# A New Class of Explanations for Classifiers with Non-Binary Features

Chunxi Ji[0000−0002−4475−1987] and Adnan Darwiche[0000−0003−3976−6735]

University of California, Los Angeles, CA 90095, USA

**Abstract.** Two types of explanations have been receiving increased attention in the literature when analyzing the decisions made by classifiers. The first type explains why a decision was made and is known as a sufficient reason for the decision, also an abductive explanation or a PI-explanation. The second type explains why some other decision was not made and is known as a necessary reason for the decision, also a contrastive or counterfactual explanation. These explanations were defined for classifiers with binary, discrete and, in some cases, continuous features. We show that these explanations can be significantly improved in the presence of non-binary features, leading to a new class of explanations that relay more information about decisions and the underlying classifiers. Necessary and sufficient reasons were also shown to be the prime implicates and implicants of the complete reason for a decision, which can be obtained using a quantification operator. We show that our improved notions of necessary and sufficient reasons are also prime implicates and implicants but for an improved notion of complete reason obtained by a new quantification operator that we also define and study.

**Keywords:** Explainable AI · Decision Graphs · Prime Implicants/Implicates.

## 1 Introduction

Explaining the decisions of classifiers has been receiving significant attention in the AI literature recently. Some explanation methods operate directly on classifiers, e.g., [43,42], while some other methods operate on symbolic encodings of their input-output behavior, e.g., [8,25,36,39], which may be compiled into tractable circuits [11,45,46,44,5,21]. When explaining the decisions of classifiers, two particular notions have been receiving increased attention in the literature: The sufficient and necessary reasons for a decision on an instance.

A *sufficient reason* for a decision [17] is a minimal subset of the instance which is guaranteed to trigger the decision. It was first introduced under the name *PI-explanation* in [45] and later called an *abductive explanation* [25].[1] Consider the classifier in Figure 1a and a patient, Susan, with the following characteristics: AGE ≥55, BTYPE=A and WEIGHT=OVER. Susan is judged as susceptible to disease by this classifier, and a sufficient reason for this decision is {AGE ≥55, BTYPE=A}.

---

[1] We will use sufficient reasons and PI/abductive explanations interchangeably.
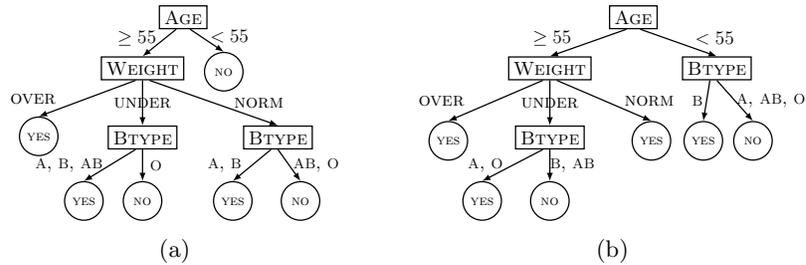
Fig. 1: Two classifiers of patients susceptible to a certain disease. The classifier in (b) will be discussed later in the paper.

Hence, the classifier will judge Susan as susceptible to disease as long as she has these two characteristics, regardless of how the feature WEIGHT is set.[2]

A *necessary reason* for a decision [18] is a minimal subset of the instance that will flip the decision if changed appropriately. It was formalized earlier in [24] under the name *contrastive explanation* which is discussed initially in [32,38].[3] Consider again the patient Susan and the classifier in Figure 1a. A necessary reason for the decision on Susan is {AGE $\geq$55}, which means that she would not be judged as susceptible to disease if she were younger than 55. The other necessary reason is {WEIGHT=OVER, BTYPE=A} so the decision on Susan can be flipped by changing these two characteristics (and this cannot be achieved by changing only one of them). Indeed, if Susan had WEIGHT=NORM and BTYPE=AB, she will not be judged as susceptible. However, since WEIGHT and BTYPE are discrete variables, there are multiple ways for changing them and some changes may not flip the decision (e.g., WEIGHT=UNDER and BTYPE=B).

The notion of a *complete reason* behind a decision was introduced in [17] and its prime implicants were shown to be the sufficient reasons for the decision. Intuitively, the complete reason is a particular condition on the instance that is both necessary and sufficient for the decision on that instance; see [16]. A declarative semantics for complete reasons was given in [19] which showed how to compute them using *universal literal quantification*. Furthermore, the prime implicates of a complete reason where shown to be the necessary reasons for the decision in [18]. Given these results, one would first use universal literal quantification to obtain the complete reason for a decision and then compute its prime implicates and implicants to obtain necessary and sufficient explanations.

---

[2] See, e.g., [13,43,48] for some approaches that can be viewed as approximating sufficient reasons and [26] for a study of the quality of some of these approximations.

[3] We will use necessary reasons and contrastive explanations interchangeably in this paper. Counterfactual explanations are related but have alternate definitions in the literature. For example, as defined in [5], they correspond to length-minimal necessary reasons; see [18]. But according to some other definitions, they include contrastive explanations (necessary reasons) as a special case; see Section 5.2 in [33]. See also [1] for counterfactual explanations that are directed towards Bayesian network classifiers and [2] for a relevant recent study and survey.

