie et al., 2009]

# Classification and Regression Trees(cont'd)

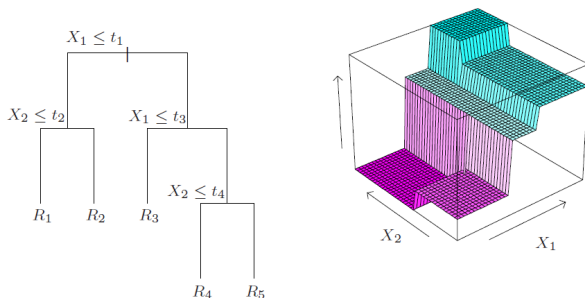


Figure: Partitioning and CART from [Hastie et al., 2009]

# Regression Trees

- For regression trees,  $c_m = \bar{y}_{R_m}$ , just sample mean at rectangle  $R_m$
- Thus, starting with all of the data, we seek the splitting variable  $j$  and split point  $s$  such that

$$\min_{j,s} \left[ \sum_{x_i \in R_1(j,s)} (y_i - \bar{y}_{R_1(j,s)})^2 + \sum_{x_i \in R_2(j,s)} (y_i - \bar{y}_{R_2(j,s)})^2 \right]$$

where  $R_1(j, s) = \{X | X_j \leq s\}$  and  $R_2(j, s) = \{X | X_j > s\}$

- Then, repeat this process on every resulting region and terminate if every leaf contains five or less observations

# Regression Trees(cont'd)

- It can be rewritten as minimizing the loss function

$$R(T) = \sum_{m=1}^{|T|} \sum_{x_i \in R_m} (y_i - \bar{y}_{R_m})^2$$

where  $|T|$  is the tree size, or the number of leaves

- Tree size  $|T|$  is a tuning parameter controlling the complexity of tree model, which should be controlled for generalization

# Cost complexity pruning

- To control the tree size, we introduce the penalized loss

$$R_\alpha(T) = R(T) + \alpha|T|$$

where  $\alpha \geq 0$  is the tuning parameter that controls the trade-off between tree size and goodness-of-fit

- Idea: Find the subtree  $T(\alpha) \subseteq T^0$  for each  $\alpha$  that minimizes  $R_\alpha(T)$  where  $T^0$  is the fully grown tree

## Cost complexity pruning(cont'd)

- For each  $\alpha$  there exists a sequence of trees  $T^0 \supseteq T^1 \supseteq \dots \supseteq T^n$  where  $T^n$  is the null tree
- The sequence  $T^s$  can be generated by replacing a subtree  $T_t$  with root node  $t$  with a leaf

$$\begin{aligned} & [R_\alpha(T - T_t) - R_\alpha(T)] \\ &= R(T - T_t) - R(T) + \alpha(|T - T_t| - |T|) \\ &= R(T) - R(T_t) + R(t) - R(T) + \alpha(|T| - |T_t| + 1 - |T|) \\ &= R(t) - R(T_t) + \alpha(1 - |T_t|) \end{aligned}$$

## Cost complexity pruning(cont'd)

- Solving  $R(T) - R(T_t) + \alpha(1 - |T_t|) = 0$  yields  $\alpha = \frac{R(t) - R(T_t)}{|T_t| - 1}$
- Starting with the full tree  $T^0$  (and  $\alpha_0 = 0$ ) in each step  $s = 1, 2, \dots$  the algorithm goes:
  - 1 Select the node  $t$  which minimizes  $\frac{R(t) - R(T_t^{s-1})}{|T_t^{s-1}| - 1}$
  - 2 Set  $T^s = T^{s-1} - T_t^{s-1}$  and  $\alpha_s = \frac{R(t) - R(T_t^{s-1})}{|T_t^{s-1}| - 1}$
  - 3 Keep this process until we get the null tree
- Hence, we get  $T^0 \supseteq T^1 \supseteq \dots \supseteq T^n$  simultaneously with  $\alpha_0 \leq \alpha_1 \leq \dots \leq \alpha_n$ . Determine the optimal  $\alpha_*$  by 5- or 10-fold cross-validation.

