

On Trustworthy Rule-Based Models and Explanations

Mohamed Siala¹, Jordi Planes², and Joao Marques-Silva³

¹LAAS-CNRS, Université de Toulouse, CNRS, INSA Toulouse,
France

²Universitat de Lleida

³ICREA & University of Lleida

July 11, 2025

Abstract

A task of interest in machine learning (ML) is that of ascribing explanations to the predictions made by ML models. Furthermore, in domains deemed high risk, the rigor of explanations is paramount. Indeed, incorrect explanations can and will mislead human decision makers. As a result, and even if interpretability is acknowledged as an elusive concept, so-called interpretable models are employed ubiquitously in high-risk uses of ML and data mining (DM). This is the case for rule-based ML models, which encompass decision trees, diagrams, sets and lists. This paper relates explanations with well-known undesired facets of rule-based ML models, which include negative overlap and several forms of redundancy. The paper develops algorithms for the analysis of these undesired facets of rule-based systems, and concludes that well-known and widely used tools for learning rule-based ML models will induce rule sets that exhibit one or more negative facets.

Keywords: Explainability Interpretability Rule-based models Formal Methods.

1 Introduction

Explainable Artificial Intelligence (XAI) is a mainstay of trustworthy AI [Adadi and Berrada, 2018, Guidotti et al., 2019, Carvalho et al., 2019, Dwivedi et al., 2023, Schwalbe and Finzel, 2024]. Furthermore, in domains that are deemed of high risk, explanations should be trustable [Rudin, 2019, Rudin et al., 2022, Huysmans et al., 2011, Freitas, 2014]. The importance of explanations and the need to trust those explanations motivated work on so-called interpretable models [Rudin, 2019, Rudin et al., 2022], even though it is generally accepted that a rigorous definition of interpretability is elusive at best [Lipton, 2018].

Rule-based models, which encompass decision trees [Breiman et al., 1984], diagrams [Hu et al., 2022, Florio et al., 2023], sets [Clark and Niblett, 1989, Lakkaraju et al., 2016, Yu et al., 2021] and lists [Rivest, 1987, Ghosh et al., 2022], epitomize interpretable models.

Work on the induction of rule-based models can be traced at least to the 1970s [Shwayder, 1971, Hyafil and Rivest, 1976], in the concrete case of decision trees.¹ Decision trees are widely used in practice and often exemplify interpretable models [Rudin, 2019, Rudin et al., 2022]. The perceived importance of interpretability has recently motivated the development of algorithms for learning optimal decision trees [Demirovic et al., 2022, Hu et al., 2020]. Decision sets (or rule sets) find a wide range of uses in different domains [Fürnkranz et al., 2012, Fürnkranz and Kliegr, 2015, Rapp et al., 2020, Hüllermeier et al., 2020, Rapp et al., 2024, Atzmueller et al., 2024]. As with decision trees, there has been recent interest in learning optimal decision sets [Lakkaraju et al., 2016]. Decision lists also find many practical uses, but claims about their interpretability are harder to justify [Marques-Silva and Ignatiev, 2023]. As a result, this paper studies decision sets, but also decision trees when viewed as a special case of decision sets.

At present, some of the best-known ML toolkits implement one or more methods of induction of rule-based models [Pedregosa et al., 2011, Rapp et al., 2020, Demsar et al., 2013]. Nevertheless, it has been argued [Marques-Silva and Ignatiev, 2023] that rule-based methods, although easier to fathom by human-decision makers, still require explanations to be computed. (Otherwise, human decision-makers would be expected to manually solve NP-hard function problems [Marques-Silva and Ignatiev, 2023].) Therefore, a key question is: *for rule-based models, when can explanations be computed trivially, such that a human decision-maker can manually produce an explanation?*

This paper shows that rigorous explanations can be found manually whenever some undesired facets of decision sets are nonexistent. Concretely, the paper relates easy-to-compute explanations with the non-existence of *negative overlap*, i.e. the existence of cases where two or more rules can fire that predict different values. Furthermore, the non-existence of redundant literals in rules is shown to be a necessary condition for minimality of explanations.

Given this state of affairs, the paper then investigates whether existing ML toolkits are able to learn rule-based models that avoid the aforementioned negative facets. As the results demonstrate, this is not the case. In addition, the paper investigates whether model-agnostic methods targeting feature selection (i.e. that produce rules as explanations) are capable of preventing negative overlap (i.e. the most worrisome negative facet). Unfortunately, as the results show, this is also not the case with the well-known explainer Anchor [Ribeiro et al., 2018].

¹Although extremely popular in ML and DM, decision trees found earlier uses in other domains, e.g. https://en.wikipedia.org/wiki/Phylogenetic_tree and https://en.wikipedia.org/wiki/Decision_tree.

Contributions. The paper studies decision sets,² concretely the problem of *negative overlap*, i.e. when two rules that predict different classes fire, but also the existence of local or global redundancies of literals in rules. The paper develops algorithms for deciding the existence of negative overlap, but also for deciding local and global redundancy. Furthermore, the results in the paper take into account possible constraints on the inputs. The paper then relates these negative facets of decision sets with the ability of human decision-makers to manually produce rigorous explanations, namely abductive explanations. In addition, the experiments confirm that implemented rule-learning algorithms in well-known toolkits exhibit the negative facets of decision sets, thus complicating (complexity-wise) the computation of rigorous explanations.

Organization. The paper is organized as follows. Section 2 introduces the notation and definitions used throughout the paper. Section 3 briefly comments on related work. Section 4 details the paper’s main contributions. Section 5 reports on the experimental results. Finally, Section 6 concludes the paper.

2 Background

The notation and definitions used throughout the paper are adapted from past works [Lakkaraju et al., 2016, Biere et al., 2021, Izza et al., 2022].

Propositional Logic and Generalizations [Biere et al., 2021]. Let $X = \{x_1, \dots, x_n\}$ be a set of Boolean variables. A literal is a Boolean variable or its negation. A clause C is a disjunction of literals and a cube L is a conjunction of literals. We use the notation $l_i \in C$ (respectively $l_i \in L$) if $C = l_1 \vee \dots \vee l_k$ (respectively $L = l_1 \wedge \dots \wedge l_k$). A conjunctive normal form (CNF) formula F is a conjunction of clauses. That is, $F = C_1 \wedge \dots \wedge C_k$ where C_j is a clause. In this case, we use the notation $C_j \in F$. Note by definition that a clause/cube is a CNF. An assignment $v = (v_1, \dots, v_n)$ is a point in $\{0, 1\}^n$. If $F = C_1 \wedge \dots \wedge C_k$ is a CNF, $v \models F$ iff $\forall C_j \in F, \exists x_i \in C_j$ such that $v_i = 1$ or $\exists \neg x_i \in C_j$ such that $v_i = 0$. If $\exists v \in \{0, 1\}^n$ such that $v \models F$ then F is said satisfiable, otherwise unsatisfiable. If F_1 and F_2 are two CNF formulas, $F_1 \models F_2$ iff $v \models F_1 \implies v \models F_2$. Note that $F_1 \models F_2$ iff $F_1 \wedge \neg F_2$ is unsatisfiable. Given a CNF formula F , the satisfiability problem (SAT) asks if F is satisfiable. SAT solvers are highly deployed in practice to answer SAT related queries, such as finding satisfying assignments or proving unsatisfiability [Biere et al., 2021]. Furthermore, extensions of propositional to more expressive logics can be handled by considering Satisfiability Modulo Theories (SMT) [Biere et al., 2021].

