

Analysis of purely random forests

Sylvain Arlot^{1,2} (joint work with Robin Genuer³)

¹Université Paris-Saclay

²Inria Saclay, Celeste project-team

³ISPED, Université de Bordeaux

Colloquium, Department of Statistics and Actuarial Science,
The University of Iowa
March 18, 2021

References: [arXiv:1407.3939](https://arxiv.org/abs/1407.3939) (v2 upcoming!)
[arXiv:1604.01515](https://arxiv.org/abs/1604.01515)

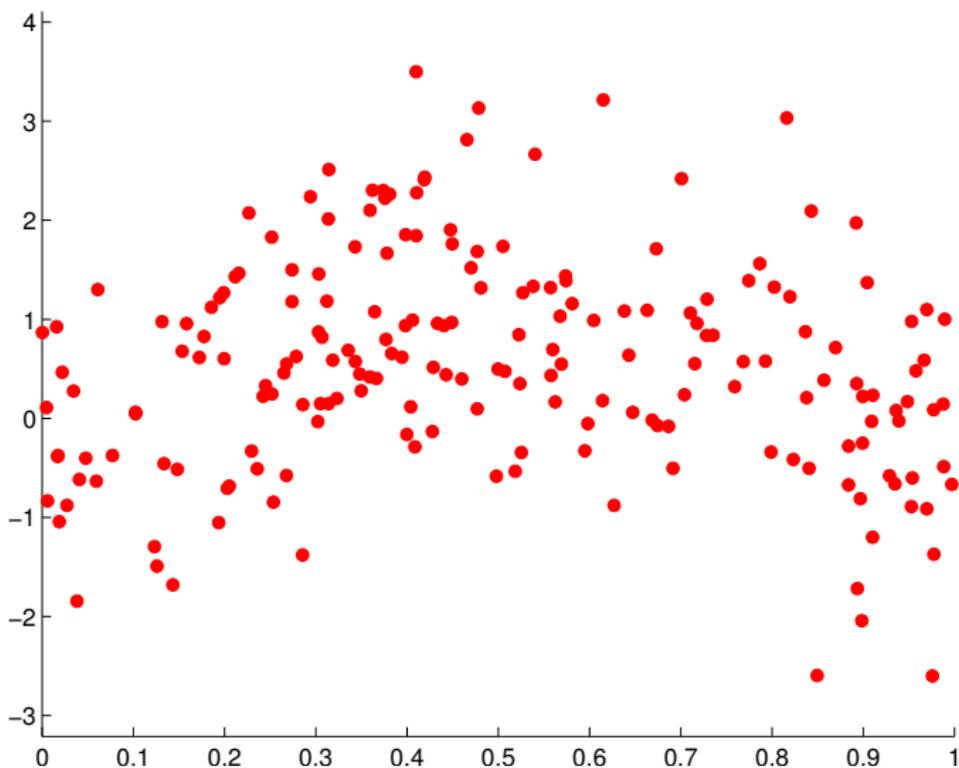
Outline

- 1 Random forests
- 2 Purely random forests
- 3 Toy forests
- 4 Hold-out random forests

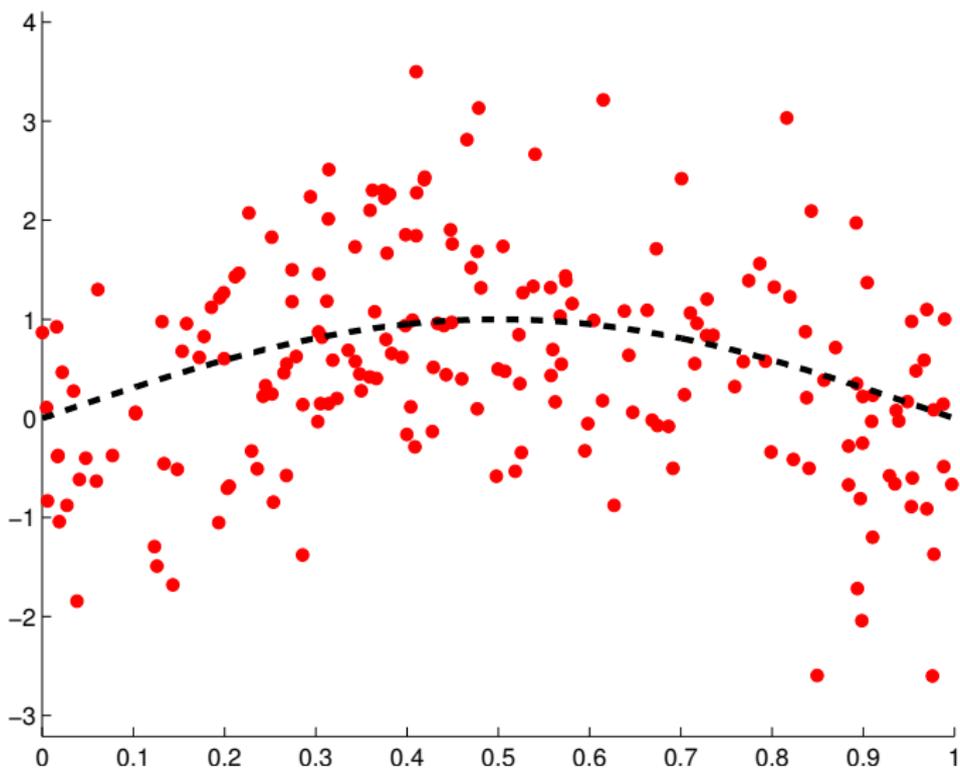
Outline

- 1 Random forests
- 2 Purely random forests
- 3 Toy forests
- 4 Hold-out random forests

Regression: data $(X_1, Y_1), \dots, (X_n, Y_n)$



Goal: find the signal (denoising)



Regression

- **Data** $D_n: (X_1, Y_1), \dots, (X_n, Y_n) \in \mathbb{R}^p \times \mathbb{R}$ (i.i.d. $\sim P$)

$$Y_i = s^*(X_i) + \varepsilon_i$$

with $s^*(X) = \mathbb{E}[Y | X]$ (regression function).

Regression

- **Data** D_n : $(X_1, Y_1), \dots, (X_n, Y_n) \in \mathbb{R}^p \times \mathbb{R}$ (i.i.d. $\sim P$)

$$Y_i = s^*(X_i) + \varepsilon_i$$

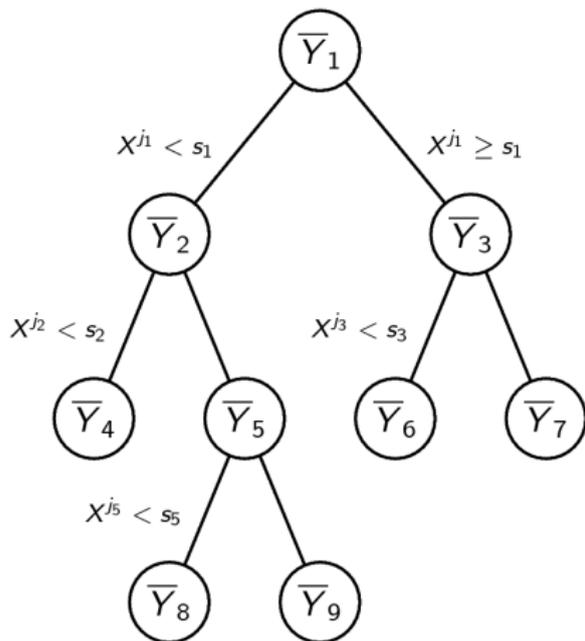
with $s^*(X) = \mathbb{E}[Y | X]$ (regression function).