## V-fold cross-validation for optimal $\alpha_*$

- 1 First, grow  $T_0$  using the whole data  $\mathcal{D}$  and get  $T^0 \supseteq T^1 \supseteq \dots \supseteq T^n$  simultaneously with  $\alpha_0 \leq \alpha_1 \leq \dots \leq \alpha_n = \infty$
- 2 Define  $\alpha'_k = \sqrt{\alpha_k \alpha_{k+1}}$ , a geometric mean
- 3 For each  $v \in \{1, 2, \dots, V\}$  and  $\alpha \in \{\alpha'_1, \alpha'_2, \dots\}$  use  $\mathcal{D}^{(v)} = \mathcal{D} - \mathcal{D}_v$  to grow  $T^{(v,0)}$  and find  $T^{(v,k')}$  with corresponding sequence of  $\alpha'_k$
- 4 Choose  $\alpha'_k$  that produces the minimum CV error then set  $\alpha_* = \alpha_k$



## V-fold cross-validation for optimal $\alpha_*$ (cont'd)

- We can use the sequence of geometric means since  $T(\alpha) = T(\alpha_k)$  for  $\alpha_k \leq \alpha < \alpha_{k+1}$
- Detailed proofs are omitted here. Ask me if you are curious how it works...!

# Classification Trees

- For classification trees,  $c_m$ , is the modal class in the rectangle  $R_m$  among the class  $\{1, 2, \dots, K\}$
- The only difference between classification trees and regression trees can be found in *loss criteria*; no more squared-loss!
- Define
  - $N_m$ : the number of samples in rectangle  $R_m$
  - $p_{mk}$ : the proportion of class  $k$  in rectangle  $R_m$

# Classification Trees(cont'd)

- Misclassification rate, or 0-1 loss:

$$R(T) = \sum_{m=1}^{|T|} \sum_{x_i \in R_m} \mathbb{I}\{y_i \neq c_m\}$$

- Gini index:

$$R(T) = \sum_{m=1}^{|T|} N_m \sum_{k=1}^K p_{mk}(1 - p_{mk})$$

- Cross-entropy:

$$R(T) = - \sum_{m=1}^{|T|} N_m \sum_{k=1}^K p_{mk} \log p_{mk}$$

## Classification Trees(cont'd)

- Consider a node  $t$  with size  $N_t$  and loss criterion  $R(t)$
- For some variable  $j$  and split point  $s$ , we split  $t$  into two nodes,  $t_R$  and  $t_L$  with size  $N_{t_R}$  and  $N_{t_L}$  and loss  $R(t_R)$  and  $R(t_L)$
- The mean decrease in loss can be defined as:

$$\Delta(j, s) = R(t) - \left( \frac{N_{t_R}}{N_t} R(t_R) + \frac{N_{t_L}}{N_t} R(t_L) \right)$$

- Find  $(j_*, s_*)$  such that

$$(j_*, s_*) = \arg \min_{j, s} \Delta(j, s)$$

# Categorical predictors?

- For an un-ordered predictor variable with  $q$  levels,  $2^q - 1$  splits are possible (Why?)
- `rpart` function in **R** deals with factors, but tree-based methods in **scikit-learn** do not.. need for dummification?

## Further applications of CART

- Multivariate Adaptive Regression Splines(MARS): [Friedman, 1991]
- Hierarchical Mixture of Experts(HME): [Jordan and Jacobs, 1994]
- And many other researches exist.. search for them if you find it interesting!