Machine Learning. We consider rule-based models for classification and regression that can be represented as a set of unordered rules. Let $\mathcal{F} = \{1, \dots, m\}$ be a set of features where each feature i takes values from a domain D_i . The

²Decision trees are a special case of a decision set, and so we also present experiments on decision trees. However, we opt not to address decision lists due to the intrinsic difficulties with their explanation [Marques-Silva and Ignatiev, 2023].

feature space is the Cartesian product of the domains $\mathbb{F} = D_1 \times \dots \times D_m$. The outcome space (i.e., classes for classification and numerical values for regression) is denoted by \mathcal{V} . A dataset is a set $\{(x, o) \mid x \in \mathbb{F} \wedge o \in \mathcal{V}\}$, and where $x = (x_1, \dots, x_m)$. A literal represents a condition on the values of a feature. We use \mathbb{L} to represent the universe of literals. A background knowledge \mathcal{B} is a propositional formula over literals from \mathbb{L} that specifies the conditions that any arbitrary point in feature space must comply with. In other words, a point in feature space x is *valid* iff $x \models \mathcal{B}$. We assume in the rest of the paper that \mathcal{B} is given as a CNF. For example, consider a dataset representing individuals and the two literals $l_1 := \text{employed}$, $l_2 := \text{salary} > 50k$. The background knowledge \mathcal{B} can contain the clause $l_1 \vee \neg l_2$ to model the fact that an unemployed individual cannot have a salary greater than 50k. Note that \mathcal{B} can be a tautology (for instance when no condition is given). In this case, any arbitrary point in feature space is a valid. A user can also miss certain constraints she is not aware of. Let $\lambda \notin \mathcal{V}$ be a dummy value. A supervised ML (classification or regression) model κ is a mapping from \mathbb{F} to $\{\lambda\} \cup \mathcal{V}$ such that $\kappa(x) = \lambda$ iff $x \not\models \mathcal{B}$.

A rule R_i is a pair (L_i, o_i) such that L_i is a conjunction of literals (i.e., cube) from $\mathcal{L} \subseteq \mathbb{L}$ and $o_i \in \mathcal{V}$. R_i fires on $x \in \mathbb{F}$ iff $x \models L_i$. With a slight abuse of notation we shall sometimes use L_i as the subset of \mathcal{L} formed by the literals in L_i . A decision set \mathcal{M} is a set of rules $\mathcal{M} = \{R_1, \dots, R_r, R_{r+1}\}$ such that $\forall i \leq r, L_i \neq \emptyset$ and $L_{r+1} = \emptyset$. R_{r+1} is called the default rule. We denote $\Delta(o)$ the set $\{R_i \mid o_i = o\}$. \mathcal{M} is used as an ML model $\kappa_{\mathcal{M}}$ as follows:

$$\kappa_{\mathcal{M}}(x) = \begin{cases} \lambda \notin \mathcal{V} & \text{if } x \not\models \mathcal{B} \\ o_{r+1} & \text{if no rule fires on } x \\ o & \text{if } \{o\} = \{o_i \mid R_i \text{ fires on } x\} \\ \text{Tie-breaking strategy otherwise} \end{cases}$$

Note that decision trees (DTs), decision diagrams (DDs), random forests (RFs) and boosted trees (BTs), can be seen as decision sets where each path represents a rule. Clearly, in such models the default rule never fires. In the case of DTs and DDs, each input fires exactly one rule (since it follows exactly one path). Thus, no tie-breaking strategy is needed. This is not the case with RFs and BTs since each input fires one rule on each tree. Therefore, a tie-breaking strategy is needed.

We extend the notion of cover and overlap from [Lakkaraju et al., 2016] by considering the background knowledge \mathcal{B} and the input space.

Definition 1 (Cover). *Given $X \subseteq \mathbb{F}$ and background knowledge \mathcal{B} , $\text{Cover}(X, B, L) = \{x \mid x \in X \wedge x \models \mathcal{B} \wedge x \models L\}$.*

Definition 2 (Overlap). *Given a background knowledge \mathcal{B} , two rules R_i and R_j such that $i, j \leq r$ overlap in $X \subseteq \mathbb{F}$ iff $\text{Cover}(X, B, L_i) \cap \text{Cover}(X, B, L_j) \neq \emptyset$.*

We say that R_i and R_j positively (respectively negatively) overlap if they overlap and $o_i = o_j$ (respectively $o_i \neq o_j$). We use the notation $R_i \ominus R_j$ if R_i and R_j negatively overlap. Observe that DTs and DDs exhibit no overlap since each input is captured by exactly one rule. This is not the case for RFs and BTs, since each input fires exactly one rule from each tree. Thus, overlaps may occur only between rules from different trees.

Formal Explanations [Izza et al., 2022, Darwiche, 2023]. Most approaches to explainability target at instance, i.e. a pair (x, c) with $x \in \mathbb{F}$ and $c \in \mathcal{V}$. We use κ throughout the paper to denote a machine learning model. Given an instance (v, c) , with $c = \kappa(v)$, a weak abductive explanation (WAXp) is a subset \mathcal{X} of the features \mathcal{F} which, if assigned the values dictated by v , is sufficient for the classifier to output prediction $c = \kappa(v)$ [Izza et al., 2022, Darwiche, 2023]:

$$\forall(x \in \mathbb{F}). \left[\bigwedge_{i \in \mathcal{X}} (x_i = v_i) \rightarrow (\kappa(x) = c) \right] \quad (1)$$

A subset-minimal WAXp is an *abductive explanation* (AXp). Recent work demonstrated the need for explaining interpretable models, including decision trees [Izza et al., 2022] and lists [Marques-Silva and Ignatiev, 2023]. To the best of our knowledge, past work did not investigate formal explanations for decision sets.

Furthermore, the definition of WAXp (see (1)) can be generalized to account for literals involving other relational operators [Izza et al., 2022] (e.g. relational operators taken from $\{\in, \geq, >, <, \leq\}$). In addition, constraints on the inputs [Gorji and Rubin, 2022, Audemard et al., 2024] can be accounted for by conjoining a set of constraints $\mathcal{C}_{\mathcal{B}}$. For example, these constraints allow capturing the background knowledge introduced earlier in this section. Concretely, we write that $\mathcal{C}_{\mathcal{B}}(x)$ holds true iff x respects the background knowledge, i.e. $x \models \mathcal{B}$.

3 Related Work

The learning of rule-based models has been the subject of research since the 1970s [Shwayder, 1971, Hyafil and Rivest, 1976]. The importance of the topic, especially given their widely accepted interpretability, has motivated recent work on learning decision sets [Rapp et al., 2020, Rapp et al., 2024, Atzmueller et al., 2024] and (optimal) trees [Demirovic et al., 2022]. These earlier works were motivated by the accepted belief that decision trees, sets and lists are interpretable [Breiman, 2001, Rudin, 2019, Rudin et al., 2022]. Accounts of methods for learning decision sets and lists include [Fürnkranz et al., 2012, Fürnkranz and Kliegr, 2015].

Motivated by the elusive nature of interpretability’s definition [Lipton, 2018], recent work [Marques-Silva and Ignatiev, 2023] uncovered practical difficulties in computing and/or using so-called interpretable models as explanations. For example, it has been shown that paths in decision trees can be arbitrarily redundant (on the number of features) when compared with an AXp [Izza et al., 2022]. Similarly, the computation of an AXp for a decision list equates with solving an NP-hard problem [Marques-Silva and Ignatiev, 2023], i.e. something that is in general beyond the capabilities of a human decision-maker. Nevertheless, past work did not address formal explanations for decision sets, arguably because of the existence of negative overlap.

Although the paper assesses rule-based methods using formal explanations, XAI is better-known by the use of model-agnostic methods [Adadi and Berrada, 2018, Guidotti et al., 2019, Carvalho et al., 2019, Minh et al., 2022, Dwivedi

et al., 2023, Schwalbe and Finzel, 2024]. Well-known examples include LIME [Ribeiro et al., 2016], SHAP [Lundberg and Lee, 2017] and Anchors [Ribeiro et al., 2018]. Since so-called interpretable models have been proposed for high-risk uses of ML, we focus on rigorous (i.e. formal) explanations.