- **Goal**: learn f measurable function $\mathcal{X} \rightarrow \mathbb{R}$ s.t. **the quadratic risk**

$$\mathbb{E}_{(X,Y) \sim P} \left[(f(X) - s^*(X))^2 \right]$$

is minimal.

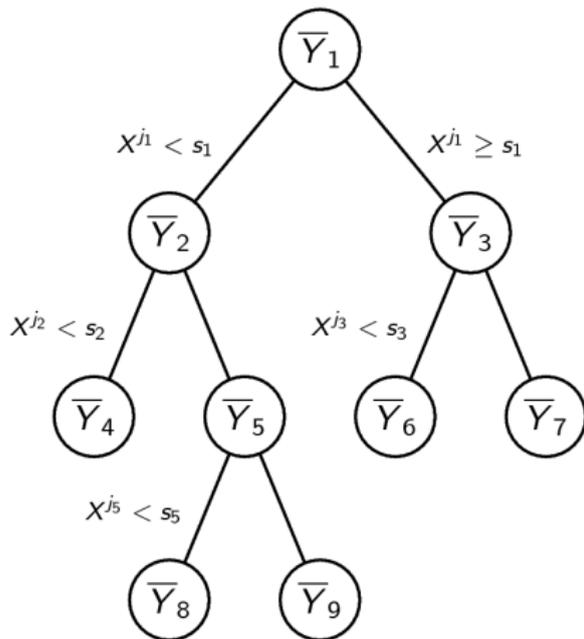
Regression tree (Breiman et al, 1984)



Tree: piecewise-constant predictor, obtained by partitioning recursively \mathbb{R}^p .

Restriction: splits parallel to the axes.

Regression tree (Breiman et al, 1984)

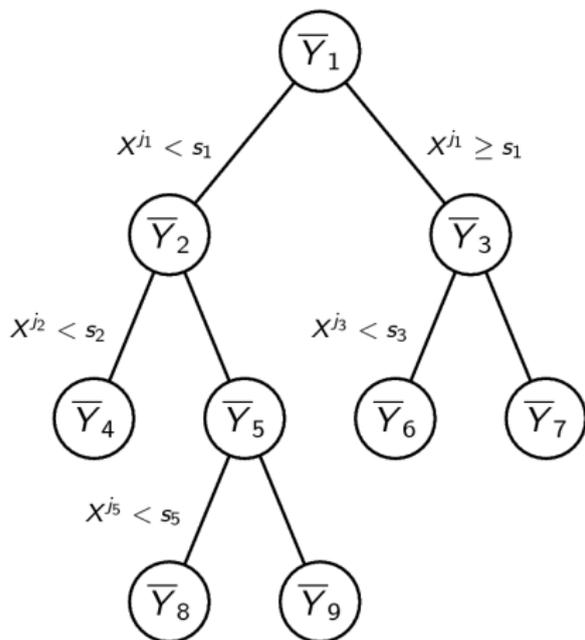


Tree: piecewise-constant predictor, obtained by partitioning recursively \mathbb{R}^p .

Restriction: splits parallel to the axes.

- 1 **Choice of the partition \mathcal{U}** (tree structure)
Usually, at each step, one looks for the best split of the data into two groups (minimize sum of within-group variances) D_n .

Regression tree (Breiman et al, 1984)



Tree: piecewise-constant predictor, obtained by partitioning recursively \mathbb{R}^p .

Restriction: splits parallel to the axes.

- 1 Choice of the partition \mathcal{U} (tree structure)
- 2 For each $\lambda \in \mathcal{U}$ (tree leaf), choice of the estimation $\hat{\beta}_\lambda$ of $s^*(x)$ when $x \in \lambda$. Here, $\hat{\beta}_\lambda = \bar{Y}_\lambda$ average of the $(Y_i)_{X_i \in \lambda}$.

Random forest (Breiman, 2001)

Definition (Random forest (Breiman, 2001))

$\{\hat{s}_{\Theta_j}, 1 \leq j \leq q\}$ collection of tree predictors, $(\Theta_j)_{1 \leq j \leq q}$ i.i.d. r.v. independent from D_n .

Random forest predictor \hat{s} obtained by **aggregating the tree collection**.

$$\hat{s}(x) = \frac{1}{q} \sum_{j=1}^q \hat{s}_{\Theta_j}(x)$$

- ensemble method (Dietterich, 1999, 2000)
- powerful **statistical learning** algorithm, for both **classification** and **regression**.

Bagging (“bootstrap aggregating”)

- **Bootstrap** (Efron, 1979): draw n i.i.d. r.v., uniform over $\{(X_i, Y_i) / i = 1, \dots, n\}$ (sampling with replacement)
 \Rightarrow **resample** D_n^b
- Bootstrapping a tree: $\hat{s}_{\text{tree}}^b = \hat{s}_{\text{tree}}(D_n^b)$
- **Bagging**: bootstrap (q independent resamples) then aggregation

$$\hat{s}_{\text{bagging}}(x) = \frac{1}{q} \sum_{j=1}^q \hat{s}_{\text{tree}}^{b,j}(x)$$

Random Forest-Random Inputs (Breiman, 2001)

Definition (RI tree)

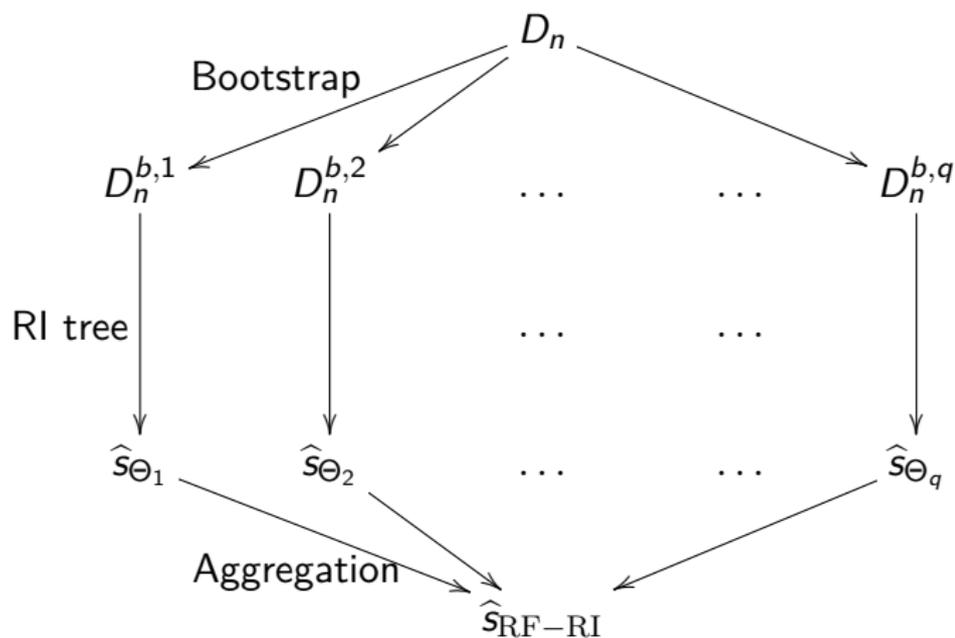
In a RI tree, at each node, **mtry** variables are randomly chosen. Then, the best cut direction is chosen only among the chosen variables.

Definition (Random forest RI)

A random forest RI (RF-RI) is obtained by **aggregating RI trees** built on independent **bootstrap resamples**.

