

On Computing Probabilistic Explanations for Decision Trees

MARCELO ARENAS, Department of Computer Science and Institute for Mathematical and Computational Engineering, Faculty of Engineering, Catholic University of Chile, Chile and IMFD, Chile

PABLO BARCELÓ, Institute for Mathematical and Computational Engineering, Faculty of Engineering and Faculty of Mathematics, Catholic University of Chile, Chile, IMFD, Chile, and CENIA, Chile

ALEXANDER KOZACHINSKIY, CENIA, Chile

MIGUEL ROMERO, Department of Computer Science, Faculty of Engineering, Catholic University of Chile, Chile and CENIA, Chile

BERNARDO SUBERCASEAUX, School of Computer Science, Carnegie Mellon University, USA

Formal XAI (explainable AI) is a growing area that focuses on computing explanations with mathematical guarantees for the decisions made by ML models. Inside formal XAI, one of the most studied cases is that of explaining the choices taken by decision trees, as they are traditionally deemed as one of the most interpretable classes of models. Recent work has focused on studying the computation of *sufficient reasons*, a kind of explanation in which given a decision tree T and an instance x , one explains the decision $T(x)$ by providing a subset y of the features of x such that for any other instance z compatible with y , it holds that $T(z) = T(x)$, intuitively meaning that the features in y are already enough to fully justify the classification of x by T . It has been argued, however, that sufficient reasons constitute a restrictive notion of explanation. For such a reason, the community has started to study their probabilistic counterpart, in which one requires that the probability of $T(z) = T(x)$ must be at least some value $\delta \in (0, 1]$, where z is a random instance that is compatible with y . Our paper settles the computational complexity of δ -sufficient-reasons over decision trees, showing that both (1) finding δ -sufficient-reasons that are minimal in size, and (2) finding δ -sufficient-reasons that are minimal inclusion-wise, are computationally intractable. By doing this, we answer two open problems originally raised by Izza et al. (2021), and extend the hardness of explanations for Boolean circuits presented by Wäldchen et al. (2021) to the more restricted case of decision trees. Furthermore, we present sharp non-approximability results under a widely believed complexity hypothesis. On the positive side, we identify structural restrictions of decision trees that make the problem tractable.

JAIR Associate Editor: Rafael Peñaloza

JAIR Reference Format:

Marcelo Arenas, Pablo Barceló, Alexander Kozachinskiy, Miguel Romero, and Bernardo Subercaseaux. 2025. On Computing Probabilistic Explanations for Decision Trees. *Journal of Artificial Intelligence Research* 83, Article 34 (August 2025), 44 pages. DOI: [10.1613/jair.1.17494](https://doi.org/10.1613/jair.1.17494)

Authors' Contact Information: Marcelo Arenas, ORCID: [0000-0003-3678-1868](https://orcid.org/0000-0003-3678-1868), marenas@uc.cl, Department of Computer Science and Institute for Mathematical and Computational Engineering, Faculty of Engineering, Catholic University of Chile, Santiago, Metropolitan Region, Chile and IMFD, Santiago, Metropolitan Region, Chile; Pablo Barceló, ORCID: [0000-0003-2293-2653](https://orcid.org/0000-0003-2293-2653), pbarcelo@uc.cl, Institute for Mathematical and Computational Engineering, Faculty of Engineering and Faculty of Mathematics, Catholic University of Chile, Santiago, Metropolitan Region, Chile and IMFD, Santiago, Metropolitan Region, Chile and CENIA, Santiago, Metropolitan Region, Chile; Alexander Kozachinskiy, ORCID: [0009-0002-8690-3640](https://orcid.org/0009-0002-8690-3640), kozmath@proton.me, CENIA, Santiago, Metropolitan Region, Chile; Miguel Romero, ORCID: [0000-0002-2615-6455](https://orcid.org/0000-0002-2615-6455), mgromero@uc.cl, Department of Computer Science, Faculty of Engineering, Catholic University of Chile, Santiago, Metropolitan Region, Chile and CENIA, Santiago, Metropolitan Region, Chile; Bernardo Subercaseaux, ORCID: [0000-0003-2295-1299](https://orcid.org/0000-0003-2295-1299), bersub@cmu.edu, School of Computer Science, Carnegie Mellon University, Pittsburgh, Pennsylvania, USA.



This work is licensed under a [Creative Commons Attribution International 4.0 License](https://creativecommons.org/licenses/by/4.0/).

© 2025 Copyright held by the owner/author(s).
DOI: [10.1613/jair.1.17494](https://doi.org/10.1613/jair.1.17494)

1 Introduction

The level of trust that AI models generate in people has been repetitively linked to our ability of *explaining* the decision of said models [1], thus motivating the area of *explainable AI* (XAI) as an important aspect of the deployment of trustworthy artificial intelligence. A sub-area of explainability that has received considerable attention over the last years, showing quick progress in theoretical and practical terms, is that of *local* explanations, that is, explanations for the outcome of a particular input to an ML model after the model has been trained. Several queries and scores have been proposed to specify explanations of this kind. These include, for example, queries based on *prime implicants* [2] or *anchors* [3], which are parts of an instance that are sufficient to explain its classification, as well as scores that intend to quantify the impact of a single feature in the output of such a classification [4, 5].

A remarkable achievement of this area of research has been the development of *formal* notions of explainability. The benefits brought about by this principled approach have been highlighted in a very thorough way by a recent survey [6]. A prime example of this kind of approach is given by “*sufficient reasons*”, which are also known as prime implicant explanations [2] or abductive explanations [7].

Given an ML model \mathcal{M} of dimension n and a Boolean input instance $\mathbf{x} \in \{0, 1\}^n$, a sufficient reason for \mathbf{x} under \mathcal{M} is a subset \mathbf{y} of the features of \mathbf{x} , such that any instance \mathbf{z} compatible with \mathbf{y} receives the same classification result as \mathbf{x} on \mathcal{M} . In more intuitive words, \mathbf{y} is a sufficient reason for \mathbf{x} under \mathcal{M} if the features in \mathbf{y} suffices to explain the output of \mathcal{M} on \mathbf{x} . In the formal explainability approach, one then aims to find sufficient reasons \mathbf{y} that satisfy one of the following optimality criteria:

- (1) they are *minimum*, that is, there are no sufficient reasons with fewer features than \mathbf{y} , or
- (2) they are *minimal*, that is, there are no sufficient reasons that are strictly contained in \mathbf{y} .

Problem. The XAI community has studied for which Boolean ML models the problem of computing (minimum or minimal) sufficient reasons is computationally tractable and for which it is computationally hard (see, e.g., [8, 9, 10]). It has been argued, however, that for practical applications sufficient reasons might be too *rigid*, as they are specified under worst-case conditions. That is, \mathbf{y} is a sufficient reason for \mathbf{x} under \mathcal{M} if *every* “completion” of \mathbf{y} is classified by \mathcal{M} in the same way as \mathbf{x} . As several authors have noted already, there is a natural way in which this notion can be relaxed in order to become more suitable for real-world explainability tasks: Instead of asking for each completion of \mathbf{y} to yield the same result as \mathbf{x} on \mathcal{M} , we could allow for a small fragment of the completions of \mathbf{y} to be classified differently than \mathbf{x} [11, 12, 13]. More precisely, we would like to ensure that a random completion of \mathbf{y} is classified as \mathbf{x} with probability at least $\delta \in (0, 1]$, a threshold that the recipient of the explanation controls. We call \mathbf{y} a δ -*sufficient reason for \mathbf{x} under \mathcal{M}* .

The study of the cost of computing minimum δ -sufficient reasons for expressive Boolean ML models based on propositional formulas was started by [11]. They show, in particular, that the decision problem of checking if \mathbf{x} admits a δ -sufficient reason of a certain size k under a model \mathcal{M} , where \mathcal{M} is specified as a CNF formula, is NP^{PP} -complete. This result shows that the problem is very difficult for complex models, at least in theoretical terms. Nonetheless, it leaves the door open for obtaining tractability results over simpler Boolean models, starting from those which are often deemed to be “easy to interpret”, for example, *decision trees* [14, 15, 16]. In particular, the study of the cost of computing both minimum and minimal δ -sufficient reasons for decision trees was initiated by [12, 17], but nothing beyond the fact that the problem lies in NP has been obtained. The work of [18] has shown that it is possible to obtain efficient algorithms that succeed with a certain probability, and that instead of finding a smallest (either cardinality or inclusion-wise) δ -sufficient reason, find δ -sufficient reasons that are small compared to the *average* size of δ -sufficient reasons for the considered model.

Decision Trees need to be explained. When looking at small decision trees (e.g., ~15 leaves, as the example depicted in Figure 1) it is tempting to think that decision trees do not require algorithms to be explained, as a

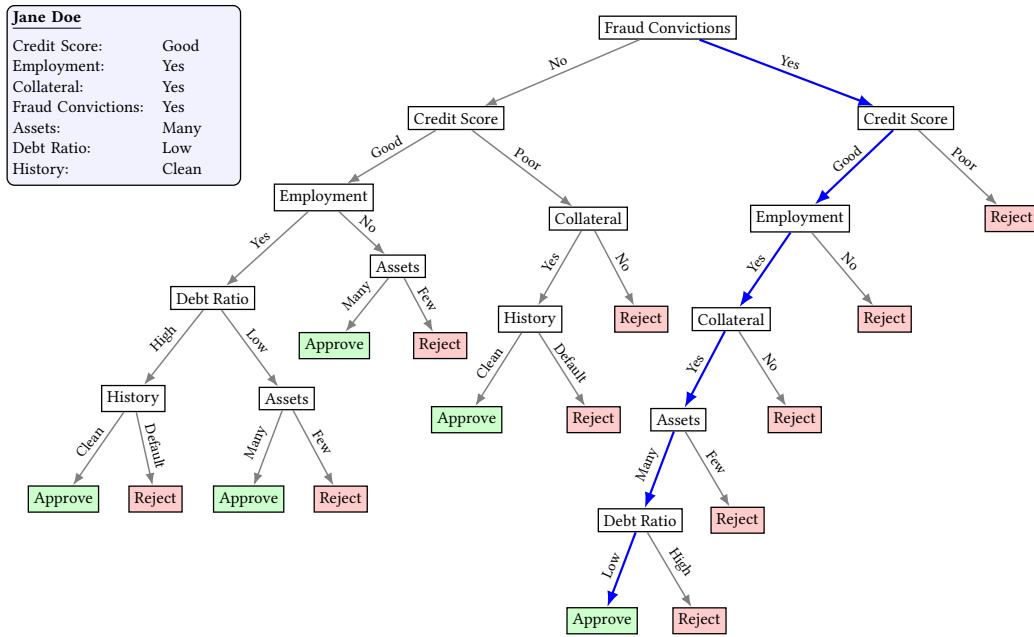


Fig. 1. Example of a decision tree for judging loan applicants with binary features. The decision path for a concrete applicant, Jane Done, is highlighted in blue.

human can readily look at the root-to-leaf path associated to any classification decision they want to understand. This view, however, has been repeatedly challenged by the formal XAI community as rather naïve [15, 19, 20]. Let us present a brief summary of some of the issues involved. First, practical training algorithms, with standard choices of hyper-parameters, often result in trees with hundreds or even thousands of nodes, which makes manual inspection significantly harder than in toy examples [21]. Second, “path explanations” (i.e., the set of features mentioned in the path that a given input instance takes along the decision tree) can be much larger than small sufficient reasons, as shown both in theory and practice [15, 19, 20]. The decision tree example presented in Figure 1 depicts a path explanation for applicant *Jane Doe* that includes 6 out of the 7 possible features. However, its subset of 5 features not including *Fraud Convictions* can be easily checked to be a (minimum) sufficient reason.¹ Not only is this reducing the size of the explanation by 1, but it also makes the explanation less misleading: we claim that *Jane Doe*’s application is not approved *because* of her fraud convictions, but *despite* them, since her other features as an applicant are very positive. Checking our claim requires inspecting not only the decision path of *Jane*’s application, but potentially the entire tree, to see that if she did not have fraud convictions, her application would have taken a different path but resulted in approval nonetheless. While in this particular toy example the feature names are suggestive enough to make rather obvious that the fraud conviction is not a good explanation for *Jane*’s approval, we claim this illustrates how in a general case, with a large number of potentially opaque features, minimal explanations are needed to understand whether a classification is happening *despite* a given a feature or *partially because* of it. As it has been claimed by others in the literature (e.g., [15]), we believe this makes a compelling argument to reject the view that decision trees are immediately interpretable by looking

¹As shown by [15], the difference in size between path explanations and minimal sufficient reasons can be arbitrarily large.

at their decision paths, and that to properly understand their decisions it might be necessary to run algorithms that obtain explanations that consider the globality of the tree.

Our results. In this paper, we provide an in-depth study of the complexity of the problem of minimum and minimal δ -sufficient reasons for decision trees.

- *Hardness results:* We first show that, assuming $P \neq NP$, neither of these problems can be solved in polynomial time. For the problem of computing minimum δ -sufficient reasons over decision trees, the hardness holds even for each fixed $\delta \in (0, 1]$.
- *Non-approximability results:* Then we study the approximation properties of these problems, and establish that neither of them can be approximated by assuming a stronger complexity conjecture; namely, that NP problems cannot be solved in quasi-polynomial time. Although this assumption is stronger than $P \neq NP$, it is still widely accepted. In fact, even the so-called ETH (*exponential time hypothesis*), stating that SAT cannot be solved just in sub-exponential time, is a standard complexity assumption now. Nevertheless, this means that our non-approximability results do not directly imply our hardness results, as the former rely on a stronger hypothesis. We mention that similar inapproximability results have been obtained by [11], but for Boolean circuits instead of decision trees, making our results much stronger (although their results are under the weaker assumption, $P \neq NP$).
- *Tractability results:* Finally, we look at structural restrictions on decision trees that, at the same time, represent meaningful practical instances and ensure that the problems studied in the paper can be solved in polynomial time. The first such restriction is having a *bounded split number*, which intuitively means that there is a constant $c \geq 1$ such that, for every node u in the decision tree T , the number of features that appear both in the subtree of T rooted at u and outside this subtree is bounded by c . We show that both problems studied in this paper, that is, computing minimum and minimal δ -sufficient reasons over decision trees, can be solved in polynomial time for decision trees of bounded split number. The second restriction studied is *monotonicity*, which refers to the class of decision trees that, when presented with an instance $x \in \{0, 1\}^n$ that is classified positively, also classify as positive every instance x' that is obtained from x by changing the value of some of its features from 0 to 1. We show that the problem of computing minimal δ -sufficient reasons can be solved in polynomial time over monotone decision trees.

Organization of the paper. We start with preliminaries in Section 2. Our main results are presented, without proof, in Section 3. Proofs of our hardness results are in Section 4, of non-approximability results in Section 5, and of tractability results in Section 6. We finalize with closing remarks and open problems in Section 7.

2 Preliminaries

Let us begin by providing the necessary definitions and notation that will be used throughout the paper.

2.1 Models and Instances

An *instance* of dimension n , with $n \geq 1$, is a tuple $x \in \{0, 1\}^n$. We use notation $x[i]$ to refer to the i -th component of this tuple, or equivalently, its i -th feature. Moreover, we consider an abstract notion of a model of dimension n , and we define it as a Boolean function $\mathcal{M} : \{0, 1\}^n \rightarrow \{0, 1\}$. That is, \mathcal{M} assigns a Boolean value to each instance of dimension n , so that we focus on binary classifiers with Boolean input features. Restricting inputs and outputs to be Boolean makes our setting cleaner while still covering several relevant practical scenarios. We use notation $\dim(\mathcal{M})$ for the dimension of a model \mathcal{M} .

2.2 Binary Decision Diagrams and Decision Trees

A *binary decision diagram* (BDD) of dimension n is a rooted directed acyclic graph \mathcal{M} with labels on edges and nodes such that:

- each leaf (a node with no outgoing edges) is labeled with **true** or **false**,
- each internal node (a node that is not a leaf) is labeled with a feature $i \in \{1, \dots, n\}$,
- each internal node has two outgoing edges, one labeled 1 and the other one labeled 0.

Every instance $\mathbf{x} \in \{0, 1\}^n$ defines a unique path $\pi_{\mathbf{x}} = u_1 \cdots u_k$ from the root u_1 to a leaf u_k of \mathcal{M} such that: if the label of u_i is $j \in \{1, \dots, n\}$, where $i \in \{1, \dots, k-1\}$, then the edge from u_i to u_{i+1} is labeled with $\mathbf{x}[j]$. Moreover, the instance \mathbf{x} is positive, denoted by $\mathcal{M}(\mathbf{x}) = 1$, if the label of u_k is **true**; otherwise the instance \mathbf{x} is negative, which is denoted by $\mathcal{M}(\mathbf{x}) = 0$. A BDD \mathcal{M} is *free* if for every path from the root to a leaf, no two nodes on that path have the same label. A *decision tree* is simply a free BDD whose underlying directed acyclic graph is a rooted tree. When a computational problem has a decision tree as part of its input, we assume a standard encoding whose number of bits is polynomial in the number of nodes of the tree.

2.3 Partial Instances

In the following, a *partial instance* of dimension n is a tuple $\mathbf{y} \in \{0, 1, \perp\}^n$. Intuitively, if $\mathbf{y}[i] = \perp$, then the value of the i -th feature is undefined. Notice that an instance is a particular case of a partial instance where all features are assigned value either 0 or 1. Given two partial instances \mathbf{x}, \mathbf{y} of dimension n , we say that \mathbf{y} is *subsumed* by \mathbf{x} , denoted by $\mathbf{y} \subseteq \mathbf{x}$, if for every $i \in \{1, \dots, n\}$ such that $\mathbf{y}[i] \neq \perp$, it holds that $\mathbf{y}[i] = \mathbf{x}[i]$. That is, \mathbf{y} is subsumed by \mathbf{x} if it is possible to obtain \mathbf{x} from \mathbf{y} by replacing some undefined values. Moreover, we say that \mathbf{y} is *properly subsumed* by \mathbf{x} , denoted by $\mathbf{y} \subsetneq \mathbf{x}$, if $\mathbf{y} \subseteq \mathbf{x}$ and $\mathbf{y} \neq \mathbf{x}$. For instance,

$$(0, \perp, 1, \perp) \subsetneq (0, 0, 1, \perp) \quad \text{and} \quad (0, \perp, 1, \perp) \not\subseteq (0, 0, 0, \perp).$$

Notice that a partial instance \mathbf{y} can be thought of as a compact representation of the set of instances \mathbf{z} such that \mathbf{y} is subsumed by \mathbf{z} , where such instances \mathbf{z} are called the *completions* of \mathbf{y} . We write $\text{COMP}(\mathbf{y})$ for the set of completions of \mathbf{y} .

2.4 Sufficient Reasons

Sufficient reasons are partial instances obtained by removing from an instance \mathbf{x} components that do not affect the final classification. Formally, fix a dimension n . Given a decision tree T , an instance \mathbf{x} , and a partial instance \mathbf{y} with $\mathbf{y} \subseteq \mathbf{x}$, we call \mathbf{y} a *sufficient reason* for \mathbf{x} under T if $T(\mathbf{x}) = T(\mathbf{z})$ for every $\mathbf{z} \in \text{COMP}(\mathbf{y})$. In other words, the features of \mathbf{y} that take value either 0 or 1 explain the decision taken by T on \mathbf{x} , as $T(\mathbf{x})$ would not change if the remaining features (i.e., those that are undefined in \mathbf{y}) were to change in \mathbf{x} , thus implying that the classification $T(\mathbf{x})$ is a consequence of the features defined in \mathbf{y} .

We say that a sufficient reason \mathbf{y} for \mathbf{x} under T is *minimal*, if it is minimal under the order induced by \subseteq , that is, if there is no sufficient reason \mathbf{y}' for \mathbf{x} under T such that $\mathbf{y}' \subsetneq \mathbf{y}$. Also, we define a *minimum* sufficient reason for \mathbf{x} under T as a sufficient reason \mathbf{y} for \mathbf{x} under T that maximizes the value $|\mathbf{y}|_{\perp} := |\{i \in \{1, \dots, n\} \mid \mathbf{y}[i] = \perp\}|$.

It turns out that minimal sufficient reasons can be computed efficiently over decision trees.

Proposition 2.1. *There is a polynomial-time² algorithm that, given a decision tree T and an instance \mathbf{x} of the same dimension, computes a minimal sufficient reason for \mathbf{x} under T .*

²Whenever we say a polynomial-time algorithm, the polynomial is with respect to the entire input size, which usually involves trees (whose size is polynomial in the number of nodes), partial instances (whose size is proportional to the dimension), and integers, which we assume are encoded in binary.

This proposition can be established by using a very simple algorithm, assuming a sub-routine to check whether a given partial instance is a sufficient reason (not necessarily minimal) of another given instance. As shown in Algorithm 2, the idea of the algorithm is as follows: start with a candidate answer \mathbf{y} which is initially equal to \mathbf{x} , the instance to explain, and maintain the invariant that \mathbf{y} is a sufficient reason for \mathbf{x} , while trying to remove defined components from \mathbf{y} until no longer possible. It is not hard to see that one can check whether a partial instance \mathbf{y} is a sufficient reason for an instance \mathbf{x} in linear time over decision trees [22], which implies that Algorithm 2 runs in polynomial time over such models.

Algorithm 2 Minimal Sufficient Reason

```

1: Input: Decision tree  $T$  and instance  $\mathbf{x}$ , both of dimension  $n$ 
2: Output: A minimal sufficient reason  $\mathbf{y}$  for  $\mathbf{x}$  under  $T$ 
3:  $\mathbf{y} \leftarrow \mathbf{x}$ 
4: for  $i \in \{1, \dots, n\}$  do
5:    $\hat{\mathbf{y}} \leftarrow \mathbf{y}$ 
6:    $\hat{\mathbf{y}}[i] \leftarrow \perp$ 
7:   if CheckSufficientReason( $T, \hat{\mathbf{y}}, \mathbf{x}$ ) then
8:      $\mathbf{y} \leftarrow \hat{\mathbf{y}}$ 
9:   end if
10: end for
11: return  $\mathbf{y}$ 

```

This algorithm is well known (see e.g., [23]), and its correctness relies on the following simple result, tracing back to [24].

