
Efficient Explanations With Relevant Sets

Yacine Izza
 Université de Toulouse, Toulouse, France
 yacine.izza@univ-toulouse.fr

Alexey Ignatiev
 Monash University, Melbourne, Australia
 alexey.ignatiev@monash.edu

Nina Narodytska
 VMware Research, CA, USA
 nnarodytska@vmware.com

Martin C. Cooper
 Université Paul Sabatier, IRIT, Toulouse, France
 martin.cooper@irit.fr

Joao Marques-Silva
 IRIT, CNRS, Toulouse, France
 joao.marques-silva@irit.fr

Abstract

Recent work proposed δ -relevant inputs (or sets) as a probabilistic explanation for the predictions made by a classifier on a given input. δ -relevant sets are significant because they serve to relate (model-agnostic) Anchors with (model-accurate) PI-explanations, among other explanation approaches. Unfortunately, the computation of smallest size δ -relevant sets is complete for NP^{PP} , rendering their computation largely infeasible in practice. This paper investigates solutions for tackling the practical limitations of δ -relevant sets. First, the paper alternatively considers the computation of subset-minimal sets. Second, the paper studies concrete families of classifiers, including decision trees among others. For these cases, the paper shows that the computation of subset-minimal δ -relevant sets is in NP, and can be solved with a polynomial number of calls to an NP oracle. The experimental evaluation compares the proposed approach with heuristic explainers for the concrete case of the classifiers studied in the paper, and confirms the advantage of the proposed solution over the state of the art.

1 Introduction

Recent work proposed δ -relevant inputs (or sets) [42], which represent probabilistic explanations for the predictions made by a classifier given some input. δ -relevant sets were shown to generalize both Anchors [33] and PI-explanations [36], thus revealing a connection between model-agnostic explanations (e.g. Anchors) and model-accurate explanations (e.g. PI-explanations). Moreover, δ -relevant sets offer a natural solution for increasing the interpretability of PI-explanations, at the cost of obtaining intuitive explanations that hold in most, but not all, points in feature space. A formidable downside of δ -relevant sets is that their computation is hard for NP^{PP} . This signifies that for most practical examples, the time for computing minimum δ -relevant sets will be prohibitive in practice. To address the computational complexity of finding minimum δ -relevant sets, a number of solutions can be envisioned. A first solution is the approximate computation of δ -relevant sets. However, for this solution, the formal guarantees offered by δ -relevant sets may no longer hold. A second solution is to identify which ML models allow for the efficient computation of δ -relevant sets. Finally, a third solution is to investigate possible ways of relaxing the original definition of δ -relevant sets [42].

This paper addresses the second and third solutions listed above. First, the paper proposes alternative definitions of δ -relevant sets. Second, the paper analyzes the computation of (different variants of) δ -relevant sets in the case of decision trees (DTs).

Preprint. Under review.

Although generally regarded as interpretable [7, 34, 25], recent work showed that DTs can exhibit *explanation redundancy* [3, 17], i.e. there exist DTs containing paths that are (possibly arbitrarily) longer than a PI-explanation [36]. Furthermore, existing experimental evidence confirms that explanation redundancy is observed in DTs learned with state of the art DT learners [17]. Thus, even in the case of DTs, the computation of δ -relevant sets is of interest when the goal is to improve the interpretability of ML models.

The main results of the paper can thus be summarized as follows. First, the paper shows that, for the decision version of computing a minimum size δ -relevant set (i.e. the problem studied in recent work [42]), is in NP in the case of DTs. The proof of this result offers a solution for computing a minimum-size δ -relevant set, in the case of DTs, by using maximum satisfiability modulo theories (MaxSMT) [4, 6]. Second, the paper shows that, in the case of DTs, a relaxed definition of δ -relevant set enables the computation of (relaxed) subset-minimal δ -relevant sets in polynomial time. Third, the paper shows that ML models based on knowledge compilation (KC) languages [11], including those studied in recent papers on XAI for KC languages [36, 37, 10, 2, 1], the computation of (relaxed) subset-minimal δ -relevant sets is also in polynomial time. Fourth, the paper shows that recently proposed duality results for explanations [15, 13], which in practice enable the enumeration of explanations, can be extended to the setting of δ -relevant sets.

Related work. The growing adoption of ML in different settings motivates the recent interest in explainability [27, 12, 35, 20, 43, 26]. Well-known approaches for explaining ML-models are model-agnostic and based on heuristic solutions [32, 21, 33]. These approaches offer no formal guarantees of rigor, and practical limitations have been reported in recent years [8, 28, 38, 19]. More recently, model-accurate non-heuristic approaches to explainability have been investigated [36, 15, 10, 13, 2, 1, 22]. These non-heuristic approaches are characterized by formal guarantees of rigor, e.g. explanations are valid in any point in feature space. Unfortunately, non-heuristic methods also exhibit a number of drawbacks, including for example scalability, explanation size, and the inability to compute explanations with probabilistic guarantees. Recent work [42] revealed ways of relating heuristic and non-heuristic explanations. Our paper builds on this recent work.

Organization. The paper is organized as follows. Section 2 introduces the notation and definitions used in the rest of the paper. Section 3 discusses δ -relevant sets and a number of alternative definitions. Section 4 delves into duality between different kinds of explanations when δ -relevant sets are considered. Section 5 discusses the computation of cardinality-minimal and subset-minimal δ -relevant sets in the case of decision trees. Section 6 presents experimental results for computing δ -relevant sets in the case of DTs. Finally, Section 7 concludes the paper.

2 Preliminaries

Classification problems & formal explanations. This paper considers classification problems, which are defined on a set of features (or attributes) $\mathcal{F} = \{1, \dots, m\}$ and a set of classes $\mathcal{K} = \{c_1, c_2, \dots, c_K\}$. Each feature $i \in \mathcal{F}$ takes values from a domain \mathbb{D}_i . In general, domains can be boolean, integer or real-valued. Feature space is defined as $\mathbb{F} = \mathbb{D}_1 \times \mathbb{D}_2 \times \dots \times \mathbb{D}_m = \{0, 1\}^m$. For boolean domains, $\mathbb{D}_i = \{0, 1\} = \mathbb{B}$, $i = 1, \dots, m$, and $\mathbb{F} = \mathbb{B}^m$. The notation $\mathbf{x} = (x_1, \dots, x_m)$ denotes an arbitrary point in feature space, where each x_i is a variable taking values from \mathbb{D}_i . The set of variables associated with features is $X = \{x_1, \dots, x_m\}$. Moreover, the notation $\mathbf{v} = (v_1, \dots, v_m)$ represents a specific point in feature space, where each v_i is a constant representing one concrete value from $\mathbb{D}_i = \{0, 1\}$. An *instance* (or example) denotes a pair (\mathbf{v}, c) , where $\mathbf{v} \in \mathbb{F}$ and $c \in \mathcal{K}$. (We also use the term *instance* to refer to \mathbf{v} , leaving c implicit.) An ML classifier \mathbb{M} is characterized by a *classification function* κ that maps feature space \mathbb{F} into the set of classes \mathcal{K} , i.e. $\kappa : \mathbb{F} \rightarrow \mathcal{K}$.