RF-RI \Leftrightarrow bagging on RI trees

Random Forest-Random Inputs



Theoretical results on RF-RI

- Few theoretical results on Breiman's original RF-RI, despite their excellent numerical performance (eg, Fernández-Delgado et al, 2014)
- Most results:
 - focus on a **specific part** of the algorithm (resampling, split criterion),
 - **modify** the algorithm (eg, subsampling instead of resampling)
 - make **strong assumptions** on s^*
- References (see **survey paper** by Biau and Scornet, 2016): Scornet, Biau & Vert (2015), Mentch & Hooker (2016), Wager & Athey (2018), Genuer & Poggi (2019), ...

Theoretical results on RF-RI

- Few theoretical results on Breiman's original RF-RI, despite their excellent numerical performance (eg, Fernández-Delgado et al, 2014)
 - Most results:
 - focus on a **specific part** of the algorithm (resampling, split criterion),
 - **modify** the algorithm (eg, subsampling instead of resampling)
 - make **strong assumptions** on s^*
 - References (see **survey paper** by Biau and Scornet, 2016): Scornet, Biau & Vert (2015), Mentch & Hooker (2016), Wager & Athey (2018), Genuer & Poggi (2019), ...
- ⇒ Here, we consider simplified RF models, for which a precise analysis is possible: **purely random forests**

Outline

- 1 Random forests
- 2 Purely random forests
- 3 Toy forests
- 4 Hold-out random forests

Purely random forests

Definition (Purely random tree)

$$\hat{s}_{\mathbb{U}}(x) = \sum_{\lambda \in \mathbb{U}} \overline{Y}_{\lambda}(D_n) \mathbb{1}_{x \in \lambda}$$

where $\overline{Y}_{\lambda}(D_n)$ is the average of $(Y_i)_{X_i \in \lambda, (X_i, Y_i) \in D_n}$ and the partition \mathbb{U} is independent from D_n .

Definition (Purely random forest)

$$\hat{s}(x) = \frac{1}{q} \sum_{j=1}^q \hat{s}_{\mathbb{U}^j}(x)$$

with $\mathbb{U}^1, \dots, \mathbb{U}^q$ i.i.d., independent from D_n .

Purely random forests

Definition (Purely random forest)

$$\hat{s}(x) = \frac{1}{q} \sum_{j=1}^q \hat{s}_{\mathbb{U}^j}(x) = \frac{1}{q} \sum_{j=1}^q \sum_{\lambda \in \mathbb{U}^j} \bar{Y}_{\lambda}(D_n) \mathbb{1}_{x \in \lambda}$$

with $\mathbb{U}^1, \dots, \mathbb{U}^q$ i.i.d., independent from D_n .

Example (“hold-out RF” model): use some **extra data** D'_n for building the trees: $\mathbb{U}^j = \mathbb{U}_{\text{RI}}(D_n^{*j})$ (can be done by splitting the sample into two subsamples D_n and D'_n).

Purely random forests

Definition (Purely random forest)

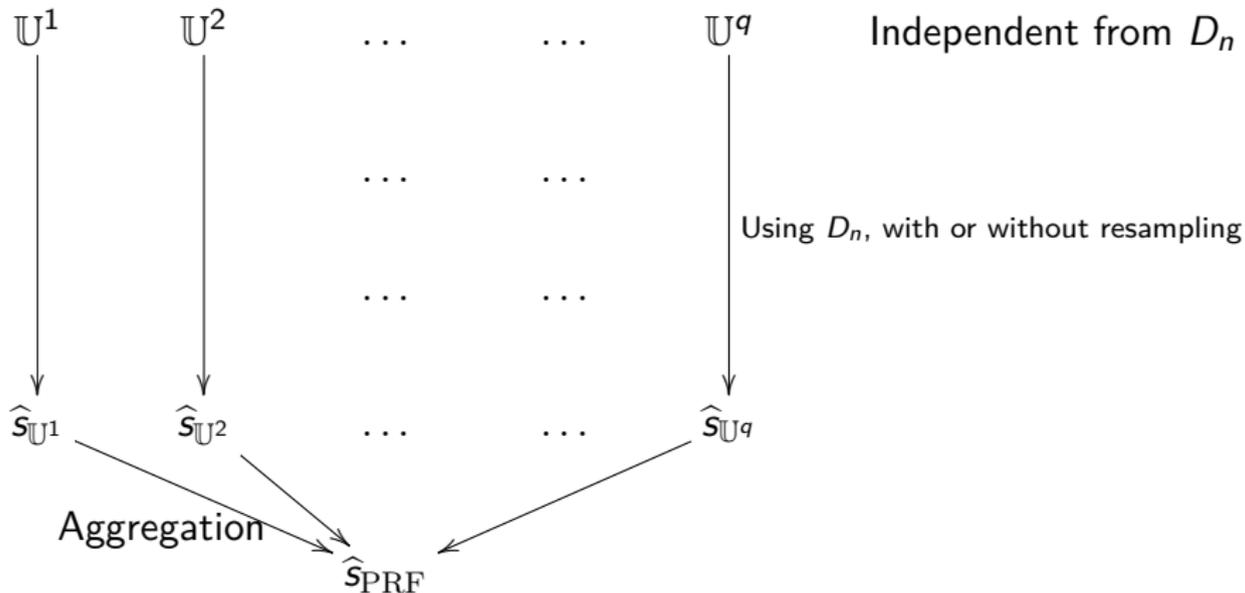
$$\hat{s}(x) = \frac{1}{q} \sum_{j=1}^q \hat{s}_{\mathbb{U}^j}(x) = \frac{1}{q} \sum_{j=1}^q \sum_{\lambda \in \mathbb{U}^j} \bar{Y}_{\lambda}(D_n) \mathbb{1}_{x \in \lambda}$$

with $\mathbb{U}^1, \dots, \mathbb{U}^q$ i.i.d., independent from D_n .

Example (“hold-out RF” model): use some **extra data** D'_n for building the trees: $\mathbb{U}^j = \mathbb{U}_{\text{RI}}(D_n^{*j})$ (can be done by splitting the sample into two subsamples D_n and D'_n).

 From now on, D_n is the sample used for computing $(Y_{\lambda}(D_n))_{\lambda \in \mathbb{U}}$, and we assume its size is n .

Purely random forests



Purely random forests: theory

- **Consistency**: Biau, Devroye & Lugosi (2008), Scornet (2014)
- **Rates of convergence**: Breiman (2004), Biau (2012), Klusowski (2018), Duroux & Scornet (2018), Mourtada, Gaiffas & Scornet (2017 & 2020)
- Some adaptivity to **dimension reduction** (sparse framework): Biau (2012), Klusowski (2018)
- Forests **decrease the estimation error** (Biau, 2012; Genuer, 2012)

Purely random forests: theory

- **Consistency**: Biau, Devroye & Lugosi (2008), Scornet (2014)
 - **Rates of convergence**: Breiman (2004), Biau (2012), Klusowski (2018), Duroux & Scornet (2018), Mourtada, Gaiffas & Scornet (2017 & 2020)
 - Some adaptivity to **dimension reduction** (sparse framework): Biau (2012), Klusowski (2018)
 - Forests **decrease the estimation error** (Biau, 2012; Genuer, 2012)
- ⇒ What about **approximation error**?
Almost the same for a forest and a tree?

