

# Probabilistic Explanations for Regression Models

---

Frederic Koriche, Jean-Marie Lagniez, and **Chi Tran**

**uai2025** — Rio de Janeiro, Brazil — 21-25 July, 2025



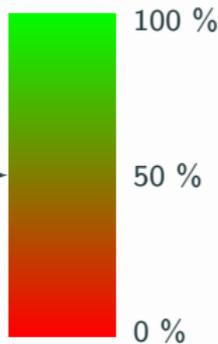
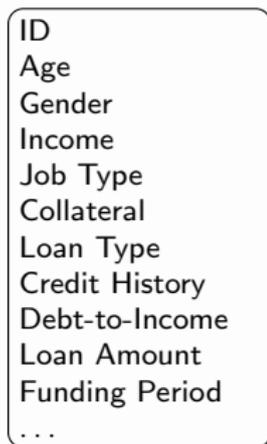
Univ. Artois, Centre de Recherche en Informatique de Lens (CRIL)  
AI Chair EXPEKCTATION of the French Agency of National Research

- 1 Regression Models
- 2 Probabilistic Explanations
- 3 Dealing with PP-Hardness
- 4 Dealing with NP-Hardness
- 5 Experiments



Conceptually, a **regression model** is a mapping  $f$  from a set  $\mathcal{X} \subseteq \mathbb{R}^d$  of *data instances* to a set  $\mathcal{Y} \subseteq \mathbb{R}$  of *outcomes*.

