

# Variable importance for random forests: a sensitivity analysis perspective

## ETICS 2021

**Clément Bénard**

**Thesis Advisors:** Gérard Biau, Sébastien da Veiga, Erwan Scornet

Safran Tech, LPSM

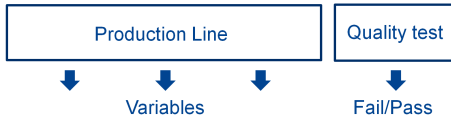
September 2021



# Industrial processes

- Context

Manufacturing process driven by controllable variables.



# Industrial processes

- Context

Manufacturing process driven by controllable variables.



- Objective

Identify production conditions generating defects: variable settings.

- Method

- 1 Fit a learning algorithm
- 2 Use variable importance to detect influential variables
- 3 Explore associated physical phenomenon with domain experts

- Random forests are an efficient approach
- MDA (Breiman, 2001): built-in variable importance algorithm for random forests

- Random forests are an efficient approach
- MDA (Breiman, 2001): built-in variable importance algorithm for random forests
  
- MDA is used intensively

- Random forests are an efficient approach
- MDA (Breiman, 2001): built-in variable importance algorithm for random forests
- MDA is used intensively
- MDA has flaws
  - Poor understanding of the MDA: what is estimated ?
  - Empirical studies show that the MDA is biased for dependent inputs (Strobl et al., 2007; Gregorutti et al., 2017; Hooker and Mentch, 2019)

- Random forests are an efficient approach
- MDA (Breiman, 2001): built-in variable importance algorithm for random forests
- MDA is used intensively
- MDA has flaws
  - Poor understanding of the MDA: what is estimated ?
  - Empirical studies show that the MDA is biased for dependent inputs (Strobl et al., 2007; Gregorutti et al., 2017; Hooker and Mentch, 2019)
- Our objective (Bénard et al., 2021)
  - Theoretical analysis of the MDA
    - First convergence result for the original MDA (Ishwaran, 2007; Zhu et al., 2015)
    - Theoretical understanding of MDA bias
  - Design of Sobol-MDA algorithm to fix the MDA flaws

- Regression setting
  - input vector  $\mathbf{X} = (X^{(1)}, \dots, X^{(p)}) \in \mathbb{R}^p$
  - output  $Y \in \mathbb{R}$
  - dataset  $\mathcal{D}_n = \{(\mathbf{X}_i, Y_i), i = 1, \dots, n\}$ ,  
where  $(\mathbf{X}_i, Y_i) \sim \mathbb{P}_{\mathbf{X}, Y}$ .



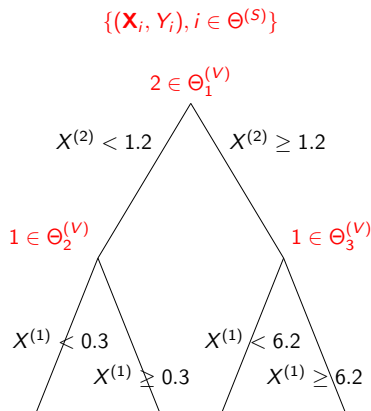
# Random forests

- Regression setting

- input vector  $\mathbf{X} = (X^{(1)}, \dots, X^{(p)}) \in \mathbb{R}^p$
- output  $Y \in \mathbb{R}$
- dataset  $\mathcal{D}_n = \{(\mathbf{X}_i, Y_i), i = 1, \dots, n\}$ ,  
where  $(\mathbf{X}_i, Y_i) \sim \mathbb{P}_{\mathbf{X}, Y}$ .

- Random forest algorithm

- Aggregation of  $\Theta$ -random trees  
 $\Theta = (\Theta^{(S)}, \Theta^{(V)})$
- $M$ : number of trees
- $m_{M,n}(\mathbf{X}, \Theta_M)$ : the forest estimate at  $\mathbf{X}$



- 1 Introduction
- 2 MDA Theoretical Limitations
  - MDA definition
  - MDA convergence
- 3 Sobol-MDA
- 4 Shapley effects

MDA principle:

decrease of accuracy of the forest when a variable is noised up

MDA principle:

decrease of accuracy of the forest when a variable is noised up

- 1 fit a random forest with  $\mathcal{D}_n$

MDA principle:

decrease of accuracy of the forest when a variable is noised up

- 1 fit a random forest with  $\mathcal{D}_n$
- 2 compute the accuracy of the forest

MDA principle:

decrease of accuracy of the forest when a variable is noised up

- 1 fit a random forest with  $\mathcal{D}_n$
- 2 compute the accuracy of the forest
- 3 permute randomly the values of a given input variable  $X^{(j)}$ :  
break the dependence between  $X^{(j)}$  and  $Y$

MDA principle:

decrease of accuracy of the forest when a variable is noised up

- 1 fit a random forest with  $\mathcal{D}_n$
- 2 compute the accuracy of the forest
- 3 permute randomly the values of a given input variable  $X^{(j)}$ :  
break the dependence between  $X^{(j)}$  and  $Y$
- 4 compute the decrease of accuracy of the forest with the permuted data

$X^{(1)}$	$X^{(2)}$	...	$X^{(j)}$	...	$X^{(p)}$	$Y$
2.1	4.3	...	0.1	...	2.6	2.3
1.7	4.1	...	9.2	...	3.8	0.4
3.4	9.2	...	3.2	...	3.6	10.2
5.6	1.2	...	8.2	...	4.2	9.1
8.9	6.8	...	6.7	...	2.9	4.5

Table: Example of the permutation of a dataset  $\mathcal{D}_n$  for  $n = 5$ .



$X^{(1)}$	$X^{(2)}$	...	$X^{(j)}$	...	$X^{(p)}$	$Y$
2.1	4.3	...	0.1	...	2.6	2.3
1.7	4.1	...	9.2	...	3.8	0.4
3.4	9.2	...	3.2	...	3.6	10.2
5.6	1.2	...	8.2	...	4.2	9.1
8.9	6.8	...	6.7	...	2.9	4.5

Table: Example of the permutation of a dataset  $\mathcal{D}_n$  for  $n = 5$ .

