

# Beyond Coefficients: Understanding Variable Importance in Modern ML

CoMeEcon  
Angel Reyero Lobo

Institut de Mathématiques de Toulouse & MIND, Inria Paris-Saclay  
Joint work with:  
Joseph Paillard

9th of June



*Inria*



# Contents

## 1 About me

## 2 Introduction

- Setting
- Explainable AI and Scientific Discovery
- Linear Models
- Random Forests
- Model-agnostic VIM
  - Permutation Feature Importance (PFI)
  - Conditional Feature Importance (CFI)
  - Leave One Covariate Out (LOCO)
- HiDimStat

## 3 Advanced Topics

- How to compare VIMs?
- Rashomon effect
- Variable importance for Conditional Independence Testing

## 4 Conclusion

## 5 References

# Index

## 1 About me

## 2 Introduction

- Setting
- Explainable AI and Scientific Discovery
- Linear Models
- Random Forests
- Model-agnostic VIM
  - Permutation Feature Importance (PFI)
  - Conditional Feature Importance (CFI)
  - Leave One Covariate Out (LOCO)
- HiDimStat

## 3 Advanced Topics

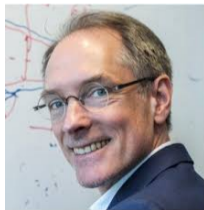
- How to compare VIMs?
- Rashomon effect
- Variable importance for Conditional Independence Testing

## 4 Conclusion

## 5 References

# About me

- I am a PhD student at the Institut de Mathématiques de Toulouse and in the MIND team at Inria Paris-Saclay
- My advisors are Pierre Neuvial (IMT) and Bertrand Thirion (Inria Paris-Saclay)
- I am an ELLIS PhD student, so I did a 6-month research visit: I came to Amsterdam to CWI to visit Peter Grünwald
- My research interests are: Interpretable Machine Learning, Conditional Independence Testing, e-values, and Missing Data



1 About me

2 Introduction

- Setting
- Explainable AI and Scientific Discovery
- Linear Models
- Random Forests
- Model-agnostic VIM
  - Permutation Feature Importance (PFI)
  - Conditional Feature Importance (CFI)
  - Leave One Covariate Out (LOCO)

● HiDimStat

3 Advanced Topics

- How to compare VIMs?
- Rashomon effect
- Variable importance for Conditional Independence Testing

4 Conclusion

5 References

# Setting: Supervised learning in tabular data

We have a dataset  $\{(x_i, y_i)\}_{i=1}^n$ , where  $x_i \in \mathbb{R}^p$  and  $y_i \in \mathbb{R}$ .

The goal is to learn a function

$$\hat{f}: \mathcal{X} \rightarrow \mathcal{Y}$$

such that

$$\hat{f}(x) \approx y$$

for unseen data points  $(x, y)$ .

To evaluate the predictive performance of a model  $f: \mathcal{X} \rightarrow \mathcal{Y}'$ , we use a loss function

$$\ell: \mathcal{Y} \times \mathcal{Y}' \rightarrow \mathbb{R}_+,$$

for example:

- Quadratic loss:  $\ell(y, \hat{y}) = (y - \hat{y})^2$
- Classification 0–1 loss:  $\ell(y, \hat{y}) = \mathbf{1}_{\{y \neq \hat{y}\}}$

## Two main goals of explainability:

- **Local Variable Importance**

- Explains a single prediction for an individual instance
  - ▶ e.g., Why did a credit scoring model reject Kayané's loan application?

## Two main goals of explainability:

- **Local Variable Importance**

- Explains a single prediction for an individual instance
  - ▶ e.g., Why did a credit scoring model reject Kayané's loan application?
- Often leads to *counterfactual explanations*:
  - ▶ What minimal changes would flip the decision?
  - ▶ e.g., Research does not pay enough, she should change the job.

## Two main goals of explainability:

### • Local Variable Importance

- Explains a single prediction for an individual instance
  - ▶ e.g., Why did a credit scoring model reject Kayané's loan application?
- Often leads to *counterfactual explanations*:
  - ▶ What minimal changes would flip the decision?
  - ▶ e.g., Research does not pay enough, she should change the job.

### • Global Variable Importance

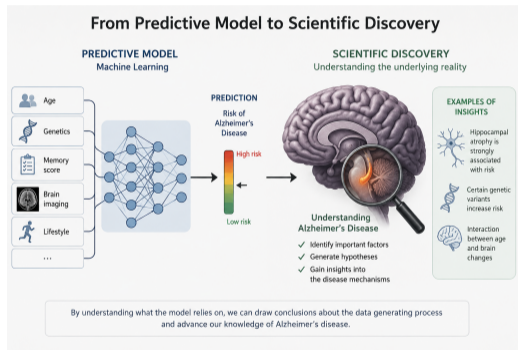
- Explains how features influence the model overall
- e.g., Is income generally important for credit approval? what about eye color?

**Focus of this talk:** Global Variable Importance

Günther et al. (2025) pessimistic about Local Variable Importance for complex models.

## Two main motivations for studying **Global Variable Importance**:

- **Explainable AI (xAI):**
  - understand the model itself, i.e., identify the features on which the model relies.
- **Scientific Discovery:**
  - use a predictive model to understand reality. Understanding which variables drive predictions can provide insights into the underlying data-generating process.



# Linear Models

**Linear Model assumption:**  $Y = \sum_{j=1}^p \beta_j X^j + \varepsilon$ , where  $\varepsilon$  is centered independent noise.

We estimate the relationship using

$$\hat{f}(X) = \sum_{j=1}^p \hat{\beta}_j X^j.$$

# Linear Models

**Linear Model assumption:**  $Y = \sum_{j=1}^p \beta_j X^j + \varepsilon$ , where  $\varepsilon$  is centered independent noise.

We estimate the relationship using

$$\hat{f}(X) = \sum_{j=1}^p \hat{\beta}_j X^j.$$

Interpretability comes naturally:

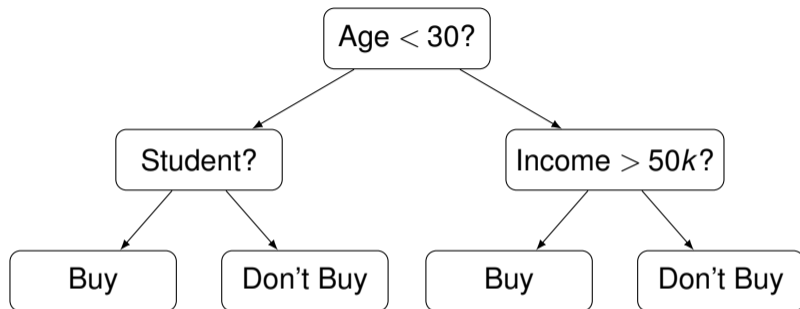
$$|\hat{\beta}_j|$$

measures the influence of feature  $X^j$  on the prediction.

