

# **A first interpretable approach: SIRUS**

---

# Why do we need interpretability?

Machine learning is used for **decision support**.

Predicting is not enough

Understanding predictions is vital

- for Machine learning to be **accepted** (sensible applications in health, justice, defense)
- To **improve algorithms** (e.g., detect unfairness and try to correct it)

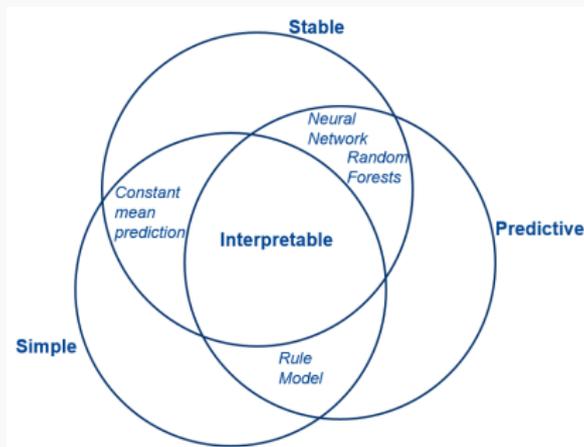
**Keywords:** trust, transparency, accountability, fairness, ethics.

**NIPS2017 debate:** Interpretability is necessary for Machine learning

<https://www.youtube.com/watch?v=93Xv8vJ2acI>

# Interpretable Models

- No agreement about a rigorous definition of interpretability [Lipton, 2016, Doshi-Velez and Kim, 2017, Murdoch et al., 2019]
- Minimum requirements for interpretability
  1. Simplicity [Murdoch et al., 2019]
  2. Stability [Yu, 2013]
  3. Predictivity [Breiman, 2001b]



# Existing Approaches

- Black-box models



E.g. Neural networks, Random forests

Combined with post-processing

E.g. variable importance  
sensitivity analysis  
local linearization

**Hard to operationalize**

# Existing Approaches

- Black-box models



E.g. Neural networks, Random forests

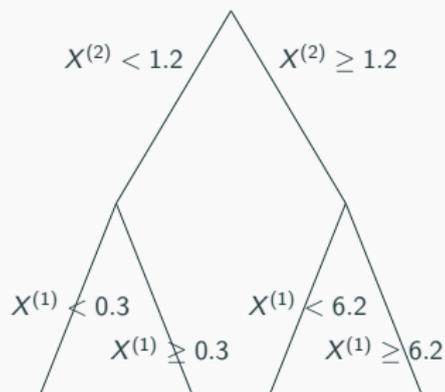
Combined with post-processing

E.g. variable importance  
sensitivity analysis  
local linearization

**Hard to operationalize**

- Interpretable models

E.g. decision trees, decision rules



**Unstable**

# SIRUS: Stable and Interpretable RULe Set

An example: SIRUS output on Titanic data set [Bénard et al., 2019]

**Average survival rate  $p_s = 39\%$ .**

|           |   |             |              |             |              |
|-----------|---|-------------|--------------|-------------|--------------|
| <b>if</b> | sex is male   | <b>then</b> | $p_s = 19\%$ | <b>else</b> | $p_s = 74\%$ |
| <b>if</b> | 1 <sup>st</sup> or 2 <sup>nd</sup> class                    | <b>then</b> | $p_s = 56\%$ | <b>else</b> | $p_s = 24\%$ |
| <b>if</b> | 1 <sup>st</sup> or 2 <sup>nd</sup> class<br>& sex is female | <b>then</b> | $p_s = 95\%$ | <b>else</b> | $p_s = 25\%$ |
| <b>if</b> | fare < 10.5£  | <b>then</b> | $p_s = 20\%$ | <b>else</b> | $p_s = 50\%$ |
| <b>if</b> | no parents or<br>children aboard                            | <b>then</b> | $p_s = 35\%$ | <b>else</b> | $p_s = 51\%$ |
| <b>if</b> | 2 <sup>st</sup> or 3 <sup>rd</sup> class<br>& sex is male   | <b>then</b> | $p_s = 14\%$ | <b>else</b> | $p_s = 64\%$ |
| <b>if</b> | sex is male<br>& age $\geq 15$                              | <b>then</b> | $p_s = 16\%$ | <b>else</b> | $p_s = 72\%$ |

## Principle

- Build a random forests and extract all decisions rules from all trees
- Select the rules that appear with a frequency larger than  $p_0$
- Aggregate the rules to obtain the final estimator.



## Principle

Frequent paths in random trees = strong and robust patterns in the data.

## Technical detail

- Preprocessing: discretize features based on their quantiles
- Random forests: building trees of depth 2

## Technical detail

- Preprocessing: discretize features based on their quantiles
- Random forests: building trees of depth 2

Probability that a  $\Theta$ -random tree contains a given path  $\mathcal{P} \in \Pi$

$$p_n(\mathcal{P}) = \mathbb{P}(\mathcal{P} \in T(\Theta, \mathcal{D}_n) | \mathcal{D}_n)$$

Selected paths

$$\hat{\mathcal{P}}_{M,n,p_0} = \{\mathcal{P} \in \Pi : \hat{p}_{M,n}(\mathcal{P}) > p_0\}$$

