# Analysis of some purely random forests

Sylvain Arlot[1] (joint work with Robin Genuer[2])

[1]UNIVERSITÉ PARIS-SUD

[2]ISPED, Université Bordeaux 2

TopData Workshop, Banyuls
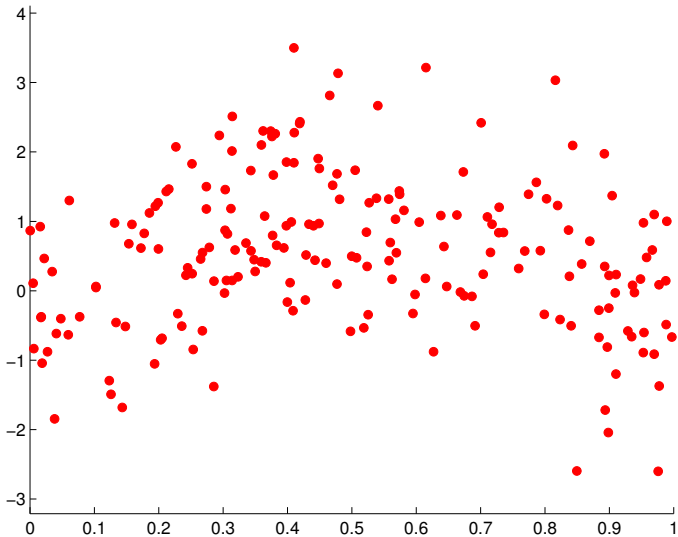13 June 2017

arXiv:1407.3939      arXiv:1604.01515

## Outline

## Outline

3/39

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

# Goal: find the signal (denoising)

## Regression

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

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

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

## Regression

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

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

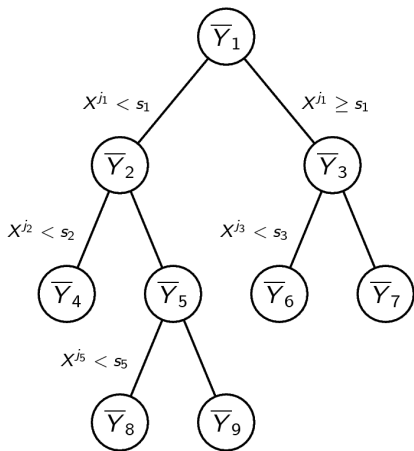with $s^\star(X) = \mathbb{E}[Y \mid X]$ (regression function).

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

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

is minimal.

## Regression tree (Breiman et al, 1984)



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

Restriction: splits parallel to the axes.

# Regression tree (Breiman et al, 1984)



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

1. **Choice of the partition $\mathbb{U}$** (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)
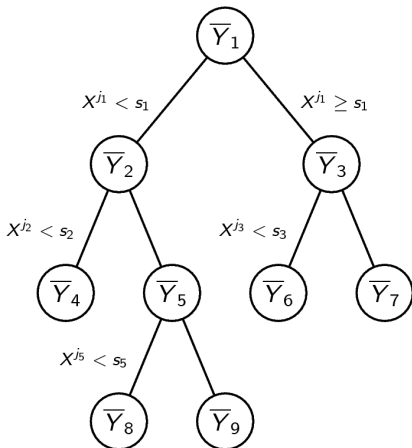


Tree: piecewise-constant predictor, obtained by partitioning recursively $\mathbb{R}^d$.

1. Choice of the partition $\mathbb{U}$ (tree structure)

2. For each $\lambda \in \mathbb{U}$ (tree leaf), choice of the estimation $\widehat{\beta}_\lambda$ of $s^\star(x)$ when $x \in \lambda$. Here, $\widehat{\beta}_\lambda = \overline{Y}_\lambda$ average of the $(Y_i)_{X_i \in \lambda}$.

# Random forest (Breiman, 2001)

## Definition (Random forest (Breiman, 2001))

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

Random forest predictor $\widehat{s}$ obtained by aggregating the tree collection.

