Econometrics October 12, 2023

Topic 20: Random Forest

by Sai Zhang

Key points: .

Disclaimer: The note is built on Prof. Jinchi Lv's lectures of the course at USC, DSO 607, High-Dimensional Statistics and Big Data Problems.

20.1 Motivation

Denote by $m(\mathbf{X})$ the measurable nonparametric regression function with p-dimensional random vector \mathbf{X} taking values in $[0,1]^p$. The Random Forest algorithm aims to learn the regression function in a nonparametric way based on the observations $\mathbf{x}_i \in [0,1]^p$, $y_i \in \mathbb{R}$, $i = 1, \dots, n$, from the model

$$y_i = m(\mathbf{x}_i) + \epsilon_i$$

where X, x_i , ε_i , $i = 1, \dots, n$ are independent, and $\{x_i\}$ and $\{\varepsilon_i\}$ are two sequences of identically distributed random variables. x_i is distributed identically as X.

Why Random Forest (RF)? RF has gained significant popularity due to its

- High accuracy: RF consistently rank among the top performer, often surpassing more complex models
- Robustness: RF are less subject to overfitting due to the ensemble nature leveraging multiple decision trees
- Interpretability: RF provide rankings of feature importance

As illustrated in Figure 20.1, in a level-2 tree, each node (cell) defines the point where the current cell split and new cells are produced. The sets of features eligible for splitting cells at level k-1 are denoted as $\Theta_k := \{\Theta_{k,1}, \cdots, \Theta_{k,2^{k-1}}\}$, where $\Theta_{k,s} \subset \{1, \cdots, p\}$.



Figure 20.1: Level-2 Tree Example

Given any T (and the associated splitting criterion) and $\Theta_{1:k}$, the tree estimate denoted as $\hat{m}_{T(\Theta_{1:k})}$ for a test

20-2 Week 20: Random Forest

point $\mathbf{c} \in [0,1]^p$ is defined as

$$\hat{m}_{T(\Theta_{1:k})}(\mathbf{c}, \mathcal{X}_n) := \sum_{(\mathbf{t}_1, \dots, \mathbf{t}_k) \in T(\Theta_{1:k})} \mathbf{1}_{\mathbf{c} \in \mathbf{t}_k} \left(\frac{\sum_{i \in \{i: \mathbf{x}_i \in \mathbf{t}_k\}} y_i}{\# \{i: \mathbf{x}_i \in \mathbf{t}_k\}} \right)$$

where $X_n := \{\mathbf{x}_i, y_i\}_{i=1}^n$ the fraction is defined as 0 when no sample is in the cell \mathbf{t}_k , and $\mathbf{1}_{\mathbf{c} \in \mathbf{t}_k}$ is an indicator function = 1 if $\mathbf{c} \in \mathbf{t}_k$ and = 0 otherwise.

20.2 Chi et al. (2022): High Dimensional RFs

Following Chi et al. (2022), for a RF model where

- a sequence of distinct $\Theta_{1:k}$ results in a distinct tree
- every set of available features $\Theta_{l,s}$, $l=1,\cdots,k$; $s=1,\cdots,2^{l-1}$

Column subsampling Define a **column subsampling** procedure: $\Theta_{l,s}$, $\forall l,s$ has $[\gamma_0 p]$ distinct integers among $1, \dots, p$, with $[\cdot]$ the ceiling function for some $0 < \gamma_0 \le 1$. γ_0 is the predetermined constant parameter of column subsampling. Introduce the boldface random mappings $\Theta_{1:k}$, which are independent and uniformly distributed over all possible $\Theta_{1:k}$ for all integer k. Then random forests estimate for \mathbf{c} with observations X_n is given by

$$\mathbb{E}\left(\hat{m}_{T(\boldsymbol{\Theta}_{1:k})}\left(\mathbf{c}, \mathcal{X}_{n}\right) \mid \mathcal{X}_{n}\right) = \sum_{\boldsymbol{\Theta}_{1:k}} \mathbb{P}\left(\bigcap_{s=1}^{k} \left\{\boldsymbol{\Theta}_{s} = \boldsymbol{\Theta}_{s}\right\}\right) \hat{m}_{T(\boldsymbol{\Theta}_{1:k})}\left(\mathbf{c}, \mathcal{X}_{n}\right)$$

The expectation is taken over sets of available features.