We now define formal explanations. Prime implicant (PI) explanations [36] denote a minimal set of literals (relating a feature value x_i and a constant $v_i \in \mathbb{D}_i$ that are sufficient for the prediction¹. Formally, given $\mathbf{v} = (v_1, \dots, v_m) \in \mathbb{F}$ with $\kappa(\mathbf{v}) = c$, a PI-explanation (AXp) is any minimal subset $\mathcal{X} \subseteq \mathcal{F}$ such that,

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

¹PI-explanations are related with abduction, and so are also referred to as abductive explanations (AXp) [14]. More recently, PI-explanations have been studied from a knowledge compilation perspective [2, 1], but also in terms of their computational complexity [3].

AXps can be viewed as answering a ‘Why?’ question, i.e. why is some prediction made given some point in feature space. A different view of explanations is a contrastive explanation [24], which answers a ‘Why Not?’ question, i.e. which features can be changed to change the prediction. A formal definition of contrastive explanation is proposed in recent work [13]. Given $\mathbf{v} = (v_1, \dots, v_m) \in \mathbb{F}$ with $\kappa(\mathbf{v}) = c$, a CXp is any minimal subset $\mathcal{Y} \subseteq \mathcal{F}$ such that,

$$\exists(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \in \mathcal{F} \setminus \mathcal{Y}} (x_j = v_j) \wedge (\kappa(\mathbf{x}) \neq c) \quad (2)$$

Building on the results of R. Reiter in model-based diagnosis [31], [13] proves a minimal hitting set (MHS, or hypergraph transversal [5]) duality relation between AXps and CXps, i.e. AXps are MHSes of CXps and vice-versa. Throughout the paper, (M)HS(\mathbb{Z}) denote the set of (minimal) hitting sets of \mathbb{Z} .

Relevant sets. δ -relevant sets were recently proposed [42] as a formalization of explanation that enables relating different types of explanation [42]. We briefly overview the definition of relevant set and associated definitions. The assumptions regarding the probabilities of logical propositions are those made in earlier work [42]. Let $\Pr_{\mathbf{x}}(A(\mathbf{x}))$ denote the probability of some proposition A defined on the vector of variables $\mathbf{x} = (x_1, \dots, x_m)$, i.e.

$$\Pr_{\mathbf{x}}(A(\mathbf{x})) = \frac{|\{\mathbf{x} \in \mathbb{F} : A(\mathbf{x})=1\}|}{|\{\mathbf{x} \in \mathbb{F}\}|}, \quad \Pr_{\mathbf{x}}(A(\mathbf{x}) \mid B(\mathbf{x})) = \frac{|\{\mathbf{x} \in \mathbb{F} : A(\mathbf{x})=1 \wedge B(\mathbf{x})=1\}|}{|\{\mathbf{x} \in \mathbb{F} : B(\mathbf{x})=1\}|} \quad (3)$$

Definition 1 (δ -relevant set [42]). Let $\kappa : \mathbb{B}^m \rightarrow \mathcal{K} = \mathbb{B}$, $\mathbf{v} \in \mathbb{B}^m$, $\kappa(\mathbf{v}) = c \in \mathbb{B}$, and $\delta \in [0, 1]$. $S \subseteq \mathcal{F}$ is a δ -relevant set for κ and \mathbf{v} if,

$$\Pr_{\mathbf{x}}(\kappa(\mathbf{x}) = c \mid \mathbf{x}_S = \mathbf{v}_S) \geq \delta \quad (4)$$

(where the restriction of \mathbf{x} to the variables with indices in S is represented by $\mathbf{x}_S = (x_i)_{i \in S}$).

(Observe that $\Pr_{\mathbf{x}}(\kappa(\mathbf{x}) = c \mid \mathbf{x}_S = \mathbf{v}_S)$ is often referred to as the *precision* of S [33, 28].) Thus, a δ -relevant set represents a set of features which, if fixed to some pre-defined value (taken from a reference vector \mathbf{v}) ensure that the probability of the prediction being the same as the one for \mathbf{v} is no less than δ .

Definition 2 (Min- δ -relevant set). Given κ , $\mathbf{v} \in \mathbb{B}^m$, and $\delta \in [0, 1]$, find the smallest k , such that there exists $S \subseteq \mathcal{F}$, with $|S| \leq k$, and S is a δ -relevant set for κ and \mathbf{v} .

With the goal of proving the computational complexity of finding a minimum-size set of features that is a δ -relevant set, earlier work [42] restricted the definition to the case where κ is represented as a boolean circuit. (Boolean circuits were restricted to propositional formulas defined using the operators \vee , \wedge and \neg , and using a set of variables representing the inputs; this explains the choice of *inputs* over *sets* in earlier work [42].) Observe that Definition 2 imposes no such restriction on the representation of the classifier, i.e. the logical representation of κ need not be a boolean circuit.

Decision trees. A decision tree \mathcal{T} is a directed acyclic graph having at most one path between every pair of nodes. \mathcal{T} has a root node, characterized by having no incoming edges. All other nodes have one incoming edge. We consider univariate decision trees where each non-terminal node is associated with a single feature x_i . Each edge is labeled with a literal, relating a feature (associated with the edge’s starting node) with some values (or range of values) from the feature’s domain. We will consider literals to be of the form $x_i \in \mathbb{E}_i$. x_i is a variable that denotes the value taken by feature i , whereas $\mathbb{E}_i \subseteq \mathbb{D}_i$ is a subset of the domain of feature i . The type of literals used to label the edges of a DT allows the representation of the DTs generated by a wide range of decision tree learners (e.g. [41]). It is assumed that for any $\mathbf{v} \in \mathbb{F}$ there exists *exactly* one path in \mathcal{T} that is consistent with \mathbf{v} . By *consistent* we mean that the literals associated with the path are satisfied (or consistent) with the feature values in \mathbf{v} .

Running example. Throughout the paper, we will consider the decision tree shown in Figure 1 as the running example².

Example 1. We consider the example DT from Figure 1. For this example and for simplicity, all features are binary with $\mathbb{D}_i = \{0, 1\}$. It is also assumed that $\Pr(x_i = 0) = \Pr(x_i = 1) = 1/2$, which we represent respectively by α and β , to allow other values to be considered. . Some of the paths

²Although the running example considers boolean features (with $\mathbb{D}_i = \{0, 1\}$) and boolean classification, similar conclusions would be obtained if we were to consider instead real-valued features, e.g. having $\mathbb{D}_i = [0, 1]$.