# Bootstrap

- Before going into bagging methods, let's have a grasp on 'bootstrap' method
- The bootstrap is one of the key resampling methods in statistics
- It was mainly used to approximately estimate the variance of an estimator  $\hat{\theta} = g(U_1, \dots, U_N)$
- Key idea: Use empirical distribution  $\hat{F}_N$  as a substitute for the true, unknown distribution  $F$

# Bootstrap(cont'd)

- For iid sample  $U_1, \dots, U_N \sim F$ , the empirical distribution can be written as

$$\hat{F}_N(A) = \frac{1}{N} \sum_{i=1}^N \mathbb{I}(U_i \in A)$$

- For univariate case, the Glivenko-Cantelli theorem says that

$$\left\| \hat{F}_N - F \right\|_{\infty} = \sup_{x \in \mathbb{R}} |\hat{F}_N(x) - F(x)| \rightarrow$$

almost surely as  $N \rightarrow \infty$



# Bootstrap(cont'd)

- Bootstrap variance estimator:
  - ① Draw  $\{U_i^b\}_{i=1}^n \sim \hat{F}_N$  randomly **with replacement** and compute  $\hat{\theta}_b = g(U_1^b, \dots, U_n^b)$
  - ② Repeat this for  $B$  times
  - ③ Compute  $\widehat{Var}(\hat{\theta}) = \frac{1}{B-1} \sum_{b=1}^B (\hat{\theta}_b - \bar{\theta})^2$  where  $\bar{\theta} = \frac{1}{B} \sum_{b=1}^B \hat{\theta}_b$
- Detailed explanations of bootstrap(constructing confidence intervals or asymptotic results, etc) are omitted hereafter

# Bootstrap Aggregation(Bagging)

- Idea: Use the bootstrap to enhance the prediction performance
- For training data  $\mathbf{Z} = \{(\mathbf{X}^{(1)}, Y^{(1)}), \dots, (\mathbf{X}^{(N)}, Y^{(N)})\}$ ,
  - ① Draw bootstrap samples  $\mathbf{Z}^b$  for  $b = 1, 2, \dots, B$
  - ② Obtain the bagging estimate

$$\hat{f}_{bag}(x) = \frac{1}{B} \sum_{b=1}^B \hat{f}^b(x)$$

where  $\hat{f}^b(\cdot) = \hat{f}(\cdot, \mathbf{Z}^b)$  is an estimator trained with  $\mathbf{Z}^b$  for  $b = 1, 2, \dots, B$

## Bagging(cont'd)

- Basically, it works well with 'unstable' estimator, i.e. an estimator with low bias and **high variance**
- Intuitively illustrating, if the variables  $U_1, \dots, U_B$  with positive correlation  $\rho$  are sampled from an identical distribution, its average has variance

$$Var(\bar{U}_B) = \rho\sigma^2 + \frac{1-\rho}{B}\sigma^2 \rightarrow \rho\sigma^2$$

as  $B \rightarrow \infty$

- For details, see [Breiman, 1996]

# Random Forests

- First suggested by [Breiman, 2001]
- The key idea of random forests is to lower the variance of bagging estimate via **reduction in correlation**
- This is achieved by random selection of inputs at each split of each bootstrapped tree

# Random Forests(cont'd)

- 1 For  $b = 1, \dots, B$ :
  - (a) Draw a bootstrap sample  $\mathbf{Z}^b$  from the training data  $\mathbf{Z}$
  - (b) Build a tree  $\hat{f}^b$  to the bootstrapped data  $\mathbf{Z}^b$  without pruning, **while choosing  $m < p$  variables at each node** (*Rule of thumb*:  $m = p/3$  for regression and  $m = \sqrt{p}$  for classification)
- 2 For the output  $\{\hat{f}^b\}_{b=1}^B$ :
  - (a) Classification:  $\hat{f}_{RF}(\mathbf{x}) = \arg \max_{j \in \{1, \dots, K\}} \frac{1}{B} \sum_{b=1}^B \mathbb{I}(\hat{f}^b(\mathbf{x}) = j)$
  - (b) Regression:  $\hat{f}_{RF}(\mathbf{x}) = \frac{1}{B} \sum_{b=1}^B \hat{f}^b(\mathbf{x})$