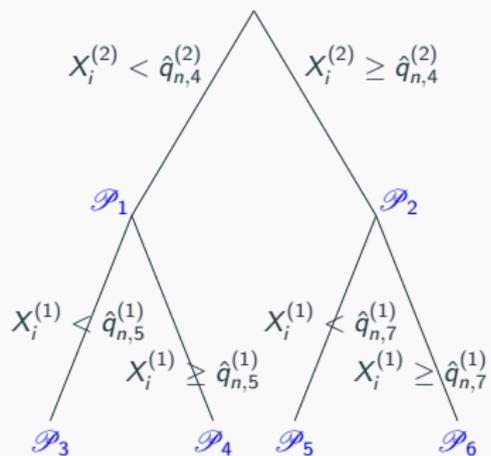
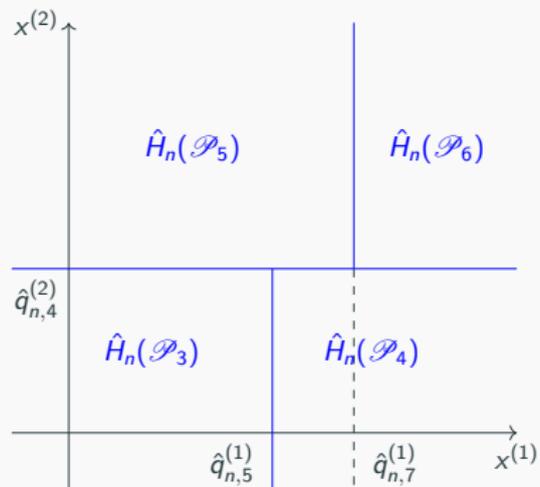
where

$$\hat{p}_{M,n}(\mathcal{P}) = \frac{1}{M} \sum_{\ell=1}^M \mathbb{1}_{\mathcal{P} \in T(\Theta_\ell, \mathcal{D}_n)}$$

is the Monte-Carlo estimate, directly computed using the random forest with  $M$  trees parametrized by  $\Theta_1, \dots, \Theta_M$ .

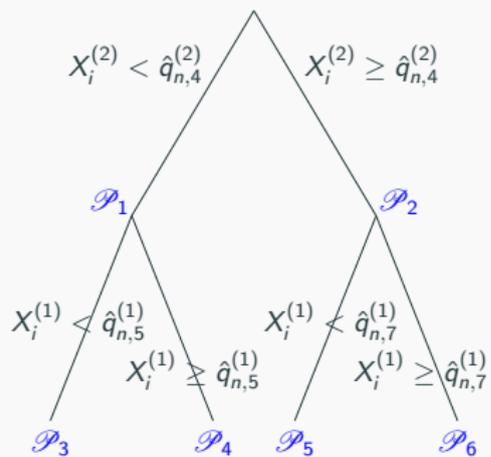
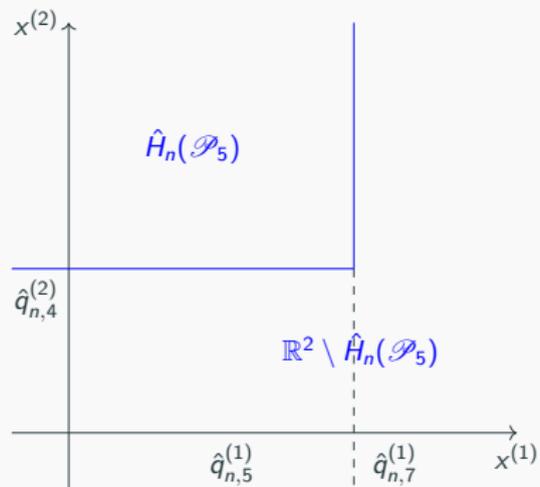
# SIRUS - Rule

How to recover a rule from a path ?

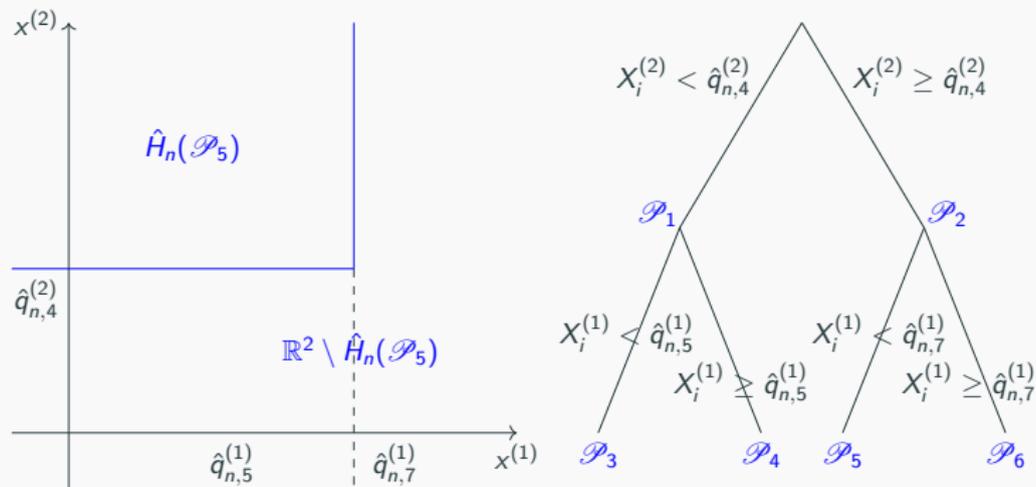


# SIRUS - Rule

How to recover a rule from a path ?

