

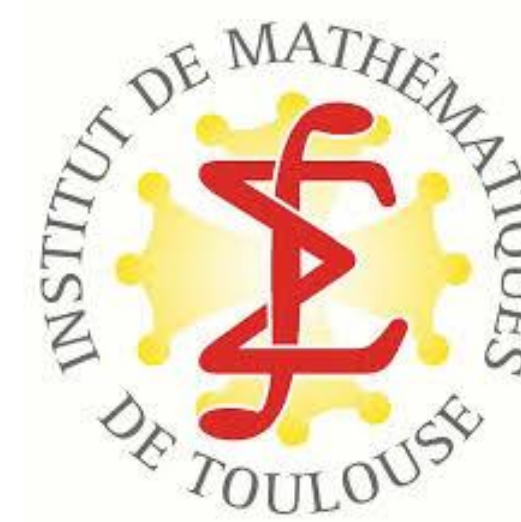
Variable Importance in High-Dimensional and Correlated Settings

Angel D Reyero Lobo^{1,2}, Pierre Neuval¹ and Bertrand Thirion²

Email: angel.reyero-lobo@inria.fr

¹ Institut de Mathématiques de Toulouse

² MIND, Inria Paris-Saclay



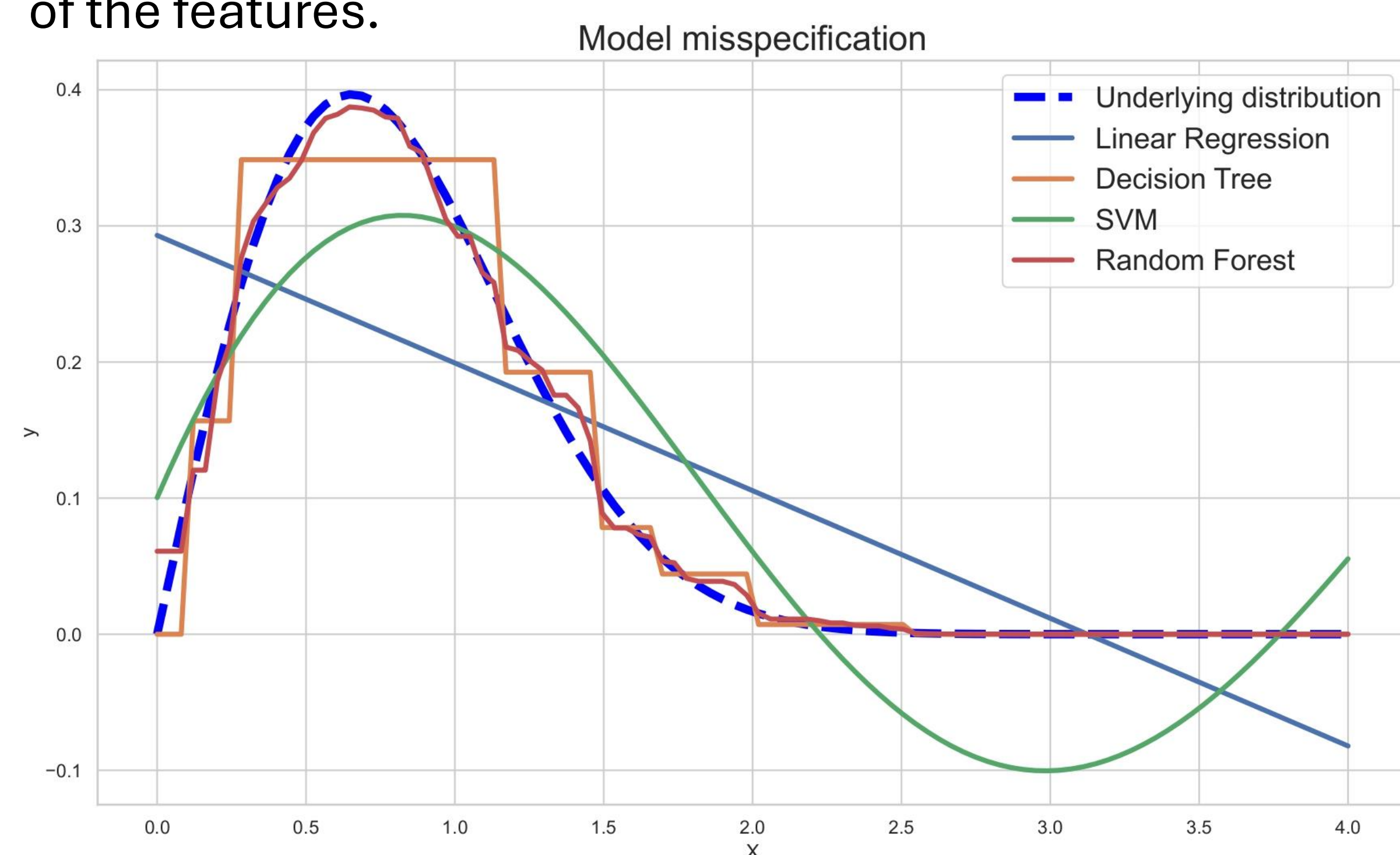
Introduction

- **Goal:** Evaluate the contribution of each feature in input X to the output variable Y .
- Their relationship can be complex, so we use a machine learning model to understand it:

$$\hat{m} \in \operatorname{argmin}_{f \in \mathcal{F}} \mathbb{E} [(f(X) - Y)^2].$$

Motivation

- Oversimplified models, while transparent, fail to capture the underlying distribution and, consequently, the intrinsic value of the features.



- Moreover, we cannot directly compare the importance of variables specific to each model between them!
- We need a **model-agnostic** approach, i.e. an interpretation method applicable to any regression technique.

Total Sobol Index

- From **Global Sensitivity Analysis**, the Sobol index is a model-agnostic statistical functional of interest that quantifies the difference in precision between the model excluding the group of covariates $s \subset \{1 \dots p\}$ and the complete model:

$$\psi_s(P) := \frac{\mathbb{E}_P [(Y - m_{P,s}(X^{-s}))^2]}{\operatorname{Var}_P(Y)} - \frac{\mathbb{E}_P [(Y - m_P(X))^2]}{\operatorname{Var}_P(Y)},$$