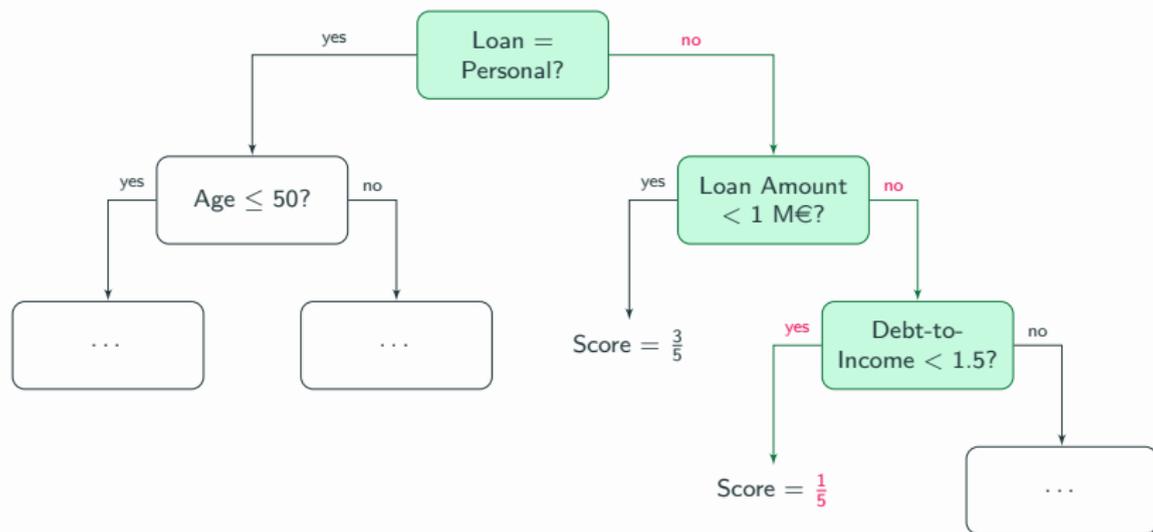
# Regression Models: Illustration



## Loan Eligibility

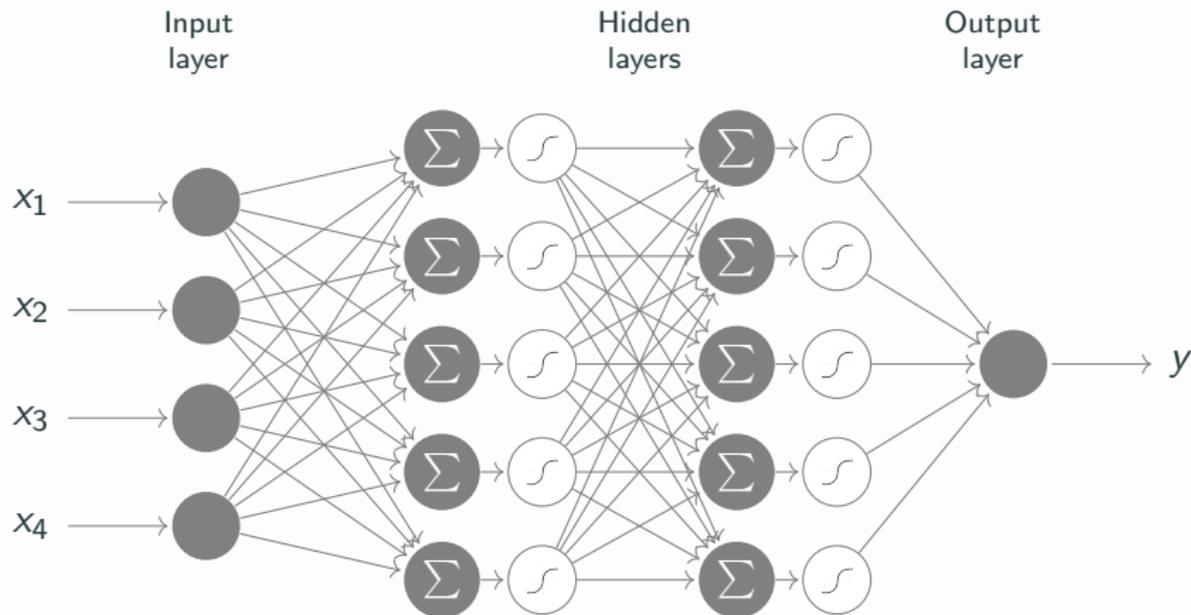
Each instance is a list of attributes about an applicant. The objective is to predict the applicant's likelihood of repaying the loan.

# Regression Models: Interpretability



Although some regression models are interpretable ...

# Regression Models: ~~Interpretability~~



Most of them are **not!**

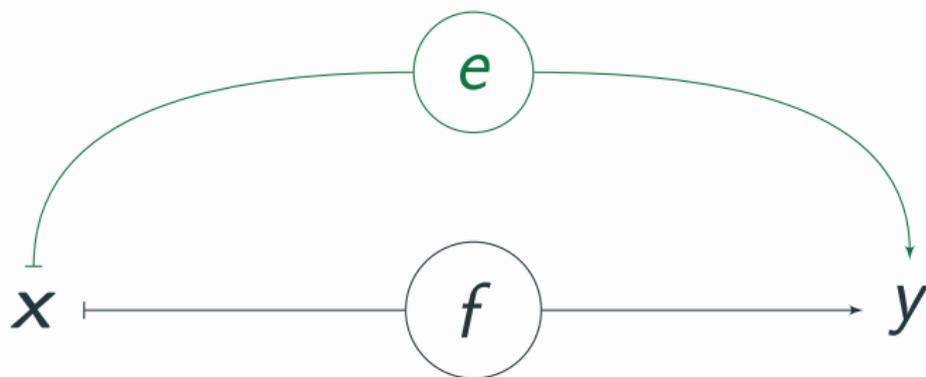
# Explaining Regressors



Why  $y$  is the outcome of  $x$ ?

Thus, a key issue is to provide answers to **why-questions**.

## Explaining Regressors



An **explanation** for an instance  $x$  with respect to a prediction model  $f$  is an interpretable surrogate model  $e$  that is consistent with  $f$  at  $x$ .

## Black-Box Model

Any regression model  $f$  is viewed as a **black box**, with only access to the outcome of any queried instance.

## Black-Box Model

Any regression model  $f$  is viewed as a **black box**, with only access to the outcome of any queried instance.

