

---

# Trading Complexity for Sparsity in Random Forest Explanations

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Random forests have long been considered as powerful model ensembles in ma-  
2 chine learning. By training multiple decision trees, whose diversity is fostered  
3 through data and feature subsampling, the resulting random forest can lead to  
4 more stable and reliable predictions than a single decision tree. This however  
5 comes at the cost of decreased interpretability: while decision trees are often easily  
6 interpretable, the predictions made by random forests are much more difficult to  
7 understand, as they involve a majority vote over hundreds of decision trees. In  
8 this paper, we examine different types of *reasons* that explain “why” an input  
9 instance is classified as positive or negative by a Boolean random forest. Notably,  
10 as an alternative to *sufficient reasons* taking the form of prime implicants of the  
11 random forest, we introduce *majoritary reasons* which are prime implicants of a  
12 strict majority of decision trees. For these different abductive explanations, the  
13 tractability of the generation problem (finding one reason) and the minimization  
14 problem (finding one shortest reason) are investigated. Experiments conducted on  
15 various datasets reveal the existence of a trade-off between runtime complexity and  
16 sparsity. Sufficient reasons - for which the identification problem is DP-complete  
17 - are slightly larger than majoritary reasons that can be generated using a simple  
18 linear-time greedy algorithm, and significantly larger than *minimal* majoritary  
19 reasons that can be approached using an anytime PARTIAL MAXSAT algorithm.

## 20 1 Introduction

21 Over the past two decades, rapid progress in statistical machine learning has led to the deployment  
22 of models endowed with remarkable predictive capabilities. Yet, as the spectrum of applications  
23 using statistical learning models becomes increasingly large, explanations for why a model is making  
24 certain predictions are ever more critical. For example, in medical diagnosis, if some model predicts  
25 that an image is malignant, then the doctor may need to know which features in the image have led to  
26 this classification. Similarly, in the banking sector, if some model predicts that a customer is a fraud,  
27 then the banker might want to know why. Therefore, having explanations for why certain predictions  
28 are made is essential for securing user confidence in machine learning technologies [21, 22].

29 This paper focuses on classifications made by *random forests*, a popular ensemble learning method  
30 that constructs multiple randomized decision trees during the training phase, and predicts by taking a  
31 majority vote over the base classifiers [8]. Since decision tree randomization is achieved by essentially  
32 coupling data subsampling (or bagging) and feature subsampling, random forests are fast and easy to  
33 implement, with few tuning parameters. Furthermore, they often make accurate and robust predictions  
34 in practice, even for small data samples and high-dimensional feature spaces [6]. For these reasons,  
35 random forests have been used in various applications including, among others, computer vision [11],  
36 crime prediction [7], ecology [12], genomics [9], and medical diagnosis [3].

37 Despite their success, random forests are much less interpretable than decision trees. Indeed, the  
 38 prediction made by a decision tree on a given data instance can be easily interpreted by reading the  
 39 unique root-to-leaf path that covers the instance. Contrastingly, there is no such *direct reason* in a  
 40 random forest, since the prediction is derived from a majority vote over multiple decision trees. So, a  
 41 key issue in random forests is to infer *abductive explanations*, that is, to explain in concise terms why  
 42 a data instance is classified as positive or negative by the model ensemble.

43 **Related Work.** Explaining random forest predictions has received increasing attention in recent  
 44 years [5, 10, 18]. Notably, in the classification setting, [10, 18] have focused on *sufficient reasons*,  
 45 which are abductive explanations involving only relevant features [13]. More specifically, if we view  
 46 any random forest classifier as a Boolean function  $f$ , then a sufficient reason for classifying a data  
 47 instance  $x$  as positive by  $f$  is a *prime implicant*  $t$  of  $f$  covering  $x$ . By construction, removing any  
 48 feature from a sufficient reason  $t$  would question the fact that  $t$  explains the way  $x$  is classified by  $f$ .  
 49 Interestingly, if  $f$  is described by a single decision tree, then generating a sufficient reason for any  
 50 input instance  $x$  can be done in linear time. Yet, in the general case where  $f$  is represented by an  
 51 arbitrary number of decision trees, the problem of identifying a sufficient reason is DP-complete.  
 52 Despite this intractability statement, the empirical results reported in [18] show that a MUS-based  
 53 algorithm for computing sufficient reasons proves quite efficient in practice.

54 In addition to “model-based” explanations investigated in [10, 18], “model-agnostic” explanations  
 55 can be applied to random forests. Notably, the LIME method [27] extrapolates a linear threshold  
 56 function  $g$  from the behavior of the random forest  $f$  around an input instance  $x$ . Yet, even if a prime  
 57 implicant of the linear threshold function can be easily computed, this explanation is *not* guaranteed  
 58 abductive since  $g$  is only an approximation of  $f$ .

59 **Contributions.** In this paper, we introduce several new notions of abductive explanations: *direct*  
 60 *reasons* extend to the case of random forests the corresponding notion defined primarily for decision  
 61 *trees*, and *majority reasons* are weak forms of abductive explanations which take into account the  
 62 averaging rule of random forests. Informally, a majoritary reason for classifying an instance  $x$  as  
 63 positive by some random forest  $f$  is a prime implicant  $t$  of a majority of decision trees in  $f$  that  
 64 covers  $x$ . Thus, any sufficient reason is a majoritary reason, but the converse is not true. For these  
 65 different reasons, we examine the tractability of both the generation (finding one explanation) and  
 66 the minimization (finding one shortest explanation) problems. To the best of our knowledge, all  
 67 complexity results related to random forest explanations are new, if we make an exception for the  
 68 intractability of generating sufficient reasons, which was recently established in [18]. Notably, direct  
 69 reasons and majoritary reasons can be derived in time polynomial in the size of the input (the instance  
 70 and the random forest used to classify it). By contrast, the identification of minimal majoritary  
 71 reasons is NP-complete, and the identification of minimal sufficient reasons is  $\Sigma_2^P$ -complete.

72 Based on these results, we provide algorithms for deriving random forest explanations, which open the  
 73 way for an empirical comparison. Our experiments made on standard benchmarks show the existence  
 74 of a trade-off between the runtime complexity of finding (possibly minimal) abductive explanations  
 75 and the sparsity of such explanations. In a nutshell, majoritary reasons and minimal majoritary  
 76 reasons offer interesting compromises in comparison to, respectively, sufficient reasons and minimal  
 77 sufficient reasons. Indeed, the size of majoritary reasons and the computational effort required to  
 78 generate them are generally smaller than those obtained for sufficient reasons. Furthermore, minimal  
 79 majoritary reasons outperform minimal sufficient reasons, since the latter are too computationally  
 80 demanding. In fact, using an *anytime* PARTIAL MAXSAT solver for minimizing majoritary reasons,  
 81 we derive sparse explanations which are typically much shorter than all other forms of abductive  
 82 explanations. Proofs and additional empirical results are provided as supplementary material.