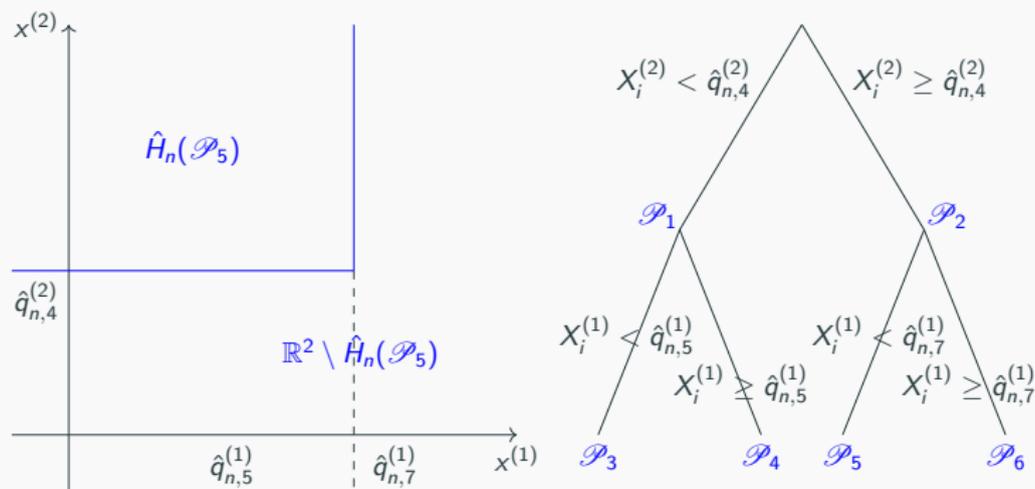


How to recover a rule from a path ?



$$\forall x \in \mathbb{R}^p, \quad \hat{g}_{n, \mathcal{P}}(x) = \begin{cases} \frac{1}{N_n(\hat{H}_n(\mathcal{P}))} \sum_{i=1}^n Y_i \mathbb{1}_{x_i \in \hat{H}_n(\mathcal{P})} & \text{if } x \in \hat{H}_n(\mathcal{P}) \\ \frac{1}{n - N_n(\hat{H}_n(\mathcal{P}))} \sum_{i=1}^n Y_i \mathbb{1}_{x_i \notin \hat{H}_n(\mathcal{P})} & \text{otherwise.} \end{cases}$$

How to recover a rule from a path ?



$$\forall x \in \mathbb{R}^p, \quad \hat{g}_{n, \mathcal{P}}(x) = \begin{cases} \frac{1}{N_n(\hat{H}_n(\mathcal{P}))} \sum_{i=1}^n Y_i \mathbb{1}_{x_i \in \hat{H}_n(\mathcal{P})} & \text{if } x \in \hat{H}_n(\mathcal{P}) \\ \frac{1}{n - N_n(\hat{H}_n(\mathcal{P}))} \sum_{i=1}^n Y_i \mathbb{1}_{x_i \notin \hat{H}_n(\mathcal{P})} & \text{otherwise.} \end{cases}$$

The final classifier corresponds to the averaging of all selected rules.

# Stability - definition

Define

- $\mathcal{D}'_n, \Theta'$  independent copies of  $\mathcal{D}_n$  and  $\Theta$
- $\hat{P}'_{M,n}(\mathcal{P}), \hat{\mathcal{P}}'_{M,n,p_0}$  built with  $\mathcal{D}'_n, \Theta'$

Dice-Sorensen index

$$\hat{S}_{M,n,p_0} = \frac{2|\hat{\mathcal{P}}_{M,n,p_0} \cap \hat{\mathcal{P}}'_{M,n,p_0}|}{|\hat{\mathcal{P}}_{M,n,p_0}| + |\hat{\mathcal{P}}'_{M,n,p_0}|}.$$

## Stability - a theoretical result

- (A1) The subsampling rate  $a_n$  satisfies  $\lim_{n \rightarrow \infty} a_n = \infty$  and  $\lim_{n \rightarrow \infty} \frac{a_n}{n} = 0$ .
- (A2) The number of trees  $M_n$  satisfies  $\lim_{n \rightarrow \infty} M_n = \infty$ .
- (A3)  $X$  has a density  $f$  with respect to the Lebesgue measure, continuous, bounded, and strictly positive.

Let  $\mathcal{U}^* = \{p^*(\mathcal{P}), \mathcal{P} \in \Pi\}$  be the set of all theoretical probabilities of appearance of all paths.

### Proposition Bénard et al. [2019]

Assume that Assumptions (A1)-(A3) are satisfied. Then, provided  $p_0 \in [0, 1] \setminus \mathcal{U}^*$ , we have

$$\lim_{n \rightarrow \infty} \hat{S}_{M_n, n, p_0} = 1, \quad \text{in probability.}$$

## Sketch of proof

The asymptotic stability of SIRUS comes from the two following points:

1. The bias of  $\hat{\rho}_{M_n, n}(\mathcal{P})$  tends to zero.
2. The variance of  $\hat{\rho}_{M_n, n}(\mathcal{P})$  tends to zero.

## Sketch of proof

The asymptotic stability of SIRUS comes from the two following points:

1. The bias of  $\hat{\rho}_{M_n, n}(\mathcal{P})$  tends to zero.
  - Prove that **CART-splitting criterion** is **consistent** and **asymptotically normal** when cuts are limited to empirical quantiles and the number of trees grows with  $n$  (A3).
2. The variance of  $\hat{\rho}_{M_n, n}(\mathcal{P})$  tends to zero.

# Sketch of proof

The asymptotic stability of SIRUS comes from the two following points:

1. The bias of  $\hat{p}_{M_n, n}(\mathcal{P})$  tends to zero.
  - Prove that **CART-splitting criterion** is **consistent** and **asymptotically normal** when cuts are limited to empirical quantiles and the number of trees grows with  $n$  (A3).
2. The variance of  $\hat{p}_{M_n, n}(\mathcal{P})$  tends to zero.

The variance can be decomposed into two terms:

- the variance generated by the Monte-Carlo randomization, which goes to 0 as the number of trees increases (A2).
- the variance of  $p_n(\mathcal{P})$ , which is a bagged estimate and thus an infinite-order U-statistic. The result comes from Mentch and Hooker [2016] since  $\lim_{n \rightarrow \infty} a_n/n = 0$  (A1).

