

A principled framework for comparing Variable Importance

Angel REYERO LOBO

Workshop: Methods for Explainable Machine Learning in Health Care
IMT & Inria Paris-Saclay

Joint work with:
Pierre NEUVIAL & Bertrand THIRION.

February 4, 2025

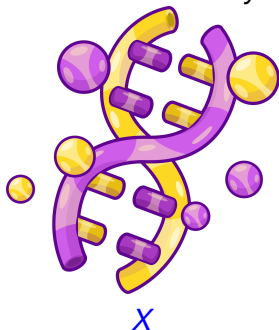


- 1 Introduction
 - Motivation
- 2 How to compare VIMs
 - State-of-the-art
 - Inconsistencies
 - General pipeline
- 3 Experiments
 - Simulated data
 - Real data
- 4 Conclusion
- 5 References

- 1 Introduction
 - Motivation
- 2 How to compare VIMs
 - State-of-the-art
 - Inconsistencies
 - General pipeline
- 3 Experiments
 - Simulated data
 - Real data
- 4 Conclusion
- 5 References

Variable Importance Measures (VIM)

How can we define / learn the importance of each covariate X^j with respect to an outcome y ?

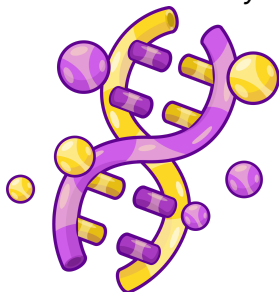


💡 Try to study their relationship using a ML model:

$$\hat{m} \in \operatorname{argmin}_{f \in \mathcal{F}} \hat{\mathbb{E}}[\mathcal{L}(f(X), y)]. \quad (1)$$

Variable Importance Measures (VIM)

How can we define / learn the importance of each covariate X^j with respect to an outcome y ?



X



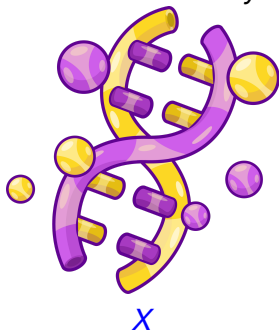
y

💡 Try to study their relationship using a ML model:

$$\hat{m} \in \operatorname{argmin}_{f \in \mathcal{F}} \hat{\mathbb{E}}[\mathcal{L}(f(X), y)]. \quad (1)$$

Variable Importance Measures (VIM)

How can we define / learn the importance of each covariate X^j with respect to an outcome y ?



💡 Try to study their relationship using a ML model:

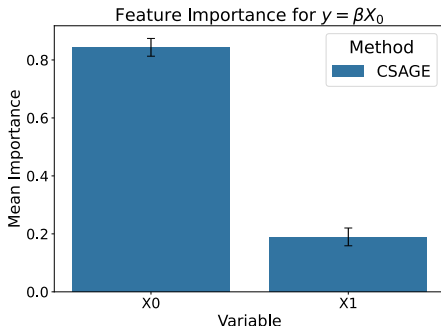
$$\hat{m} \in \operatorname{argmin}_{f \in \mathcal{F}} \hat{\mathbb{E}}[\mathcal{L}(f(X), y)]. \quad (1)$$

Assumption (Identifiability): X^j is not a function of X^{-j} .

What does it mean to be important?

“Feature importance as how much predictive power it provides to the model. We can then define “important” features as those whose absence degrades m ’s performance.”

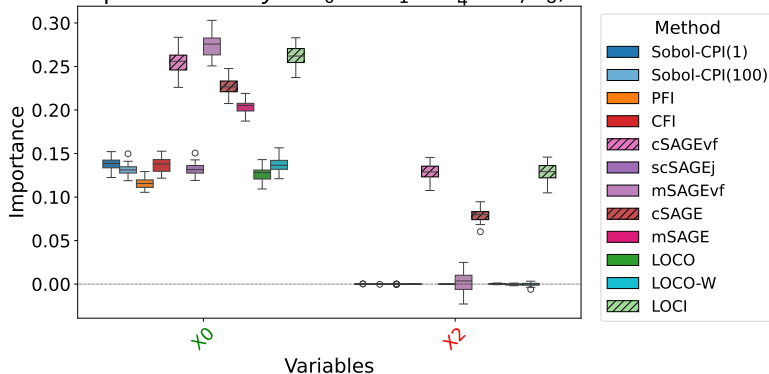
Covert et al. (2020) NeurIPS



⚠ Gap between variable importance and variable selection.