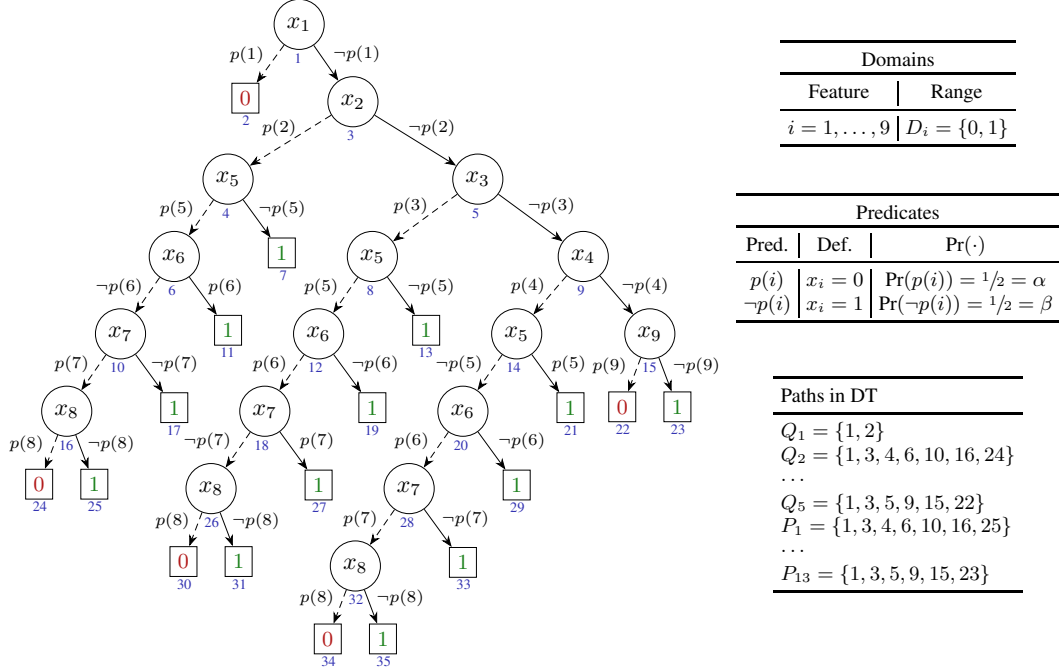


Figure 1: DT used as running example

in the DT are also shown. Moreover, let the instance be $\mathbf{v} = (v_1, v_2, v_3, v_4, v_5, v_6, v_7, v_8, v_9) = (1, 1, 1, 1, 0, 0, 0, 0, 1)$ with prediction $c = 1$. Since \mathbf{v} is consistent with the path ending at node 23, by inspection, we can conclude that a possible explanation is $\mathcal{X} = \{1, 2, 3, 4, 9\}$, i.e. the features listed in the path. It can be shown that this corresponds to a PI-explanation.

3 Complementary Definitions of Relevant Sets

Given the prohibitive complexity of solving the Min- δ -Relevant-Set problem, this section proposes alternative definitions of minimal relevant sets, which are shown to yield efficient algorithms for some concrete ML models. First, we consider subset-minimal relevant sets instead of cardinality-minimal sets. However, we relax the restrictions that features are boolean and the classification problem is restricted to two classes.

Min-C δ -Relevant-Sets. Following earlier work on PI-explanations [36], we consider subset-minimal relevant sets.

Definition 3 (C δ -relevant set). Let $\kappa : \mathbb{F} \rightarrow \mathcal{K}$, $\delta \in [0, 1]$, and instance $(\mathbf{v} \in \mathbb{F}, c \in \mathcal{K})$. $\mathcal{S} \subseteq \mathcal{F}$ is a C δ -relevant set for the classifier-instance pair, κ and (\mathbf{v}, c) , if (4) holds.

(The difference of C δ to plain δ relevant sets is that \mathbb{F} and \mathcal{K} become unrestricted.) As noted in earlier work, a (smallest) PI-explanation is a 1-relevant set for a given pair κ and (\mathbf{v}, c) . Furthermore, the main difference with respect to Anchors [33] is the assumptions made with respect to sampling. As noted in earlier work [42], δ -relevant sets can be related with different efforts for computing explanations [33, 36, 18].

Definition 4 (Min-C δ -Relevant-Set). Let $\kappa : \mathbb{F} \rightarrow \mathcal{K}$, $\delta \in [0, 1]$, and instance $(\mathbf{v} \in \mathbb{F}, c \in \mathcal{K})$. A Min-C δ -Relevant-Set is a (subset-)minimal subset $\mathcal{S} \subseteq \mathcal{F}$ that is C δ -relevant for κ and (\mathbf{v}, c) .

(Observe that, in contrast with the definition of Min- δ -Relevant-Set [42], where the objective is to compute a cardinality-minimal set, the objective of the definition of Min-C δ -Relevant-Set it to compute a subset-minimal set.) For the case where κ is implemented as a boolean circuit (propositional formula defined on the operators \vee , \wedge and \neg), Min- δ -Relevant-Set is hard for NP^{PP} , with the decision problem in NP^{PP} [42]. Although the complexity of Min-C δ -Relevant-Set is unknown, we conjecture that it is similar to the one of Min- δ -Relevant-Set. Moreover, we have the following result,

Proposition 1. Deciding whether a set $\mathcal{S} \in \mathcal{F}$ is a C δ -relevant set is PP-hard.

(The proof is included in the supplementary materials.) It should be underlined that the high complexity of exactly solving Min- δ -Relevant-Set (and Min-C δ -Relevant-Set) in the general case represents a key practical limitation. One additional difficulty with computing a subset-minimal C δ -relevant set is that its definition is non-monotone. (4) can be written as follows,

$$\Pr_{\mathbf{x}}(\kappa(\mathbf{x}) = c \mid \mathbf{x}_{\mathcal{S}} = \mathbf{v}_{\mathcal{S}}) = \frac{|\{\mathbf{x} \in \mathbb{F} : \kappa(\mathbf{x}) = c \wedge (\mathbf{x}_{\mathcal{S}} = \mathbf{v}_{\mathcal{S}})\}|}{|\{\mathbf{x} \in \mathbb{F} : (\mathbf{x}_{\mathcal{S}} = \mathbf{v}_{\mathcal{S}})\}|}$$

As the size of set \mathcal{S} is reduced (e.g. as we search for a minimal set), both the numerator and the denominator can change. Hence, the value of $\Pr_{\mathbf{x}}(\kappa(\mathbf{x}) = c \mid \mathbf{x}_{\mathcal{S}} = \mathbf{v}_{\mathcal{S}})$ is not guaranteed to shrink, and in some settings this value can grow.

Min-I δ -Relevant-Sets. We show later in the paper that, by considering the probability of the conjunction of two events instead of the conditional probability, the resulting monotone definition of relevant set enables more efficient computation of subset-minimal relevant sets. One has four possible outcomes when considering two events. In our case that means we can have: $[\kappa(\mathbf{x}) = \kappa(\mathbf{v}), \mathbf{x}_{\mathcal{S}} = \mathbf{v}_{\mathcal{S}}]$, $[\kappa(\mathbf{x}) = \kappa(\mathbf{v}), \mathbf{x}_{\mathcal{S}} \neq \mathbf{v}_{\mathcal{S}}]$, $[\kappa(\mathbf{x}) \neq \kappa(\mathbf{v}), \mathbf{x}_{\mathcal{S}} = \mathbf{v}_{\mathcal{S}}]$, and $[\kappa(\mathbf{x}) \neq \kappa(\mathbf{v}), \mathbf{x}_{\mathcal{S}} \neq \mathbf{v}_{\mathcal{S}}]$. We are interested in picking sets \mathcal{S} that minimize the odds of picking an assignment consistent with \mathcal{S} and obtaining a different prediction. Hence, our concern will be to identify sets \mathcal{S} that *minimize* $\Pr(\kappa(\mathbf{x}) \neq \kappa(\mathbf{v}), \mathbf{x}_{\mathcal{S}} = \mathbf{v}_{\mathcal{S}})$.