Risk of a single tree (regressogram)

Given the partition \mathbb{U} , regressogram estimator

$$\hat{s}_{\mathbb{U}}(x) := \sum_{\lambda \in \mathbb{U}} \bar{Y}_{\lambda} \mathbb{1}_{x \in \lambda}$$

where \bar{Y}_{λ} is the average of $(Y_i)_{X_i \in \lambda}$.

$$\hat{s}_{\mathbb{U}} \in \operatorname{argmin}_{f \in \mathcal{S}_{\mathbb{U}}} \left\{ \frac{1}{n} \sum_{i=1}^n (Y_i - f(X_i))^2 \right\}$$

where $\mathcal{S}_{\mathbb{U}}$ is the vector space of functions which are constant over each $\lambda \in \mathbb{U}$.

Risk of a single tree (regressogram)

Given the partition \mathbb{U} , regressogram estimator

$$\hat{s}_{\mathbb{U}}(x) := \sum_{\lambda \in \mathbb{U}} \bar{Y}_{\lambda} \mathbb{1}_{x \in \lambda}$$

where \bar{Y}_{λ} is the average of $(Y_i)_{X_i \in \lambda}$.

$$\hat{s}_{\mathbb{U}} \in \operatorname{argmin}_{f \in \mathcal{S}_{\mathbb{U}}} \left\{ \frac{1}{n} \sum_{i=1}^n (Y_i - f(X_i))^2 \right\}$$

where $\mathcal{S}_{\mathbb{U}}$ is the vector space of functions which are constant over each $\lambda \in \mathbb{U}$.

Define:

$$\tilde{s}_{\mathbb{U}}(x) := \sum_{\lambda \in \mathbb{U}} \beta_{\lambda} \mathbb{1}_{x \in \lambda} \quad \text{where } \beta_{\lambda} := \mathbb{E}[s^*(X) \mid X \in \lambda] .$$

$$\Rightarrow \tilde{s}_{\mathbb{U}} \in \operatorname{argmin}_{f \in \mathcal{S}_{\mathbb{U}}} \mathbb{E} \left[(f(X) - s^*(X))^2 \right] \quad \text{and} \quad \tilde{s}_{\mathbb{U}}(x) = \mathbb{E}[\hat{s}_{\mathbb{U}}(x) \mid \mathbb{U}]$$

Risk decomposition: single tree

$$\begin{aligned} & \mathbb{E}\left[(s^*(X) - \widehat{s}_U(X))^2\right] \\ &= \mathbb{E}\left[(s^*(X) - \widetilde{s}_U(X))^2\right] + \mathbb{E}\left[(\widetilde{s}_U(X) - \widehat{s}_U(X))^2\right] \\ &= \text{Approximation error} + \text{Estimation error} \end{aligned}$$

Risk decomposition: single tree

$$\begin{aligned} & \mathbb{E} \left[(s^*(X) - \widehat{s}_{\mathbb{U}}(X))^2 \right] \\ &= \mathbb{E} \left[(s^*(X) - \widetilde{s}_{\mathbb{U}}(X))^2 \right] + \mathbb{E} \left[(\widetilde{s}_{\mathbb{U}}(X) - \widehat{s}_{\mathbb{U}}(X))^2 \right] \\ &= \text{Approximation error} + \text{Estimation error} \end{aligned}$$

If s^* is smooth, $X \sim \mathcal{U}([0, 1]^p)$ and \mathbb{U} regular partition into D pieces, then

$$\mathbb{E} \left[(s^*(X) - \widetilde{s}_{\mathbb{U}}(X))^2 \right] \propto \frac{1}{D^{2/p}}$$

Risk decomposition: single tree

$$\begin{aligned} & \mathbb{E} \left[(s^*(X) - \widehat{s}_{\mathbb{U}}(X))^2 \right] \\ &= \mathbb{E} \left[(s^*(X) - \widetilde{s}_{\mathbb{U}}(X))^2 \right] + \mathbb{E} \left[(\widetilde{s}_{\mathbb{U}}(X) - \widehat{s}_{\mathbb{U}}(X))^2 \right] \\ &= \text{Approximation error} + \text{Estimation error} \end{aligned}$$

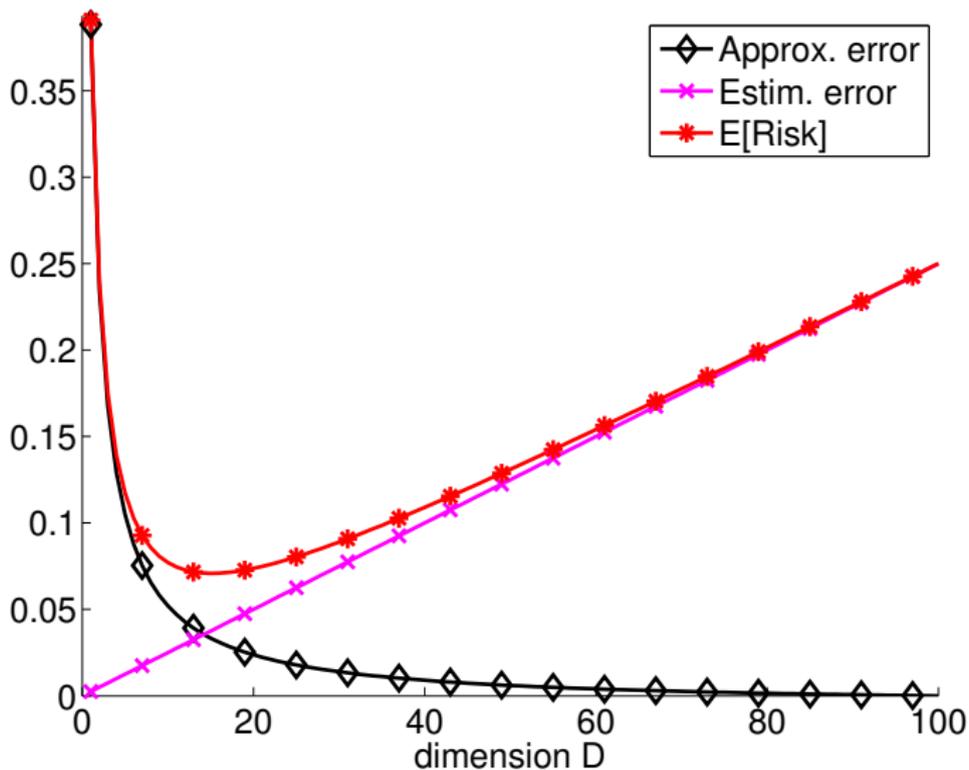
If s^* is smooth, $X \sim \mathcal{U}([0, 1]^p)$ and \mathbb{U} regular partition into D pieces, then

$$\mathbb{E} \left[(s^*(X) - \widetilde{s}_{\mathbb{U}}(X))^2 \right] \propto \frac{1}{D^{2/p}}$$

If $\text{var}(Y | X) = \sigma^2$ does not depend on X , then

$$\mathbb{E} \left[(\widetilde{s}_{\mathbb{U}}(X) - \widehat{s}_{\mathbb{U}}(X))^2 \right] \approx \frac{\sigma^2 D}{n}$$

Approximation and estimation errors, $p = 1$



Risk decomposition: purely random forest

$(\mathbb{U}^j)_{1 \leq j \leq q}$ finite partitions, i.i.d. $\sim \mathcal{U}$

Estimator (forest): $\hat{s}_{\mathbb{U}^{1 \dots q}}(x) := \frac{1}{q} \sum_{j=1}^q \hat{s}_{\mathbb{U}^j}(x)$

Ideal forest: $\tilde{s}_{\mathbb{U}^{1 \dots q}}(x) := \frac{1}{q} \sum_{j=1}^q \tilde{s}_{\mathbb{U}^j}(x) = \mathbb{E}[\hat{s}_{\mathbb{U}^{1 \dots q}}(x) \mid \mathbb{U}^{1 \dots q}]$