Proposition 2.2. *For any class of models \mathfrak{C} , if a partial instance \mathbf{y} of dimension n is a sufficient reason for an instance \mathbf{x} under a model $\mathcal{M} \in \mathfrak{C}$, but not a minimal sufficient reason, then there is a partial instance $\hat{\mathbf{y}}$ which is equal to \mathbf{y} except that for some $i \in \{1, \dots, n\}$ we have that $\hat{\mathbf{y}}[i] = \perp$ and $\mathbf{y}[i] \neq \perp$, and such that $\hat{\mathbf{y}}$ is also a sufficient reason for \mathbf{x} under \mathcal{M} .*

The following theorem shows a stark contrast between the complexity of computing minimal and minimum sufficient reasons over decision trees.

THEOREM 2.3 ([9]). *Assuming $P \neq NP$, there is no polynomial-time algorithm that, given a decision tree T and an instance \mathbf{x} of the same dimension, computes a minimum sufficient reason for \mathbf{x} under T .*

2.5 Probabilistic Sufficient Reasons

Arguably, the notion of sufficient reason is a natural notion of explanation for the result of a classifier. However, such a concept imposes a severe restriction by asking all completions of a partial instance to be classified in the same way. To overcome this limitation, a probabilistic generalization of sufficient reasons was proposed by [11] and [12]. More precisely, this notion allows to settle a confidence $\delta \in (0, 1]$ ³ on the fraction of completions of a partial instance that yield the same classification.

Definition 2.4 (Probabilistic sufficient reasons). *Given a value $\delta \in (0, 1]$, a δ -sufficient reason (δ -SR for short) for an instance \mathbf{x} under a decision tree T is a partial instance \mathbf{y} such that $\mathbf{y} \subseteq \mathbf{x}$ and*

$$\Pr_{\mathbf{z}}[T(\mathbf{z}) = T(\mathbf{x}) \mid \mathbf{z} \in \text{COMP}(\mathbf{y})] := \frac{|\{\mathbf{z} \in \text{COMP}(\mathbf{y}) \mid T(\mathbf{z}) = T(\mathbf{x})\}|}{2^{|\mathbf{y}|_{\perp}}} \geq \delta.$$

³Given rational numbers a, b with $a \leq b$, we use notation $[a, b]$ for the set of rational numbers x such that $a \leq x \leq b$ (and $[a, b)$ for the set $\{x \in \mathbb{Q} \mid a \leq x < b\}$).

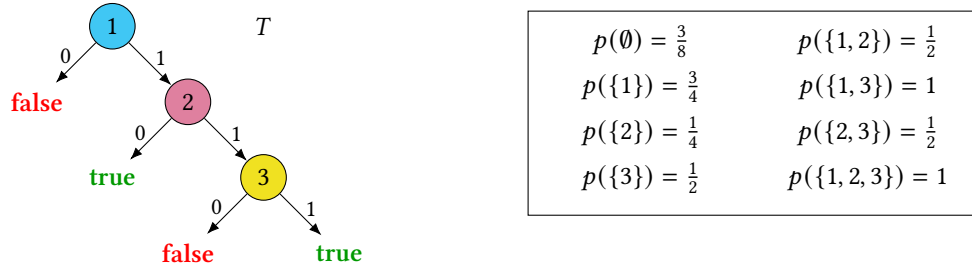


Fig. 2. The decision tree T and the values $p(X)$ from Example 2.5.

Minimal and minimum δ -sufficient reasons are defined analogously as in the case of minimal and minimum sufficient reasons.

Example 2.5. Consider the decision tree T over features $\{1, 2, 3\}$ shown in Figure 2 and the input instance $\mathbf{x} = (1, 1, 1)$. Notice that $T(\mathbf{x}) = 1$. For each $X \subseteq \{1, 2, 3\}$, we show the probability $p(X) := \Pr_z[T(z) = 1 \mid z \in \text{COMP}(\mathbf{y}_X)]$, where \mathbf{y}_X is the partial instance that is obtained from \mathbf{x} by fixing $\mathbf{y}_X[i] = \perp$ for each $i \notin X$. We observe, for instance, that \mathbf{x} itself is neither a minimum nor a minimal 1-SR for \mathbf{x} under T , as $\mathbf{y}_{\{1,3\}} = (1, \perp, 1)$ is also a 1-SR. In turn, $\mathbf{y}_{\{1,3\}}$ is both a minimal and a minimum 1-SR for \mathbf{x} under T . The partial instance $\mathbf{y}_{\{1,3\}}$ is not, however, a minimal or a minimum $3/4$ -SR for \mathbf{x} under T , as $\mathbf{y}_{\{1\}} = (1, \perp, \perp)$ is also a $3/4$ -SR. \square

Example 2.6. Consider the decision tree T over features $\{1, 2, 3\}$ shown in Figure 3 and the input instance $\mathbf{x} = (1, 1, 1)$. Notice that $T(\mathbf{x}) = 1$. Exactly as in Example 2.5, we display as well the probabilities $p(X)$ for each $X \subseteq \{1, 2, 3\}$. Interestingly, this example illustrates that Proposition 2.2 does not hold when $\delta < 1$. Indeed, consider that $\mathbf{y}_{\{1,2,3\}}$ is a $5/8$ -SR which is not minimal, as \mathbf{y}_\emptyset is also a $5/8$ -SR, but if we remove any single feature from $\mathbf{y}_{\{1,2,3\}}$, we obtain a partial instance which is not a $5/8$ -SR. \square

As illustrated on Example 2.6, it is not true in general that if $\mathbf{y}' \subseteq \mathbf{y}$ then

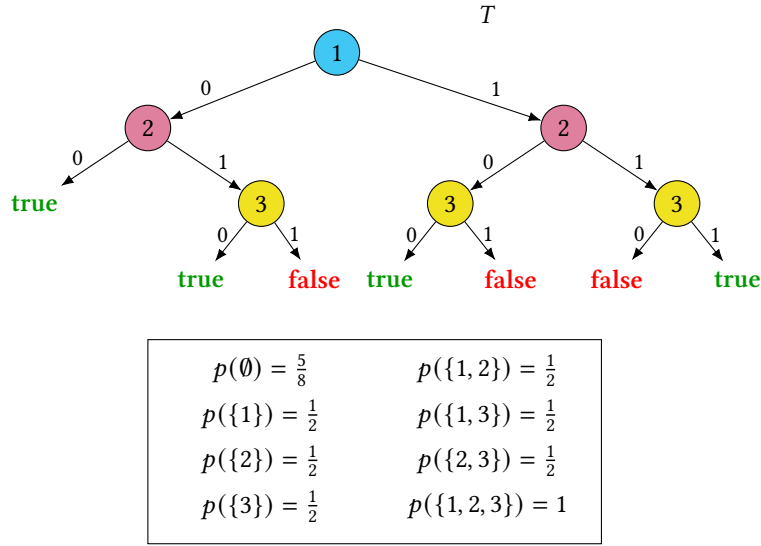
$$\Pr_z[T(z) = T(\mathbf{x}) \mid z \in \text{COMP}(\mathbf{y}')] \leq \Pr_z[T(z) = T(\mathbf{x}) \mid z \in \text{COMP}(\mathbf{y})],$$

which means that standard algorithms for finding minimal sets holding monotone predicates (such as Algorithm 2) cannot be used to compute minimal δ -SRs. In fact, we will show later that no polynomial time algorithm exists for computing minimal δ -SRs (at least under complexity assumptions).

2.6 Computing Probabilistic Sufficient Reasons

The problems of computing minimum and minimal δ -SR on decision trees were defined and left open by [12, 17]. These problems are formally defined as follows.

PROBLEM: : Compute-Minimum-SR INPUT : A decision tree T of dimension n , an instance \mathbf{x} of dimension n and $\delta \in (0, 1]$ OUTPUT : A minimum δ -SR for \mathbf{x} under T
PROBLEM: : Compute-Minimal-SR INPUT : A decision tree T of dimension n , an instance \mathbf{x} of dimension n and $\delta \in (0, 1]$ OUTPUT : A minimal δ -SR for \mathbf{x} under T

Fig. 3. The decision tree T and the values $p(X)$ from Example 2.6.

3 Main Results

In this section, we summarize the main results presented in this paper regarding the computation of probabilistic sufficient reasons over decision trees.

3.1 Hardness Results

We first establish that, if $P \neq NP$, neither Compute-Minimum-SR nor Compute-Minimal-SR can be solved in polynomial time. We first consider the problem Compute-Minimum-SR, and in fact, prove a stronger result by considering the family of problems δ -Compute-Minimum-SR which are obtained from Compute-Minimum-SR when $\delta \in (0, 1]$ is assumed to be fixed. More precisely, we obtain as a corollary of Theorem 2.3 that 1-Compute-Minimum-SR cannot be solved efficiently. Moreover, a non-trivial modification of the proof of this theorem shows that this negative result continues to hold for every fixed $\delta \in (0, 1]$.

THEOREM 3.1. *Fix any value for $\delta \in (0, 1]$. Then, assuming that $P \neq NP$, there is no polynomial-time algorithm for the δ -Compute-Minimum-SR problem.*

Let us now look at the problem Compute-Minimal-SR. When $\delta = 1$, this problem can be solved in polynomial time as stated in Theorem 2.3. However, it was conjectured by [12] that assuming $P \neq NP$, this positive behavior does not extend to the general problem Compute-Minimal-SR, in which δ is an input confidence parameter. Our main result confirms that this conjecture is correct.

THEOREM 3.2. *Assuming that $P \neq NP$, the Compute-Minimal-SR problem cannot be solved in polynomial time.*

The proof of this theorem heavily relies on the following result, which we believe to be of independent interest. The result establishes that the following problem over formulas in CNF, which we call Minimal-Expected-Clauses, is NP-hard. The problem Minimal-Expected-Clauses is defined as follows. Let φ be a CNF formula over variables $X = \{x_1, \dots, x_n\}$. Partial assignments of the variables in X , as well as the notions of subsumption and completions over them, are defined in the same way as for partial instances over features. For a partial assignment μ over

X , we denote by $E(\varphi, \mu)$ the expected number of clauses of φ satisfied by a random completion of μ . We then consider the following problem for fixed $k \geq 2$ (recall that a k -CNF formula is a CNF formula where each clause has at most k literals):

PROBLEM : k -Minimal-Expected-Clauses
 INPUT : (φ, σ) , for φ a k -CNF formula and σ a partial assignment
 OUTPUT : Yes, if there is a $\mu \subseteq \sigma$ such that $E(\varphi, \mu) \geq E(\varphi, \sigma)$
 and No otherwise

THEOREM 3.3. k -Minimal-Expected-Clauses is NP-hard for every $k \geq 2$.

3.2 Non-Approximability Results

By assuming a stronger complexity assumption, namely, $\text{NP} \not\subseteq \text{QP}$, where QP consists of all problems that have quasi-polynomial time algorithms, we can actually show stronger negative results that imply the non-approximability of the problems studied in this paper. Intuitively, we show that it is hard to distinguish cases when a very good sufficient reason with almost all variables left unfixed exists from cases when in all mildly good sufficient reasons just a few variables are left unfixed.

More precisely, for Compute-Minimum-SR we show the following:

THEOREM 3.4. Assuming that $\text{NP} \not\subseteq \text{QP}$, for every $\varepsilon > 0$ there is no polynomial-time algorithm that, given a decision tree T of dimension n and an input $\mathbf{x} \in \{0, 1\}^n$, distinguishes between the following two cases:

- there exists an $(1 - \varepsilon)$ -sufficient reason \mathbf{y} for \mathbf{x} under T with $|\mathbf{y}|_{\perp} \geq n - n^{\varepsilon}$.
- there exists no ε -sufficient reason \mathbf{y} for \mathbf{x} under T with $|\mathbf{y}|_{\perp} \geq n^{\varepsilon}$.

We now address the Compute-Minimal-SR problem. We obtain the following inapproximability result.

THEOREM 3.5. Assuming that $\text{NP} \not\subseteq \text{QP}$, for every $\varepsilon \in (0, 1)$, there exists no polynomial-time algorithm that, given a decision tree T of dimension n and an input $\mathbf{x} \in \{0, 1\}^n$, distinguishes between the following two cases:

- the only $(1/2 + \varepsilon)$ -sufficient reason for \mathbf{x} under T is \mathbf{x} itself.
- there exists a $(1 - \varepsilon)$ -sufficient reason for \mathbf{x} under T other than \mathbf{x} .

3.3 Tractability Results

We finally study restrictions on decision trees that lead to polynomial time algorithms for Compute-Minimum-SR or Compute-Minimal-SR, hence avoiding the general intractability results shown in the previous section. We identify two such restrictions: *bounded split number* and *monotonicity*.

3.3.1 *Bounded Split Number*. Let T be a decision tree of dimension n . For a set U of nodes of T , we denote by $\mathcal{F}(U)$ the set of features from $\{1, \dots, n\}$ labeling the nodes in U . For each node u of T , we denote by N_u^{\downarrow} the set of nodes appearing in T_u , that is, the subtree of T rooted at u . On the other hand, we denote by N_u^{\uparrow} the set of nodes of T minus the set of nodes N_u^{\downarrow} . We define the *split number* of the decision tree T to be

$$\max_{\text{node } u \text{ in } T} \left| \mathcal{F}(N_u^{\downarrow}) \cap \mathcal{F}(N_u^{\uparrow}) \right|.$$

Intuitively, the split number of a decision tree T is a measure of the interaction (number of common features) between the subtrees of the form T_u and their exterior. A small split number allows us to essentially treat each subtree T_u independently (in particular, the left and right subtrees below any node), which in turn leads to efficient algorithms for the problems Compute-Minimum-SR and Compute-Minimal-SR.

THEOREM 3.6. Both problems Compute-Minimum-SR and Compute-Minimal-SR can be solved in polynomial time for decision trees with split number at most c , where $c \geq 1$ is any fixed integer.

3.3.2 Monotonicity. Classifiers that are *monotone* have been studied in the context of XAI as they often present tractable cases for different explanations, as shown in [10]. The computation of minimal sufficient reasons for monotone models was known to be in polynomial time [24]. We show that this is also the case for computing minimal δ -SRs under a mild assumption on the class of models.

Let us define the ordering \preceq for instances in $\{0, 1\}^n$ as follows:

$$\mathbf{x} \preceq \mathbf{z} \quad \text{iff} \quad \mathbf{x}[i] \leq \mathbf{z}[i], \text{ for all } i \in \{1, \dots, n\}.$$

We can now define monotonicity as follows. A model \mathcal{M} of dimension n is said to be *monotone* if for every pair of instances $\mathbf{x}, \mathbf{z} \in \{0, 1\}^n$, it holds that

$$\mathbf{x} \preceq \mathbf{z} \implies \mathcal{M}(\mathbf{x}) \leq \mathcal{M}(\mathbf{z}).$$

We now prove that the problem of computing minimal probabilistic sufficient reasons can be solved in polynomial time for any class \mathfrak{C} of monotone Boolean models for which the problem of counting positive completions can be solved efficiently. Formally, the latter problem is defined as follows: given a model $\mathcal{M} \in \mathfrak{C}$ of dimension n and a partial instance $\mathbf{y} \in \{0, 1, \perp\}^n$, compute

$$|\{\mathbf{x} \in \text{COMP}(\mathbf{y}) \mid \mathcal{M}(\mathbf{x}) = 1\}|.$$

We call this problem \mathfrak{C} -#Positive-Completions.

THEOREM 3.7. *Let \mathfrak{C} be a class of monotone Boolean models such that \mathfrak{C} -#Positive-Completions can be solved in polynomial time. Then the problem Compute-Minimal-SR can be solved in polynomial time over \mathfrak{C} .*

As a corollary, the computation of minimal probabilistic sufficient reasons can be carried out in polynomial time not only over monotone decision trees, but also over monotone and free BDDs.

Corollary 3.8. *The problem of computing minimal δ -SRs can be solved in polynomial time over the class of monotone and free BDDs.*

4 Proofs of the Hardness Results

We start by stating an auxiliary lemma that we will need for Theorem 3.1.

Lemma 4.1. *Fix $\delta \in (0, 1)$. Given as input an integer n one can build in $n^{O(1)}$ time a decision tree T_δ of dimension n , such that*

$$\delta - \frac{1}{2^n} \leq \Pr_z [T_\delta(\mathbf{z}) = 1 \mid \mathbf{z} \in \text{COMP}(\perp^n)] \leq \delta,$$

and moreover, there exists an instance \mathbf{x}^\dagger for T_δ such that every partial instance $\mathbf{y} \subseteq \mathbf{x}^\dagger$ holds

$$\Pr_z [T_\delta(\mathbf{z}) = 1 \mid \mathbf{z} \in \text{COMP}(\mathbf{y})] \leq \Pr_z [T_\delta(\mathbf{z}) = 1 \mid \mathbf{z} \in \text{COMP}(\perp^n)] \leq \delta.$$

Let us now restate and prove Theorem 3.1 assuming Lemma 4.1, and after that we present the proof of Lemma 4.1.

THEOREM 3.1 (RESTATED). *Fix any value for $\delta \in (0, 1]$. Then, assuming that $P \neq NP$, there is no polynomial-time algorithm for the δ -Compute-Minimum-SR problem.*

PROOF OF THEOREM 3.1. We prove that deciding whether a δ -SR of size k exists is NP-hard. We reduce from the case $\delta = 1$, proved NP-hard by [9]. We assume of course that $\delta < 1$, as otherwise the result is already known.

Let (T, \mathbf{x}, k) be an input of the Minimum Sufficient Reason problem (i.e., $\delta = 1$), and let n be the dimension of T and \mathbf{x} . Assume without loss of generality that $T(\mathbf{x}) = 1$. If the given input of Minimum Sufficient Reason is

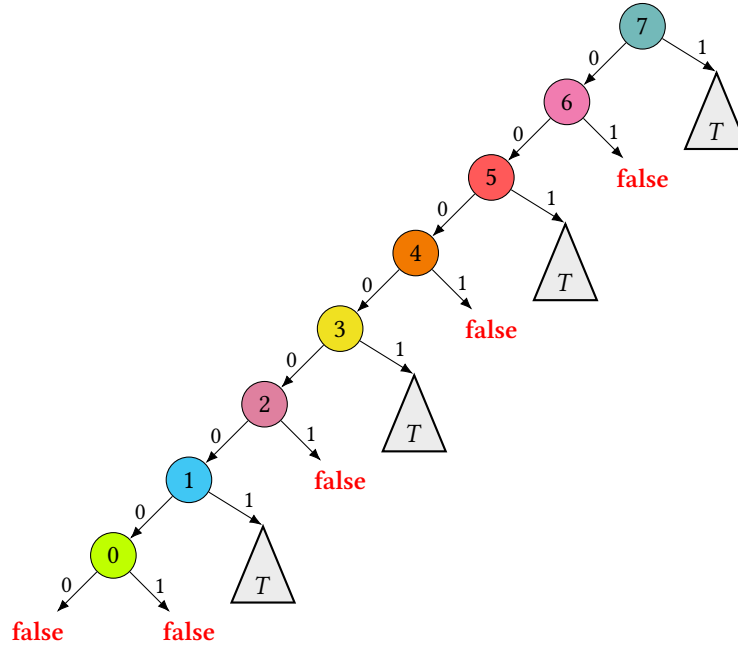


Fig. 4. Illustration of the construction of F_δ for $\delta = \frac{2}{3}$ and $n = 2$. Thus $2n + 3 + \lceil \log(1/\delta) \rceil = 8$ and $c = \lfloor \frac{2}{3} \cdot 2^8 \rfloor = 170 = 2^7 + 2^5 + 2^3 + 2^1$.

positive, then there is a partial instance $\mathbf{y} \subseteq \mathbf{x}$ with $|\mathbf{y}|_\perp \geq n - k$ such that $\Pr_z[T(z) = 1 \mid z \in \text{COMP}(\mathbf{y})] = 1$, and otherwise for every partial instance $\mathbf{y} \subseteq \mathbf{x}$ with $|\mathbf{y}|_\perp \geq n - k$ it holds that

$$\Pr_z [T(z) = 0 \mid z \in \text{COMP}(\mathbf{y})] \geq \frac{1}{2^n}.$$

Let us build a tree F_δ with $3n + 3 + \lceil \log(1/\delta) \rceil^4$ variables as follows. First build T_δ of dimension $2n + 3 + \lceil \log(1/\delta) \rceil$ by using Lemma 4.1, and then replace every true leaf of T_δ by a copy of T . Assume the $2n + 3 + \lceil \log(1/\delta) \rceil$ variables of T_δ are disjoint from the n variables that appear in T , and thus F_δ has the proposed number of variables. An example of the construction of F_δ is illustrated in Figure 4.

Define

$$\delta' := \Pr_z \left[T_\delta(z) = 1 \mid z \in \text{COMP} \left(\perp^{2n+2+\lceil \log(1/\delta) \rceil} \right) \right],$$

and recall that $|\delta' - \delta| \leq \frac{1}{2^{2n+3+\lceil \log(1/\delta) \rceil}}$. Now, let us build a final decision tree T^* with $4n + k + 4 + \lceil \log(1/\delta') \rceil + \lceil \log(1/\delta) \rceil$ variables as follows. Create $\ell := n + k + 1 + \lceil \log(1/\delta') \rceil$ vertices, labeled r_i for $i \in \{1, \dots, \ell\}$, and assume these labels are disjoint from the ones used in F_δ . Let r_1 be the root of T^* , and for each $i \in \{1, \dots, \ell - 1\}$, connect vertex labeled with r_i to vertex labeled with r_{i+1} using a 0-edge. The 0-edge from vertex labeled r_ℓ goes towards a leaf labeled with **true**. The 1-edge from every vertex r_i goes towards the root of a different copy of F_δ . Note that this construction, illustrated in Figure 5, takes polynomial time. Now, consider the instance \mathbf{x}^* that is defined (i) exactly as \mathbf{x} for the variables of T , (ii) exactly as in the instance \mathbf{x}^\dagger coming from Lemma 4.1 for the variables of T_δ in F_δ , and (iii) with all variables r_i set to 0. Note that $T^*(\mathbf{x}^*) = 1$. Now we prove both directions of the reduction separately. Assume first that the instance (T, \mathbf{x}, k) is a positive instance for Minimum Sufficient

⁴We use notation $\log(x)$ to refer to the logarithm in base 2 of x

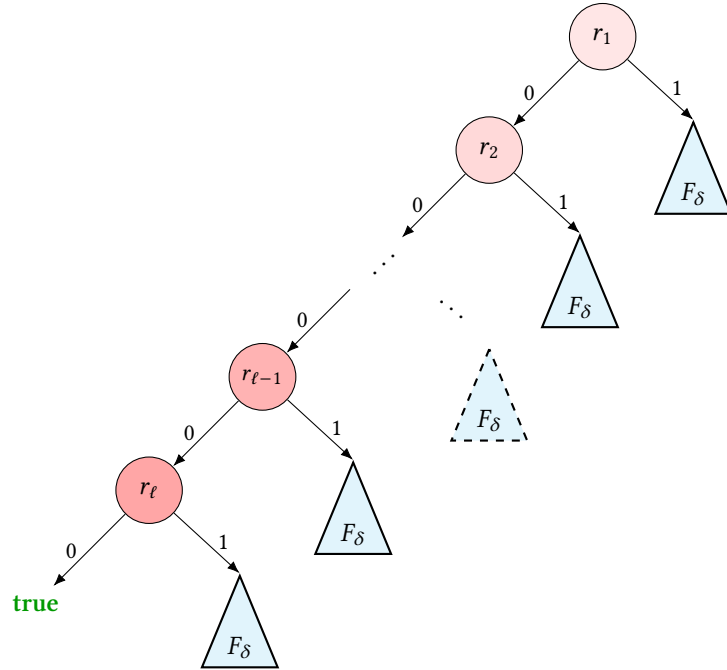


Fig. 5. Illustration of the construction of T^* . Recall that $\ell = n + k + 1 + \lceil \log(1/\delta') \rceil$.