The same idea extends to:

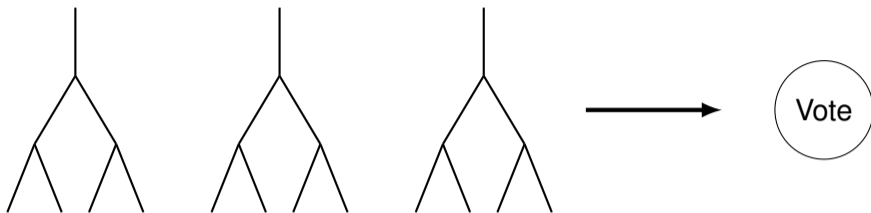
- Generalized Linear Models (e.g. Logistic Regression),
- Regularized linear models (Ridge, Lasso).

**Limitation:** linear models struggle to capture nonlinear interactions.



- Easy to visualize and explain,
- Naturally capture interactions and nonlinearities,
- But often unstable and less predictive than ensemble methods.

# Random Forests



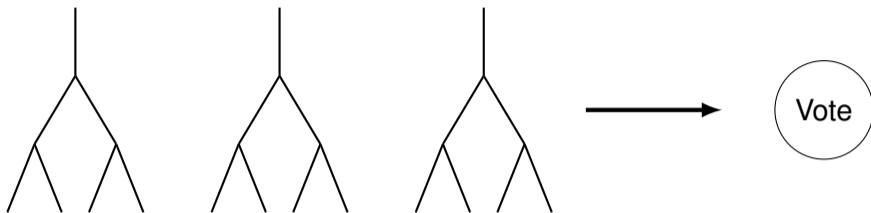
A random forest aggregates many decision trees trained on random subsets of:

- observations (bootstrap sampling),
- features (feature subsampling).

Advantages:

- Strong predictive performance,
- Captures complex interactions,
- More stable than a single tree.

# Random Forests



A random forest aggregates many decision trees trained on random subsets of:

- observations (bootstrap sampling),
- features (feature subsampling).

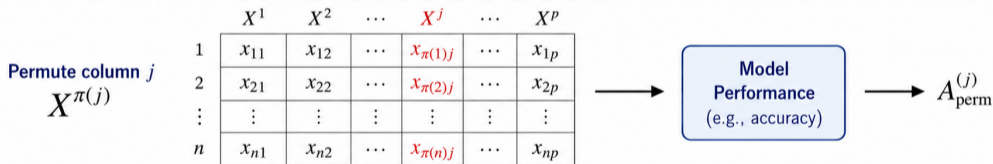
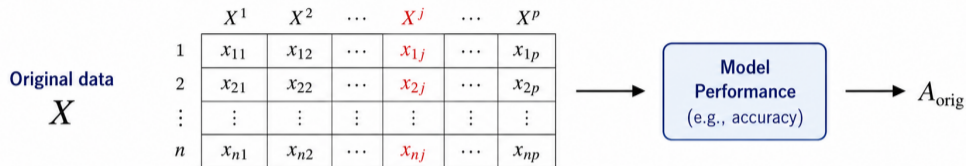
Advantages:

- Strong predictive performance,
- Captures complex interactions,
- More stable than a single tree.

? How can we assign an importance to  $X^j$  in the RF model?

# Mean Decrease Accuracy (MDA)

## Mean Decrease Accuracy (MDA)



$$\text{MDA}_j = A_{\text{orig}} - A_{\text{perm}}^{(j)}$$

If permuting feature  $j$  decreases the performance (i.e.,  $A_{\text{perm}}^{(j)} < A_{\text{orig}}$ ), then  $X^j$  is **important**.

# Mean Decrease Accuracy (MDA)

**Idea:** a feature is important if shuffling its values decreases the model accuracy.

For a feature  $X^j$ :

- 1 Compute the prediction error on the test set,

$$\text{Err}_{\text{orig}} = \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} \ell(\hat{f}(x_i), y_i).$$

- 2 Randomly permute the values of the feature  $X^j$  among individuals:  $X^{\pi(j)}$ .
- 3 Recompute the prediction error,

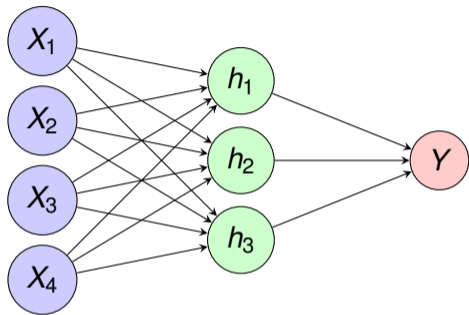
$$\text{Err}_{\text{perm}}^{(j)} = \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} \ell(\hat{f}(x_i^{\pi(j)}), y_i).$$

- 4 Define the importance score:  $\text{MDA}_j = \text{Err}_{\text{perm}}^{(j)} - \text{Err}_{\text{orig}}$ .

Large increase in error  $\Rightarrow$  important feature.

# Model-agnostic Variable Importance

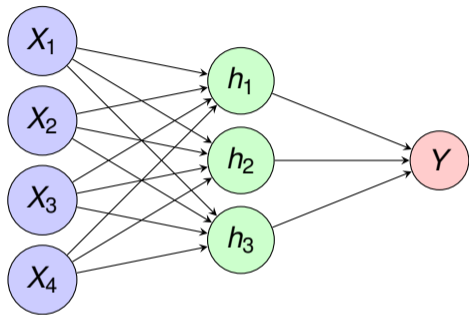
Modern machine learning models include:  
Gradient Boosting, Neural Networks,  
SuperLearners, foundation models, etc.



**Question:** how do we interpret such complex black-box models?

# Model-agnostic Variable Importance

Modern machine learning models include:  
Gradient Boosting, Neural Networks,  
SuperLearners, foundation models, etc.



**Question:** how do we interpret such complex black-box models?

A key idea is to generalize Mean Decrease Accuracy (MDA) to *any* predictor  $\hat{f}$ .

This leads to **Permutation Feature Importance (PFI)**.

# Permutation Feature Importance (PFI)

PFI measures the increase in prediction loss when a feature is randomly permuted.

Let  $\mathcal{D}_{\text{test}} = \{(x_i, y_i)\}_{i=1}^{n_{\text{test}}}$ .