## 83 2 Preliminaries

84 For an integer  $n$ , let  $[n] = \{1, \dots, n\}$ . By  $\mathcal{F}_n$  we denote the class of all Boolean functions from  
 85  $\{0, 1\}^n$  to  $\{0, 1\}$ , and we use  $X_n = \{x_1, \dots, x_n\}$  to denote the set of input Boolean variables. Any  
 86 Boolean vector  $x \in \{0, 1\}^n$  is called an *instance*. For any function  $f \in \mathcal{F}_n$ , an instance  $x \in \{0, 1\}^n$   
 87 is called a *positive example* of  $f$  if  $f(x) = 1$ , and a *negative example* otherwise.

88 We refer to  $f$  as a propositional formula when it is described using the Boolean connectives  $\wedge$   
 89 (conjunction),  $\vee$  (disjunction) and  $\neg$  (negation), together with the constants 1 (true) and 0 (false). As

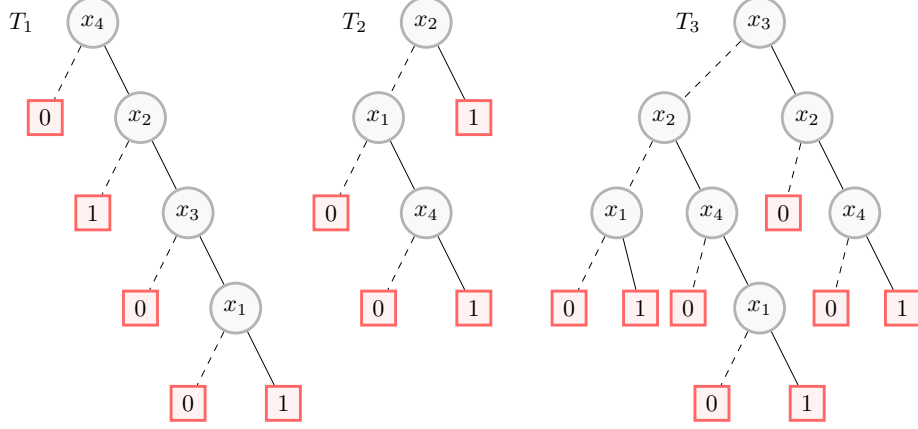


Figure 1: A random forest  $F = \{T_1, T_2, T_3\}$  for recognizing *Cattleya* orchids. The left (resp. right) child of any decision node labelled by  $x_i$  corresponds to the assignment of  $x_i$  to 0 (resp. 1).

90 usual, a *literal*  $l_i$  is a variable  $x_i$  or its negation  $\neg x_i$ , also denoted  $\bar{x}_i$ . A *term* (or *monomial*)  $t$  is a  
 91 conjunction of literals, and a *clause*  $c$  is a disjunction of literals. A DNF *formula* is a disjunction of  
 92 terms and a CNF *formula* is a conjunction of clauses. The set of variables occurring in a formula  $f$  is  
 93 denoted  $\text{Var}(f)$ . In the rest of the paper, we shall often treat instances as terms, and terms as sets of  
 94 literals. Given an assignment  $\mathbf{z} \in \{0, 1\}^n$ , the corresponding term is defined as

$$t_{\mathbf{z}} = \bigwedge_{i=1}^n x_i^{z_i} \text{ where } x_i^0 = \bar{x}_i \text{ and } x_i^1 = x_i$$

95 A term  $t$  *covers* an assignment  $\mathbf{z}$  if  $t \subseteq t_{\mathbf{z}}$ . An *implicant* of a Boolean function  $f$  is a term that  
 96 implies  $f$ , that is, a term  $t$  such that  $f(\mathbf{z}) = 1$  for every assignment  $\mathbf{z}$  covered by  $t$ . A *prime implicant*  
 97 of  $f$  is an implicant  $t$  of  $f$  such that no proper subset of  $t$  is an implicant of  $f$ .

98 With these basic notions in hand, a (Boolean) *decision tree* on  $X_n$  is a binary tree  $T$ , each of whose  
 99 internal nodes is labeled with one of  $n$  input variables, and whose leaves are labeled 0 or 1. Every  
 100 variable is supposed (w.l.o.g.) to occur at most once on any root-to-leaf path (read-once property).  
 101 The value  $T(\mathbf{x}) \in \{0, 1\}$  of  $T$  on an input instance  $\mathbf{x}$  is given by the label of the leaf reached from  
 102 the root as follows: at each node go to the left or right child depending on whether the input value of  
 103 the corresponding variable is 0 or 1, respectively. A (Boolean) *random forest* on  $X_n$  is an ensemble  
 104  $F = \{T_1, \dots, T_m\}$ , where each  $T_i$  ( $i \in [m]$ ) is a decision tree on  $X_n$ , and such that the value  
 105  $F(\mathbf{x}) \in \{0, 1\}$  on an input instance  $\mathbf{x}$  is given by

$$F(\mathbf{x}) = \begin{cases} 1 & \text{if } \frac{1}{m} \sum_{i=1}^m T_i(\mathbf{x}) > \frac{1}{2} \\ 0 & \text{otherwise.} \end{cases}$$

106 The size of  $F$  is given by  $|F| = \sum_{i=1}^m |T_i|$ , where  $|T_i|$  is the number of nodes occurring in  $T_i$ . The  
 107 class of decision trees on  $X_n$  is denoted  $\text{DT}_n$ , and the class of random forests with at most  $m$  decision  
 108 trees (with  $m \geq 1$ ) over  $\text{DT}_n$  is denoted  $\text{RF}_{n,m}$ .  $\text{RF}_n$  is the union of all  $\text{RF}_{n,m}$  for  $m \in \mathbb{N}$ .

109 **Example 1.** The random forest  $F = \{T_1, T_2, T_3\}$  in Figure 1 is composed of three decision trees.  
 110 It separates *Cattleya* orchids from other orchids using the following features:  $x_1$ : “has fragrant  
 111 flowers”,  $x_2$ : “has one or two leaves”,  $x_3$ : “has large flowers”, and  $x_4$ : “is sympodial”.

112 It is well-known that any decision tree  $T$  can be transformed into its negation  $\neg T \in \text{DT}_n$ , by simply  
 113 reverting the label of leaves. Negating a random forest can also be achieved in polynomial time:

114 **Proposition 1.** There exists a linear-time algorithm that computes a random forest  $\neg F \in \text{RF}_{n,m}$   
 115 equivalent to the negation of a given random forest  $F \in \text{RF}_{n,m}$ .

116 Another important property of decision trees is that any  $T \in \text{DT}_n$  can be transformed in linear time  
 117 into an equivalent disjunction of terms  $\text{DNF}(T)$ , where each term coincides with a 1-path (i.e., a path  
 118 from the root to a leaf labeled with 1), or a conjunction of clauses  $\text{CNF}(T)$ , where each clause is the  
 119 negation of term describing a 0-path. When switching to random forests, the picture is quite different:

120 **Proposition 2.** Any CNF or DNF formula can be converted in linear time into an equivalent random  
 121 forest, but there is no polynomial-space translation from RF to CNF or to DNF.