Risk decomposition: purely random forest

$(\mathbb{U}^j)_{1 \leq j \leq q}$ finite partitions, i.i.d. $\sim \mathcal{U}$

Estimator (forest): $\hat{s}_{\mathbb{U}^{1 \dots q}}(x) := \frac{1}{q} \sum_{j=1}^q \hat{s}_{\mathbb{U}^j}(x)$

Ideal forest: $\tilde{s}_{\mathbb{U}^{1 \dots q}}(x) := \frac{1}{q} \sum_{j=1}^q \tilde{s}_{\mathbb{U}^j}(x) = \mathbb{E}[\hat{s}_{\mathbb{U}^{1 \dots q}}(x) \mid \mathbb{U}^{1 \dots q}]$

Quadratic risk decomposition (given $X = x$)

$$\begin{aligned} \mathbb{E} \left[(s^*(x) - \hat{s}_{\mathbb{U}^{1 \dots q}}(x))^2 \right] &= \mathbb{E} \left[(s^*(x) - \tilde{s}_{\mathbb{U}^{1 \dots q}}(x))^2 \right] \\ &\quad + \mathbb{E} \left[(\tilde{s}_{\mathbb{U}^{1 \dots q}}(x) - \hat{s}_{\mathbb{U}^{1 \dots q}}(x))^2 \right] + \delta_{\mathbb{U}^{1 \dots q}}(x) \end{aligned}$$

Approximation error: $B_{\mathcal{U}, q}(x) := \mathbb{E} \left[(s^*(x) - \tilde{s}_{\mathbb{U}^{1 \dots q}}(x))^2 \right]$

Approximation error decomposition (given $X = x$)

$$\mathcal{B}_{U,q}(x) = \mathcal{B}_{U,\infty}(x) + \frac{\mathcal{V}_U(x)}{q}$$

$$\text{where } \mathcal{B}_{U,\infty}(x) := \left(\mathbb{E}[s^*(x) - \tilde{s}_U(x)] \right)^2$$

$$\text{and } \mathcal{V}_U(x) := \text{var}(\tilde{s}_U(x))$$

$\mathcal{B}_{U,\infty}(x)$ is the **approx. error of the infinite forest**: $\tilde{s}_{U,\infty}(x) := \mathbb{E}[\tilde{s}_U(x)]$

to be compared with the **approximation error of a single tree**

$$\mathcal{B}_{U,1}(x) = \mathcal{B}_{U,\infty}(x) + \mathcal{V}_U(x)$$

Outline

- 1 Random forests
- 2 Purely random forests
- 3 Toy forests**
- 4 Hold-out random forests

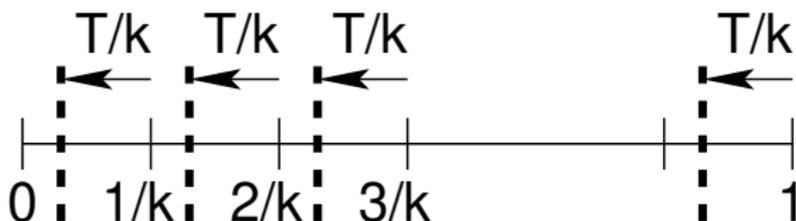
Toy forests

Assume: $\mathcal{X} = [0, 1)^p$ and X uniform over $[0, 1)^p$

If $p = 1$, $\mathbb{U} \sim \mathcal{U}_k^{\text{toy}}$ defined by:

$$\mathbb{U} = \left\{ \left[0, \frac{1-T}{k} \right), \left[\frac{1-T}{k}, \frac{2-T}{k} \right), \dots, \left[\frac{k-T}{k}, 1 \right) \right\}$$

where T has uniform distribution over $[0, 1]$.



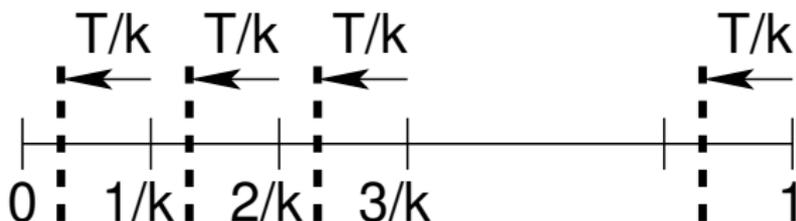
Toy forests

Assume: $\mathcal{X} = [0, 1]^p$ and X uniform over $[0, 1]^p$

If $p = 1$, $\mathbb{U} \sim \mathcal{U}_k^{\text{toy}}$ defined by:

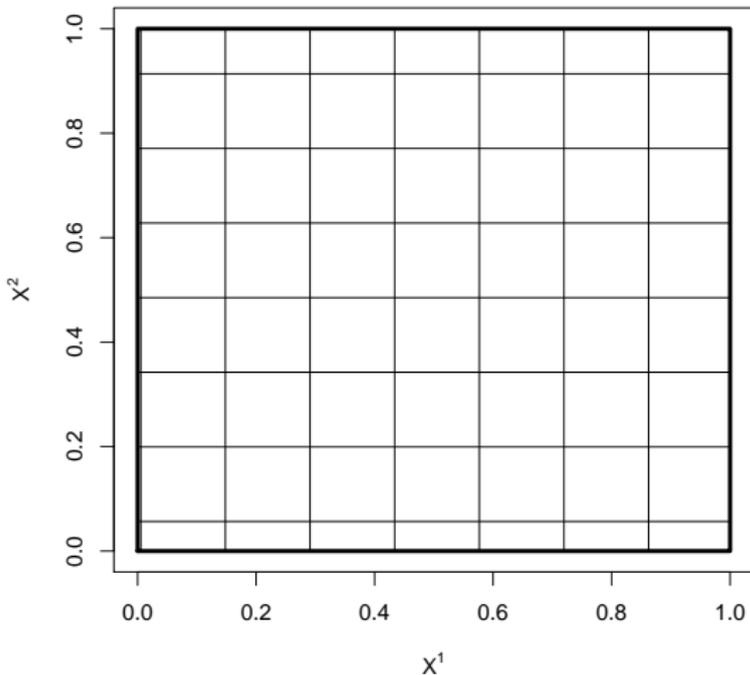
$$\mathbb{U} = \left\{ \left[0, \frac{1-T}{k} \right), \left[\frac{1-T}{k}, \frac{2-T}{k} \right), \dots, \left[\frac{k-T}{k}, 1 \right) \right\}$$

where T has uniform distribution over $[0, 1]$.