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

- 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 $\{(X_i, Y_i) \,/\, i = 1, \ldots, n\}$ (sampling with replacement)
  $\Rightarrow$ resample $D_n^b$

# Bagging ("bootstrap aggregating")

- Bootstrap (Efron, 1979): draw $n$ i.i.d. r.v., uniform over $\{(X_i, Y_i) \,/\, i = 1, \ldots, n\}$ (sampling with replacement)
  $\Rightarrow$ resample $D_n^b$

- Bootstrapping a tree: $\widehat{s}_{\mathrm{tree}}^b = \widehat{s}_{\mathrm{tree}}(D_n^b)$

## Bagging ("bootstrap aggregating")

- Bootstrap (Efron, 1979): draw $n$ i.i.d. r.v., uniform over $\{(X_i, Y_i) / i = 1, \ldots, n\}$ (sampling with replacement)
  $\Rightarrow$ resample $D_n^b$

- Bootstrapping a tree: $\widehat{s}_{\mathrm{tree}}^b = \widehat{s}_{\mathrm{tree}}(D_n^b)$

- Bagging: bootstrap ($q$ independent resamples) then aggregation

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

9/39

# 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.

# 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.

# 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.

$$RF\text{-}RI \quad \Leftrightarrow \quad \text{bagging on RI trees}$$

## Random Forest-Random Inputs



The figure shows $D_n$ at the top. From $D_n$, via **Bootstrap**, arrows point to $D_n^{b,1}$, $D_n^{b,2}$, ..., ..., $D_n^{b,q}$.

From each, via **RI tree**, arrows point down to $\widehat{s}_{\Theta_1}$, $\widehat{s}_{\Theta_2}$, ..., ..., $\widehat{s}_{\Theta_q}$.

Via **Aggregation**, these combine into $\widehat{s}_{\mathrm{RF-RI}}$.

11/39

# Example of application of random forests: Kinect

Depth image  $\Rightarrow$ depth comparison features at each pixel



$\Rightarrow$ body part at each pixel  $\Rightarrow$ body part positions  $\Rightarrow \cdots$

Figures from Shotton et al (2011)    12/39

# Theoretical results on RF-RI

- Few theoretical results on Breiman's original RF-RI

- 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 $s^\star$

- References (see survey paper by Biau and Scornet, 2016): Mentch & Hooker (2014), Scornet, Biau & Vert (2015), Wager & Athey (2015).

## Theoretical results on RF-RI

- Few theoretical results on Breiman's original RF-RI

- 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 $s^\star$

- References (see survey paper by Biau and Scornet, 2016): Mentch & Hooker (2014), Scornet, Biau & Vert (2015), Wager & Athey (2015).

$\Rightarrow$ Here, we consider simplified RF models, for which a precise analysis is possible: purely random forests

13/39

## Outline

14/39

## Purely random forests

---

**Definition (Purely random tree)**

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

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)**

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

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

## Purely random forests

---

**Definition (Purely random forest)**

$$\widehat{s}(x) = \frac{1}{q} \sum_{j=1}^{q} \widehat{s}_{\mathbb{U}^j}(x) = \frac{1}{q} \sum_{j=1}^{q} \sum_{\lambda \in \mathbb{U}^j} \overline{Y_\lambda}(D_n) \mathbb{1}_{x \in \lambda}$$

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

---

Example ("hold-out RF" model): (random) split of the sample into $D_n$ (used for defining the labels $\overline{Y_\lambda}$) and $D'_n$ (used for building the trees $\mathbb{U}^j = \mathbb{U}_{\mathrm{RI}}(D_n'^{\star j})$).

## Purely random forests

> ### Definition (Purely random forest)
>
> $$\widehat{s}(x) = \frac{1}{q} \sum_{j=1}^{q} \widehat{s}_{\mathbb{U}^j}(x) = \frac{1}{q} \sum_{j=1}^{q} \sum_{\lambda \in \mathbb{U}^j} \overline{Y_\lambda}(D_n) \mathbb{1}_{x \in \lambda}$$
>
> with $\mathbb{U}^1, \ldots, \mathbb{U}^q$ i.i.d., independent from $D_n$.

Example ("hold-out RF" model): (random) split of the sample into $D_n$ (used for defining the labels $\overline{Y_\lambda}$) and $D_n'$ (used for building the trees $\mathbb{U}^j = \mathbb{U}_{\mathrm{RI}}(D_n'^{\star j})$).

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

## Purely random forests

$\mathbb{U}^1$ $\mathbb{U}^2$ ... ... $\mathbb{U}^q$     Independent from $D_n$

...    ...

Using $D_n$, with or without resampling

...    ...