### 122 3 Random Forest Explanations

123 The key focus of this study is to explain *why* a given (Boolean) random forest classifies some incoming  
 124 data instance as positive or negative. This calls for a notion of abductive explanation<sup>1</sup>. Formally,  
 125 given a Boolean function  $f \in \mathcal{F}_n$  and an instance  $\mathbf{x} \in \{0, 1\}^n$ , an *abductive explanation* for  $\mathbf{x}$   
 126 given  $f$  is an implicant  $t$  of  $f$  (resp.  $\neg f$ ) if  $f(\mathbf{x}) = 1$  (resp.  $f(\mathbf{x}) = 0$ ) that covers  $\mathbf{x}$ . An abductive  
 127 explanation  $t$  for  $\mathbf{x}$  given  $f$  always exists, since  $t = t_{\mathbf{x}}$  is such a (trivial) explanation. So, in the rest  
 128 of this section, we shall mainly concentrate on *sparse* forms of abductive explanations.

129 Before delving into details, it is worth mentioning that if  $f$  is represented by a random forest then,  
 130 without loss of generality, we can focus on the case where  $\mathbf{x}$  is a positive example of  $f$ , because  $\neg f$   
 131 can be computed in linear time (by Proposition 1). Nevertheless, for the sake of clarity, we shall  
 132 consider both cases  $f(\mathbf{x}) = 1$  and  $f(\mathbf{x}) = 0$  in our definitions.

#### 133 3.1 Direct Reasons

134 For a decision tree  $T \in \text{DT}_n$  and a data instance  $\mathbf{x} \in \{0, 1\}^n$ , the *direct reason* of  $\mathbf{x}$  given  $T$  is the  
 135 term  $t_{\mathbf{x}}^T$  corresponding to the unique root-to-leaf path of  $T$  that covers  $\mathbf{x}$ . We can extend this simple  
 136 form of abductive explanation to random forests as follows:

137 **Definition 1.** Let  $F = \{T_1, \dots, T_m\}$  be a random forest in  $\text{RF}_{n,m}$ , and  $\mathbf{x} \in \{0, 1\}^n$  be an instance.  
 138 Then, the direct reason for  $\mathbf{x}$  given  $F$  is the term  $t_{\mathbf{x}}^F$  defined by

$$t_{\mathbf{x}}^F = \begin{cases} \bigwedge_{T_i \in F: T_i(\mathbf{x})=1} t_{\mathbf{x}}^{T_i} & \text{if } F(\mathbf{x}) = 1 \\ \bigwedge_{T_i \in F: T_i(\mathbf{x})=0} t_{\mathbf{x}}^{T_i} & \text{if } F(\mathbf{x}) = 0 \end{cases}$$

139 By construction,  $t_{\mathbf{x}}^F$  is an abductive explanation which can be computed in  $\mathcal{O}(|F|)$  time.

140 **Example 2.** Considering Example 1 again, the instance  $\mathbf{x} = (1, 1, 1, 1)$  is recognized as a Cattleya  
 141 orchid, since  $F(\mathbf{x}) = 1$ . The direct reason for  $\mathbf{x}$  given  $F$  is  $t_{\mathbf{x}}^F = x_1 \wedge x_2 \wedge x_3 \wedge x_4$ . It coincides  
 142 with  $t_{\mathbf{x}}$ . Consider now the instance  $\mathbf{x}' = (0, 1, 0, 0)$ ; it is not recognized as a Cattleya orchid, since  
 143  $F(\mathbf{x}') = 0$ . The direct reason for  $\mathbf{x}'$  given  $F$  is  $t_{\mathbf{x}'}^F = x_2 \wedge \bar{x}_3 \wedge \bar{x}_4$ . It is a better abductive explanation  
 144 than  $t_{\mathbf{x}'}$  itself since it does not contain  $\bar{x}_1$ , which is locally irrelevant.

#### 145 3.2 Sufficient Reasons

146 Another valuable notion of abductive explanation is the one of *sufficient reason*<sup>2</sup>, defined for any  
 147 Boolean classifier [13]. In the setting of random forests, such explanations can be defined as follows:

148 **Definition 2.** Let  $F \in \text{RF}_n$  be a random forest and  $\mathbf{x} \in \{0, 1\}^n$  be an instance. A sufficient reason  
 149 for  $\mathbf{x}$  given  $F$  is a prime implicant  $t$  of  $F$  (resp.  $\neg F$ ) if  $F(\mathbf{x}) = 1$  (resp.  $F(\mathbf{x}) = 0$ ) that covers  $\mathbf{x}$ .

150 **Example 3.** For our running example,  $x_1 \wedge x_2 \wedge x_4$  and  $x_3 \wedge x_4$  are the sufficient reasons for  $\mathbf{x}$   
 151 given  $F$ .  $\bar{x}_4$  and  $\bar{x}_1 \wedge x_2 \wedge \bar{x}_3$  are the sufficient reasons for  $\mathbf{x}'$  given  $F$ .

152 Unlike arbitrary abductive explanations, all features occurring in a sufficient reason  $t$  are *relevant*.  
 153 Indeed, removing any literal from  $t$  would question the fact that  $t$  implies  $F$ . To this point, the direct  
 154 reason  $t_{\mathbf{x}}^F$  for  $\mathbf{x}$  given  $F$  may contain arbitrarily many more features than a sufficient reason for  $\mathbf{x}$   
 155 given  $F$ , since this was already shown in the case where  $F$  consists in a single decision tree [17].

156 The problem of finding a sufficient reason  $t$  for an input instance  $\mathbf{x} \in \{0, 1\}^n$  with respect to a given  
 157 random forest  $F \in \text{RF}_n$ , has recently been shown DP-complete [18]. In fact, even the apparently  
 158 simple task of *checking* whether  $t$  is an implicant of  $F$  is already hard:

159 **Proposition 3.** Let  $F$  be a random forest in  $\text{RF}_n$  and  $t$  be a term over  $X_n$ . Then, deciding whether  $t$   
 160 is an implicant of  $F$  is coNP-complete.

161 The above result is in stark contrast with the computational complexity of checking whether a term  $t$   
 162 is an implicant of a decision tree  $T$ . This task can be solved in polynomial time, using the fact that

<sup>1</sup>Unlike [15], we do not require those explanations to be minimal w.r.t. set inclusion, in order to keep the concept distinct (and actually more general) than the one of sufficient reasons.

<sup>2</sup>Sufficient reasons are also known as prime-implicant explanations [29].

163  $T$  can be converted (in linear time) into its clausal form  $\text{CNF}(T)$ , together with the fact that testing  
 164 whether  $t$  implies  $\text{CNF}(T)$  can be done in  $\mathcal{O}(|T|)$  time. That mentioned, in the case of random forests,  
 165 the implicant test can be achieved via a call to a SAT oracle:

166 **Proposition 4.** *Let  $F = \{T_1, \dots, T_m\}$  be a random forest of  $\text{RF}_{n,m}$ , and  $t$  be a (satisfiable) term  
 167 over  $X_n$ . Let  $H$  be the CNF formula*

$$\{(\bar{y}_i \vee c) : i \in [m], c \in \text{CNF}(\neg T_i)\} \cup \text{CNF}\left(\sum_{i=1}^m y_i > \frac{m}{2}\right)$$

168 where  $\{y_1, \dots, y_m\}$  are fresh variables and  $\text{CNF}\left(\sum_{i=1}^m y_i > \frac{m}{2}\right)$  is a CNF encoding of the cardinality  
 169 constraint  $\sum_{i=1}^m y_i > \frac{m}{2}$ . Then,  $t$  is an implicant of  $F$  if and only if  $H \wedge t$  is unsatisfiable.