The main results of this paper, namely the direct relationship between easy-to-compute explanations and the non-existence of well-known negative facets of rule-based models, are novel. The observation that rule-based models, obtained with well-known toolkits, exhibit those negative facets, is also a novel result, to the best of our knowledge.

4 Overlap and Redundancy

In this section, we let \mathcal{B} be a background knowledge and $\mathcal{M} = \{R_1, \dots, R_r, R_{r+1}\}$ be a decision set where R_{r+1} is the default rule such that each rule $R_{i \leq r}$ fires on at least one valid input (w.r.t. \mathcal{B}). As mentioned in the introduction, we provide a formal framework to address the following questions: (i) How can we generate all (negative) overlap?; (ii) Is rule R_i redundant in \mathcal{M} ?; and (iii) Is literal l redundant in a given rule?.

We use Example 1 throughout the paper to illustrate the different concepts.

Example 1. \mathcal{B} is background knowledge that encodes the following constraints (in a CNF): $(salary > 0) \leftrightarrow (age \geq 18)$; $(size = 140) \rightarrow (size > 120)$; $(weight > 90) \rightarrow (weight \geq 85)$; and $(weight \geq 85) \rightarrow (weight > 80)$. The decision set contains the following rules:

- $R_1 = ((salary > 0) \wedge (size \neq 140) \wedge (age > 10) \wedge (color = blue) \wedge (weight > 80), 1)$
- $R_2 = ((salary > 0) \wedge (size = 140), 1)$
- $R_3 = ((salary > 0) \wedge (weight > 90), 1)$
- $R_4 = ((size > 120) \wedge (weight < 85), 0)$

4.1 Overlap

We start by giving a sufficient and necessary condition to check if two rules negatively overlap.

Lemma 1 (Overlap Check). *Two rules R_i and R_j overlap iff $\mathcal{B} \wedge L_i \wedge L_j$ is satisfiable.*

Proof. $\mathcal{B} \wedge L_i \wedge L_j$ is satisfiable iff $\exists x \in \mathbb{F}, x \models \mathcal{B} \wedge L_i$ and $x \models \mathcal{B} \wedge L_j$. This is equivalent to $\exists x \in \mathbb{F}, \{x\} \in Cover(\mathbb{F}, \mathcal{B}, L_i) \cap Cover(\mathbb{F}, \mathcal{B}, L_j)$. The latter means that R_i and R_j overlap. \square \square

In Example 1, one can use Lemma 1 to show that R_3 and R_4 do not overlap, in contrast to R_1 and R_4 , which do.

Algorithm 1 Negative Overlap Pairs

```
1: Function: Pairs
2: Input:  $\mathbb{F}, O, \mathcal{M} = \{R_1, \dots, R_r\}, \mathcal{B}$ 
3: Output:  $\Pi = \{(i, j) \mid R_i \ominus R_j\}$ 
4:  $\Pi = \emptyset$ 
5:  $\Psi = \text{GetList}(o_1, o_2, \dots, o_r)$ 
6:  $g = |\Psi|$ 
7: for  $a$  in  $1, \dots, g - 1$  do
8:   for  $b$  in  $a + 1, \dots, g$  do
9:     for  $R_i$  in  $\Delta(\Psi(a))$  do
10:      for  $R_j$  in  $\Delta(\Psi(b))$  do
11:        if  $\mathcal{B} \wedge L_i \wedge L_j$  is SATISFIABLE then
12:           $\Pi \leftarrow \Pi \cup \{(i, j)\}$ 
13:        end if
14:      end for
15:    end for
16:  end for
17: end for
18: Return  $\Pi$ 
```

We consider now the question of generating all negative overlap. Algorithm 1 finds all pairs of rules that exhibit a negative overlap. We use $\text{GetList}(o_1, o_2, \dots, o_r)$ as a function that computes a list that contains the distinct values in $\{o_1, \dots, o_r\}$.

Algorithm1 terminates because each pair of rules will be visited at most once. The correctness of Algorithm1 follows from the fact that each pair (R_i, R_j) such that $o_i \neq o_j$ is visited exactly once in Line 12. The complexity of Algorithm 1 is $O(|\mathcal{M}|^2 \times f(\mathcal{M}))$ where $f(M)$ is the worst complexity of $\mathcal{B} \wedge L_i \wedge L_j$ for an arbitrary pair of rules (R_i, R_j) . This observation follows from the fact that computing GetList can be naturally be done in $O(|\mathcal{M}|)$ and the fact that the satisfiability check in Line 12 is called at most once for each pair (R_i, R_j) .

Finally, one might ask whether the default rule can be triggered. Proposition 1 shows that this can be achieved with one SAT call.

Proposition 1 (Default Rule Application). *The default rule is triggered iff $\mathcal{B} \wedge \neg L_1 \dots \wedge \neg L_r$ is satisfiable.*

Sketch. No rule fires on a solution to $\mathcal{B} \wedge \neg L_1 \dots \wedge \neg L_r$. □ □

4.2 Redundancy

In order to study rule and literal redundancy, we provide a formal definition of decision sets equivalence. We denote by $S_{\mathcal{M}}(o) = \cup_{R_i \in \Delta(o)} \{x \in \mathbb{F} \mid x \models \mathcal{B} \wedge L_i\}$.

Definition 3 (Decision Set Equivalence). *Let \mathcal{M}_1 and \mathcal{M}_2 be two decision sets defined over the same feature space \mathbb{F} and output \mathcal{V} and having the same default rule. \mathcal{M}_1 is equivalent to \mathcal{M}_2 iff $\forall o \in \mathcal{V}, S_{\mathcal{M}_1}(o) = S_{\mathcal{M}_2}(o)$.*

The following lemma is an immediate consequence of Definition 3.

Lemma 2 (Lemma Decision Set Equivalence). *Let \mathcal{M}_1 and \mathcal{M}_2 be two equivalent decision sets that exhibit no negative overlap and let \mathcal{B} be a background knowledge. Then $\forall x \models \mathcal{B}, \kappa_{\mathcal{M}_1}(x) = \kappa_{\mathcal{M}_2}(x)$.*

We introduce the notion of rule redundancy to capture the fact that removing a given rule from a decision set leads to an equivalent decision set.

Definition 4 (Rule Redundancy). *A rule R_i is redundant in \mathcal{M} iff $\mathcal{M} \setminus R_i$ is equivalent to \mathcal{M}*

Let $G_i = \Delta(o_i) \setminus \{R_i\} = \{R_{i_1}, \dots, R_{i_z}\}$ where $R_{i_m} = (L_{i_m}, o_{i_m})$.

Proposition 2 (Rule Redundancy Check). *A rule R_i is redundant in \mathcal{M} iff $\mathcal{B} \wedge L_i \models L_{i_1} \vee \dots \vee L_{i_z}$.*

Proof. Let $\mathcal{M}^* = \mathcal{M} \setminus R_i$. Clearly R_i is redundant in \mathcal{M} iff $S_{\mathcal{M}}(o_i) = S_{\mathcal{M}^*}(o_i)$. In other words, iff $\cup_{R_j \in \Delta(o_i)} \{x \in \mathbb{F} \mid x \models \mathcal{B} \wedge L_j\} = \cup_{R_j \in \Delta(o_i) \setminus R_i} \{x \in \mathbb{F} \mid x \models \mathcal{B} \wedge L_j\}$. The latter is true iff $\mathcal{B} \wedge L_i \models L_{i_1} \vee \dots \vee L_{i_z}$. \square \square

Following Proposition 2, one can check if a rule is redundant with one SAT oracle since $\mathcal{B} \wedge L_i \models L_{i_1} \vee \dots \vee L_{i_z}$ iff $\mathcal{B} \wedge L_i \wedge \neg L_{i_1} \wedge \dots \wedge \neg L_{i_z}$ is unsatisfiable. For instance, in Example 1, this allows to show that R_3 is redundant.

