

514 Proofs

515 Proof of Proposition 1

516 *Proof.* Let T be the complete binary tree of depth k , formed by $n = 2^k - 1$ internal nodes and 2^k
 517 leaves. We assume a breadth-first ordering of internal nodes, such that the root is labeled by x_1 , the
 518 nodes of depth 1 are labeled by x_2 and x_3 , and so on. Each internal node at depth $k - 1$ from the
 519 root of T has two children, one of it is a 0-leaf and the other one is a 1-leaf. For an arbitrary instance
 520 $\mathbf{x} \in \{0, 1\}^n$ and any complete subtree T' of T of depth d , let $s(\mathbf{x}, T')$ denote the set of sufficient
 521 reasons of \mathbf{x} given T' , and let $\sigma(\mathbf{x}, d) = |s(\mathbf{x}, T')|$ denote the number of those sufficient reasons.
 522 We show by induction on d that:

$$\sigma(\mathbf{x}, 1) = 1 \quad (1)$$

$$\sigma(\mathbf{x}, d + 1) = \sigma(\mathbf{x}, d)(\sigma(\mathbf{x}, d) + 1) \quad (2)$$

523 For the base case (1), any complete subtree T' of T of depth $d = 1$ has a single internal node, say x_i ,
 524 with two leaves labeled by 0 and 1, respectively. Therefore, the unique sufficient reason for \mathbf{x} given
 525 T' is either x_i or \bar{x}_i , and hence, $\sigma(\mathbf{x}, 1) = 1$. Now, consider any complete subtree T' of T of depth
 526 $d + 1$ rooted at a node x_i . Let $T'_l(x_i)$ and $T'_r(x_i)$ denote the subtrees of depth d , respectively rooted
 527 at the left child of x_i and the right child of x_i . Suppose without loss of generality that the unique path
 528 leading to $T'(\mathbf{x}) = 1$ includes the left child of x_i (i.e. $T'_l(\mathbf{x}) = 1$). By construction,

$$s(\mathbf{x}, T') = \{t_l \wedge t_r : t_l \in s(\mathbf{x}, T'_l), t_r \in s(\mathbf{x}, T'_r)\} \\ \cup \{l_i \wedge t_l : t_l \in s(\mathbf{x}, T'_l)\}$$

529 where $l_i = \bar{x}_i$ if $x_i = 0$ in \mathbf{x} , and $l_i = x_i$ otherwise. Since by induction hypothesis $s(\mathbf{x}, T'_l) =$
 530 $s(\mathbf{x}, T'_r) = \sigma(\mathbf{x}, d)$, it follows that $\sigma(\mathbf{x}, d + 1) = \sigma(\mathbf{x}, d)^2 + \sigma(\mathbf{x}, d)$. Finally, since the doubly
 531 exponential sequence² given by $a(1) = 1$ and $a(d + 1) = a(d)^2 + a(d)$ satisfies $a(d) = \lfloor c^{2^{d-1}} \rfloor$,
 532 where $c \sim 1.59791$, it follows that $\sigma(\mathbf{x}, k) \geq \lfloor (3/2)^{2^{k-1}} \rfloor$. Using $2^{k-1} = (n + 1)/2$, we get the
 533 desired result. \square

534 Proof of Proposition 2

535 *Proof.* One first need the following lemma that gives a recursive characterization of the set of
 536 sufficient reasons for an instance given a Boolean classifier:

537 **Lemma 1.** For any Boolean function $f \in \mathcal{F}_n$ and any instance $\mathbf{x} \in \{0, 1\}^n$, the following inductive
 538 characterization of $sr(\mathbf{x}, f)$ holds:

$$\begin{aligned} sr(\mathbf{x}, 1) &= \{1\} \\ sr(\mathbf{x}, 0) &= \{\} \\ 539 \quad sr(\mathbf{x}, f) &= sr(\mathbf{x}, (f \mid \ell) \wedge (f \mid \bar{\ell})) \cup \{\ell \wedge t_\ell : t_\ell \in sr(\mathbf{x}, f \mid \ell) \text{ s.t. } t_\ell \not\models f \mid \bar{\ell}\} \\ &\quad \text{where } Var(\ell) \subseteq Var(f) \text{ and } t_x \models \ell \end{aligned}$$

and

$$sr(\mathbf{x}, (f \mid \ell) \wedge (f \mid \bar{\ell})) = \max(\{t_\ell \wedge t_{\bar{\ell}} : t_\ell \in sr(\mathbf{x}, f \mid \ell), t_{\bar{\ell}} \in sr(\mathbf{x}, f \mid \bar{\ell})\}, \models).$$