$X^{(1)}$	$X^{(2)}$	...	$X^{(j)}$	...	$X^{(p)}$	$Y$
2.1	4.3	...	0.1	...	2.6	2.3
1.7	4.1	...	9.2	...	3.8	0.4
3.4	9.2	...	3.2	...	3.6	10.2
5.6	1.2	...	8.2	...	4.2	9.1
8.9	6.8	...	6.7	...	2.9	4.5

$X^{(1)}$	$X^{(2)}$	...	$X^{(j)}$	...	$X^{(p)}$	$Y$
2.1	4.3	...	6.7	...	2.6	2.3
1.7	4.1	...	3.2	...	3.8	0.4
3.4	9.2	...	9.2	...	3.6	10.2
5.6	1.2	...	0.1	...	4.2	9.1
8.9	6.8	...	8.2	...	2.9	4.5

**Table:** Example of the permutation of a dataset  $\mathcal{D}_n$  for  $n = 5$ .

$X^{(1)}$	$X^{(2)}$	...	$X^{(j)}$	...	$X^{(p)}$	Y
2.1	4.3	...	0.1	...	2.6	2.3
1.7	4.1	...	9.2	...	3.8	0.4
3.4	9.2	...	3.2	...	3.6	10.2
5.6	1.2	...	8.2	...	4.2	9.1
8.9	6.8	...	6.7	...	2.9	4.5

$X^{(1)}$	$X^{(2)}$	...	$X^{(j)}$	...	$X^{(p)}$	Y
2.1	4.3	...	6.7	...	2.6	2.3
1.7	4.1	...	3.2	...	3.8	0.4
3.4	9.2	...	9.2	...	3.6	10.2
5.6	1.2	...	0.1	...	4.2	9.1
8.9	6.8	...	8.2	...	2.9	4.5

**Table:** Example of the permutation of a dataset  $\mathcal{D}_n$  for  $n = 5$ .

Explained variance of Y = 16.4

Explained variance of Y = 13.7

$$\text{MDA}(X^{(j)}) = 16.4 - 13.7 = 2.7$$

$X^{(1)}$	$X^{(2)}$	...	$X^{(j)}$	...	$X^{(p)}$	$Y$
2.1	4.3	...	0.1	...	2.6	2.3
1.7	4.1	...	9.2	...	3.8	0.4
3.4	9.2	...	3.2	...	3.6	10.2
5.6	1.2	...	8.2	...	4.2	9.1
8.9	6.8	...	6.7	...	2.9	4.5

$X^{(1)}$	$X^{(2)}$	...	$X^{(j)}$	...	$X^{(p)}$	$Y$
2.1	4.3	...	6.7	...	2.6	2.3
1.7	4.1	...	3.2	...	3.8	0.4
3.4	9.2	...	9.2	...	3.6	10.2
5.6	1.2	...	0.1	...	4.2	9.1
8.9	6.8	...	8.2	...	2.9	4.5

Table: Example of the permutation of a dataset  $\mathcal{D}_n$  for  $n = 5$ .

Question: Can I use  $\mathcal{D}_n$  to both fit the forest and compute accuracy ?

No: overfitting and inflated accuracy.

How to handle this in practice?

The explained variance estimate of MDA algorithms differ across implementations

**Train-Test MDA:** train data to fit the forest, and test data for accuracy

The explained variance estimate of MDA algorithms differ across implementations

**Train-Test MDA:** train data to fit the forest, and test data for accuracy

**Out-of-bag (OOB) samples:**  $\mathcal{D}_n$  is bootstrap prior to the construction of each tree, leaving aside a portion of  $\mathcal{D}_n$ , which is not involved in the tree growing and defines the “out-of-bag” sample.

$X^{(1)}$	$X^{(2)}$	...	$X^{(j)}$	...	$X^{(p)}$	$Y$
2.1	4.3	...	0.1	...	2.6	2.3
1.7	4.1	...	9.2	...	3.8	0.4
3.4	9.2	...	3.2	...	3.6	10.2
5.6	1.2	...	8.2	...	4.2	9.1
8.9	6.8	...	6.7	...	2.9	4.5

Selected samples:  $\Theta_\ell^{(S)} = \{1, 3, 4\}$

The explained variance estimate of MDA algorithms differ across implementations

**Train-Test MDA:** train data to fit the forest, and test data for accuracy

**Out-of-bag (OOB) samples:**  $\mathcal{D}_n$  is bootstrap prior to the construction of each tree, leaving aside a portion of  $\mathcal{D}_n$ , which is not involved in the tree growing and defines the “out-of-bag” sample.

$X^{(1)}$	$X^{(2)}$	...	$X^{(j)}$	...	$X^{(p)}$	$Y$
2.1	4.3	...	0.1	...	2.6	2.3
1.7	4.1	...	9.2	...	3.8	0.4
3.4	9.2	...	3.2	...	3.6	10.2
5.6	1.2	...	8.2	...	4.2	9.1
8.9	6.8	...	6.7	...	2.9	4.5

OOB samples:  $\{1, \dots, n\} \setminus \Theta_\ell^{(s)} = \{2, 5\}$

The explained variance estimate of MDA algorithms differ across implementations

**Train-Test MDA:** train data to fit the forest, and test data for accuracy

**Out-of-bag (OOB) samples:**  $\mathcal{D}_n$  is bootstrap prior to the construction of each tree, leaving aside a portion of  $\mathcal{D}_n$ , which is not involved in the tree growing and defines the “out-of-bag” sample.

MDA Version	Package	Error	Data
Train-Test	scikit-learn randomForestSRC	Forest	Testing dataset
Breiman-Cutler	randomForest (normalized) ranger / randomForestSRC	Tree	OOB sample
Ishwaran-Kogalur	randomForestSRC	Forest	OOB sample

Table: Summary of the different MDA algorithms.



- $i \in \{1, \dots, n\} \setminus \Theta_\ell^{(S)} = \{2, 5\}$ : OOB sample of the  $\ell$ -th tree

$X^{(1)}$	$X^{(2)}$	...	$X^{(j)}$	...	$X^{(p)}$	$Y$
2.1	4.3	...	0.1	...	2.6	2.3
1.7	4.1	...	9.2	...	3.8	0.4
3.4	9.2	...	3.2	...	3.6	10.2
5.6	1.2	...	8.2	...	4.2	9.1
8.9	6.8	...	6.7	...	2.9	4.5

# Breiman-Cutler MDA

- $i \in \{1, \dots, n\} \setminus \Theta_\ell^{(S)} = \{2, 5\}$ : OOB sample of the  $\ell$ -th tree
- $N_{n,\ell} = \sum_{i=1}^n \mathbb{1}_{i \notin \Theta_\ell^{(S)}} = 2$ : size of the OOB sample of the  $\ell$ -th tree

$X^{(1)}$	$X^{(2)}$	...	$X^{(j)}$	...	$X^{(p)}$	$Y$
2.1	4.3	...	0.1	...	2.6	2.3
1.7	4.1	...	9.2	...	3.8	0.4
3.4	9.2	...	3.2	...	3.6	10.2
5.6	1.2	...	8.2	...	4.2	9.1
8.9	6.8	...	6.7	...	2.9	4.5