$\widehat{s}_{\mathbb{U}^1}$ $\widehat{s}_{\mathbb{U}^2}$ ... ... $\widehat{s}_{\mathbb{U}^q}$

Aggregation

$\widehat{s}_{\mathrm{RF-RI}}$

16/39

# Purely random forests: theory

- Consistency: Biau, Devroye & Lugosi (2008), Scornet (2014)

- Rates of convergence: Breiman (2004), Biau (2012)
- Some adaptivity to dimension reduction (sparse framework): Biau (2012)

- 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)
- Some adaptivity to dimension reduction (sparse framework): Biau (2012)

- Forests decrease the estimation error (Biau, 2012; Genuer, 2012)

$\Rightarrow$ 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

$$\widehat{s}_{\mathbb{U}}(x) := \sum_{\lambda \in \mathbb{U}} \overline{Y_\lambda} \mathbb{1}_{x \in \lambda}$$

where $\overline{Y_\lambda}$ is the average of $(Y_i)_{X_i \in \lambda}$.

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

where $S_{\mathbb{U}}$ 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

$$\widehat{s}_{\mathbb{U}}(x) := \sum_{\lambda \in \mathbb{U}} \overline{Y_\lambda} \mathbb{1}_{x \in \lambda}$$

where $\overline{Y_\lambda}$ is the average of $(Y_i)_{X_i \in \lambda}$.

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

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

Define:

$$\widetilde{s}_{\mathbb{U}}(x) := \sum_{\lambda \in \mathbb{U}} \beta_\lambda \mathbb{1}_{x \in \lambda} \quad \text{where } \beta_\lambda := \mathbb{E}[s^\star(X) \,|\, X \in \lambda] \ .$$

$\Rightarrow \widetilde{s}_{\mathbb{U}} \in \operatorname{argmin}_{f \in S_{\mathbb{U}}} \mathbb{E}\big[(f(X) - s^\star(X))^2\big]$ and $\widetilde{s}_{\mathbb{U}}(x) = \mathbb{E}\big[\widehat{s}_{\mathbb{U}}(x) \,|\, \mathbb{U}\big]$

## Risk decomposition: single tree

$$\mathbb{E}\Big[\big(\widehat{s}_{\mathbb{U}}(X) - s^{\star}(X)\big)^2\Big]$$
$$= \mathbb{E}\Big[\big(\tilde{s}_{\mathbb{U}}(X) - s^{\star}(X)\big)^2\Big] + \mathbb{E}\Big[\big(\widehat{s}_{\mathbb{U}}(X) - \tilde{s}_{\mathbb{U}}(X)\big)^2\Big]$$
$$= \text{Approximation error} + \text{Estimation error}$$

## Risk decomposition: single tree

$$\mathbb{E}\Big[\big(\widehat{s}_{\mathbb{U}}(X) - s^{\star}(X)\big)^2\Big]$$
$$= \mathbb{E}\Big[\big(\tilde{s}_{\mathbb{U}}(X) - s^{\star}(X)\big)^2\Big] + \mathbb{E}\Big[\big(\widehat{s}_{\mathbb{U}}(X) - \tilde{s}_{\mathbb{U}}(X)\big)^2\Big]$$
$$= \text{ Approximation error } + \text{ Estimation error}$$

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

$$\mathbb{E}\Big[\big(\tilde{s}_{\mathbb{U}}(X) - s^{\star}(X)\big)^2\Big] \propto \frac{1}{K^2}$$

## Risk decomposition: single tree

$$\mathbb{E}\Big[\big(\widehat{s}_{\mathbb{U}}(X) - s^{\star}(X)\big)^2\Big]$$
$$= \mathbb{E}\Big[\big(\widetilde{s}_{\mathbb{U}}(X) - s^{\star}(X)\big)^2\Big] + \mathbb{E}\Big[\big(\widehat{s}_{\mathbb{U}}(X) - \widetilde{s}_{\mathbb{U}}(X)\big)^2\Big]$$
$$= \text{ Approximation error } + \text{ Estimation error}$$

If $s^{\star}$ is smooth, $X \sim \mathcal{U}([0,1])$ and $\mathbb{U}$ regular partition into $K$ pieces, then
$$\mathbb{E}\Big[\big(\widetilde{s}_{\mathbb{U}}(X) - s^{\star}(X)\big)^2\Big] \propto \frac{1}{K^2}$$