540 *Proof.* Let us recall first the following inductive characterization of $pi(f)$, the set of prime implicants
 541 of $f \in \mathcal{F}_n$, based on the Shannon decomposition of f over any of its variables x (see e.g., [2]):

$$\begin{aligned} pi(1) &= \{1\} \\ pi(0) &= \{\} \\ 542 \quad pi(f) &= pi((f \mid \bar{x}) \wedge (f \mid x)) \\ &\quad \cup \{\bar{x} \wedge t_{\bar{x}} : t_{\bar{x}} \in pi(f \mid \bar{x}) \text{ s.t. } \nexists t \in pi((f \mid \bar{x}) \wedge (f \mid x)), t_{\bar{x}} \models t\} \\ &\quad \cup \{x \wedge t_x : t_x \in pi(f \mid x) \text{ s.t. } \nexists t \in pi((f \mid \bar{x}) \wedge (f \mid x)), t_x \models t\} \\ &\quad \text{where } x \in Var(f) \end{aligned}$$

and

$$pi((f \mid \bar{x}) \wedge (f \mid x)) = \max(\{t_{\bar{x}} \wedge t_x : t_{\bar{x}} \in pi(f \mid \bar{x}), t_x \in pi(f \mid x)\}, \models).$$

²See <https://oeis.org/A007018>.

543 For the base cases $sr(\mathbf{x}, 1) = \{1\}$ and $sr(\mathbf{x}, 0) = \{\}$, the result is obvious. For the general case,
 544 taking $x \in Var(\ell)$, we have:

$$\begin{aligned}
 sr(\mathbf{x}, f) &= \{t \in pi(f) : t_{\mathbf{x}} \models t\} \\
 &= \{t \in pi((f \mid \bar{x}) \wedge (f \mid x)) : t_{\mathbf{x}} \models t\} \\
 &\cup \{t \in \{\bar{x} \wedge t_{\bar{x}} : t_{\bar{x}} \in pi(f \mid \bar{x}) \text{ s.t. } \nexists t' \in pi((f \mid \bar{x}) \wedge (f \mid x)), t_{\bar{x}} \models t'\}, \text{ and } t_{\mathbf{x}} \models t\} \\
 545 &\cup \{t \in \{x \wedge t_x : t_x \in pi(f \mid x) \text{ s.t. } \nexists t' \in pi((f \mid \bar{x}) \wedge (f \mid x)), t_x \models t'\}, \text{ and } t_{\mathbf{x}} \models t\} \\
 &= sr(\mathbf{x}, (f \mid \bar{x}) \wedge (f \mid x)) \\
 &\cup \{t \in \{\bar{x} \wedge t_{\bar{x}} : t_{\bar{x}} \in pi(f \mid \bar{x}) \text{ s.t. } \nexists t' \in pi((f \mid \bar{x}) \wedge (f \mid x)), t_{\bar{x}} \models t'\}, \text{ and } t_{\mathbf{x}} \models t\} \\
 &\cup \{t \in \{x \wedge t_x : t_x \in pi(f \mid x) \text{ s.t. } \nexists t' \in pi((f \mid \bar{x}) \wedge (f \mid x)), t_x \models t'\}, \text{ and } t_{\mathbf{x}} \models t\}
 \end{aligned}$$

546 Now, since \mathbf{x} is an instance, whatever ℓ , it cannot be the case that $t_{\mathbf{x}} \models \ell$ and $t_{\mathbf{x}} \models \bar{\ell}$. Suppose that
 547 $\ell = x$ (the case $\ell = \bar{x}$ is similar). In this situation, no element of $\{\bar{x} \wedge t_{\bar{x}} : t_{\bar{x}} \in pi(f \mid \bar{x}) \text{ s.t. } \nexists t \in$
 548 $pi((f \mid \bar{x}) \wedge (f \mid x)), t_{\bar{x}} \models t\}$, and $t_{\mathbf{x}} \models t\}$ can belong to $sr(\mathbf{x}, f)$. As a consequence, we get that:

$$\begin{aligned}
 sr(\mathbf{x}, f) &= sr(\mathbf{x}, (f \mid \bar{x}) \wedge (f \mid x)) \\
 549 &\cup \{t \in \{\ell \wedge t_{\ell} : t_{\ell} \in pi(f \mid \ell) \text{ s.t. } \nexists t' \in pi((f \mid \bar{\ell}) \wedge (f \mid \ell)), t_{\ell} \models t'\}, \text{ and } t_{\mathbf{x}} \models t\} \\
 &\text{where } Var(\ell) \subseteq Var(f) \text{ and } t_{\mathbf{x}} \models \ell
 \end{aligned}$$