# Breiman-Cutler MDA

- $i \in \{1, \dots, n\} \setminus \Theta_\ell^{(S)} = \{2, 5\}$ : OOB sample of the  $\ell$ -th tree
- $N_{n,\ell} = \sum_{i=1}^n \mathbb{1}_{i \notin \Theta_\ell^{(S)}} = 2$ : size of the OOB sample of the  $\ell$ -th tree
- $\mathbf{X}_{i,\pi_{j\ell}}$ :  $i$ -th observation where the  $j$ -th component is permuted across the OOB sample of the  $\ell$ -th tree

$X^{(1)}$	$X^{(2)}$	...	$X^{(j)}$	...	$X^{(p)}$	$Y$	$X^{(1)}$	$X^{(2)}$	...	$X^{(j)}$	...	$X^{(p)}$	$Y$
2.1	4.3	...	0.1	...	2.6	2.3	2.1	4.3	...	0.1	...	2.6	2.3
1.7	4.1	...	9.2	...	3.8	0.4	1.7	4.1	...	6.7	...	3.8	0.4
3.4	9.2	...	3.2	...	3.6	10.2	3.4	9.2	...	3.2	...	3.6	10.2
5.6	1.2	...	8.2	...	4.2	9.1	5.6	1.2	...	8.2	...	4.2	9.1
8.9	6.8	...	6.7	...	2.9	4.5	8.9	6.8	...	9.2	...	2.9	4.5

$\mathbf{X}_i$

$\mathbf{X}_{i,\pi_{j\ell}}$

- $i \in \{1, \dots, n\} \setminus \Theta_\ell^{(S)} = \{2, 5\}$ : OOB sample of the  $\ell$ -th tree
- $N_{n,\ell} = \sum_{i=1}^n \mathbb{1}_{i \notin \Theta_\ell^{(S)}} = 2$ : size of the OOB sample of the  $\ell$ -th tree
- $\mathbf{X}_{i,\pi_{j\ell}}$ :  $i$ -th observation where the  $j$ -th component is permuted across the OOB sample of the  $\ell$ -th tree

$$\widehat{\text{MDA}}_{M,n}^{(BC)}(\mathbf{X}^{(j)}) = \frac{1}{M} \sum_{\ell=1}^M \frac{1}{N_{n,\ell}} \sum_{i=1}^n [(Y_i - m_n(\mathbf{X}_{i,\pi_{j\ell}}, \Theta_\ell))^2 - (Y_i - m_n(\mathbf{X}_i, \Theta_\ell))^2] \mathbb{1}_{i \notin \Theta_\ell^{(S)}}$$

- $i \in \{1, \dots, n\} \setminus \Theta_\ell^{(S)} = \{2, 5\}$ : OOB sample of the  $\ell$ -th tree
- $N_{n,\ell} = \sum_{i=1}^n \mathbb{1}_{i \notin \Theta_\ell^{(S)}} = 2$ : size of the OOB sample of the  $\ell$ -th tree
- $\mathbf{X}_{i,\pi_{j\ell}}$ :  $i$ -th observation where the  $j$ -th component is permuted across the OOB sample of the  $\ell$ -th tree

$$\widehat{\text{MDA}}_{M,n}^{(BC)}(X^{(j)}) = \frac{1}{M} \sum_{\ell=1}^M \frac{1}{N_{n,\ell}} \sum_{i=1}^n [(Y_i - m_n(\mathbf{X}_{i,\pi_{j\ell}}, \Theta_\ell))^2 - (Y_i - m_n(\mathbf{X}_i, \Theta_\ell))^2] \mathbb{1}_{i \notin \Theta_\ell^{(S)}}$$

Quadratic risk of the  $\ell$ -th tree

- $i \in \{1, \dots, n\} \setminus \Theta_\ell^{(S)} = \{2, 5\}$ : OOB sample of the  $\ell$ -th tree
- $N_{n,\ell} = \sum_{i=1}^n \mathbb{1}_{i \notin \Theta_\ell^{(S)}} = 2$ : size of the OOB sample of the  $\ell$ -th tree
- $\mathbf{X}_{i,\pi_{j\ell}}$ :  $i$ -th observation where the  $j$ -th component is permuted across the OOB sample of the  $\ell$ -th tree

$$\widehat{\text{MDA}}_{M,n}^{(BC)}(\mathbf{X}^{(j)}) = \frac{1}{M} \sum_{\ell=1}^M \frac{1}{N_{n,\ell}} \sum_{i=1}^n \left[ (Y_i - m_n(\mathbf{X}_{i,\pi_{j\ell}}, \Theta_\ell))^2 - (Y_i - m_n(\mathbf{X}_i, \Theta_\ell))^2 \right] \mathbb{1}_{i \notin \Theta_\ell^{(S)}}$$

Inflated quadratic risk of the  $\ell$ -th tree where  $\mathbf{X}^{(j)}$  is permuted

- $i \in \{1, \dots, n\} \setminus \Theta_\ell^{(S)} = \{2, 5\}$ : OOB sample of the  $\ell$ -th tree
- $N_{n,\ell} = \sum_{i=1}^n \mathbb{1}_{i \notin \Theta_\ell^{(S)}} = 2$ : size of the OOB sample of the  $\ell$ -th tree
- $\mathbf{X}_{i,\pi_{j\ell}}$ :  $i$ -th observation where the  $j$ -th component is permuted across the OOB sample of the  $\ell$ -th tree

$$\widehat{\text{MDA}}_{M,n}^{(BC)}(\mathbf{X}^{(j)}) = \frac{1}{M} \sum_{\ell=1}^M \frac{1}{N_{n,\ell}} \sum_{i=1}^n [(Y_i - m_n(\mathbf{X}_{i,\pi_{j\ell}}, \Theta_\ell))^2 - (Y_i - m_n(\mathbf{X}_i, \Theta_\ell))^2] \mathbb{1}_{i \notin \Theta_\ell^{(S)}}$$