One can also build an equivalent decision set with no redundant rules by checking and removing redundant rules iteratively. Note that the order in which the redundant rules are removed matters as it might return a different decision set at each execution.

We assume in the rest of this section that no rule is redundant. Suppose that L_i contains at least two literals and that $l \in L_i$. We denote by $\mathcal{M}_l^i = \mathcal{M} \cup (L_i \setminus l, o_i) \setminus R_i$ the decision set identical to \mathcal{M} except that l is removed from L_i . We give a formal definition of literal redundancy.

Definition 5 (Literal Redundancy). *A literal l is redundant in L_i iff $l \in L_i$ and \mathcal{M}_l^i is equivalent to \mathcal{M} .*

Informally speaking, a literal is redundant in L_i iff its removal from L_i leads to an equivalent decision set. In the following we prove that there are only two cases of redundancies that we call local and global redundancies, and we show sufficient and necessary conditions to find (and remove) them. When using L_i , we suppose that it contains at least two literals.

We denote by $L_i^l = L_i \cup \{\neg l\} \setminus \{l\}$. We define the following sets to address literal redundancy: $\Omega_i = \cup_{R_j \in G_i} \{x \in \mathbb{F} \mid x \models \mathcal{B} \wedge L_j\}$, $\Theta_i = \{x \in \mathbb{F} \mid x \models \mathcal{B} \wedge L_i \setminus \{l\}\}$, $\Xi_i = \{x \in \mathbb{F} \mid x \models \mathcal{B} \wedge L_i\}$, $\Upsilon_i = \{x \in \mathbb{F} \mid x \models \mathcal{B} \wedge L_i^l\}$. By construction, we have:

- $\Theta_i = \Xi_i \cup \Upsilon_i$
- $S_{\mathcal{M}}(o_i) = \Omega_i \cup \Xi_i$
- $S_{\mathcal{M}_l^i}(o_i) = \Omega_i \cup \Theta_i = \Omega_i \cup \Xi_i \cup \Upsilon_i$

Proposition 3 (Literal Redundancy (1)). *A literal l is redundant in L_i iff $l \in L_i$ and $\Omega_i \cup \Xi_i = \Omega_i \cup \Theta_i = \Omega_i \cup \Xi_i \cup \Upsilon_i$*

Proof. Observe first that \mathcal{M}_l^i is equivalent to \mathcal{M} iff $S_{\mathcal{M}}(o_i) = S_{\mathcal{M}_l^i}(o_i)$. Therefore, l is redundant in L_i iff $\Omega_i \cup \Xi_i = \Omega_i \cup \Theta_i = \Omega_i \cup \Xi_i \cup \Upsilon_i$. \square \square

Lemma 3 (Local Redundancy). *If $l \in L_i$ and $\mathcal{B} \wedge L_i \setminus \{l\} \models l$ then l is redundant in L_i . This is called local redundancy.*

Proof. If $\mathcal{B} \wedge L_i \setminus \{l\} \models l$ then $\Xi_i = \Theta_i$ and thus $S_{\mathcal{M}_l^i}(o_i) = \Omega_i \cup \Theta_i = \Omega_i \cup \Xi_i = S_{\mathcal{M}}(o_i)$. Therefore, by Proposition 3, l is redundant in L_i . \square \square

In Example 1, (*age* > 10) is locally redundant in R_1 .

Recall that $G_i = \Delta(o_i) \setminus \{R_i\} = \{R_{i_1}, \dots, R_{i_z}\}$ and $L_i^{\bar{l}} = L_i \cup \{\neg l\} \setminus \{l\}$.

Lemma 4 (Global Redundancy). *If l is not locally redundant in L_i and $\mathcal{B} \wedge L_i^{\bar{l}} \models L_{i_1} \vee \dots \vee L_{i_z}$, then l is redundant in L_i . This is called global redundancy.*

Proof. If $\mathcal{B} \wedge L_i^{\bar{l}} \models L_{i_1} \vee \dots \vee L_{i_z}$ then $\Upsilon_i \subseteq \Omega_i$. Thus, since $\Theta_i = \Xi_i \cup \Upsilon_i$, we have $S_{\mathcal{M}_l^i}(o_i) = \Omega_i \cup \Theta_i = \Omega_i \cup \Xi_i \cup \Upsilon_i = \Omega_i \cup \Xi_i = S_{\mathcal{M}}(o_i)$. Therefore, by Proposition 3, l is redundant in L_i . \square \square

In Example 1, (*size* \neq 140) is globally redundant in R_1 .

Theorem 1 (Literal Redundancy (2)). *A literal $l \in L_i$ is redundant iff it is locally redundant or globally redundant.*

Proof. \implies : If l is redundant, then by Proposition 3 we have $\Omega_i \cup \Xi_i = \Omega_i \cup \Xi_i \cup \Upsilon_i$. Observe that $\Upsilon_i \cap \Xi_i = \emptyset$. This is because if $x \in \Upsilon_i \cap \Xi_i$, then $x \models \mathcal{B} \wedge L_i \wedge L_i^{\bar{l}}$ which is false because $L_i \wedge L_i^{\bar{l}}$ contains l and $\neg l$. Therefore, there are only two cases for $\Omega_i \cup \Xi_i = \Omega_i \cup \Xi_i \cup \Upsilon_i$ to hold. Either $\Upsilon_i = \emptyset$ or $\Upsilon_i \neq \emptyset$ and $\Upsilon_i \subseteq \Omega_i$. The first case is true iff $\mathcal{B} \wedge L_i \setminus \{l\} \models l$, that is, l is locally redundant. The second case is true iff $\mathcal{B} \wedge L_i^{\bar{l}} \models L_{i_1} \vee \dots \vee L_{i_z}$, that is, l is globally redundant \Leftarrow : trivial. \square \square

Corollary 1 (Assessing Literal Redundancy). *A literal $l \in L_i$ is redundant iff one of the following conditions holds:*

1. **Local redundancy:**

$$\mathcal{B} \wedge (L_i \setminus \{l\}) \wedge \neg l \text{ is unsatisfiable.}$$

2. **Global redundancy:** (1) does not hold, and

$$\mathcal{B} \wedge L_i^{\bar{l}} \wedge \neg L_{i_1} \wedge \dots \wedge \neg L_{i_z} \text{ is unsatisfiable.}$$

Proof. Immediate from Theorem 1 and Lemmas 3 and 4. \square \square

Corollary 1 can be used to iteratively remove redundant literals, thus building decision sets with no rules/literal redundancies. It should be noted that different removal orders might lead to different decision sets.

Example 2. Suppose that $\mathcal{B} = (b \vee w) \wedge (\neg d \vee f)$ and $\mathcal{M} = \{R_1, R_2, R_3\}$ where $R_1 : (L_1 = a \wedge b, o_1)$, $R_2 : (L_2 = a \wedge w, o_1)$, and $R_3 : (L_3 = c \wedge d \wedge f, o_2)$.

- $\mathcal{B} \wedge L_3 \setminus \{f\} \models f$. Therefore, f is locally redundant in L_3 .
- $L_1^{\bar{b}} = a \wedge \neg b$, and $\mathcal{B} \wedge L_1^{\bar{b}} = (b \vee w) \wedge a \wedge \neg b \equiv a \wedge \neg b \wedge w$. Thus $\mathcal{B} \wedge L_1^{\bar{b}} \models R_2$ and therefore b is globally redundant in L_1 .

4.2.1 Relation to Abductive Explanations.

Proposition 4. Suppose that $L_k \subseteq \{x_i = v_i^j \mid i \in [1, m], v_i^j \in D_i\}$. If $R_k = (L_k, o_k)$ fires on x , and there is no negative overlap involving R_k , then the features used in L_k represent a WAXp.