Motivation

Feature importance for $y = X_0 + 2X_1 - X_4^2 + X_7X_8$, $R^2 = 0.99$

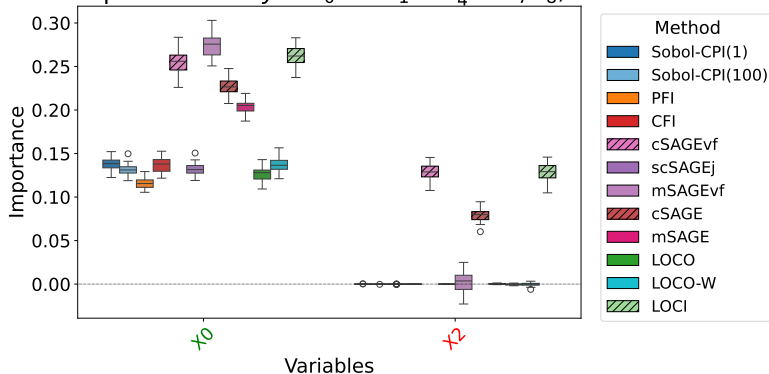


1 From X_0 : How to compare VIMs?

2 From X_2 : What's the minimum for a VIM?

Motivation

Feature importance for $y = X_0 + 2X_1 - X_4^2 + X_7X_8$, $R^2 = 0.99$



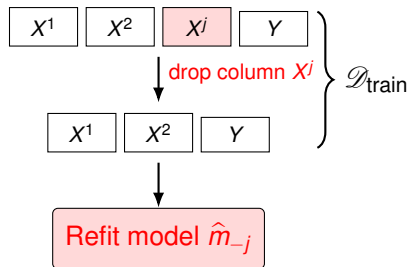
1 From X_0 : How to compare VIMs?

2 From X_2 : What's the minimum for a VIM?

- Minimal axiom: $\psi(j, P) = 0$ if and only if $X^j \perp\!\!\!\perp y \mid X^{-j}$.

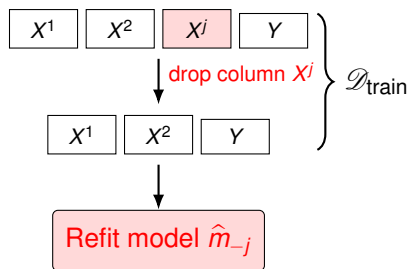
- 1 Introduction
 - Motivation
- 2 How to compare VIMs
 - State-of-the-art
 - Inconsistencies
 - General pipeline
- 3 Experiments
 - Simulated data
 - Real data
- 4 Conclusion
- 5 References

LOCO



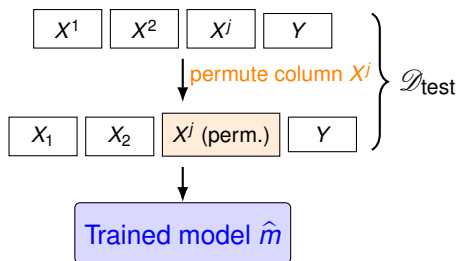
$$\Delta_{\text{LOCO}} = \text{perf}(\text{refit } \hat{m}_{-j}) - \text{perf}(\text{orig } \hat{m})$$

LOCO



$$\Delta_{\text{LOCO}} = \text{perf}(\text{refit } \hat{m}_{-j}) - \text{perf}(\text{orig } \hat{m})$$

PFI



$$\Delta_{\text{PFI}} = \text{perf}(\hat{m}(\text{perm } X^{(j)})) - \text{perf}(\hat{m}(\text{orig } X))$$

How can we compare VIMs?