170 Based on such an encoding, the sufficient reasons for an instance  $\mathbf{x}$  given a random forest  $F$  can  
 171 be characterized in terms of MUS (minimal unsatisfiable subsets), as suggested in [18]. This  
 172 characterization is useful because many SAT-based algorithms for computing a MUS (or even all  
 173 MUSes) of a CNF formula have been pointed out for the past decade [2, 19, 20], and hence, one can  
 174 take advantage of them for computing sufficient reasons.

175 Going one step further, a natural way for improving the clarity of sufficient reasons is to focus on  
 176 those of minimal size. Specifically, given  $F \in \text{RF}_n$  and  $\mathbf{x} \in \{0, 1\}^n$ , a *minimal sufficient reason* for  
 177  $\mathbf{x}$  with respect to  $F$  is a sufficient reason for  $\mathbf{x}$  given  $F$  of minimal size.<sup>3</sup>

178 **Example 4.** *For our running example,  $x_3 \wedge x_4$  is the unique minimal sufficient reason for  $\mathbf{x}$  given  $F$ ,  
 179 and  $\bar{x}_4$  is the unique minimal reason for  $\mathbf{x}'$  given  $F$ .*

180 As a by-product of the characterization of a sufficient reason in terms of MUS [18], a minimal  
 181 sufficient reason for  $\mathbf{x}$  given  $f$  can be viewed as a *minimal MUS*. Thus, we can exploit algorithms for  
 182 computing minimal MUSes (see e.g., [16]) in order to derive minimal sufficient reasons. However,  
 183 deriving a minimal sufficient reason is computationally harder than deriving a sufficient reason:

184 **Proposition 5.** *Let  $F \in \text{RF}_n$ ,  $\mathbf{x} \in \{0, 1\}^n$ , and  $k \in \mathbb{N}$ . Then, deciding whether there exists a  
 185 minimal sufficient reason  $t$  for  $\mathbf{x}$  given  $F$  containing at most  $k$  features is  $\Sigma_2^P$ -complete.*

### 186 3.3 Majoritary Reasons

187 Based on the above considerations, a natural question arises: does there exist a middle ground  
 188 between direct reasons, which main contain many irrelevant features but are easy to calculate, and  
 189 sufficient reasons, which only contain relevant features but are potentially much harder to generate?  
 190 Inspired by the way prime implicants can be computed when dealing with decision trees, we can  
 191 reply in the affirmative using the notion of *majoritary reasons*, defined as follows.

192 **Definition 3.** *Let  $F = \{T_1, \dots, T_m\}$  be a random forest in  $\text{RF}_{n,m}$  and  $\mathbf{x} \in \{0, 1\}^n$  be an instance.  
 193 Then, a *majoritary reason* for  $\mathbf{x}$  given  $F$  is a term  $t$  covering  $\mathbf{x}$ , such that  $t$  is an implicant of at least  
 194  $\lfloor \frac{m}{2} \rfloor + 1$  decision trees  $T_i$  (resp.  $\neg T_i$ ) if  $F(\mathbf{x}) = 1$  (resp.  $F(\mathbf{x}) = 0$ ), and for every  $l \in t$ ,  $t \setminus \{l\}$   
 195 does not satisfy this last condition.*

196 **Example 5.** *For our running example,  $\mathbf{x}$  has three majoritary reasons given  $F$ :  $x_1 \wedge x_2 \wedge x_4$ ,  
 197  $x_1 \wedge x_3 \wedge x_4$ , and  $x_2 \wedge x_3 \wedge x_4$ . Those reasons are better than  $t_{\mathbf{x}}^F$  in the sense that they are shorter  
 198 than this direct reason. Contrastingly,  $\mathbf{x}'$  has four majoritary reasons given  $F$ :  $\bar{x}_1 \wedge \bar{x}_4$ ,  $x_2 \wedge \bar{x}_4$ ,  
 199  $\bar{x}_3 \wedge \bar{x}_4$ , and  $\bar{x}_1 \wedge x_2 \wedge \bar{x}_3$ . Each of the two majoritary reasons  $x_2 \wedge \bar{x}_4$ ,  $\bar{x}_3 \wedge \bar{x}_4$  show that  $t_{\mathbf{x}'}^F$   
 200 contains some irrelevant literals for the task of classifying  $\mathbf{x}'$  using  $F$ .*

201 In general, the notions of majoritary reasons and of sufficient reasons do not coincide. Indeed, a  
 202 sufficient reason  $t$  is a prime implicant (covering  $\mathbf{x}$ ) of the forest  $F$ , while a majoritary reason  $t'$  is an  
 203 implicant (covering  $\mathbf{x}$ ) of a strict majority of decision trees in the forest  $F$  satisfying the additional  
 204 condition that  $t'$  is a prime implicant of at least one of these decision trees. Viewing majoritary  
 205 reasons as “weak” forms of sufficient reasons, they can include irrelevant features:

206 **Proposition 6.** *Let  $F = \{T_1, \dots, T_m\}$  be a random forest of  $\text{RF}_{n,m}$  and  $\mathbf{x} \in \{0, 1\}^n$  such that  
 207  $F(\mathbf{x}) = 1$ . Unless  $m < 3$ , it can be the case that every majoritary reason for  $\mathbf{x}$  given  $F$  contains  
 208 arbitrarily many more features than any sufficient reason for  $\mathbf{x}$  given  $F$ .*

<sup>3</sup>Minimal sufficient reasons should not to be confused with *minimum-cardinality explanations* [29], where the minimality condition bears on the features set to 1 in the data instance  $\mathbf{x}$ .

209 What makes majoritary reasons valuable is that they are abductive and can be generated in linear time.  
 210 The evidence that any majoritary reason  $t$  for  $\mathbf{x}$  given  $F$  is an abductive explanation for  $\mathbf{x}$  given  $F$   
 211 comes directly from the fact that if  $t$  implies a majority of decision trees in  $F$ , then it is an implicant  
 212 of  $F$  (note that the converse implication does not hold in general).

213 The tractability of generating majoritary reasons lies in the fact that they can be found using a simple  
 214 greedy algorithm. For the case where  $F(\mathbf{x}) = 1$ , start with  $t = t_{\mathbf{x}}$ , and iterate over the literals  $l$  of  $t$   
 215 by checking whether  $t$  deprived of  $l$  is an implicant of at least  $\lfloor \frac{m}{2} \rfloor + 1$  decision trees of  $F$ . If so,  
 216 remove  $l$  from  $t$  and proceed to the next literal. Once all literals in  $t_{\mathbf{x}}$  have been examined, the final  
 217 term  $t$  is by construction an implicant of a strict majority of decision trees in  $F$ , such that removing  
 218 any literal from it would lead to a term that is no longer an implicant of this majority. So,  $t$  is by  
 219 construction a majoritary reason. The case where  $F(\mathbf{x}) = 0$  is similar, by simply replacing each  
 220  $T_i$  with its negation in  $F$ . This greedy algorithm runs in  $\mathcal{O}(n|F|)$  time, using the fact that, on each  
 221 iteration, checking whether  $t$  is an implicant of  $T_i$  (for each  $i \in [m]$ ) can be done in  $\mathcal{O}(|T_i|)$  time.