Proof. By construction. \square

Proposition 5. Suppose that $L_k \subseteq \{x_i = v_i^j \mid i \in [1, m], v_i^j \in D_i\}$. If $R_k = (L_k, o_k)$ fires on v , there is no negative overlap, and L_k contains no (global or local) redundant literal, then the features from L_k represent a AXp. \square

Proof. Suppose by contradiction that the features from L_k do not define an AXp. Then there is a literal $l \in L_k$ such that $\forall x \in \mathbb{F}$ such that $x \models \mathcal{B}$, if $x \models L_k \setminus \{l\}$, then $\kappa_{\mathcal{M}(x)} = o_k$. In this case, \mathcal{M}_l^k (i.e., the decision set identical to \mathcal{M} except for L_k which is replaced with $L_k \setminus \{l\}$) is equivalent to \mathcal{M} . Then, by Definition 5, l is redundant, thus the contradiction. \square

Observe that, if the conditions of Proposition 5 hold, then the literals in L_k represent an AXp, and so can be identified manually by a human decision-maker. Otherwise, as proved in earlier work for the concrete case of decision lists [Marques-Silva and Ignatiev, 2023], finding an AXp is computationally hard.

5 Experiments

We evaluate the different desired properties on different use cases including decision sets, decision trees, and anchor explanations. All SAT calls are performed using the PySAT toolkit³ with its default configuration [Ignatiev et al., 2024, Ignatiev et al., 2018]. All experiments run on AppleM1 Pro that has 32G memory and 8 cores.

Prediction Models & Datasets. In order to make our evaluation as broad and as unbiased as possible, we selected datasets from the UCI machine learning repository⁴ with the parameters: $\mathcal{P} = (Task, Min, Max, Nb, Types)$ on each use case (whenever relevant) where:

- $Task \subseteq \{classification, regression\}$ is the prediction task
- Min (respectively Max) is the minimum (respectively *maximum*) size of the dataset.
- Nb is the minimum number of inputs of each class present in the dataset in case of classification.

³<https://pysathq.github.io/>.

⁴<https://archive.ics.uci.edu>

- $Types \subseteq \{numerical, binary\}$ is the type of features.

We describe the different prediction models along with their tailored setting.

- **Orange (v3)**⁵: a library to learn decision sets for classification. The datasets are selected using the parameters $\mathcal{P} = (\{classification\}, 100, 10^6, 20, \{binary, numerical\})$.
- **Boomer [Rapp et al., 2020]**⁶: A library for learning gradient boosted multi-label classification rules. We use the default Boomer datasets⁷.
- **scikit-learn (v1.6.1)**⁸ and **Interpretable AI (IAI)**⁹ to learn decision trees (DTs) for classification and regression. scikit-learn learns trees in a greedy way with no guarantee of optimality whereas IAI learns optimal decision trees. The parameters used for the datasets are $\mathcal{P} = (\{classification, regression\}, 100, 4 * 10^6, 20, \{binary, numerical\})$.

Background Knowledge. In our empirical study, the Boolean variables that are used in the different decision sets represent a domain relation of the form $(x_f \bowtie v_f)$ where $\bowtie \in \{=, >, \geq, \leq, <\}$ for some $f \in \mathbb{F}$ and $v_f \in D_f$. We implemented a general purpose procedure to generate a background knowledge \mathcal{B} for each use case that maintains domain coherence. For instance, if $(length > 30)$ and $(length = 17)$ appear in a decision set, then \mathcal{B} contains the clause $\neg(length = 17) \vee \neg(length > 30)$.

Given a set of rules, for each feature f , we first compute a list, called Val_f , that contains all distinct values from the domain of f that are used in the decision set (or Anchor explanations). Val_f is increasingly ordered if the values are numerical. We also collect the set of unary relations used for f , denoted by $Relations_f$, which can be any subset of $\{=, >, \geq, \leq, <\}$. The background knowledge \mathcal{B} is then constructed using Algorithm 2 as a CNF. Note that we need only the three operators $>, \geq$, and $=$, since $(x > v)$ is equivalent to $(x \geq v - 1)$ and $(x < v)$ is equivalent to $(x \leq v - 1)$. Algorithm 2 follows standard procedures for encoding finite domains into SAT [Ohrimenko et al., 2009].

Learning Setting. For Orange, scikit-learn, and IAI, a grid search is used to select the best values for the maximum rule length, the minimum covered examples per rule, among others. Each dataset used with Orange, scikit-learn, and IAI is split into 80% for training and 20% for testing. Boomer’s learning parameters are the default ones except for the maximum number of rules that we fix to 100 with one label classification. The detailed grid search parameters are given in Table 1. Cross validation is performed with 5 folds for all experiments using stratified sampling and each execution is randomly repeated 4 times.

Experimental Pipeline. All decision sets that have only one output are discarded. For each decision set, we first remove duplicate rules and rules that never fire. After this preprocessing, we run Algorithm 1 to find all overlap. Next, we look for redundant rules then remove them. Finally, we compute

⁵<https://orangedatamining.com>

⁶<https://github.com/mrapp-ke/MLRL-Boomer>

⁷<https://github.com/mrapp-ke/Boomer-Datasets>

⁸<https://scikit-learn.org/stable/>

⁹<https://www.interpretable.ai/>

Algorithm 2 Domain Constraints As a Background Knowledge

```
1: Function: Build  $\mathcal{B}$ 
2: Input:  $Relations_1, \dots, Relations_m, Val_1, \dots, Val_m$ 
3: Output: A background knowlegde  $\mathcal{B}$  as a CNF
4:  $\mathcal{B} = \emptyset$ 
5: for  $f \in \{1, \dots, m\}$  do
6:   if  $'=' \in Relations_f$  then
7:      $\mathcal{B} \leftarrow \mathcal{B} \cup \{(f = Val_f[i]) \implies \neg(f = Val_f[j]) \mid i < j \in [1, |Val_f|]\}$ 
8:   end if
9:   if  $'>' \in Relations_f$  then
10:     $\mathcal{B} \leftarrow \mathcal{B} \cup \{(f > Val_f[i+1]) \implies (f > Val_f[i]) \mid i \in [1, |Val_f| - 1]\}$ 
11:   end if
12:   if  $'\geq' \in Relations_f$  then
13:     $\mathcal{B} \leftarrow \mathcal{B} \cup \{(f \geq Val_f[i+1]) \implies (f \geq Val_f[i]) \mid i \in [1, |Val_f| - 1] : \}$ 
14:   end if
15:   if  $\{'=', '\geq'\} \subseteq Relations_f$  then
16:     $\mathcal{B} \leftarrow \mathcal{B} \cup \{(x = Val_f[i]) \implies (x \geq Val_f[i]) \mid i \in [1, |Val_f|]\}$ 
17:     $\mathcal{B} \leftarrow \mathcal{B} \cup \{(x = Val_f[i]) \implies \neg(x \geq Val_f[i+1]) \mid i \in [1, |Val_f| - 1]\}$ 
18:   end if
19:   if  $\{'=', '>'\} \subseteq Relations_f$  then
20:     $\mathcal{B} \leftarrow \mathcal{B} \cup \{(x = Val_f[i]) \implies \neg(x > Val_f[i]) \mid i \in [1, |Val_f|]\}$ 
21:     $\mathcal{B} \leftarrow \mathcal{B} \cup \{(x = Val_f[i+1]) \implies (x > Val_f[i]) \mid i \in [1, |Val_f| - 1]\}$ 
22:   end if
23:   if  $\{'\geq', '>'\} \subseteq Relations_f$  then
24:     $\mathcal{B} \leftarrow \mathcal{B} \cup \{(x > Val_f[i]) \implies (x \geq Val_f[i]) \mid i \in [1, |Val_f|]\}$ 
25:   end if
26:   if  $\{'=', '\geq', '>'\} \subseteq Relations_f$  then
27:     $\mathcal{B} \leftarrow \mathcal{B} \cup \{(x \geq Val_f[i]) \implies (x = Val_f[i]) \vee (x > Val_f[i]) \mid i \in [1, |Val_f|]\}$ 
28:   end if
29: end for
30: Return  $\mathcal{B}$ 
```

all locally/globally redundant literals. We use a timeout of one hour on each decision set to find all pairs negative overlap and rule/literal redundancies.