Definition 5 (I δ -relevant set). Let $\kappa : \mathbb{F} \rightarrow \mathcal{K}$, $\delta \in [0, 1]$, and instance $(\mathbf{v} \in \mathbb{F}, c \in \mathcal{K})$. $\mathcal{S} \subseteq \mathcal{F}$ is a I δ -relevant set for the classifier-instance pair, κ and (\mathbf{v}, c) , if (5) holds:

$$\Pr_{\mathbf{x}}(\kappa(\mathbf{x}) \neq \kappa(\mathbf{v}), \mathbf{x}_{\mathcal{S}} = \mathbf{v}_{\mathcal{S}}) \leq \delta \quad (5)$$

From the definition of conditional probability (see above in this section), it is immediate to observe that,

$$\Pr_{\mathbf{x}}(\kappa(\mathbf{x}) \neq \kappa(\mathbf{v}), \mathbf{x}_{\mathcal{S}} = \mathbf{v}_{\mathcal{S}}) = \frac{|\{\mathbf{x} \in \mathbb{F} : \kappa(\mathbf{x}) \neq c \wedge (\mathbf{x}_{\mathcal{S}} = \mathbf{v}_{\mathcal{S}})\}|}{|\{\mathbf{x} \in \mathbb{F}\}|}$$

Definition 6 (Min-I δ -Relevant-Set). Let $\kappa : \mathbb{F} \rightarrow \mathcal{K}$, $\delta \in [0, 1]$, and instance $(\mathbf{v} \in \mathbb{F}, c \in \mathcal{K})$. A Min-I δ -Relevant-Set is a minimal subset $\mathcal{S} \subseteq \mathcal{F}$ that is I δ -relevant for κ and (\mathbf{v}, c) .

By observing that for larger sets we can only increase the likelihood of the function differing from the value of $\kappa(\mathbf{v})$, we have the following result.

Proposition 2. I δ -relevant sets are monotone, i.e. for $\mathcal{S}_1 \supseteq \mathcal{S}_2$, it is the case that,

$$\Pr_{\mathbf{x}}(\kappa(\mathbf{x}) \neq \kappa(\mathbf{v}), \mathbf{x}_{\mathcal{S}_1} = \mathbf{v}_{\mathcal{S}_1}) \leq \Pr_{\mathbf{x}}(\kappa(\mathbf{x}) \neq \kappa(\mathbf{v}), \mathbf{x}_{\mathcal{S}_2} = \mathbf{v}_{\mathcal{S}_2})$$

4 Duality Results for Relevant Sets

Duality results between different types of explanation enable the implementation of explanation enumeration algorithms [15, 13]³ This section proves one initial duality result between δ -relevant sets. Given earlier work [15, 13], additional results can be envisioned.

Let C be a predicate, $C : 2^{\mathbb{F}} \rightarrow \{0, 1\}$, such that,

$$C(\mathcal{S}) = [\Pr_{\mathbf{x}}(\kappa(\mathbf{x}) \neq c, \mathbf{x}_{\mathcal{S}} = \mathbf{v}_{\mathcal{S}}) \leq \delta]$$

We associate with C a set of subsets of \mathbb{F} ,

$$\mathbb{C} = \{\mathcal{S} \subseteq \mathbb{F} \mid C(\mathcal{S})\}$$

In addition, we define a set of minimal sets,

$$\mathbb{C}_{\min} = \{\mathcal{S} \subseteq \mathbb{F} \mid C(\mathcal{S}) \wedge \forall (\mathcal{U} \subsetneq \mathcal{S}). \neg C(\mathcal{U})\}$$

Next, we introduce the dual predicate D , $D : 2^{\mathbb{F}} \rightarrow \{0, 1\}$, such that,

$$D(\mathcal{T}) = \neg C(\mathcal{F} \setminus \mathcal{T}) = [\Pr_{\mathbf{x}}(\kappa(\mathbf{x}) \neq c, \mathbf{x}_{\mathcal{F} \setminus \mathcal{T}} = \mathbf{v}_{\mathcal{F} \setminus \mathcal{T}}) > \delta]$$

The dual of δ -relevant sets are sets \mathcal{T} of features which if changed entail a change of class with a probability $> \delta$ and are thus probabilistic analogues of contrastive explanations [13]. As done above,

³These recent duality results about explanations build on the work of Reiter [31]. In this section, we follow loosely a recent overview [39].

Path R_j	Q_1	Q_2	Q_3	Q_4	Q_5	P_1	P_2	P_3	P_4	P_5	P_6	P_7	P_8	P_9	P_{10}	P_{11}	P_{12}	P_{13}
$\Pr(R_j)$	α^1	$\alpha^4\beta^2$	$\alpha^4\beta^3$	$\alpha^4\beta^4$	$\alpha^1\beta^4$	$\alpha^3\beta^3$	$\alpha^2\beta^3$	$\alpha^3\beta^1$	$\alpha^1\beta^2$	$\alpha^4\beta^3$	$\alpha^4\beta^2$	$\alpha^3\beta^2$	$\alpha^1\beta^3$	$\alpha^3\beta^5$	$\alpha^2\beta^4$	$\alpha^1\beta^5$	$\alpha^2\beta^3$	β^5

Table 1: Path probabilities for running example

we can define the following sets:

$$\begin{aligned}\mathbb{D} &= \{\mathcal{T} \subseteq \mathbb{F} \mid D(\mathcal{T})\} \\ \mathbb{D}_{\min} &= \{\mathcal{T} \subseteq \mathbb{F} \mid D(\mathcal{T}) \wedge \forall (\mathcal{V} \subsetneq \mathcal{T}). \neg D(\mathcal{V})\}\end{aligned}$$

Given the above, monotonicity of the predicates C and D (see Proposition 2), allows us to prove the following results,

Proposition 3. $\mathbb{C} = \text{HS}(\mathbb{D})$, $\mathbb{D} = \text{HS}(\mathbb{C})$, $\mathbb{C}_{\min} = \text{MHS}(\mathbb{D})$, and $\mathbb{D}_{\min} = \text{MHS}(\mathbb{C})$.

Proof. First, we consider $\mathbb{C} = \text{HS}(\mathbb{D})$. Suppose, there exists \mathcal{S} such that it is not a hitting set of sets in \mathbb{D} . Namely, \mathcal{S} does not hit some set $\mathcal{T} \in \mathbb{D}$, $\mathcal{S} \cap \mathcal{T} = \emptyset$. By definition, $\mathcal{S} \subset \mathcal{F} \setminus \mathcal{T}$. We recall that a predicate P is *monotone* if for all $S, S' \subseteq \mathcal{F}$,

$$S \subset S' \wedge P(S) \rightarrow P(S').$$

Hence, as $\mathcal{S} \subset \mathcal{F} \setminus \mathcal{T}$ and $C(\mathcal{S})$ holds, $C(\mathcal{F} \setminus \mathcal{T})$ must hold. This leads to a contradiction as $D(\mathcal{T}) = \neg C(\mathcal{F} \setminus \mathcal{T})$ by definition. The reverse direction, $\mathbb{D} = \text{HS}(\mathbb{C})$, is similar.