222 By analogy with minimal sufficient reasons, a natural way of improving the quality of majoritary  
 223 reasons is to seek for shortest ones. Let  $F \in \text{RF}_n$  be a random forest and  $\mathbf{x} \in \{0, 1\}^n$  be an instance.  
 224 Then, a *minimal majoritary reason* for  $\mathbf{x}$  given  $F$  is a minimal-size majoritary reason for  $\mathbf{x}$  given  $F$ .

225 **Example 6.** For our running example, the three majoritary reasons for  $\mathbf{x}$  given  $F$  are its minimal  
 226 majoritary reasons. Contrastingly, among the majoritary reasons for  $\mathbf{x}'$  given  $F$ , only  $\bar{x}_1 \wedge \bar{x}_4$ ,  
 227  $x_2 \wedge \bar{x}_4$ , and  $\bar{x}_3 \wedge \bar{x}_4$  are minimal majoritary reasons.

228 Unsurprisingly, the optimization task for majoritary reasons is more demanding than the generation  
 229 task. Yet, minimal majoritary reasons are easier to find than minimal sufficient reasons. Specifically:

230 **Proposition 7.** Let  $F \in \text{RF}_n$ ,  $\mathbf{x} \in \{0, 1\}^n$ , and  $k \in \mathbb{N}$ . Then, deciding whether there exists a  
 231 minimal majoritary reason  $t$  for  $\mathbf{x}$  given  $F$  containing at most  $k$  features is NP-complete.

232 A common approach for handling NP-optimization problems is to rely on modern constraint solvers.  
 233 From this perspective, recall that a PARTIAL MAXSAT problem consists of a pair  $(C_{\text{soft}}, C_{\text{hard}})$   
 234 where  $C_{\text{soft}}$  and  $C_{\text{hard}}$  are (finite) sets of clauses. The goal is to find a Boolean assignment that  
 235 maximizes the number of clauses  $c$  in  $C_{\text{soft}}$  that are satisfied, while satisfying all clauses in  $C_{\text{hard}}$ .

236 **Proposition 8.** Let  $F \in \text{RF}_{n,m}$  and  $\mathbf{x} \in \{0, 1\}^n$  be an instance such that  $F(\mathbf{x}) = 1$ . Let  
 237  $(C_{\text{soft}}, C_{\text{hard}})$  be an instance of the PARTIAL MAXSAT problem such that:

$$C_{\text{soft}} = \{\bar{x}_i : x_i \in t_{\mathbf{x}}\} \cup \{x_i : \bar{x}_i \in t_{\mathbf{x}}\}$$

$$C_{\text{hard}} = \{(\bar{y}_i \vee c_{|\mathbf{x}}) : i \in [m], c \in \text{CNF}(T_i)\} \cup \text{CNF}\left(\sum_{i=1}^m y_i > \frac{m}{2}\right)$$

238 where  $c_{|\mathbf{x}} = c \cap t_{\mathbf{x}}$  is the restriction of  $c$  to the literals in  $t_{\mathbf{x}}$ ,  $\{y_1, \dots, y_m\}$  are fresh variables and  
 239  $\text{CNF}(\sum_{i=1}^m y_i > \frac{m}{2})$  is a CNF encoding of the constraint  $\sum_{i=1}^m y_i > \frac{m}{2}$ . The intersection of  $t_{\mathbf{x}}$  with  
 240  $t_{\mathbf{z}^*}$ , where  $\mathbf{z}^*$  is an optimal solution of  $(C_{\text{soft}}, C_{\text{hard}})$ , is a minimal majoritary reason for  $\mathbf{x}$  given  $F$ .

241 Clearly, in the case where  $F(\mathbf{x}) = 0$ , it is enough to consider the same instance of PARTIAL MAXSAT  
 242 as above, except that  $C_{\text{hard}} = \{(\bar{y}_i \vee c_{|\mathbf{x}}) : i \in [m], c \in \text{CNF}(\neg T_i)\} \cup \text{CNF}(\sum_{i=1}^m y_i > \frac{m}{2})$ .

243 Thanks to this characterization result, one can leverage the numerous algorithms that have been  
 244 developed so far for PARTIAL MAXSAT (see e.g. [1, 23, 24, 28]) in order to compute minimal  
 245 majoritary reasons. We took advantage of it to achieve some of the experiments reported in Section 4.

## 246 4 Experiments

247 **Empirical setting.** The empirical protocol was as follows. We have considered 15 datasets, which  
 248 are standard benchmarks from the well-known repositories Kaggle (www.kaggle.com), OpenML  
 249 (www.openml.org), and UCI (archive.ics.uci.edu/ml/). These datasets are *compas*, *placement*,  
 250 *recidivism*, *adult*, *ad\_data*, *mnist38*, *mnist49*, *gisette*, *dexter*, *dorothea*, *farm-ads*, *higgs\_boson*,  
 251 *christine*, *gina*, and *bank*. *mnist38* and *mnist49* are subsets of the *mnist* dataset, restricted to the  
 252 instances of 3 and 8 (resp. 4 and 9) digits. Due to space constraints, additional information about  
 253 the datasets (especially the numbers and types of features, the number of instances), and about the  
 254 random forests that have been trained (especially, the number of Boolean features used, the number

255 of trees, the depth of the trees, the mean accuracy) are reported as a supplementary material. We  
256 used only datasets for binary classification, which is a very common kind of dataset. Categorical  
257 features have been treated as arbitrary numbers (the scale is nominal). As to numeric features, no data  
258 preprocessing has taken place: these features have been binarized on-the-fly by the random forest  
259 learning algorithm that has been used.

260 For every benchmark  $b$ , a 10-fold cross validation process has been achieved. Namely, a set of 10  
261 random forest  $F_b$  have been computed and evaluated from the labelled instances of  $b$ , partitioned  
262 into 10 parts. One part was used as the test set and the remaining 9 parts as the training set for  
263 generating a random forest. The classification performance for  $b$  was measured as the mean accuracy  
264 obtained over the 10 random forests generated from  $b$ . As to the random forest learner, we have used  
265 the implementation provided by the Scikit-Learn [26] library in his version 0.23.2. The maximal  
266 depth of any decision tree in a forest has been bounded at 8. All other hyper-parameters of the  
267 learning algorithm have been set to their default value except the number of trees. We made some  
268 preliminary tests for tuning this parameter in order to ensure that the accuracy is good enough. For  
269 each benchmark  $b$ , each random forest  $F$ , and a subset of 25 instances  $x$  picked up at random in the  
270 corresponding test set (leading to 250 instances per dataset) we have run the algorithms described in  
271 Section 3 for deriving the direct reason for  $x$  given  $F$ , a sufficient reason for  $x$  given  $F$ , a majoritary  
272 reason  $x$  given  $F$ , a minimal majoritary reason for  $x$  given  $F$ , and a minimal sufficient reason for  $x$   
273 given  $F$ .

274 For computing sufficient reasons and minimal majoritary reasons, we took advantage of the Pysat  
275 library [14] (version 0.1.6.dev15) which provides the implementation of the RC2 PARTIAL MAXSAT  
276 solver and an interface to MUSER [4]. When deriving majoritary reasons, we picked up uniformly at  
277 random 50 permutations of the literals describing the instance and tried to eliminate those literals  
278 (within the greedy algorithm) following the ordering corresponding to the permutation. As a majori-  
279 tary reason for the instance, we kept a smallest reason among those that have been derived (of course,  
280 the corresponding computation time that has been measured is the cumulated time over the 50 tries).  
281 Sufficient reasons have been computed as MUSes, as explained before.