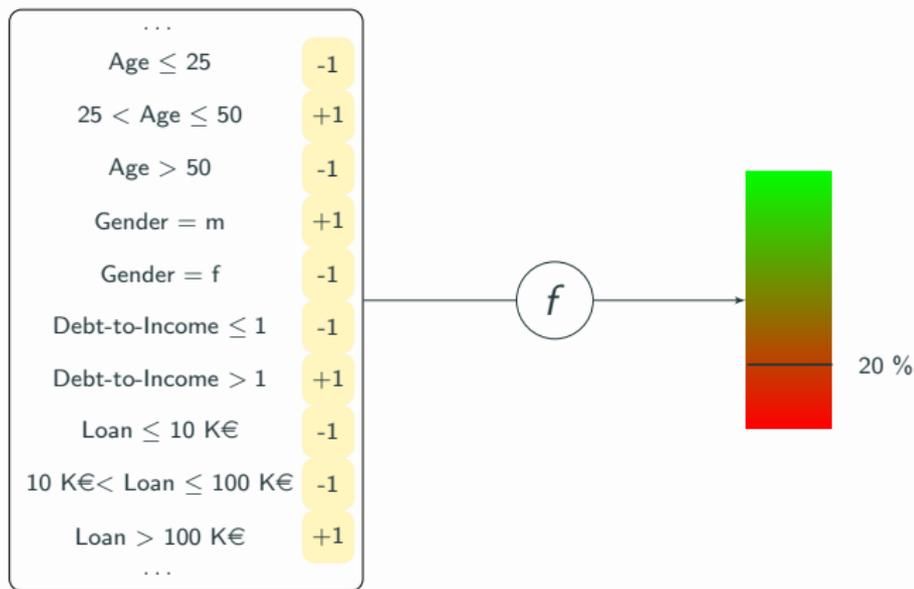
## Main Assumption [Ribeiro et al., 2016]

The input space  $\mathcal{X}$  is a Boolean hypercube, where each dimension is **interpretable**.

In other words, each instance is a vector  $\mathbf{x} \in \{\pm 1\}^d$ , where  $[d] = \{1, \dots, d\}$  is a set of interpretable features\*.

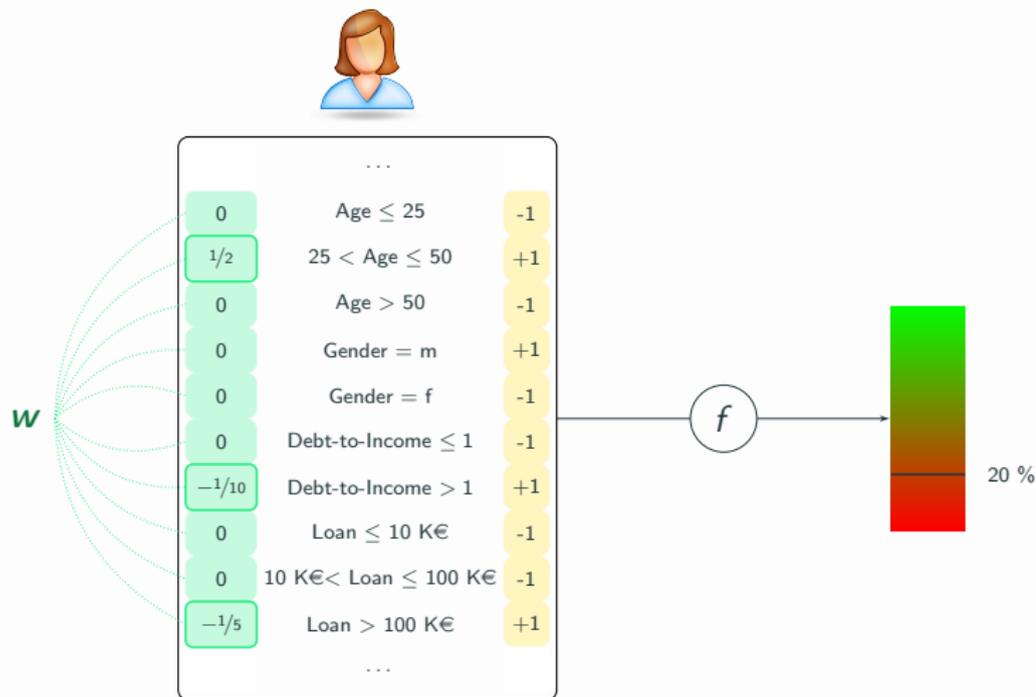
\*Data instances over raw attributes are first transformed into Boolean vectors, using discretization techniques or interpretable latent spaces.