where $m_P(x) := \mathbb{E}_P [Y|X = x]$ and $m_{P,s}(x^{-s}) := \mathbb{E}_P [Y|X_{-s} = x^{-s}]$.

- See [Bénard et al. \(2022\)](#).

Estimation

- Naive idea: **plug-in** estimate.

$$\hat{\psi}_s(P) := \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} \frac{(y_i - \hat{m}_{P,s}(x_i^{-s}))^2}{\widehat{\operatorname{Var}}_P(Y)} - \frac{(y_i - \hat{m}_P(x_i))^2}{\widehat{\operatorname{Var}}_P(Y)}.$$

- It needs a bias correction ([Williamson et al. \(2021\)](#)).
- It is **computationally intensive**: we need to refit a $m_{P_0,s}$ for each group s !

- **Refitting** is expensive but **predicting** is cheap.

Selected references:

- Clément Bénard, Sébastien da Veiga, and Erwan Scornet. Mda for random forests: inconsistency, and a practical solution via the sobol-mda, 2022.
- Williamson BD, Gilbert PB, Carone M, Simon N. Nonparametric variable importance assessment using machine learning techniques. *Biometrics*. 2021; 77: 9–22.
- Chamma, A., Engemann, D. A., and Thirion, B. (2023). Statistically valid variable importance assessment through conditional permutations.

(Conditional) Permutation Feature Importance

- **Permutation Feature Importance (PFI):** It measures the predictive capability of the j -th covariate by breaking its relationship with the output.
- In practice, the j -th column of the test matrix is permuted:

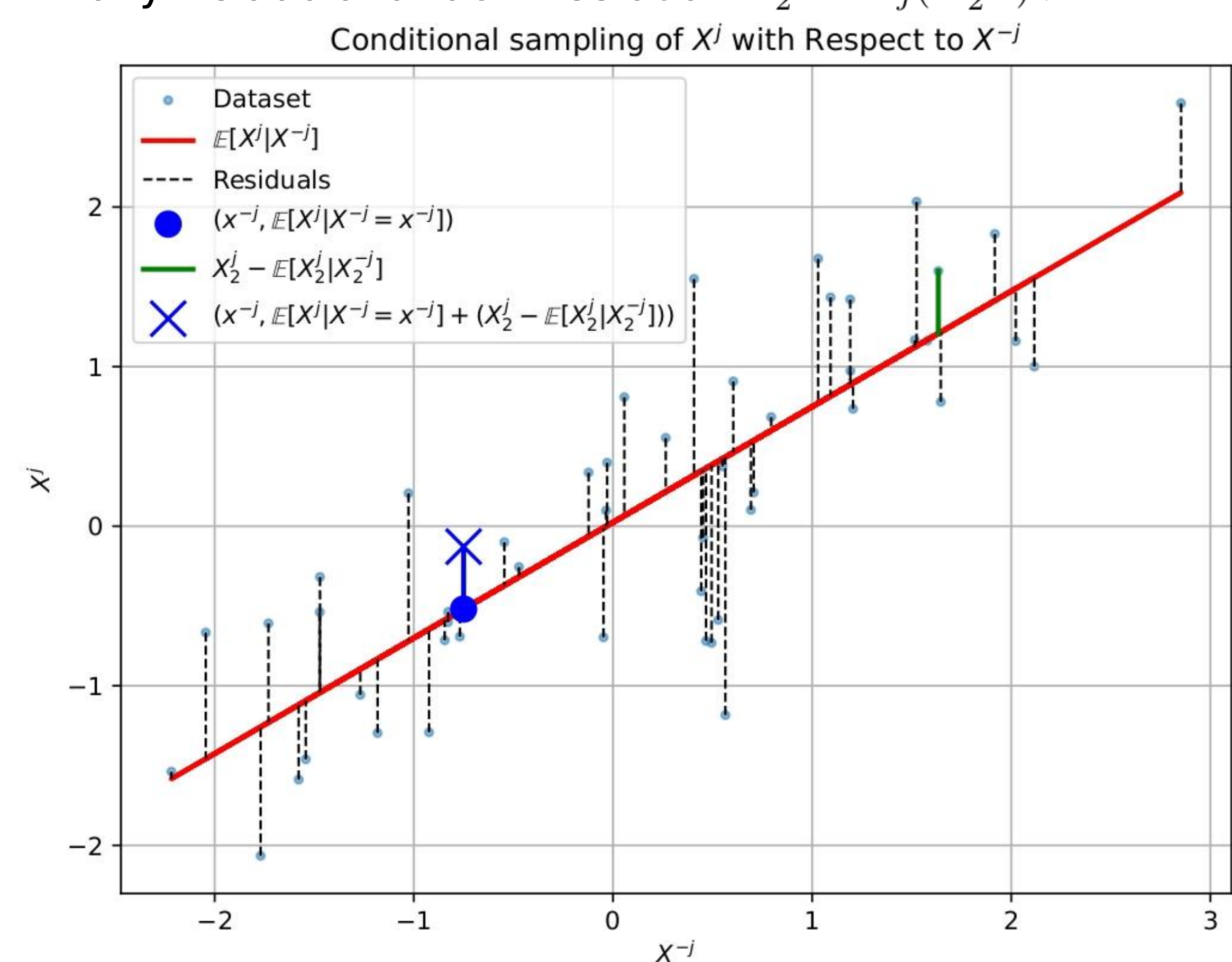
$$\hat{\chi}_j^{\text{PFI}} := \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} \left((y_i - \hat{m}(x_i^{(j)}))^2 - (y_i - \hat{m}(x_i))^2 \right).$$

- It also breaks the relationship with the rest of the input covariates, inducing **extrapolation bias**.
- **Conditional Permutation Importance (CPI):** introduced by [Chamma et al. \(NeurIPS 2023\)](#), the j -th covariate is *conditionally* permuted on the rest of the input covariates, breaking only its relationship with the output variable:

$$\hat{\chi}_j^{\text{CPI}} := \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} \left((y_i - \hat{m}(\tilde{x}_i^{(j)}))^2 - (y_i - \hat{m}(x_i))^2 \right).$$

An idea of conditional sampling/ permutation

- Under Gaussianity assumption, the conditional distribution $X^j|X^{-j} = x^{-j}$ is i.i.d. to $\tilde{X}^{(j)}|X^{-j} = x^{-j} \sim \nu_{-j}(x^{-j}) + (X_2^j - \nu_{-j}(X_2^{-j}))$ where $\nu_{-j}(x^{-j}) := \mathbb{E}[X^j|X^{-j} = x^{-j}]$ and $X_2 \stackrel{\text{iid}}{\sim} X$.
- In practice, we regress X^j given X^{-j} to obtain $\hat{\nu}_{-j}$. Then, to sample from $\tilde{X}^{(j)}|X^{-j} = x^{-j}$ we first predict $\hat{\nu}_{-j}(x^{-j})$ and finally we add a random residual $X_2^j - \hat{\nu}_{-j}(X_2^{-j})$.



Discussion

- **Aim:** a **computationally feasible**, **statistically valid** and **model-agnostic** variable importance measure that works in **highly-dimensional** and **correlated** settings.
- How can we relate the total Sobol index with CPI?
- How can we generalize the conditional sampling?
- What assumptions are necessary about the model estimate \hat{m} to provide information about the intrinsic variable importance?
- Does the total Sobol index capture all the interactions between the covariates, or do we need to move forward with Shapley values?