# Numerical experiments

## Competitors:

- CART [Breiman et al., 1984]
- Classical rule learning: RIPPER [Cohen, 1995]
- Frequent pattern mining: CBA [Classification Based on Association Rules, Liu et al., 1998], BRL [Bayesian Rule List, Letham et al., 2015]
- Tree ensemble: RuleFit [Friedman and Popescu, 2008], Node Harvest [Meinshausen, 2010].

## Metrics:

- Accuracy/Error: 1-AUC
- Stability: Dice-Sorensen index
- Simplicity: Number of rules output by the procedure

# Accuracy

|                  | Black box     | Decision tree | Classical rule learning | Frequent pattern mining |                               | Tree ensemble                 |              |             |
|------------------|---------------|---------------|-------------------------|-------------------------|-------------------------------|-------------------------------|--------------|-------------|
| Dataset          | Random Forest | CART          | RIPPER                  | CBA                     | BRL                           | RuleFit                       | Node harvest | SIRUS       |
| Authentication   | $10^{-4}$     | 0.02          | 0.02                    | 0.14                    | 0.009                         | <b><math>9.10^{-4}</math></b> | 0.02         | 0.03        |
| Breast Wisconsin | 0.009         | 0.06          | 0.07                    | 0.05                    | 0.02                          | <b>0.01</b>                   | <b>0.01</b>  | <b>0.01</b> |
| Credit Approval  | 0.07          | 0.14          | 0.15                    | 0.14                    | 0.11                          | <b>0.08</b>                   | <b>0.07</b>  | 0.09        |
| Credit German    | 0.20          | 0.29          | 0.38                    | 0.40                    | 0.33                          | <b>0.23</b>                   | <b>0.26</b>  | <b>0.25</b> |
| Diabetes         | 0.17          | 0.25          | 0.29                    | 0.30                    | 0.25                          | <b>0.18</b>                   | <b>0.19</b>  | <b>0.19</b> |
| Haberman         | 0.31          | 0.48          | 0.39                    | 0.50                    | 0.43                          | <b>0.37</b>                   | <b>0.34</b>  | <b>0.35</b> |
| Heart C2         | 0.10          | 0.19          | 0.23                    | 0.17                    | 0.23                          | 0.12                          | 0.12         | <b>0.10</b> |
| Heart H2         | 0.11          | 0.23          | 0.24                    | 0.24                    | 0.16                          | <b>0.11</b>                   | <b>0.11</b>  | <b>0.12</b> |
| Heart Statlog    | 0.10          | 0.20          | 0.21                    | 0.17                    | 0.22                          | 0.12                          | 0.12         | <b>0.10</b> |
| Hepatitis        | 0.12          | 0.48          | 0.39                    | 0.36                    | 0.33                          | 0.20                          | 0.23         | <b>0.17</b> |
| Ionosphere       | 0.02          | 0.11          | 0.12                    | 0.13                    | 0.10                          | <b>0.04</b>                   | 0.07         | 0.07        |
| Kr vs Kp         | $9.10^{-4}$   | 0.02          | <b>0.009</b>            | 0.05                    | 0.01                          | <b>0.009</b>                  | 0.04         | 0.04        |
| Liver Disorders  | 0.23          | 0.33          | 0.35                    | 0.48                    | 0.44                          | <b>0.27</b>                   | 0.30         | 0.35        |
| Mushrooms        | 0             | 0.007         | $3.10^{-5}$             | $5.10^{-4}$             | <b><math>2.10^{-5}</math></b> | $5.10^{-4}$                   | 0.002        | $6.10^{-4}$ |
| Sonar            | 0.07          | 0.27          | 0.26                    | 0.25                    | 0.44                          | <b>0.12</b>                   | 0.16         | 0.2         |
| Spambase         | 0.01          | 0.11          | 0.08                    | 0.12                    | 0.05                          | <b>0.02</b>                   | 0.04         | 0.07        |

**Figure 1:** Model error (1-AUC) over a 10-fold cross-validation for UCI datasets. Results are averaged over 10 repetitions of the cross-validation. Values within 10% of the minimum are displayed in bold, random forest is put aside.

# Simplicity

|                  | Decision tree | Classical rule learning | Frequent pattern mining |     | Tree ensemble |              |       |
|------------------|---------------|-------------------------|-------------------------|-----|---------------|--------------|-------|
| Dataset          | CART          | RIPPER                  | CBA                     | BRL | RuleFit       | Node harvest | SIRUS |
| Authentication   | 21            | 7                       | 7                       | 17  | 49            | 30           | 13    |
| Breast Wisconsin | 7             | 12                      | 24                      | 7   | 24            | 32           | 24    |
| Credit Approval  | 5             | 4                       | 55                      | 4   | 15            | 27           | 16    |
| Credit German    | 18            | 3                       | 69                      | 4   | 33            | 33           | 20    |
| Diabetes         | 13            | 3                       | 17                      | 6   | 26            | 31           | 8     |
| Haberman         | 2             | 1                       | 2                       | 2   | 3             | 17           | 5     |
| Heart C2         | 10            | 3                       | 34                      | 4   | 23            | 36           | 20    |
| Heart H2         | 5             | 2                       | 29                      | 3   | 12            | 24           | 12    |
| Heart Statlog    | 10            | 3                       | 27                      | 4   | 22            | 35           | 16    |
| Hepatitis        | 2             | 2                       | 14                      | 2   | 8             | 14           | 12    |
| Ionosphere       | 4             | 4                       | 38                      | 4   | 20            | 35           | 15    |
| Kr vs Kp         | 16            | 15                      | 29                      | 9   | 18            | 13           | 24    |
| Liver Disorders  | 15            | 3                       | 2                       | 3   | 19            | 33           | 17    |
| Mushrooms        | 4             | 8                       | 25                      | 11  | 10            | 22           | 23    |
| Sonar            | 6             | 4                       | 33                      | 2   | 32            | 83           | 19    |
| Spambase         | 14            | 16                      | 126                     | 16  | 68            | 60           | 21    |

**Figure 2:** Mean model size over a 10-fold cross-validation for UCI datasets. Results are averaged over 10 repetitions of the cross-validation.