282 We also derived a “LIME explanation” for each instance. Such an explanation has been generated  
283 thanks to the following approach. For any  $x$  under consideration, one first used LIME [27] to generate  
284 an associated linear model  $w_x$  where  $w_x \in \mathbb{R}^n$ . This linear model  $w_x$  classifies any instance  $x'$  as a  
285 positive instance if and only if  $w_x \cdot x' > 0$ . Furthermore,  $w_x$  classifies the instance to be explained  $x$   
286 in the same way as the black box model considered at start (in our case, the random forest  $F$ ). We ran  
287 the LIME implementation linked to [27] in its latest version. Interestingly, a minimal sufficient reason  
288  $t$  for  $x$  given  $w_x$  can be generated in polynomial time from  $w_x$ . We call it a LIME explanation for  $x$ .  
289 The computation of  $t$  is as follows. If  $x$  is classified positively by  $w_x$ , in order to derive  $t$ , it is enough  
290 to sum in a decreasing way the positive weights  $w_i$  occurring in  $w_x$  until this sum exceeds the sum  
291 of the opposites of all the negative weights occurring in  $w_x$ . The term  $t$  composed of the variables  $x_i$   
292 corresponding to the positive weights that have been selected is by construction a minimal sufficient  
293 reason for  $x$  given  $w_x$  since for every  $x'$  covered by  $t$ , the inequation  $w_x \cdot x' > 0$  necessarily holds;  
294 indeed, it holds in the worst situation where all the variables associated with a positive weight in  $w_x$   
295 and not belonging to  $t$  are set to 0, whilst all the variables associated with a negative weight in  $w_x$   
296 are set to 1. Similarly, if  $x$  is classified negatively by  $w_x$ , in order to derive  $t$ , it is enough to sum in  
297 an increasing way the negative weights  $w_i$  occurring in  $w_x$  until this sum is lower than or equal to  
298 the opposite of the sum of all the positive weights occurring in  $w_x$ . This time, the term  $t$  composed  
299 of the variables  $x_i$  corresponding to the negative weights that have been selected is by construction a  
300 minimal sufficient reason for  $x$  given  $w_x$ .

301 All the experiments have been conducted on a computer equipped with Intel(R) XEON E5-2637 CPU  
302 @ 3.5 GHz and 128 Gib of memory. A time-out (TO) of 600s has been considered for each instance  
303 and each type of explanation, except LIME explanations.

304 **Results.** A first conclusion that can be drawn from our experiments is the intractability of computing  
305 in practice minimal sufficient reasons (this is not surprising, since this coheres with the complexity  
306 result given by Proposition 5). Indeed, we have been able to compute within the time limit of 600s a  
307 minimal reason for only 10 instances and a single dataset (*compas*).

308 Due to space limitations, we report hereafter empirical results about two datasets only, namely  
309 *placement* and *gissette* (the results obtained on the other datasets are similar and available as a

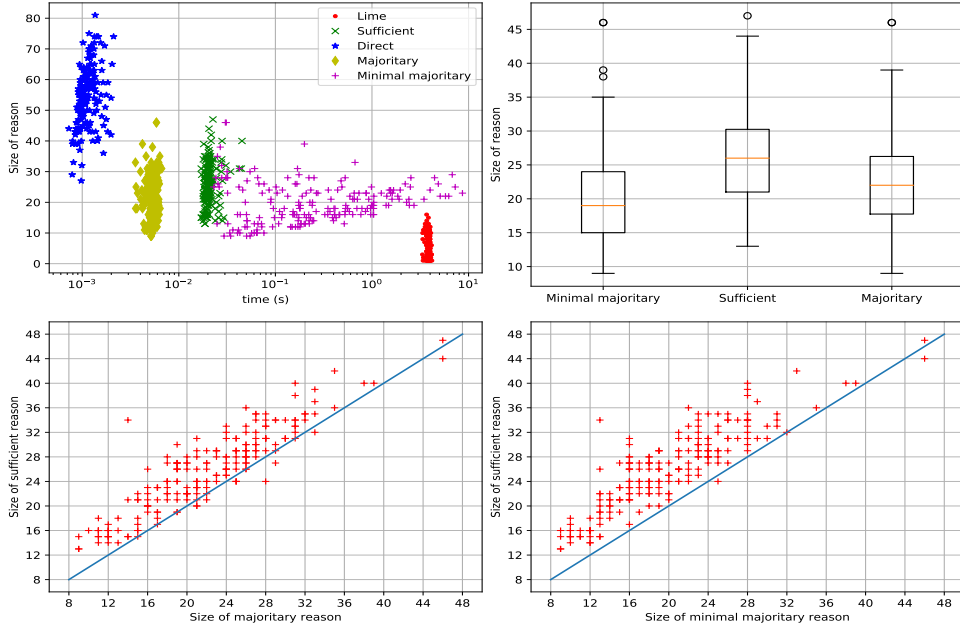


Figure 2: Empirical results for the *placement* dataset.

310 supplementary material). The *placement* data set is about the placement of students in a campus. It  
 311 consists of 215 labelled instances. Students are described using 13 features, related to their curricula,  
 312 the type and work experience and the salary. An instance is labelled as positive when the student  
 313 gets a job. The random forest that has been generated consists of 25 trees, and its mean accuracy  
 314 is 97.6%. *gisette* is a much larger dataset, based on 5000 features and containing 7000 labelled  
 315 instances. Features correspond to pixels. The problem is to separate the highly confusable digits 4  
 316 and 9. An instance is labelled as positive whenever the picture represents a 9. The random forest that  
 317 has been generated consists of 85 trees, and its mean accuracy is 96%.

318 Figure 2 provides the results obtained for *placement*, using four plots. Each dot represents an instance.  
 319 The first plot shows the time needed to compute a reason on the x-axis, and the size of this reason on  
 320 the y-axis. On this plot, no dot corresponds to a minimal sufficient reason because their computation  
 321 did not terminate before the time-out. The plot also highlights that all the other reasons have been  
 322 computed within the time limit, and in general using a small amount of time. In particular, it shows  
 323 that the direct reason can be quite large, that the computation of LIME explanations is usually more  
 324 expensive than the ones of the other explanations, and that LIME explanations can be very short (but  
 325 one must keep in mind that they are not abductive explanations in general<sup>4</sup>). A box plot about the  
 326 sizes of all the explanations is reported (the LIME ones and the direct reasons are not presented for  
 327 the sake of readability). The figure also provides two scatter plots, aiming to compare the size of  
 328 majoritary reasons with the size of sufficient reasons, as well as the size of the minimal majoritary  
 329 reasons with the size of sufficient reasons. These plots clearly show the benefits that can be offered  
 330 by considering majoritary reasons and minimal majoritary reasons instead of sufficient reasons.