Decision Set Statistics.

- Train: Training accuracy (classification) or Training MSE (regression)
- Test: Testing accuracy (classification) or Training MSE (regression)
- NR: Number of rules
- NP: Cardinality of $\{o_i \mid \mathcal{M} = \cup_i (L_i, o_i)\}$
- TO: CPU time (s) to find all negative overlap
- TB: CPU time (s) to generate the background knowledge
- TC: CPU time (s) to find all redundant rules
- TR: CPU time (s) to find all redundant literals
- BS: Size of the background knowledge
- RS: Sum of the sizes of the rules

Table 1: Grid Search Parameters

	Orange	Sklearn Class.	Sklearn Reg.	IAI Class.	IAI Reg.
Beam Width	10,30	-	-	-	-
Min Covered	5,15	-	-	-	-
Max Rule Length	3,5	-	-	-	-
Criterion	-	gini, entropy	sqr err, fried mse	-	mse
Max Depth	-	3,5,7,9	3,5,7,9,11	3,5,7	3,5,7
Min Sample Leaf	-	5,15,25	5,15,25	-	-
Min Bucket	-	- -	-	5,15	5,15

- RM: Maximum rule size
- NO: Number of negative overlap
- $PO = \frac{NO}{Total}$: Percentage of negative overlap where Total is the total number of pairs of rules associated to different predictions

Model Statistics. We report for each prediction model the following:

- DS: The total number of decision sets
- EX: The total number of decision sets that timed out
- IR: Number of instances that admit at least one redundant rule
- IL: Number of instances that admit at least one locally redundant literal
- PL: Percentage of locally redundant literals for instances that admit at least one locally redundant literal
- IG: Number of instances that admit at least one globally redundant literal
- PG: Percentage of globally redundant literals for instances that admit at least one globally redundant literal

In the rest of the section, we focus on the most important observations.

Summary. Tables 2 and Tables 3 give the full statistics for each learning model¹⁰. Instance-related statistics are averaged for each model. Decision sets that are worse than random guess are ignored. Instance statistics are averaged for each prediction model. Only the results of the experiments that did not reach the timeout are reported. The time to generate the background knowledge (TB) is often less than a second. The time to find redundant rules (TC) is often few seconds, except for some decision sets where it took about a minute. The runtime to find all literal redundancies (TR), however, is much longer. To observe this more accurately, we present in Figure 1 its box plot across all models. The x-axis is the time in seconds and the y-axis is the TR value for each model. This is expected because every literal is checked for redundancy by application of Corollary 1.

5.1 Rule and Literal Redundancy

We are interested in this section in the evaluation of the presence of redundant rules and locally/globally redundant literals, their correlations with other characteristics, as well as the efficiency of our approach.

¹⁰The detailed results can be found at <https://siala.github.io/data/2025-ecml/>

Table 2: Summary of the Results (1)

	DS	EX	NR	NP	TO	TB	TC
sklearn classification	196	21	35	3	0	0	0
sklearn regression	28	8	70	69	0	0	8
IAI classification	177	0	17	4	0	0	0
IAI regression	28	0	56	53	0	0	3
Orange	127	12	175	2	76	0	2
Boomer	180	16	97	49	0	0	0

Table 3: Summary of the Results (2)

	TR	BS	RS	RM	IR	IL	PL	IG	PG
sklearn classification	94	23	233	5	0	123	7	175	33
sklearn regression	879	85	541	7	0	20	18	0	0
IAI classification	19	13	96	4	0	82	3	135	10
IAI regression	333	77	367	5	0	25	15	2	0
Orange	10	12139	405	3	16	1	0	13	0
Boomer	21	30	180	11	42	6	0	20	0

Redundancy. We note first that rule redundancy does not occur often as we can see in column IR in Table 3 except for Boomer. Figure 2 represents a box plot of the percentage of local (respectively global) redundancies PL (respectively PG) for all learning models. Orange and Boomer barely exhibit literal redundancies (see columns IL and IG in Table 3). Regression models did not show any global literal redundancy except for 2 cases with IAI regression trees. This is expected because for a literal to be globally redundant, there should be at least two rules predicting the exact same value, which is rare in regression. Classification trees, however, exhibit a noticeable presence of global redundancy (‘PG IAI classification’ and ‘PG sklearn classification’). Figure 2 shows a significant presence of local redundancy in all tree models. We note that for each prediction task (regression, classification), IAI trees have fewer local/global redundancies than sklearn trees (in terms of the median and the maximum values). This suggests that optimal trees tend to reduce redundancy.

Correlations. We looked into different correlations between local/global redundancy and other statistics. We report the results only for models where at least 30% of its decision trees/sets exhibit local/global redundancy. There was a moderate negative correlation of global redundancy with the number of prediction outcomes (i.e., size of \mathcal{V}) with scikit-learn and IAI classification trees. Figure 3 shows the most important correlations of local redundancy with the statistics mentioned earlier. For instance, on the x-axis, with NR we show the correlation of the local redundancy values found by each model with the number of rules. Clearly local redundancy with scikit-learn regression trees highly correlates with NR, NP, BS, RM. IAI regression trees has the same tendency.

Figure 1: Box Plot of TR.

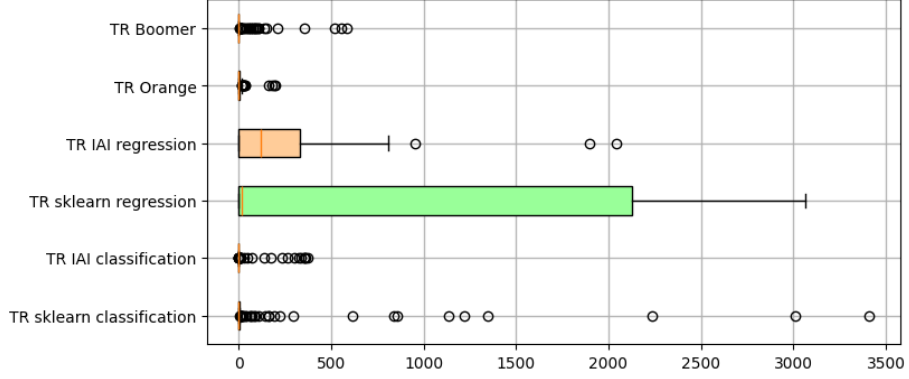
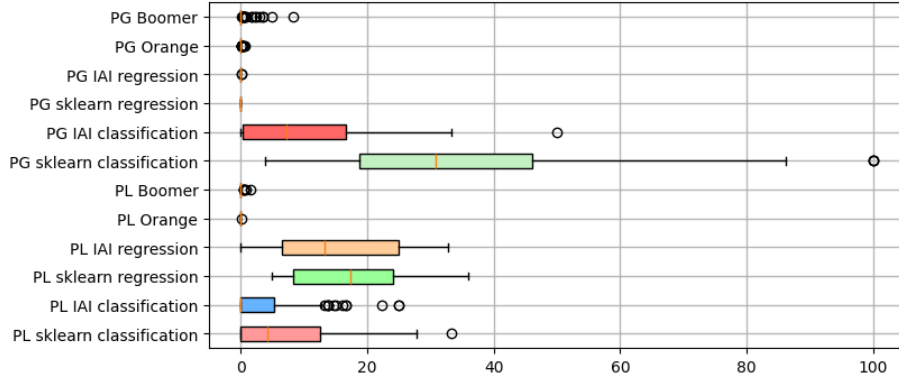


Figure 2: Local/Global Redundancies. The X-axis represents the values of PL/PG



5.2 Negative Overlap

We evaluate the presence of negative overlap on Orange and Boomer and their relationship with relevant statistics. Boomer timed out on 4 datasets (emotions, image, scene, yeast) after the one hour time limit. The results are summarized in Table 4 for instances that did not timeout. The most important observation is the high percentage of negative overlap (column PO). Indeed, with Boomer decision sets, almost every pair of rules with different predictions overlap. Such an observation is worth reporting to the user. The results are less spectacular for orange with an average close to 50% but still worth noting. The runtime to find all negative overlap per instance is not negligible.