Second, we consider $\mathbb{C}_{\min} = \text{MHS}(\mathbb{D})$. The proof that $\mathcal{S} \in \mathbb{C}_{\min}$ is a hitting set of \mathbb{D} follows from the argument above as $\mathbb{C}_{\min} \subseteq \mathbb{C}$. Next, we suppose $\mathcal{S} \in \mathbb{C}_{\min}$ is *not a minimal* hitting set of \mathbb{D} . Let $\mathcal{G} \subset \mathcal{S}$ be the minimal hitting set. By definition of minimality of \mathcal{S} , $\neg C(\mathcal{G})$ must hold. Consider \mathcal{W} such that $\mathcal{G} = \mathcal{F} \setminus \mathcal{W}$. We have that $\neg C(\mathcal{G}) = \neg C(\mathcal{F} \setminus \mathcal{W}) = D(\mathcal{W})$. Therefore, $\mathcal{W} \in \mathbb{D}$. As $\mathcal{G} \cap \mathcal{W} = \emptyset$ by construction, \mathcal{G} does not hit $\mathcal{W} \in \mathbb{D}$ and it is not a hitting set. The reverse direction, $\mathbb{D}_{\min} = \text{MHS}(\mathbb{C})$, is similar. \square

5 Relevant Sets for DTs & Other Classifiers

This section shows that the decision problem for δ -relevant (and so $C\delta$ -relevant) sets is in NP when κ is represented by a decision tree⁴. Thus, Min- $C\delta$ -Relevant-Set can be solved with at most a logarithmic number of calls to an NP oracle. (This is true since we minimize on the number of features.) This section also shows the decision problem for $I\delta$ -relevant sets is in P. Thus, the Min- $I\delta$ -Relevant-Set can be solved in polynomial time in the case of DTs. Furthermore, the section extends the previous results to the case of knowledge compilation (KC) languages [11].

Path probabilities for DTs. Let $\mathbf{v} \in \mathbb{F}$ and suppose that $\kappa(\mathbf{v}) = c \in \mathcal{K}$. For a DT \mathcal{T} , let $\mathcal{P} = \{P_1, \dots, P_M\}$ denote the paths with prediction c , and let $\mathcal{Q} = \{Q_1, \dots, Q_N\}$ denote the paths with a prediction in $\mathcal{K} \setminus \{c\}$. Let $\mathcal{R} = \mathcal{P} \cup \mathcal{Q}$ denote the set of all paths in the DT \mathcal{T} . For $R_j \in \mathcal{R}$, let $\|R_j\|$ denote the number of points in \mathbb{F} consistent with R_j . Thus, the path probability of any path $R_j \in \mathcal{R}$ is, $\Pr(R_j) = \|R_j\| / \|\mathbb{F}\|$. (The path probability of some tree path R_j is the empirical probability of a point in feature space chosen at random being consistent with the path R_j .) As a result, we get,

$$\sum_{R_j \in \mathcal{P}} \Pr(R_j) + \sum_{R_j \in \mathcal{Q}} \Pr(R_j) = 1$$

Moreover, $\|\mathbb{F}\| = \|\mathbb{D}_1\| \times \dots \times \|\mathbb{D}_m\|$. For each path R_j , let d_i denote the number of values in \mathbb{D}_i that is consistent with the literals defined on x_i in path R_j . Thus, $\|R_j\| = d_1 \times \dots \times d_m$.

Example 2. For the example in Figure 1, Table 1 shows the path probability for each path in the DT, computed using the above definition of path probability.

Min- $C\delta$ -Relevant-Sets for DTs. Whereas in the general case, deciding whether there exists a δ -relevant set of size no greater than k is complete for NP^{PP} [42], in the the case of DTs, one can prove that this decision problem is in NP (and the same applies in the case of a subset-minimal $C\delta$ -relevant set).

Proposition 4. For DTs, given $\mathbf{v} \in \mathbb{F}$, with $\kappa(\mathbf{v}) = c \in \mathcal{K}$, deciding the existence of min- δ -relevant set of size at most k is in NP.

⁴For simplicity, the paper considers the case of non-continuous features. However, in the case of DTs, the results generalize to continuous features.

Proof. We reduce the problem of deciding the existence of a min- δ -relevant set of size at most k to the decision version of the maximum satisfiability modulo theories (SMT) problem [4, 6] (assuming a suitable quantifier-free theory).

Let s_i be a boolean variable such that $s_i = 1$ iff $i \in \mathcal{F}$ is included in the δ -relevant set. Moreover, let t_j be a boolean variable, such that $t_j = 1$ iff path $R_j \in \mathcal{P} \cup \mathcal{Q}$ is inconsistent, i.e. some feature i added to the δ -relevant set makes R_j inconsistent. Thus, if path R_j is inconsistent with the value given to feature i , then it must be the case that,

$$s_i \rightarrow t_j$$

Also, if R_j is deemed inconsistent, then it must be the case that,

$$t_j \rightarrow \bigvee_{i \in I_j} s_i$$

where i ranges over the set of features I_j that make R_j inconsistent, given \mathbf{v} .

The set of picked features \mathcal{S} , (i.e. \mathcal{S} is the set of features having $s_i = 1$), ensures that

$$\Pr_{\mathbf{x}}(\kappa(\mathbf{x}) = \kappa(\mathbf{v}) | \mathbf{x}_{\mathcal{S}} = \mathbf{v}_{\mathcal{S}}) \geq \delta$$

From which we get,

$$\begin{aligned} \Pr_{\mathbf{x}}(\kappa(\mathbf{x}) = \kappa(\mathbf{v}) \wedge \mathbf{x}_{\mathcal{S}} = \mathbf{v}_{\mathcal{S}}) &\geq \delta \times \Pr_{\mathbf{x}}(\mathbf{x}_{\mathcal{S}} = \mathbf{v}_{\mathcal{S}}) \Leftrightarrow \\ \sum_{j, R_j \in \mathcal{P}} \neg t_j \times \Pr(R_j) &\geq \delta \times \sum_{j, R_j \in \mathcal{P} \cup \mathcal{Q}} \neg t_j \times \Pr(R_j) \end{aligned}$$

which is a linear inequality on the t_j variables, since the probabilities are constant. An additional constraint is that the number of s_i variables assigned value 1 cannot exceed k , i.e. the bound on the size of the relevant set \mathcal{S} :

$$\sum_{i \in \mathcal{F}} s_i \leq k$$

which is another linear inequality, this one on the s_i variables. By conjoining all the constraints, and assignment to the s_i and t_j variables that satisfies the constraints picks a δ -relevant set whose size does not exceed k . \square

Clearly, since the decision problem is in NP, it is immediate how to compute a cardinality minimal δ -relevant set by binary search on the number of s_i variables included in the set. Since the number of variables equals the size of \mathcal{F} , then we are guaranteed to need (in the worst-case) a logarithmic number of calls to an NP oracle. Furthermore, since the decision problem for the min- δ -relevant problem is in NP, it is also the case that the decision problem for the min-C δ -relevant problem is also in NP. Finally, we conjecture that the min- δ -relevant set, but also the min-C δ -relevant problem are both hard for NP. These conjectures are further justified below.

Min-I δ -Relevant-Sets for DTs. One apparent reason to the conjectured complexity is the fact that the conditional probability used for defining δ -relevant and C δ -relevant sets is non-monotone. As a result, earlier in the paper we introduced I δ -relevant sets, which were shown to be monotone in Proposition 2. We now show that, in the case of DTs, computing a subset-minimal I δ -relevant set is in P. The criterion for a set of features to be I δ -relevant is:

$$\Pr_{\mathbf{x}}(\kappa(\mathbf{x}) \neq \kappa(\mathbf{v}), \mathbf{x}_{\mathcal{S}} = \mathbf{v}_{\mathcal{S}}) \leq \delta$$

Observe that this constraint holds when $\mathcal{S} = \mathcal{F}$ and, by Proposition 2 I δ -relevant sets are monotone. As a result, we can compute a subset minimal I δ -relevant set, as proposed in Algorithm 1. (The algorithm is standard, and can be traced to at least the work of Chinneck & Dravnieks [9]. The novelty is its use for computing min-I δ -relevant sets.) The algorithm maintains an invariant representing the assertion that the set \mathcal{S} is a I δ -relevant set. Initially, all features are included in set \mathcal{S} , i.e. $\mathcal{S} = \mathcal{F}$, and so \mathcal{S} is a I δ -relevant set. The algorithm then iteratively removes one feature at a time, and checks whether the invariant is broken. If it is, then the feature is added back to set \mathcal{S} . Otherwise, we are guaranteed, by monotonicity, that for any superset of set \mathcal{S} , the invariant holds.