550 If $t = \ell \wedge t_{\ell}$ is such that $t_{\mathbf{x}} \models t$ holds, then we have $t_{\mathbf{x}} \models t_{\ell}$. Hence, we have:

$$\begin{aligned}
 sr(\mathbf{x}, f) &= sr(\mathbf{x}, (f \mid \bar{x}) \wedge (f \mid x)) \\
 551 &\cup \{\ell \wedge t_{\ell} : t_{\ell} \in sr(\mathbf{x}, f \mid \ell) \text{ s.t. } \nexists t' \in pi((f \mid \bar{\ell}) \wedge (f \mid \ell)), t_{\ell} \models t'\} \\
 &\text{where } Var(\ell) \subseteq Var(f) \text{ and } t_{\mathbf{x}} \models \ell
 \end{aligned}$$

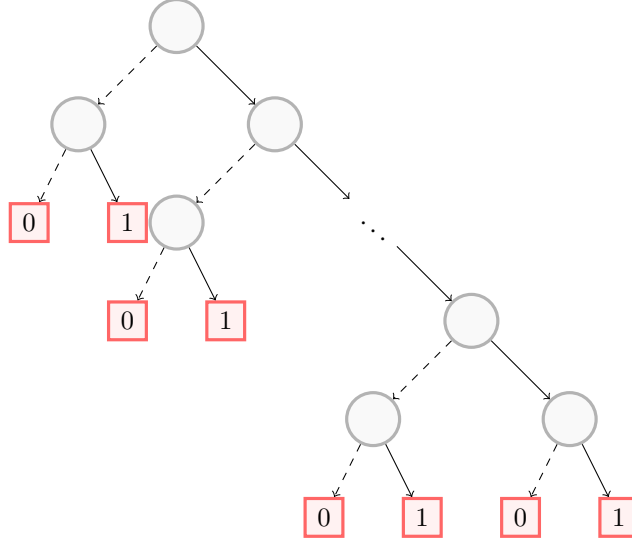
552 Consider now the condition $\exists t' \in pi((f \mid \bar{\ell}) \wedge (f \mid \ell)), t_{\ell} \models t'$ and suppose that it is satisfied. Since
 553 $pi((f \mid \bar{\ell}) \wedge (f \mid \ell)) = max(\{t'_{\bar{\ell}} \wedge t'_{\ell} : t'_{\bar{\ell}} \in pi(f \mid \bar{\ell}), t'_{\ell} \in pi(f \mid \ell)\}, \models)$, there exist $t'_{\bar{\ell}} \in pi(f \mid \bar{\ell})$
 554 and $t'_{\ell} \in pi(f \mid \ell)$ such that $t' = t'_{\bar{\ell}} \wedge t'_{\ell}$. Thus, we have $t_{\ell} \models t'_{\bar{\ell}} \wedge t'_{\ell}$, and in particular $t_{\ell} \models t'_{\ell}$ holds.
 555 But since t_{ℓ} and t'_{ℓ} are prime implicants of $f \mid \ell$, this implies that $t_{\ell} \equiv t'_{\ell}$ holds. Furthermore, from
 556 $t_{\ell} \models t'_{\bar{\ell}} \wedge t'_{\ell}$ we get that $t_{\ell} \models t'_{\bar{\ell}}$. In addition, a prime implicant $t'_{\bar{\ell}}$ of $f \mid \bar{\ell}$ such that $t_{\ell} \models t'_{\bar{\ell}}$ exists if
 557 and only if $t_{\ell} \models f \mid \bar{\ell}$. Altogether, the condition $\exists t' \in pi((f \mid \bar{\ell}) \wedge (f \mid \ell)), t_{\ell} \models t'$ is equivalent to
 558 $t_{\ell} \models f \mid \bar{\ell}$. Thus, we get that:

$$\begin{aligned}
 sr(\mathbf{x}, f) &= sr(\mathbf{x}, (f \mid \ell) \wedge (f \mid \bar{\ell})) \\
 559 &\cup \{\ell \wedge t_{\ell} : t_{\ell} \in sr(\mathbf{x}, f \mid \ell) \text{ s.t. } t_{\ell} \not\models f \mid \bar{\ell}\} \\
 &\text{where } Var(\ell) \subseteq Var(f) \text{ and } t_{\mathbf{x}} \models \ell
 \end{aligned}$$