Reason. Then we claim that there is a δ -SR for T^* of size at most k . Indeed, let $\mathbf{y} \subseteq \mathbf{x}$ be a sufficient reason for \mathbf{x} under T with at most k defined components. Then consider the partial instance $\mathbf{y}^* \subseteq \mathbf{x}^*$, that is only defined in the components where \mathbf{y} is defined. Now let us study $\Pr_z[T^*(z') = 1 \mid z \in \text{COMP}(\mathbf{y}^*)]$. The probability that z ends up in the true leaf on the 0-edge from vertex r_ℓ is $\frac{1}{2^\ell}$. In any other case, z takes a path that goes into a copy of F_δ , where its probability of acceptance is $\delta' \geq \delta - \frac{1}{2^{2n+3+\lceil \log(1/\delta) \rceil}}$ because of Lemma 4.1 and using that \mathbf{y}^* is undefined for all the variables of T_c . These two facts imply that

$$\Pr_z[T^*(z) = 1 \mid z \in \text{COMP}(\mathbf{y}^*)] \geq \frac{1}{2^\ell} + \delta - \frac{1}{2^{2n+3+\lceil \log(1/\delta) \rceil}}.$$

Now consider that

$$\begin{aligned} \delta &\leq \delta' + \frac{1}{2^{2n+3+\lceil \log(1/\delta) \rceil}} = \delta' + \frac{1}{2^{2n+3}} \cdot \frac{1}{2^{\lceil \log(1/\delta) \rceil}} \\ &\leq \delta' + \frac{1}{2^{2n+3}} \cdot \frac{1}{2^{\log(1/\delta)}} = \delta' + \frac{\delta}{2^{2n+3}}, \end{aligned}$$

from where $\delta \left(1 - \frac{1}{2^{2n+3}}\right) \leq \delta'$, and thus

$$\begin{aligned}
 \log(1/\delta) &\geq \log(1/\delta') + \log\left(1 - \frac{1}{2^{2n+3}}\right) \\
 &= \log(1/\delta') - \log\left(\frac{2^{2n+3}}{2^{2n+3} - 1}\right) \\
 &= \log(1/\delta') - 2n - 3 + \log(2^{2n+3} - 1) \\
 &\geq \log(1/\delta') - 2n - 3 + 2n + 2 && \text{(using } 2^{2n+3} - 1 \geq 2^{2n+2}\text{)} \\
 &= \log(1/\delta') - 1.
 \end{aligned}$$

From this we obtain that

$$\begin{aligned}
 \ell &= n + k + 1 + \lceil \log(1/\delta') \rceil \\
 &\leq 2n + 1 + \lceil \log(1/\delta') \rceil \\
 &\leq 2n + 1 + \log(1/\delta') + 1 \\
 &\leq 2n + 3 + \log(1/\delta) \\
 &\leq 2n + 3 + \lceil \log(1/\delta) \rceil,
 \end{aligned}$$

which allows us to conclude that

$$\Pr_z[T^*(z) = 1 \mid z \in \text{COMP}(\mathbf{y}^*)] \geq \frac{1}{2^\ell} + \delta - \frac{1}{2^{2n+3+\lceil \log(1/\delta) \rceil}} \geq \delta.$$

On the other hand, if (T, \mathbf{x}, k) is a negative instance for Minimum Sufficient Reason, consider any partial instance $\mathbf{y}^* \subseteq \mathbf{x}^*$ with at most k defined components, and note that by hypothesis we have that $\Pr_z[T(z) = 1 \mid z \in \text{COMP}(\mathbf{y}^*)] \leq 1 - \frac{1}{2^n}$. This implies, together with the second part of Lemma 4.1, that

$$\Pr_z[F_\delta(z) = 1 \mid z \in \text{COMP}(\mathbf{y}^*)] \leq \delta \left(1 - \frac{1}{2^n}\right),$$

and thus subsequently

$$\Pr_z[T^*(z) = 1 \mid z \in \text{COMP}(\mathbf{y}^*)] \leq \delta \left(1 - \frac{1}{2^n}\right) + \frac{1}{2^{\ell-k}},$$

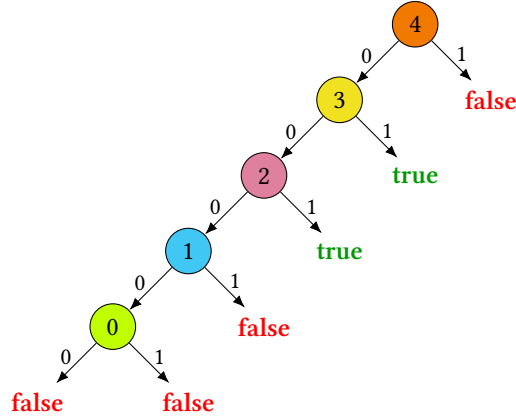


Fig. 6. Example of T_c , the tree constructed in the proof of Lemma 4.1, for $n = 5$ and $\delta = \frac{2}{5}$. In this case $c = 12 = 2^2 + 2^3$. Note that $\Pr_z[T_c(z) = 1 \mid z \in \text{COMP}(\perp^n)] = \frac{c}{2^n} = \frac{12}{32}$, and $|\frac{2}{5} - \frac{12}{32}| = \frac{1}{40} < \frac{1}{32}$.

by using that with at most k defined components in T^* , the probability of reaching the true leaf across the 0-edge from r_ℓ is at most $\frac{1}{2^{\ell-k}}$. To conclude, note that

$$\begin{aligned}
 \Pr_z[T^*(z) = 1 \mid z \in \text{COMP}(\mathbf{y}^*)] &\leq \delta \left(1 - \frac{1}{2^n}\right) + \frac{1}{2^{\ell-k}} \\
 &= \delta - \frac{\delta}{2^n} + \frac{1}{2^{n+1+\lceil \log(1/\delta') \rceil}} \\
 &\leq \delta - \frac{\delta}{2^n} + \frac{1}{2^{n+1+\log(1/\delta')}} \\
 &= \delta + \frac{(\delta' - \delta) - \delta}{2^{n+1}} \\
 &\leq \delta + \frac{\frac{1}{2^{2n+3+\lceil \log(1/\delta) \rceil}} - \delta}{2^{n+1}} \\
 &\leq \delta + \frac{\frac{\delta}{2^{2n+3}} - \delta}{2^{n+1}} \\
 &= \delta + \delta \left(\frac{-1 + \frac{1}{2^{2n+3}}}{2^{n+1}}\right) \\
 &< \delta.
 \end{aligned}$$

We have thus concluded that \mathbf{y}^* is a δ -SR for \mathbf{x}^* over T^* if and only if (T, \mathbf{x}, k) is a positive instance of Minimum Sufficient Reason, which completes our reduction. \square

Now we settle our debt by proving Lemma 4.1, the only part missing from the previous proof.

PROOF OF LEMMA 4.1. Let $c = \lfloor \delta 2^n \rfloor$, and note that $\delta - \frac{1}{2^n} \leq \frac{c}{2^n} \leq \delta$, from where $|\delta - \frac{c}{2^n}| \leq \frac{1}{2^n}$. This implies that we can prove the first part of the lemma by building in polynomial time a tree T_c over n variables, that has exactly c different positive instances, as then its probability of accepting a random completion of \perp^n will be exactly $\frac{c}{2^n}$. Note as well that $c < 2^n$ as $\delta < 1$.

As a first step, let us write c in binary, obtaining

$$c = \alpha_0 2^0 + \alpha_1 2^1 + \dots + \alpha_{n-1} 2^{n-1},$$

with $\alpha_i \in \{0, 1\}$ for each i . Then to build T_c start by creating n vertices, labeled 0 through $n - 1$. These n labels are the variables of T_c . For each $i \in \{1, \dots, n - 1\}$, connect vertex labeled i to vertex labeled $i - 1$ with a 0-edge, making vertex labeled $n - 1$ the root of T_c . Then, for each vertex with label $i \in \{0, \dots, n - 1\}$, set its 1-edge towards a leaf with label **true** if $\alpha_i = 1$, and towards a leaf with label **false** if $\alpha_i = 0$. The 0-edge of vertex labeled 0 goes towards a leaf with label **false**. Now let us count how many different positive instances does T_c have. We can do this by summing over all true leaves of T_c . Each true leaf comes from a 1-edge from a vertex labeled $i \in \{0, \dots, n - 1\}$. For every $i \in \{0, \dots, n - 1\}$, if the vertex labeled i has a true leaf when following its 1-edge, then the number of instances reaching that true leaf is exactly 2^i , as the variables whose value is not determined by the path to that leaf are those with labels less than i , which are exactly i variables. Therefore, the number of different positive instances of T_c along a 1-edge is the sum of 2^i for every i such that $\alpha_i = 1$, which is exactly c . An example is given in Figure 6. This concludes the proof of the first part of the lemma as the construction is clearly polynomial in n . For the second part, let us build \mathbf{x}^\dagger by setting $\mathbf{x}^\dagger[i] = 1 - \alpha_i$ for every $i \in \{0, \dots, n - 1\}$. In the example presented in Figure 6, we would build

$$\mathbf{x}^\dagger = (1, 1, 0, 0, 1).$$

We now prove that for any $\mathbf{y} \subseteq \mathbf{x}^\dagger$, it holds that

$$\Pr_z [T_c(\mathbf{z}) = 1 \mid \mathbf{z} \in \text{COMP}(\mathbf{y})] \leq \Pr_z [T_c(\mathbf{z}) = 1 \mid \mathbf{z} \in \text{COMP}(\perp^n)].$$

We do this via a finite induction argument by strengthening our induction hypothesis; for $i \in \{0, \dots, n - 1\}$, let T_c^i be the sub-tree of T_c rooted at the vertex labeled i , and let us claim that for every $i \in \{0, \dots, n - 1\}$ we have that

$$\Pr_z [T_c^i(\mathbf{z}) = 1 \mid \mathbf{z} \in \text{COMP}(\mathbf{y})] \leq \Pr_z [T_c^i(\mathbf{z}) = 1 \mid \mathbf{z} \in \text{COMP}(\perp^n)],$$

which implies what we want to show when taking $i = n - 1$. The base case of the induction is when $i = 0$, in which case the claim trivially holds as if $\mathbf{y}[0] = \perp$ we have equality, and if $\mathbf{y}[0] \neq \perp$ then by construction

$$\Pr_z [T_c^0(\mathbf{z}) = 1 \mid \mathbf{z} \in \text{COMP}(\mathbf{y})] = 0.$$

For the inductive case, let $i > 0$, and proceed by cases; if $\mathbf{y}[i] = \perp$, then by letting $t_i \in \{0, 1\}$ be an indicator variable for whether the leaf across the 1-edge from vertex i is labeled **true** we have that

$$\begin{aligned} \Pr_z [T_c^i(\mathbf{z}) = 1 \mid \mathbf{z} \in \text{COMP}(\mathbf{y})] &= \frac{1}{2} t_i + \frac{1}{2} \Pr_z [T_c^{i-1}(\mathbf{z}) = 1 \mid \mathbf{z} \in \text{COMP}(\mathbf{y})] \\ &\leq \frac{1}{2} t_i + \frac{1}{2} \Pr_z [T_c^{i-1}(\mathbf{z}) = 1 \mid \mathbf{z} \in \text{COMP}(\perp^n)] \\ &= \Pr_z [T_c^i(\mathbf{z}) = 1 \mid \mathbf{z} \in \text{COMP}(\perp^n)], \end{aligned}$$

where the inequality has used the inductive hypothesis. On the other hand, if $\mathbf{y}[i] = 1$, that implies $\mathbf{x}^\dagger[i] = 1$ and thus $\alpha_i = 0$, which means the leaf across the 1-edge from vertex i is labeled with **false**, and thus

$$\Pr_z [T_c^i(\mathbf{z}) = 1 \mid \mathbf{z} \in \text{COMP}(\mathbf{y})] = 0,$$

which trivially satisfies the claim. For the last case, if $\mathbf{y}[i] = 0$, then $\mathbf{x}^\dagger[i] = 0$ and thus $\alpha_i = 1$, which means the leaf across the 1-edge from vertex i is labeled with **true**. Therefore, we have

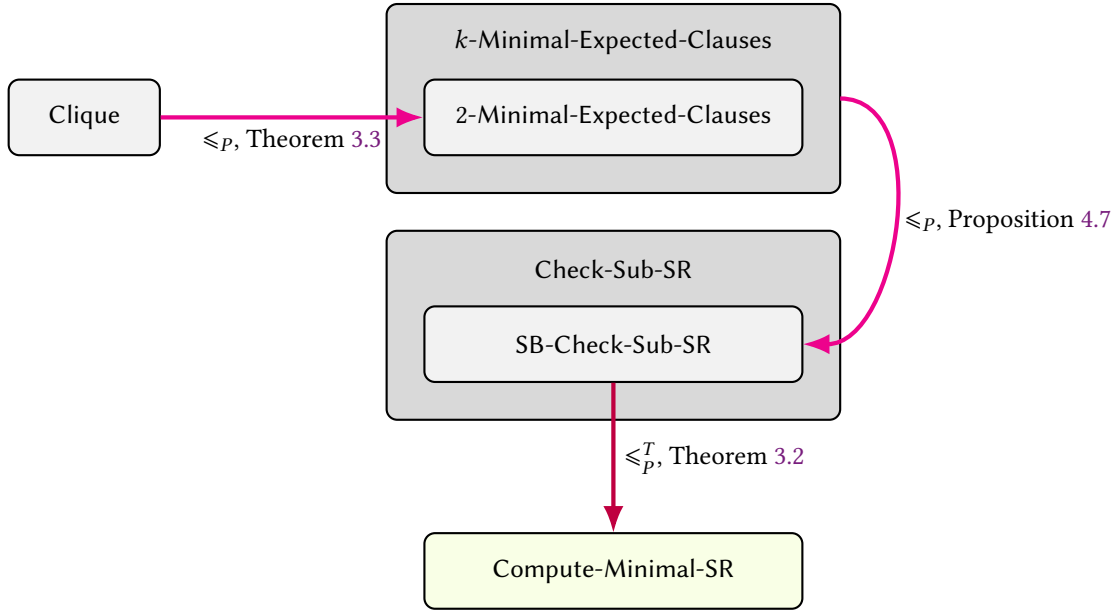


Fig. 7. Roadmap for the proof of Theorem 3.2. Arrows denote polynomial-time many-to-one reductions (\leq_p) or Turing reductions (\leq_p^T).

$$\begin{aligned}
 \Pr_z [T_c^i(z) = 1 \mid z \in \text{COMP}(\mathbf{y})] &= \Pr_z [T_c^{i-1}(z) = 1 \mid z \in \text{COMP}(\mathbf{y})] \\
 &\leq \Pr_z [T_c^{i-1}(z') = 1 \mid z \in \text{COMP}(\perp^n)] \\
 &\leq \frac{1}{2} + \frac{1}{2} \Pr_z [T_c^{i-1}(z) = 1 \mid z \in \text{COMP}(\perp^n)] \quad (\text{as } \Pr[\cdot] \leq 1) \\
 &= \Pr_z [T_c^i(z) = 1 \mid z \in \text{COMP}(\perp^n)].
 \end{aligned}$$

This completes the induction argument, and thus we conclude the proof of the lemma. □

Our goal now is to prove Theorem 3.2. Figure 7 presents a roadmap of the proof with its corresponding chain of reductions. Following the roadmap, we need first Theorem 3.3, which we restate and prove next.

THEOREM 3.3 (RESTATED). *k*-Minimal-Expected-Clauses is NP-hard for every $k \geq 2$.

PROOF OF THEOREM 3.3. We start with some simple observations regarding $E(\varphi, \mu)$, for a CNF formula φ and a partial assignment μ . By linearity of expectation, we can write $E(\varphi, \mu)$ as the sum

$$E(\varphi, \mu) = \sum_{C \text{ clause of } \varphi} \text{Prob}_{C,\mu}, \tag{4.2}$$

where $\text{Prob}_{C,\mu}$ is the probability that a random completion of μ satisfies the clause C . In turn, the probabilities $\text{Prob}_{C,\mu}$ can be easily computed as:

- $\text{Prob}_{C,\mu} = 1$, if there is a positive literal x in C with $\mu(x) = 1$ or there a negative literal $\neg x$ in C with $\mu(x) = 0$.

- $\text{Prob}_{C,\mu} = 1 - \frac{1}{2^\eta}$, where η is the number of literals in C of the form x or $\neg x$ with $\mu(x) = \perp$.

Note that for an assignment σ , $E(\varphi, \sigma)$ is simply the number of clauses of φ satisfied by σ .

We reduce from the Clique problem. Recall this problem asks, given a graph G and an integer $k \geq 3$, whether there is a *clique* of size k , that is, a set K of k vertices such that there is an edge between any pair of distinct vertices from K . Let G be a graph and $k \geq 3$. We can assume without loss of generality that k is odd and the degree of every vertex x of G , denoted by $\deg(x)$, is at least $k - 1$; if k is even we can consider the equivalent instance given by the graph G' that extends G with a fresh node connected via an edge with all the other nodes and $k' = k + 1$. On the other hand, we can iteratively remove vertices of degree less than $k - 1$ as those cannot be part of any clique of size k . Starting from the Clique instance (G, k) we will construct an instance (φ, σ) for 2-Minimal-Expected-Clauses as follows. The variables of φ are the nodes of G . For each variable x we have the following clauses in φ :

- A clause $A_x = (\neg x)$. This clause A_x is repeated $\deg(x) - \frac{k-1}{2}$ times in φ . Note this quantity is always a positive integer since k is odd.
- A set of clauses $\mathcal{B}_x = \{(x \vee \neg y_1), \dots, (x \vee \neg y_{\deg(x)})\}$, where $y_1, \dots, y_{\deg(x)}$ are the neighbors of x in G . Each clause in \mathcal{B}_x appears only once in φ .

Additionally, for each set $\{x, y\}$ where $x \neq y$ and $\{x, y\}$ is *not* an edge in G , we have a clause $Z_{x,y} = (x \vee y)$ repeated $4e$ times in φ , where e is the number of edges in G .

We define the assignment σ such that $\sigma(x) = 1$, for all variables x of φ .

For an arbitrary partial assignment μ to the variables of φ , with $\mu \subseteq \sigma$, we define

$$\text{utility}_{\varphi,\sigma}(\mu) := E(\varphi, \mu) - E(\varphi, \sigma).$$

In particular, the instance (φ, σ) is a Yes-instance for 2-Minimal-Expected-Clauses if and only if there is $\mu \subseteq \sigma$ with $\text{utility}_{\varphi,\sigma}(\mu) \geq 0$. By equation (4.2), we can write

$$\text{utility}_{\varphi,\sigma}(\mu) = \sum_{C \text{ clause of } \varphi} \text{utility}_{\varphi,\sigma}(\mu, C),$$

where $\text{utility}_{\varphi,\sigma}(\mu, C)$ is defined as:

$$\text{utility}_{\varphi,\sigma}(\mu, C) := \text{Prob}_{C,\mu} - \text{Prob}_{C,\sigma}.$$

Now observe that we have:

$$\text{Prob}_{C,\sigma} = \begin{cases} 0 & \text{if } C = A_x \text{ for some variable } x \\ 1 & \text{if } C \in \mathcal{B}_x \text{ for some variable } x \text{ or } C = Z_{x,y} \text{ for some set } \{x, y\} \end{cases}$$

On the other hand, for the probability $\text{Prob}_{C,\mu}$ we have the following:

- (1) Assume $C = A_x = (\neg x)$ for some variable x . Then $\text{Prob}_{C,\mu}$ is
 - (a) $\frac{1}{2}$, if $\mu(x) = \perp$ (and hence $\text{utility}_{\varphi,\sigma}(\mu, C) = \frac{1}{2}$), and
 - (b) 0 otherwise (then $\text{utility}_{\varphi,\sigma}(\mu, C) = 0$).
- (2) Suppose $C = (x \vee \neg y) \in \mathcal{B}_x$ for some variable x . Then $\text{Prob}_{C,\mu}$ is
 - (a) $\frac{3}{4}$ if $\mu(x) = \perp$ and $\mu(y) = \perp$ (and hence $\text{utility}_{\varphi,\sigma}(\mu, C) = -\frac{1}{4}$),
 - (b) $\frac{1}{2}$ if $\mu(x) = \perp$ and $\mu(y) = 1$ (and then $\text{utility}_{\varphi,\sigma}(\mu, C) = -\frac{1}{2}$), and
 - (c) 1 otherwise (then $\text{utility}_{\varphi,\sigma}(\mu, C) = 0$).
- (3) Suppose $C = Z_{x,y} = (x \vee y)$ for some set $\{x, y\}$. Then $\text{Prob}_{C,\mu}$ is
 - (a) $\frac{3}{4}$ if $\mu(x) = \perp$ and $\mu(y) = \perp$ (and hence $\text{utility}_{\varphi,\sigma}(\mu, C) = -\frac{1}{4}$), and
 - (b) 1 otherwise (then $\text{utility}_{\varphi,\sigma}(\mu, C) = 0$).

We now show the correctness of our construction, that is, G has a clique of size $k \geq 3$ if and only if (φ, σ) is a Yes-instance for the problem 2-Minimal-Expected-Clauses.