Observation resampling Let $A = \{a_1, \dots, a_B\}$ be a set of subsamples with each a_i consisting of $\lceil bn \rceil$ observations (indices) drawn without replacement from $\{1, \dots, n\}$ for some positive integer B and $0 < b \le 1$; in addition, each a_i is independent of model training. The default values of B and B are 500 and 0.632 1 . Then the tree estimate using subsample B is define as

$$\hat{m}_{T(\Theta_{1:k}),a}\left(\mathbf{c},\mathcal{X}_{n}\right) \coloneqq \sum_{\left(\mathbf{t}_{1},\cdots,\mathbf{t}_{k}\right)\in T(\Theta_{1:k})} \mathbf{1}_{\mathbf{c}\in\mathbf{t}_{k}} \left(\frac{\sum_{i\in a\cap\left\{i:\mathbf{x}_{i}\in\mathbf{t}_{k}\right\}} y_{i}}{\#\left(a\cap\left\{i:\mathbf{x}_{i}\in\mathbf{t}_{k}\right\}\right)}\right)$$

the random forests estimate given A is then

$$B^{-1} \sum_{a \in A} \mathbb{E} \left[\hat{m}_{T,a} \left(\mathbf{\Theta}_{1:k}, \mathbf{c}, \mathcal{X}_n \right) \mid \mathcal{X}_n \right] := B^{-1} \sum_{a \in A} \mathbb{E} \left[\hat{m}_{T(\mathbf{\Theta}_{1:k}),a} \left(\mathbf{c}, \mathcal{X}_n \right) \mid \mathcal{X}_n \right]$$

CART-split criterion Given a cell t, a subset of observation indices a and a set of available features $\Theta \subset \{1, \dots, p\}$, the CART-split is defined as

$$(\hat{j}, \hat{c}) = \arg \min_{j \in \Theta, c \in \{x_{ij}: \mathbf{x}_i \in \mathbf{t}, i \in a\}} \left[\sum_{i \in a \cap P_L} (\overline{y}_L - y_i)^2 + \sum_{i \in a \cap P_R} (\overline{y}_R - y_i)^2 \right]$$
(20.1)

 $^{^{1}}$ Or, b = 1 but observations are drawn with replacement.

Week 20: Random Forest 20-3

where

$$P_L := \left\{ i : \mathbf{x}_i \in \mathbf{t}, x_{ij} < c \right\}$$

$$\overline{y}_L := \sum_{i \in a \cap P_L} \frac{y_i}{\#(a \cap P_L)}$$

$$P_R := \left\{ i : \mathbf{x}_i \in \mathbf{t}, x_{ij} \ge c \right\}$$

$$\overline{y}_R := \sum_{i \in a \cap P_R} \frac{y_i}{\#(a \cap P_R)}$$

The CART-split criterion conditional on the sample is a deterministic splitting criterion; conditioning on another sample leads to another deterministic splitting criterion. Define \hat{T}_a as the sample tree growing rule that is associated with a splitting criterion following Eq. (20.1), the tree estimates using \hat{T}_a can be similarly defined as

$$\hat{m}_{\hat{T}_a(\Theta_{1:k})}(\mathbf{c}, \mathcal{X}_n) := \sum_{(\mathbf{t}_1, \dots, \mathbf{t}_k) \in \hat{T}_a(\Theta_{1:k})} \mathbf{1}_{\mathbf{c} \in \mathbf{t}_k} \left(\frac{\sum_{i \in \{i: \mathbf{x}_i \in \mathbf{t}_k\}} y_i}{\# \{i: \mathbf{x}_i \in \mathbf{t}_k\}} \right)$$

the definition is the same for $\hat{m}_{\hat{T}_{a,d}}$. Then the random forests estimate for a test point $\mathbf{c} \in [0,1]^p$ is given by

$$B^{-1} \sum_{a \in A} \mathbb{E} \left(\hat{m}_{\hat{T}_a, a} \left(\mathbf{\Theta}_{1:k}, \mathbf{c}, \mathcal{X}_n \right) \mid \mathcal{X}_n \right)$$

where the average and conditional expectation correspond to the sample and column subsamplings respectively, and they are interchangeable.

Bias-variance decomposition For a tree growing rule T and $\Theta_{1:k}$, the population version is defined as

$$m_{T(\Theta_{1:k})}^{*}(\mathbf{c}) := \sum_{(\mathbf{t}_{1}, \dots, \mathbf{t}_{k}) \in T(\Theta_{1:k})} \mathbf{1}_{\mathbf{c} \in \mathbf{t}_{k}} \mathbb{E}\left(m(\mathbf{X}) \mid \mathbf{X} \in \mathbf{t}_{k}\right)$$
(20.2)

for each test point $\mathbf{c} \in [0,1]^p$. And the \mathbb{L}^2 prediction loss for random forests is defined as

$$\mathbb{E}\left[m(\mathbf{X}) - B^{-1} \sum_{a \in A} \mathbb{E}\left(\hat{m}_{\hat{T}_a, a}\left(\mathbf{\Theta}_{1:k}, \mathbf{X}, \mathcal{X}_n\right) \mid \mathbf{X}, \mathcal{X}_n\right)\right]^2$$
(20.3)

if we use the full sample $a = \{1, \dots, n\}$, and denote \hat{T}_a and $\hat{m}_{\hat{T}_a, a}$ as \hat{T} and $\hat{m}_{\hat{T}}$, the sample subsampling and average $B^{-1} \sum_{a \in A} (\cdot)$ in the random forests estimate are no longer needed, then Eq.(20.3) can be simplified as

$$\mathbb{E}\left[m(\mathbf{X}) - \mathbb{E}\left(\hat{m}_{\hat{T}}\left(\mathbf{\Theta}_{1:k}, \mathbf{X}, \mathcal{X}_{n}\right) \mid \mathbf{X}, \mathcal{X}_{n}\right)\right]^{2}$$