Define the importance of feature  $j$  as:

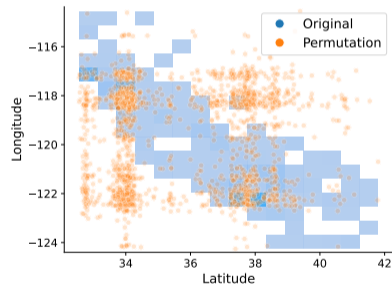
$$\text{PFI}_j = \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} \ell(\hat{f}(x_i^{\pi(j)}), y_i) - \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} \ell(\hat{f}(x_i), y_i).$$

where  $x_i^{\pi(j)}$  is obtained by replacing the  $j$ -th feature value of  $x_i$  with the value of feature  $j$  from another randomly chosen observation.

# Permutation Feature Importance (PFI): Limitations

Although PFI is model-agnostic and widely used, it has several issues:

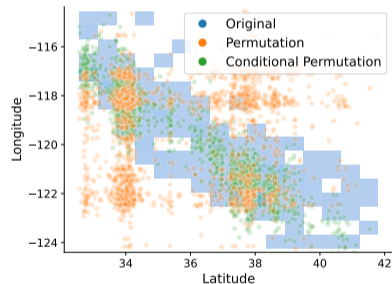
- **Theoretical limitations:** PFI does not correspond to an interpretable estimand (Bénard et al. (2022)).
- **Extrapolation issues:** permuting features may create unrealistic samples outside the data distribution (Strobl et al. (2008); Hooker et al. (2021)).



# Permutation Feature Importance (PFI): Limitations

Although PFI is model-agnostic and widely used, it has several issues:

- **Theoretical limitations:** PFI does not correspond to an interpretable estimand (Bénard et al. (2022)).
- **Extrapolation issues:** permuting features may create unrealistic samples outside the data distribution (Strobl et al. (2008); Hooker et al. (2021)).



**Key idea:** permute while keeping realistic observations.

## PFI

$$P_{(X^{\pi(j)}, Y)} = P_{(Y, X^{-j})} P_{X^j}$$

- × Breaks  $X^j - Y$
- × Breaks  $X^j - X^{-j}$
- × Unrealistic samples

## PFI

$$P_{(X^{\pi(j)}, Y)} = P_{(Y, X^{-j})} P_{X^j}$$

- × Breaks  $X^j - Y$
- × Breaks  $X^j - X^{-j}$
- × Unrealistic samples

## CFI

$$P_{(\tilde{X}^{(j)}, Y)} = P_{(Y, X^{-j})} P_{X^j | X^{-j}}$$

- ✓ Preserves  $X^j - X^{-j}$
- ✓ Realistic samples
- ✓ Breaks  $X^j - Y | X^{-j}$

Key idea: condition on the remaining features  $X^{-j}$  before resampling  $X^j$ .

# Conditional Feature Importance (CFI)

CFI measures the increase in prediction loss when a feature is *conditionally* permuted.

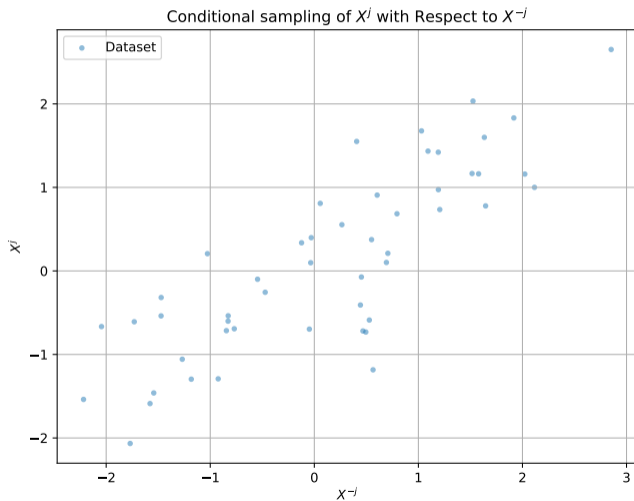
Let  $\mathcal{D}_{\text{test}} = \{(x_i, y_i)\}_{i=1}^{n_{\text{test}}}$ .

Define the importance of feature  $j$  as:

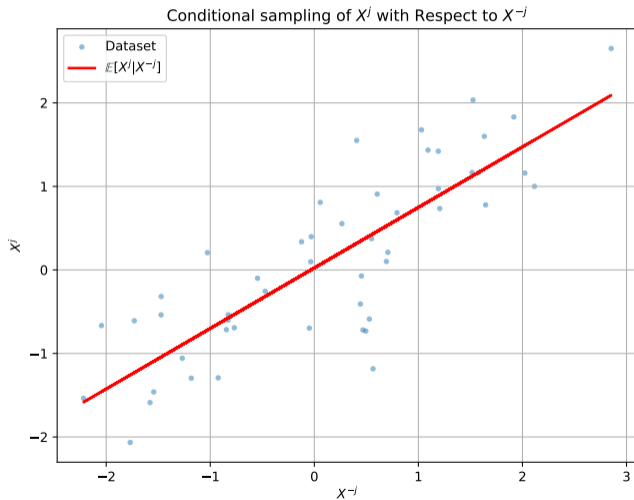
$$\psi_{\text{CFI}}(j) = \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} \ell(\hat{f}(\tilde{x}_i^{(j)}), y_i) - \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} \ell(\hat{f}(x_i), y_i).$$

where  $\tilde{x}_i^{(j)}$  replaces the  $j$ -th feature value of  $x_i$  with a new sample from  $\mathcal{L}(X^j | X^{-j})$ .