Example 3. We consider the running example (see Figure 1, with instance (\mathbf{v}, c) given by $\mathbf{v} = (v_1, v_2, v_3, v_4, v_5, v_6, v_7, v_8, v_9) = (1, 1, 1, 1, 0, 0, 0, 0, 1)$ and $c = \kappa(\mathbf{v}) = 1$. As argued earlier, by setting $\delta = 0$, one AXP is $\mathcal{X} = \{1, 2, 3, 4, 9\}$. Let $\epsilon(\mathcal{S}) = \Pr_{\mathbf{x}}(\kappa(\mathbf{x}) = \kappa(\mathbf{v}), \mathbf{x}_{\mathcal{S}} = \mathbf{x}_{\mathcal{S}})$, denoting the error associated with some set of features \mathcal{S} . With the purpose of improving the interpretability of the explanation, we set $\delta = 0.03$, and work towards finding an explanation with fewer literals. To illustrate the execution of the algorithm, we assume that the features are analyzed in the order $\{1, 2, 3, 4, 9\}$. Table 2 summarizes the execution of the algorithm. The algorithm first analyzes

Algorithm 1 Finding one min- Id -relevant set (IDRS)**Input:** Classifier κ , instance \mathbf{v} , value δ **Output:** IDRS \mathcal{S}

```

1: procedure findIDRS( $\kappa, \mathbf{v}, \delta$ )
2:    $\mathcal{S} \leftarrow \{1, \dots, m\}$ 
3:   for  $i \in \{1, \dots, m\}$  do  $\triangleright$  Invariant:  $\Pr_{\mathbf{x}}(\kappa(\mathbf{x}) \neq \kappa(\mathbf{v}), \mathbf{x}_{\mathcal{S}} = \mathbf{v}_{\mathcal{S}}) \leq \delta$ 
4:      $\mathcal{S} \leftarrow \mathcal{S} \setminus \{i\}$ 
5:     if  $\Pr_{\mathbf{x}}(\kappa(\mathbf{x}) \neq \kappa(\mathbf{v}), \mathbf{x}_{\mathcal{S}} = \mathbf{v}_{\mathcal{S}}) > \delta$  then
6:        $\mathcal{S} \leftarrow \mathcal{S} \cup \{i\}$ 
7:   return  $\mathcal{S}$ 

```

\mathcal{S}	i	$\mathcal{R} = \mathcal{S} \setminus \{i\}$	\mathcal{Q} paths consistent with \mathcal{R}	$\epsilon(\mathcal{R})$	Decision
$\{1, 2, 3, 4, 9\}$	1	$\{2, 3, 4, 9\}$	Q_1	0.5	Keep 1
$\{1, 2, 3, 4, 9\}$	2	$\{1, 3, 4, 9\}$	Q_2	0.0157	Drop 2
$\{1, 3, 4, 9\}$	3	$\{1, 4, 9\}$	Q_2, Q_3	0.0234	Drop 3
$\{1, 4, 9\}$	4	$\{1, 9\}$	Q_2, Q_3, Q_4	0.0273	Drop 4
$\{1, 9\}$	9	$\{1\}$	Q_2, Q_3, Q_4, Q_5	> 0.03	Keep 9

Table 2: Execution of Algorithm 1

dropping feature 1 from the explanation \mathcal{X} . In this case, only path Q_1 can be made consistent. Given that the probability of picking an assignment consistent with Q_1 is 0.5, then feature 1 cannot be dropped from the explanation, as that would put the error about the target threshold. Next, the algorithm considers dropping feature 2 from the explanation. In this case, only path Q_2 can be made consistent. Given that the probability of picking an assignment consistent with Q_2 is $(1/2)^6 = 0.015625$, then are still below the target absolute fraction of error of $\delta = 0.03$. Hence, we remove feature 2 from the explanation. For feature 3, and since feature 2 is already dropped, then both paths Q_2 and Q_3 can be made consistent. In this case, the error raises to 0.0234, but it is still below 0.3, and so feature 3 is also dropped from the explanation. A similar analysis allows concluding that feature 4 can also be dropped from the explanation. In contrast, after removing features 2, 3 and 4, feature 9 cannot be dropped from the explanation. The resulting approximate explanation (i.e. a Id -relevant set) is thus $\{1, 9\}$. Moreover, the explanation $\{1, 9\}$ ensures that, in more than 97% of the points in feature space consistent with the values of features 1 and 9, the prediction will be the intended one, i.e. 1.

KC languages [11]. Knowledge compilation (KC) languages [11] aim at simplifying queries and transformations in knowledge bases, and have recently been used as ML models. Concrete examples include binary decision diagrams [36, 37], among others [2, 3, 1]. By noting that the explanation algorithm proposed for DTs exploits counting of models after conditioning (i.e. fixing to a value) a set of selected features, then we can conclude that, for KC languages that implement conditioning and model counting in polynomial time, one min- Id -relevant set can also be computed in polynomial time.

6 Experimental Results

This section summarizes the experimental results, which aim at demonstrating the efficiency of min- Id -relevant sets if computed as explanations for DT models trained for various well-known datasets, over heuristic explanations of Anchor [33], both in terms of runtime performance and explanation precision.

Implementation and benchmarks. Min- Id -relevant sets are computed following the ideas of Section 5 and utilizing the polynomial-time Algorithm 1. The prototype implementation of the algorithm (IDRS) is written in Perl while the overall experiment is set up and run in Python.⁵ The precision of the resulting explanations is then assessed using the generic (and non-monotone) precision metric of Anchor [33]. The experiments are conducted on a MacBook Pro with a Dual-Core Intel Core i5 2.3GHz CPU with 8GByte RAM running macOS Catalina.

The benchmarks used in the experiment include publicly available and widely used datasets. The datasets originate from UCI ML Repository [40] and Penn ML Benchmarks [29]. The number of

⁵The prototype implementation, benchmarks, instructions and log files of the experiment will be made publicly available in the final version of the paper.