# Stability

| Dataset          | Decision tree | Classical rule learning | Frequent pattern mining |             | Tree ensemble |              |             |
|------------------|---------------|-------------------------|-------------------------|-------------|---------------|--------------|-------------|
|                  | CART          | RIPPER                  | CBA                     | BRL         | RuleFit       | Node harvest | SIRUS       |
| Authentication   | 0.41          | 0.36                    | <b>0.87</b>             | <b>0.86</b> | 0.48          | 0.59         | <b>0.81</b> |
| Breast Wisconsin | 0.21          | 0.55                    | <b>0.80</b>             | <b>0.78</b> | 0.34          | 0.71         | 0.70        |
| Credit Approval  | 0.52          | 0.32                    | 0.43                    | 0.52        | 0.25          | 0.23         | <b>0.75</b> |
| Credit German    | 0.46          | 0.22                    | 0.51                    | 0.41        | 0.24          | 0.48         | <b>0.66</b> |
| Diabetes         | 0.29          | 0.21                    | 0.46                    | <b>0.73</b> | 0.39          | 0.45         | <b>0.81</b> |
| Haberman         | <b>0.83</b>   | 0.09                    | <b>0.79</b>             | 0.50        | 0.46          | 0.52         | 0.65        |
| Heart C2         | 0.25          | 0.35                    | 0.38                    | 0.60        | 0.39          | 0.49         | <b>0.71</b> |
| Heart H2         | 0.46          | 0.27                    | 0.52                    | <b>0.73</b> | 0.29          | 0.29         | <b>0.65</b> |
| Heart Statlog    | 0.30          | 0.41                    | 0.41                    | <b>0.75</b> | 0.35          | 0.48         | <b>0.83</b> |
| Hepatitis        | 0.26          | 0.16                    | 0.24                    | 0.34        | 0.26          | 0.49         | <b>0.68</b> |
| Ionosphere       | <b>0.96</b>   | 0.39                    | 0.13                    | 0.70        | 0.17          | 0.33         | 0.69        |
| Kr vs Kp         | 0.71          | 0.74                    | <b>0.84</b>             | <b>0.80</b> | 0.19          | 0.27         | <b>0.87</b> |
| Liver Disorders  | 0.23          | 0.10                    | <b>0.91</b>             | 0.50        | 0.24          | 0.31         | 0.58        |
| Mushrooms        | <b>1</b>      | 0.84                    | <b>0.98</b>             | 0.80        | 0.69          | 0.48         | 0.86        |
| Sonar            | 0.34          | 0.04                    | 0.09                    | 0.19        | 0.09          | 0.20         | <b>0.55</b> |
| Spambase         | 0.49          | 0.10                    | 0.46                    | <b>0.86</b> | 0.28          | 0.66         | <b>0.78</b> |

**Figure 3:** Mean stability over a 10-fold cross-validation for UCI datasets.

Results are averaged over 10 repetitions of the cross-validation. Values within 10% of the maximum are displayed in bold.

# Conclusion on SIRUS

- Output a small, stable, and predictive set of rules
  - Predictive performances are on par with RF
  - Stability and number of rules improved over state-of-the-art algorithms
- Theoretical guarantees of stability for SIRUS
- Relies heavily on quantile discretization and a limited tree depth

## **Post-hoc methods: Sobol indices and Shapley effects**

---

## Introduction - Industrial Context

A first interpretable approach: SIRUS

Post-hoc methods: Sobol indices and Shapley effects

### Introduction

#### MDA Theoretical Limitations

MDA definition

MDA convergence

#### Sobol-MDA

#### Shapley effects

- MDA [Breiman, 2001a]: built-in variable importance algorithm for random forests

- MDA [Breiman, 2001a]: built-in variable importance algorithm for random forests
- MDA is used intensively

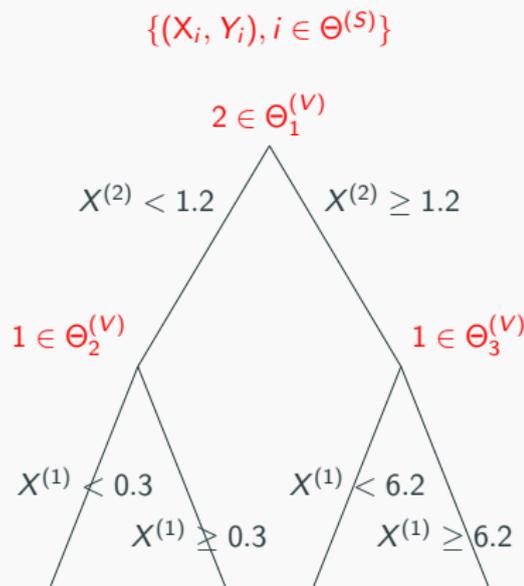
- MDA [Breiman, 2001a]: built-in variable importance algorithm for random forests
- MDA is used intensively
- MDA has flaws
  - Poor understanding of the MDA: what is estimated ?
  - Empirical studies show that the MDA is biased for dependent inputs [Strobl et al., 2007, Gregorutti et al., 2017, Hooker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

- MDA [Breiman, 2001a]: 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  $\mathcal{D}_n = \{(X_i, Y_i), i = 1, \dots, n\}$ ,  
where  $(X_i, Y_i) \sim \mathbb{P}_{\mathbf{X}, Y}$ .
- 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  $\mathbf{X}$



## Introduction - Industrial Context

A first interpretable approach: SIRUS

Post-hoc methods: Sobol indices and Shapley effects

Introduction

MDA Theoretical Limitations

MDA definition

MDA convergence

Sobol-MDA

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

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 1:** 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 1:** 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 1:** Example of the permutation of a dataset  $\mathcal{D}_n$  for  $n = 5$ .