# How random forests decorrelates

- We can express  $\hat{f}^b(\mathbf{x}) = \hat{f}(\mathbf{x}, \Theta_b)$  where  $\Theta_b$  characterizes the  $b$ -th tree in terms of split variable, splitting points, terminal node values and  $\mathbf{Z}^b$  itself
- From that point of view, we can consider  $\Theta_b$  as a random variable
- However, distribution of  $\Theta_b$  or  $\Theta$  need not be specified since it would be only used theoretically

# How random forests decorrelates(classification)

- By [Breiman, 2001], almost surely for all  $\Theta$

$$\frac{1}{B} \sum_{b=1}^B \mathbb{I} \left( \hat{f}(\mathbf{x}, \Theta) = j \right) \rightarrow P_{\Theta} \left( \hat{f}(\mathbf{x}, \Theta) = j \right)$$

- Define a *margin function*

$$\begin{aligned} mr(\mathbf{X}, Y) &= P_{\Theta} \left( \hat{f}(\mathbf{x}, \Theta) = Y \right) - \max_{j \neq Y} P_{\Theta} \left( \hat{f}(\mathbf{x}, \Theta) = j \right) \\ &= P_{\Theta} \left( \hat{f}(\mathbf{x}, \Theta) = Y \right) - P_{\Theta} \left( \hat{f}(\mathbf{x}, \Theta) = \hat{j}(\mathbf{X}, Y) \right) \end{aligned}$$

where  $\hat{j}(\mathbf{X}, Y) = \arg \max_{j \neq Y} P_{\Theta} \left( \hat{f}(\mathbf{x}, \Theta) = j \right)$

## How random forests decorrelates(classification, cont'd)

- Furthermore, define a *raw margin function* by

$$rmg(\mathbf{X}, Y, \Theta) = \mathbb{I}(\hat{f}(\mathbf{x}, \Theta) = Y) - \mathbb{I}(\hat{f}(\mathbf{x}, \Theta) = \hat{j}(\mathbf{X}, Y))$$

so that we can rewrite a margin function as

$$mg(\mathbf{X}, Y) = E_{\Theta} rmg(\mathbf{X}, Y, \Theta)$$



## How random forests decorrelates(classification, cont'd)

- Lastly, define the generalization error by

$$PE^* = P_{\mathbf{X},Y} (mr(\mathbf{X}, Y) < 0)$$

- By the Chebyshev inequality, we can easily show that

$$PE^* \leq \frac{Var(mr(\mathbf{X}, Y))}{s^2}$$

where  $s = E_{\mathbf{X},Y} mr(\mathbf{X}, Y)$ , strength of a classifier  $\hat{f}(\mathbf{X}, \Theta)$

# How random forests decorrelates(classification, cont'd)

## Theorem

*Assume that  $s > 0$ . Then*

$$PE^* \leq \bar{\rho} \frac{1 - s^2}{s^2}$$

*where  $\bar{\rho} = E_{\Theta} E_{\Theta'} [\rho(\Theta, \Theta') sd(\Theta) sd(\Theta')] / (E_{\Theta} sd(\Theta))^2$ , the average correlation between  $rmg(\mathbf{X}, Y, \Theta)$  and  $rmg(\mathbf{X}, Y, \Theta')$ .*

## How random forests decorrelates(classification, cont'd)

Proof.

$$\begin{aligned} \text{Var}(mr(\mathbf{X}, Y)) &= E[E_{\Theta}rmg(\mathbf{X}, Y, \Theta)]^2 - [EE_{\Theta}rmg(\mathbf{X}, Y, \Theta)]^2 \\ &= E_{\Theta, \Theta'}Cov(rmg(\mathbf{X}, Y, \Theta), rmg(\mathbf{X}, Y, \Theta')) \end{aligned}$$

which can be implied from the fact that

$$[E_{\Theta}f(\Theta)]^2 = E_{\Theta, \Theta'}f(\Theta)f(\Theta')$$

for independent  $\Theta, \Theta'$  and any  $f$ .

## How random forests decorrelates(classification, cont'd)

Proof (cont'd).