Risks are computed over the OOB sample of each tree

- $i \in \{1, \dots, n\} \setminus \Theta_\ell^{(S)} = \{2, 5\}$ : OOB sample of the  $\ell$ -th tree
- $N_{n,\ell} = \sum_{i=1}^n \mathbb{1}_{i \notin \Theta_\ell^{(S)}} = 2$ : size of the OOB sample of the  $\ell$ -th tree
- $\mathbf{X}_{i,\pi_{j\ell}}$ :  $i$ -th observation where the  $j$ -th component is permuted across the OOB sample of the  $\ell$ -th tree

$$\widehat{\text{MDA}}_{M,n}^{(BC)}(\mathbf{X}^{(j)}) = \frac{1}{M} \sum_{\ell=1}^M \frac{1}{N_{n,\ell}} \sum_{i=1}^n [(Y_i - m_n(\mathbf{X}_{i,\pi_{j\ell}}, \Theta_\ell))^2 - (Y_i - m_n(\mathbf{X}_i, \Theta_\ell))^2] \mathbb{1}_{i \notin \Theta_\ell^{(S)}}$$

Average over all trees



- 1 Introduction
- 2 MDA Theoretical Limitations
  - MDA definition
  - MDA convergence
- 3 Sobol-MDA
- 4 Shapley effects

(A1)

The response  $Y \in \mathbb{R}$  follows

$$Y = m(\mathbf{X}) + \varepsilon$$

where

- $\mathbf{X} = (X^{(1)}, \dots, X^{(p)}) \in [0, 1]^p$
- $\mathbf{X}$  admits a density  $f$  such that  $c_1 < f(\mathbf{x}) < c_2$ , with constants  $c_1, c_2 > 0$
- $m$  is continuous
- the noise  $\varepsilon$  is sub-Gaussian and centered

# Assumptions

(A2): the theoretical tree is consistent  
(always true with slight modifications of the forest algorithm)

(A2): the theoretical tree is consistent  
(always true with slight modifications of the forest algorithm)

(A2)

*The randomized theoretical CART tree built with the distribution of  $(\mathbf{X}, Y)$  is consistent, that is, for all  $\mathbf{x} \in [0, 1]^p$ , almost surely,*

$$\lim_{k \rightarrow \infty} \Delta(m, A_k^*(\mathbf{x}, \Theta)) = 0.$$

# Assumptions

(A2): the theoretical tree is consistent  
(always true with slight modifications of the forest algorithm)

(A2)

*The randomized theoretical CART tree built with the distribution of  $(\mathbf{X}, Y)$  is consistent, that is, for all  $\mathbf{x} \in [0, 1]^p$ , almost surely,*

$$\lim_{k \rightarrow \infty} \Delta(m, A_k^*(\mathbf{x}, \Theta)) = 0.$$

(A3): tree partition is not too complex with respect to  $n$

# Assumptions

(A2): the theoretical tree is consistent  
(always true with slight modifications of the forest algorithm)

(A2)

*The randomized theoretical CART tree built with the distribution of  $(\mathbf{X}, Y)$  is consistent, that is, for all  $\mathbf{x} \in [0, 1]^p$ , almost surely,*

$$\lim_{k \rightarrow \infty} \Delta(m, A_k^*(\mathbf{x}, \Theta)) = 0.$$

(A3): tree partition is not too complex with respect to  $n$

(A3)

*The asymptotic regime of  $a_n$ , the size of the subsampling without replacement, and the number of terminal leaves  $t_n$  is such that  $a_n \leq n - 2$ ,  $a_n/n < 1 - \kappa$  for a fixed  $\kappa > 0$ ,  $\lim_{n \rightarrow \infty} a_n = \infty$ ,  $\lim_{n \rightarrow \infty} t_n = \infty$ , and  $\lim_{n \rightarrow \infty} t_n \frac{(\log(a_n))^9}{a_n} = 0$ .*

## Theorem (Bénard et al. (2021))

If Assumptions (A1), (A2), and (A3) are satisfied, then, for all  $M \in \mathbb{N}^*$  and  $j \in \{1, \dots, p\}$  we have

$$\widehat{MDA}_{M,n}^{(BC)}(X^{(j)}) \xrightarrow{\mathbb{L}^1} \mathbb{E}[(m(\mathbf{X}) - m(\mathbf{X}_{\pi_j}))^2]$$

$\mathbf{X}_{\pi_j}$ :  $\mathbf{X}$  where the  $j$ -th component is replaced by an independent copy, i.e.

$$\mathbf{X}_{\pi_j} = (X^{(1)}, \dots, X^{(j)}, \dots, X^{(p)})$$

**Limit interpretation?**

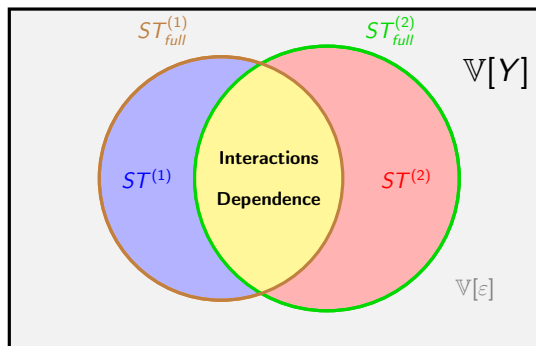


Figure: Standard and full total Sobol indices for  $Y = m(\mathbf{X}^{(1)}, \mathbf{X}^{(2)}) + \varepsilon$ .

**Total Sobol index** (Sobol, 1993)

$$ST^{(1)} = \frac{\mathbb{E}[V(m(\mathbf{X})|\mathbf{X}^{(-1)})]}{V(Y)}$$

**Full total Sobol index** (Mara et al., 2015; Benoumechiara, 2019)

$$ST_{full}^{(1)} = \frac{\mathbb{E}[V(m(\mathbf{X}_{\pi_j})|\mathbf{X}^{(-1)})]}{V(Y)}$$



## Proposition (Bénard et al. (2021))

If Assumptions (A1), (A2) and (A3) are satisfied, then for all  $M \in \mathbb{N}^*$  and  $j \in \{1, \dots, p\}$  we have

$$\widehat{MDA}_{M,n}^{(BC)}(\mathbf{X}^{(j)}) \xrightarrow{\mathbb{L}^1} \mathbb{V}[Y] \times ST^{(j)} + \mathbb{V}[Y] \times ST_{full}^{(j)} + MDA_3^{*(j)}.$$

The term  $MDA_3^{*(j)}$  is not an importance measure and is defined by

$$MDA_3^{*(j)} = \mathbb{E}[(\mathbb{E}[m(\mathbf{X})|\mathbf{X}^{(-j)}] - \mathbb{E}[m(\mathbf{X}_{\pi_j})|\mathbf{X}^{(-j)}])^2].$$