# Explaining Regressors: Illustration



Using a set of interpretable features, let us explain why our applicant is not eligible for a loan ...

# Explaining Regressors: Illustration



A clear way to explain the outcome of an instance is to identify a **sparse** vector of weights  $w$ .

## Explaining Regressors: Illustration



$$\frac{1}{2} [25 < \text{Age} \leq 50] - \frac{1}{10} [\text{DTI} > 1] - \frac{1}{5} [\text{Loan} > 100 \text{ K€}] \longrightarrow \text{Score} = \frac{1}{5}$$

This vector  $\mathbf{w}$  can be seen as a **weighted decision rule**, where the head is determined by the sum of activated weights.

The explanation model  $e$  is therefore given by:

$$e(\mathbf{z}) = \mathbf{w} \cdot \mathbf{z} \quad \text{for all } \mathbf{z} \in \{\pm 1\}^d$$

- 1 Regression Models
- 2 Probabilistic Explanations**
- 3 Dealing with PP-Hardness
- 4 Dealing with NP-Hardness
- 5 Experiments

# What is a Good Explanation?

Two Main Criteria [Molnar, 2025]:

## Conciseness

Is the explanation short enough to be understandable?

⇒ Measured using  $\|\mathbf{w}\|_0$ , which is the number of nonzero weights of  $\mathbf{w}$ .

# What is a Good Explanation?

Two Main Criteria [Molnar, 2025]:

## Conciseness

Is the explanation short enough to be understandable?

⇒ Measured using  $\|\mathbf{w}\|_0$ , which is the number of nonzero weights of  $\mathbf{w}$ .

## Precision

Does the explanation predict the outcome as truthfully as possible?

⇒ An explanation  $\mathbf{w}$  is *sufficient* if  $\mathbf{w} \cdot \mathbf{z} = \mathbf{w} \cdot \mathbf{x}$  implies  $f(\mathbf{z}) = f(\mathbf{x})$  for all  $\mathbf{z} \in \{\pm 1\}^d$ .

# What is a Good Explanation?

Two Main Criteria [Molnar, 2025]:

## Conciseness

Is the explanation short enough to be understandable?

⇒ Measured using  $\|\mathbf{w}\|_0$ , which is the number of nonzero weights of  $\mathbf{w}$ .

## Precision

Does the explanation predict the outcome as truthfully as possible?

⇒ An explanation  $\mathbf{w}$  is *sufficient* if  $\mathbf{w} \cdot \mathbf{z} = \mathbf{w} \cdot \mathbf{x}$  implies  $f(\mathbf{z}) = f(\mathbf{x})$  for all  $\mathbf{z} \in \{\pm 1\}^d$ .

Unfortunately, both criteria are **clashing**: sufficient explanations may require too many weights to be understandable!

## Definition

Given a probability distribution  $\mathcal{D}$  over  $\{\pm 1\}^d$ , the **precision error** of an explanation  $\mathbf{w}$  for  $\mathbf{x}$  is defined as

$$\mathbb{E}_{\mathbf{z} \sim \mathcal{D}}[|f(\mathbf{z}) - f(\mathbf{x})| \mid \mathbf{w} \cdot \mathbf{z} = \mathbf{w} \cdot \mathbf{x}]$$

In other words, the error of  $\mathbf{w}$  is the expected gap between  $f(\mathbf{z})$  and  $f(\mathbf{x})$ , when the projections of  $\mathbf{z}$  and  $\mathbf{x}$  to  $\mathbf{w}$  are the same.

