**UCLouvain**

**FACULTÉ DES SCIENCES**

**Faculté des sciences**

# XAI : A Comparative Study of Post-hoc and Intrinsic Explainability Methods for Categorical Tabular Data

Author: **Siméon GODFRIN**
Supervisors: **Olivier CAELEN, Corentin VANDE KERCKHOVE, Robin VAN OIRBEEK**
Reader: **Marco SAERENS**
Academic year 2025–2026
Master [120] en science des données, orientation statistique

Ceci est la version corrigée de mon mémoire

# Contents

# Chapter 1

# Introduction

## 1.1    Contextualization

This Master's thesis focuses on eXplainable Artificial Intelligence (XAI) and aims to compare several explainability methods within the framework of a multi-class prediction task applied to an almost fully categorical dataset. The dataset consists of a survey on the habits of Quebec residents conducted in 2022. One of the primary objectives of collecting this dataset was to enable the development of a tool allowing individuals to position themselves on the political spectrum solely based on their lifestyle and demographic features (e.g., age, income). This represents a challenging and relatively uncommon case on which to test both the explainability methods and the comparison benchmark on a real-world dataset.

A central challenge of this task lies in defining criteria that allow for a meaningful comparison between fundamentally different explainability techniques. This difficulty arises from several factors, including diversity in explanation techniques, variability in explanation outputs, dependence on data type, and the lack of consensus on what constitutes a good explanation in the Machine Learning (ML) field. To address this challenge, a structured comparison framework is proposed in Section 5, largely inspired by the functionally grounded evaluation approaches described in [14] and [35]. Those scientific articles provide answers to the question of how to evaluate and compare XAI methods in a systematic way, based on the Ethics Guidelines for Trustworthy AI published by the European Union's High Level Expert Group (AI HLEG) [20].

This work begins by introducing XAI and the accuracy and interpretability tradeoff in section 1.2. Section 2 presents the datasets and the predictive model, a Neural Network (NN), used as a black-box model to be explained. Section 3 defines the theoretical foundation of the three post-hoc explainability techniques used — SHAP, LIME and Counterfactuals — while Section 4 focuses on the intrinsically interpretable model used — Optimal Sparse Decision Tree (OSDT). Those techniques were chosen because they suit the nature of the data and the task at hand. Moreover, the three post-hoc techniques are commonly used and I was interested in understanding them more deeply and to see how they compare to an intrinsically interpretable model and to each other's. The section 5 defines the characteristics providing the formal definition of a XAI's method desired properties, on which we are going to compare the 4 approaches. The section 6 discusses the findings and concludes the thesis.

## 1.2    Organization

### 1.2.1    Defintion of XAI

Following DARPA (Defense Advanced Research Projects Agency), one of the earliest actor of the field, XAI can be defined at aiming "to produce more explainable models, while maintaining a high level of learning performance (prediction accuracy) [. . . ].". This field resurfaced in response to the prohiminent use of Neural Network (and DNN), the need to explain those models and the increasingly more complex datasets. This is still an ongoing and recent research field without a well established consensus about the taxonomy. The term "explainable" and "interpretable" are either used interchangeably or defined as two different concepts. In the following work, one will use them interchangeably but the term "interpretable/interpretability" will mostly be used following the current ruling taxonomy in the scientific and ML community. [5] Attempts have been done to establish a taxonomy but one finds it hard to bend to since most online papers do not follow it yet.

**The interpretability and accuracy trade-off**

One common idea in the ML world is that there exists a tradeoff between interpretability and accuracy [6, 5, 22] :

- On one hand, simple models— such as linear and logistic regression (LR), decision tree (DT) — are intrinsicaly interpretable, given their construction, but are deemed less accurate and less fit to model complex real-life problems.