## Proposition (Bénard et al. (2021))

If Assumptions (A1), (A2) and (A3) are satisfied, then for all  $M \in \mathbb{N}^*$  and  $j \in \{1, \dots, p\}$  we have

$$(i) \quad \widehat{MDA}_{M,n}^{(TT)}(X^{(j)}) \xrightarrow{\mathbb{L}^1} \mathbb{V}[Y] \times ST^{(j)} + \mathbb{V}[Y] \times ST_{full}^{(j)} + MDA_3^{*(j)}$$

$$(ii) \quad \widehat{MDA}_{M,n}^{(BC)}(X^{(j)}) \xrightarrow{\mathbb{L}^1} \mathbb{V}[Y] \times ST^{(j)} + \mathbb{V}[Y] \times ST_{full}^{(j)} + MDA_3^{*(j)}.$$

If additionally  $M \rightarrow \infty$ , then

$$(iii) \quad \widehat{MDA}_{M,n}^{(IK)}(X^{(j)}) \xrightarrow{\mathbb{L}^1} \mathbb{V}[Y] \times ST^{(j)} + MDA_3^{*(j)}.$$

**If inputs  $\mathbf{X}$  are independent:**  $MDA_3^{*(j)} = 0$  and  $ST^{(j)} = ST_{full}^{(j)}$ .

Corollary (Bénard et al. (2021))

*If  $\mathbf{X}$  has independent components, and if Assumptions (A1)-(A3) are satisfied, for all  $M \in \mathbb{N}^*$  and  $j \in \{1, \dots, p\}$  we have*

$$\widehat{MDA}_{M,n}^{(TT)}(X^{(j)}) \xrightarrow{\mathbb{L}^1} 2\mathbb{V}[Y] \times ST^{(j)}$$
$$\widehat{MDA}_{M,n}^{(BC)}(X^{(j)}) \xrightarrow{\mathbb{L}^1} 2\mathbb{V}[Y] \times ST^{(j)}.$$

*If additionally  $M \rightarrow \infty$ , then*

$$\widehat{MDA}_{M,n}^{(IK)}(X^{(j)}) \xrightarrow{\mathbb{L}^1} \mathbb{V}[Y] \times ST^{(j)}.$$

This Corollary completes the result from (Gregorutti, 2015).

# Additive regression function

If  $m$  is additive:  $\text{MDA}_3^{*(j)} = 0$ .

Corollary (Bénard et al. (2021))

If the regression function  $m$  is additive, and if Assumptions (A1)-(A3) are satisfied, for all  $M \in \mathbb{N}^*$  and  $j \in \{1, \dots, p\}$  we have

$$\begin{aligned}\widehat{\text{MDA}}_{M,n}^{(TT)}(X^{(j)}) &\xrightarrow{\mathbb{L}^1} \mathbb{V}[Y] \times ST^{(j)} + \mathbb{V}[Y] \times ST_{full}^{(j)} \\ \widehat{\text{MDA}}_{M,n}^{(BC)}(X^{(j)}) &\xrightarrow{\mathbb{L}^1} \mathbb{V}[Y] \times ST^{(j)} + \mathbb{V}[Y] \times ST_{full}^{(j)}.\end{aligned}$$

If additionally  $M \rightarrow \infty$ , then

$$\widehat{\text{MDA}}_{M,n}^{(IK)}(X^{(j)}) \xrightarrow{\mathbb{L}^1} \mathbb{V}[Y] \times ST^{(j)}.$$

- When inputs  $\mathbf{X}$  are dependent and have interactions, the MDA is artificially inflated by the term  $MDA_3$  and is therefore misleading.

- When inputs  $\mathbf{X}$  are dependent and have interactions, the MDA is artificially inflated by the term  $MDA_3$  and is therefore misleading.
- MDA versions have different theoretical counterparts, and thus different meanings: be careful when using forest packages !

- When inputs  $\mathbf{X}$  are dependent and have interactions, the MDA is artificially inflated by the term  $MDA_3$  and is therefore misleading.
- MDA versions have different theoretical counterparts, and thus different meanings: be careful when using forest packages !
- For variable selection, the total Sobol index is the relevant component

$$\mathbb{V}[Y] \times ST^{(j)} + \cancel{\mathbb{V}[Y] \times ST_{full}^{(j)}} + \cancel{MDA_3^{+(j)}}$$

- When inputs  $\mathbf{X}$  are dependent and have interactions, the MDA is artificially inflated by the term  $MDA_3$  and is therefore misleading.
- MDA versions have different theoretical counterparts, and thus different meanings: be careful when using forest packages !
- For variable selection, the total Sobol index is the relevant component

$$\mathbb{V}[Y] \times ST^{(j)} + \cancel{\mathbb{V}[Y] \times ST_{full}^{(j)}} + \cancel{MDA_3^{+(j)}}$$

- We develop the Sobol-MDA: a fast and consistent estimate of  $ST^{(j)}$  for random forests



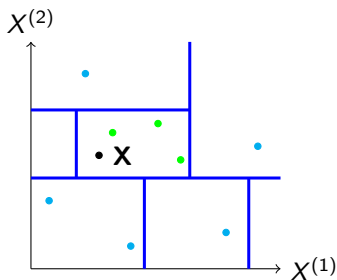
- 1 Introduction
- 2 MDA Theoretical Limitations
  - MDA definition
  - MDA convergence
- 3 Sobol-MDA
- 4 Shapley effects

Principle: **project** the partition of each tree along the  $j$ -th direction to remove  $X^{(j)}$  from the prediction process.

Principle: **project** the partition of each tree along the  $j$ -th direction to remove  $X^{(j)}$  from the prediction process.

$$\widehat{\text{S-MDA}}_{M,n}(X^{(j)}) = \frac{1}{\hat{\sigma}_Y^2} \frac{1}{n} \sum_{i=1}^n \left[ Y_i - m_{M,n}^{(-j, OOB)}(\mathbf{x}_i^{(-j)}, \Theta_M) \right]^2 - \left[ Y_i - m_{M,n}^{(OOB)}(\mathbf{x}_i, \Theta_M) \right]^2$$

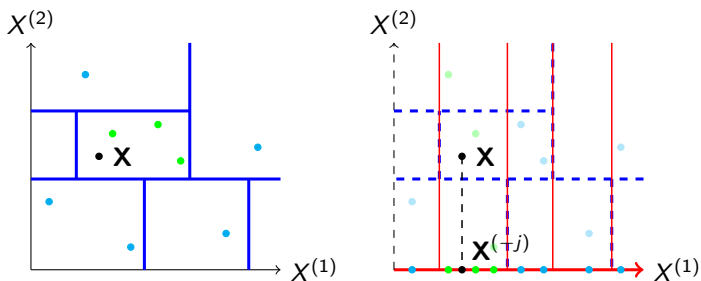
Principle: **project** the partition of each tree along the  $j$ -th direction to remove  $X^{(j)}$  from the prediction process.