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

$$\mathbb{E}\Big[\big(\widetilde{s}_{\mathbb{U}}(X) - \widehat{s}_{\mathbb{U}}(X)\big)^2\Big] \approx \frac{\sigma^2 K}{n}$$

19/39

# Approximation and estimation errors

## Risk decomposition: purely random forest

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

Estimator (forest): $\quad \widehat{s}_{\mathbb{U}^{1\cdots q}}(x) := \dfrac{1}{q} \sum_{j=1}^{q} \widehat{s}_{\mathbb{U}^j}(x)$

Ideal forest: $\quad \tilde{s}_{\mathbb{U}^{1\cdots q}}(x) := \dfrac{1}{q} \sum_{j=1}^{q} \tilde{s}_{\mathbb{U}^j}(x) = \mathbb{E}\big[\widehat{s}_{\mathbb{U}^{1\cdots q}}(x)\,|\,\mathbb{U}^{1\cdots q}\big]$

# Risk decomposition: purely random forest

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

Estimator (forest):     $\widehat{s}_{\mathbb{U}^{1\cdots q}}(x) := \dfrac{1}{q} \sum_{j=1}^{q} \widehat{s}_{\mathbb{U}^j}(x)$

Ideal forest:     $\tilde{s}_{\mathbb{U}^{1\cdots q}}(x) := \dfrac{1}{q} \sum_{j=1}^{q} \tilde{s}_{\mathbb{U}^j}(x) = \mathbb{E}\big[\widehat{s}_{\mathbb{U}^{1\cdots q}}(x) \,|\, \mathbb{U}^{1\cdots q}\big]$

**Quadratic risk decomposition (given $X = x$)**

$$\mathbb{E}\Big[\big(\widehat{s}_{\mathbb{U}^{1\cdots q}}(x) - s^\star(x)\big)^2\Big] = \mathbb{E}\Big[\big(\tilde{s}_{\mathbb{U}^{1\cdots q}}(x) - s^\star(x)\big)^2\Big]$$
$$+ \mathbb{E}\Big[\big(\widehat{s}_{\mathbb{U}^{1\cdots q}}(x) - \tilde{s}_{\mathbb{U}^{1\cdots q}}(x)\big)^2\Big]$$

# Risk decomposition: purely random forest

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

Estimator (forest):     $\widehat{s}_{\mathbb{U}^{1 \cdots q}}(x) := \dfrac{1}{q} \displaystyle\sum_{j=1}^{q} \widehat{s}_{\mathbb{U}^j}(x)$

Ideal forest:     $\widetilde{s}_{\mathbb{U}^{1 \cdots q}}(x) := \dfrac{1}{q} \displaystyle\sum_{j=1}^{q} \widetilde{s}_{\mathbb{U}^j}(x) = \mathbb{E}\big[\widehat{s}_{\mathbb{U}^{1 \cdots q}}(x) \,|\, \mathbb{U}^{1 \cdots q}\big]$

**Quadratic risk decomposition (given $X = x$)**

$$\mathbb{E}\Big[\big(\widehat{s}_{\mathbb{U}^{1 \cdots q}}(x) - s^{\star}(x)\big)^2\Big] = \mathbb{E}\Big[\big(\widetilde{s}_{\mathbb{U}^{1 \cdots q}}(x) - s^{\star}(x)\big)^2\Big]$$
$$+ \mathbb{E}\Big[\big(\widehat{s}_{\mathbb{U}^{1 \cdots q}}(x) - \widetilde{s}_{\mathbb{U}^{1 \cdots q}}(x)\big)^2\Big]$$

Bias term (approximation error):
$$\mathcal{B}_{\mathcal{U},q}(x) := \mathbb{E}\Big[\big(\widetilde{s}_{\mathbb{U}^{1 \cdots q}}(x) - s^{\star}(x)\big)^2\Big]$$

# Bias decomposition (given $X = x$)

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

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

$$\text{and} \quad \mathcal{V}_{\mathcal{U}}(x) := \text{var}(\tilde{s}_{\mathbb{U}}(x))$$