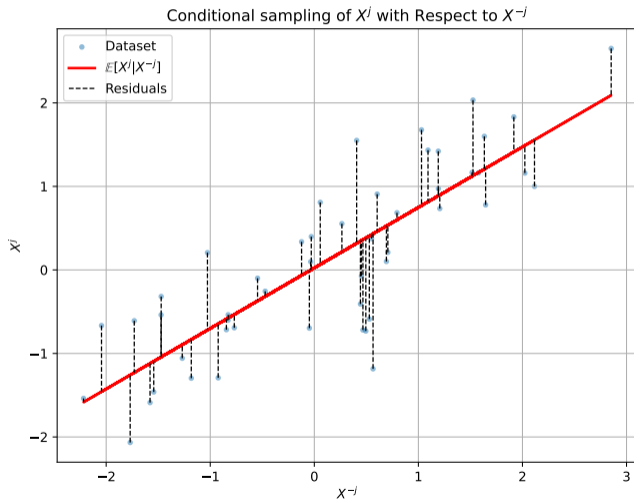
# Conditional Permutation



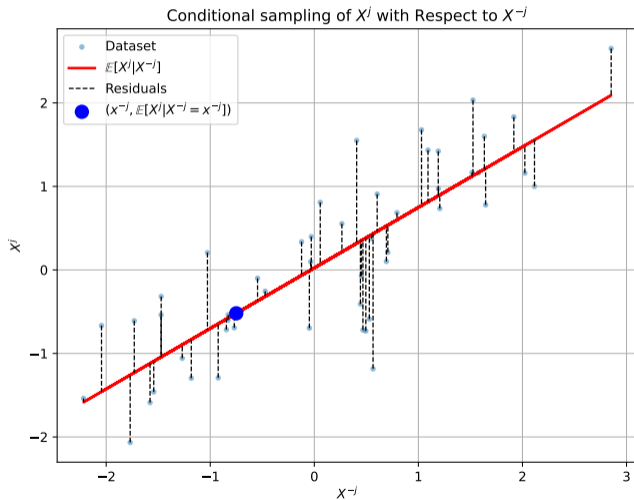
# Conditional Permutation



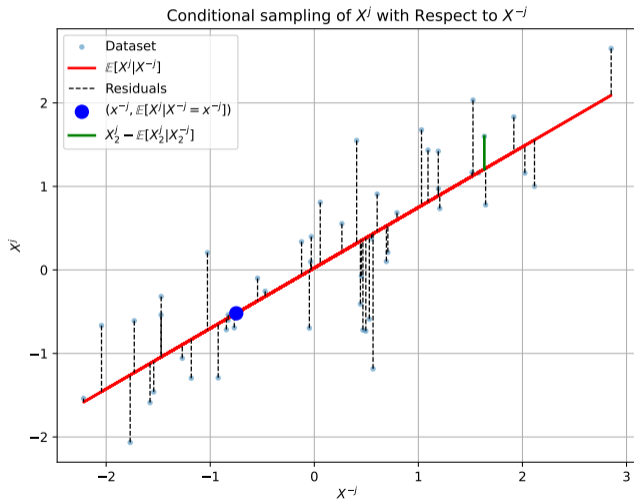
# Conditional Permutation



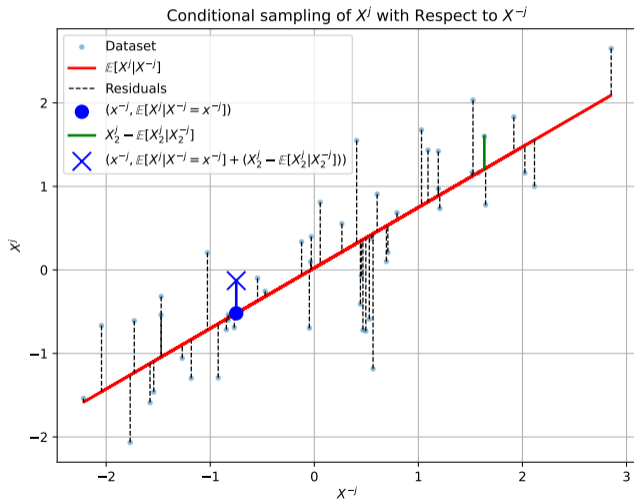
# Conditional Permutation



# Conditional Permutation



# Conditional Permutation



## Step 1: Learn the conditional mean

$$\hat{v}_j(X^{-j}) \approx \mathbb{E}[X^j | X^{-j}]$$



## Step 2: Permute residuals

$$\hat{\varepsilon}_j = X^j - \hat{v}_j(X^{-j}), \quad \tilde{X}^{(j)} = \hat{v}_j(X^{-j}) + \hat{\varepsilon}_j^\pi$$



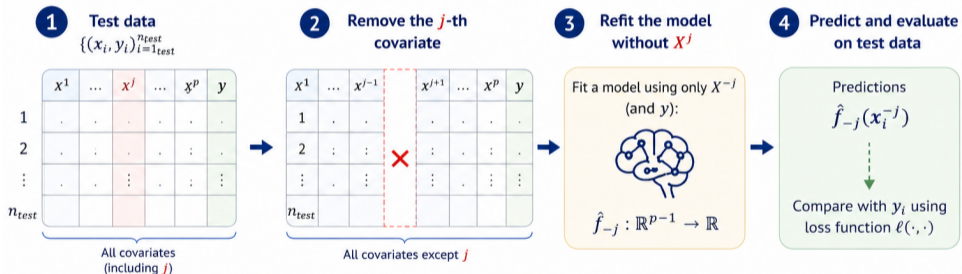
## Step 3: Measure loss increase

$$\psi_{\text{CFI}}(j) = \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} \ell(\hat{f}(\tilde{x}_i^{(j)}), y_i) - \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} \ell(\hat{f}(x_i), y_i).$$

 **Permute residuals, not features.**

## Leave One Covariate Out (LOCO)

LOCO measures the importance of a feature by evaluating a model that has been refit **without** that feature.



LOCO importance of covariate  $j$

$$\psi_{\text{LOCO}}(j) = \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} \left[ \ell(\hat{f}_{-j}(x_i^{-j}), y_i) - \ell(\hat{f}(x_i), y_i) \right]$$

- $x_i^{-j}$ : all covariates except  $j$  for observation  $i$
- $\hat{f}_{-j}$ : model refit without  $X^j$
- $\hat{f}$ : full model using all covariates
- $\ell(\cdot, \cdot)$ : loss function

We measure feature importance by comparing model performance with and without  $X^j$ :

$$\psi_{\text{LOCO}}(j) = \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} \left[ \ell(\hat{f}_{-j}(\mathbf{x}_i^{-j}), y_i) - \ell(\hat{f}(\mathbf{x}_i), y_i) \right]$$

- It is computationally expensive, since we must retrain a model for each feature.
- It converges to the total Sobol' index, a standard estimand in sensitivity analysis.
- Due to optimization error, it is more unstable than CFI (Reyero-Lobo et al. (2025b)).



## HiDimStat: High-dimensional statistical inference tool for Python

The HiDimStat package provides statistical inference methods to solve the problem of support recovery in the context of high-dimensional and spatially structured data.

### Installation

HiDimStat working only with Python 3, ideally Python 3.10+. For installation, run the following from terminal:

```
pip install hidimstat
```

Or if you want the latest version available (for example to contribute to the development of this project):

```
git clone https://github.com/mind-inria/hidimstat.git
cd hidimstat
pip install -e .
```

#### ☰ On this page

[Installation](#)

[Dependencies](#)

[Documentation & Examples](#)

[Build the documentation](#)

[References](#)

[References](#)

[📄 Show Source](#)

# Index

1

About me

2

Introduction

- Setting
- Explainable AI and Scientific Discovery
- Linear Models
- Random Forests
- Model-agnostic VIM
  - Permutation Feature Importance (PFI)
  - Conditional Feature Importance (CFI)
  - Leave One Covariate Out (LOCO)
- HiDimStat

3

Advanced Topics

- How to compare VIMs?
- Rashomon effect
- Variable importance for Conditional Independence Testing

4

Conclusion

5

References

# How Do We Compare Variable Importance Measures?

We have encountered several notions of importance:

- Linear model coefficients,
- Permutation-based importance (PFI/CFI),
- Leave-One-Covariate-Out importance (LOCO).