**Figure:** Partition of  $[0, 1]^2$  by a random tree (left side) projected on the subspace span by  $\mathbf{X}^{(-2)} = X^{(1)}$  (right side), for  $p = 2$  and  $j = 2$ .

$$\widehat{\text{S-MDA}}_{M,n}(X^{(j)}) = \frac{1}{\hat{\sigma}_Y^2} \frac{1}{n} \sum_{i=1}^n \left[ Y_i - m_{M,n}^{(-j, \text{OOB})}(\mathbf{X}_i^{(-j)}, \Theta_M) \right]^2 - \left[ Y_i - m_{M,n}^{(\text{OOB})}(\mathbf{X}_i, \Theta_M) \right]^2$$

Principle: **project** the partition of each tree along the  $j$ -th direction to remove  $X^{(j)}$  from the prediction process.



**Figure:** Partition of  $[0, 1]^2$  by a random tree (left side) projected on the subspace span by  $\mathbf{X}^{(-2)} = X^{(1)}$  (right side), for  $p = 2$  and  $j = 2$ .

$$\widehat{\text{S-MDA}}_{M,n}(X^{(j)}) = \frac{1}{\hat{\sigma}_Y^2} \frac{1}{n} \sum_{i=1}^n \left[ Y_i - m_{M,n}^{(-j, \text{OOB})}(\mathbf{X}_i^{(-j)}, \Theta_M) \right]^2 - \left[ Y_i - m_{M,n}^{(\text{OOB})}(\mathbf{X}_i, \Theta_M) \right]^2$$

The Sobol-MDA recovers the appropriate theoretical counterpart for variable selection: the total Sobol index

Theorem (Bénard et al. (2021))

*If Assumptions (A1), (A2'), and (A3') are satisfied, for all  $M \in \mathbb{N}^*$  and  $j \in \{1, \dots, p\}$*

$$\widehat{S\text{-MDA}}_{M,n}(X^{(j)}) \xrightarrow{P} ST^{(j)}.$$

Settings (Archer and Kimes, 2008; Gregorutti et al., 2017)

- $p = 200$  input variables
- 5 independent groups of 40 variables
- each group is a Gaussian vector, strongly correlated

Settings (Archer and Kimes, 2008; Gregorutti et al., 2017)

- $p = 200$  input variables
- 5 independent groups of 40 variables
- each group is a Gaussian vector, strongly correlated
- 1 variable from each group involved in  $m$

$$m(\mathbf{X}) = 2X^{(1)} + X^{(41)} + X^{(81)} + X^{(121)} + X^{(161)}.$$

- independent Gaussian noise with  $\mathbb{V}[\varepsilon] = 10\% \mathbb{V}[Y]$

$$Y = m(\mathbf{X}) + \varepsilon$$



Settings (Archer and Kimes, 2008; Gregorutti et al., 2017)

- $p = 200$  input variables
- 5 independent groups of 40 variables
- each group is a Gaussian vector, strongly correlated
- 1 variable from each group involved in  $m$

$$m(\mathbf{X}) = 2X^{(1)} + X^{(41)} + X^{(81)} + X^{(121)} + X^{(161)}.$$

- independent Gaussian noise with  $\mathbb{V}[\varepsilon] = 10\% \mathbb{V}[Y]$

$$Y = m(\mathbf{X}) + \varepsilon$$

- $n = 1000$  observations
- $M = 300$  trees

# Sobol-MDA Experiments

$\widehat{\text{S-MDA}}$		$\widehat{\text{BC-MDA}}/2\mathbb{V}[Y]$		$\widehat{\text{IK-MDA}}/\mathbb{V}[Y]$	
$\mathbf{X}^{(1)}$	0.035	$\mathbf{X}^{(1)}$	0.048	$\mathbf{X}^{(1)}$	0.056
$\mathbf{X}^{(161)}$	0.005	$\mathbf{X}^{(25)}$	0.010	$\mathbf{X}^{(5)}$	0.009
$\mathbf{X}^{(81)}$	0.004	$\mathbf{X}^{(31)}$	0.008	$\mathbf{X}^{(81)}$	0.007
$\mathbf{X}^{(121)}$	0.004	$\mathbf{X}^{(14)}$	0.008	$\mathbf{X}^{(41)}$	0.005
$\mathbf{X}^{(41)}$	0.002	$\mathbf{X}^{(40)}$	0.007	$\mathbf{X}^{(161)}$	0.005
$\mathbf{X}^{(179)}$	0.002	$\mathbf{X}^{(3)}$	0.007	$\mathbf{X}^{(15)}$	0.005
$\mathbf{X}^{(13)}$	0.001	$\mathbf{X}^{(17)}$	0.006	$\mathbf{X}^{(121)}$	0.005
$\mathbf{X}^{(25)}$	0.001	$\mathbf{X}^{(26)}$	0.006	$\mathbf{X}^{(7)}$	0.005
$\mathbf{X}^{(73)}$	0.001	$\mathbf{X}^{(41)}$	0.006	$\mathbf{X}^{(4)}$	0.004
$\mathbf{X}^{(155)}$	0.001	$\mathbf{X}^{(121)}$	0.006	$\mathbf{X}^{(28)}$	0.004

Table: Sobol-MDA, normalized BC-MDA, and normalized IK-MDA estimates with influential variables in blue.

Additional experiments are available in B nard et al. (2021)  
(non-linear data with interactions and dependence)

- analytical example
- backward variable selection with real data

Sobol-MDA can be associated with any black-box algorithm

- fit a black box  $\hat{f}$  on  $\mathcal{D}_n$
- generate a large sample  $\mathcal{D}'_N$  with  $\hat{f}$
- run the Sobol-MDA with  $\mathcal{D}'_N$

- 1 Introduction
- 2 MDA Theoretical Limitations
  - MDA definition
  - MDA convergence
- 3 Sobol-MDA
- 4 Shapley effects

# Definition of Shapley effects

- Originally defined in economics and game theory (Shapley, 1953)

# Definition of Shapley effects

- Originally defined in economics and game theory (Shapley, 1953)
- Attribute the value produced by a joint team to its individual members

# Definition of Shapley effects

- Originally defined in economics and game theory (Shapley, 1953)
- Attribute the value produced by a joint team to its individual members
- Difference of produced value between a subset of the team and the same subteam with an additional member (averaged over all possible subteams).