Dataset	#F	#I	δ	IDRS						Anchor							
				Length		Precision (%)		Runtime (s)		Length		Precision (%)		Runtime (s)			
				M	avg	avg	dev	m	M	avg	M	avg	avg	dev	m	M	avg
adult	12	1766	0.0	10	5.1	100	0.0	0.04	0.07	0.05	12	5.3	87.8	16.7	0.14	2.99	1.20
			0.01	6	3.3	85.7	20.8	0.04	0.08	0.04							
			0.02	6	2.8	83.0	16.4	0.04	0.08	0.05							
			0.05	5	1.9	77.7	21.0	0.04	0.11	0.06							
allhyper	29	1113	0.0	7	4.4	100	0.0	0.05	0.05	0.05	29	1.2	89.5	7.0	0.28	5.75	0.35
			0.01	6	3.0	98.4	4.3	0.04	0.08	0.05							
			0.02	6	1.0	97.7	6.3	0.05	0.07	0.05							
			0.05	4	1.0	97.7	6.3	0.04	0.10	0.05							
ann-thyroid	21	2139	0.0	10	3.9	100	0.0	0.08	0.10	0.08	21	1.3	96.4	8.7	0.22	8.63	0.48
			0.01	6	1.4	96.9	11.4	0.07	0.13	0.08							
			0.02	6	1.0	96.8	11.2	0.08	0.12	0.08							
			0.05	5	0.1	95.2	9.9	0.07	0.17	0.10							
fars	29	2790	0.0	15	5.9	100	0.0	0.67	0.92	0.69	29	9.0	73.5	40.3	0.30	57.43	7.54
			0.01	6	2.0	75.2	30.9	0.58	0.81	0.69							
			0.02	6	2.1	70.2	35.5	0.67	0.98	0.71							
			0.05	5	1.7	58.6	38.0	0.63	0.89	0.70							
kddcup	41	4368	0.0	14	11.4	100	0.0	0.44	4.14	0.46	39	2.6	23.1	19.0	0.42	137.3	10.59
			0.01	8	4.4	53.7	42.9	0.42	0.84	0.45							
			0.02	7	4.2	51.8	22.0	0.45	0.61	0.46							
			0.05	6	2.8	38.7	22.0	0.41	0.54	0.44							

Table 3: Assessing explanations of $I\delta$ -relevant sets (IDRS) and comparison with Anchor’s explanations. Columns #F and #I report, resp., the number of features and the number of tested instances, in the dataset. (Note that for a dataset containing less (resp. more) than 10,000 instances, 30% (resp. 3%) of its instances, randomly selected, are used to be explained. Moreover, duplicate rows in the datasets are filtered.) Sub-Columns M and avg of column Length show, resp., the maximum and average size of an explanation. Sub-columns avg and dev of column Precision show, resp., the average and standard deviation of the explanation’s precision. Sub-columns m, M and avg of column Runtime report, resp., minimal, maximal and average time in seconds to compute an explanation.

training instances (resp. features) in the target datasets varies from 3710 to 145585 (resp. 12 to 41). All the decision tree models are trained using the learning tool *ITI (Incremental Tree Induction)* [41, 16]. Note that the accuracy of all the models is above 73%, the maximum depth of the trees varies from 14 to 60 and the total number of nodes varies from 49 to 9969.

The experiment was set to iterate over some of the *unique* (see below) instances of a dataset and to compute an explanation for each such instance: either a min- $I\delta$ -relevant set or an anchor. As the baseline, we ran Anchor with the default explanation precision of 0.95. The prototype implementation IDRS was run for the values of δ from $\{0.05, 0.02, 0.01, 0.0\}$. It should be observed that the proposed experiment gives an advantage to Anchor, as Anchor is allowed to compute explanations guided by its own metric, whereas $I\delta$ -relevant sets *know nothing* about this metric (which they will be assessed with).

Results. Table 3 details the results of our experiment. First of all, observe that $I\delta$ -relevant sets are extremely simple to compute. Concretely, the runtime of our prototype explainer IDRS normally takes just *a fraction of a second* per data instance (and never exceeds a second) to get a subset-minimal $I\delta$ -relevant set. This means that it is at least 1–2 orders of magnitude faster than Anchor, which can take up to 138 seconds to get a single explanation, with the average explanation time being up to 10 seconds. Also observe that the runtime of the proposed approach is not affected by the value of δ and tends to be negligible overall.

Second, length-wise $I\delta$ -relevant sets also outperform Anchor. In particular, it is not surprising that the largest $I\delta$ -relevant sets correspond to $\delta=0$ and these on average include up to 11.4 features. Explanation size gets further improved when δ is 0.01, 0.02 or 0.05. Concretely, it is reduced to *a few* literals per explanation (on average below 5). (Also, please refer to the value of standard deviation shown in the tables.) On the contrary, Anchor’s explanations utilize up to 39 literals, with the average explanation containing 9 literals. These results show an important difference between IDRS and Anchor in terms of interpretability [23].

Finally and somewhat unexpectedly, IDRS outperforms Anchor in terms of explanation precision. Clearly, the precision of $I\delta$ -relevant sets (i.e. $\delta=0$) is 100%, which demonstrates a significant improvement over anchors. What is more interesting, however, is that the average precision of IDRS does not significantly drop down in case of $\delta \in \{0.05, 0.02, 0.01\}$. In particular, its precision is on par with (or better than) the explanations provided by Anchor. All the points above confirm that $I\delta$ -relevant sets if computed for DT models provide a viable alternative to Anchor’s explanations from all the considered perspectives, including runtime performance, explanation size, and precision.

7 Conclusions

δ -relevant sets [42] enable the computation of approximate (i.e. non-universally true) explanations, and reveal connections between heuristic explanations and non-heuristic explanations. A major drawback of δ -relevant sets is their computational complexity. This paper shows that for DTs, deciding whether there exists a set of at most k features that δ -relevant is in NP. Furthermore, the paper proposes relaxed alternative definitions of δ -relevant sets, and shows that such alternative definitions enable the computation of minimal approximate explanations in polynomial time. The paper also derives a first result on the duality between sets of features representing different kinds of (relaxed) δ relevant sets. The experimental results, obtained on large DTs learned with a state of the art tree learner, confirm the practical efficiency and the quality of explanations when compared with the well-known Anchor heuristic explainer [33].

Acknowledgments. This work was supported by the AI Interdisciplinary Institute ANITI, funded by the French program “Investing for the Future – PIA3” under Grant agreement no. ANR-19-PI3A-0004, and by the H2020-ICT38 project COALA “Cognitive Assisted agile manufacturing for a Labor force supported by trustworthy Artificial intelligence”.