$\mathcal{B}_{\mathcal{U},\infty}(x)$ is the bias of the infinite forest: $\tilde{s}_{\mathbb{U},\infty}(x) := \mathbb{E}\big[\tilde{s}_{\mathbb{U}}(x)\big]$

to be compared with the bias of a single tree

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

22/39

## Outline

1. Random forests

2. Purely random forests

3. Toy forests in one dimension

4. Hold-out random forests

23/39

## Toy forests in one dimension

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

$\mathbb{U} \sim \mathcal{U}_k^{\text{toy}}$ defined by:

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

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

## Interpretation of the ideal infinite forest

### Proposition (A. & Genuer, 2014)

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

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

where

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

## Interpretation of the ideal infinite forest: proof

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

$$\tilde{s}_{\mathbb{U}}(x) = \frac{1}{|I_{\mathbb{U}}(x)|} \int_{I_{\mathbb{U}}(x)} s^{\star}(t)\, \mathrm{d}t$$

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: proof

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

$$\tilde{s}_{\mathbb{U}}(x) = \frac{1}{|I_{\mathbb{U}}(x)|} \int_{I_{\mathbb{U}}(x)} s^{\star}(t) \, \mathrm{d}t$$

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

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

$$\begin{aligned}
\tilde{s}_{\mathbb{U}, \infty}(x) &= \mathbb{E}_{\mathbb{U}}\left[\tilde{s}_{\mathbb{U}}(x)\right] \\
&= k \int_0^1 s^{\star}(t) \, \mathbb{P}\left(x + \frac{V_x - 1}{k} \leqslant t < x + \frac{V_x}{k}\right) \mathrm{d}t \\
&= k \int_0^1 s^{\star}(t) \, \mathbb{P}\left(k(t - x) < V_x \leqslant k(t - x) + 1\right) \mathrm{d}t
\end{aligned}$$

## Analysis of the approximation error

(H2)    $s^\star$ twice differentiable over $(0,1)$ and $s^{\star\prime\prime}$ bounded

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

$$s^\star(t) - s^\star(x) = s^{\star\prime}(x)(t-x) + \frac{1}{2}s^{\star\prime\prime}(c_{t,x})(t-x)^2$$

## Analysis of the approximation error

(H2)    $s^\star$ twice differentiable over $(0,1)$ and $s^{\star\prime\prime}$ bounded

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

$$s^\star(t) - s^\star(x) = s^{\star\prime}(x)(t-x) + \frac{1}{2}s^{\star\prime\prime}(c_{t,x})(t-x)^2$$

Therefore,

$$\tilde{s}_{\mathbb{U}}(x) - s^\star(x) = k \int_{x+\frac{V_x-1}{k}}^{x+\frac{V_x}{k}} (s^\star(t) - s^\star(x))\,\mathrm{d}t$$

$$= k\, s^{\star\prime}(x) \int_{x+\frac{V_x-1}{k}}^{x+\frac{V_x}{k}} (t-x)\,\mathrm{d}t + R_1(x)$$

$$= \frac{s^{\star\prime}(x)}{k}\left(V_x - \frac{1}{2}\right) + R_1(x)$$

where $R_1(x) = \frac{k}{2}\int_{x+\frac{V_x-1}{k}}^{x+\frac{V_x}{k}} s^{\star\prime\prime}(c_{t,x})(t-x)^2\,\mathrm{d}t$

27/39

# Analysis of the approximation error

$$\left(\mathbb{E}_{\mathbb{U}}\big[\tilde{s}_{\mathbb{U}}(x) - s^\star(x)\big]\right)^2 \leqslant \frac{\Box}{k^4} \qquad \mathcal{V}_{\mathcal{U}}(x) \underset{k \to +\infty}{\sim} \frac{\Box}{k^2}$$

---

**Proposition (A. & Genuer, 2014)**

*Assuming (H2), for every $x \in \left[\frac{1}{k}, 1 - \frac{1}{k}\right]$,*

$$\mathcal{B}_{\mathcal{U}_k^{\mathrm{toy}},1}(x) \underset{k \to +\infty}{\sim} \frac{\Box}{k^2} \qquad \mathcal{B}_{\mathcal{U}_k^{\mathrm{toy}},\infty}(x) \leqslant \frac{\Box}{k^4}$$

$$\int_{\frac{1}{k}}^{1-\frac{1}{k}} \mathcal{B}_{\mathcal{U}_k^{\mathrm{toy}},1}(x)\,\mathrm{d}x \underset{k \to +\infty}{\sim} \frac{\Box}{k^2} \qquad \int_{\frac{1}{k}}^{1-\frac{1}{k}} \mathcal{B}_{\mathcal{U}_k^{\mathrm{toy}},\infty}(x)\,\mathrm{d}x \leqslant \frac{\Box}{k^4}$$