For any random variables  $U$  and  $V$ ,  $\rho(U, V) = \frac{\text{Cov}(U, V)}{\sqrt{\text{Var}(U)}\sqrt{\text{Var}(V)}}$ . Thus

$$\text{Var}(mr(\mathbf{X}, Y)) = E_{\Theta, \Theta'} \rho(\Theta, \Theta') sd(\Theta) sd(\Theta')$$

where  $\rho(\Theta, \Theta')$  is the correlation btw  $rmg(\mathbf{X}, Y, \Theta)$  and  $rmg(\mathbf{X}, Y, \Theta')$  holding  $\Theta$  and  $\Theta'$  fixed and  $sd(\Theta)$  is the standard dev. of  $rmg(\mathbf{X}, Y, \Theta)$  for fixed  $\Theta$ .

## How random forests decorrelates(classification, cont'd)

Proof (cont'd).

By the definition of  $\bar{\rho}$ , we have

$$\text{Var}(mr(\mathbf{X}, Y)) = \bar{\rho} \{E_{\Theta} sd(\Theta)\}^2 \leq \bar{\rho} E_{\Theta} sd(\Theta)^2$$

by the Jensen's inequality. Now consider

$$E_{\Theta} sd(\Theta)^2 = E_{\Theta} \left[ E_{\mathbf{X}, Y} rmg(\mathbf{X}, Y, \Theta)^2 - (E_{\mathbf{X}, Y} rmg(\mathbf{X}, Y, \Theta))^2 \right]$$

## How random forests decorrelates(classification, cont'd)

Proof.

Since

$$\begin{aligned}s^2 &= [E_{\mathbf{X},Y}mr(\mathbf{X},Y)]^2 = [E_{\mathbf{X},Y}E_{\Theta}rmg(\mathbf{X},Y)]^2 \\ &= [E_{\Theta}E_{\mathbf{X},Y}rmg(\mathbf{X},Y)]^2 \leq E_{\Theta} [E_{\mathbf{X},Y}rmg(\mathbf{X},Y)]^2\end{aligned}$$

and

$$E_{\Theta}E_{\mathbf{X},Y}rmg(\mathbf{X},Y,\Theta)^2 \leq 1.$$

Therefore  $E_{\Theta}sd(\Theta)^2 \leq 1 - s^2$  and this completes the proof. □

## How random forests decorrelate(regression, cont'd)

- Also by [Breiman, 2001], we have that almost surely for  $\Theta$

$$\hat{f}_{RF}(\mathbf{x}) = \frac{1}{B} \sum_{b=1}^B \hat{f}(\mathbf{x}, \Theta_b) \rightarrow E_{\Theta} \hat{f}(\mathbf{x}, \Theta)$$

- Define generalization error of forests and each tree by

$$PE^*(forest) = E_{\mathbf{X}, Y} \left( Y - E_{\Theta} \hat{f}(\mathbf{x}, \Theta) \right)^2$$

and

$$PE^*(tree) = E_{\Theta} E_{\mathbf{X}, Y} \left( Y - \hat{f}(\mathbf{x}, \Theta) \right)^2$$

# How random forests decorrelate(regression, cont'd)

## Theorem

Assume for all  $\Theta$  that  $E_{\mathbf{X}, Y} [Y - \hat{f}(\mathbf{x}, \Theta)] = 0$ . Then

$$PE^*(\text{forest}) \leq \bar{\rho} PE^*(\text{tree})$$

where  $\bar{\rho} = E_{\Theta} E_{\Theta'} [\rho(\Theta, \Theta') sd(\Theta) sd(\Theta')] / (E_{\Theta} sd(\Theta))^2$ , the average correlation between  $Y - \hat{f}(\mathbf{X}, \Theta)$  and  $Y - \hat{f}(\mathbf{X}, \Theta')$ .