References

- [1] G. Audemard, S. Bellart, L. Bounia, F. Koriche, J. Lagniez, and P. Marquis. On the computational intelligibility of boolean classifiers. *CoRR*, abs/2104.06172, 2021.
- [2] G. Audemard, F. Koriche, and P. Marquis. On tractable XAI queries based on compiled representations. In *KR*, pages 838–849, 2020.
- [3] P. Barceló, M. Monet, J. Pérez, and B. Subercaux. Model interpretability through the lens of computational complexity. In *NeurIPS*, 2020.
- [4] C. W. Barrett and C. Tinelli. Satisfiability modulo theories. In E. M. Clarke, T. A. Henzinger, H. Veith, and R. Bloem, editors, *Handbook of Model Checking*, pages 305–343. Springer, 2018.
- [5] C. Berge. *Hypergraphs: combinatorics of finite sets*. Elsevier, 1984.
- [6] N. Bjørner, A. Phan, and L. Fleckenstein. νz - an optimizing SMT solver. In C. Baier and C. Tinelli, editors, *TACAS*, pages 194–199, 2015.
- [7] L. Breiman. Statistical modeling: The two cultures. *Statistical science*, 16(3):199–231, 2001.
- [8] O. Camburu, E. Giunchiglia, J. Foerster, T. Lukasiewicz, and P. Blunsom. Can I trust the explainer? verifying post-hoc explanatory methods. *CoRR*, abs/1910.02065, 2019.
- [9] J. W. Chinneck and E. W. Dravnieks. Locating minimal infeasible constraint sets in linear programs. *INFORMS J. Comput.*, 3(2):157–168, 1991.
- [10] A. Darwiche and A. Hirth. On the reasons behind decisions. In *ECAI*, pages 712–720, 2020.
- [11] A. Darwiche and P. Marquis. A knowledge compilation map. *J. Artif. Intell. Res.*, 17:229–264, 2002.
- [12] R. Guidotti, A. Monreale, S. Ruggieri, F. Turini, F. Giannotti, and D. Pedreschi. A survey of methods for explaining black box models. *ACM Comput. Surv.*, 51(5):93:1–93:42, 2019.
- [13] A. Ignatiev, N. Narodytska, N. Asher, and J. Marques-Silva. From contrastive to abductive explanations and back again. In *AI*IA*, 2020.
- [14] A. Ignatiev, N. Narodytska, and J. Marques-Silva. Abduction-based explanations for machine learning models. In *AAAI*, pages 1511–1519, 2019.
- [15] A. Ignatiev, N. Narodytska, and J. Marques-Silva. On relating explanations and adversarial examples. In *NeurIPS*, pages 15857–15867, 2019.
- [16] Incremental Decision Tree Induction. <https://www-lrn.cs.umass.edu/iti/>, 2020.
- [17] Y. Izza, A. Ignatiev, and J. Marques-Silva. On explaining decision trees. *CoRR*, abs/2010.11034, 2020.
- [18] P. Khosravi, Y. Liang, Y. Choi, and G. Van den Broeck. What to expect of classifiers? reasoning about logistic regression with missing features. In *IJCAI*, pages 2716–2724, 2019.
- [19] H. Lakkaraju and O. Bastani. "how do I fool you?": Manipulating user trust via misleading black box explanations. In *AIES*, pages 79–85, 2020.
- [20] Z. C. Lipton. The mythos of model interpretability. *Commun. ACM*, 61(10):36–43, 2018.

- [21] S. M. Lundberg and S. Lee. A unified approach to interpreting model predictions. In *NIPS*, pages 4765–4774, 2017.
- [22] E. L. Malfa, A. Zbrzezny, R. Michelmores, N. Paoletti, and M. Kwiatkowska. On guaranteed optimal robust explanations for NLP models. In *IJCAI*, 2021. In press, available from <https://arxiv.org/abs/2105.03640>.
- [23] G. A. Miller. The magical number seven, plus or minus two: Some limits on our capacity for processing information. *Psychological review*, 63(2):81, 1956.
- [24] T. Miller. Explanation in artificial intelligence: Insights from the social sciences. *Artif. Intell.*, 267:1–38, 2019.
- [25] C. Molnar. *Interpretable Machine Learning*. 2020. <http://tiny.cc/6c76tz>.
- [26] D. Monroe. Deceiving AI. *Commun. ACM*, 64, 2021.
- [27] G. Montavon, W. Samek, and K. Müller. Methods for interpreting and understanding deep neural networks. *Digital Signal Processing*, 73:1–15, 2018.
- [28] N. Narodytska, A. A. Shrotri, K. S. Meel, A. Ignatiev, and J. Marques-Silva. Assessing heuristic machine learning explanations with model counting. In *SAT*, pages 267–278, 2019.
- [29] R. S. Olson, W. La Cava, P. Orzechowski, R. J. Urbanowicz, and J. H. Moore. PMLB: a large benchmark suite for machine learning evaluation and comparison. *BioData Mining*, 10(1):36, 2017.
- [30] C. H. Papadimitriou. *Computational complexity*. Addison-Wesley, 1994.
- [31] R. Reiter. A theory of diagnosis from first principles. *Artif. Intell.*, 32(1):57–95, 1987.
- [32] M. T. Ribeiro, S. Singh, and C. Guestrin. "Why should I trust you?": Explaining the predictions of any classifier. In *KDD*, pages 1135–1144, 2016.
- [33] M. T. Ribeiro, S. Singh, and C. Guestrin. Anchors: High-precision model-agnostic explanations. In *AAAI*, 2018.
- [34] C. Rudin. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence*, 1(5):206–215, 2019.
- [35] W. Samek, G. Montavon, A. Vedaldi, L. K. Hansen, and K. Müller, editors. *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*. Springer, 2019.
- [36] A. Shih, A. Choi, and A. Darwiche. A symbolic approach to explaining bayesian network classifiers. In *IJCAI*, pages 5103–5111, 2018.
- [37] A. Shih, A. Choi, and A. Darwiche. Compiling bayesian network classifiers into decision graphs. In *AAAI*, pages 7966–7974, 2019.
- [38] D. Slack, S. Hilgard, E. Jia, S. Singh, and H. Lakkaraju. Fooling LIME and SHAP: adversarial attacks on post hoc explanation methods. In *AIES*, pages 180–186, 2020.
- [39] J. Slaney. Set-theoretic duality: A fundamental feature of combinatorial optimisation. In *ECAI*, pages 843–848, 2014.
- [40] UCI Machine Learning Repository. <https://archive.ics.uci.edu/ml>, 2020.
- [41] P. E. Utgoff, N. C. Berkman, and J. A. Clouse. Decision tree induction based on efficient tree restructuring. *Mach. Learn.*, 29(1):5–44, 1997.
- [42] S. Wäldchen, J. MacDonald, S. Hauch, and G. Kutyniok. The computational complexity of understanding binary classifier decisions. *J. Artif. Intell. Res.*, 70:351–387, 2021.
- [43] D. S. Weld and G. Bansal. The challenge of crafting intelligible intelligence. *Commun. ACM*, 62(6):70–79, 2019.

Supplementary Material

Proofs

Deciding δ -relevancy.

Definition 7 (MajSAT[30]). *Given a boolean function $f : \{0, 1\}^n \rightarrow \{0, 1\}$, the MajSAT problem is to decide whether the number of points \mathbf{x} with $f(\mathbf{x}) = 1$ exceeds the number of points with $f(\mathbf{x}) = 0$.*

It is well-known that MajSAT is PP-complete [30].

Proposition 5. *Deciding whether a set S is a $C\delta$ -relevant set is PP-hard.*

Proof. [Sketch]

We reduce MajSAT to deciding $C\delta$ -relevancy.

Let $f : \{0, 1\}^n \rightarrow \{0, 1\}$ be a boolean function. The variables of f are $X = \{x_1, \dots, x_n\}$.

We want to decide whether $\Pr(f(\mathbf{x}) = 1) > \Pr(f(\mathbf{x}) = 0)$. We create another function $F : \{0, 1\}^n \times \{0, 1\}^2 \rightarrow \{0, 1\}$, such that the variables of F are X and $P = \{p_1, p_2\}$. Moreover, F is defined as follows:

$$F(\mathbf{x}, \mathbf{p}) = \begin{cases} 1 & \text{if } p_1 = p_2 = 1 \\ f(\mathbf{x}) & \text{otherwise} \end{cases}$$

Set $(\mathbf{x}_a, \mathbf{p}_a) = ((0, \dots, 0), (1, 1))$. Clearly, $F(\mathbf{x}_a, \mathbf{p}_a) = 1$.

Moreover, set $\delta = 0.75$ and pick $S = \{p_1\}$.

Now, if $\Pr(F(\mathbf{x}_b, \mathbf{p}_b) = 1 | (\mathbf{x}_b, \mathbf{p}_b)_S = (\mathbf{x}_a, \mathbf{p}_a)_S) > \delta$ iff the number of points with $f(\mathbf{x}) = 1$ exceeds the number of points with $f(\mathbf{x}) = 0$. \square