- **Leave One Covariate Out(LOCO):**

$$\hat{\psi}_{\text{LOCO}}^j = \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} \mathcal{L}(y_i, \hat{m}_{-j}(\mathbf{x}_i^{-j})) - \mathcal{L}(y_i, \hat{m}(\mathbf{x}_i)).$$

- **Permutation Feature Importance(PFI):**

$$\hat{\psi}_{\text{PFI}}^j = \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} \mathcal{L}(y_i, \hat{m}(\mathbf{x}_i^{(j)})) - \mathcal{L}(y_i, \hat{m}(\mathbf{x}_i)).$$

where the j -th covariate is **permuted**.

➔ LOCO uses **refitting** and PFI uses **perturbation**.

How can we compare VIMs?

“LOCO differs from the other methods [...] since most of the other methods don’t require retraining the model. However, due to retraining the model, the interpretation shifts from only interpreting that one single model to interpreting the learner and how model training reacts to changes in the features.”

Molnar (2025), Interpretable Machine Learning, 3rd Edition

How can we compare VIMs?

“LOCO differs from the other methods [...] since most of the other methods don’t require retraining the model. However, due to retraining the model, the interpretation shifts from only interpreting that one single model to interpreting the learner and how model training reacts to changes in the features.”

Molnar (2025), Interpretable Machine Learning, 3rd Edition

Basically, compare VIMs based on the **inference** procedure used.

Inconsistencies in comparison by inference

The Total Sobol Index can be estimated in many different ways!

$$\psi_{\text{TSI}} := \mathbb{E} \left[\text{Var}(y \mid X^{-j}) \right]$$

Inconsistencies in comparison by inference

The Total Sobol Index can be estimated in many different ways!

$$\begin{aligned}\psi_{\text{TSI}} &:= \mathbb{E} \left[\text{Var}(y \mid X^{-j}) \right] \\ &= \mathbb{E} \left[\left(m_{-j}(X^{-j}) - y \right)^2 \right] - \mathbb{E} \left[(m(X) - y)^2 \right] && \text{refitting} \\ &= \frac{1}{2} \left(\mathbb{E} \left[\left(m(\tilde{X}^{(j)}) - y \right)^2 \right] - \mathbb{E} \left[(m(X) - y)^2 \right] \right) && \text{perturbation}\end{aligned}$$

Inconsistencies in comparison by inference

The Total Sobol Index can be estimated in many different ways!

$$\begin{aligned}\psi_{\text{TSI}} &:= \mathbb{E} \left[\text{Var}(y \mid X^{-j}) \right] \\&= \mathbb{E} \left[\left(m_{-j}(X^{-j}) - y \right)^2 \right] - \mathbb{E} \left[(m(X) - y)^2 \right] && \text{refitting} \\&= \frac{1}{2} \left(\mathbb{E} \left[\left(m(\tilde{X}^{(j)}) - y \right)^2 \right] - \mathbb{E} \left[(m(X) - y)^2 \right] \right) && \text{perturbation} \\&= \mathbb{E} \left[\left(\mathbb{E} \left[m(X) \mid X^{-j} \right] - y \right)^2 \right] - \mathbb{E} \left[(m(X) - y)^2 \right] && \text{marginalization} \\&= \mathbb{E} \left[(m_{-j}(X^{-j}) - m(X))^2 \right] && \text{variance} \\&= \sigma^2(R_{-j}^2 - R^2).\end{aligned}$$



1 Theoretical Index:

- Define goals and choose a matching theoretical quantity.
- Verify if it satisfies the minimal axiom.



1 Theoretical Index:

- Define goals and choose a matching theoretical quantity.
- Verify if it satisfies the minimal axiom.

2 Estimation:

- Select a procedure aligned with your desired inference properties:
 - ▶ E.g., double robustness, computational feasibility, extrapolation issues, benefit from unlabeled data, or simpler relationships between inputs than between inputs and outputs, ...



1 Theoretical Index:

- Define goals and choose a matching theoretical quantity.
- Verify if it satisfies the minimal axiom.