## How random forests decorrelate(regression, cont'd)

Proof.

$$\begin{aligned} PE^*(forest) &= E_{\mathbf{X}, Y} \left( Y - E_{\Theta} \hat{f}(\mathbf{x}, \Theta) \right)^2 \\ &= E_{\mathbf{X}, Y} \left[ E_{\Theta} \left( Y - \hat{f}(\mathbf{x}, \Theta) \right) \right]^2 \\ &= E_{\mathbf{X}, Y} \left[ E_{\Theta} \left( Y - \hat{f}(\mathbf{x}, \Theta) \right) \cdot E_{\Theta'} \left( Y - \hat{f}(\mathbf{x}, \Theta') \right) \right] \\ &= E_{\Theta, \Theta'} \left[ E_{\mathbf{X}, Y} \left( Y - \hat{f}(\mathbf{x}, \Theta) \right) \cdot \left( Y - \hat{f}(\mathbf{x}, \Theta') \right) \right] \\ &= E_{\Theta, \Theta'} Cov \left( Y - \hat{f}(\mathbf{x}, \Theta), Y - \hat{f}(\mathbf{x}, \Theta') \right) \\ &= E_{\Theta, \Theta'} \left[ \rho(\Theta, \Theta') sd(\Theta) sd(\Theta') \right] \end{aligned}$$

## How random forests decorrelate(regression, cont'd)

Proof (cont'd).

By the definition of average correlation  $\bar{\rho}$ , we have

$$\begin{aligned} PE^*(forest) &= \bar{\rho} \cdot (E_{\Theta} sd(\Theta))^2 \\ &\leq \bar{\rho} \cdot E_{\Theta} sd(\Theta)^2 \\ &= \bar{\rho} \cdot E_{\Theta} E_{\mathbf{X}, Y} \left( Y - \hat{f}(\mathbf{x}, \Theta) \right)^2 \\ &= \bar{\rho} \cdot PE^*(tree) \end{aligned}$$

which completes the proof. □

# Out-of-bag(OOB) estimates

- By nature of the bootstrap, some of observations may not be included in  $b$ -th bootstrap sample, i.e.

$$P\left(i \notin \mathbf{Z}^b\right) = \left(1 - \frac{1}{N}\right)^N \approx 0.367$$

- Denote the set of such samples  $OOB(b)$ , which we can use as a *test set*
- Thus, in principle, random forests do not require cross-validation!

# Out-of-bag(OOB) estimates(cont'd)

- Define

$$Q(\mathbf{X}, j) = \frac{\sum_{b=1}^B \mathbb{I}(\hat{f}(\mathbf{X}, \Theta_b) = j : (\mathbf{X}, Y) \in OOB(b))}{\sum_{b=1}^B \mathbb{I}((\mathbf{X}, Y) \in OOB(b))}$$

the out-of-bag proportion of votes cast at  $\mathbf{X}$  for class  $j$ , which is an estimator of  $P_{\Theta}(\hat{f}(\mathbf{X}, \Theta) = j)$

- Then an estimator of  $mr(\mathbf{X}, Y)$  is  $Q(\mathbf{X}, Y) - \max_{j \neq Y} Q(\mathbf{X}, j)$

## Out-of-bag(OOB) estimates(cont'd)

- Now, an estimate of generalization error can be written as

$$\frac{1}{N} \sum_{i=1}^N \mathbb{I} \left( Q \left( \mathbf{X}^{(i)}, Y^{(i)} \right) - \max_{j \neq Y^{(i)}} Q \left( \mathbf{X}^{(i)}, j \right) < 0 \right)$$

- OOB estimates of strength and correlation can be calculated in similar ways, where details can be found in [Breiman, 2001]

# Variable importance

- One of the advantages of decision tree model is easy interpretability
- Such an advantage disappears in random forests due to bootstrap and random selection of features
- To explain which predictors contribute to forests more than the others, we introduce the *variable importance*

## Variable importance(cont'd)

- [Breiman, 2001] suggests a *permutation importance*
- The idea is that if the variable  $X_j$  is important, the error using randomly shuffled data (permutation) should be much worse than that using the correct one
- The permutation importance by OOB set measures how much the accuracy deteriorates due to permutation in OOB sets