331 Figure 3 synthesizes the results obtained for *gisette*, using four plots again. Three of them are of the  
 332 same kind as the plots used for *placement*. Conclusions similar to those drawn for *placement* can  
 333 be derived for *gisette*, with some exceptions. First of all, this time, no dot corresponds to a minimal  
 334 majoritary reason because their computation did not terminate before the time-out. Furthermore,  
 335 LIME explanations are very long here. This can be explained by the fact that the computation  
 336 achieved by LIME relies on a binary representation of the instance that is quite different (and possibly  
 337 much larger) than the one considered in the representation of the random forest. Indeed, each decision  
 338 tree of the forest focuses only on a subset of most important features (in the sense of Gini criterion)  
 339 found during the learning phase. In our experiments, the size of LIME explanations was typically  
 340 high for datasets based on many features.

<sup>4</sup>See also [25] that reports some experiments about ANCHOR (the successor of LIME), assessing the quality of the explanations computed using ANCHOR.



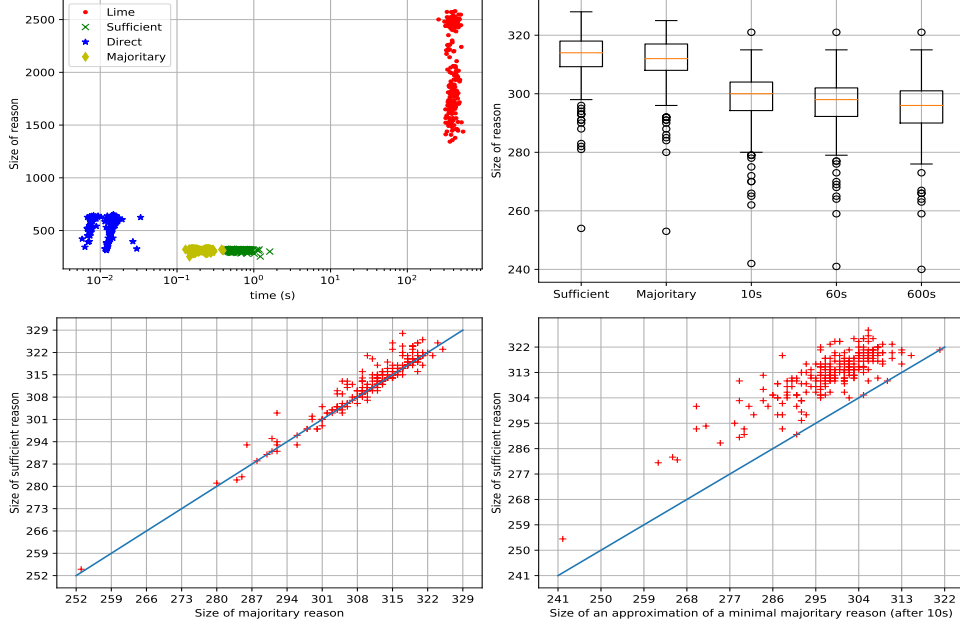


Figure 3: Empirical results for the *gisette* dataset.

341 When minimal majoritary reasons are hard to be computed (as it is the case for *gisette*), an approach  
342 consists in approximating them. Interestingly, one can take advantage of an incremental PARTIAL  
343 MAXSAT ALGORITHM, like LMHS [28], to do the job. Specifically, the result given in Proposition  
344 8 provides a way to derive abductive explanations for an instance  $x$  given a random forest  $F$  in an  
345 *anytime* fashion. Basically, using LMHS, a Boolean assignment  $z$  satisfying all the hard constraints of  
346  $C_{\text{hard}}$  and a given number, say  $k$ , of soft constraints from  $C_{\text{soft}}$  is looked for ( $k$  is set to 0 at start).  
347 If such an assignment is found, then one looks for an assignment satisfying  $k + 1$  soft constraint,  
348 and so on, until an optimal solution is found or a preset time bound is reached. In many cases, the  
349 most demanding step from a computational standpoint is the one for which  $k$  is the optimal value  
350 (but one ignores it) and one looks for an assignment that satisfies  $k + 1$  soft constraint (and such an  
351 assignment does not exist). By construction, every  $z$  that is generated that way is such that  $t_x \cap t_z$   
352 is an implicant of  $F$  that covers  $x$  (and hence, an abductive explanation). The approximation  $z$  of  
353 a minimal majoritary reason for  $x$  given  $F$ , which is obtained when the time limit is met, can be  
354 significantly shorter than the sufficient reason for  $x$  given  $F$  that has been derived. In our experiments,  
355 we used three time limits: 10s, 60s, 600s. As the box plot and the dedicated scatter plot given in  
356 Figure 3 show it, the sizes of the approximations  $z$  which are derived gently decrease with time.  
357 Interestingly, the size savings that are achieved in comparison to sufficient reasons are significant,  
358 even for the smallest time bound of 10s that has been considered.

## 359 5 Conclusion

360 In this paper, we have introduced, analyzed and evaluated some new notions of abductive explanations  
361 suited to random forest classifiers, namely majoritary reasons and minimal majoritary reasons.  
362 Our investigation reveals the existence of a trade-off between runtime complexity and sparsity for  
363 abductive explanations. Unlike sufficient reasons, majoritary reasons and minimal majoritary reasons  
364 may contain irrelevant features. Despite this evidence, majoritary reasons and minimal majoritary  
365 reasons appear as valuable alternative to sufficient reasons. Indeed, majoritary reasons can be  
366 computed in polynomial time while sufficient reasons cannot (unless  $P = NP$ ). In addition, most of  
367 the time in our experiments, majoritary reasons appear as slightly smaller than sufficient reasons.  
368 Minimal majoritary reasons can be looked for when majoritary reasons are too large, but this is at  
369 the cost of an extra computation time that can be important, and even prohibitive in some cases.  
370 However, minimal majoritary reasons can be approximated using an *anytime* PARTIAL MAXSAT  
371 algorithm. Empirically, approximations can be derived within a small amount of time and their sizes  
372 are significantly smaller than the ones of sufficient reasons.