By Jensen's inequality and Cauchy-Schwarz inequality,

$$\frac{1}{2}\mathbb{E}\left[m(\mathbf{X}) - \mathbb{E}\left(\hat{m}_{\hat{T}}\left(\mathbf{\Theta}_{1:k}, \mathbf{X}, \mathcal{X}_{n}\right) \mid \mathbf{X}, \mathcal{X}_{n}\right)\right]^{2}$$

$$\leq \mathbb{E}\left[m(\mathbf{X}) - m_{\hat{T}}^{*}\left(\mathbf{\Theta}_{1:k}, \mathbf{X}\right)\right]^{2} + \mathbb{E}\left[m_{\hat{T}}^{*}\left(\mathbf{\Theta}_{1:k}, \mathbf{X}\right) - \hat{m}_{\hat{T}}\left(\mathbf{\Theta}_{1:k}, \mathbf{X}, \mathcal{X}_{n}\right)\right]^{2}$$
approximation error (squared bias) estimation variance

Consistency of RF Models

For a cell t and its two daughter cells t' and t", define

$$\begin{split} (\mathbb{I})_{\mathbf{t},\mathbf{t}'} &:= \mathbb{P}\left(X \in \mathbf{t}' \mid X \in \mathbf{t}\right) \operatorname{Var}\left(m(X) \mid X \in \mathbf{t}'\right) + \mathbb{P}\left(X \in \mathbf{t}'' \mid X \in \mathbf{t}\right) \operatorname{Var}\left(m(X) \mid X \in \mathbf{t}''\right) \\ (\mathbb{II})_{\mathbf{t},\mathbf{t}'} &:= \mathbb{P}\left(X \in \mathbf{t}' \mid X \in \mathbf{t}\right) \left[\mathbb{E}(m(X) \mid X \in \mathbf{t}') - \mathbb{E}(m(X) \mid X \in \mathbf{t})\right]^2 \\ &+ \mathbb{P}\left(X \in \mathbf{t}'' \mid X \in \mathbf{t}\right) \left[\mathbb{E}(m(X) \mid X \in \mathbf{t}'') - \mathbb{E}(m(X) \mid X \in \mathbf{t})\right]^2 \end{split}$$

20-4 Week 20: Random Forest

and $(\mathbb{I})_{t,t''}$ and $(\mathbb{II})_{t,t''}$ are defined similarly.

in this context, we assume the following regularity conditions:

- **Absolutely continuous distribution**: f, the density function of X, is bounded away from 0 and ∞
- Covariates and model errors: assume $p = O(n^{K_0})$ for $K_0 > 0$, and there is a symmetric distribution around 0 for ϵ_1 , s.t. $\mathbb{E} |\epsilon_1|^q < \infty$ for sufficiently large q > 0
- **Bounded regression functions**: $\sup_{\mathbf{c} \in [0,1]^p} |m(\mathbf{c})| \le M_0$, for some $M_0 > 0$
- **Sufficient impurity decrease**: $\exists \alpha_1 \ge 1 \text{ s.t. } \forall \mathbf{t} = t_1 \times \cdots \times t_p$,

$$\operatorname{Var}\left[m(\mathbf{X}) \mid \mathbf{X} \in \mathbf{t}\right] \leq \alpha_1 \sup_{j \in \{1, \dots, p\}, c \in t_j} (\mathbb{II})_{\mathbf{t}, \mathbf{t}(j, c)}$$

where

- $-(\mathbb{I})_{t,t'}$: conditional bias decrease (or conditional impurity decrease)
- $Var[m(X) \mid X \in t]$: conditional *total* bias, $Var[m(X) \mid X \in t] = (\mathbb{I})_{t,t'} + (\mathbb{I})_{t,t'}$
- Intuition: having large conditional bias decrease on each cell is a desired property for achieving a good control of the squared bias of random forests estimate

Sufficient impurity decrease (SID) Define the functional class

$$SID(\alpha) := \{m(\mathbf{X}) : m(\mathbf{X}) \text{ satisfies SID with } \alpha_1 \leq \alpha \}$$

the size of $SID(\alpha)$ is **non-decreasing** in $\alpha \ge 1$: if $m(\mathbf{X}) \in SID(\alpha - c)$ for some $\alpha - c \ge 1$ and c > 0, then $m(\mathbf{X}) \in SID(\alpha)^2$.

²Many popular regression functions belong to this functional class.

Week 20: Random Forest 20-5

References

Chien-Ming Chi, Patrick Vossler, Yingying Fan, and Jinchi Lv. Asymptotic properties of high-dimensional random forests. *The Annals of Statistics*, 50(6):3415–3438, 2022.