A natural question arises:

How can we compare these different notions of importance?

# How Do We Compare Variable Importance Measures?

We have encountered several notions of importance:

- Linear model coefficients,
- Permutation-based importance (PFI/CFI),
- Leave-One-Covariate-Out importance (LOCO).

A natural question arises:

How can we compare these different notions of importance?

Some authors argue that LOCO and PFI/CFI are fundamentally different because

- LOCO compares *two predictive models*,
- PFI/CFI evaluates a *single model* under perturbed inputs.

# How Do We Compare Variable Importance Measures?

We have encountered several notions of importance:

- Linear model coefficients,
- Permutation-based importance (PFI/CFI),
- Leave-One-Covariate-Out importance (LOCO).

A natural question arises:

How can we compare these different notions of importance?

Some authors argue that LOCO and PFI/CFI are fundamentally different because

- LOCO compares *two predictive models*,
- PFI/CFI evaluates a *single model* under perturbed inputs.

**Exercise.** Assume

$$Y = f_*(X) + \varepsilon, \quad \varepsilon \perp\!\!\!\perp X, \quad \mathbb{E}[\varepsilon] = 0, \quad \mathbb{E}[\varepsilon^2] < \infty.$$

Show that under squared loss,

$$\boxed{\psi_{\text{CFI}}(j) = 2 \psi_{\text{LOCO}}(j)}.$$

# CFI and $2\times$ LOCO are Equivalent

Under squared loss,

$$\psi_{\text{CFI}}(j) = \mathbb{E} \left[ (f_{\star}(\tilde{X}^{(j)}) - Y)^2 - (f_{\star}(X) - Y)^2 \right].$$

Using

$$Y = f_{\star}(X) + \varepsilon, \quad \mathbb{E}[\varepsilon] = 0,$$

we obtain

$$\psi_{\text{CFI}}(j) = \mathbb{E} \left[ (f_{\star}(\tilde{X}^{(j)}) - f_{\star}(X))^2 \right].$$

Since  $\tilde{X}^j \stackrel{d}{=} X^j$ ,  $\tilde{X}^j \perp\!\!\!\perp X^j \mid X^{-j}$ , the variables  $f_{\star}(\tilde{X}^{(j)})$ ,  $f_{\star}(X)$  are i.i.d. given  $X^{-j}$ . Therefore, using Tower's property, we obtain that

$$\psi_{\text{CFI}}(j) = 2 \mathbb{E} \left[ \text{Var}(f_{\star}(X) \mid X^{-j}) \right]$$

Recall

$$\psi_{\text{LOCO}}(j) = \mathbb{E} \left[ (f_{-j,\star}(X^{-j}) - Y)^2 - (f_{\star}(X) - Y)^2 \right].$$

Since

$$f_{-j,\star}(X^{-j}) = \mathbb{E} [f_{\star}(X) \mid X^{-j}],$$

the additive independent centered noise yields

$$\psi_{\text{LOCO}}(j) = \mathbb{E} \left[ (f_{-j,\star}(X^{-j}) - f_{\star}(X))^2 \right] = \mathbb{E} \left[ \text{Var}(f_{\star}(X) \mid X^{-j}) \right].$$

Hence,

$$\boxed{\psi_{\text{CFI}}(j) = 2 \psi_{\text{LOCO}}(j)}.$$

# How Should We Compare VIMs?

## Exercise.

Assume the linear model  $Y = \sum_{k=1}^p \beta_k X^k + \varepsilon$ . Show that

$$\psi_{\text{LOCO}}(j) = \beta_j^2 \mathbb{E} \left[ \text{Var}(X^j \mid X^{-j}) \right].$$

# How Should We Compare VIMs?

## Exercise.

Assume the linear model  $Y = \sum_{k=1}^p \beta_k X^k + \varepsilon$ . Show that

$$\psi_{\text{LOCO}}(j) = \beta_j^2 \mathbb{E} \left[ \text{Var}(X^j \mid X^{-j}) \right].$$

Indeed,

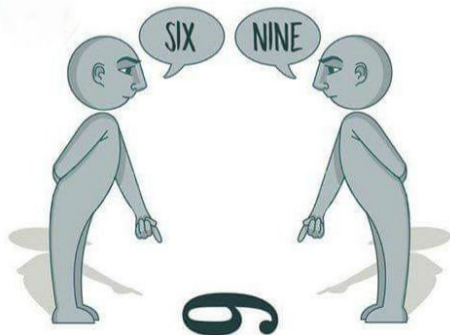
$$\psi_{\text{LOCO}}(j) = \mathbb{E} \left[ \text{Var}(f_{\star}(X) \mid X^{-j}) \right] = \mathbb{E} \left[ \text{Var} \left( \sum_{k=1}^p \beta_k X^k \mid X^{-j} \right) \right] = \beta_j^2 \mathbb{E} \left[ \text{Var}(X^j \mid X^{-j}) \right].$$

- Stronger correlations imply smaller LOCO importance.
- Regression coefficients do *not* account for redundancy among predictors.
- Different VIMs quantify fundamentally different notions of importance.

For a general framework to compare VIMs, see Reyero-Lobo et al. (2025a).

# The Rashomon Effect

- Many different models can explain the same dataset **equally well**.
- However, these models may rely on **different patterns or features** in the data.
- This phenomenon is known as the **Rashomon effect**.



# Aggregate Models or Aggregate Importances?

For each trained model  $\{\hat{f}_b\}_{b=1}^B$ , we obtain an importance estimate

$$\hat{\psi}_b := \psi(\hat{f}_b).$$

This raises an important question:

- **Aggregate the importances**

$$\hat{\psi} = \frac{1}{B} \sum_{b=1}^B \hat{\psi}_b$$

- **Aggregate the models first, then compute the importance**

$$\hat{\psi} = \psi\left(\frac{1}{B} \sum_{b=1}^B \hat{f}_b\right)$$

*Should we aggregate the models or the importance measures?*

# Aggregate Models, Not Explanations

We proved that the excess risk satisfies

$$\frac{1}{n} \sum \ell(\hat{f}(x_i), y_i) - \mathbb{E}[\ell(f_*(X), Y)] = \varepsilon + O_P(n^{-1/2}) = \mathbb{E}[\|\hat{f} - f_*\|^2] + O_P(n^{-1/2}).$$

- The term  $O_P(n^{-1/2})$  corresponds to the **test/estimation error**.
- The dominant term  $\varepsilon$  comes from the **training procedure**.
- The importance bias comes mainly from model bias!  
We need model-agnostic measures; otherwise, we can't explain the distribution.