# Analysis of the approximation error

$$\left(\mathbb{E}_{\mathbb{U}}\left[\tilde{s}_{\mathbb{U}}(x) - s^\star(x)\right]\right)^2 \leqslant \frac{\square}{k^4} \qquad \mathcal{V}_{\mathcal{U}}(x) \underset{k \to +\infty}{\sim} \frac{\square}{k^2}$$

---

**Proposition (A. & Genuer, 2014)**

*Assuming (H2), for every $x \in \left[\frac{1}{k}, 1 - \frac{1}{k}\right]$,*

$$\mathcal{B}_{\mathcal{U}_k^{\mathrm{toy}},1}(x) \underset{k \to +\infty}{\sim} \frac{\square}{k^2} \qquad \mathcal{B}_{\mathcal{U}_k^{\mathrm{toy}},\infty}(x) \leqslant \frac{\square}{k^4}$$

$$\int_{\frac{1}{k}}^{1-\frac{1}{k}} \mathcal{B}_{\mathcal{U}_k^{\mathrm{toy}},1}(x)\,\mathrm{d}x \underset{k \to +\infty}{\sim} \frac{\square}{k^2} \qquad \int_{\frac{1}{k}}^{1-\frac{1}{k}} \mathcal{B}_{\mathcal{U}_k^{\mathrm{toy}},\infty}(x)\,\mathrm{d}x \leqslant \frac{\square}{k^4}$$

---

Rate $k^{-4}$ is tight assuming:

(H3)   $s^\star$ three times differentiable over $(0,1)$ and $s^{\star\prime\prime\prime}$ bounded

## Estimation error

General fact (Jensen's inequality):

$$\mathbb{E}\Big[\big(\widehat{s}_{\mathbb{U},\,\infty}(X) - \tilde{s}_{\mathbb{U},\,\infty}(X)\big)^2\Big] \leqslant \mathbb{E}\Big[\big(\widehat{s}_{\mathbb{U}}(X) - \tilde{s}_{\mathbb{U}}(X)\big)^2\Big]$$

## Estimation error

General fact (Jensen's inequality):

$$\mathbb{E}\Big[\big(\widehat{s}_{\mathbb{U},\infty}(X) - \tilde{s}_{\mathbb{U},\infty}(X)\big)^2\Big] \leqslant \mathbb{E}\Big[\big(\widehat{s}_{\mathbb{U}}(X) - \tilde{s}_{\mathbb{U}}(X)\big)^2\Big]$$

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

$$\mathbb{E}\Big[\big(\widehat{s}_{\mathbb{U}}(X) - \tilde{s}_{\mathbb{U}}(X)\big)^2\Big] \approx \frac{\sigma^2 k}{n}$$

$$\mathbb{E}\Big[\big(\widehat{s}_{\mathbb{U},\infty}(X) - \tilde{s}_{\mathbb{U},\infty}(X)\big)^2\Big] \approx \frac{2}{3}\frac{\sigma^2 k}{n}$$

(A. & Genuer, 2016)

## Summary: risk analysis

Single tree $(q=1)$     Infinite forest $(q=\infty)$

$$\mathbb{E}\left[\left(\widehat{s}_{\mathbb{U}^{1\cdots q}}(x) - s^{\star}(x)\right)^2\right] \approx \quad \frac{c_1(s^{\star}, x)}{k^2} + \frac{\sigma^2 k}{n} \quad \frac{c_2(s^{\star}, x)}{k^4} + \frac{2\sigma^2 k}{3n}$$

where $\quad c_1(s^{\star}, x) = \dfrac{s^{\star\prime}(x)^2}{12} \quad$ and $\quad c_2(s^{\star}, x) = \dfrac{s^{\star\prime\prime}(x)^2}{144} \ .$

Assumptions:

- $x \in (0, 1)$ far from boundary
- (H3) $s^{\star}$ three times differentiable over $(0, 1)$ and $s^{\star\prime\prime\prime}$ bounded
- $\mathcal{X}$ uniform over $[0, 1]$
- $\mathrm{var}(Y|X) = \sigma^2$
- no resampling for computing labels