Negative Overlap in Boomer. As the results in this section confirm (see Table 4), Boomer [Rapp et al., 2020] exhibits extensive negative overlap. This is to be expected. In contrast with the approach outlined in this paper, where neg-

Figure 3: Pearson Correlations of Local Redundancies (PL)

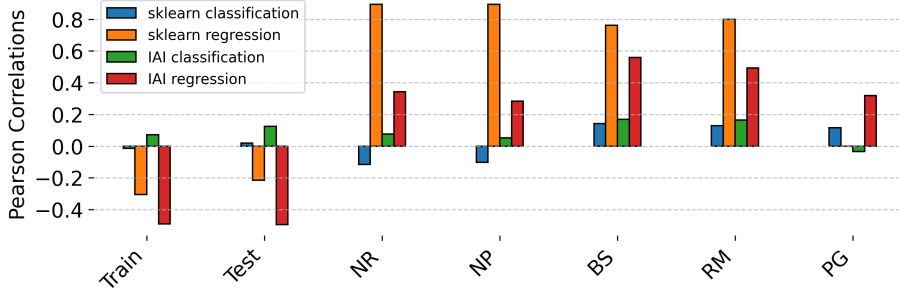


Table 4: Summary of the Negative Overlap Results

	DS	EX	Train	Test	NR	NP	IR	TO	TR	BS	RS	RM	PO
Orange	127	12	70	71	175	2	16	76	10	12139	405	3	50
Boomer	180	16	95	95	97	49	42	0	21	30	180	11	99

ative overlap is targeted as a reason for non-interpretability, Boomer exploits boosting (and as a result negative overlap) to build high-accuracy rule ensembles. The theoretical and practical advantages of boosting are well-known [Freund, 1995, Collins et al., 2002], namely to allow the learning of strong classifiers. As argued in this paper, a downside of negative overlap (and so of rule ensembles) is that finding explanations becomes a computationally-hard challenge. Our experiments are reported for completeness, and confirm the previous remarks.

5.3 Application to Anchor Explanations

Anchors are well-known model-agnostic explanations representing local, “sufficient” conditions for predictions [Ribeiro et al., 2018]. The question we ask here addresses precisely one of the open questions in [Ribeiro et al., 2018]: How to find potentially conflicting anchors? To answer this question, we generate anchors for different inputs, then apply our approach to find negative overlap between anchors.

We reproduced the exact experiments in [Ribeiro et al., 2018] with the three datasets: *adult* for predicting whether a person makes $> 50K$ annually; *rcdv* for predicting recidivism for individuals released from prison; and *lending* for predicting whether a loan on the Lending Club website will turn out bad. For each dataset, four models are used for prediction: boosted trees with xgboost, random forest, logistic regression, and neural networks. Each model is built using the exact configuration in the original paper [Ribeiro et al., 2018]. For each dataset and each model, we generate all anchors of the validation set and look for all negative overlap.

Table 5 presents the results for each dataset and each model. As we can see, negative overlap in Anchor explanations is present in all use cases. Often,

Table 5: Anchor Experiments

Learner	Dataset	Train	Test	NR	TO	NO	PO	RM
xgboost	recidivism	92.39	74.33	333	0	87	0.31	17
randomforest	recidivism	93.52	75.46	321	0	65	0.25	17
logistic	recidivism	62.59	60.00	196	0	735	7.81	12
nn	recidivism	87.47	71.49	341	1	150	0.52	17
xgboost	lending	90.10	82.89	260	0	384	2.47	15
randomforest	lending	91.25	83.60	278	0	207	1.18	15
logistic	lending	82.56	83.51	50	0	54	9.38	14
nn	lending	88.00	82.54	159	0	66	1.07	16
xgboost	adult	90.35	84.26	565	8	3195	4.03	14
randomforest	adult	93.52	85.60	558	7	2534	3.27	13
logistic	adult	83.00	82.98	378	3	2788	7.86	13
nn	adult	92.47	83.62	597	11	3212	3.61	14

anchors of random forests exhibit the lowest percentage of negative overlap, whereas those of logistic regression have the highest percentage. We also observe that the best (and respectively, worst) models in terms of prediction quality tend to have the lowest (respectively, highest) percentages of negative overlap. These observations suggest that the quality of Anchor explanations depends on the prediction quality of the learner/model.

6 Conclusions

This paper investigates the occurrence of negative facets of decision sets, namely negative overlap and (global or local) literal redundancy. Dedicated algorithms for their identification are proposed. Furthermore, the paper reveals the tight relationship between decision sets for which manual explanations can be devised, and the non-existence of the aforementioned negative facets. A first set of experiments confirms that these negative facets occur ubiquitously in existing implementations of decision sets, thus rendering unrealistic the manual identification of explanations. A second set of experiments confirms that the explanations obtained with the well-known explainer Anchors will also exhibit the same negative facets.

Acknowledgements. Mohamed Siala would like to thank INSA Toulouse for funding his research visit to the University of Lleida. This work was supported in part by the MCIN/AEI/10.13039/501100011033/FEDER, UE under the project PID2022-139835NB-C22. This work was supported in part by the Spanish Government under grant PID 2023-152814OB-I00. The authors at University of Lleida would like to thank the Catalan Government for the quality accreditation given to their research group GREiA (2021 SGR 1615).

References

- [Adadi and Berrada, 2018] Adadi, A. and Berrada, M. (2018). Peeking inside the black-box: A survey on explainable artificial intelligence (XAI). *IEEE Access*, 6:52138–52160.
- [Atzmueller et al., 2024] Atzmueller, M., Fürnkranz, J., Kliegr, T., and Schmid, U. (2024). Explainable and interpretable machine learning and data mining. *Data Min. Knowl. Discov.*, 38(5):2571–2595.
- [Audemard et al., 2024] Audemard, G., Lagniez, J., Marquis, P., and Szczepanski, N. (2024). Deriving provably correct explanations for decision trees: The impact of domain theories. In *IJCAI*, pages 3688–3696.
- [Biere et al., 2021] Biere, A., Heule, M., van Maaren, H., and Walsh, T., editors (2021). *Handbook of Satisfiability - Second Edition*, volume 336 of *Frontiers in Artificial Intelligence and Applications*. IOS Press.
- [Breiman, 2001] Breiman, L. (2001). Statistical modeling: The two cultures. *Statistical science*, 16(3):199–231.
- [Breiman et al., 1984] Breiman, L., Friedman, J. H., Olshen, R. A., and Stone, C. J. (1984). *Classification and Regression Trees*. Wadsworth.
- [Carvalho et al., 2019] Carvalho, D. V., Pereira, E. M., and Cardoso, J. S. (2019). Machine learning interpretability: A survey on methods and metrics. *Electronics*, 8(8):832.
- [Clark and Niblett, 1989] Clark, P. and Niblett, T. (1989). The CN2 induction algorithm. *Machine Learning*, 3(4):261–283.
- [Collins et al., 2002] Collins, M., Schapire, R. E., and Singer, Y. (2002). Logistic regression, AdaBoost and Bregman distances. *Mach. Learn.*, 48(1-3):253–285.
- [Darwiche, 2023] Darwiche, A. (2023). Logic for explainable AI. In *LICS*, pages 1–11.
- [Demirovic et al., 2022] Demirovic, E., Lukina, A., Hebrard, E., Chan, J., Bailey, J., Leckie, C., Ramamohanarao, K., and Stuckey, P. J. (2022). MurTree: optimal decision trees via dynamic programming and search. *J. Mach. Learn. Res.*, 23:26:1–26:47.
- [Demsar et al., 2013] Demsar, J. et al. (2013). Orange: data mining toolbox in python. *J. Mach. Learn. Res.*, 14(1):2349–2353.
- [Dwivedi et al., 2023] Dwivedi, R. et al. (2023). Explainable AI (XAI): core ideas, techniques, and solutions. *ACM Comput. Surv.*, 55(9):194:1–194:33.
- [Florio et al., 2023] Florio, A. M., Martins, P., Schiffer, M., Serra, T., and Vidal, T. (2023). Optimal decision diagrams for classification. pages 7577–7585. AAAI Press.