560 Finally, if $t \in max(\{t_{\bar{x}} \wedge t_x : t_{\bar{x}} \in pi(f \mid \bar{x}), t_x \in pi(f \mid x)\}, \models)$, then by construction t is
 561 such that there exist $t_{\bar{x}} \in pi(f \mid \bar{x})$ and $t_x \in pi(f \mid x)$ satisfying $t = t_{\bar{x}} \wedge t_x$. If $t_{\mathbf{x}} \models t$ holds,
 562 then $t_{\mathbf{x}} \models t_{\bar{x}}$ and $t_{\mathbf{x}} \models t_x$ hold. Hence $t_{\bar{x}} \in sr(\mathbf{x}, f \mid \bar{x})$ and $t_x \in sr(\mathbf{x}, f \mid x)$. Consequently,
 563 $t \in max(\{t_{\bar{x}} \wedge t_x : t_{\bar{x}} \in sr(\mathbf{x}, f \mid \bar{x}), t_x \in sr(\mathbf{x}, f \mid x)\}, \models)$. \square

564 From the inductive characterization of $sr(\mathbf{x}, f)$ given by the previous proposition, we can easily
 565 derive a bottom-up algorithm allowing to derive $sr(\mathbf{x}, f)$ when f is represented by a decision tree.

566 Consider now a decision tree T of depth $k \geq 1$ having the following form:



567

568 T has $2k - 1$ decision nodes and $2k$ leaves. Suppose that the variables associated with the decision
 569 nodes are in one-to-one correspondence with the decision nodes (i.e., they are all distinct). The
 570 number of variables occurring in T is thus $n = 2k - 1$, therefore T has $2n + 1$ nodes. Consider now
 571 the instance $\mathbf{x} \in \{0, 1\}^n$ such that $x_i = 1$ for every $i \in [n]$. We are going to prove by induction on
 572 the depth k of such a tree T that \mathbf{x} has 2^{k-1} minimal reasons given T , each of them containing k
 573 literals. The proof takes advantage of the recursive characterization of the set of all sufficient reasons
 574 for an instance given a decision tree, as made precise by Lemma 1.

- 575 • Base case $k = 1$. We have $n = 1$. T consists of a decision node labelled by the single
 576 variable of X_n , say x , a left child that is a 0-leaf and a right child that is a 1-leaf. T is
 577 equivalent to x and x is implied by $t_{\mathbf{x}}$. Hence, x is the unique sufficient reason for \mathbf{x} given T ,
 578 so it is also the unique minimal reason for \mathbf{x} given T . As expected, the number of minimal
 579 reasons for \mathbf{x} given T is equal to 2^{k-1} . The size of the unique minimal reason is $k = 1$.
- 580 • Inductive step $k > 1$. Let x be the variable of X_n labelling the root node of T . By
 581 construction, the left child T_l of T is equivalent to a single variable, say x_l , that is the unique
 582 minimal reason for \mathbf{x} given T_l . The right child T_r of T has the same form as T , but with
 583 depth $k - 1$. By induction hypothesis, we know that \mathbf{x} has 2^{k-2} minimal reasons given T_r ,
 584 each of them containing $k - 1$ literals. As shown by Lemma 1, provided that the variables
 585 labelling the decision nodes are pairwise distinct, the minimal reasons for \mathbf{x} given T are
 586 obtained by extending every minimal reason for \mathbf{x} given T_r with x_l and by extending every
 587 minimal reason for \mathbf{x} given T_r with x . Accordingly, \mathbf{x} has $2 \times (2^{k-2}) = 2^{k-1}$ minimal
 588 reasons given T and each of them contains $k - 1 + 1 = k$ literals.