# MDA illustration

| $X^{(1)}$ | $X^{(2)}$ | ... | $X^{(j)}$ | ... | $X^{(p)}$ | $Y$  |
|-----------|-----------|-----|-----------|-----|-----------|------|
| 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 1:** Example of the permutation of a dataset  $\mathcal{D}_n$  for  $n = 5$ .

Explained variance of  $Y = 16.4$

Explained variance of  $Y = 13.7$

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

## MDA illustration

| $X^{(1)}$ | $X^{(2)}$ | ... | $X^{(j)}$ | ... | $X^{(p)}$ | $Y$  |
|-----------|-----------|-----|-----------|-----|-----------|------|
| 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 1:** 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?

## MDA versions

The explained variance estimate of MDA algorithms differ across implementations

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

## MDA versions

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

## MDA versions

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

# MDA versions

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<br>randomForestSRC                       | Forest | Testing dataset |
| Breiman-Cutler   | randomForest (normalized)<br>ranger / randomForestSRC | Tree   | OOB sample      |
| Ishwaran-Kogalur | randomForestSRC                                       | Forest | OOB sample      |

## Breiman-Cutler MDA

- $i \in \{1, \dots, n\} \setminus \Theta_\ell^{(S)} = \{2, 5\}$ : 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
- $X_{i,\pi_{j\ell}}$ :  $i$ -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  |

$X_i$

$X_{i,\pi_{j\ell}}$

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

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

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

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

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
- $X_{i,\pi_{j\ell}}$ :  $i$ -th observation where the  $j$ -th component is permuted across the OOB sample of the  $\ell$ -th tree

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

Inflated quadratic risk of the  $\ell$ -th tree where  $X^{(j)}$  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
- $X_{i,\pi_{j\ell}}$ :  $i$ -th observation where the  $j$ -th component is permuted across the OOB sample of the  $\ell$ -th tree

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

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
- $X_{i,\pi_{j\ell}}$ :  $i$ -th observation where the  $j$ -th component is permuted across the OOB sample of the  $\ell$ -th tree

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

Average over all trees

## Introduction - Industrial Context

A first interpretable approach: SIRUS

Post-hoc methods: Sobol indices and Shapley effects

Introduction

MDA Theoretical Limitations

MDA definition

MDA convergence

Sobol-MDA

Shapley effects

# Assumptions

(A1)

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

$$Y = m(X) + \varepsilon$$

where

- $X = (X^{(1)}, \dots, X^{(p)}) \in [0, 1]^p$
- $X$  admits a density  $f$  such that  $c_1 < f(x) < c_2$ , 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)

# 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  $(X, Y)$  is consistent, that is, for all  $x \in [0, 1]^p$ , almost surely,*

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

# 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  $(X, Y)$  is consistent, that is, for all  $x \in [0, 1]^p$ , almost surely,*

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

(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  $(X, Y)$  is consistent, that is, for all  $x \in [0, 1]^p$ , almost surely,*

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

(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

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

$X_{\pi_j}$ :  $X$  where the  $j$ -th component is replaced by an independent copy, i.e.  
 $X_{\pi_j} = (X^{(1)}, \dots, X^{(j)}, \dots, X^{(p)})$

**Limit interpretation?**

# Sensitivity analysis

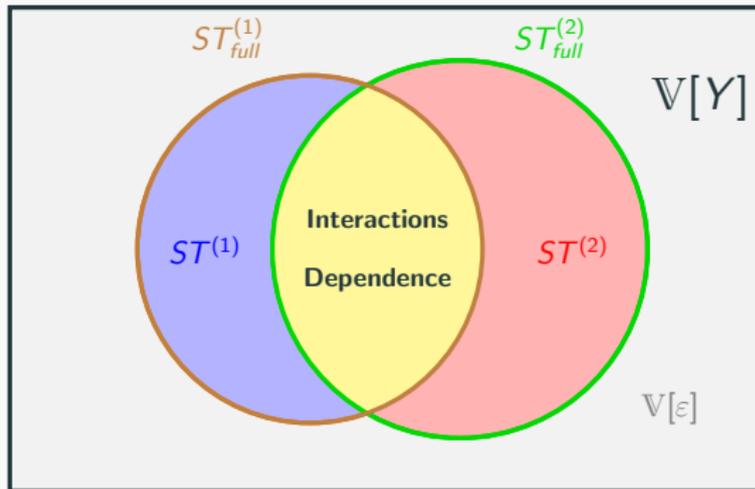


Figure 4: Standard and full total Sobol indices for  $Y = m(X^{(1)}, X^{(2)}) + \varepsilon$ .

**Total Sobol index** [Sobol, 1993]

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

**Full total Sobol index** [Mara et al., 2015, Benoumechiara, 2019]

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

## 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

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

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

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

## 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

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

# Independent inputs

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

## Corollary (Bénard et al. [2021])

If  $X$  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

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

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

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  $X$  are dependent and have interactions, the MDA is artificially inflated by the term  $MDA_3$  and is therefore misleading.