# Computing Probabilistic Explanations

Given

- a black-box model  $f$ ,
- an instance  $\mathbf{x}$  to explain,
- a probability distribution  $\mathcal{D}$  over instances,
- a conciseness parameter  $k$ ,

our problem is to

$$\text{Minimize} \quad \mathbb{E}_{\mathbf{z} \sim \mathcal{D}} [ |f(\mathbf{z}) - f(\mathbf{x})| \mid \mathbf{w} \cdot \mathbf{z} = \mathbf{w} \cdot \mathbf{x} ]$$

$$\text{Subject to} \quad \mathbf{w} \cdot \mathbf{x} = f(\mathbf{x})$$

$$\|\mathbf{w}\|_0 \leq k$$

# Computing Probabilistic Explanations

Given

- a black-box model  $f$ ,
- an instance  $\mathbf{x}$  to explain,
- a probability distribution  $\mathcal{D}$  over instances,
- a conciseness parameter  $k$ ,

our problem is to

Objective Function

Minimize  $\mathbb{E}_{\mathbf{z} \sim \mathcal{D}}[|f(\mathbf{z}) - f(\mathbf{x})| \mid \mathbf{w} \cdot \mathbf{z} = \mathbf{w} \cdot \mathbf{x}]$

Subject to  $\mathbf{w} \cdot \mathbf{x} = f(\mathbf{x})$

$$\|\mathbf{w}\|_0 \leq k$$

# Computing Probabilistic Explanations

Given

- a black-box model  $f$ ,
- an instance  $\mathbf{x}$  to explain,
- a probability distribution  $\mathcal{D}$  over instances,
- a conciseness parameter  $k$ ,

our problem is to

Objective Function

Minimize  $\mathbb{E}_{\mathbf{z} \sim \mathcal{D}}[|f(\mathbf{z}) - f(\mathbf{x})| \mid \mathbf{w} \cdot \mathbf{z} = \mathbf{w} \cdot \mathbf{x}]$

Subject to  $\mathbf{w} \cdot \mathbf{x} = f(\mathbf{x})$  Consistency Constraint

$\|\mathbf{w}\|_0 \leq k$  Conciseness Constraint

## Computational Complexity

The problem of finding an explanation  $\mathbf{w}$  of size at most  $k$  and precision error at most  $\varepsilon$  is **NP<sup>PP</sup>-hard**.

Thus, the problem involves *two independent* sources of complexity:

- Evaluating the precision error of a given vector  $\mathbf{w}$  is **PP-hard**,
- Finding a vector  $\mathbf{w}$  of size at most  $k$  that achieves minimal precision error is **NP-hard**.

- 1 Regression Models
- 2 Probabilistic Explanations
- 3 Dealing with PP-Hardness**
- 4 Dealing with NP-Hardness
- 5 Experiments

# Optimizing Precision

Consider again our optimization problem:

Minimize  $\mathbb{E}_{z \sim \mathcal{D}}[|f(z) - f(x)| \mid \mathbf{w} \cdot z = \mathbf{w} \cdot x]$   $P(\mathbf{w})$

Subject to  $\mathbf{w} \cdot \mathbf{x} = f(\mathbf{x})$

$$\|\mathbf{w}\|_0 \leq k$$

The objective function  $P(\mathbf{w})$  involves a **conditional** expectation, which is very difficult to evaluate.

# Optimizing Fidelity

Now, consider the following variant:

$$\text{Minimize } \mathbb{E}_{\mathbf{z} \sim \mathcal{D}}[(\mathbf{w} \cdot \mathbf{z} - f(\mathbf{z}))^2] \quad F(\mathbf{w})$$

$$\text{Subject to } \mathbf{w} \cdot \mathbf{x} = f(\mathbf{x})$$

$$\|\mathbf{w}\|_0 \leq k$$

By substituting the precision error  $P(\mathbf{w})$  with **fidelity error**  $F(\mathbf{w})$  [Ribeiro et al., 2016], the objective function is an **unconditional** expectation that can be approximated through sampling.