(\Rightarrow) Suppose G has a clique K of size $k \geq 3$. Let μ be the partial assignment that sets $\mu(x) = \perp$ if $x \in K$ and $\mu(x) = 1$ if $x \notin K$. Note that $\mu \subseteq \sigma$. We claim that $\text{utility}_{\varphi, \sigma}(\mu) = 0$ and hence (φ, σ) is a Yes-instance. Let C be a clause in φ . If C is of the form $Z_{x,y}$, then $\text{utility}_{\varphi, \sigma}(\mu, C) = 0$. Indeed, by construction, $\{x, y\}$ is not an edge, and since K is a clique, then $\mu(x) = 1$ or $\mu(y) = 1$. This means we are always in case 3(b) above. If $x \notin K$ and C is of the form A_x or belongs to \mathcal{B}_x , then $\text{utility}_{\varphi, \sigma}(\mu, C) = 0$, since $\mu(x) = 1$ and hence we fall either in case 1(b) or 2(c) above. It follows that $\text{utility}_{\varphi, \sigma}(\mu)$ is the sum of the utilities of all the clauses involved with variables $x \in K$. That is:

$$\text{utility}_{\varphi, \sigma}(\mu) = \sum_{x \in K} \left(\deg(x) - \frac{k-1}{2} \right) \text{utility}_{\varphi, \sigma}(\mu, A_x) + \sum_{x \in K} \sum_{C \in \mathcal{B}_x} \text{utility}_{\varphi, \sigma}(\mu, C). \quad (4.3)$$

Take $x \in K$. Then $\text{utility}_{\varphi, \sigma}(\mu, A_x) = \frac{1}{2}$ as $\mu(x) = \perp$, and then case 1(a) applies. On the other hand, for a clause $C \in \mathcal{B}_x$ we have two cases:

- $C = (x \vee \neg y)$ for $y \in K$. In this case, $\text{utility}_{\varphi, \sigma}(\mu, C) = -\frac{1}{4}$ as we are in case 2(a) above.
- $C = (x \vee \neg y)$ for $y \notin K$. In this case, $\text{utility}_{\varphi, \sigma}(\mu, C) = -\frac{1}{2}$ as we are in case 2(b) above.

Moreover, note that the first case occurs exactly for $k-1$ clauses in \mathcal{B}_x , as x has precisely $k-1$ neighbors in the clique K . The second case occurs exactly for $\deg(x) - (k-1) \geq 0$ clauses in \mathcal{B}_x . Replacing in equation (4.3), we obtain:

$$\begin{aligned} \text{utility}_{\varphi, \sigma}(\mu) &= \sum_{x \in K} \left(\frac{k-1}{4} + \frac{\deg(x)}{2} - \frac{k-1}{2} \right) + \left(-\frac{k-1}{4} - \frac{\deg(x)}{2} + \frac{k-1}{2} \right) \\ &= 0. \end{aligned}$$

We conclude that (φ, σ) is a Yes-instance.

(\Leftarrow) Suppose now that there is a partial assignment μ , with $\mu \subseteq \sigma$ and $\text{utility}_{\varphi, \sigma}(\mu) \geq 0$. Let K be the set of variables x such that $\mu(x) = \perp$. For $x \notin K$ and $C = A_x$ or $C \in \mathcal{B}_x$, we have $\text{utility}_{\varphi, \sigma}(\mu, C) = 0$, as we are in cases 1(b) or 2(c) above. Then we can write:

$$\begin{aligned} \text{utility}_{\varphi, \sigma}(\mu) &= \sum_{x \in K} \left[\left(\deg(x) - \frac{k-1}{2} \right) \text{utility}_{\varphi, \sigma}(\mu, A_x) + \sum_{C \in \mathcal{B}_x} \text{utility}_{\varphi, \sigma}(\mu, C) \right] \\ &\quad + \sum_{\{x, y\} \text{ non-edge}} 4e(\text{utility}_{\varphi, \sigma}(\mu, Z_{x,y})). \end{aligned} \quad (4.4)$$

We claim that $|K| \geq k$. Towards a contradiction, suppose $|K| = \ell < k$. Since $\text{utility}_{\varphi, \sigma}(\mu, Z_{x,y}) \leq 0$ for every pair $\{x, y\}$, the last term in equation (4.4) is ≤ 0 , and then:

$$\text{utility}_{\varphi, \sigma}(\mu) \leq \sum_{x \in K} \left(\frac{k-1}{2} + \deg(x) - (k-1) \right) \text{utility}_{\varphi, \sigma}(\mu, A_x) + \sum_{x \in K} \sum_{C \in \mathcal{B}_x} \text{utility}_{\varphi, \sigma}(\mu, C). \quad (4.5)$$

Take $x \in K$. Following the same argument as before, we have that

$$\text{utility}_{\varphi, \sigma}(\mu, A_x) = \frac{1}{2},$$

and for a clause $C \in \mathcal{B}_x$ we have the two cases:

- $C = (x \vee \neg y)$ for $y \in K$, and $\text{utility}_{\varphi, \sigma}(\mu, C) = -\frac{1}{4}$.
- $C = (x \vee \neg y)$ for $y \notin K$, and $\text{utility}_{\varphi, \sigma}(\mu, C) = -\frac{1}{2}$.

Let say the first case occurs precisely for r clauses from \mathcal{B}_x . Then:

$$\sum_{C \in \mathcal{B}_x} \text{utility}_{\varphi, \sigma}(\mu, C) = -\frac{r}{4} - \frac{\deg(x) - r}{2} = \frac{r}{4} - \frac{\deg(x)}{2}. \quad (4.6)$$

Note that $r \leq \ell - 1$ and from equation (4.6) we obtain (recall $\ell < k$):

$$\sum_{C \in \mathcal{B}_x} \text{utility}_{\varphi, \sigma}(\mu, C) \leq \frac{\ell - 1}{4} - \frac{\deg(x)}{2} < \frac{k - 1}{4} - \frac{\deg(x)}{2}.$$

Replacing in equation (4.5), we obtain:

$$\begin{aligned} \text{utility}_{\varphi, \sigma}(\mu) &< \sum_{x \in K} \left[\left(\frac{k - 1}{4} + \frac{\deg(x)}{2} - \frac{k - 1}{2} \right) + \frac{k - 1}{4} - \frac{\deg(x)}{2} \right] \\ &= \sum_{x \in K} \left[\frac{\deg(x)}{2} - \frac{k - 1}{4} + \frac{k - 1}{4} - \frac{\deg(x)}{2} \right] \\ &= 0. \end{aligned}$$

We conclude that $\text{utility}_{\varphi, \sigma}(\mu) < 0$ which is a contradiction. Hence $|K| \geq k$.

Finally, we show that K is a clique. By contradiction, assume there is a pair $\{\tilde{x}, \tilde{y}\}$ such that $\tilde{x} \neq \tilde{y}$, $\tilde{x}, \tilde{y} \in K$ and $\{\tilde{x}, \tilde{y}\}$ is not an edge in G . Then there is a clause $Z_{\tilde{x}, \tilde{y}}$ which is repeated e times in φ . Since $\mu(\tilde{x}) = \perp$ and $\mu(\tilde{y}) = \perp$, we have $\text{utility}_{\varphi, \sigma}(\mu, Z_{\tilde{x}, \tilde{y}}) = -\frac{1}{4}$, as we are in case 3(a) above. As $\text{utility}_{\varphi, \sigma}(\mu, Z_{x, y}) \leq 0$ for all pairs $\{x, y\}$, we obtain:

$$\sum_{\{x, y\} \text{ non-edge}} 4e (\text{utility}_{\varphi, \sigma}(\mu, Z_{x, y})) \leq 4e (\text{utility}_{\varphi, \sigma}(\mu, Z_{\tilde{x}, \tilde{y}})) \leq -e$$

For $x \in K$, since $\text{utility}_{\varphi, \sigma}(\mu, C) \leq 0$, for all $C \in \mathcal{B}_x$, we have $\sum_{C \in \mathcal{B}_x} \text{utility}_{\varphi, \sigma}(\mu, C) \leq 0$ and thus

$$\sum_{x \in K} \sum_{C \in \mathcal{B}_x} \text{utility}_{\varphi, \sigma}(\mu, C) \leq 0.$$

On the other hand, for $x \in K$, we have $\text{utility}_{\varphi, \sigma}(\mu, A_x) = \frac{1}{2}$. Combining all this with equation (4.4) we obtain:

$$\begin{aligned} \text{utility}_{\varphi, \sigma}(\mu) &\leq \sum_{x \in K} \left(\frac{\deg(x)}{2} - \frac{k - 1}{4} \right) - e \\ &< \sum_{x \in K} \frac{\deg(x)}{2} - e \quad (\text{since } k \geq 3) \\ &\leq \sum_{x \text{ in } G} \frac{\deg(x)}{2} - e \\ &= 0. \end{aligned}$$

We conclude that $\text{utility}_{\varphi, \sigma}(\mu) < 0$, and thus obtain a contradiction. Therefore, G contains a clique of size k . \square

We now continue towards the proof of Theorem 3.2. The next step is to show that the following decision problem is NP-hard.

PROBLEM : Check-Sub-SR
INPUT : (T, \mathbf{y}) , for a decision tree T of dimension n , and partial instance $\mathbf{y} \in \{0, 1, \perp\}^n$
OUTPUT : Yes, if there is a partial instance $\mathbf{y}' \subseteq \mathbf{y}$ such that $\Pr_z[T(z) = 1 \mid z \in \text{COMP}(\mathbf{y}')] \geq \Pr_z[T(z) = 1 \mid z \in \text{COMP}(\mathbf{y})]$, and No otherwise

Proposition 4.7. Check-Sub-SR is NP-hard.

PROOF. We reduce from 2-Minimal-Expected-Clauses. Let (φ, σ) be an instance of 2-Minimal-Expected-Clauses. Let m be the number of clauses of φ and assume that φ has n variables x_1, \dots, x_n . Without loss of generality we assume that m is a power of 2. Define a decision tree T of dimension $n + m - 1$ as follows. Start with a *perfect* binary tree S of depth $\log_2 m$, that is, each internal node has two children, and each leaf is at depth $\log_2 m$. In particular, S has m leaves and $m - 1$ internal nodes. All the internal nodes of S are labeled with a different fresh feature from $\{n + 1, \dots, n + m - 1\}$. For each clause C in φ , pick a different leaf ℓ_C of S . It is easy to see that for each clause C we can define a decision tree S_C over the features $\{1, \dots, n\}$ such that for every assignment $\mu : \{x_1, \dots, x_n\} \rightarrow \{0, 1\}$ to the variables of φ , the corresponding instance $\mathbf{x} \in \{0, 1\}^n$ where $x[i] = \mu(x_i)$ satisfies that $S_C(\mathbf{x}) = 1$ if and only if μ satisfies C . The decision tree T is obtained from S by identifying for each clause C , the leaf ℓ_C with the root of the decision tree S_C .

For any partial assignment $\mu : \{x_1, \dots, x_n\} \rightarrow \{0, 1, \perp\}$ for φ , we denote by \mathbf{y}_μ the partial instance of dimension $n + m - 1$ such that $\mathbf{y}_\mu[i] = \mu(x_i)$ for every $i \in \{1, \dots, n\}$ and $\mathbf{y}_\mu[i] = \perp$ for every $i \in \{n + 1, \dots, n + m - 1\}$. The output of the reduction is (T, \mathbf{y}_σ) . Observe that the transformation from μ to \mathbf{y}_μ is a bijection between the sets $\{\mu \mid \mu \subseteq \sigma\}$ and $\{\mathbf{y}_\mu \mid \mathbf{y}_\mu \subseteq \mathbf{y}_\sigma\}$. By construction, for any partial assignment $\mu \subseteq \sigma$, we have that

$$\Pr_z[T(z) = 1 \mid z \in \text{COMP}(\mathbf{y}_\mu)] = \frac{E(\varphi, \mu)}{m}.$$

Hence (φ, σ) is a Yes-instance of 2-Minimal-Expected-Clauses if and only if (T, \mathbf{y}_σ) is a Yes-instance for the Check-Sub-SR problem. \square

Remark. We can assume that the instance (T, \mathbf{y}_σ) constructed in the proof of Proposition 4.7 satisfies

$$\Pr_z[T(z) = 1 \mid z \in \text{COMP}(\mathbf{y}_\sigma)] > \frac{1}{2}.$$

Indeed, the above probability is simply $\frac{E(\varphi, \sigma)}{m}$, where m is the number of clauses of φ . On the other hand, from the proof of Theorem 3.3, we can choose σ such that $\sigma(x_i) = 1$ for every variable $x_i \in \{x_1, \dots, x_n\}$ of φ . It follows that $E(\varphi, \sigma)$ is simply the number of clauses satisfied by σ , which are all the clauses in \mathcal{B}_x for some variable x , and all the clauses of the form $Z_{x,y}$. Note that the total number of clauses from the sets \mathcal{B}_x is greater than the total number of clauses of the form A_x , and hence $\frac{E(\varphi, \sigma)}{m} > \frac{1}{2}$. Indeed, there are $\deg(x)$ clauses in \mathcal{B}_x , and summing over all the variables x , we obtain $2e$, where e are the number of edges in the graph G . On the other hand, each clause A_x is repeated $\frac{k-1}{2} + \deg(x) - (k-1) = \deg(x) - \frac{k-1}{2}$ times. Taking the sum over all the variables x , we obtain $2e - n \left(\frac{k-1}{2} \right) < 2e$. This property will be useful for the proof of Theorem 3.2, which we recall next.

THEOREM 3.2 (RESTATED). *Assuming that $P \neq NP$, the Compute-Minimal-SR problem cannot be solved in polynomial time.*

We show that if Compute-Minimal-SR admits a polynomial time algorithm, then Check-Sub-SR is in P, which contradicts the assumption that $P \neq NP$. In particular, we show a Turing-reduction from a variant of

Check-Sub-SR to Compute-Minimal-SR, thus establishing that the latter cannot be solved in polynomial time unless $P = NP$.

For the sake of readability, given a partial instance \mathbf{y} , in this proof we use notation $z \sim \mathbb{U}(\mathbf{y})$ to indicate that z is generated uniformly at random from the set $\text{COMP}(\mathbf{y})$. For instance, we obtain the following simplification by using this terminology:

$$\Pr_z [T(z) = 1 \mid z \in \text{COMP}(\mathbf{y})] = \Pr_{z \sim \mathbb{U}(\mathbf{y})} [T(z) = 1]$$

We require a particular kind of hard instances for the Check-Sub-SR to make our reduction work. In particular, we now define the notion of *strongly-balanced* inputs, which intuitively captures the idea that defined features in a partial instance \mathbf{y} appear at the same depth in different branches of the decision tree T . To make this definition precise, consider that every path π from the root to a leaf in a decision tree can be identified with a sequence of labels s_π corresponding to the labels of the nodes of π , where the last label of π is either **true** or **false**. We use notation $s_\pi[i]$ for the i -th label in the sequence s_π . With this notation, we can introduce the following definition.

Definition 4.8. *Given a decision tree T of dimension d and $\mathbf{y} \in \{0, 1, \perp\}^n$ a partial instance, we say that the pair (T, \mathbf{y}) is strongly-balanced if*

$$\Pr_{z \sim \mathbb{U}(\mathbf{y})} [T(z) = 1] > \frac{1}{2},$$

and there exists $k \in \mathbb{N}$ such that for every root-to-leaf path π in T , the sequence s_π satisfies

$$\mathbf{y}[s_\pi[i]] = \perp \iff i \leq k.$$

If (T, \mathbf{y}) is strongly-balanced, then there exists a unique value $k \in \mathbb{N}$ that satisfies the second condition of the definition. We denote this value by $u(T, \mathbf{y})$. In particular, if $\mathbf{y} \in \{0, 1\}^n$, then (T, \mathbf{y}) is strongly-balanced and $u(T, \mathbf{y}) = 0$.

Now let us define the following problem.

PROBLEM : SB-Check-SUB-SR
 INPUT : (T, \mathbf{y}) , for T a decision tree of dimension n and $\mathbf{y} \in \{0, 1, \perp\}^n$ a partial instance, where (T, \mathbf{y}) is strongly-balanced.
 OUTPUT : Yes, if there is a partial instance $\mathbf{y}' \subseteq \mathbf{y}$ such that $\Pr_{z \sim \mathbb{U}(\mathbf{y}')} [T(z) = 1] \geq \Pr_{z \sim \mathbb{U}(\mathbf{y})} [T(z) = 1]$, and No otherwise.

One can now check that the proof of Proposition 4.7 directly proves NP-hardness for this problem, and thus we can reduce from it to prove hardness for the computation variant. Indeed, the first part of the definition of strongly-balanced follows from the remark given after the proof of Proposition 4.7. The second part follows from the fact that the construction in the proof of Proposition 4.7 starts with a perfect binary tree.

We can now start proving Theorem 3.2, which requires a few auxiliary claims that for the sake of readability will be assumed within this proof, and then we will conclude this section by proving each claim individually.

PROOF OF THEOREM 3.2. Assume an instance (T, \mathbf{y}) for the SB-Check-Sub-SR problem. Let us enumerate the features in T as $1, \dots, n$. Also, let S be the set of features defined in \mathbf{y} , that is, $\mathbf{y}[i] \neq \perp \iff i \in S$. We first build a decision tree T' of dimension $2n - |S|$, with the following features:

- (1) Create a feature i for $i \in S$.
- (2) For every $i \in \{1, \dots, n\} \setminus S$ create features i and i' .

Note that this amounts to the promised number of features. We build T' through a recursive process \mathcal{R} defined next. First, note that any decision tree can be described inductively as either a **true/false** leaf, or a tuple (r, L, R) , where r is the root node, L is a decision tree whose root is the left child of r , and R is a decision tree whose root

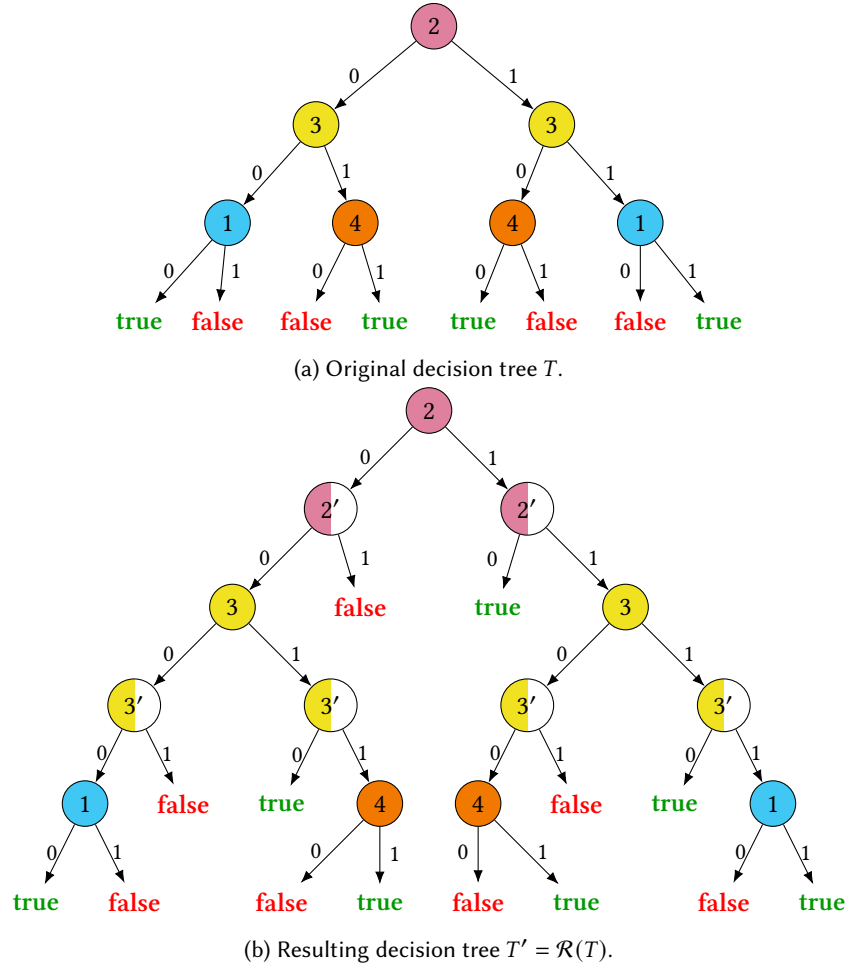


Fig. 8. Illustration of the recursive process \mathcal{R} over an example where $\mathbf{y} = (0, \perp, \perp, 0)$. Note that the pair (T, \mathbf{y}) is strongly-balanced.

is the right child of r . We can now define \mathcal{R} as a recursive procedure that when called with argument τ proceed as follows:

- (1) If τ is a leaf, then simply return τ .
- (2) If $\tau = (r, L, R)$, and node r is labeled with feature $i \in S$, then simply return

$$(r, \mathcal{R}(L), \mathcal{R}(R)).$$

- (3) If $\tau = (r, L, R)$, and node r is labeled with feature $i \in \{1, \dots, n\} \setminus S$, then return the following decision tree:

$$(r, (u, \mathcal{R}(L), \mathbf{false}), (v, \mathbf{true}, \mathcal{R}(R))),$$

where nodes u and v are new nodes, both labeled with feature i' .

As anticipated, $T' = \mathcal{R}(T)$. An example illustrating the process is depicted in Figure 8. Now we create a tree T'' of dimension $2n - |S| + m$ that on top of the previous features incorporates features b_j for each $j \in \{1, \dots, m\}$, where m is an integer we specify later on. In order to construct T'' , we start by defining \mathbf{y}_0 as the partial instance of dimension $2n - |S| + m$ such that $\mathbf{y}_0[i] = \mathbf{y}[i]$ for every $i \in S$ and

$$\mathbf{y}_0[b_j] = 0, \quad \forall j \in \{1, \dots, m\},$$

with the remaining components of \mathbf{y}_0 being left undefined. Let $T_{\mathbf{y}_0}$ be a tree of dimension $2n - |S| + m$ that accepts exactly the completions of \mathbf{y}_0 ; this can be trivially done by creating a tree that accepts exactly the features that are defined in \mathbf{y}_0 , and then observing that when running an instance whose feature space is a superset of this, then the instance is accepted if and only if it is a completion of \mathbf{y}_0 . Now let T_1 be a tree of dimension $2n + |S| + m$ that implements the following Boolean formula:

$$\phi = \sum_{j=1}^m b_j \geq 2.$$

Claim 4.9. *Decision tree T_1 , implementing the function ϕ , can be constructed in polynomial time.*

Now, let us build an instance \mathbf{x} of dimension $2n - |S| + m$ as follows. For each $i \in S$ let $\mathbf{x}[i] = \mathbf{y}[i]$, thus ensuring \mathbf{x} is a completion of \mathbf{y} . Then for each $j \in \{1, \dots, m\}$ let $\mathbf{x}[b_j] = 0$, and finally for each $i \in \{1, \dots, n\} \setminus S$, let $\mathbf{x}[i] = 0$ and $\mathbf{x}[i'] = 1$. For example, if $\mathbf{y} = (0, \perp, \perp, 1)$, and $m = 3$ then

$$\mathbf{x} = (0, 0, 1, 0, 1, 1, 0, 0, 0),$$

where the features b_j have been placed at the end of the vector.