589 Finally, since $n = 2k - 1$, we have $k = \frac{n+1}{2}$ and the number of minimal reasons for \mathbf{x} given T is
 590 equal to $2^{k-1} = 2^{\frac{n-1}{2}}$. □

591 Proof of Proposition 3

592 *Proof.* The algorithms to compute $Nec_s(\mathbf{x}, T)$, $Rel_s(\mathbf{x}, T)$, and $Irr_s(\mathbf{x}, T)$ are as follows: first
 593 compute $CNF(T)$ and then remove from this set of clauses every literal that does not belong to $t_{\mathbf{x}}$.
 594 This can be done in $\mathcal{O}(n|T|)$ time. By construction, the resulting CNF formula f is monotone: every
 595 literal in it occurs with the same polarity as the one it has in $t_{\mathbf{x}}$. Furthermore, the size of f cannot
 596 exceed the size of $CNF(T)$, thus the size of T .

597 Since f is a monotone CNF formula, its prime implicants can be computed by removing from f every
 598 clause that is a strict superset of another clause of f . This can be achieved in quadratic time in the size
 599 of f , thus in the size of T . Let g be the resulting formula in prime implicants form and equivalent to
 600 f . g is equivalent to the complete reason for \mathbf{x} given T . Since it is in prime implicants form, g is

601 Lit-dependent on every literal occurring in it (i.e., g is Lit-simplified, see Proposition 8 in [23] for
602 details), hence so is the complete reason for x given T .

603 This means that for every literal ℓ occurring in g , there exists a sufficient reason for x given T that
604 contains ℓ , so that $Rel_s(x, T)$ is the set of literals occurring in g and $Irr_s(x, T)$ is the complement
605 of $Rel_s(x, T)$ in the set of all literals over X_n . Finally, since by definition the literals of $Nec_s(x, T)$
606 must belong to every sufficient reason for x given T , they are given by the unit clauses that belong to
607 g . \square

608 **Proof of Proposition 4**

609 *Proof.* We call MINIMAL REASON the problem that asks, given $T \in \text{DT}_n$, $x \in \{0, 1\}^n$ with
610 $T(x) = 1$ and $k \in \mathbb{N}$, whether there is an implicant t of T of size at most k that covers x .

611 Our objective is to prove that MINIMAL REASON is NP-hard. To this end, let us first recall that
612 a *vertex cover* of an undirected graph $G = (X, E)$ is a subset $V \subseteq X$ of vertices such that
613 $\{y, z\} \cap V \neq \emptyset$ for every edge $e = \{y, z\}$ in E . In the MIN VERTEX COVER problem, we are given
614 a graph G together with an integer $k \in \mathbb{N}$, and the task is to find a vertex cover V of G of size at most
615 k . MIN VERTEX COVER is a well-known NP-hard problem [20], and we now show that it can be
616 reduced in polynomial time to MINIMAL REASON.

617 Suppose that we are given a graph $G = (X, E)$ and assume, without loss of generality, that G does
618 not include isolated vertices. For any $y \in X$, let $E_y = \{e \in E : y \in e\}$ denote the set of edges in
619 G that are adjacent to y , and let $N_y = \{z \in X : \{y, z\} \in E\}$ denote the set of neighbors of y in
620 G . By $G \setminus y$, we denote the deletion of y from G , obtained by removing y and its adjacent edges,
621 i.e., $G \setminus y = (X \setminus \{y\}, E \setminus E_y)$. We associate with G a decision tree $T(G)$ over $X_n = X$ using the
622 following recursive algorithm. If G is the empty graph (i.e. $E = \emptyset$), then return the decision tree
623 rooted at a 1-leaf. Otherwise, pick a node $y \in X$ and generate a decision tree $T(G)$ such that:

- 624 (1) the root is labeled by y ;
- 625 (2) the left child is the decision tree encoding the monomial $\bigwedge N_y$;
- 626 (3) the right child is the decision tree $T(G')$ returned by calling the algorithm on $G' = G \setminus y$.

627 By construction, $T(G)$ is a complete backtrack search tree of the formula $\text{CNF}(E) = \bigwedge \{(y \vee z) :$
628 $\{y, z\} \in E\}$, which implies that $T(G)$ and $\text{CNF}(E)$ are logically equivalent. Furthermore, $T(G)$ is a
629 comb-shaped tree since recursion only on the rightmost branch. In particular, the algorithm runs in
630 $\mathcal{O}(n|E|)$ time, since step (1) takes $\mathcal{O}(1)$ time, step (2) takes $\mathcal{O}(n)$ time, and step (3) is called at most
631 $|E|$ times.

