

# Local vs. Global Interpretability: A Computational Complexity Perspective

What is computationally harder - interpreting an ML model *locally* or *globally*?  
The answer: It depends...

## Introduction

- Can we *mathematically* assess how interpretable an ML model is?
- Previous work addresses this using *computational complexity theory* – the harder it is to generate an explanation, the less interpretable a model is.
- For example, some explanation forms can be efficiently obtained for decision trees and linear models but are intractable for neural networks.

## Local vs. Global Interpretability

- The problem: previous work mainly analyzed the complexity of obtaining *local* explanation forms.
- Interpretability can be *local* (understanding specific decisions) or *global* (understand general behaviors)
- For example, Molnar et. al. claim that while the weights of a linear model help interpret local decisions, they don't explain global behaviors.

We analyze the computational complexity of obtaining local explanations (for a specific instance  $x$ ) compared to their global variants (for *all*  $x$  in a given domain)

### Local Sufficient Reason

$$\forall(\mathbf{z} \in \mathbb{F}). [f(\mathbf{x}_S; \mathbf{z}_{\bar{S}}) = f(\mathbf{x})]$$

$$f \left( \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \right) = \text{"Class 1"}$$

$$f \left( \begin{array}{|c|c|c|} \hline \text{Free} & \text{Free} & \\ \hline \text{Free} & \text{Free} & \\ \hline & & \\ \hline \end{array} \right) = \text{"Class 1"}$$

### Global Sufficient Reason

$$\forall(\mathbf{x}, \mathbf{z} \in \mathbb{F}). [f(\mathbf{x}_S; \mathbf{z}_{\bar{S}}) = f(\mathbf{x})]$$

$$f \left( \begin{array}{|c|c|c|} \hline \text{Free} & \text{Free} & \\ \hline \text{Free} & & \\ \hline & & \\ \hline \end{array} \right) = \text{"Class 1"}$$

$$f \left( \begin{array}{|c|c|c|} \hline \text{Free} & \text{Free} & \\ \hline \text{Free} & & \\ \hline & & \\ \hline \end{array} \right) = \text{"Class 2"}$$

$$f \left( \begin{array}{|c|c|c|} \hline \text{Free} & \text{Free} & \\ \hline \text{Free} & & \\ \hline & & \\ \hline \end{array} \right) = \text{"Class 3"}$$

### Explanation forms we analyze:

- Feature selection – selecting minimum local vs. global sufficient reasons
  - Identifying local vs. global redundant or necessary features
  - A probabilistic notion of sufficiency – whether local or global

In **linear models**, feature selection is *harder to perform globally than locally*:

Explanation Type	Local	Global
Check Sufficient Reason	PTIME	coNP-Complete
Minimum Sufficient Reason	PTIME	coNP-Complete

Less globally interpretable

In **neural networks and decision trees** – the *local task is harder than the global one*.

Explanation Type	Decision Trees		Neural Networks	
	Local	Global	Local	Global
Minimum Sufficient Reason	NP-Complete	PTIME	$\Sigma_2^P$ -Complete	coNP-Complete
Feature Redundancy	coNP-Complete	PTIME	$\Pi_2^P$ -Complete	coNP-Complete

Less locally interpretable

### Why does this happen?

There can be an exponential number of minimal *local* sufficient reasons:

$$f \left( \begin{array}{|c|c|c|} \hline \text{Free} & \text{Free} & \\ \hline \text{Free} & & \\ \hline & & \text{Free} \\ \hline \end{array} \right) = \text{"Class 1"} \quad f \left( \begin{array}{|c|c|c|} \hline & & \text{Free} \\ \hline & \text{Free} & \\ \hline \text{Free} & & \\ \hline \end{array} \right) = \text{"Class 1"}$$

$$f \left( \begin{array}{|c|c|c|} \hline \text{Free} & \text{Free} & \\ \hline \text{Free} & & \\ \hline & & \\ \hline \end{array} \right) = \text{"Class 1"}$$

There is one unique minimal *global* sufficient reason:

$$f \left( \begin{array}{|c|c|c|} \hline \text{Free} & \text{Free} & \\ \hline \text{Free} & & \\ \hline & & \\ \hline \end{array} \right) = \text{"Class 1"} \quad f \left( \begin{array}{|c|c|c|} \hline \text{Free} & \text{Free} & \\ \hline \text{Free} & & \\ \hline & & \\ \hline \end{array} \right) = \text{"Class 2"}$$

$$f \left( \begin{array}{|c|c|c|} \hline \text{Free} & \text{Free} & \\ \hline \text{Free} & & \\ \hline & & \\ \hline \end{array} \right) = \text{"Class 3"}$$

This property simplifies the global task compared to the local one. However, for linear models, the weights enable the development of efficient algorithms, which is not evident in neural networks and decision trees.

## Conclusion

- We prove that there is often a *strict complexity gap* between obtaining an explanation locally vs. globally.
- In some cases, such gaps justify folklore claims (e.g., linear models can be interpreted locally using their weights but not globally).
- In other cases, they yield unexpected outcomes (e.g., neural networks and decision trees are easier to interpret globally than locally in some contexts).