Let \mathbf{y}^* be the partial instance of dimension $2n - |S| + m$ such that $\mathbf{y}^*[i] = \mathbf{y}[i]$ for every $i \in S$, and undefined otherwise. Let us abuse notation and assume now that T' has dimension $2n - |S| + m$, even though it only explicitly uses the first $2n - |S|$ features, as this would make it compatible with other decision trees and instance we have constructed. Finally, let T^* be the decision tree defined as

$$T^* = T_{\mathbf{y}_0} \cup (T' \cap T_1),$$

and note that the union and intersection of decision trees can be computed in polynomial time through a standard algorithm (see e.g., [25]).

Let us now define

$$\delta = \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [(T' \cap T_1)(z) = 1].$$

We now claim that the result of $\text{Compute-Minimal-SR}(T^*, \mathbf{x}, \delta)$ is different from \mathbf{y}^* if and only if (T, \mathbf{y}) is a positive instance of SB-Check-SUB-SR. But before we can prove this, we need some intermediary tools and claims that we develop next.

Let us start by distinguishing two kinds of leaves of T' . Let us say that the leaves of T' introduced in step 1 of the recursive procedure \mathcal{R} are *natural*, while those introduced in step 3 are *artificial*. We denote by \mathcal{N} the set of natural leaves of T' and by \mathcal{A} the set of artificial leaves of T' . Moreover, let $\mathcal{N}_t, \mathcal{N}_f$ represent the **true** and **false** natural leaves, and define $\mathcal{A}_t, \mathcal{A}_f$ analogously. We also use $T'_{\downarrow z}$ to denote the leaf where instance z ends when evaluated over tree T' . With this notation, $T'(z) = 1$ is equivalent to $T'_{\downarrow z} \in \mathcal{A}_t \cup \mathcal{N}_t$. We now make the following claims.

Claim 4.10. *For every partial instance $\mathbf{y}'^* \subseteq \mathbf{y}^*$, it holds that*

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}'^*)} [T'_{\downarrow z} \in \mathcal{A}_t \mid T'_{\downarrow z} \in \mathcal{A}] = \Pr_{z \sim \mathbb{U}(\mathbf{y}'^*)} [T'_{\downarrow z} \in \mathcal{A}_f \mid T'_{\downarrow z} \in \mathcal{A}] = \frac{1}{2}.$$

Claim 4.11. Given a partial instance $\mathbf{y}' \subseteq \mathbf{y}$, we define \mathbf{y}'^* as the partial instance of dimension $2n - |S| + m$ that matches \mathbf{y}' on its defined features, and holds $\mathbf{y}'^*[i] = \mathbf{y}'^*[i'] = \perp$ for every feature i such that $\mathbf{y}'[i] = \perp$. Then it holds that

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}')} [T(z) = 1] = \Pr_{z \sim \mathbb{U}(\mathbf{y}'^*)} \left[T'_{\downarrow z} \in \mathcal{N}_t \mid T'_{\downarrow z} \in \mathcal{N} \right].$$

Claim 4.12. By choosing $m \geq \max\{2u(T, \mathbf{y}) + 2n, 9\}$, we have that

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [(T' \cap T_1)(z) = 1] > \frac{1}{2}.$$

Claim 4.13. For every partial instance $\mathbf{y}'^* \subseteq \mathbf{y}^*$, it holds that

$$\begin{aligned} \Pr_{z \sim \mathbb{U}(\mathbf{y}'^*)} \left[T'_{\downarrow z} \in \mathcal{A} \right] &= \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} \left[T'_{\downarrow z} \in \mathcal{A} \right], \\ \Pr_{z \sim \mathbb{U}(\mathbf{y}'^*)} \left[T'_{\downarrow z} \in \mathcal{N} \right] &= \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} \left[T'_{\downarrow z} \in \mathcal{N} \right]. \end{aligned}$$

With these claims we can finally prove the reduction is correct. That is, we show that the result of the procedure `Compute-Minimal-SR`(T^*, \mathbf{x}, δ) is different from \mathbf{y}^* if and only if (T, \mathbf{y}) is a positive instance of `SB-Check-SUB-SR`.

(\Rightarrow) Assume the result of `Compute-Minimal-SR`(T^*, \mathbf{x}, δ) is some partial instance \mathbf{y}'^* different from \mathbf{y}^* . Note immediately that it is not possible that $\mathbf{y}^* \subseteq \mathbf{y}'^*$ as by definition of δ , we have that

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T^*(z) = 1] \geq \delta,$$

which would contradict the minimality of \mathbf{y}'^* . Let us first prove that $\mathbf{y}'^* \subseteq \mathbf{y}^*$. As a first step, we show that $\mathbf{y}'^*[i] = \mathbf{y}'^*[i'] = \perp$ for every $i \notin S$. We do this by way of contradiction, assuming first that $\mathbf{y}'^*[i] = 0$ or $\mathbf{y}'^*[i'] = 1$ for some $i \notin S$, and exposing how either case generates a contradiction.⁵

- (1) If there is some feature $i \notin S$ such that $\mathbf{y}'^*[i] = 0$, let us define \mathbf{y}^\dagger to be equal to \mathbf{y}'^* except that $\mathbf{y}^\dagger[i] = \perp$. This means that $\mathbf{y}^\dagger \subseteq \mathbf{y}'^*$. Moreover, let \mathbf{y}_1^\dagger be equal to \mathbf{y}'^* except that $\mathbf{y}_1^\dagger[i] = 1$. We now show that \mathbf{y}^\dagger is also a valid output of the computation problem, which contradicts the minimality of \mathbf{y}'^* . Indeed, given that $(T_{\mathbf{y}_0} \cap T_1)(z) = 0$ for every instance z , it holds that

$$\begin{aligned} \Pr_{z \sim \mathbb{U}(\mathbf{y}^\dagger)} [T^*(z) = 1] &= \Pr_{z \sim \mathbb{U}(\mathbf{y}^\dagger)} [T_{\mathbf{y}_0}(z) = 1] + \Pr_{z \sim \mathbb{U}(\mathbf{y}^\dagger)} [(T' \cap T_1)(z) = 1] \\ &= \Pr_{z \sim \mathbb{U}(\mathbf{y}'^*)} [T_{\mathbf{y}_0}(z) = 1] + \Pr_{z \sim \mathbb{U}(\mathbf{y}^\dagger)} [(T' \cap T_1)(z) = 1], \end{aligned}$$

and we can then observe that the events $T'(z) = 1$ and $T_1(z) = 1$ are independent, thus implying that

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}^\dagger)} [(T' \cap T_1)(z) = 1] = \Pr_{z \sim \mathbb{U}(\mathbf{y}^\dagger)} [T'(z) = 1] \cdot \Pr_{z \sim \mathbb{U}(\mathbf{y}^\dagger)} [T_1(z) = 1].$$

By construction of \mathbf{y}^\dagger , we also have that

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}^\dagger)} [T_1(z) = 1] = \Pr_{z \sim \mathbb{U}(\mathbf{y}'^*)} [T_1(z) = 1].$$

⁵Recall that $\mathbf{x}[i] = 0$ and $\mathbf{x}[i'] = 1$, so if $\mathbf{y}'^*[i] \neq \perp$ or $\mathbf{y}'^*[i'] \neq \perp$, then $\mathbf{y}'^*[i] = 0$ or $\mathbf{y}'^*[i'] = 1$ as $\mathbf{y}'^* \subseteq \mathbf{x}$.

Thus, it now suffices to show that

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}_1^\dagger)} [T'(z) = 1] \geq \Pr_{z \sim \mathbb{U}(\mathbf{y}'^*)} [T'(z) = 1],$$

as this implies by definition of \mathbf{y}'^* , \mathbf{y}^\dagger and \mathbf{y}_1^\dagger that

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}^\dagger)} [T'(z) = 1] \geq \Pr_{z \sim \mathbb{U}(\mathbf{y}'^*)} [T'(z) = 1],$$

which in turn implies by the previous discussion that

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}^\dagger)} [T^*(z) = 1] \geq \Pr_{z \sim \mathbb{U}(\mathbf{y}'^*)} [T^*(z) = 1] \geq \delta,$$

and leads to a contradiction.

In order to prove that $\Pr_{z \sim \mathbb{U}(\mathbf{y}_1^\dagger)} [T'(z) = 1] \geq \Pr_{z \sim \mathbb{U}(\mathbf{y}'^*)} [T'(z) = 1]$, we consider two cases, either $\mathbf{y}'^*[i'] = 1$ or $\mathbf{y}'^*[i'] = \perp$.

(a) If $\mathbf{y}'^*[i'] = 1$, then we have that

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}_1^\dagger)} [T'(z) = 1] \geq \Pr_{z \sim \mathbb{U}(\mathbf{y}'^*)} [T'(z) = 1],$$

as for every node u in T' labeled with i , any completion z of \mathbf{y}'^* that goes through u is rejected by construction (landing in an artificial **false** leaf), and for paths of T' that do not go through a node labeled u there is no difference between completions of \mathbf{y}'^* and those for \mathbf{y}_1^\dagger .

(b) If $\mathbf{y}'^*[i'] = \perp$, then we have that for every node u in T' which corresponds to (u, L, R) according to the recursive definition of a decision tree, and is labeled with i , the probability of acceptance of a random completion z conditioned on its path going through u , which is denoted by $z \rightsquigarrow u$, is as follows:

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}'^*)} [T'(z) = 1 \mid z \rightsquigarrow u] = \frac{1}{2} \cdot \Pr_{z \sim \mathbb{U}(\mathbf{y}'^*)} [L(z) = 1 \mid z \rightsquigarrow u],$$

while on the other hand we have

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}_1^\dagger)} [T'(z) = 1 \mid z \rightsquigarrow u] = \frac{1}{2} + \frac{1}{2} \cdot \Pr_{z \sim \mathbb{U}(\mathbf{y}_1^\dagger)} [R(z) = 1 \mid z \rightsquigarrow u].$$

By considering that

$$\frac{1}{2} \geq \frac{1}{2} \cdot \Pr_{z \sim \mathbb{U}(\mathbf{y}'^*)} [L(z) = 1 \mid z \rightsquigarrow u],$$

we obtain

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}_1^\dagger)} [T'(z) = 1 \mid z \rightsquigarrow u] \geq \frac{1}{2} \cdot \Pr_{z \sim \mathbb{U}(\mathbf{y}'^*)} [L(z) = 1 \mid z \rightsquigarrow u],$$

from where

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}_1^\dagger)} [T'(z) = 1 \mid z \rightsquigarrow u] \geq \Pr_{z \sim \mathbb{U}(\mathbf{y}'^*)} [T'(z) = 1 \mid z \rightsquigarrow u].$$

Therefore, by noting that $z \rightsquigarrow u_1$ and $z \rightsquigarrow u_2$ are disjoint events if u_1, u_2 are different nodes of T' with label i , and defining $N(T', i)$ as the set of nodes of T' with label i , we have that

$$\begin{aligned} \Pr_{z \sim \mathbb{U}(\mathbf{y}_1^\dagger)} [T'(z) = 1] &= \sum_{u \in N(T', i)} \Pr_{z \sim \mathbb{U}(\mathbf{y}_1^\dagger)} [T'(z) = 1 \mid z \rightsquigarrow u] \cdot \Pr_{z \sim \mathbb{U}(\mathbf{y}_1^\dagger)} [z \rightsquigarrow u] \\ &\geq \sum_{u \in N(T', i)} \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'(z) = 1 \mid z \rightsquigarrow u] \cdot \Pr_{z \sim \mathbb{U}(\mathbf{y}_1^\dagger)} [z \rightsquigarrow u] \\ &= \sum_{u \in N(T', i)} \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'(z) = 1 \mid z \rightsquigarrow u] \cdot \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [z \rightsquigarrow u] \\ &= \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'(z) = 1] \end{aligned}$$

This concludes the proof of this case.

- (2) It remains to analyze the case when $\mathbf{y}^*[i] = \perp$ and $\mathbf{y}^*[i'] = 1$, which can be proved in the same way as the previous case $\mathbf{y}^*[i] = 0$ and $\mathbf{y}^*[i'] = \perp$.

After this case analysis, we can safely assume that $\mathbf{y}^*[i] = \mathbf{y}^*[i'] = \perp$ for every $i \notin S$. We now show that $\mathbf{y}^*[b_j] = \perp$ for every $j \in \{1, \dots, m\}$. To see this, consider that in general it could be that \mathbf{y}^* forces a certain number k of features b_j to get value 0, meaning that there is a set $K \subseteq \{1, \dots, m\}$ with $|K| = k$ such that $\mathbf{y}^*[b_j] = 0$ for $j \in K$. We argue that $k = 0$. Let us start by arguing that $k \leq m - 2$. Indeed, assume expecting a contradiction that $k > m - 2$, then by definition

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T_1(z) = 1] = 0,$$

and thus

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T^*(z) = 1] = \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T_{y_0}(z) = 1],$$

but

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T_{y_0}(z) = 1] \leq \frac{1}{2},$$

as \mathbf{y}^* cannot be a superset of \mathbf{y}^* , and thus at least one feature i of \mathbf{y}^* is undefined in \mathbf{y}^* , and the event $z[i] = \mathbf{y}^*[i]$, which happens with probability $\frac{1}{2}$, is required for $T_{y_0}(z) = 1$. But by definition of δ , if \mathbf{y}^* is the output of the computational problem, then its probability of acceptance is at least

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [(T' \cap T_1)(z) = 1],$$

and this probability is greater than $\frac{1}{2}$ (Claim 4.12), and thus we have a contradiction. We now safely assume $k \leq m - 2$ and thus $m - k \geq 2$. Observe that as at least one component of \mathbf{y}^* is undefined in \mathbf{y}^* , we have

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T_{y_0}(z) = 1] \leq \frac{1}{2} \cdot \left(\frac{1}{2}\right)^{m-k},$$

and thus

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [(T' \cap T_1)(z) = 1] \geq \delta - \frac{1}{2} \cdot \left(\frac{1}{2}\right)^{m-k},$$

which, considering that

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T_1(z) = 1] = 1 - \left(\frac{1}{2}\right)^{m-k} - (m-k) \left(\frac{1}{2}\right)^{m-k},$$

implies that

$$\Pr_{\mathbf{y} \sim \bar{\mathbb{U}}(\mathbf{z}^*)} [T'(\mathbf{z}) = 1] \geq \frac{\delta - \frac{1}{2} \cdot \left(\frac{1}{2}\right)^{m-k}}{1 - \left(\frac{1}{2}\right)^{m-k} - (m-k) \left(\frac{1}{2}\right)^{m-k}},$$

as $T'(\mathbf{z}) = 1$ and $T_1(\mathbf{z}) = 1$ are independent events. We now show that the right side of the previous equation is greater than δ . Indeed, for ease of notation set $r := (m - k + 1)$ and note that the right side can be rewritten as

$$\frac{\delta - \left(\frac{1}{2}\right)^r}{1 - r \left(\frac{1}{2}\right)^{r-1}}.$$

Now consider that as $m - k \geq 2$ we have that $r \geq 3$ and thus $2^{r-1} > r$, which implies $r \left(\frac{1}{2}\right)^{r-1} < 1$, and thus the denominator of the previous equation is positive, implying that what we want to show is equivalent to

$$\delta - \left(\frac{1}{2}\right)^r > \delta \left(1 - r \left(\frac{1}{2}\right)^{r-1}\right),$$

which is in turn equivalent to

$$\delta r \left(\frac{1}{2}\right)^{r-1} > \left(\frac{1}{2}\right)^r,$$

but as by Claim 4.12 we have $\delta > \frac{1}{2}$, and $r \geq 3$, the previous equation is trivially true. We have therefore showed that

$$\Pr_{\mathbf{z} \sim \bar{\mathbb{U}}(\mathbf{y}^*)} [T'(\mathbf{z}) = 1] > \delta.$$

Now let \mathbf{y}^\ominus be the partial instance such that $\mathbf{y}^\ominus[i] = \mathbf{y}^*[i]$ for every $i \in S$, and is undefined in all other features. Note that $\mathbf{y}^\ominus \subseteq \mathbf{y}^*$ and also $\mathbf{y}^\ominus \subseteq \mathbf{y}^*$. If $\mathbf{y}^\ominus = \mathbf{y}^*$, then $\mathbf{y}^* \subseteq \mathbf{y}^*$ which is what we are hoping to prove. So we now assume $\mathbf{y}^\ominus \subsetneq \mathbf{y}^*$ expecting a contradiction. Note that T' does not use the b_j features at all, and therefore we have that

$$\Pr_{\mathbf{z} \sim \bar{\mathbb{U}}(\mathbf{y}^\ominus)} [T'(\mathbf{z}) = 1] = \Pr_{\mathbf{z} \sim \bar{\mathbb{U}}(\mathbf{y}^*)} [T'(\mathbf{z}) = 1] > \delta.$$

We use this to prove that \mathbf{y}^\ominus would have been a valid outcome of the computing problem, thus contradicting the minimality of \mathbf{y}^* . Indeed,

$$\Pr_{\mathbf{z} \sim \bar{\mathbb{U}}(\mathbf{y}^\ominus)} [T^*(\mathbf{z}) = 1] > \Pr_{\mathbf{z} \sim \bar{\mathbb{U}}(\mathbf{y}^\ominus)} [T'(\mathbf{z}) = 1] \cdot \Pr_{\mathbf{z} \sim \bar{\mathbb{U}}(\mathbf{y}^\ominus)} [T_1(\mathbf{z}) = 1],$$

and note that as

$$\Pr_{\mathbf{z} \sim \bar{\mathbb{U}}(\mathbf{y}^\ominus)} [T'(\mathbf{z}) = 1] > \delta,$$

it must be the case that

$$\Pr_{\mathbf{z} \sim \bar{\mathbb{U}}(\mathbf{y}^\ominus)} [T'(\mathbf{z}) = 1] \geq \delta + \left(\frac{1}{2}\right)^{2n-|S|},$$

as only $2n - |S|$ features appear as labels in T' and, thus, the completion probability of any partial instance over T' must be an integer multiple of $\left(\frac{1}{2}\right)^{2n-|S|}$. Now let us abbreviate $2n - |S|$ as ℓ and choose $m \geq 2\ell$. We thus have

that

$$\begin{aligned}
\Pr_{z \sim \mathbb{U}(\mathbf{y}^\ominus)} [T^\star(z) = 1] &\geq \Pr_{z \sim \mathbb{U}(\mathbf{y}^\ominus)} [T'(z) = 1] \cdot \Pr_{z \sim \mathbb{U}(\mathbf{y}^\ominus)} [T_1(z) = 1] \\
&\geq \left(\delta + \left(\frac{1}{2} \right)^\ell \right) \left(1 - (m+1) \left(\frac{1}{2} \right)^m \right) \\
&\geq \left(\delta + \left(\frac{1}{2} \right)^\ell \right) \left(1 - (2\ell+1) \left(\frac{1}{2} \right)^{2\ell} \right) \\
&= \delta - \delta(2\ell+1) \left(\frac{1}{2} \right)^{2\ell} + \left(\frac{1}{2} \right)^\ell - (2\ell+1) \left(\frac{1}{2} \right)^{3\ell} \\
&\geq \delta - (2\ell+1) \left(\frac{1}{2} \right)^{2\ell} + \left(\frac{1}{2} \right)^\ell - (2\ell+1) \left(\frac{1}{2} \right)^{2\ell} \\
&= \delta - (4\ell+2) \left(\frac{1}{2} \right)^{2\ell} + \left(\frac{1}{2} \right)^\ell \\
&= \delta + \left(\frac{1}{2} \right)^\ell \left(1 - (4\ell+2) \left(\frac{1}{2} \right)^\ell \right),
\end{aligned}$$

where the last parenthesis is positive for $\ell \geq 5$, which can be assumed without loss of generality as $2n - |S| < 5$ requires $n < 5$ (since $|S| \leq n$), and thus the original instance of the decision problem would have constant size. We have thus concluded that

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}^\ominus)} [T^\star(z) = 1] \geq \Pr_{z \sim \mathbb{U}(\mathbf{y}^\ominus)} [(T' \cap T_1)(z) = 1] \geq \delta, \quad (4.14)$$

thus showing that \mathbf{y}^\ominus is a valid outcome for the computing problem, which contradicts the minimality of \mathbf{y}^\star . This in turn implies that $\mathbf{y}^{\star\star} = \mathbf{y}^\ominus$, and thus subsequently that $\mathbf{y}^{\star\star} \subseteq \mathbf{y}^\star$. Let us now show how by combining Claims 4.10, 4.11 and 4.13, we can conclude the forward direction entirely. Indeed, note that the trivial equality

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}^\star)} [T_1(z) = 1] = \Pr_{z \sim \mathbb{U}(\mathbf{y}^{\star\star})} [T_1(z) = 1]$$

implies that

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}^\star)} [T'(z) = 1] \leq \Pr_{z \sim \mathbb{U}(\mathbf{y}^{\star\star})} [T'(z) = 1], \quad (4.15)$$

as we already have proved that $\Pr_{z \sim \mathbb{U}(\mathbf{y}^{\star\star})} [(T' \cap T_1)(z) = 1] \geq \delta$ by Equation 4.14 and the fact that $\mathbf{y}^\ominus = \mathbf{y}^{\star\star}$, and we have that $\delta = \Pr_{z \sim \mathbb{U}(\mathbf{y}^\star)} [(T' \cap T_1)(z) = 1]$. We can now use Claims 4.10, 4.11 and 4.13 together

with Equation 4.15 to conclude that

$$\begin{aligned}
 \Pr_{z \sim \bar{U}(\mathbf{y})} [T(z) = 1] &= \Pr_{z \sim \bar{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{N}_t \mid T'_{\downarrow z} \in \mathcal{N}] && \text{(Claim 4.11)} \\
 &= \frac{\Pr_{z \sim \bar{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{N}_t]}{\Pr_{z \sim \bar{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{N}]} \\
 &= \frac{\Pr_{z \sim \bar{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{N}_t]}{\Pr_{z \sim \bar{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{N}]} && \text{(Claim 4.13)} \\
 &= \frac{\Pr_{z \sim \bar{U}(\mathbf{y}^*)} [T'(z) = 1] - \Pr_{z \sim \bar{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{A}_t]}{\Pr_{z \sim \bar{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{N}]} \\
 &= \frac{\Pr_{z \sim \bar{U}(\mathbf{y}^*)} [T'(z) = 1] - \Pr_{z \sim \bar{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{A}_t \mid T'_{\downarrow z} \in \mathcal{A}] \cdot \Pr_{z \sim \bar{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{A}]}{\Pr_{z \sim \bar{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{N}]} \\
 &= \frac{\Pr_{z \sim \bar{U}(\mathbf{y}^*)} [T'(z) = 1] - \Pr_{z \sim \bar{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{A}_t \mid T'_{\downarrow z} \in \mathcal{A}] \cdot \Pr_{z \sim \bar{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{A}]}{\Pr_{z \sim \bar{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{N}]} && \text{(Claims 4.10 and 4.13)} \\
 &\leq \frac{\Pr_{z \sim \bar{U}(\mathbf{y}^*)} [T'(z) = 1] - \Pr_{z \sim \bar{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{A}_t \mid T'_{\downarrow z} \in \mathcal{A}] \cdot \Pr_{z \sim \bar{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{A}]}{\Pr_{z \sim \bar{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{N}]} && \text{(Equation 4.15)} \\
 &= \frac{\Pr_{z \sim \bar{U}(\mathbf{y}^*)} [T'(z) = 1] - \Pr_{z \sim \bar{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{A}_t]}{\Pr_{z \sim \bar{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{N}]} \\
 &= \frac{\Pr_{z \sim \bar{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{N}_t]}{\Pr_{z \sim \bar{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{N}]} \\
 &= \Pr_{z \sim \bar{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{N}_t \mid T'_{\downarrow z} \in \mathcal{N}] \\
 &= \Pr_{z \sim \bar{U}(\mathbf{y}')} [T(z) = 1], && \text{(Claim 4.11)}
 \end{aligned}$$

where \mathbf{y}' is the partial instance of dimension n such that $\mathbf{y}'[i] = \mathbf{y}^*[i]$ for every i such that $\mathbf{y}^*[i] \neq \perp$, and \mathbf{y}' is undefined in all other features. By this definition, $\mathbf{y}' \subseteq \mathbf{y}$ as we had $\mathbf{y}^* \subseteq \mathbf{y}^*$ (because by assumption $\mathbf{y}^* \neq \mathbf{y}^*$), and thus we have effectively proved that the instance (T, z) is a positive instance of SB-Check-SUB-SR. This concludes the proof of the forward direction.