## Rates of convergence

Corollary: risk convergence rates (far from boundaries, with $k = k_n^\star$ optimal):

$$\text{Tree} \geqslant \square\, n^{-2/3}$$

$$\text{Infinite forest} \leqslant \square\, n^{-4/5} \quad \Rightarrow \quad \text{minimax } \mathcal{C}^2$$

## Rates of convergence

Corollary: risk convergence rates (far from boundaries, with $k = k_n^\star$ optimal):

$$\text{Tree} \geqslant \Box \, n^{-2/3}$$
$$\text{Infinite forest} \leqslant \Box \, n^{-4/5} \quad \Rightarrow \quad \text{minimax } \mathcal{C}^2$$

Remarks:

- $q \geqslant \Box \, (k_n^\star)^2$ is sufficient to get an "infinite" forest

## Rates of convergence

Corollary: risk convergence rates (far from boundaries, with $k = k_n^\star$ optimal):

$$\text{Tree} \geqslant \Box \, n^{-2/3}$$

$$\text{Infinite forest} \leqslant \Box \, n^{-4/5} \quad \Rightarrow \quad \text{minimax } \mathcal{C}^2$$

Remarks:

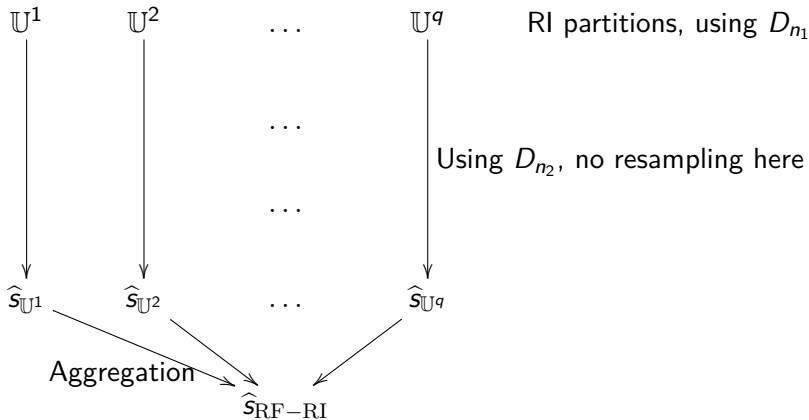- $q \geqslant \Box \, (k_n^\star)^2$ is sufficient to get an "infinite" forest

- 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

31/39

## Outline

## Definition (Biau, 2012)

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

$\mathbb{U}^1 \qquad \mathbb{U}^2 \qquad \ldots \qquad \mathbb{U}^q$  RI partitions, using $D_{n_1}$

. . .

Using $D_{n_2}$, no resampling here

. . .

$\widehat{s}_{\mathbb{U}^1} \qquad \widehat{s}_{\mathbb{U}^2} \qquad \ldots \qquad \widehat{s}_{\mathbb{U}^q}$

Aggregation

$\widehat{s}_{\mathrm{RF-RI}}$

$\Rightarrow$ purely random forest

## Numerical experiments: framework

- Data generation:
  $X_i \sim \mathcal{U}([0,1]^d)$      $Y_i = s^\star(X_i) + \varepsilon_i$
  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$      $\sigma^2 = 1/16$

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

- Data split: $n_1 = 1\,280$     $n_2 = 25\,600$

- Forests definition:
  $\texttt{nodesize} = 1$
  $k \in \{2^5, 2^6, 2^7, 2^8\}$
  "Large" forests are made of $q = k$ trees.