373 **References**

- 374 [1] C. Ansótegui, M. L. Bonet, and J. Levy. SAT-based MaxSAT algorithms. *Artificial Intelligence*,  
375 196:77–105, 2013.
- 376 [2] G. Audemard, J-M. Lagniez, and L. Simon. Improving glucose for incremental SAT solving  
377 with assumptions: Application to MUS extraction. In *Proceedings of the 16th International*  
378 *Conference on Theory and Applications of Satisfiability Testing (SAT'13)*, pages 309–317, 2013.
- 379 [3] A. T. Azar, H. I. Elshazly, A. E. Hassanien, and A. M. Elkorany. A random forest classifier for  
380 lymph diseases. *Computer Methods and Programs in Biomedicine*, 113(2):465–473, 2014.
- 381 [4] Anton Belov and João Marques-Silva. Muser2: An efficient MUS extractor. *J. Satisf. Boolean*  
382 *Model. Comput.*, 8(3/4):123–128, 2012.
- 383 [5] C. Bénard, G. Biau, S. Da Veiga, and E. Scornet. Interpretable random forests via rule extraction.  
384 In *Proceedings of the 24th International Conference on Artificial Intelligence and Statistics*,  
385 *AISTATS'21*, pages 937–945, 2021.
- 386 [6] G. Biau. Analysis of a random forests model. *Journal of Machine Learning Research*, 13:1063–  
387 1095, 2012.
- 388 [7] A. Bogomolov, B. Lepri, J. Staiano, N. Oliver, F. Pianesi, and A. Pentland. Once upon a crime:  
389 Towards crime prediction from demographics and mobile data. In *Proceedings of the 16th*  
390 *International Conference on Multimodal Interaction, ICMI'14*, pages 427–434. ACM, 2014.
- 391 [8] L. Breiman. Random forests. *Machine Learning*, 45(1):5–32, 2001.
- 392 [9] X. Chen and H. Ishwaran. Random forests for genomic data analysis. *Genomics*, 99(6):323–329,  
393 2012.
- 394 [10] A. Choi, A. Shih, A. Goyanka, and A. Darwiche. On symbolically encoding the behavior  
395 of random forests. In *Proceedings of the 3rd Workshop on Formal Methods for ML-Enabled*  
396 *Autonomous Systems (FoMLAS)*, 2020.
- 397 [11] A. Criminisi and J. Shotton. *Decision Forests for Computer Vision and Medical Image Analysis*.  
398 Advances in Computer Vision and Pattern Recognition. Springer, 2013.
- 399 [12] R. Cutler, C. E. Jr. Thomas, K. H. Beard, A. Cutler, K. T. Hess, J. Gibson, and J. J. Lawler.  
400 Random forests for classification in ecology. *Ecology*, 88(11):2783–2792, 2007.
- 401 [13] A. Darwiche and A. Hirth. On the reasons behind decisions. In *Proceedings of the 24th*  
402 *European Conference on Artificial Intelligence (ECAI'20)*, pages 712–720, 2020.
- 403 [14] A. Ignatiev, A. Morgado, and J. Marques-Silva. PySAT: A python toolkit for prototyping with  
404 SAT oracles. In *Proceedings of the 21st International Conference on Theory and Applications*  
405 *of Satisfiability Testing (SAT'2018)*, pages 428–437, 2018.
- 406 [15] A. Ignatiev, N. Narodytska, and J. Marques-Silva. Abduction-based explanations for ma-  
407 chine learning models. In *Proceedings of the 23rd AAAI Conference on Artificial Intelligence*  
408 *(AAAI'19)*, pages 1511–1519, 2019.
- 409 [16] A. Ignatiev, A. Previti, M. Liffiton, and J. Marques-Silva. Smallest MUS extraction with minimal  
410 hitting set dualization. In *Proceedings of the 21st International Conference on Principles and*  
411 *Practice of Constraint Programming (CP'15)*, pages 173–182, 2015.
- 412 [17] Y. Izza, A. Ignatiev, and J. Marques-Silva. On explaining decision trees. *CoRR*, abs/2010.11034,  
413 2020.
- 414 [18] Y. Izza and J. Marques-Silva. On explaining random forests with SAT. In *Proceedings of the*  
415 *30th International Joint Conference on Artificial Intelligence (IJCAI'21)*, page to appear, 2021.
- 416 [19] M. Liffiton, A. Previti, A. Malik, and J. Marques-Silva. Fast, flexible MUS enumeration.  
417 *Constraints An Int. J.*, 21(2):223–250, 2016.
- 418 [20] J. Marques-Silva, M. Janota, and C. Mencía. Minimal sets on propositional formulae. Problems  
419 and reductions. *Artificial Intelligence*, 252:22–50, 2017.
- 420 [21] T. Miller. Explanation in artificial intelligence: Insights from the social sciences. *Artificial*  
421 *Intelligence*, 267:1–38, 2019.
- 422 [22] Ch. Molnar. *Interpretable Machine Learning - A Guide for Making Black Box Models Explain-*  
423 *able*. Leanpub, 2019.

- 424 [23] A. Morgado, A. Ignatiev, and J. Marques-Silva. MSCG: robust core-guided MaxSAT solving. *J.*  
425 *Satisf. Boolean Model. Comput.*, 9(1):129–134, 2014.
- 426 [24] N. Narodytska and F. Bacchus. Maximum satisfiability using core-guided MaxSAT resolution.  
427 In *Proceedings of the 28th AAAI Conference on Artificial Intelligence*, pages 2717–2723, 2014.
- 428 [25] N. Narodytska, A. Shrotri, K. Meel, A. Ignatiev, and J. Marques-Silva. Assessing heuristic  
429 machine learning explanations with model counting. In *Proceedings of 22nd International*  
430 *Conference on the Theory and Applications of Satisfiability Testing (SAT'19)*, pages 267–278,  
431 2019.
- 432 [26] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel,  
433 P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher,  
434 M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine*  
435 *Learning Research*, 12:2825–2830, 2011.
- 436 [27] M. T. Ribeiro, S. Singh, and C. Guestrin. "why should I trust you?": Explaining the predictions  
437 of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on*  
438 *Knowledge Discovery and Data Mining*, pages 1135–1144. ACM, 2016.
- 439 [28] P. Saikko, J. Berg, and M. Järvisalo. LMHS: A SAT-IP hybrid MaxSAT solver. In *Proceedings of*  
440 *the 19th International Conference of Theory and Applications of Satisfiability Testing (SAT'16)*,  
441 pages 539–546, 2016.
- 442 [29] A. Shih, A. Choi, and A. Darwiche. A symbolic approach to explaining bayesian network  
443 classifiers. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial*  
444 *Intelligence (IJCAI'18)*, pages 5103–5111, 2018.

445 **Checklist**

- 446 1. For all authors...
- 447 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
448 contributions and scope? [Yes]
- 449 (b) Did you describe the limitations of your work? [Yes]
- 450 (c) Did you discuss any potential negative societal impacts of your work? [No] One cannot  
451 expect any negative impact (the paper is about explaining predictions).
- 452 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
453 them? [Yes]
- 454 2. If you are including theoretical results...
- 455 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 456 (b) Did you include complete proofs of all theoretical results? [Yes] As a supplementary  
457 material.
- 458 3. If you ran experiments...
- 459 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
460 mental results (either in the supplemental material or as a URL)? [Yes]
- 461 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
462 were chosen)? [Yes]
- 463 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
464 ments multiple times)? [No] But the results we obtained have been averaged over a  
465 number of trials.
- 466 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
467 of GPUs, internal cluster, or cloud provider)? [Yes]
- 468 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 469 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 470 (b) Did you mention the license of the assets? [Yes]
- 471 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]  
472 The pieces of software we used are furnished as a supplementary material.
- 473 (d) Did you discuss whether and how consent was obtained from people whose data you’re  
474 using/curating? [No] This issue is irrelevant for this paper.
- 475 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
476 information or offensive content? [No] The datasets we used are anonymized and do  
477 not contain personally identifiable information or offensive content.
- 478 5. If you used crowdsourcing or conducted research with human subjects...
- 479 (a) Did you include the full text of instructions given to participants and screenshots, if  
480 applicable? [No] We did not use crowdsourcing or conducted research with human  
481 subjects.
- 482 (b) Did you describe any potential participant risks, with links to Institutional Review  
483 Board (IRB) approvals, if applicable? [No] We did not use crowdsourcing or conducted  
484 research with human subjects.
- 485 (c) Did you include the estimated hourly wage paid to participants and the total amount  
486 spent on participant compensation? [No] We did not use crowdsourcing or conducted  
487 research with human subjects.