**Key message:** Most of the error in ML models comes from **training the model**, not from test estimation.

# Aggregating Before vs After

We compare two strategies:

- **Aggregate models first (Ensemble)**

$$\psi_n^{\text{ens}} = \psi(\hat{f}_{\text{ens}})$$

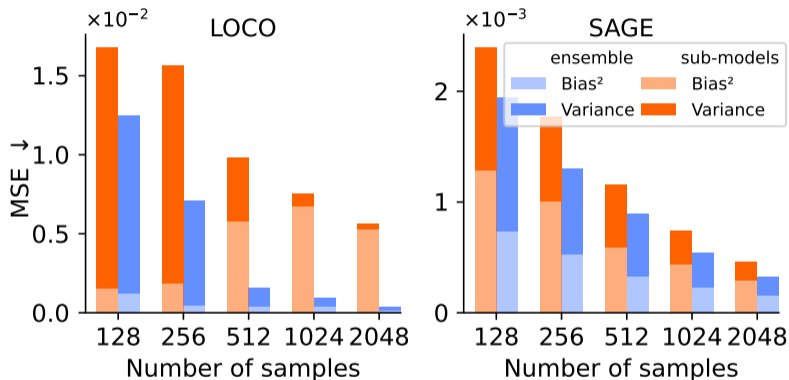
- **Aggregate explanations after training (Sub-models)**

$$\psi_n^{\text{sub}} = \frac{1}{B} \sum_{b=1}^B \psi(\hat{f}_b)$$

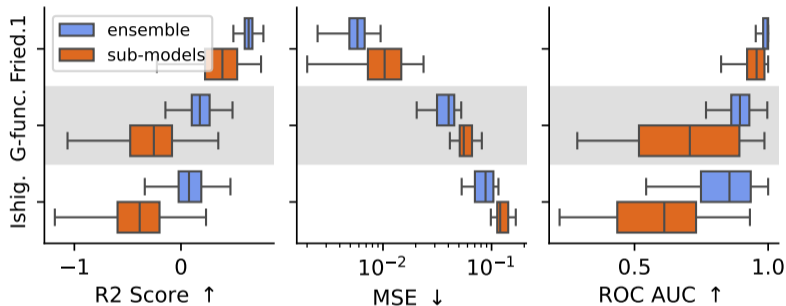
We obtain the relation

$$\psi_n^{\text{ens}} - \psi_{\star} = (\psi_n^{\text{sub}} - \psi_{\star}) \left( \rho + \frac{1-\rho}{B} \right) + O_P(n^{-1/2})$$

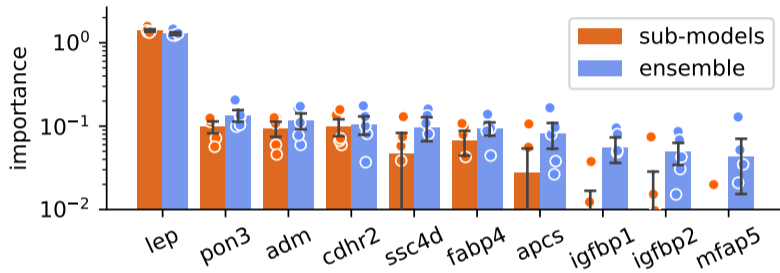
**Conclusion:** Aggregating models first (ensemble) leads to a smaller error.



*Aggregation of models (ensemble) vs aggregation of explanations (sub-models).*



*Simulation results comparing aggregation of models (ensemble) with aggregation of explanations(sub-models).*



*UKBiobank comparing aggregation of models (ensemble) with aggregation of explanations (sub-models).*

# Statistical Guarantees for Scientific Discovery

We often use machine learning models for **scientific discovery**:

Which variables truly influence the outcome?

Marginal independence,  $X^j \perp\!\!\!\perp Y$ , is not enough:

- Variables spuriously correlated with important ones are considered important.

# Statistical Guarantees for Scientific Discovery

We often use machine learning models for **scientific discovery**:

Which variables truly influence the outcome?

Marginal independence,  $X^j \perp\!\!\!\perp Y$ , is not enough:

- Variables spuriously correlated with important ones are considered important.

Instead, we test **Conditional Independence Testing (CIT)**

$$H_0: X^j \perp\!\!\!\perp Y \mid X^{-j}.$$

- $X^j$  is important if it provides information about  $Y$  that cannot be recovered from  $X^{-j}$ .
- If features are too correlated, consider **grouping** them.

To make scientific claims, we need **statistical guarantees** (e.g. type-I or FDR control).

Many examples to do so in HiDimStat:

<https://hidimstat.github.io/stable/generated/gallery/examples>

# The Model-X Principle

Most modern CIT procedures rely on the **Model-X assumption**:

$$\tilde{X}^{(j)} \sim \mathcal{L}(X^j \mid X^{-j}).$$

We can sample copies of  $X^j$  while preserving its dependence with the rest.

Most modern CIT procedures rely on the **Model-X assumption**:

$$\tilde{X}^{(j)} \sim \mathcal{L}(X^j | X^{-j}).$$

We can sample copies of  $X^j$  while preserving its dependence with the rest.

Given an importance statistic  $T(X^j, Y, X^{-j})$ , we compare it with  $T(\tilde{X}^{(j)}, Y, X^{-j})$ .

Under the null hypothesis ( $H_0: X^j \perp\!\!\!\perp Y | X^{-j}$ ), the two quantities are exchangeable:

$$T(X^j, Y, X^{-j}) \stackrel{\mathcal{L}}{=} T(\tilde{X}^{(j)}, Y, X^{-j}).$$

This simple observation is the foundation of

- Conditional Randomization Tests (CRT),
- Model-X Knockoffs.

# What should the test statistic $T(X^j, Y, X^{-j})$ be?

The power of the procedure depends on the choice of  $T(X^j, Y, X^{-j})$ .

Traditionally, one often uses the absolute value of a regression coefficient:

$$T = |\hat{\beta}_j|.$$

# What should the test statistic $T(X^j, Y, X^{-j})$ be?

The power of the procedure depends on the choice of  $T(X^j, Y, X^{-j})$ .

Traditionally, one often uses the absolute value of a regression coefficient:

$$T = |\hat{\beta}_j|.$$

But isn't this just an  
Interpretability / Variable Importance measure?

For linear models,  $T = |\hat{\beta}_j|$  is a natural importance score.

## **Question:**

How can we construct powerful CIT using importance measures from modern ML?

Many questions arise when using interpretable ML for conditional independence testing:

- How can we test whether  $\psi_{\text{CFI}}(j) = 0$ ?
- How can we make the Conditional Randomization Test model-agnostic?
- How can we make knockoffs model-agnostic?
- How can we obtain sequentially valid procedures?

# Conditional Independence Testing with Interpretable ML

Many questions arise when using interpretable ML for conditional independence testing:

- How can we test whether  $\psi_{\text{CFI}}(j) = 0$ ?
  - The CLT is not valid due to the vanishing variance under the null.
  - We can use variance corrections or nonparametric tests (Reyero-Lobo et al. (2025b)).
- How can we make the Conditional Randomization Test model-agnostic?
  