632 Now, with an instance $P_1 = (G, k)$ of MIN VERTEX COVER, we associate the instance $P_2 =$
633 $(T(G), x, k)$ of MINIMAL REASON, where $x = (1, \dots, 1)$. Based on the above algorithm, P_2 can
634 be constructed in time polynomial in the size of P_1 .

635 Let V be a solution of P_1 . Since V is a vertex cover of G , the term $t_V = \bigwedge V$ is an implicant of the
636 formula $\text{CNF}(E)$. Since $t_V \subseteq t_x$ and $|t_V| \leq k$, it follows from the fact that $\text{CNF}(E)$ and $T(G)$ are
637 logically equivalent that t_V is a solution of P_2 .

638 Conversely, let t be a solution of P_2 . Since t is an implicant of $T(G)$, it follows that t is an implicant
639 of $\text{CNF}(E)$. This together with the fact that $t \subseteq t_x$ implies that the subset of vertices $V \subseteq X_n$,
640 satisfying $\bigwedge V = t$, is a vertex cover of G . Since $|V| \leq k$, it is therefore a solution of P_1 . \square

642 **Proof of Proposition 5**

643 *Proof.* Let x^* be any solution of $(C_{\text{soft}}, C_{\text{hard}})$. Observe that the set of all hard clauses $c \cap t_x$ (where
644 c is a clause of $\text{CNF}(T)$) corresponds to a *monotone* CNF formula. Therefore, in order to satisfy such
645 a clause $c \cap t_x$, x^* must set a literal ℓ of t_x to 1. Thus, x^* satisfies all the hard clauses of C_{hard} if
646 and only if the term consisting of the literals that are shared by $t_x = \bigwedge_{i=1}^n \ell_i$ and t_{x^*} is an implicant
647 of T and is implied by x .

648 The soft clauses of C_{soft} are used to select among the assignments that satisfy all the hard clauses,
 649 the ones that correspond to minimal sufficient reasons. Soft clauses are given by literals ℓ_i , which
 650 are precisely the complementary literals to those occurring in t_x . Having a soft clause ℓ_i violated by
 651 x^* means that the literal $\bar{\ell}$ of t_x is necessary to get an implicant of T given the assignment of the
 652 other variables in x^* . Whenever a soft clause ℓ_i is violated by x^* a penalty of 1 incurs. This ensures
 653 that the term consisting of the literals that are shared by $t_x = \bigwedge_{i=1}^n \ell_i$ and t_{x^*} is a minimal sufficient
 654 reason for x given T . \square

655 Proof of Proposition 6

656 *Proof.* By definition, the sufficient reasons t for x given f are the prime implicants of f that covers
 657 x . Thus, they are precisely the prime implicants of the (conjunctively-interpreted) set of clauses
 658 $\{c \cap t_x : c \in \text{CNF}(f)\}$ where $\text{CNF}(f)$ is any CNF formula equivalent to f . Furthermore, the complete
 659 reason for x given f (equivalent to the disjunction of all the sufficient reasons for x given f [8]) is
 660 a monotone Boolean function because every sufficient reason covers x which assigns in a unique
 661 way every variable from X_n . The prime implicants of such a monotone function are precisely the
 662 minimal hitting sets of the prime implicants of the function. Because of the minimal hitting set
 663 duality between sufficient reasons and contrastive explanations for x given f [16], the contrastive
 664 explanations for x given f are thus the sets of literals corresponding to the prime implicants of
 665 $\{c \cap t_x : c \in \text{CNF}(f)\}$. Now, since the (conjunctively-interpreted) set of clauses $\{c \cap t_x : c \in \text{CNF}(f)\}$
 666 is equivalent to the complete reason for x given f , it is a monotone function, and as a consequence,
 667 its prime implicants are its minimal elements w.r.t. \subseteq . This comes from the correctness of any
 668 resolution-based algorithm for generating prime implicants (see e.g., [26]). Finally, when f is a
 669 decision tree T , $\{c \cap t_x : c \in \text{CNF}(T)\}$ can be computed in time polynomial in $n + |T|$ because
 670 $\text{CNF}(T)$ can be computed in time linear in $|T|$. Using an extra quadratic time in the size of this
 671 set $\{c \cap t_x : c \in \text{CNF}(T)\}$, its minimal elements w.r.t. \subseteq can be selected. The resulting set is by
 672 construction the set of all the contrastive explanations for x given T , and this set has been computed
 673 in time polynomial in $n + |T|$. \square