- When inputs  $X$  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  $X$  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

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

## MDA summary

- When inputs  $X$  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

## Introduction - Industrial Context

A first interpretable approach: SIRUS

Post-hoc methods: Sobol indices and Shapley effects

Introduction

MDA Theoretical Limitations

MDA definition

MDA convergence

Sobol-MDA

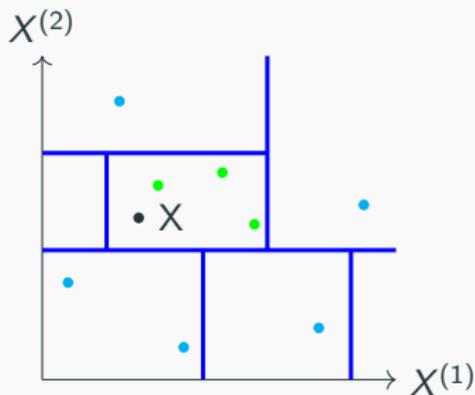
Shapley effects

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.

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

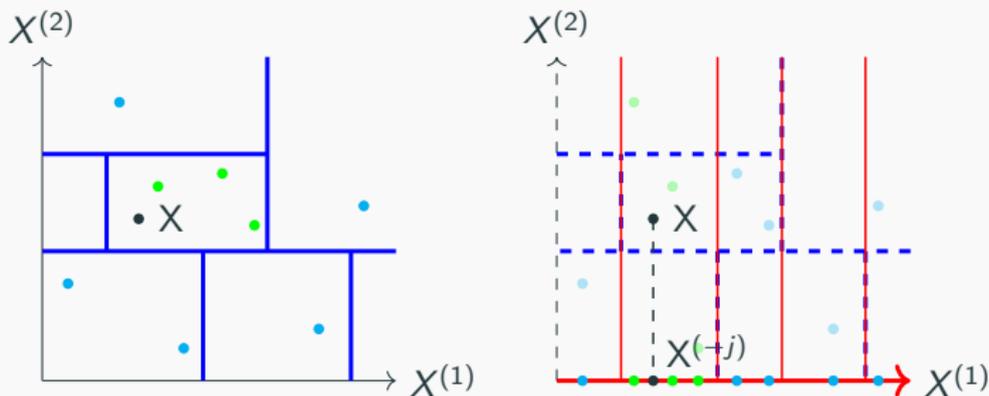
Principle: **project** the partition of each tree along the  $j$ -th direction to remove  $X^{(j)}$  from the prediction process.



**Figure 5:** Partition of  $[0, 1]^2$  by a random tree (left side) projected on the subspace span by  $X^{(-2)} = X^{(1)}$  (right side), for  $p = 2$  and  $j = 2$ .

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

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



**Figure 5:** Partition of  $[0, 1]^2$  by a random tree (left side) projected on the subspace span by  $X^{(-2)} = X^{(1)}$  (right side), for  $p = 2$  and  $j = 2$ .

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

# Consistency of the Sobol-MDA

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  $j \in \{1, \dots, p\}$*

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

# Sobol-MDA Experiments

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

# Sobol-MDA Experiments

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  $\mathbb{V}[\varepsilon] = 10\% \mathbb{V}[Y]$

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

# Sobol-MDA Experiments

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

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

**Table 3:** Sobol-MDA, normalized BC-MDA, and normalized IK-MDA estimates with influential variables in blue.

# Additional Experiments

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  $\mathcal{D}'_N$

## Introduction - Industrial Context

A first interpretable approach: SIRUS

Post-hoc methods: Sobol indices and Shapley effects

Introduction

MDA Theoretical Limitations

MDA definition

MDA convergence

Sobol-MDA

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

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

## Definition of Shapley effects

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

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

**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

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

**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

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

**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.  $\mathbb{V}[\mathbb{E}[Y|X^{(U)}]]$  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

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

**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.  $\mathbb{V}[\mathbb{E}[Y|X^{(U)}]]$  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

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

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

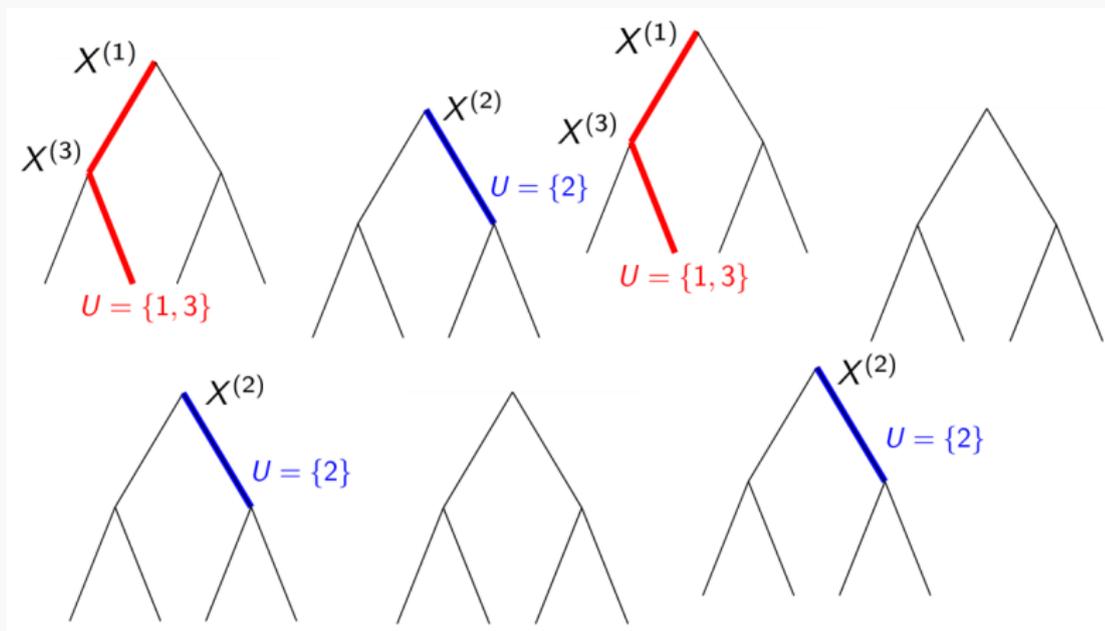
Two obstacles arise to estimate Shapley effects:

1. the computational complexity is exponential with the dimension  $p$   
**Literature:** Monte-Carlo methods
2.  $\mathbb{V}[\mathbb{E}[Y|X^{(U)}]]$  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:

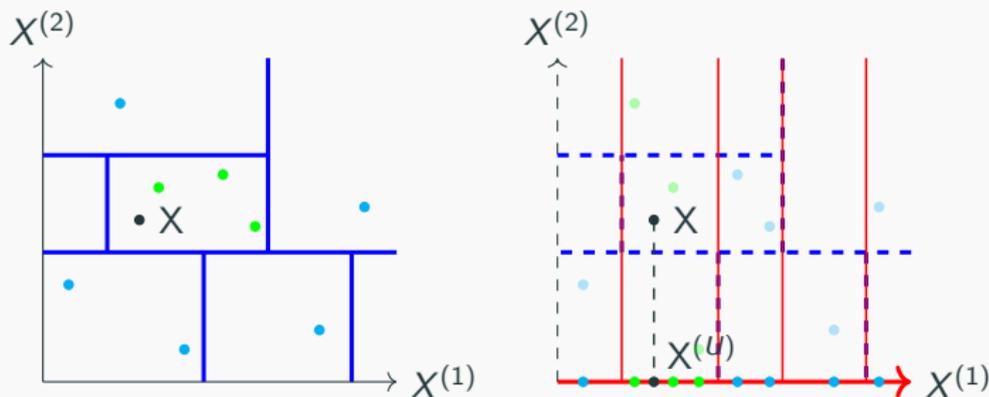
1. sample many subsets  $U$ , typically a few hundreds, based on their occurrence frequency  $\hat{p}_{M,n}(U)$  in the random forest



# 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  $\mathbb{V}[\mathbb{E}[Y|X^{(U)}]]$  with the projected forest algorithm for all selected  $U$  and their complementary sets  $\{1, \dots, p\} \setminus U$ :  $\hat{v}_{M,n}(U)$



**Figure 6:** Partition of  $[0, 1]^2$  by a random tree (left side) projected on the subspace span by  $X^{(U)} = X^{(1)}$  (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  $\mathbb{V}[\mathbb{E}[Y|X^{(U)}]]$  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  $\hat{Sh}_{M,n}$  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*

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

# Conclusion

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

# Conclusion

- 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`

# Conclusion

- 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 ?



## References

---

- 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, G. Biau, S. Da Veiga, and E. Scornet. Sirius: making random forests interpretable. *arXiv preprint arXiv:1908.06852*, 2019.
- 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, 2001a.

- L. Breiman. Statistical modeling: The two cultures (with comments and a rejoinder by the author). *Statistical Science*, 16:199–231, 2001b.
- L. Breiman, J. Friedman, R. Olshen, and C. Stone. *Classification and Regression Trees*. Chapman & Hall, New York, 1984.
- William W Cohen. Fast effective rule induction. In *Machine learning proceedings 1995*, pages 115–123. Elsevier, 1995.
- F. Doshi-Velez and B. Kim. Towards a rigorous science of interpretable machine learning. *arXiv:1702.08608*, 2017.
- Jerome H Friedman and Bogdan E Popescu. Predictive learning via rule ensembles. *The Annals of Applied Statistics*, 2(3):916–954, 2008.
- 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.
- Benjamin Letham, Cynthia Rudin, Tyler H McCormick, and David Madigan. Interpretable classifiers using rules and bayesian analysis: Building a better stroke prediction model. *The Annals of Applied Statistics*, 9(3):1350–1371, 2015.
- Z.C. Lipton. The mythos of model interpretability. *arXiv:1606.03490*, 2016.
- Bing Liu, Wynne Hsu, Yiming Ma, et al. Integrating classification and association rule mining. In *Kdd*, volume 98, pages 80–86, 1998.
- 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.

- Nicolai Meinshausen. Node harvest. *The Annals of Applied Statistics*, pages 2049–2072, 2010.
- L. Mentch and G. Hooker. Quantifying uncertainty in random forests via confidence intervals and hypothesis tests. *Journal of Machine Learning Research*, 17:841–881, 2016.
- W.J. Murdoch, C. Singh, K. Kumbier, R. Abbasi-Asl, and B. Yu. Interpretable machine learning: Definitions, methods, and applications. *arXiv:1901.04592*, 2019.
- 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(1):25, 2007.
- B. Yu. Stability. *Bernoulli*, 19:1484–1500, 2013.
- R. Zhu, D. Zeng, and M. R. Kosorok. Reinforcement learning trees. *Journal of the American Statistical Association*, 110:1770–1784, 2015.

# Sparsity / interaction order

Interpretable without sparsity:

- Philosophical question?
- RF do not work well in the additive nonsparse context
- Therefore SIRUS does not work either

Sparsity

- SIRUS inherits sparsity properties of RF
- Sparsity has an impact on the required number of rules (which is automatically chosen based on an accuracy value).

Interaction order:

- RF can miss complex signals
- Sirius can detect only interactions of order two (see IRF and signed iterative RF)
- Can we adapt SIRUS to handle high-level interactions (they are masked by low-level interactions in the current version)? Modify the probability to encourage high-level interactions.

RF modification:

- Can we adapt RF to predict only when the output is larger than a threshold?
- Can we adapt the splitting criterion to focus on these regions (by adapting the pinball loss used for conditional quantiles)?