- How can we make knockoffs model-agnostic?
  
- How can we obtain sequentially valid procedures?

# Conditional Independence Testing with Interpretable ML

Many questions arise when using interpretable ML for conditional independence testing:

- How can we test whether  $\psi_{\text{CFI}}(j) = 0$ ?
  - The CLT is not valid due to the vanishing variance under the null.
  - We can use variance corrections or nonparametric tests (Reyero-Lobo et al. (2025b)).
- How can we make the Conditional Randomization Test model-agnostic?
  - The Holdout Randomization Test (HRT) is model-agnostic but requires a train–test split.
  - Semi-knockoffs (Reyero-Lobo et al. (2026)) avoid the need for a train–test split.
- How can we make knockoffs model-agnostic?
  
- How can we obtain sequentially valid procedures?

# Conditional Independence Testing with Interpretable ML

Many questions arise when using interpretable ML for conditional independence testing:

- How can we test whether  $\psi_{\text{CFI}}(j) = 0$ ?
  - The CLT is not valid due to the vanishing variance under the null.
  - We can use variance corrections or nonparametric tests (Reyero-Lobo et al. (2025b)).
- How can we make the Conditional Randomization Test model-agnostic?
  - The Holdout Randomization Test (HRT) is model-agnostic but requires a train–test split.
  - Semi-knockoffs (Reyero-Lobo et al. (2026)) avoid the need for a train–test split.
- How can we make knockoffs model-agnostic?
  - **Semi-knockoffs** can control the FDR through the knockoff threshold.
  - They do not need conditional densities but conditional expectations.
- How can we obtain sequentially valid procedures?

Many questions arise when using interpretable ML for conditional independence testing:

- How can we test whether  $\psi_{\text{CFI}}(j) = 0$ ?
  - The CLT is not valid due to the vanishing variance under the null.
  - We can use variance corrections or nonparametric tests (Reyero-Lobo et al. (2025b)).
- How can we make the Conditional Randomization Test model-agnostic?
  - The Holdout Randomization Test (HRT) is model-agnostic but requires a train–test split.
  - Semi-knockoffs (Reyero-Lobo et al. (2026)) avoid the need for a train–test split.
- How can we make knockoffs model-agnostic?
  - **Semi-knockoffs** can control the FDR through the knockoff threshold.
  - They do not need conditional densities but conditional expectations.
- How can we obtain sequentially valid procedures?
  - We can use **e-values**.
  - Joint work with: Sebastian Arias, Michele Meziu, and Peter Grünwald.

# A Quick Introduction to the E-World

**E-value.** A random variable  $S$  satisfying  $S \geq 0$ ,  $\mathbb{E}_P[S] \leq 1 \quad \forall P \in \mathcal{H}_0$ .

We can obtain type-I error control: by Markov's inequality,

$$\mathbb{P}_P(S \geq 1/\alpha) \leq \frac{\mathbb{E}_P[S]}{1/\alpha} \leq \alpha,$$

so we reject  $H_0$  whenever  $S \geq 1/\alpha$ .

# A Quick Introduction to the E-World

**E-value.** A random variable  $S$  satisfying  $S \geq 0$ ,  $\mathbb{E}_P[S] \leq 1 \quad \forall P \in \mathcal{H}_0$ .

We can obtain type-I error control: by Markov's inequality,

$$\mathbb{P}_P(S \geq 1/\alpha) \leq \frac{\mathbb{E}_P[S]}{1/\alpha} \leq \alpha,$$

so we reject  $H_0$  whenever  $S \geq 1/\alpha$ .

**E-process.** A sequence  $(S_t)_{t \geq 1}$  such that

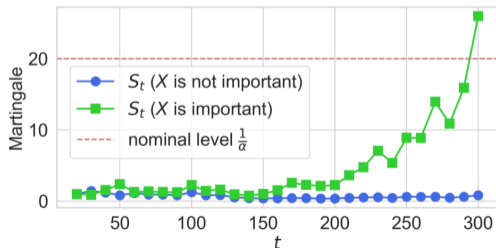
$$S_t \geq 0, \quad \mathbb{E}_P[S_\tau] \leq 1$$

for every stopping time  $\tau$  and every  $P \in \mathcal{H}_0$ .

By Ville's inequality,

$$\mathbb{P}_P\left(\sup_{t \geq 1} S_t \geq \frac{1}{\alpha}\right) \leq \alpha.$$

Reject as soon as  $S_t \geq 1/\alpha$ .



# Sequential Conditional Independence Testing with ML

**Goal:** Test  $Y \perp\!\!\!\perp X \mid Z$  assuming Model-X ( $P_{X|Z}$  is known and  $\tilde{X} \sim P_{X|Z}$ ).

**Existing approaches:**

- **Model-X e-value** Grünwald et al. (2024): based on likelihood ratio

$$\frac{f_{Y|X,Z}}{f_{Y|Z}}.$$

# Sequential Conditional Independence Testing with ML

**Goal:** Test  $Y \perp\!\!\!\perp X \mid Z$  assuming Model-X ( $P_{X|Z}$  is known and  $\tilde{X} \sim P_{X|Z}$ ).

**Existing approaches:**

- **Model-X e-value** Grünwald et al. (2024): based on likelihood ratio

$$\frac{f_{Y|X,Z}}{f_{Y|Z}}.$$

- **e-CRT** Shaer et al. (2023): based on importance statistics

$$q = T(X, Y, Z), \quad \tilde{q} = T(\tilde{X}, Y, Z),$$

and an antisymmetric score  $g(\tilde{q}, q) = -g(q, \tilde{q})$  giving an e-value  $1 + g(q, \tilde{q})$ .

# Sequential Conditional Independence Testing with ML

**Goal:** Test  $Y \perp\!\!\!\perp X \mid Z$  assuming Model-X ( $P_{X|Z}$  is known and  $\tilde{X} \sim P_{X|Z}$ ).

**Existing approaches:**

- **Model-X e-value** Grünwald et al. (2024): based on likelihood ratio

$$\frac{f_{Y|X,Z}}{f_{Y|Z}}.$$

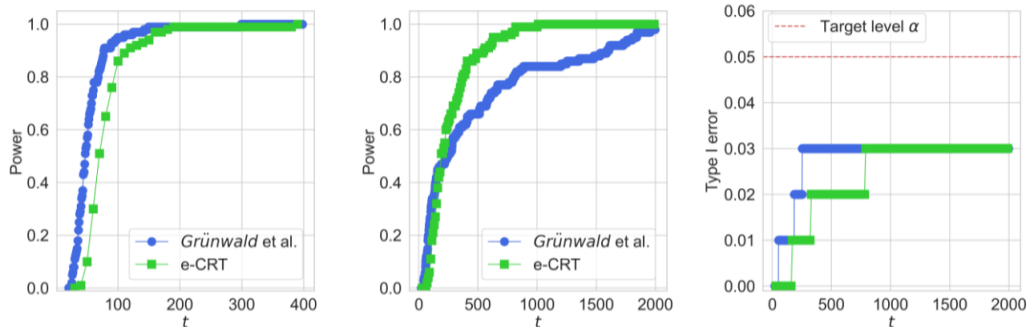
- **e-CRT** Shaer et al. (2023): based on importance statistics

$$q = T(X, Y, Z), \quad \tilde{q} = T(\tilde{X}, Y, Z),$$

and an antisymmetric score  $g(\tilde{q}, q) = -g(q, \tilde{q})$  giving an e-value  $1 + g(q, \tilde{q})$ .