- Compute integrated approximation/estimation errors

## Numerical experiments: results ($d = 5$)

|  | Single tree | | Large forest | |
|---|---|---|---|---|
| No bootstrap $\mathtt{mtry} = d$ | $\dfrac{0.13}{k^{0.17}}$ | $+ \dfrac{1.04\sigma^2 k}{n_2}$ | $\dfrac{0.13}{k^{0.17}}$ | $+ \dfrac{1.04\sigma^2 k}{n_2}$ |
| Bootstrap $\mathtt{mtry} = d$ | $\dfrac{0.14}{k^{0.17}}$ | $+ \dfrac{1.06\sigma^2 k}{n_2}$ | $\dfrac{0.15}{k^{0.29}}$ | $+ \dfrac{0.08\sigma^2 k}{n_2}$ |
| No bootstrap $\mathtt{mtry} = \lfloor d/3 \rfloor$ | $\dfrac{0.23}{k^{0.19}}$ | $+ \dfrac{1.01\sigma^2 k}{n_2}$ | $\dfrac{0.06}{k^{0.31}}$ | $+ \dfrac{0.06\sigma^2 k}{n_2}$ |
| Bootstrap $\mathtt{mtry} = \lfloor d/3 \rfloor$ | $\dfrac{0.25}{k^{0.20}}$ | $+ \dfrac{1.02\sigma^2 k}{n_2}$ | $\dfrac{0.06}{k^{0.34}}$ | $+ \dfrac{0.05\sigma^2 k}{n_2}$ |

35/39

# Numerical experiments: results ($d = 10$)

| | Single tree | | Large forest | |
|---|---|---|---|---|
| No bootstrap $\mathtt{mtry} = d$ | $\dfrac{0.11}{k^{0.12}}$ | $+ \dfrac{1.03\sigma^2 k}{n_2}$ | $\dfrac{0.11}{k^{0.12}}$ | $+ \dfrac{1.03\sigma^2 k}{n_2}$ |
| Bootstrap $\mathtt{mtry} = d$ | $\dfrac{0.11}{k^{0.11}}$ | $+ \dfrac{1.05\sigma^2 k}{n_2}$ | $\dfrac{0.10}{k^{0.19}}$ | $+ \dfrac{0.04\sigma^2 k}{n_2}$ |
| No bootstrap $\mathtt{mtry} = \lfloor d/3 \rfloor$ | $\dfrac{0.21}{k^{0.18}}$ | $+ \dfrac{1.08\sigma^2 k}{n_2}$ | $\dfrac{0.08}{k^{0.25}}$ | $+ \dfrac{0.04\sigma^2 k}{n_2}$ |
| Bootstrap $\mathtt{mtry} = \lfloor d/3 \rfloor$ | $\dfrac{0.20}{k^{0.16}}$ | $+ \dfrac{1.05\sigma^2 k}{n_2}$ | $\dfrac{0.07}{k^{0.26}}$ | $+ \dfrac{0.03\sigma^2 k}{n_2}$ |

Random forests
○○○○○○○○○○○

Purely random forests
○○○○○○○○

Toy forests
○○○○○○○○

Hold-out random forests
○○○○

**Conclusion**

## 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

- General result on the approximation error under (H2)/(H3): 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 ($\propto k^{-2}$ for toy forests)

Random forests
○○○○○○○○○○

Purely random forests
○○○○○○○○

Toy forests
○○○○○○○○

Hold-out random forests
○○○○

**Conclusion**
○○○○

## Approximation error: generalization

- General result on the approximation error under (H2)/(H3): 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 ($\propto k^{-2}$ for toy forests)

- toy forests in dimension $d$: approximation error $\propto k^{-2/d}$ vs. $k^{-4/d}$ (infinite forest reaches minimax $\mathcal{C}^2$ rates)

## Approximation error: generalization

- General result on the approximation error under (H2)/(H3):
  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 ($\propto k^{-2}$ for toy forests)

- toy forests in dimension $d$: approximation error $\propto k^{-2/d}$ vs. $k^{-4/d}$ (infinite forest reaches minimax $\mathcal{C}^2$ rates)

- purely uniformly random forests in dimension 1 (split a random cell, chosen with probability equal to its volume): rates similar to toy forests

38/39

# Approximation error: generalization

- General result on the approximation error under (H2)/(H3): 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 ($\propto k^{-2}$ for toy forests)

- toy forests in dimension $d$: approximation error $\propto k^{-2/d}$ vs. $k^{-4/d}$ (infinite forest reaches minimax $\mathcal{C}^2$ rates)

- purely uniformly random forests in dimension 1 (split a random cell, chosen with probability equal to its volume): rates similar to toy forests

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

# Open problems / future work

- Extensive numerical experiments? (other functions $s^\star$, ...)

- Theory on approximation error of hold-out RF?
  $\Rightarrow$ understand the typical shape of a cell of a RI tree
  ($x$ centered on average? square distance to boundary?)

- Theory on estimation error of other models (beyond toy)?
  of hold-out RF?