If $p > 1$, T_j for each coordinate $j = 1, \dots, p$, independent

Toy forest, $p = 2$: example



Interpretation of the ideal infinite forest ($p = 1$)

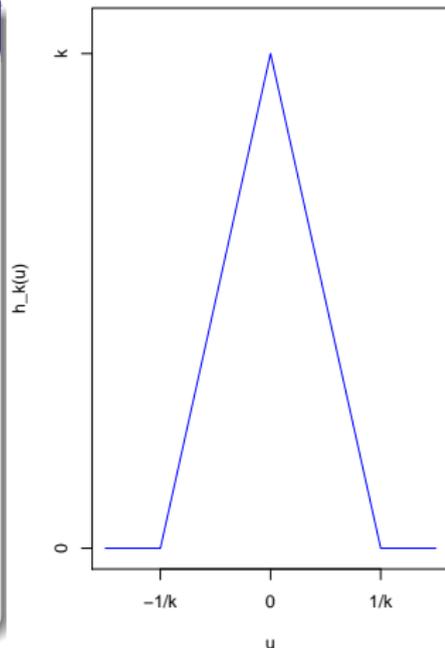
Proposition (A. & Genuer, 2014–2020)

For any $x \in \left[\frac{1}{k}, 1 - \frac{1}{k}\right]$, the ideal infinite forest at x satisfies:

$$\tilde{s}_{\mathbb{U}, \infty}(x) = (s^* * h_k)(x) = \int_0^1 s^*(t) h_k(x - t) dt$$

where

$$h_k(u) = \begin{cases} k(1 - ku) & \text{if } 0 \leq u \leq \frac{1}{k} \\ k(1 + ku) & \text{if } -\frac{1}{k} \leq u \leq 0 \\ 0 & \text{if } |u| \geq \frac{1}{k} \end{cases}$$



Interpretation of the ideal infinite forest ($p = 1$): proof

$I_{\mathbb{U}}(x)$:= the interval of \mathbb{U} to which x belongs

$$\tilde{s}_{\mathbb{U}}(x) = \frac{1}{|I_{\mathbb{U}}(x)|} \int_{I_{\mathbb{U}}(x)} s^*(t) dt$$

If $x \in \left[\frac{1}{k}, 1 - \frac{1}{k}\right]$, $I_{\mathbb{U}}(x) = \left[x + \frac{V_x - 1}{k}, x + \frac{V_x}{k}\right)$

where V_x has uniform distribution over $[0, 1]$.

Interpretation of the ideal infinite forest ($p = 1$): proof

$I_{\mathbb{U}}(x) :=$ the interval of \mathbb{U} to which x belongs

$$\tilde{s}_{\mathbb{U}}(x) = \frac{1}{|I_{\mathbb{U}}(x)|} \int_{I_{\mathbb{U}}(x)} s^*(t) dt$$

If $x \in \left[\frac{1}{k}, 1 - \frac{1}{k}\right]$, $I_{\mathbb{U}}(x) = \left[x + \frac{V_x - 1}{k}, x + \frac{V_x}{k}\right)$

where V_x has uniform distribution over $[0, 1]$.

$$\begin{aligned} \tilde{s}_{\mathbb{U}, \infty}(x) &= \mathbb{E}_{\mathbb{U}}[\tilde{s}_{\mathbb{U}}(x)] \\ &= k \int_0^1 s^*(t) \mathbb{P}\left(x + \frac{V_x - 1}{k} \leq t < x + \frac{V_x}{k}\right) dt \\ &= k \int_0^1 s^*(t) \underbrace{\mathbb{P}(k(t - x) < V_x \leq k(t - x) + 1)}_{=h_k(x-t) \text{ if } 1/k \leq x \leq 1-1/k} dt \end{aligned}$$

Analysis of the approximation error, $p = 1$, $x \in \left[\frac{1}{k}, 1 - \frac{1}{k}\right]$

(H2) s^* twice differentiable over $(0, 1)$ and $s^{*''}$ bounded

Taylor-Lagrange formula: for every $t \in (0, 1)$, some $c_{t,x} \in (0, 1)$ exists such that

$$s^*(t) - s^*(x) = s^{*'}(x)(t - x) + \frac{1}{2}s^{*''}(c_{t,x})(t - x)^2$$

Analysis of the approximation error, $p = 1$, $x \in \left[\frac{1}{k}, 1 - \frac{1}{k}\right]$

(H2) s^* twice differentiable over $(0, 1)$ and $s^{*''}$ bounded

Taylor-Lagrange formula: for every $t \in (0, 1)$, some $c_{t,x} \in (0, 1)$ exists such that

$$s^*(t) - s^*(x) = s^{*'}(x)(t - x) + \frac{1}{2}s^{*''}(c_{t,x})(t - x)^2$$

Therefore,

$$\begin{aligned} \tilde{s}_{\cup}(x) - s^*(x) &= k \int_{x + \frac{V_x - 1}{k}}^{x + \frac{V_x}{k}} (s^*(t) - s^*(x)) dt \\ &= k s^{*'}(x) \int_{x + \frac{V_x - 1}{k}}^{x + \frac{V_x}{k}} (t - x) dt + R_1(x) \\ &= \frac{s^{*'}(x)}{k} \left(V_x - \frac{1}{2} \right) + R_1(x) \end{aligned}$$

where $R_1(x) = \frac{k}{2} \int_{x + \frac{V_x - 1}{k}}^{x + \frac{V_x}{k}} s^{*''}(c_{t,x})(t - x)^2 dt.$

Analysis of the approximation error, $p = 1$, $x \in \left[\frac{1}{k}, 1 - \frac{1}{k}\right]$

(H2) s^* twice differentiable over $(0, 1)$ and $s^{*''}$ bounded

$$\tilde{s}_{\mathbb{U}}(x) - s^*(x) = \frac{s^{*'}(x)}{k} \left(V_x - \frac{1}{2} \right) + R_1(x)$$

where $R_1(x) = \frac{k}{2} \int_{x+\frac{V_x-1}{k}}^{x+\frac{V_x}{k}} s^{*''}(c_{t,x})(t-x)^2 dt$.

Hence,

$$\mathcal{B}_{\mathcal{U}_k^{\text{toy}}, \infty}(x) = \left(\mathbb{E}_{\mathbb{U}}[s^*(x) - \tilde{s}_{\mathbb{U}}(x)] \right)^2 = \left(\mathbb{E}_{\mathbb{U}}[R_1(x)] \right)^2 \leq \frac{\square}{k^4}$$

and

$$\mathcal{V}_{\mathcal{U}_k^{\text{toy}}}(x) = \text{var} \left(\frac{s^{*'}(x)}{k} \left(V_x - \frac{1}{2} \right) + R_1(x) \right) \underset{k \rightarrow +\infty}{\sim} \frac{s^{*'}(x)^2 \text{var}(V_x)}{k^2}.$$

Analysis of the approximation error, $p \geq 1$

(H θ) $s^* \in \mathcal{C}^{1,(\theta-1)}(\mathcal{X})$ Hölder space, $\theta \in [1, 2]$

$$\mathcal{B}_{\mathcal{U}_k^{\text{toy}}, \infty}(x) = \left(\mathbb{E}_{\mathbb{U}}[s^*(x) - \tilde{s}_{\mathbb{U}}(x)] \right)^2 \leq \frac{\square}{k^{2\theta/p}} \quad \mathcal{V}_{\mathcal{U}_k^{\text{toy}}}(x) \underset{k \rightarrow +\infty}{\sim} \frac{\square}{k^{2/p}}$$

Proposition (A. & Genuer, 2014–2020)

Assuming (H θ), $\theta \in [1, 2]$, $\forall x \in [\frac{1}{k}, 1 - \frac{1}{k}]^p$,

$$\mathcal{B}_{\mathcal{U}_k^{\text{toy}}, 1}(x) \underset{k \rightarrow +\infty}{\sim} \frac{\square}{k^{2/p}} \quad \mathcal{B}_{\mathcal{U}_k^{\text{toy}}, \infty}(x) \leq \frac{\square}{k^{2\theta/p}}$$

$$\int_{[\frac{1}{k}, 1 - \frac{1}{k}]^p} \mathcal{B}_{\mathcal{U}_k^{\text{toy}}, 1}(x) dx \underset{k \rightarrow +\infty}{\sim} \frac{\square}{k^{2/p}} \quad \int_{[\frac{1}{k}, 1 - \frac{1}{k}]^p} \mathcal{B}_{\mathcal{U}_k^{\text{toy}}, \infty}(x) dx \leq \frac{\square}{k^{2\theta/p}}$$

Rate $k^{-4/p}$ is tight assuming θ -Hölder smoothness, $\theta > 2$.