- [Freitas, 2014] Freitas, A. A. (2014). Comprehensible classification models: a position paper. *SIGKDD Explor. Newsl.*, 15(1):1–10.
- [Freund, 1995] Freund, Y. (1995). Boosting a weak learning algorithm by majority. *Inf. Comput.*, 121(2):256–285.
- [Fürnkranz et al., 2012] Fürnkranz, J., Gamberger, D., and Lavrac, N. (2012). *Foundations of Rule Learning*. Cognitive Technologies. Springer.
- [Fürnkranz and Kliegr, 2015] Fürnkranz, J. and Kliegr, T. (2015). A brief overview of rule learning. In *RuleML*, pages 54–69.
- [Ghosh et al., 2022] Ghosh, B., Malioutov, D., and Meel, K. S. (2022). Efficient learning of interpretable classification rules. *J. Artif. Intell. Res.*, 74:1823–1863.
- [Gorji and Rubin, 2022] Gorji, N. and Rubin, S. (2022). Sufficient reasons for classifier decisions in the presence of domain constraints. In *AAAI*, pages 5660–5667.
- [Guidotti et al., 2019] Guidotti, R., Monreale, A., Ruggieri, S., Turini, F., Giannotti, F., and Pedreschi, D. (2019). A survey of methods for explaining black box models. *ACM Comput. Surv.*, 51(5):93:1–93:42.
- [Hu et al., 2022] Hu, H., Huguet, M., and Siala, M. (2022). Optimizing binary decision diagrams with maxsat for classification. In *AAAI*, pages 3767–3775. AAAI Press.
- [Hu et al., 2020] Hu, H., Siala, M., Hebrard, E., and Huguet, M. (2020). Learning optimal decision trees with maxsat and its integration in adaboost. In *IJCAI*, pages 1170–1176.
- [Hüllermeier et al., 2020] Hüllermeier, E., Fürnkranz, J., Mencía, E. L., Nguyen, V., and Rapp, M. (2020). Rule-based multi-label classification: Challenges and opportunities. In *RuleML*, pages 3–19.
- [Huysmans et al., 2011] Huysmans, J., Dejaeger, K., Mues, C., Vanthienen, J., and Baesens, B. (2011). An empirical evaluation of the comprehensibility of decision table, tree and rule based predictive models. *Decision Support Systems*, 51(1):141–154.
- [Hyafil and Rivest, 1976] Hyafil, L. and Rivest, R. L. (1976). Constructing optimal binary decision trees is NP-complete. *Inf. Process. Lett.*, 5(1):15–17.
- [Ignatiev et al., 2018] Ignatiev, A., Morgado, A., and Marques-Silva, J. (2018). PySAT: A Python toolkit for prototyping with SAT oracles. In *SAT*, pages 428–437.
- [Ignatiev et al., 2024] Ignatiev, A., Tan, Z. L., and Karamanos, C. (2024). Towards universally accessible SAT technology. In *SAT*, pages 4:1–4:11.

- [Izza et al., 2022] Izza, Y., Ignatiev, A., and Marques-Silva, J. (2022). On tackling explanation redundancy in decision trees. *J. Artif. Intell. Res.*, 75:261–321.
- [Lakkaraju et al., 2016] Lakkaraju, H., Bach, S. H., and Leskovec, J. (2016). Interpretable decision sets: A joint framework for description and prediction. In *KDD*, pages 1675–1684.
- [Lipton, 2018] Lipton, Z. C. (2018). The mythos of model interpretability. *Commun. ACM*, 61(10):36–43.
- [Lundberg and Lee, 2017] Lundberg, S. M. and Lee, S. (2017). A unified approach to interpreting model predictions. In *NeurIPS*, pages 4765–4774.
- [Marques-Silva and Ignatiev, 2023] Marques-Silva, J. and Ignatiev, A. (2023). No silver bullet: interpretable ML models must be explained. *Frontiers Artif. Intell.*, 6.
- [Minh et al., 2022] Minh, D., Wang, H. X., Li, Y. F., and Nguyen, T. N. (2022). Explainable artificial intelligence: a comprehensive review. *Artif. Intell. Rev.*, 55(5):3503–3568.
- [Ohrimenko et al., 2009] Ohrimenko, O., Stuckey, P. J., and Codish, M. (2009). Propagation via lazy clause generation. *Constraints An Int. J.*, 14(3):357–391.
- [Pedregosa et al., 2011] Pedregosa, F. et al. (2011). Scikit-learn: Machine learning in python. *J. Mach. Learn. Res.*, 12:2825–2830.
- [Rapp et al., 2024] Rapp, M., Fürnkranz, J., and Hüllermeier, E. (2024). On the efficient implementation of classification rule learning. *Adv. Data Anal. Classif.*, 18(4):851–892.
- [Rapp et al., 2020] Rapp, M., Mencía, E. L., Fürnkranz, J., Nguyen, V., and Hüllermeier, E. (2020). Learning gradient boosted multi-label classification rules. In *ECML*, pages 124–140.
- [Ribeiro et al., 2016] Ribeiro, M. T., Singh, S., and Guestrin, C. (2016). "why should I trust you?": Explaining the predictions of any classifier. In *KDD*, pages 1135–1144.
- [Ribeiro et al., 2018] Ribeiro, M. T., Singh, S., and Guestrin, C. (2018). Anchors: High-precision model-agnostic explanations. In *AAAI*, pages 1527–1535.
- [Rivest, 1987] Rivest, R. L. (1987). Learning decision lists. *Mach. Learn.*, 2(3):229–246.
- [Rudin, 2019] Rudin, C. (2019). Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature machine intelligence*, 1(5):206–215.

- [Rudin et al., 2022] Rudin, C., Chen, C., Chen, Z., Huang, H., Semenova, L., and Zhong, C. (2022). Interpretable machine learning: Fundamental principles and 10 grand challenges. *Statistic Surveys*, 16:1–85.
- [Schwalbe and Finzel, 2024] Schwalbe, G. and Finzel, B. (2024). A comprehensive taxonomy for explainable artificial intelligence: a systematic survey of surveys on methods and concepts. *Data Min. Knowl. Discov.*, 38(5):3043–3101.
- [Shwayder, 1971] Shwayder, K. (1971). Conversion of limited-entry decision tables to computer programs - A proposed modification to Pollack’s algorithm. *Commun. ACM*, 14(2):69–73.
- [Yu et al., 2021] Yu, J., Ignatiev, A., Stuckey, P. J., and Bodic, P. L. (2021). Learning optimal decision sets and lists with SAT. *J. Artif. Intell. Res.*, 72:1251–1279.