# Optimizing Empirical Fidelity

Finally, given a sample set  $\{(\mathbf{z}_1, f(\mathbf{z}_1)), \dots, (\mathbf{z}_m, f(\mathbf{z}_m))\}$  drawn from the distribution  $\mathcal{D}$  and labeled by the predictor  $f$ , consider the problem:

$$\begin{array}{ll} \text{Minimize} & \frac{1}{m} \sum_{i=1}^m (\mathbf{w} \cdot \mathbf{z}_i - f(\mathbf{z}_i))^2 \\ \text{Subject to} & \mathbf{w} \cdot \mathbf{x} = f(\mathbf{x}) \\ & \|\mathbf{w}\|_0 \leq k \end{array}$$

$\hat{F}(\mathbf{w})$

By approximating the fidelity  $F(\mathbf{w})$  through **empirical fidelity**  $\hat{F}(\mathbf{w})$ , the objective function is now easy to evaluate.

## Approximation Guarantees

Let  $\mathbf{w}^*$  be an optimal explanation for the precision. Then, using a number of samples  $m$  that is logarithmic in  $d$  and quadratic in  $k$ , any explanation  $\mathbf{w}$  that is optimal for the empirical fidelity satisfies with high probability:

$$P(\mathbf{w}) \leq \sqrt{\hat{F}(\mathbf{w}^*)} + o(1)$$

## Approximation Guarantees

Let  $\mathbf{w}^*$  be an optimal explanation for the precision. Then, using a number of samples  $m$  that is logarithmic in  $d$  and quadratic in  $k$ , any explanation  $\mathbf{w}$  that is optimal for the empirical fidelity satisfies with high probability:

$$P(\mathbf{w}) \leq \sqrt{\hat{F}(\mathbf{w}^*)} + o(1)$$

Computing probabilistic explanations of optimal empirical fidelity can be solved via [Mixed Integer Programming](#).

- 1 Regression Models
- 2 Probabilistic Explanations
- 3 Dealing with PP-Hardness
- 4 Dealing with NP-Hardness**
- 5 Experiments

# Optimizing Empirical Fidelity Remains Hard

Even with approximation guarantees, the final problem:

$$\text{Minimize} \quad \frac{1}{m} \sum_{i=1}^m (\mathbf{w} \cdot \mathbf{z}_i - f(\mathbf{z}_i))^2$$

$$\text{Subject to} \quad \mathbf{w} \cdot \mathbf{x} = f(\mathbf{x})$$

$$\|\mathbf{w}\|_0 \leq k$$

Conciseness Constraint

is still NP-hard because the conciseness constraint is **not** convex!

# Optimizing Empirical Fidelity Remains Hard

Even with approximation guarantees, the final problem:

$$\text{Minimize } \frac{1}{m} \sum_{i=1}^m (\mathbf{w} \cdot \mathbf{z}_i - f(\mathbf{z}_i))^2$$

$$\text{Subject to } \mathbf{w} \cdot \mathbf{x} = f(\mathbf{x})$$

$$\|\mathbf{w}\|_0 \leq k$$

Conciseness Constraint

is still NP-hard because the conciseness constraint is **not** convex!

⇒ We thus need additional assumptions to achieve polynomial-time efficiency.

## Concentration Inequality [Achlioptas, 2001]

Let  $\mathcal{D}$  be the uniform distribution over  $\{\pm 1\}^d$ . Then, for any  $\mathbf{w} \in \mathbb{R}^d$  and any  $\varepsilon \in (0, 1)$ , we have

$$\mathbb{P}_{\mathbf{Z} \sim \mathcal{D}^m} \left[ \left| \frac{1}{m} \|\mathbf{Z}\mathbf{w}\|_2^2 - \|\mathbf{w}\|_2^2 \right| > \varepsilon \right] \leq 2e^{-\Omega(m)}$$

## Concentration Inequality [Achlioptas, 2001]

Let  $\mathcal{D}$  be the uniform distribution over  $\{\pm 1\}^d$ . Then, for any  $\mathbf{w} \in \mathbb{R}^d$  and any  $\varepsilon \in (0, 1)$ , we have

$$\mathbb{P}_{\mathbf{Z} \sim \mathcal{D}^m} \left[ \left| \frac{1}{m} \|\mathbf{Z}\mathbf{w}\|_2^2 - \|\mathbf{w}\|_2^2 \right| > \varepsilon \right] \leq 2e^{-\Omega(m)}$$