Estimation error

General fact (Jensen's inequality):

$$\mathbb{E}\left[\left(\tilde{s}_{U,\infty}(X) - \hat{s}_{U,\infty}(X)\right)^2\right] \leq \mathbb{E}\left[\left(\tilde{s}_U(X) - \hat{s}_U(X)\right)^2\right]$$

Estimation error

General fact (Jensen's inequality):

$$\mathbb{E}\left[\left(\tilde{s}_{U,\infty}(X) - \hat{s}_{U,\infty}(X)\right)^2\right] \leq \mathbb{E}\left[\left(\tilde{s}_U(X) - \hat{s}_U(X)\right)^2\right]$$

For the toy forest, without any resampling for computing labels and assuming that $\text{var}(Y|X) = \sigma^2$:

$$\begin{aligned}\mathbb{E}\left[\left(\tilde{s}_U(X) - \hat{s}_U(X)\right)^2\right] &\approx \frac{\sigma^2 k}{n} \\ \mathbb{E}\left[\left(\tilde{s}_{U,\infty}(X) - \hat{s}_{U,\infty}(X)\right)^2\right] &\approx \frac{2}{3} \frac{\sigma^2 k}{n}\end{aligned}$$

(A. & Genuer, 2016)

Summary: risk analysis

$$\mathbb{E}\left[(s^*(x) - \widehat{s}_{\cup_{1 \dots q}}(x))^2\right] \approx \frac{c_1(s^*, x)}{k^{2/p}} + \frac{\sigma^2 k}{n} \leq \frac{c'_\theta(s^*, x)}{k^{2\theta/p}} + \frac{2\sigma^2 k}{3n}$$

Single tree

$(q = 1)$

Infinite forest

$(q = \infty)$

where $c_1(s^*, x) = \frac{s^{*\prime}(x)^2}{12}$.

Assumptions:

- $x \in (0, 1)^p$ far from boundary
- $(H_\theta) s^* \in \mathcal{C}^{1,(\theta-1)}(\mathcal{X})$, $\theta \in [1, 2]$
- X uniform over $[0, 1]^p$
- $\text{var}(Y|X) = \sigma^2$
- no resampling for computing labels

Rates of convergence

Corollary: risk convergence rates (far from boundaries, with $k = k_n^*$ optimal), under $(H\theta)$, $\theta \in [1, 2]$:

Tree risk $\geq \square n^{-2/(2+p)}$ if s^* not constant, $\theta > 1$

Infinite forest risk $\leq \square n^{-2\theta/(2\theta+p)} \Rightarrow \text{minimax } C^\theta, \theta \in [1, 2]$

Rates of convergence

Corollary: risk convergence rates (far from boundaries, with $k = k_n^*$ optimal), under $(H\theta)$, $\theta \in [1, 2]$:

Tree risk $\geq \square n^{-2/(2+p)}$ if s^* not constant, $\theta > 1$

Infinite forest risk $\leq \square n^{-2\theta/(2\theta+p)} \Rightarrow \text{minimax } \mathcal{C}^\theta, \theta \in [1, 2]$

Remarks:

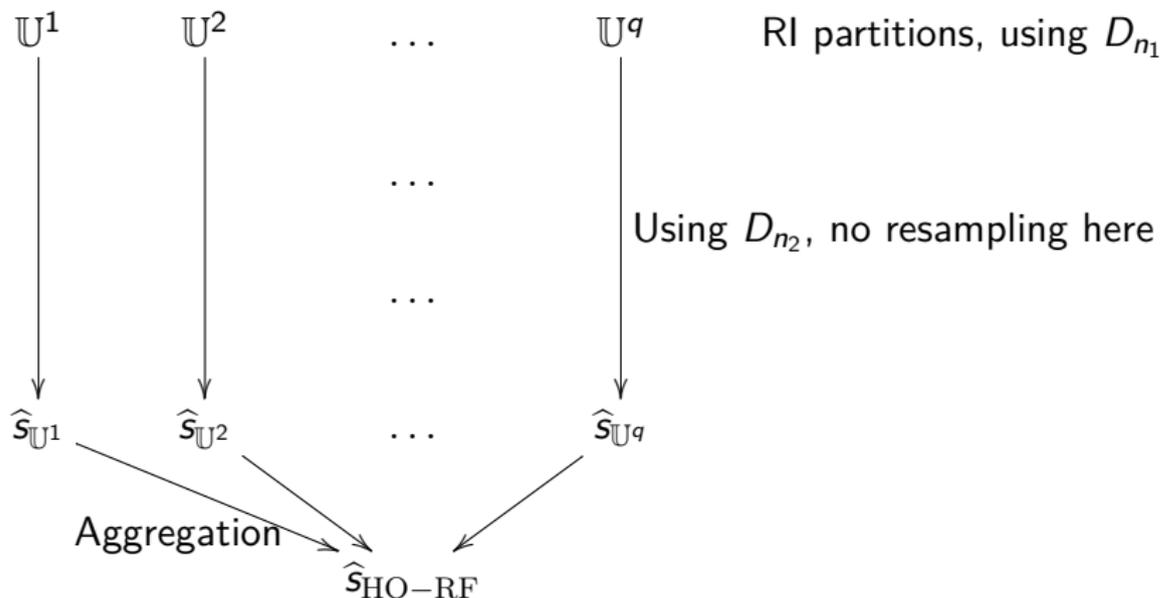
- $q \geq \square (k_n^*)^2$ is sufficient to get an “infinite” forest
- with subsampling a out of n for computing labels: estimation error of a single tree $\frac{\sigma^2 k}{a}$ instead of $\frac{\sigma^2 k}{n}$; no change for infinite forest

Outline

- 1 Random forests
- 2 Purely random forests
- 3 Toy forests
- 4 Hold-out random forests**

Definition (Biau, 2012)

Split D_n into D_{n_1} and D_{n_2}



\Rightarrow purely random forest
 closely related to double-sample trees (Wager & Athey, 2018)

Numerical experiments: framework

- Data generation:

$$\begin{aligned} X_i &\sim \mathcal{U}([0, 1]^p) & Y_i &= s^*(X_i) + \varepsilon_i \\ \varepsilon_i &\sim \mathcal{N}(0, \sigma^2) & \sigma^2 &= 1/16 \end{aligned}$$

$$s^* : \mathbf{x} \in [0, 1]^p \mapsto \frac{1}{10} \times \left[10 \sin(\pi x_1 x_2) + 20(x_3 - 0.5)^2 + 10x_4 + 5x_5 \right].$$

- Data split: $n_1 = 1\,280$ $n_2 = 25\,600$
- Forests definition:
 nodesize = 1
 $k \in \{2^5, 2^6, 2^7, 2^8, 2^9\}$
 “Large” forests are made of $q = k$ trees.
- Compute integrated approximation/estimation errors

Numerical experiments: results ($p = 5$)

	Single tree		Large forest	
No bootstrap $m_{\text{try}} = p$	0.13	$\frac{1.04\sigma^2 k}{n_2}$	0.13	$\frac{1.04\sigma^2 k}{n_2}$
Bootstrap $m_{\text{try}} = p$	0.14	$\frac{1.06\sigma^2 k}{n_2}$	0.15	$\frac{0.08\sigma^2 k}{n_2}$
No bootstrap $m_{\text{try}} = \lfloor p/3 \rfloor$	0.23	$\frac{1.01\sigma^2 k}{n_2}$	0.06	$\frac{0.06\sigma^2 k}{n_2}$
Bootstrap $m_{\text{try}} = \lfloor p/3 \rfloor$	0.25	$\frac{1.02\sigma^2 k}{n_2}$	0.06	$\frac{0.05\sigma^2 k}{n_2}$

$$\frac{2}{2+p} \approx 0.286$$

$$\frac{4}{4+p} \approx 0.444$$

Numerical experiments: results ($p = 10$)