2 Estimation:

- Select a procedure aligned with your desired inference properties:
 - ▶ E.g., double robustness, computational feasibility, extrapolation issues, benefit from unlabeled data, or simpler relationships between inputs than between inputs and outputs, ...

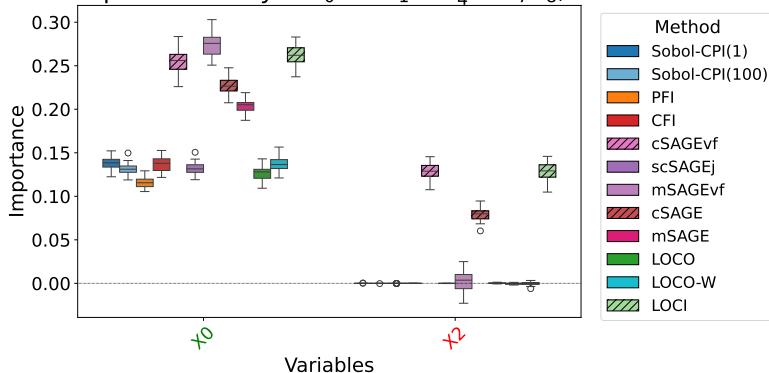
3 Type-I error:

- Provide statistical guarantees for the important covariates.

- 1 Introduction
 - Motivation
- 2 How to compare VIMs
 - State-of-the-art
 - Inconsistencies
 - General pipeline
- 3 Experiments
 - Simulated data
 - Real data
- 4 Conclusion
- 5 References

Simulated data

Feature importance for $y = X_0 + 2X_1 - X_4^2 + X_7X_8$, $R^2 = 0.99$



- From X_0 :
 - Sobol-CPI(1) \simeq Sobol-CPI(100) \simeq CFI \simeq scSAGEj \simeq LOCO(W).
 - cSAGEvf \simeq LOCI.
- From X_2 :
 - cSAGEvf, cSAGE and LOCI do not satisfy the minimal axiom!

Bike dataset

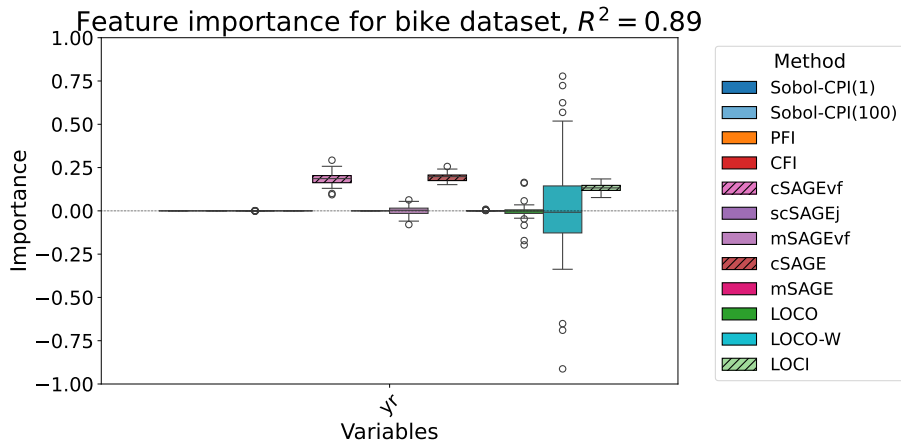


Figure 1: Boxplots of the VIMs for the feature *year*: Methods satisfying the minimal axiom assign no importance to this variable.