- On the other hand, Deep Neural Networks (DNNs) models yield a higher accuracy "at the cost" of interpretability. (By increasing the model's complexity).
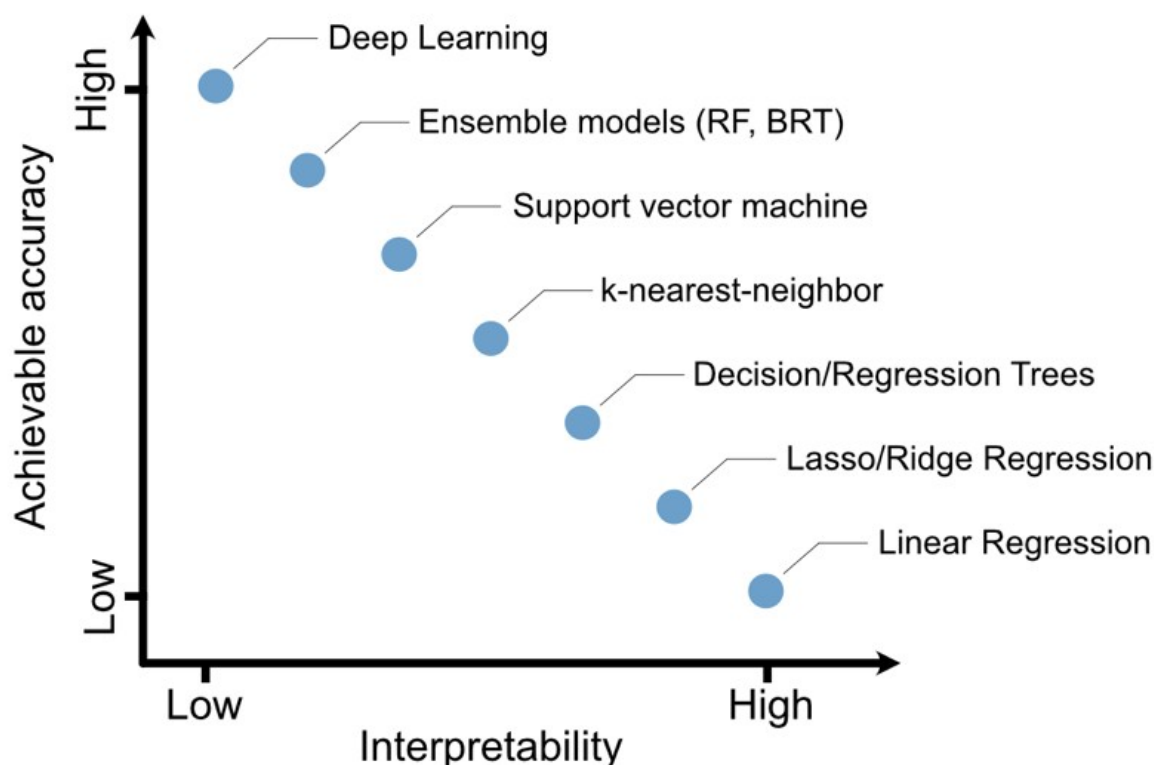


Figure 1.1: The accuracy-interpretability tradeoff : as the accuracy increases, the interpretability decreases. Source : [40]

The belief in this tradeoff partly emerged from older sources such as Leo Breiman, who introduced the Rashomon effect in statistics [13].

He asserted, more than 20 years ago, the superiority of accuracy over interpretability for statistical models.

This belief isn't simply an old idea, it is still present in the recent literature. Johansson et al.(2011) [28] write "Models exhibiting [high accuracy] are many times more complex and opaque, while interpretable models may lack the necessary accuracy." In their study of 16 biomedical classification tasks, black-box methods (like ensembles or SVMs) generally obtained higher accuracies than transparent models. This resonates with the intuition that a sufficiently complex model can fit intricate, non-linear relationships in data that a simpler model cannot capture. Those two examples might give a piece of explanation of the current inertia and tendency in the ML community **to reach for complexity first**, motivated by the desire of maximizing accuracy in a world where data is increasingly complex.

Nonetheless, the belief into this tradeoff, or not, should not limit one to consider intrinsically interpretable models when the context allows it.

The existence of this tradeoff as general rule is debated, as it appears to be based more on intuition and heuristic than in solid empirical evidence. Most recent research tends to align with a circumstantial existence of the trade-off. Some qualify it as a "myth" [46], while others found that when taking the user's interpretation into account that "there are contexts where black-box model can be more explainable and less confusing". [11]

The performance gap between complex and simple models exists but is not always large. Johansson et al.(2011) [28] found that, although black-box models tended to be best, the accuracy difference was typically modest — suggesting the trade-off might be a gentle slope rather than a steep cliff.