(\Leftarrow) Assume the instance (T, \mathbf{y}) is a positive instance of SB-Check-SUB-SR and, thus, there exists some $\mathbf{y}' \subseteq \mathbf{y}$ such that

$$\Pr_{z \sim \bar{U}(\mathbf{y}')} [T(z) = 1] \geq \Pr_{z \sim \bar{U}(\mathbf{y})} [T(z) = 1].$$

Define \mathbf{y}^* of the dimension of T^* based on \mathbf{y}' by setting $\mathbf{y}^*[i] = \mathbf{y}'[i]$ for every i such that $\mathbf{y}'[i] \neq \perp$, and leave the rest of \mathbf{y}^* undefined. Note that this definition immediately implies $\mathbf{y}^* \subseteq \mathbf{y}^*$. By Claim 4.11 the previous

equation implies that

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}'^*)} \left[T'_{\downarrow z} \in \mathcal{N}_t \mid T'_{\downarrow z} \in \mathcal{N} \right] \geq \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} \left[T'_{\downarrow z} \in \mathcal{N}_t \mid T'_{\downarrow z} \in \mathcal{N} \right],$$

which implies in turn that

$$\frac{\Pr_{z \sim \mathbb{U}(\mathbf{y}'^*)} [T'_{\downarrow z} \in \mathcal{N}_t]}{\Pr_{z \sim \mathbb{U}(\mathbf{y}'^*)} [T'_{\downarrow z} \in \mathcal{N}]} \geq \frac{\Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{N}_t]}{\Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{N}]}.$$

By Claim 4.13 the denominators of the previous inequality are equal and, thus,

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}'^*)} [T'_{\downarrow z} \in \mathcal{N}_t] \geq \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{N}_t],$$

from where

$$\begin{aligned} \Pr_{z \sim \mathbb{U}(\mathbf{y}'^*)} [T'_{\downarrow z} \in \mathcal{N}_t] + \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{A}_t] &\geq \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{N}_t] + \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{A}_t] \\ &= \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'(z) = 1]. \end{aligned}$$

But combining Claims 4.10 and 4.13 we have that

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{A}_t] = \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{A}_t],$$

which when combined with the previous equation gives us

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}'^*)} [T'(z) = 1] \geq \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'(z) = 1],$$

and using again that

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}'^*)} [T_1(z) = 1] = \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T_1(z) = 1],$$

we obtain that

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}'^*)} [(T' \cap T_1)(z) = 1] \geq \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [(T' \cap T_1)(z) = 1] = \delta.$$

Finally, by observing that

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}'^*)} [T^*(z) = 1] \geq \Pr_{z \sim \mathbb{U}(\mathbf{y}'^*)} [(T' \cap T_1)(z) = 1] \geq \delta,$$

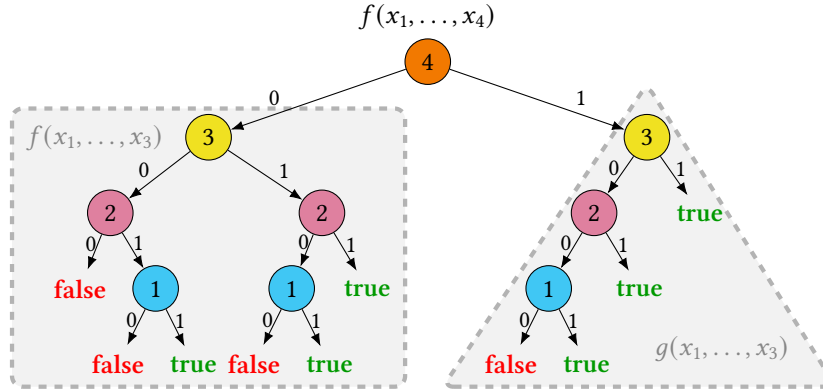
we have that \mathbf{y}'^* is a valid output for the computational problem, and given that it is a strict subset of \mathbf{y}^* , the result of $\text{Compute-Minimal-SR}(T^*, \mathbf{x}, \delta)$ cannot be equal to \mathbf{y}^* . This concludes the backward direction, and with it the entire proof is complete. \square

We now restate and prove each of the claims used in the preceding theorem.

Claim 4.9 (Restated). *Decision tree T_1 , implementing the function ϕ , can be constructed in polynomial time.*

PROOF OF CLAIM 4.9. This proof can be easily done by a direct construction. Indeed, consider the following Boolean formulas:

$$\begin{aligned} f(x_1, \dots, x_n) &\equiv \sum_{i=1}^n x_i \geq 2, \\ g(x_1, \dots, x_n) &\equiv \sum_{i=1}^n x_i \geq 1 = \bigvee_{i=1}^n x_i. \end{aligned}$$


 Fig. 9. Illustration of T_1 following the construction for Claim 4.9.

We then note that

$$f(x_1, \dots, x_n) = [\neg x_n \wedge f(x_1, \dots, x_{n-1})] \vee [x_n \wedge g(x_1, \dots, x_{n-1})],$$

and thus we can build a decision tree for f recursively as illustrated in Figure 9. Note that $g(x_1, \dots, x_n)$ can be trivially implemented by a decision tree of size $O(n)$. Thus, the recursive equation characterizing the size $\alpha(n)$ of a decision tree for $f(x_1, \dots, x_n)$ according to our construction is simply $\alpha(n) = 1 + \alpha(n-1) + O(n)$, from where we get $\alpha(n) \in O(n^2)$, thus concluding the proof of the claim. \square

Claim 4.10 (Restated). *For every partial instance $\mathbf{y}^* \subseteq \mathbf{y}^*$, it holds that*

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{A}_t \mid T'_{\downarrow z} \in \mathcal{A}] = \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{A}_f \mid T'_{\downarrow z} \in \mathcal{A}] = \frac{1}{2}.$$

PROOF OF CLAIM 4.10. Observe that every leaf $\ell \in \mathcal{A}$ has a parent node v in T' labeled with some feature i' whose parent node u in T' is labeled with feature i . Let $G(\ell) = u$ be said the grandparent of ℓ , and assume that $G^{-1}(u) = \{\ell' \mid G(\ell') = u\}$. With this notation, we can split the set \mathcal{A} as follows:

$$\mathcal{A} = \bigcup_{\text{node } u \text{ with label } i \notin S} \{T'_{\downarrow z} \mid G(T'_{\downarrow z}) = u\},$$

where the union is actually disjoint. Thus, we have for every partial instance $\mathbf{y}^* \subseteq \mathbf{y}^*$:

$$\begin{aligned} & \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{A}_t \mid T'_{\downarrow z} \in \mathcal{A}] \\ &= \sum_{\substack{\text{node } u \\ \text{with label } i \notin S}} \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{A}_t \mid T'_{\downarrow z} \in \mathcal{A} \cap G^{-1}(u)] \cdot \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{A} \cap G^{-1}(u) \mid T'_{\downarrow z} \in \mathcal{A}]. \end{aligned} \quad (4.16)$$

Observe that

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{A}_t \mid T'_{\downarrow w} \in \mathcal{A} \cap G^{-1}(u)] = \frac{1}{2},$$

since the event is equivalent to $z[i] = 1, z[i'] = 0$, and this is equally likely to the complement event $z[i] = 0, z[i'] = 1$, given that $\mathbf{y}^*[i] = \mathbf{y}^*[i'] = \perp$. Therefore,

$$\begin{aligned} \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} \left[T'_{\downarrow z} \in \mathcal{A}_t \mid T'_{\downarrow z} \in \mathcal{A} \right] &= \sum_{\text{node } u \text{ with label } i \notin S} \frac{1}{2} \cdot \Pr_{z \in \mathbb{U}(\mathbf{y}^*)} \left[T'_{\downarrow z} \in \mathcal{A} \cap G^{-1}(u) \mid T'_{\downarrow z} \in \mathcal{A} \right] \\ &= \frac{1}{2} \cdot \sum_{\text{node } u \text{ with label } i \notin S} \Pr_{z \in \mathbb{U}(\mathbf{y}^*)} \left[T'_{\downarrow z} \in \mathcal{A} \cap G^{-1}(u) \mid T'_{\downarrow z} \in \mathcal{A} \right] \\ &= \frac{1}{2}. \end{aligned} \quad \square$$

Claim 4.11 (Restated). *Given a partial instance $\mathbf{y}' \subseteq \mathbf{y}$, we define \mathbf{y}^* as the partial instance of dimension $2n - |S| + m$ that matches \mathbf{y}' on its defined features, and holds $\mathbf{y}^*[i] = \mathbf{y}^*[i'] = \perp$ for every feature i such that $\mathbf{y}'[i] = \perp$. Then it holds that*

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}')} [T(z) = 1] = \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} \left[T'_{\downarrow z} \in \mathcal{N}_t \mid T'_{\downarrow z} \in \mathcal{N} \right].$$

PROOF OF CLAIM 4.11. Given that the resulting leaf is natural, for every node u of T' such that u is labeled with feature $i \notin S$, the tuple (u, L, R) was considered when constructing T' and the path of w in T' goes through u , it must be the case that $z[i] = z[i'] = 0$ or $z[i] = z[i'] = 1$, as otherwise $T'_{\downarrow z} \in \mathcal{A}$. But these two alternatives are equally likely by definition of T' . Thus, by using a simple induction argument, for every leaf ℓ of T with a corresponding natural leaf ℓ' of T' , it holds that

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}')} [T_{\downarrow z} = \ell] = \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} \left[T'_{\downarrow z} = \ell' \mid T'_{\downarrow z} \in \mathcal{N} \right],$$

from where the claim immediately follows. \square

Claim 4.12 (Restated). *By choosing $m \geq \max\{2u(T, \mathbf{y}) + 2n, 9\}$, we have that*

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [(T' \cap T_1)(z) = 1] > \frac{1}{2}.$$

PROOF OF CLAIM 4.12. First, consider that for T_1 we have

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T_1(z) = 1] = 1 - \left(\frac{1}{2}\right)^m - m \left(\frac{1}{2}\right)^m = 1 - (m+1) \left(\frac{1}{2}\right)^m,$$

while on the other hand

$$\begin{aligned} \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'(z) = 1] &= \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{A}_t] + \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{N}_t] \\ &= \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{A}_t \mid T'_{\downarrow z} \in \mathcal{A}] \cdot \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{A}] \\ &\quad + \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{N}_t \mid T'_{\downarrow z} \in \mathcal{N}] \cdot \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{N}] \\ &= \frac{1}{2} \cdot \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{A}] + \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{N}_t \mid T'_{\downarrow z} \in \mathcal{N}] \cdot \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{N}], \end{aligned}$$

where the last equality uses Claim 4.10. Let us now show that

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{N}_t \mid T'_{\downarrow z} \in \mathcal{N}] \geq \frac{1}{2} + \left(\frac{1}{2}\right)^n.$$

By Claim 4.11, this is the same as showing that

$$\Pr_{z \sim \mathbb{U}(\mathbf{y})} [T(z) = 1] \geq \frac{1}{2} + \left(\frac{1}{2}\right)^n,$$

which is guaranteed by the definition of the SB-Check-SUB-SR problem, as we know that

$$\Pr_{z \sim \mathbb{U}(\mathbf{y})} [T(z) = 1] > \frac{1}{2},$$

and also that $\Pr_{z \sim \mathbb{U}(\mathbf{y})} [T(z) = 1]$ must be of the form $\left(\frac{k}{2^n}\right)$ with $k \in \mathbb{N}$, given that n is the dimension of T .

Now, consider that

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{A}] = 1 - \left(\frac{1}{2}\right)^{u(T, \mathbf{y})}.$$

Notice that this holds because T is strongly-balanced, so falling into a natural leaf in T' requires going through $u(T, \mathbf{y})$ layers without choosing an artificial leaf, which happens with probability $\frac{1}{2}$ at each layer. Thus, we have that

$$\begin{aligned} \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'(z) = 1] &\geq \frac{1}{2} \left(1 - \left(\frac{1}{2}\right)^{u(T, \mathbf{y})}\right) + \left(\frac{1}{2} + \left(\frac{1}{2}\right)^n\right) \left(\frac{1}{2}\right)^{u(T, \mathbf{y})} \\ &= \frac{1}{2} + \left(\frac{1}{2}\right)^{n+u(T, \mathbf{y})}. \end{aligned}$$

Moreover, given that T' and T_1 impose restrictions over disjoint sets of features, we have that

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'(z) = 1 \mid T_1(z) = 1] = \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'(z) = 1].$$

Putting together all the previous results, we obtain that

$$\begin{aligned} \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [(T' \cap T_1)(z) = 1] &= \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'(z) = 1 \mid T_1(z) = 1] \cdot \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T_1(z) = 1] \\ &= \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'(z) = 1] \cdot \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T_1(z) = 1] \\ &\geq \left(\frac{1}{2} + \left(\frac{1}{2}\right)^{n+u(T, \mathbf{y})}\right) \left(1 - (m+1) \left(\frac{1}{2}\right)^m\right) \\ &= \frac{1}{2} - (m+1) \left(\frac{1}{2}\right)^{m+1} + \left(\frac{1}{2}\right)^{n+u(T, \mathbf{y})} - (m+1) \left(\frac{1}{2}\right)^{m+n+u(T, \mathbf{y})} \\ &\geq \frac{1}{2} - (m+1) \left(\frac{1}{2}\right)^m + \left(\frac{1}{2}\right)^{n+u(T, \mathbf{y})} \\ &\geq \frac{1}{2} - \left(\frac{1}{2}\right)^{m - \lceil \log(m+1) \rceil} + \left(\frac{1}{2}\right)^{n+u(T, \mathbf{y})}. \end{aligned}$$

But we are assuming $m \geq \max\{2n + 2u(T, \mathbf{y}), 9\}$, which implies that $m \geq 2n + 2u(T, \mathbf{y})$. Hence, we have that

$$m - \lceil \log(m+1) \rceil > n + u(T, \mathbf{y}),$$

as $m - \lceil \log(m+1) \rceil > \frac{m}{2}$ since $m \geq 9$. We conclude that

$$\frac{1}{2} - \left(\frac{1}{2}\right)^{m - \lceil \log(m+1) \rceil} + \left(\frac{1}{2}\right)^{n+u(T, \mathbf{y})} > \frac{1}{2},$$

from which the claim follows. \square

Claim 4.13 (Restated). *For every partial instance $\mathbf{y}'^* \subseteq \mathbf{y}^*$, it holds that*

$$\begin{aligned} \Pr_{z \sim \mathbb{U}(\mathbf{y}'^*)} [T'_{\downarrow z} \in \mathcal{A}] &= \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{A}], \\ \Pr_{z \sim \mathbb{U}(\mathbf{y}'^*)} [T'_{\downarrow z} \in \mathcal{N}] &= \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{N}]. \end{aligned}$$

PROOF OF CLAIM 4.13. We only need to prove that

$$\Pr_{z \sim \mathbb{U}(\mathbf{y}'^*)} [T'_{\downarrow z} \in \mathcal{A}] = \Pr_{z \sim \mathbb{U}(\mathbf{y}^*)} [T'_{\downarrow z} \in \mathcal{A}].$$

As shown in the proof of Claim 4.12, this follows from the strongly-balanced property of T . \square

5 Proofs of Non-approximability Results

Both Theorem 3.4 and Theorem 3.5 (whose statements we will recall shortly) are consequences of the hardness of a related problem: we are given a decision tree T , and we are allowed to fix an arbitrary subset of the variables to 1. Our goal is to maximize the probability that T outputs 1 under this fixation. Note that fixing to 0 is not allowed, otherwise, we could just find an arbitrary input on which T is 1 (unless T is always 0) and fix all the variables according to this input. If T is equal to 1 on the input of all 1s, the optimal solution is fixing all variables to 1, but of course this will not happen in hard instances. For this problem, we obtain the following hardness result for its promise version.

THEOREM 5.1. *Unless $\text{QP} = \text{NP}$, for every $\kappa > 0$ there exists no polynomial-time algorithm that, given a decision tree T over n variables, distinguishes between the following two cases:*

- there exists a partial instance $\mathbf{y} \in \{1, \perp\}^n$ such that $\Pr_z[T(z) = 1 \mid z \in \text{COMP}(\mathbf{y})] \geq 1 - \kappa$;
- there exists no partial instance $\mathbf{y} \in \{1, \perp\}^n$ such that $\Pr_z[T(z) = 1 \mid z \in \text{COMP}(\mathbf{y})] \geq \kappa$.

We first derive Theorem 3.4 and Theorem 3.5 from this result, and then we prove it separately.

THEOREM 3.4 (RESTATED). *Assuming that $\text{NP} \not\subseteq \text{QP}$, for every $\varepsilon > 0$ there is no polynomial-time algorithm that, given a decision tree T of dimension n and an input $\mathbf{x} \in \{0, 1\}^n$, distinguishes between the following two cases:*

- there exists an $(1 - \varepsilon)$ -sufficient reason \mathbf{y} for \mathbf{x} under T with $|\mathbf{y}|_{\perp} \geq n - n^{\varepsilon}$.
- there exists no ε -sufficient reason \mathbf{y} for \mathbf{x} under T with $|\mathbf{y}|_{\perp} \geq n^{\varepsilon}$.

PROOF OF THEOREM 3.4. For every $\varepsilon > 0$, we provide a polynomial-time reduction from the problem from Theorem 5.1 with $\kappa = \varepsilon/2$ to the problem from Theorem 3.4 with parameter ε . Let Δ be any positive integer larger than $1/\varepsilon$. Define $m = (n + 2\Delta)^{\Delta}$. Given a decision tree T over n variables, consider a decision tree T_1 over $n + m$ variables, defined as follows:

$$T_1(x_1, \dots, x_{n+m}) = T(x_1, \dots, x_n) \vee (x_{n+1} \wedge x_{n+2} \wedge \dots \wedge x_{n+m}).$$

Note that T_1 can be constructed in polynomial time from T (to every leaf of T that outputs 0 attach a tree computing the conjunction $x_{n+1} \wedge x_{n+2} \wedge \dots \wedge x_{n+m}$). Let e^k denote a Boolean vector of dimension k where all coordinates are 1s. It remains to establish two claims.

- If there exists a partial instance $\mathbf{y} \in \{1, \perp\}^n$ such that $\Pr_z[T(z) = 1 \mid z \in \text{COMP}(\mathbf{y})] \geq 1 - \kappa$, then there exists a $(1 - \varepsilon)$ -sufficient reason \mathbf{y}_1 for e^{n+m} under T_1 such that $|\mathbf{y}_1|_{\perp} \geq (n + m) - (n + m)^{\varepsilon}$.
- If there exists no partial instance $\mathbf{y} \in \{1, \perp\}^n$ such that $\Pr_z[T(z) = 1 \mid z \in \text{COMP}(\mathbf{y})] \geq \kappa$, then there exists no ε -sufficient reason \mathbf{y}_1 for e^{n+m} under T_1 such that $|\mathbf{y}_1|_{\perp} \geq (n + m)^{\varepsilon}$.

Let us start with the first claim. First, note that $T(e^{n+m}) = 1$. Let \mathbf{y}_1 be simply an extension of \mathbf{y} to the variables x_{n+1}, \dots, x_{n+m} by letting them be undefined. Thus,

$$|\mathbf{y}_1|_{\perp} \geq m = (n+m) - n = (n+m) - (n^{1/\varepsilon})^\varepsilon \geq (n+m) - (n+m)^\varepsilon$$

where the last inequality is because $m = (n+2\Delta)^\Delta \geq n^{1/\varepsilon}$. It remains to show that \mathbf{y}_1 is a $(1-\varepsilon)$ -sufficient reason for e^{n+m} under T_1 . This is because by sampling a random completion of \mathbf{y}_1 just in variables x_1, \dots, x_n , we already make T to be equal to 1 with probability $1 - \kappa = 1 - \varepsilon/2 > 1 - \varepsilon$. But if T is equal to 1, then T_1 is also equal to 1.