# Definition of Shapley effects

- Originally defined in economics and game theory (Shapley, 1953)
- Attribute the value produced by a joint team to its individual members
- Difference of produced value between a subset of the team and the same subteam with an additional member (averaged over all possible subteams).
- Adapted by Owen (2014) to variable importance in machine learning:

# Definition of Shapley effects

- Originally defined in economics and game theory (Shapley, 1953)
- Attribute the value produced by a joint team to its individual members
- Difference of produced value between a subset of the team and the same subteam with an additional member (averaged over all possible subteams).
- Adapted by Owen (2014) to variable importance in machine learning:
  - member of the team = input variable

# Definition of Shapley effects

- Originally defined in economics and game theory (Shapley, 1953)
- Attribute the value produced by a joint team to its individual members
- Difference of produced value between a subset of the team and the same subteam with an additional member (averaged over all possible subteams).
- Adapted by Owen (2014) to variable importance in machine learning:
  - member of the team = input variable
  - value function = explained output variance

# Definition of Shapley effects

Formally, the Shapley effect of the  $j$ -th variable is defined by

$$Sh^*(X^{(j)}) = \sum_{U \subset \{1, \dots, p\} \setminus \{j\}} \frac{1}{p} \binom{p-1}{|U|}^{-1} \frac{\mathbb{V}[\mathbb{E}[Y | \mathbf{X}^{(U \cup \{j\})}]] - \mathbb{V}[\mathbb{E}[Y | \mathbf{X}^{(U)}]]}{\mathbb{V}[Y]}.$$

# Definition of Shapley effects

Formally, the Shapley effect of the  $j$ -th variable is defined by

$$Sh^*(X^{(j)}) = \sum_{U \subset \{1, \dots, p\} \setminus \{j\}} \frac{1}{p} \binom{p-1}{|U|}^{-1} \frac{\mathbb{V}[\mathbb{E}[Y | \mathbf{X}^{(U \cup \{j\})}]] - \mathbb{V}[\mathbb{E}[Y | \mathbf{X}^{(U)}]]}{\mathbb{V}[Y]}.$$

**Main property:** equitably allocate contributions due to dependence and interactions across input variables

# Definition of Shapley effects

Formally, the Shapley effect of the  $j$ -th variable is defined by

$$Sh^*(X^{(j)}) = \sum_{U \subset \{1, \dots, p\} \setminus \{j\}} \frac{1}{p} \binom{p-1}{|U|}^{-1} \frac{\mathbb{V}[\mathbb{E}[Y | \mathbf{X}^{(U \cup \{j\})}]] - \mathbb{V}[\mathbb{E}[Y | \mathbf{X}^{(U)}]]}{\mathbb{V}[Y]}.$$

**Main property:** equitably allocate contributions due to dependence and interactions across input variables

Two obstacles arise to estimate Shapley effects:

- 1 the computational complexity is exponential with the dimension  $p$

# Definition of Shapley effects

Formally, the Shapley effect of the  $j$ -th variable is defined by

$$Sh^*(X^{(j)}) = \sum_{U \subset \{1, \dots, p\} \setminus \{j\}} \frac{1}{p} \binom{p-1}{|U|}^{-1} \frac{\mathbb{V}[\mathbb{E}[Y|\mathbf{X}^{(U \cup \{j\})}]] - \mathbb{V}[\mathbb{E}[Y|\mathbf{X}^{(U)}]]}{\mathbb{V}[Y]}.$$

**Main property:** equitably allocate contributions due to dependence and interactions across input variables

Two obstacles arise to estimate Shapley effects:

- 1 the computational complexity is exponential with the dimension  $p$
- 2  $\mathbb{V}[\mathbb{E}[Y|\mathbf{X}^{(U)}]]$  requires a fast and accurate estimate for all variable subsets  $U \subset \{1, \dots, p\}$

# Definition of Shapley effects

Formally, the Shapley effect of the  $j$ -th variable is defined by

$$Sh^*(X^{(j)}) = \sum_{U \subset \{1, \dots, p\} \setminus \{j\}} \frac{1}{p} \binom{p-1}{|U|}^{-1} \frac{\mathbb{V}[\mathbb{E}[Y|\mathbf{X}^{(U \cup \{j\})}]] - \mathbb{V}[\mathbb{E}[Y|\mathbf{X}^{(U)}]]}{\mathbb{V}[Y]}.$$

**Main property:** equitably allocate contributions due to dependence and interactions across input variables

Two obstacles arise to estimate Shapley effects:

- 1 the computational complexity is exponential with the dimension  $p$   
**Literature: Monte-Carlo methods**
- 2  $\mathbb{V}[\mathbb{E}[Y|\mathbf{X}^{(U)}]]$  requires a fast and accurate estimate for all variable subsets  $U \subset \{1, \dots, p\}$



# Definition of Shapley effects

Formally, the Shapley effect of the  $j$ -th variable is defined by

$$Sh^*(X^{(j)}) = \sum_{U \subset \{1, \dots, p\} \setminus \{j\}} \frac{1}{p} \binom{p-1}{|U|}^{-1} \frac{\mathbb{V}[\mathbb{E}[Y|\mathbf{X}^{(U \cup \{j\})}]] - \mathbb{V}[\mathbb{E}[Y|\mathbf{X}^{(U)}]]}{\mathbb{V}[Y]}.$$

**Main property:** equitably allocate contributions due to dependence and interactions across input variables

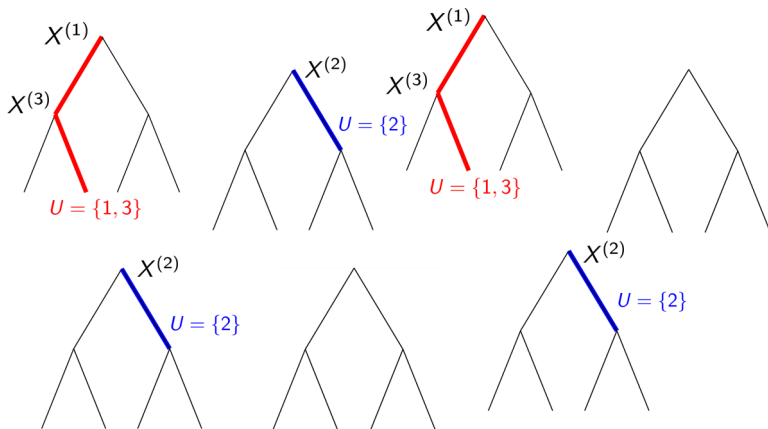
Two obstacles arise to estimate Shapley effects:

- 1 the computational complexity is exponential with the dimension  $p$   
**Literature:** Monte-Carlo methods
- 2  $\mathbb{V}[\mathbb{E}[Y|\mathbf{X}^{(U)}]]$  requires a fast and accurate estimate for all variable subsets  $U \subset \{1, \dots, p\}$   
**Literature:** strong approximation of the conditional distributions

# SHAFF: SHAPley effects via random Forests

SHAFF proceeds in three steps:

- 1 sample many subsets  $U$ , typically a few hundreds, based on their occurrence frequency  $\hat{p}_{M,n}(U)$  in the random forest



# SHAFF: SHAPley effects via random Forests

SHAFF proceeds in three steps:

- 1 sample many subsets  $U$ , typically a few hundreds, based on their occurrence frequency  $\hat{p}_{M,n}(U)$  in the random forest
- 2 estimate  $\mathbb{V}[\mathbb{E}[Y|\mathbf{X}^{(U)}]]$  with the projected forest algorithm for all selected  $U$  and their complementary sets  $\{1, \dots, p\} \setminus U$ :  $\hat{v}_{M,n}(U)$

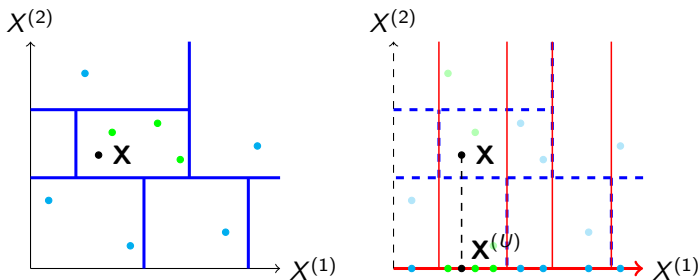


Figure: Partition of  $[0, 1]^2$  by a random tree (left side) projected on the subspace span by  $\mathbf{X}^{(U)} = X^{(1)}$  (right side), for  $p = 2$  and  $U = \{1\}$ .

# SHAFF: SHAPley effects via random Forests

**SHAFF** proceeds in three steps:

- 1 sample many subsets  $U$ , typically a few hundreds, based on their occurrence frequency  $\hat{p}_{M,n}(U)$  in the random forest
- 2 estimate  $\mathbb{V}[\mathbb{E}[Y|\mathbf{X}^{(U)}]]$  with the projected forest algorithm for all selected  $U$  and their complementary sets  $\{1, \dots, p\} \setminus U$ :  $\hat{v}_{M,n}(U)$
- 3 solve a weighted linear regression problem to recover Shapley effects  $\text{Sh}_{M,n}$  by minimizing in  $\beta$

$$\ell_{M,n}(\beta) = \frac{1}{K} \sum_{U \in \mathcal{U}_{n,K}} \frac{w(U)}{\hat{p}_{M,n}(U)} (\hat{v}_{M,n}(U) - \beta^T I(U))^2,$$

where  $w(U) = \frac{p-1}{\binom{p}{|U|} |U| (p-|U|)}$  and  $I(U)$  is the binary vector of dimension  $p$  where the  $j$ -th component takes the value 1 if  $j \in U$  and 0 otherwise.

(A4)

The number of Monte-Carlo sampling  $K_n$  and the number of trees  $M_n$  grow with  $n$ , such that  $M_n \rightarrow \infty$  and  $n.M_n/K_n \rightarrow 0$ .

Theorem

If Assumptions (A1), (A2'), (A3'), and (A4) are satisfied, then **SHAFF** is consistent, that is

$$\hat{\text{Sh}}_{M_n, n} \xrightarrow{P} \text{Sh}^*.$$

- Strong connections between the MDA and Sobol indices
- MDA does not target the appropriate quantity

- Strong connections between the MDA and Sobol indices
- MDA does not target the appropriate quantity
- Sobol-MDA fixes the flaws of original MDA
- R/C++ package `Sobo1MDA`, available online on Gitlab (<https://gitlab.com/drti/sobolmda>), and based on the package `ranger`

- Strong connections between the MDA and Sobol indices
- MDA does not target the appropriate quantity
- Sobol-MDA fixes the flaws of original MDA
- R/C++ package `Sobo1MDA`, available online on Gitlab (<https://gitlab.com/drti/sobolmda>), and based on the package `ranger`
- SHAFF: generalization of projected random forests to Shapley effects
- R/C++ package `shaff`, available online on Gitlab (<https://gitlab.com/drti/shaff>), and based on the package `ranger`



# Questions ?



- K.J. Archer and R.V. Kimes. Empirical characterization of random forest variable importance measures. *Computational Statistics & Data Analysis*, 52:2249–2260, 2008.
- C. Bénard, S. Da Veiga, and E. Scornet. Mda for random forests: inconsistency, and a practical solution via the sobol-mda. *arXiv preprint arXiv:2102.13347*, 2021.
- N. Benoumechiara. *Treatment of dependency in sensitivity analysis for industrial reliability*. PhD thesis, Sorbonne Université ; EDF R&D, 2019.
- L. Breiman. Random forests. *Machine Learning*, 45:5–32, 2001.
- B. Gregorutti. *Random forests and variable selection : analysis of the flight data recorders for aviation safety*. PhD thesis, Université Pierre et Marie Curie - Paris VI, 2015.
- B. Gregorutti, B. Michel, and P. Saint-Pierre. Correlation and variable importance in random forests. *Statistics and Computing*, 27:659–678, 2017.
- G. Hooker and L. Mentch. Please stop permuting features: an explanation and alternatives. *arXiv preprint arXiv:1905.03151*, 2019.
- H. Ishwaran. Variable importance in binary regression trees and forests. *Electronic Journal of Statistics*, 1:519–537, 2007.
- T. A Mara, S. Tarantola, and P. Annoni. Non-parametric methods for global sensitivity analysis of model output with dependent inputs. *Environmental Modelling & Software*, 72: 173–183, 2015.
- A.B. Owen. Sobol’indices and shapley value. *SIAM/ASA Journal on Uncertainty Quantification*, 2:245–251, 2014.
- L.S. Shapley. A value for n-person games. *Contributions to the Theory of Games*, 2:307–317, 1953.
- I.M. Sobol. Sensitivity estimates for nonlinear mathematical models. *Mathematical Modelling and Computational Experiments*, 1:407–414, 1993.
- C. Strobl, A.-L. Boulesteix, A. Zeileis, and T. Hothorn. Bias in random forest variable importance measures: illustrations, sources and a solution. *BMC bioinformatics*, 8:25, 2007.
- R. Zhu, D. Zeng, and M. R. Kosorok. Reinforcement learning trees. *Journal of the American Statistical Association*, 110:1770–1784, 2015.