	Single tree		Large forest	
No bootstrap $m_{\text{try}} = p$	$\frac{0.11}{k^{0.12}}$	$+\frac{1.03\sigma^2 k}{n_2}$	$\frac{0.11}{k^{0.12}}$	$+\frac{1.03\sigma^2 k}{n_2}$
Bootstrap $m_{\text{try}} = p$	$\frac{0.11}{k^{0.11}}$	$+\frac{1.05\sigma^2 k}{n_2}$	$\frac{0.10}{k^{0.19}}$	$+\frac{0.04\sigma^2 k}{n_2}$
No bootstrap $m_{\text{try}} = \lfloor p/3 \rfloor$	$\frac{0.21}{k^{0.18}}$	$+\frac{1.08\sigma^2 k}{n_2}$	$\frac{0.08}{k^{0.25}}$	$+\frac{0.04\sigma^2 k}{n_2}$
Bootstrap $m_{\text{try}} = \lfloor p/3 \rfloor$	$\frac{0.20}{k^{0.16}}$	$+\frac{1.05\sigma^2 k}{n_2}$	$\frac{0.07}{k^{0.26}}$	$+\frac{0.03\sigma^2 k}{n_2}$

$$\frac{2}{2+p} \approx 0.167$$

$$\frac{4}{4+p} \approx 0.286$$

Conclusion

- Forests improve the **order of magnitude** of the **approximation error**, compared to a single tree
- **Estimation error** seems to change only by a **constant factor** (at least for toy forests);
not contradictory with literature: here, we fix k ; different picture if `nodesize` is fixed (+subsampling)

Conclusion

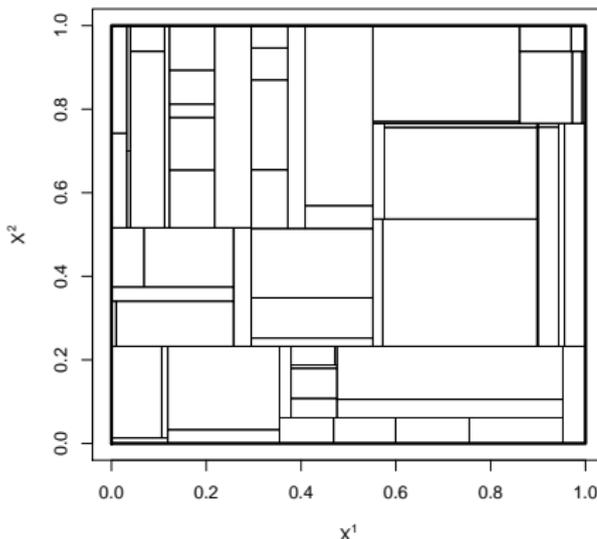
- Forests improve the **order of magnitude** of the **approximation error**, compared to a single tree
- **Estimation error** seems to change only by a **constant factor** (at least for toy forests);
not contradictory with literature: here, we fix k ; different picture if `nodesize` is fixed (+subsampling)
- Randomization:
randomization of labels seems to have no impact;
strong impact of **randomization of partitions** (hold-out RF: both bootstrap and `mtry`)

Approximation error: generalization: purf

Purely uniformly random forests:

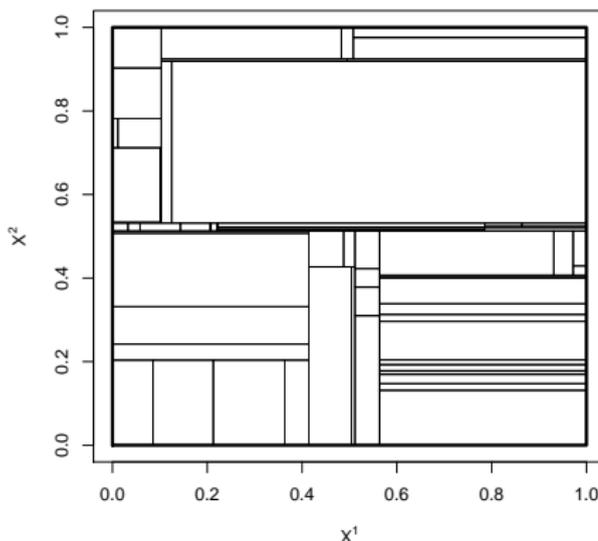
split a random cell, chosen with probability equal to its volume

⇒ in dimension $p = 1$, rates similar to toy forests

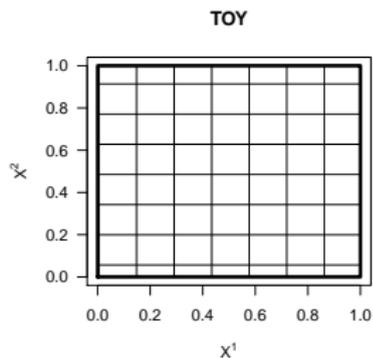
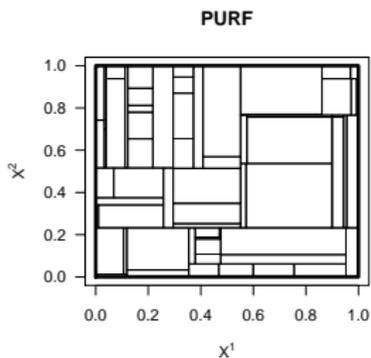
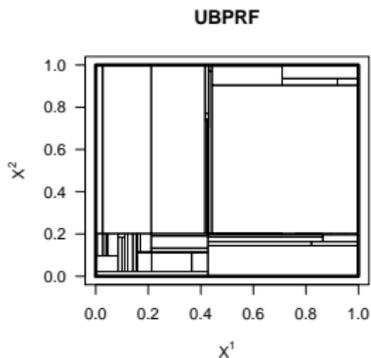


Approximation error: generalization: bprf

Balanced purely random forests in dimension p :
full binary tree, uniform splits $\Rightarrow k^{-\alpha}$ (tree) vs. $k^{-2\alpha}$ (forest)
where $\alpha = -\log_2\left(1 - \frac{1}{2p}\right) \Rightarrow$ not minimax rates!



Approximation error: 4 PRFs, different rates



Approximation error: generalization

- General result on the **approximation error** under $(H\theta)$:
e.g., roughly, if x is **centered in its cell** (on average over \mathbb{U}),
tree approx. error $\propto \mathcal{M}_2$ **infinite forest** approx. error $\propto \mathcal{M}_2^2$
where $\mathcal{M}_2 \approx$ average **square distance from x to the boundary**
of its cell ($\propto k^{-2/p}$ for toy forests)
- other PRF studied in the literature: Mondrian forests (Mourtada, Gaïffas & Scornet 2017 & 2020), centered random forests (Biau, 2012; Klusowski, 2018), ...

Open problems / future work

- Theory on **approximation error of hold-out RF?**
⇒ understand the typical shape of the cell that contains x , for a RI tree
(x centered on average? square distance to boundary?)
- Theory on **estimation error** of other PRF (beyond toy and PURF), with **lower bounds?** of hold-out RF?
- Extensive numerical experiments? (other functions s^* , ...)