Now, let us show the second claim. Assume that \mathbf{y}_1 is an ε -sufficient reason for e^{n+m} under T_1 . Our goal is to show that $|\mathbf{y}_1|_{\perp} < (n+m)^\varepsilon$. Consider a uniformly random completion \mathbf{z} of \mathbf{y}_1 . On one hand, \mathbf{z} makes T_1 equal to 1 with probability at least ε . On the other hand, variables that are defined in \mathbf{y}_1 are all equal to 1 (as \mathbf{y}_1 is a sufficient reason for the input of all 1s), meaning that \mathbf{z} makes T equal to 1 with probability less than $\kappa = \varepsilon/2$. Hence, a uniformly random completion of \mathbf{y}_1 makes the conjunction $x_{n+1} \wedge x_{n+2} \wedge \dots \wedge x_{n+m}$ true with probability larger than $\varepsilon/2$. On the other hand, the probability of making this conjunction true is 2 to the power of (minus the number of variables among x_{n+1}, \dots, x_{n+m} that are undefined in \mathbf{y}_1). The number of such variables is at least $|\mathbf{y}_1|_{\perp} - n$, which gives us:

$$\varepsilon/2 < 2^{n-|\mathbf{y}_1|_{\perp}},$$

from where we get an upper bound on $|\mathbf{y}_1|_{\perp}$:

$$|\mathbf{y}_1|_{\perp} < n + \log_2(2/\varepsilon) \leq n + \log_2(2\Delta) \leq n + 2\Delta = m^{1/\Delta} \leq (n+m)^\varepsilon.$$

This finishes the proof. \square

THEOREM 3.5 (RESTATED). *Assuming that $\text{NP} \not\subseteq \text{QP}$, for every $\varepsilon \in (0, 1)$, there exists no polynomial-time algorithm that, given a decision tree T of dimension n and an input $\mathbf{x} \in \{0, 1\}^n$, distinguishes between the following two cases:*

- *the only $(1/2 + \varepsilon)$ -sufficient reason for \mathbf{x} under T is \mathbf{x} itself.*
- *there exists a $(1 - \varepsilon)$ -sufficient reason for \mathbf{x} under T other than \mathbf{x} .*

PROOF OF THEOREM 3.5. Given $\varepsilon > 0$, we establish a polynomial-time reduction from the problem given in Theorem 5.1 with $\kappa = \varepsilon$ to the problem from Theorem 3.5 with parameter ε . Suppose we are given a decision tree T of dimension n , and we want to distinguish the case when there exists a fixation of some variables to 1s under which the probability of having 1 in the output of T is at least $1 - \varepsilon$, from the case when all such fixations give probability less than ε . First, we check if $T(1, \dots, 1) = 1$, and if so, we do not need any reduction. From now on, assume that $T(1, \dots, 1) = 0$. We define a tree

$$T_1(x_1, \dots, x_n) = T(x_1, \dots, x_n) \vee (x_1 \wedge x_2 \wedge \dots \wedge x_n)$$

and we consider the input $e^n = (1, \dots, 1)$ to it, on which T_1 is equal to 1 because of the conjunction. If we can fix some variables to 1s such that T is equal to 1 with probability at least $1 - \varepsilon$, then this fixation is a $(1 - \varepsilon)$ -sufficient reason for \mathbf{x} under T_1 , just because T_1 is 1 whenever T is. This sufficient reason cannot be e^n itself because T is 0 on the input of all 1s. Now, assume that under any fixation of some variables to 1s, the probability that T is equal to 1 is lower than ε . Then for any fixation where at least one variable is not fixed, the probability of T_1 to be equal to 1 is lower than $1/2 + \varepsilon$, because the conjunction $(x_1 \wedge x_2 \wedge \dots \wedge x_n)$ is equal to 1 with probability at most $1/2$ then. This concludes the proof as then there is no $(1/2 + \varepsilon)$ -sufficient reason for e^n under T_1 other than e^n . \square

PROOF OF THEOREM 5.1. To simplify the presentation, within this proof we employ the following notation:

$$T(\mathbf{y}) = \Pr_{\mathbf{z}}[T(\mathbf{z}) = 1 \mid \mathbf{z} \in \text{COMP}(\mathbf{y})]$$

for a decision tree T of dimension n and for a partial instance $\mathbf{y} \in \{1, \perp\}^n$.

We use a standard probability concentration inequality due to Hoeffding [26].

Proposition 5.2. *Let us assume that X_1, X_2, \dots, X_n are independent Bernoulli random variables. Denote $S_n = X_1 + \dots + X_n$. Then for every $\delta > 0$ we have:*

$$\Pr \left[\frac{S_n}{n} \geq \frac{\mathbb{E}S_n}{n} + \delta \right] \leq \exp\{-2\delta^2 n\}.$$

We reduce from the *1-in- k exact hitting set problem* (1-in-EkHS) problem. It is a version of the SAT problem where each clause contains exactly k variables (without negations) and a clause is satisfied if and only if there is exactly one variable in it which is set to 1. We use the following result of Guruswami and Trevisan [27, Lemma 13]:

THEOREM 5.3. *For every $\delta > 0$ there exists k such that it is NP-hard to distinguish satisfiable instances of 1-in-EkHS from instances of 1-in-EkHS for which it is impossible to satisfy more than a fraction $(1/e + \delta)$ of clauses.*

We fix any k such that it is NP-hard to distinguish satisfiable instances of 1-in-EkHS from instances of 1-in-EkHS for which it is impossible to satisfy more than a fraction $1/2$ of clauses. Now, fix $\kappa > 0$. To establish Theorem 5.1, we take any instance φ of 1-in-EkHS with n variables and m clauses and in *quasi-polynomial* time construct a decision tree T over N variables such that

- if φ is satisfiable, then there exists $\mathbf{y} \in \{1, \perp\}^N$ such that $T(\mathbf{y}) \geq 1 - \kappa$;
- if it is impossible to satisfy more than a fraction $1/2$ of clauses of φ , then there exists no $\mathbf{y} \in \{1, \perp\}^N$ such that $T(\mathbf{y}) \geq \kappa$.

Let l be the smallest integral number such that $m \leq 2^{2l}$. Note that $l = \Theta(\log m)$. The most technical part of the proof is to construct in polynomial time an $O(\log m)$ -depth decision tree L over $n + 2l + 2$ variables such that for some absolute constant $c > 0$ we have:

- if φ is satisfiable, then there exists $\mathbf{y} \in \{1, \perp\}^{n+2l+2}$ such that $L(\mathbf{y}) \geq 7/8$;
- if it is impossible to satisfy more than a fraction $1/2$ of clauses of φ , then there for every $\mathbf{y} \in \{1, \perp\}^{n+2l+2}$ we have $L(\mathbf{y}) \leq 7/8 - \delta$, where $\delta = \frac{c}{\sqrt{\ln(m)}}$.

We first describe how to “increase the gap” from $(7/8 - \delta$ vs. $7/8)$ to $(\kappa$ vs. $1 - \kappa)$. The construction of L is given afterwards.

Finishing the proof modulo the construction of L . We define a decision tree T as

$$K = \frac{2 \ln(2/\kappa)}{\delta^2} = O(\log m)$$

“independent copies” of L . More specifically, we let T be over $N = K(n + 2l + 2)$ variables. First, T runs L on the first $n + 2l + 2$ variables, then on the second $n + 2l + 2$ variables, and so on (in total, K runs). In the end, it outputs 1 if and only if for at least a fraction $7/8 - \delta/2$ of runs the output of L was 1.

Assume first that φ is satisfiable. Then, by definition of L , there exists $\mathbf{y} \in \{1, \perp\}^{n+2l+2}$ such that $L(\mathbf{y}) \geq 7/8$. Repeat it K times to obtain a partial input $Y \in \{1, \perp\}^N$ for T . We claim that $T(Y) \geq 1 - \kappa$. Indeed, a random completion of Y can be generated as K independent samples of a random completion of \mathbf{y} . The tree T outputs 1 if and only if L outputs 1 for at least a fraction $7/8 - \delta/2$ of the samples. On the other hand, the probability that L outputs 1 on one of the samples is $L(\mathbf{y}) \geq 7/8$. Hence, the average fraction of samples on which L outputs 1 is at least $7/8$. By Hoeffding’s inequality (Proposition 5.2), the probability that T outputs 1 is at least $1 - \exp\{-2(\delta/2)^2 \cdot K\} = 1 - \kappa/2 > 1 - \kappa$.

Now, assume that it is impossible to satisfy more than a fraction $1/2$ of clauses of φ . Then $L(\mathbf{y}) \leq 7/8 - \delta$ for all $\mathbf{y} \in \{1, \perp\}^{n+2l+2}$. We claim that $T(Y) < \kappa$ for every $Y \in \{1, \perp\}^N$. Let \mathbf{y}_1 be a restriction of Y to the first $n + 2l + 2$ variables, \mathbf{y}_2 be a restriction of Y to the second $n + 2l + 2$ variables, and so on. Generating a random completion of Y is the same as independently sampling random completions of $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_K$. For each $i \in \{1, \dots, K\}$, the

probability that L outputs 1 on a random completion of \mathbf{y}_i is $L(\mathbf{y}_i) \leq 7/8 - \delta$. Hence, the average fraction of $i \in \{1, \dots, K\}$ such that L outputs 1 on a random completion of \mathbf{y}_i is at most $7/8 - \delta$. In turn, $T(Y)$ is the probability that the fraction of $i \in \{1, \dots, K\}$ such that L outputs 1 on a random completion of \mathbf{y}_i is at least $7/8$. Hence, by Hoeffding's inequality, $T(Y)$ is bounded by $\exp\{-2(\delta/2)^2 \cdot K\} = \kappa/2 < \kappa$.

The depth of T is (depth of L) $\times K = O(\log^2 m)$. Hence, the size of T is quasi-polynomial. The time to construct T from L is also quasi-polynomial because it is polynomial in the size of T .

Construction of L . Let x_1, \dots, x_n denote variables of φ . The variables of L will be denoted by $x_1, \dots, x_n, y_1, \dots, y_{2l+1}, z$.

First, L asks the values of y_1, \dots, y_{2l+1} . Now, we will call a binary word “fat” if it has more 1s than 0s, and “thin” otherwise. If $y_1 \dots y_{2l+1}$ is thin, L outputs 1. To describe what L does when $y_1 \dots y_{2l+1}$ is fat, we first notice that in $\{0, 1\}^{2k+1}$ exactly half of the words are fat and exactly half are thin. In particular, the number of fat words is 2^{2l} . We fix any surjective mapping from the set of fat words in $\{0, 1\}^{2l+1}$ to the set of clauses of φ . Such mapping exists because $m \leq 2^{2l}$ by our choice of l .

So, when $y_1 \dots y_{2l+1}$ is fat, L takes the clause C , assigned to $y_1 \dots y_{2l+1}$, and runs a decision tree L_C , defined as follows. If C is over variables x_{i_1}, \dots, x_{i_k} , then L_C is over x_{i_1}, \dots, x_{i_k} and z . The tree L_C asks the values of all its variables. It outputs 1 if and only if either there is exactly one 0 among x_{i_1}, \dots, x_{i_k} or ($x_{i_1} = \dots = x_{i_k} = 1$ and $z = 0$).

The depth of L is $2l + 1 + k + 1 = O(l) = O(\log m)$, and it can be constructed from φ in polynomial time.

Consider any partial input \mathbf{p} to L_C in which no coordinate is equal to 0 (i.e. we can only have variables that are fixed to 1 and undefined variables). We call \mathbf{p} good for C if z is undefined in \mathbf{p} and there is also exactly one undefined variable among x_{i_1}, \dots, x_{i_k} in \mathbf{p} . Otherwise, we call \mathbf{p} bad for C .

Lemma 5.4. *If \mathbf{p} is good for C , we have $L_C(\mathbf{p}) = 3/4$, and if \mathbf{p} is bad for C , $L_C(\mathbf{p}) \leq 5/8$.*

PROOF OF LEMMA 5.4. Assume first that \mathbf{p} is good for C . Let x_{i_j} be a variable among x_{i_1}, \dots, x_{i_k} which is undefined in \mathbf{p} . Consider a random completion of \mathbf{p} . If x_{i_j} is 0 in this completion, then there is exactly one 0 among x_{i_1}, \dots, x_{i_k} , and hence L_C outputs 1. If x_{i_j} is 1, L_C outputs 1 with probability $1/2$, depending on whether $z = 0$. Overall, we get $L_C(\mathbf{p}) = (1/2) \cdot 1 + (1/2) \cdot (1/2) = 3/4$.

Now, consider the case when \mathbf{p} is bad for C . Assume first that there are $t \geq 2$ undefined variables among x_{i_1}, \dots, x_{i_k} in \mathbf{p} . Consider a random completion of \mathbf{p} . By definition, L_C outputs 1 on it only in the following two cases:

- (a) There is exactly 1 undefined variable among x_{i_1}, \dots, x_{i_k} which is equal to 0 in our random completion;
- (b) all undefined variables among x_{i_1}, \dots, x_{i_k} are equal to 1 and $z = 0$ in our random completion.

The probability of (a) is $t2^{-t}$. The probability of (b) is 2^{-t-1} if z is undefined and 0 otherwise. Overall, we get $L_C(\mathbf{p}) \leq t2^{-t} + 2^{-t-1}$. This expression decreases in t , and for $t = 2$ it is equal to $5/8$.

Now, assume that the number of undefined variables among x_{i_1}, \dots, x_{i_k} is at most 1. If it is 1, then z has to be fixed to 1 in \mathbf{p} (otherwise, \mathbf{p} is good for C). Then L_C can only output 1 if the unique undefined variable among x_{i_1}, \dots, x_{i_k} is 0, and this happens with probability $1/2$ in a random completion of \mathbf{p} . That is, in this case, $L_C(\mathbf{p}) = 1/2$. Now, if all variables x_{i_1}, \dots, x_{i_k} are fixed to 1s in \mathbf{p} , then L_C outputs 1 only if $z = 0$. This either has probability $1/2$ (if z is undefined) or 0 (if z is fixed to 1). Hence, in this case, $L_C(\mathbf{p}) \leq 1/2$. \square

Assume first that φ is satisfiable. Fix any satisfying assignment α to φ . We construct a partial input $\mathbf{p} \in \{1, \perp\}^{n+2l+2}$ to L such that $L(\mathbf{p}) = 7/8$. We define \mathbf{p} on x_1, \dots, x_n as follows:

$$x_i = 0 \text{ in } \alpha \implies x_i = 1 \text{ in } \mathbf{p} \quad (5.5)$$

$$x_i = 1 \text{ in } \alpha \implies x_i = \perp \text{ in } \mathbf{p} \quad (5.6)$$

Variables y_1, \dots, y_{2l+1}, z are undefined in \mathbf{p} . Since α is a satisfying assignment to φ , for every clause C there is exactly one variable in C which is undefined in \mathbf{p} . This means that the restriction of \mathbf{p} to variables of L_C is

good in the sense of Lemma 5.4, for every clause of C . Thus, conditioned on the event that $y_1 \dots y_{2l+1}$ is fat, L outputs 1 with probability $3/4$ on a random completion of \mathbf{p} . Now, variables y_1, \dots, y_{2l+1} are undefined in \mathbf{p} , which means that they are all sampled independently uniformly at random in our completion. The probability of the event “ $y_1 \dots y_{2l+1}$ is fat” is $1/2$, which contributes $(1/2) \cdot (3/4) = 3/8$ to $L(\mathbf{p})$. In turn, conditioned on the event “ $y_1 \dots y_{2l+1}$ is thin” (which also happens with probability $1/2$), our tree always outputs 1. Overall, we get $L(\mathbf{p}) = 3/8 + 1/2 = 7/8$.

We now show that if it is impossible to satisfy more than a half of clauses of φ , then for every $\mathbf{p} \in \{1, \perp\}^{n+2l+2}$ we have $L(\mathbf{p}) \leq 7/8 - \delta$, where $\delta = c/\sqrt{\ln(m)}$ and $c > 0$ is some absolute constant. Again, we do so by taking a random completion of \mathbf{p} and showing that the probability that L outputs 1 on it is at most $7/8 - \delta$.

Assume first that y_1, \dots, y_{2l+1} are all undefined in \mathbf{p} . We turn \mathbf{p} into an assignment α to φ by reversing (5.5–5.6):

$$\begin{aligned} x_i = 1 \text{ in } \mathbf{p} &\implies x_i = 0 \text{ in } \alpha \\ x_i = \perp \text{ in } \mathbf{p} &\implies x_i = 1 \text{ in } \alpha \end{aligned}$$

Observe that if α does not satisfy a clause C , then the restriction of \mathbf{p} to variables of L_C is bad for C in the sense of Lemma 5.4. Now, the number of unsatisfied clauses is at least $m/2$. Each of these clauses is assigned to some fat $y_1 \dots y_{2l+1}$. That is, for at least $m/2$ fat $y_1 \dots y_{2k+1}$, the probability that L outputs 1 (conditioned on this fixation of $y_1 \dots y_{2k+1}$) is at most $5/8$. For any other fat $y_1 \dots y_{2k+1}$, this conditional probability is at most $3/4$. Overall, the contribution of fat $y_1 \dots y_{2k+1}$ to $L(\mathbf{p})$ does not exceed

$$\frac{m/2}{2^{2l+1}} \cdot (5/8) + \frac{2^{2l} - m/2}{2^{2l+1}} \cdot (3/4),$$

and the contribution of thin $y_1 \dots y_{2k+1}$ to $L(\mathbf{p})$ is, as before, $1/2$. Overall, we get

$$L(\mathbf{p}) \leq \frac{m/2}{2^{2l+1}} \cdot (5/8) + \frac{2^{2l} - m/2}{2^{2l+1}} \cdot (3/4) + 1/2 = 7/8 - \frac{m}{2^{2l+5}}.$$

It remains to recall that l was chosen as the smallest integral number such that $m \leq 2^{2l}$. This means that $m \geq 2^{2(l-1)}$. This gives us $L(\mathbf{p}) \leq 7/8 - 1/128 < 7/8 - \delta$.

It remains to consider the case when in \mathbf{p} at least one variable among y_1, \dots, y_{2l+1} is fixed to 1. As before, the probability that L outputs 1, conditioned on the event “ y_1, \dots, y_{2k+1} is fat”, is at most $3/4$. And this probability is 1 conditioned on the event “ y_1, \dots, y_{2k+1} is thin”. Now, however, since at least one variable among y_1, \dots, y_{2k+1} is fixed to 1, the probability to be thin is slightly less than $1/2$, which makes $L(\mathbf{p})$ slightly less than $7/8$. More specifically,

$$\begin{aligned} L(\mathbf{p}) &\leq \Pr[y_1 \dots y_{2l+1} \text{ is thin}] + (3/4) \Pr[y_1 \dots y_{2l+1} \text{ is fat}] \\ &= 1 - (1/4) \Pr[y_1 \dots y_{2l+1} \text{ is fat}], \end{aligned}$$

where y_1, \dots, y_{2l+1} are chosen from a random completion of \mathbf{p} . That is, now our goal is to lower bound $\Pr[y_1 \dots y_{2l+1} \text{ is fat}]$, assuming that at least one of y_1, \dots, y_{2l+1} is fixed to 1, and the rest of them are chosen independently uniformly at random. The more variables we fix to 1, the larger becomes $\Pr[y_1 \dots y_{2l+1} \text{ is fat}]$. Hence, without loss of generality we may assume that exactly one variable is fixed to 1, say, y_{2l+1} . Then $y_1 \dots y_{2l+1}$

is fat if and only if there are at least l ones among y_1, \dots, y_{2l} . In conclusion, we get

$$\begin{aligned}
 \Pr[y_1 \dots y_{2l+1} \text{ is fat}] &\geq \frac{\binom{2l}{2l} + \dots + \binom{2l}{l+1} + \binom{2l}{l}}{2^{2l}} \\
 &= \frac{\frac{\binom{2l}{2l} + \binom{2l}{0}}{2} + \dots + \frac{\binom{2l}{l+1} + \binom{2l}{l-1}}{2} + \frac{\binom{2l}{l} + \binom{2l}{l}}{2}}{2^{2l}} \\
 &= \frac{\binom{2l}{2l} + \dots + \binom{2l}{0}}{2 \cdot 2^{2l}} + \frac{\binom{2l}{l}}{2 \cdot 2^{2l}} = \frac{1}{2} + \Omega(1/\sqrt{l}) \\
 &= \frac{1}{2} + \Omega(1/\sqrt{\ln m}).
 \end{aligned}$$

From this, we get $L(\mathbf{p}) \leq 1 - (1/4) \cdot (1/2 + \Omega(1/\sqrt{\ln m})) = 7/8 - \Omega(1/\sqrt{\ln m})$, as required. \square

6 Proofs of the Tractability Results

We now restate and proof our main positive results.

THEOREM 3.6 (RESTATED). *Both problems Compute-Minimum-SR and Compute-Minimal-SR can be solved in polynomial time for decision trees with split number at most c , where $c \geq 1$ is any fixed integer.*

PROOF OF THEOREM 3.6. It suffices to provide a polynomial time algorithm for Compute-Minimum-SR. (The same algorithm works for Compute-Minimal-SR as a minimum δ -SR is in particular minimal.) In turn, using standard arguments, it is enough to provide a polynomial time algorithm for the following decision problem Check-Minimum-SR: Given a tuple $(T, \mathbf{y}, \delta, k)$, where T is a decision tree of dimension n , $\mathbf{y} \in \{0, 1, \perp\}^n$ is a partial instance, $\delta \in (0, 1]$, and $k \geq 0$ is an integer, decide whether there is a partial instance $\mathbf{y}' \subseteq \mathbf{y}$ such that $n - |\mathbf{y}'|_{\perp} \leq k$ (i.e., \mathbf{y}' has at most k defined components) and

$$\Pr_z[T(z) = 1 \mid z \in \text{COMP}(\mathbf{y}')] \geq \delta.$$

To solve Check-Minimum-SR over an instance $(T, \mathbf{y}, \delta, k)$, where T has split number at most c , we apply dynamic programming over T in a bottom-up manner. Let $Z \subseteq \{1, \dots, n\}$ be the set of features defined in \mathbf{y} , that is, features i with $\mathbf{y}[i] \neq \perp$. Those are the features we could eventually remove when looking for \mathbf{y}' . For each node u in T , we solve a polynomial number of subproblems over the subtree T_u . We define

$$\text{Int}(u) := \mathcal{F}(N_u^{\downarrow}) \cap \mathcal{F}(N_u^{\uparrow}) \cap Z \quad \text{New}(u) := \left(\mathcal{F}(N_u^{\downarrow}) \setminus \text{Int}(u)\right) \cap Z.$$

In other words, $\text{Int}(u)$ are the features appearing both inside and outside T_u , while $\text{New}(u)$ are the features only inside T_u , that is, the new features introduced below u . Both sets are restricted to Z as features not in Z play no role in the process.