**Limitations:** Model-X is growth-rate optimal under correct specification, but can be suboptimal under misspecification. e-CRT is flexible but leaves both  $T$  and  $g$  unspecified.

# e-CRT vs Model-X



**Figure 1: From Shaer et al. (2023):** In the well-specified setting (left), Model-X is more powerful. However, in the misspecified setting (center), where the conditional density cannot be learned exactly, e-CRT can achieve higher power.

## Our proposal:

$$\frac{2\hat{q}(t, \tilde{t})}{\hat{q}(t, \tilde{t}) + \hat{q}(\tilde{t}, t)},$$

where

$$t = \ell(\hat{m}(X, Z), Y), \quad \tilde{t} = \ell(\hat{m}(\tilde{X}, Z), Y),$$

and  $\hat{q}$  is an estimate of the joint density of  $(t, \tilde{t})$ .

- We estimate a 2D density, not a conditional density  $f_{Y|X,Z}$ .
- We can derandomize it:

$$\mathbb{E}_{\tilde{X}|X,Y,Z} \left[ \frac{2\hat{q}(t, \tilde{t})}{\hat{q}(t, \tilde{t}) + \hat{q}(\tilde{t}, t)} \right].$$

**Joint work with:** Sebastian Arias, Michele Meziu, Peter Grünwald.

# Index

1

About me

2

Introduction

- Setting
  - Explainable AI and Scientific Discovery
  - Linear Models
  - Random Forests
  - Model-agnostic VIM
    - Permutation Feature Importance (PFI)
    - Conditional Feature Importance (CFI)
    - Leave One Covariate Out (LOCO)
  - HiDimStat
- 3
- Advanced Topics
- How to compare VIMs?
  - Rashomon effect
  - Variable importance for Conditional Independence Testing

4

Conclusion

5

References

- Machine learning has become a powerful tool for **scientific discovery**:
  - It can capture complex patterns that may not have been identified previously.
- Modern models are highly complex, motivating **model-agnostic** interpretability.
- The general pipeline of IML should be:
  - 1 define what we mean by *importance*, which determines the **target quantity** of interest
  - 2 compare different **estimators** (e.g., LOCO and CFI) and their properties, e.g. stability, computational cost, robustness to misspecification.
  - 3 derive rigorous **statistical guarantees**.

Conditional independence testing uses interpretability to construct powerful tests.

Everything is already implemented in <https://hidimstat.github.io>!

# Index

1

About me

2

Introduction

- Setting
- Explainable AI and Scientific Discovery
- Linear Models
- Random Forests
- Model-agnostic VIM
  - Permutation Feature Importance (PFI)
  - Conditional Feature Importance (CFI)
  - Leave One Covariate Out (LOCO)
- HiDimStat

3

Advanced Topics

- How to compare VIMs?
- Rashomon effect
- Variable importance for Conditional Independence Testing

4

Conclusion

5

References

- Clément Bénéard, Sébastien Da Veiga, and Erwan Scornet. Mean decrease accuracy for random forests: inconsistency, and a practical solution via the sobol-mda. *Biometrika*, 109(4):881–900, 02 2022.
- Peter Grünwald, Alexander Henzi, and Tyron Lardy. Anytime-valid tests of conditional independence under model-x. *Journal of the American Statistical Association*, 119(546):1554–1565, 2024. doi: 10.1080/01621459.2023.2205607. URL <https://doi.org/10.1080/01621459.2023.2205607>.
- Eric Günther, Balázs Szabados, Robi Bhattacharjee, Sebastian Bordt, and Ulrike von Luxburg. Informative post-hoc explanations only exist for simple functions, 2025. URL <https://arxiv.org/abs/2508.11441>.

- Giles Hooker, Lucas Mentch, and Siyu Zhou. Unrestricted permutation forces extrapolation: variable importance requires at least one more model, or there is no free variable importance. *Statistics and Computing*, 31(82):1–16, 2021. doi: 10.1007/s11222-021-10057-z. URL <https://doi.org/10.1007/s11222-021-10057-z>.
- Joseph Paillard, Angel Reyero Lobo, Denis A. Engemann, and Bertrand Thirion. Aggregate models, not explanations: Improving feature importance estimation. In *International Conference on Machine Learning (ICML)*, 2026.
- Angel Reyero-Lobo, Pierre Neuvial, and Bertrand Thirion. A principled approach for comparing variable importance, 2025a. URL <https://arxiv.org/abs/2507.17306>.
- Angel Reyero-Lobo, Pierre Neuvial, and Bertrand Thirion. Conditional feature importance revisited: Double robustness, efficiency and inference, 2025b. URL <https://arxiv.org/abs/2501.17520>.

- Angel Reyero-Lobo, Bertrand Thirion, and Pierre Neuvial. Semi-knockoffs: a model-agnostic conditional independence testing method with finite-sample guarantees. In *International Conference on Machine Learning (ICML)*, 2026.
- Shalev Shaer, Gal Maman, and Yaniv Romano. Model-x sequential testing for conditional independence via testing by betting. In Francisco Ruiz, Jennifer Dy, and Jan-Willem van de Meent, editors, *Proceedings of The 26th International Conference on Artificial Intelligence and Statistics*, volume 206 of *Proceedings of Machine Learning Research*, pages 2054–2086. PMLR, 25–27 Apr 2023. URL <https://proceedings.mlr.press/v206/shaer23a.html>.
- Carolin Strobl, Anne-Laure Boulesteix, Thomas Kneib, Thomas Augustin, and Achim Zeileis. Conditional variable importance for random forests. *BMC Bioinformatics*, 9(1): 307, 2008. ISSN 1471-2105. doi: 10.1186/1471-2105-9-307. URL <https://doi.org/10.1186/1471-2105-9-307>.

# Thank you!

Questions?

angel.reyero-lobo@inria.fr  
joseph.paillard@roche.com