**In case of tabular data**, the gap in performance between interpretable and black-box model tends to be lower.[22], [46] Neural networks, in this framework, generally find no decisive advantage. It has been known for a very long time that very simple models perform surprisingly well for tabular data [25]. And this finding resurfaces in recent research.

For example, research in finance and medicine have noted that simpler models (like logistic regression or decision trees) often come within a few percentage points of the accuracy of state-of-the-art algorithms. [46]

Goethals, Martens & Evgeniou (2022) [22] analysed 90 benchmark classification tabular datasets and finds that the "trade-off exists for most (69%) of the datasets, but that somewhat surprisingly for the majority of cases it is rather small while for only a few it [the trade-off] is very large". They found that the gap in accuracy between black box and white box models is correlated to dataset's characteristics such as how difficult it is to linearly separate it. The finding reinforces our belief in the superiority of black box models in capturing the non-linearities. It also opens the door to a better identification of the conditions as to *when* using which type of models.

**In the case of *raw/unstructured data***, such as images, neural networks currently have an advantage over other approaches. "The difference in performance between comprehensible models compared to a black box ones such as DL is considered "unbridgeable." Goethals, Martens & Evgeniou (2022) [22]

Yet, even in these domains, hybrid approaches are emerging (e.g. prototype-based convolutional networks [15])

"These two data extremes show that in ML, the dichotomy between the accurate black box and the less-accurate interpretable model is false" [47], or at the very least, is much more nuanced. In certain contexts, i.e. tabular data, depending on the task at hand, we must not choose accuracy *or* interpretability — we can have *both*.

Therefore, in case of tabular data, a well-constructed transparent model *could* yield both an high accuracy and be intrinsically interpretable. [48] This idea is central, to the "Interpretable Machine Learning" ("IML" in short), which one will frame as sub-field within XAI.

Those findings serve to remind the ML practitioners to explore the space of simpler models before resorting to complex black-box models if the context suits it.

## The Rashomon effect

In the context of tabular data, I have previously mentioned that interpretable models can match the performance of black-box models. In principle, if one searches the space of simpler models thoroughly, one may discover a model that performs comparably to its more complex counterpart.

This observation shifts the focus of model selection: if multiple models achieve near-identical accuracy, it becomes reasonable to prioritize other properties such as interpretability or fairness. In the direct continuity of this idea, the Rashomon set refers to the ML models achieving similar performance above a given threshold. [48] Research on the Rashomon set offers theoretical support for a practice often used heuristically in applied machine learning: selecting the simplest model that performs well. For a more formal treatment of the Rashomon set and its implications, see [50].

In conclusion, this trade-off is a central idea, a knot in the field, embodied by the Post-hoc vs Intrinsically interpretable models dichotomy. The belief in a strict trade-off between accuracy and interpretability may lead practitioners to assume a forced compromise. However, empirical studies [17, 25, 48] have shown that this trade-off is often circumstantial and depends on the specific data and task at hand. The Rashomon set offers a more flexible perspective: if many models perform equally well, at least one of them may be interpretable. This insight softens the perceived trade-off and aligns with what has been observed in practice.

### 1.2.2   The three axis

This work will revolve around the comparison of 4 explainability approaches — three post-hoc and one intrinsically interpretable model— examined across three dimensions inspired by [6] : i) Exploratory Data Analysis, ii) Model-Level analysis and iii) Comparisons of the approaches.
The comprehension and analysis of the data will be considered as an entire part of the explainability process since there is no understanding of the model if one does not understand the task at hand.
The three axis are the following :

- (I) Exploratory Data Analysis (EDA) : the comprehension and exploration of the raw data.

- (II) Models analysis : Definition of the post-hoc interpretability methods — SHAP, LIME, Counterfactuals— and of intrinsically interpretable model —Optimal Sparse Decision Tree (OSDT).

- (III) Comparison between the methods based on functionally-grounded evaluation, i.e. quantitative metrics independent of explainees.