Each particular subproblem is indexed by a possible size $s \in \{0, \dots, k\}$ and a possible set $J \subseteq \text{Int}(u)$ with $|J| \leq s$, and the goal is to compute the quantity:

$$p_{u,s,J} := \max_{\mathbf{y}' \in C_{u,s,J}} \Pr_z[T_u(z) = 1 \mid z \in \text{COMP}(\mathbf{y}')],$$

where $C_{u,s,J}$ is the space of partial instances $\mathbf{y}' \subseteq \mathbf{y}$ with $n - |\mathbf{y}'|_{\perp} \leq s$ and such that $\mathbf{y}'[i] = \mathbf{y}[i]$ for $i \in J$ and $\mathbf{y}'[i] = \perp$ for $i \in \text{Int}(u) \setminus J$. In other words, the set J fixes the behavior on $\text{Int}(u)$ (keep features in J , remove features in $\text{Int}(u) \setminus J$) and hence the maximization occurs over choices on the set $\text{New}(u)$ (which features are kept and which features are removed). The key idea is that $p_{u,s,J}$ can be computed inductively using the information already computed for the children u_1 and u_2 of u . Intuitively, this holds since the common features between T_{u_1} and T_{u_2} are at most c , which is a fixed constant, and hence we can efficiently synchronize the information stored for u_1 and u_2 . Finally, to solve the instance $(T, \mathbf{y}, \delta, k)$ we simply check whether $p_{r,k,0} \geq \delta$, for the root r of T .

Formally, let us define for a set $H \subseteq Z$, the partial instance $\mathbf{y}_H \in \{0, 1, \perp\}^n$ such that $\mathbf{y}_H[i] = \mathbf{y}[i]$ for every $i \in H$, and $\mathbf{y}_H[i] = \perp$ for every $i \notin H$. In particular, $\mathbf{y}_H \subseteq \mathbf{y}$. Then we can write $p_{u,s,J}$ as

$$p_{u,s,J} = \max_{\substack{K \subseteq \text{New}(u) \\ |K| \leq s - |J|}} \Pr_z[T_u(z) = 1 \mid z \in \text{COMP}(\mathbf{y}_{J \cup K})].$$

Let u_1 and u_2 be the children of u . We have that $\text{New}(u)$ is the disjoint union of:

$$\text{New}(u) = \text{New}(u_1) \cup \text{New}(u_2) \cup \text{Sync}(u),$$

where $\text{Sync}(u) := \text{New}(u) \cap \left(\mathcal{F}(N_{u_1}^\perp) \cap \mathcal{F}(N_{u_2}^\perp) \right)$. In other words, the features in $\text{Sync}(u)$ are the features that are in both T_{u_1} and T_{u_2} but not outside T_u . We conclude by explaining the computation of $p_{u,s,J}$. We consider the following cases:

- (1) The feature i labeling u is in J . This means we have to keep feature i . If $\mathbf{y}[i] = 0$, then to compute $p_{u,s,J}$ we can simply look at u_1 (the left child). Note that $\text{Int}(u_1)$ is the disjoint union of $\text{Int}(u_1) \cap \text{Int}(u)$ and $\text{Sync}(u)$. Then

$$p_{u,s,J} = \max_{\substack{J' \subseteq \text{Sync}(u) \\ |J'| \leq s - |\text{Int}(u_1) \cap J|}} p_{u_1,s,(\text{Int}(u_1) \cap J) \cup J'}.$$

This computation can be done in polynomial time as $\text{Sync}(u) \leq c$ and then there are a constant number of possible $J' \subseteq \text{Sync}(u)$. The case when $\mathbf{y}[i] = 1$ is analogous, taking u_2 instead of u_1 .

- (2) The feature i labeling u is either outside Z or belongs to $\text{Int}(u) \setminus J$. This means feature i is undefined. Again, we have that $\text{Int}(u_1)$ is the disjoint union of $\text{Int}(u_1) \cap \text{Int}(u)$ and $\text{Sync}(u)$. Similarly, $\text{Int}(u_2)$ is the disjoint union of $\text{Int}(u_2) \cap \text{Int}(u)$ and $\text{Sync}(u)$. Then

$$p_{u,s,J} = \max_{\substack{J' \subseteq \text{Sync}(u) \\ |J'| \leq s - |J|}} \max_{\substack{0 \leq s_1, s_2 \leq s \\ s_1 + s_2 \leq s - |J| - |J'|}} \frac{1}{2} \cdot p_{u_1,s_1,(\text{Int}(u_1) \cap J) \cup J'} + \frac{1}{2} \cdot p_{u_2,s_2,(\text{Int}(u_2) \cap J) \cup J'}.$$

Again, this can be done in polynomial time as $\text{Sync}(u) \leq c$.

- (3) Finally, the remaining case is that the feature i labeling u is in $\text{New}(u)$. In that case we have the two possibilities: either we keep feature i or we remove it. If $s - |J| = 0$, then the only possible choice is to remove the feature i , and hence $p_{u,s,J}$ is computed exactly as in case (2). If $s - |J| > 0$. Then we take the maximum between the cases when we keep feature i and the case when we remove feature i . For the latter, $p_{u,s,J}$ is computed exactly as in case (2). For the former, we compute $p_{u,s,J}$ in a similar way as in case (1). More precisely, if $\mathbf{y}[i] = 0$, then:

$$p_{u,s,J} = \max_{\substack{J' \subseteq \text{Sync}(u) \\ |J'| \leq s - 1 - |\text{Int}(u_1) \cap J|}} p_{u_1,s-1,(\text{Int}(u_1) \cap J) \cup J'}.$$

The case $\mathbf{y}[i] = 1$ is analogous.

This finishes the proof of the theorem. □

THEOREM 3.7 (RESTATED). *Let \mathfrak{C} be a class of monotone Boolean models such that \mathfrak{C} -#Positive-Completions can be solved in polynomial time. Then the problem Compute-Minimal-SR can be solved in polynomial time over \mathfrak{C} .*

PROOF OF THEOREM 3.7. Let us first introduce some notation. For any partial instance \mathbf{y} of dimension n and integer $i \in \{1, \dots, n\}$, if $\mathbf{y}[i] \neq \perp$ then we use $\mathbf{y} \setminus \{i\}$ to denote the partial instance \mathbf{y}' that is exactly equal to \mathbf{y} except for $\mathbf{y}'[i] = \perp$. Now, to prove Theorem 3.7, we will use the following lemma, which is a probabilistic counterpart to Proposition 2.2.

Lemma 6.1. *Let \mathfrak{C} be a class of monotone models, $\mathcal{M} \in \mathfrak{C}$ a model of dimension n , and $\mathbf{x} \in \{0, 1\}^n$ an instance. Consider any $\delta \in (0, 1]$. Then if $\mathbf{y} \subseteq \mathbf{x}$ is a δ -SR for \mathbf{x} under \mathcal{M} which is not minimal, then there is a partial instance $\mathbf{y} \setminus \{i\}$, for some $i \in \{1, \dots, n\}$, that is also a δ -SR for \mathbf{x} under \mathcal{M} .*

After proving Lemma 6.1, we can easily prove Theorem 3.7. Indeed, it suffices to make a slight modification to Algorithm 2: instead of “CheckSufficientReason($T, \hat{\mathbf{y}}, \mathbf{x}$)”, we need to “CheckProbSufficientReason($T, \hat{\mathbf{y}}, \mathbf{x}, \delta$)”, which is easy since computing the probability associated to a partial instance can be done in polynomial time [12].⁶ Therefore, one can compute in polynomial time a minimal δ -SR for \mathbf{x} under \mathcal{M} . This finishes the proof of Theorem 3.7. \square

We now prove Lemma 6.1, the key for Theorem 3.7.

PROOF OF LEMMA 6.1. Note that, if $\mathcal{M}(\mathbf{x}) = 1$ then we can safely assume that for every i where $\mathbf{y}[i] \neq \perp$ it holds that $\mathbf{y}[i] = 1$, as otherwise if $\mathbf{y}[i^*] = 0$ for some i^* , then the lemma trivially holds by setting $\mathbf{y}' = \mathbf{y} \setminus \{i^*\}$ because of monotonicity. Similarly, if $\mathcal{M}(\mathbf{x}) = 0$ then we can safely assume that for every i where $\mathbf{y}[i] \neq \perp$ it holds that $\mathbf{y}[i] = 0$.

As by hypothesis \mathbf{y} is not minimal, there exists a δ -SR $\mathbf{y}^* \subsetneq \mathbf{y}$ that minimizes $|\mathbf{y}^*|_{\perp}$. We prove that $|\mathbf{y}^*|_{\perp} = |\mathbf{y}|_{\perp} + 1$, from where the lemma immediately follows.

Assume for the sake of a contradiction that $|\mathbf{y}^*|_{\perp} > |\mathbf{y}|_{\perp} + 1$. Then, there must exist a feature i^* that $\mathbf{y}^*[i^*] = \perp \neq \mathbf{y}[i^*]$, and such that $\mathbf{y}^* \cup \{i^*\} \neq \mathbf{y}$, where $\mathbf{y}^* \cup \{i^*\}$ is defined as

$$(\mathbf{y}^* \cup \{i^*\})[i] = \begin{cases} \mathbf{y}[i^*] & \text{if } i = i^* \\ \mathbf{y}^*[i^*] & \text{otherwise.} \end{cases}$$

Similarly, we denote $\mathbf{y}^* \cup (i^* \rightarrow \alpha)$, with $\alpha \in \{0, 1\}$, the partial instance defined as

$$(\mathbf{y}^* \cup (i^* \rightarrow \alpha))[i] = \begin{cases} \alpha & \text{if } i = i^* \\ \mathbf{y}^*[i^*] & \text{otherwise.} \end{cases}$$

We now claim that $\mathbf{y}^* \cup \{i^*\}$ is also a δ -SR for \mathbf{x} under \mathcal{M} , which contradicts the minimality of \mathbf{y}^* , as $|\mathbf{y}^* \cup \{i^*\}|_{\perp} < |\mathbf{y}^*|_{\perp}$. Let us denote by $C(\mathcal{M}, \mathbf{y})$ the number of completions $\mathbf{z} \in \text{COMP}(\mathbf{y})$ such that $\mathcal{M}(\mathbf{z}) = 1$. Now there are two cases, if $\mathcal{M}(\mathbf{x}) = 1$ then

$$\begin{aligned} C(\mathcal{M}, \mathbf{y}^*) &= C(\mathcal{M}, \mathbf{y}^* \cup (i^* \rightarrow 0)) + C(\mathcal{M}, \mathbf{y}^* \cup (i^* \rightarrow 1)) \\ &\leq 2C(\mathcal{M}, \mathbf{y}^* \cup (i^* \rightarrow 1)), \end{aligned}$$

where the inequality holds because of monotonicity. This implies that

$$\frac{C(\mathcal{M}, \mathbf{y}^* \cup (i^* \rightarrow 1))}{2^{|\mathbf{y}^* \cup (i^* \rightarrow 1)|_{\perp}}} = \frac{C(\mathcal{M}, \mathbf{y}^* \cup (i^* \rightarrow 1))}{2^{|\mathbf{y}^*|_{\perp} - 1}} \geq \frac{C(\mathcal{M}, \mathbf{y}^*)}{2 \cdot 2^{|\mathbf{y}^*|_{\perp} - 1}} \geq \delta,$$

which implies that $\mathbf{y}^* \cup (i^* \rightarrow 1)$ is also a δ -SR (note that $\mathbf{y}^* \cup (i^* \rightarrow 1) = \mathbf{y}^* \cup \{i^*\}$ because of the initial observation), contradicting the minimality of \mathbf{y}^* . Similarly, if $\mathcal{M}(\mathbf{x}) = 0$, then

$$\begin{aligned} C(\mathcal{M}, \mathbf{y}^*) &= C(\mathcal{M}, \mathbf{y}^* \cup (i^* \rightarrow 0)) + C(\mathcal{M}, \mathbf{y}^* \cup (i^* \rightarrow 1)) \\ &\geq 2C(\mathcal{M}, \mathbf{y}^* \cup (i^* \rightarrow 0)), \end{aligned}$$

⁶This can be easily seen by observing that with dynamic programming we can easily count the number of “completions” of a partial instance with a given classification [9].

from where

$$\begin{aligned} \frac{2^{|\mathbf{y}^* \cup (i^* \rightarrow 0)|_{\perp}} - C(\mathcal{M}, \mathbf{y}^* \cup (i^* \rightarrow 0))}{2^{|\mathbf{y}^* \cup (i^* \rightarrow 0)|_{\perp}}} &= \frac{2^{|\mathbf{y}^*|_{\perp}-1} - C(\mathcal{M}, \mathbf{y}^* \cup (i^* \rightarrow 0))}{2^{|\mathbf{y}^*|_{\perp}-1}} \\ &\geq 1 - \frac{C(\mathcal{M}, \mathbf{y}^*)}{2 \cdot 2^{|\mathbf{y}^*|_{\perp}-1}} \geq \delta, \end{aligned}$$

thus implying that $\mathbf{y}^* \cup (i^* \rightarrow 0)$ is also a δ -SR for x under \mathcal{M} , which again contradicts the minimality of \mathbf{y}^* . \square

7 Final Remarks and Open Problems

We have settled the complexity of two open problems in formal XAI, proving that both minimal and minimum size probabilistic explanations can be hard to compute, and even approximate, for decision trees. These results further support the idea that decision trees may not always be interpretable [9, 14, 28, 19], while being the first results to do so even in a probabilistic setting. Our study has focused on Boolean models, and moreover it is based on the assumption that features are independent and identically distributed, with probability $1/2$. Naturally, our hardness results in this restricted setting imply hardness in more general settings, but we leave as future work the study of tractable cases (e.g., bounded split-number trees, monotone classes that allow counting) under general, potentially correlated, distributions.

The results proven in this paper suggest that minimal (or minimum) explanations might be hard to obtain in practice even for decision trees, especially in problems where the feature space has large dimension. A way to circumvent the limitations proven in our work is to relax the guarantee of minimality, or introduce some probability of error, as done in the work of [18]. More recently, the concept of (δ, ε) -sufficient reasons has been proposed [29], but such a relaxation will generally not be helpful for tractability over decision trees given the inapproximability results from Theorems 3.4 and 3.5. A promising direction of future research is to better understand the kind of guarantees and settings in which is still possible to obtain tractability.

Finally, our work leaves open some interesting technical questions:

- (1) What is the *parameterized* complexity of computing minimum δ -SRs over decision trees, assuming that the parameter is the size of the explanation one is looking for? It is not hard to see that W[2] hardness follows from our proofs (i.e., they are parameterized reductions all the way to Set Cover), but membership in any class is fully open. Parameterized complexity might be of particular interest for this class of problems as explanations might reasonably be expected to be very small in practice.
- (2) Does Theorem 3.1 continue to hold for monotone decision trees? That is, is it the case that computing minimum δ -SRs over monotone decision trees is hard for every fixed $\delta \in (0, 1]$?
- (3) Is it the case that the hardness of computation for minimal sufficient reasons holds for a *fixed* $\delta \in (0, 1]$? If so, does it hold for every fixed such a δ , or only for some?
- (4) Is it possible to show the non-approximability results in Theorems 3.4 and 3.5 under the more standard assumption that $P \neq NP$?
- (5) Is it possible to extend the positive behavior of decision trees with bounded split number to more powerful Boolean ML models such as free BDDs?

Acknowledgements

Arenas and Barceló are funded by ANID - Millennium Science Initiative Program - Code ICN17002. Barceló, Kozachinskiy and Romero are supported by the National Center for Artificial Intelligence CENIA FB210017, Basal ANID. Kozachinskiy is supported by the ANID Fondecyt Iniciación grant 11250060. Subercaseaux is supported by the U.S. National Science Foundation under grant DMS-2434625.

References

- [1] A. B. Arrieta et al. 2020. Explainable artificial intelligence (XAI): concepts, taxonomies, opportunities and challenges toward responsible AI. *Information Fusion*, 58, (June 2020), 82–115. doi: [10.1016/j.inffus.2019.12.012](https://doi.org/10.1016/j.inffus.2019.12.012).
- [2] A. Shih, A. Choi, and A. Darwiche. 2018. A symbolic approach to explaining bayesian network classifiers. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence (IJCAI'18)*. AAAI Press, Stockholm, Sweden, 5103–5111. ISBN: 9780999241127.
- [3] M. T. Ribeiro, S. Singh, and C. Guestrin. 2018. Anchors: high-precision model-agnostic explanations. In *AAAI*, 1527–1535.
- [4] S. M. Lundberg and S.-I. Lee. 2017. A unified approach to interpreting model predictions. In *NeurIPS*, 4765–4774.
- [5] T. Yan and A. D. Procaccia. 2021. If you like shapley then you'll love the core. In *AAAI*, 5751–5759.
- [6] J. Marques-Silva and A. Ignatiev. 2022. Delivering trustworthy AI through Formal XAI. In *AAAI*.
- [7] A. Ignatiev, N. Narodytska, N. Asher, and J. Marques-Silva. 2021. From contrastive to abductive explanations and back again. In *AIxIA 2020 – Advances in Artificial Intelligence*. Springer International Publishing, 335–355. doi: [10.1007/978-3-030-77091-4_21](https://doi.org/10.1007/978-3-030-77091-4_21).
- [8] J. Marques-Silva, T. Gerspacher, M. C. Cooper, A. Ignatiev, and N. Narodytska. 2020. Explaining naive bayes and other linear classifiers with polynomial time and delay. In *NeurIPS*.
- [9] P. Barceló, M. Monet, J. Pérez, and B. Subercaseaux. 2020. Model interpretability through the lens of computational complexity. In *NeurIPS*, 15487–15498.
- [10] J. Marques-Silva, T. Gerspacher, M. C. Cooper, A. Ignatiev, and N. Narodytska. 2021. Explanations for monotonic classifiers. In *ICML*, 7469–7479.
- [11] S. Wäldchen, J. MacDonald, S. Hauch, and G. Kutyniok. 2021. The computational complexity of understanding binary classifier decisions. *J. Artif. Intell. Res.*, 70, 351–387.
- [12] Y. Izza, A. Ignatiev, N. Narodytska, M. C. Cooper, and J. Marques-Silva. 2021. Efficient explanations with relevant sets. *ArXiv*, abs/2106.00546.
- [13] E. Wang, P. Khosravi, and G. V. den Broeck. 2021. Probabilistic sufficient explanations. In *IJCAI*. Z. Zhou, (Ed.), 3082–3088.
- [14] Z. C. Lipton. 2018. The mythos of model interpretability. *Queue*, 16, 3, 31–57.
- [15] Y. Izza, A. Ignatiev, and J. Marques-Silva. 2020. On explaining decision trees. (2020). doi: [10.48550/arXiv.2010.11034](https://doi.org/10.48550/arXiv.2010.11034).
- [16] L. H. Gilpin, D. Bau, B. Z. Yuan, A. Bajwa, M. A. Specter, and L. Kagal. 2018. Explaining explanations: an overview of interpretability of machine learning. In *DSAA*. F. Bonchi, F. J. Provost, T. Eliassi-Rad, W. Wang, C. Cattuto, and R. Ghani, (Eds.), 80–89.
- [17] Y. Izza, A. Ignatiev, N. Narodytska, M. C. Cooper, and J. Marques-Silva. 2022. Provably precise, succinct and efficient explanations for decision trees. (2022). doi: [10.48550/arXiv.2205.09569](https://doi.org/10.48550/arXiv.2205.09569).
- [18] G. Blanc, J. Lange, and L.-Y. Tan. 2021. Provably efficient, succinct, and precise explanations. In *NeurIPS*. A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan, (Eds.)
- [19] G. Audemard, S. Bellart, L. Bounia, F. Koriche, J.-M. Lagniez, and P. Marquis. 2022. On the explanatory power of Boolean decision trees. *Data & Knowledge Engineering*, 142, (Nov. 2022), 102088. doi: [10.1016/j.datak.2022.102088](https://doi.org/10.1016/j.datak.2022.102088).
- [20] Y. Izza, A. Ignatiev, and J. Marques-Silva. 2022. On tackling explanation redundancy in decision trees. *Journal of Artificial Intelligence Research*, 75, (Sept. 2022), 261–321. doi: [10.1613/jair.1.13575](https://doi.org/10.1613/jair.1.13575).
- [21] R. Gomes Mantovani, T. Horváth, A. L. D. Rossi, R. Cerri, S. Barbon Junior, J. Vanschoren, and A. C. P. L. F. d. Carvalho. 2024. Better trees: an empirical study on hyperparameter tuning of classification decision tree induction algorithms. *Data Mining and Knowledge Discovery*, 38, 3, (Jan. 2024), 1364–1416. doi: [10.1007/s10618-024-01002-5](https://doi.org/10.1007/s10618-024-01002-5).
- [22] G. Audemard, S. Bellart, L. Bounia, F. Koriche, J. Lagniez, and P. Marquis. 2021. On the computational intelligibility of boolean classifiers. In *KR*, 74–86.
- [23] X. Huang, Y. Izza, A. Ignatiev, and J. Marques-Silva. 2021. On efficiently explaining graph-based classifiers. In *KR*, 356–367.
- [24] J. Goldsmith, M. Hagen, and M. Mundhenk. 2005. Complexity of DNF and isomorphism of monotone formulas. In *MFCS*. Springer Berlin Heidelberg, 410–421. doi: [10.1007/11549345_36](https://doi.org/10.1007/11549345_36).
- [25] I. Wegener. 2000. *Branching Programs and Binary Decision Diagrams*. Society for Industrial and Applied Mathematics, (Jan. 2000). doi: [10.1137/1.9780898719789](https://doi.org/10.1137/1.9780898719789).
- [26] W. Hoeffding. 1963. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58, 301, 13–30.
- [27] V. Guruswami and L. Trevisan. 2005. The complexity of making unique choices: approximating 1-in-k sat. In *Proceedings of the 8th international workshop on Approximation, Randomization and Combinatorial Optimization Problems, and Proceedings of the 9th international conference on Randomization and Computation: algorithms and techniques*, 99–110.
- [28] M. Arenas, D. Báez, P. Barceló, J. Pérez, and B. Subercaseaux. 2021. Foundations of symbolic languages for model interpretability. In *NeurIPS*, 11690–11701.
- [29] B. Subercaseaux, M. Arenas, and K. S. Meel. 2024. Probabilistic explanations for linear models. (2024). <https://arxiv.org/abs/2501.00154> arXiv: 2501.00154 [cs.AI].

Received 27 October 2024; accepted 13 July 2025