Thus, if the number of samples  $m$  is sufficiently large, any matrix  $\mathbf{Z} \sim \mathcal{D}^m$  satisfies, with high probability, the **Restricted Isometry Property** for all  $k$ -sparse vectors  $\mathbf{w}$  [Baraniuk et al., 2008]:

$$(1 - \beta_k) \|\mathbf{w}\|_2^2 \leq \frac{1}{m} \|\mathbf{Z}\mathbf{w}\|_2^2 \leq (1 + \beta_k) \|\mathbf{w}\|_2^2$$

# Iterative Hard Thresholding

Input      Query  $(\mathbf{x}, f(\mathbf{x}))$ , sparsity level  $k$ , samples  $(\mathbf{Z}, \mathbf{y})$

Initialize    $\mathbf{w}_0 = 0$

For each     $t = 1, 2, \dots$  do

$$\mathbf{v}_t = \mathbf{w}_{t-1} - \frac{1}{m} \mathbf{Z}^T (\mathbf{Z} \mathbf{w}_{t-1} - \mathbf{y})$$

Gradient Descent

$$\mathbf{w}_t = \operatorname{argmin}_{\mathbf{w} \in \mathcal{W}} \|\mathbf{v}_t - \mathbf{w}\|_2$$

Projection onto feasible explanations

$$\mathcal{W} = \{\mathbf{w} \in \mathbb{R}^d : \mathbf{w} \cdot \mathbf{x} = f(\mathbf{x}) \text{ and } \|\mathbf{w}\|_0 \leq k\}$$

## Efficiency Guarantees

Let  $\mathbf{w}^*$  be an optimal explanation for the precision. Then, using a polynomial number of samples drawn uniformly at random, the IHT algorithm is guaranteed to find, with high probability, a  $k$ -sparse explanation  $\mathbf{w}_t$  that achieves

$$P(\mathbf{w}_t) \leq 7\sqrt{\hat{F}(\mathbf{w}^*)} + o(1)$$

Furthermore,  $\mathbf{w}_t$  can be computed in polynomial time with respect to  $d$ ,  $k$ , and  $\log_2 \lceil 1/\hat{F}(\mathbf{w}^*) \rceil$ .

- 1 Regression Models
- 2 Probabilistic Explanations
- 3 Dealing with PP-Hardness
- 4 Dealing with NP-Hardness
- 5 Experiments**

## Benchmarks

- Prediction Tasks: 20 regression datasets from OpenML
- Black-box  $f$ : neural networks (MLP) learned from train set
- Instance  $x$ : selected uniformly at random from the test set
- Distributions:  $\mathcal{D}$  parameterized by spread  $\sigma \in [0, 1]$
- Explanation Sizes:  $k \in \{1, \dots, 10\}$
- Number of samples:  $m = 1000$
- Timeout: 60 seconds

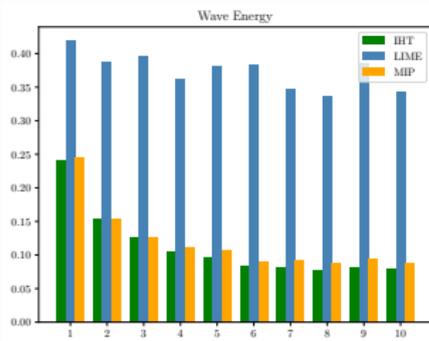
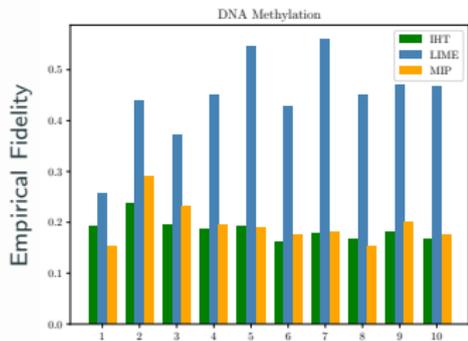
## Competitors

- IHT, MIP (Gurobi solver) versus CVX (convex relaxation), LIME, and MAPLE.