## Variable importance(cont'd)

- Let  $\pi_b^j$  be a permutation on the  $j$ -th variable in  $b$ -th OOB set. Define

$$VI(j, b) = \frac{\sum_{i \in OOB(b)} \left[ \mathbb{I} \left( Y^{(i)} = \hat{f}^b \left( \mathbf{X}^{(i)} \right) \right) - \mathbb{I} \left( Y^{(i)} = \hat{f}^b \left( \mathbf{X}_{\pi_b^j}^{(i)} \right) \right) \right]}{|OOB(b)|}$$

where  $\mathbf{X}_{\pi_b^j}^{(i)}$  is the value of permuted  $j$ -th variable for observation  $i$ .

- Now the permutation variable importance is written as

$$VI(j) = \frac{1}{B} \sum_{b=1}^B VI(j, b)$$



# Asymptotic properties of random forests

- Simplified proof given by [Breiman, 2004]
- Further researches from [Lin et al., 2006], [Biau, 2012], [Wager and Athey, 2018] or etc

# Further applications of random forests

- Survival data: [Ishwaran et al., 2008]
- Causal inference: [Wager and Athey, 2018]
- Generalized random forests: [Athey et al., 2019]

# Homeworks

- ① By using Chebyshev's inequality, prove that

$$PE^* \leq \frac{Var(mr(\mathbf{X}, Y))}{s^2}.$$

- ② Show that

$$P(i \notin \mathbf{Z}^b) = \left(1 - \frac{1}{N}\right)^N \approx 0.367$$

holds. Why does the approximation hold?

- ③ Are there any methods for quantifying variable importances other than permutation? If then, suggest briefly how you can!

# References I

 Athey, S., Tibshirani, J., Wager, S., et al. (2019).

Generalized random forests.

*The Annals of Statistics*, 47(2):1148–1178.

 Biau, G. (2012).

Analysis of a random forests model.

*Journal of Machine Learning Research*, 13(Apr):1063–1095.

 Bishop, C. M. (2006).

*Pattern recognition and machine learning*.

springer.

# References II



Breiman, L. (1996).

Bagging predictors.

*Machine learning*, 24(2):123–140.



Breiman, L. (2001).

Random forests.




*Machine learning*, 45(1):5–32.



Breiman, L. (2004).

Consistency for a simple model of random forests.

# References III

-  Breiman, L., Friedman, J., Stone, C. J., and Olshen, R. A. (1984).  
*Classification and regression trees*.  
CRC press.
-  Friedman, J. H. (1991).  
Multivariate adaptive regression splines.  
*The annals of statistics*, pages 1–67.
-  Hastie, T. and Tibshirani, R. (1987).  
Generalized additive models: some applications.  
*Journal of the American Statistical Association*, 82(398):371–386.

# References IV



Hastie, T., Tibshirani, R., and Friedman, J. (2009).

*The elements of statistical learning: data mining, inference, and prediction.*

Springer Science & Business Media.



Ishwaran, H., Kogalur, U. B., Blackstone, E. H., Lauer, M. S., et al. (2008).

Random survival forests.

*The annals of applied statistics*, 2(3):841–860.

# References V



Jordan, M. I. and Jacobs, R. A. (1994).

Hierarchical mixtures of experts and the em algorithm.

*Neural computation*, 6(2):181–214.



Lin, Y., Zhang, H. H., et al. (2006).

Component selection and smoothing in multivariate nonparametric regression.

*The Annals of Statistics*, 34(5):2272–2297.



Murphy, K. P. (2012).

*Machine learning: a probabilistic perspective*.

MIT press.



# References VI



Ravikumar, P., Lafferty, J., Liu, H., and Wasserman, L. (2009).

Sparse additive models.

*Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 71(5):1009–1030.



Wager, S. and Athey, S. (2018).

Estimation and inference of heterogeneous treatment effects using random forests.

*Journal of the American Statistical Association*, 113(523):1228–1242.