Bike dataset

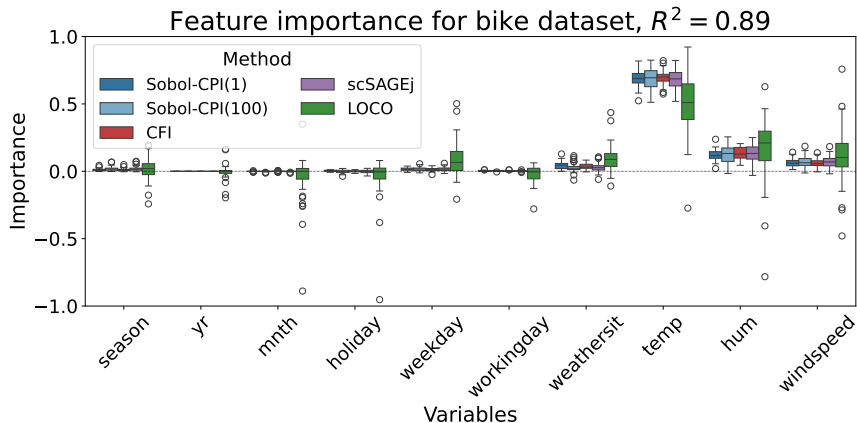


Figure 2: Boxplots of the VIMs estimating ψ_{TSI} for all features: Refitting approaches exhibit poorer inference properties.

Breast Cancer data with a correlated artificial null

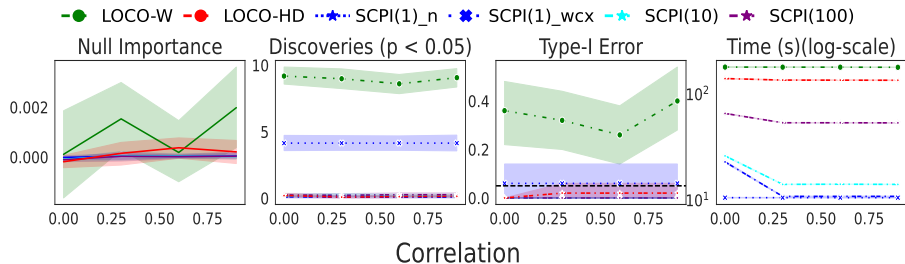


Figure 3: Double robustness and inference: Sobol-CPI assigns zero importance to the null. Sobol-CPI(1)-Wilcoxon makes discoveries while controlling error and staying efficient.

- 1 Introduction
 - Motivation
- 2 How to compare VIMs
 - State-of-the-art
 - Inconsistencies
 - General pipeline
- 3 Experiments
 - Simulated data
 - Real data
- 4 Conclusion
- 5 References

1 How to compare VIMs?



- 1 Conceptual comparison in the **theoretical index**.
 - 2 Inference comparison in the **estimation**.
 - 3 **Statistical guarantees** for the selected features.
- ✓ We provide a guide to help practitioners select a meaningful VIM.

2 What's the minimum for a VIM?

- **Minimal axiom:** $\psi(j, P) = 0$ if and only if $X^j \perp\!\!\!\perp y \mid X^{-j}$.
- ✓ Intuitive: Important if its absence degrades the model.
- ✓ Link between **variable selection** and **variable importance**!

- 1 Introduction
 - Motivation
- 2 How to compare VIMs
 - State-of-the-art
 - Inconsistencies
 - General pipeline
- 3 Experiments
 - Simulated data
 - Real data
- 4 Conclusion
- 5 References

References

- Ian Covert, Scott M Lundberg, and Su-In Lee. Understanding global feature contributions with additive importance measures. *Advances in Neural Information Processing Systems*, 33:17212–17223, 2020.
- Fiona Katharina Ewald, Ludwig Bothmann, Marvin N. Wright, Bernd Bischl, Giuseppe Casalicchio, and Gunnar König. *A Guide to Feature Importance Methods for Scientific Inference*, page 440–464. Springer Nature Switzerland, 2024. ISBN 9783031637971. doi: 10.1007/978-3-031-63797-1_22. URL http://dx.doi.org/10.1007/978-3-031-63797-1_22.
- Jing Lei, Max G'Sell, Alessandro Rinaldo, Ryan J. Tibshirani, and Larry Wasserman. Distribution-free predictive inference for regression. *Journal of the American Statistical Association*, 113(523): 1094–1111, 2018. doi: 10.1080/01621459.2017.1307116. URL <https://doi.org/10.1080/01621459.2017.1307116>.

Christoph Molnar. *Interpretable Machine Learning*. 3 edition, 2025.