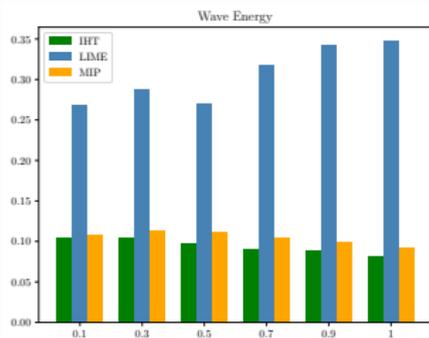
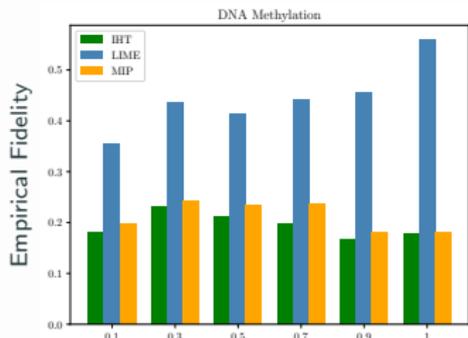
# Results

Benchmark Name	Empirical Fidelity				
	CVX	IHT	LIME	MAPLE	MIP
Airfoil Self Noise	0.040	0.055	0.321	0.218	<b>0.049</b>
Auto MPG	0.031	0.069	0.338	0.122	<b>0.039</b>
Liver Disorders	0.059	0.091	0.209	0.147	<b>0.068</b>
Medical Charges	0.040	<b>0.049</b>	0.408	0.204	<b>0.049</b>
Ailerons	0.050	0.201	0.647	0.113	<b>0.085</b>
Auto Imports	0.067	0.232	0.528	0.148	<b>0.107</b>
DNA Methylation	0.121	<b>0.192</b>	0.582	0.168	<b>0.191</b>
NCI 60 Thioguanine	0.062	0.235	0.534	0.108	<b>0.132</b>
Student Performance	0.074	0.143	0.454	0.169	<b>0.105</b>
Wave Energy	0.017	<b>0.080</b>	0.301	0.128	0.091

Results on 4 low-dimensional benchmarks and 6 medium-dimensional benchmarks, using  $k = 7$ ,  $m = 1000$ , and  $\sigma = 1$ . Entries highlighted in green indicate that all generated explanations were  $k$ -sparse.



Increasing Explanation Sizes  $k$



Increasing Spread  $\sigma$

Comparisons among IHT, Lime and MIP for several parameters.

# Conclusion

- Probabilistic explanations achieve a balance between *conciseness* and *precision*.
- However, computing these explanations is very challenging ( $\text{NP}^{\text{PP}}$ -hard).
- By replacing precision with fidelity, they can be approached with **Mixed Integer Programming**, while offering approximation guarantees.
- Under the uniform distribution, these explanations can be efficiently approached using **Iterative Hard Thresholding**.
- Empirical results on real-world benchmarks support our theoretical findings.

# Conclusion

- Probabilistic explanations achieve a balance between *conciseness* and *precision*.
- However, computing these explanations is very challenging ( $\text{NP}^{\text{PP}}$ -hard).
- By replacing precision with fidelity, they can be approached with **Mixed Integer Programming**, while offering approximation guarantees.
- Under the uniform distribution, these explanations can be efficiently approached using **Iterative Hard Thresholding**.
- Empirical results on real-world benchmarks support our theoretical findings.

Thank you for your attention!

- Dimitris Achlioptas. Database-friendly random projections. In *Proceedings of the 20th ACM Symposium on Principles of Database Systems (PODS)*, page 274–281, 2001.
- Richard Baraniuk, Mark Davenport, Ronald DeVore, and Michael Wakin. A simple proof of the restricted isometry property for random matrices. *Constructive Approximation*, 28:253–263, 2008.
- Christoph Molnar. *Interpretable Machine Learning*. leanpub.com, 3rd edition, 2025. URL <https://christophm.github.io/interpretable-ml-book>.
- Marco Túlio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should I trust you?": Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1135–1144, 2016.