ISBN 978-3-911578-03-5. URL [https:](https://christophm.github.io/interpretable-ml-book)

[//christophm.github.io/interpretable-ml-book](https://christophm.github.io/interpretable-ml-book).

Angel Reyero-Lobo, Pierre Neuvial, and Bertrand Thirion. Conditional feature importance revisited: Double robustness, efficiency and inference, 2025a. URL <https://arxiv.org/abs/2501.17520>.

Angel Reyero-Lobo, Pierre Neuvial, and Bertrand Thirion. A principled approach for comparing variable importance, 2025b. URL <https://arxiv.org/abs/2507.17306>.

Thank you — Questions?



hidimstat package



Article

Theoretical indices

Index	Definition	MA
ψ_{TSI}	$\mathbb{E} [\ell(m_{-j}(X^{-j}), y)] - \mathbb{E} [\ell(m(X), y)]$	Yes
ψ_{SAGE}	$\sum_{S \subset -\{j\}} w_S \left(\mathbb{E} [\ell(y, \mathbb{E} [m(X) \mid X^S])] - \mathbb{E} [\ell(y, \mathbb{E} [m(X) \mid X^{S \cup \{j\}}])] \right)$	No
ψ_{LOCI}	$\mathbb{E} [\ell(y, \mathbb{E} [y])] - \mathbb{E} [\ell(y, m_j(X^j))]$	No
ψ_{mSAGEvf}	$\mathbb{E} [\ell(y, \mathbb{E} [y])] - \mathbb{E} [\ell(y, \mathbb{E} [m(X^{(-j)})])]$	Yes
ψ_{mSAGE}	$\sum_{S \subset -\{j\}} w_S \left(\mathbb{E} [\ell(y, \mathbb{E} [m(X^{(-S)}) \mid X^S])] - \mathbb{E} [\ell(y, \mathbb{E} [m(X^{(-S)}) \mid X^{S \cup \{j\}}])] \right)$	Yes
ψ_{PFI}	$\mathbb{E} [\ell(m(X^{(j)}), y)] - \mathbb{E} [\ell(m(X), y)]$	Yes

Method	Theoretical quantity	Index	Estim
cSAGE	$\sum_{S \subset -\{j\}} w_S (v(S \cup \{j\}) - v(S))$	ψ_{SAGE}	M
cSAGEvf	$v(\{j\})$	ψ_{LOCI}	M
mSAGEvf	$v^m(\{j\})$	ψ_{mSAGEvf}	M
mSAGE	$\sum_{S \subset -\{j\}} w_S (v^m(S \cup \{j\}) - v^m(S))$	ψ_{mSAGE}	M
scSAGEvf	$v(-\{j\} \cup \{j\}) - v(-\{j\})$	ψ_{TSI}	M
LOCO	$\mathbb{E} [\ell(m_{-j}(X^{-j}), y)] - \mathbb{E} [\ell(m(X), y)]$	ψ_{TSI}	R
LOCO-W	$\mathbb{E} [\ell(m_{-j}(X^{-j}), y)] - \mathbb{E} [\ell(m(X), y)]$	ψ_{TSI}	R
LOCI	$\mathbb{E} [\ell(m_j(X^j), y)] - \mathbb{E} [\ell(m(X), y)]$	ψ_{LOCI}	R
PFI	$\mathbb{E} [\ell(m(X^{(j)}), y)] - \mathbb{E} [\ell(m(X), y)]$	ψ_{PFI}	P
CFI	$\mathbb{E} [\ell(m(\tilde{X}^{(j)}), y)] - \mathbb{E} [\ell(m(X), y)]$	ψ_{TSI}	P
Sobol-CPI(n-cal)	$\frac{n_{\text{cal}}}{n_{\text{cal}}+1} \left(\mathbb{E} \left[\ell \left(\frac{1}{n_{\text{cal}}} \sum_{k=1}^{n_{\text{cal}}} m(\tilde{X}_k^{(j)}), y \right) \right] - \mathbb{E} [\ell(m(X), y)] \right)$	ψ_{TSI}	P/M