Necessary and sufficient reasons are *subsets* of the instance being explained so each reason corresponds to a set of variable settings (Feature=Value), like WEIGHT=UNDER and BTYPE=B, which we shall call *simple literals.* Since necessary and sufficient reasons correspond to sets of simple literals, we will refer to them as *simple* or *classical* explanations. We will show next that these simple explanations can be significantly improved if the classifier has non-binary features, leading to more general notions of necessary, sufficient and complete reasons that provide more informative explanations of decisions.

Consider again the decision on Susan discussed above which had the sufficient reason {AGE $\geq$55, BTYPE=A}. Such an explanation can be viewed as a *property* of the instance which guarantees the decision. The property has a specific form: a conjunction of feature settings (i.e., instance characteristics) which leaves out characteristics of the instance that are irrelevant to the decision (WEIGHT=OVER). However, the following is a weaker property of the instance which will also trigger the decision: {AGE $\geq$55, BTYPE $\in$ {A, B}}. This property tells us that not only is WEIGHT=OVER irrelevant to the decision, but also that BTYPE=A is not particularly relevant since BTYPE could have been B and the decision would have still been triggered. In other words, what is really relevant is that BTYPE $\in$ {A, B} or, alternatively, BTYPE $\notin$ {AB, O}. Clearly, this kind of explanation reveals more information about why the classifier made its decision. We will later formalize and study a new class of explanations for this purpose, called *general sufficient reasons,* which arise only when the classifier has non-binary features.

A necessary reason for a decision can also be understood as a property of the instance, but one that will flip the decision if violated in a *certain* manner [18]. As mentioned earlier, {WEIGHT=OVER, BTYPE=A} is a necessary reason for the decision on Susan. This reason corresponds to the property (WEIGHT=OVER or BTYPE=A). We can flip the decision by violating this property through changing the values of WEIGHT and BTYPE in the instance. Since these variables are non-binary, there are multiple changes (six total) that will violate the property. Some violations will flip the decision, others will not (we are only guaranteed that at least one violation will flip the decision). For example, WEIGHT=NORM, BTYPE=O and WEIGHT=UNDER, BTYPE=AB will both violate the property but only the first one will flip the decision. However, the following weaker property is guaranteed to flip the decision regardless of how it is violated: (WEIGHT=OVER or BTYPE $\in$ {A, B, AB}). We can violate this property using two different settings of WEIGHT and BTYPE, both of which will flip the decision. This property corresponds to the *general necessary reason* {WEIGHT=OVER, BTYPE $\in$ {A, B, AB}}, a new notion that we introduce and study later. Similar to general sufficient reasons, general necessary reasons provide more information about the behavior of a classifier and arise only when the classifier has non-binary features.

We stress here that using simple explanations in the presence of non-binary features is quite prevalent in the literature; see, e.g., [4,6,8,18,23,28,35]. Two notable exceptions are [12,27] which we discuss in more detail later.[4]

---

[4] Interestingly, the axiomatic study of explanations in [3] allows non-binary features, yet Axiom 4 (*feasibility*) implies that explanations must be simple.

Our study of general necessary and sufficient reasons follows a similar structure to recent developments on classical necessary and sufficient reasons. In particular, we define a new quantification operator like the one defined in [19] and show how it can be used to compute the *general reason* of a decision, and that its prime implicates and implicants contain the general necessary and sufficient reasons. Complete reasons are known to be monotone formulas. We show that general reasons are *fixated formulas* which include monotone ones. We introduce the fixation property and discuss some of its (computational) implications.

This paper is structured as follows. We start in Section 2 by discussing the syntax and semantics of formulas with discrete variables which are needed to capture the input-output behavior of classifiers with non-binary features. We then introduce the new quantification operator in Section 3 where we study its properties and show how it can be used to formulate the new notion of general reason. The study of general necessary and sufficient reasons is conducted in Section 4 where we also relate them to their classical counterparts and argue further for their utility. Section 5 provides closed-form general reasons for a broad class of classifiers and Section 6 discusses the computation of general necessary and sufficient reasons based on general reasons. We finally close with some remarks in Section 7. Proofs of all results can be found in Appendix A.

## 2   Representing Classifiers using Class Formulas

We now discuss the syntax and semantics of *discrete formulas,* which we use to represent the input-output behavior of classifiers. Such symbolic formulas can be automatically compiled from certain classifiers, like Bayesian networks, random forests and some types of neural networks; see [16] for a summary.

We assume a finite set of variables $\Sigma$ which represent classifier features. Each variable $X \in \Sigma$ has a finite number of *states* $x_1, \ldots, x_n$, $n > 1$. A *literal* $\ell$ for variable $X$, called $X$-literal, is a set of states such that $\emptyset \subset \ell \subset \{x_1, \ldots, x_n\}$. We will often denote a literal such as $\{x_1, x_3, x_4\}$ by $x_{134}$ which reads: the state of variable $X$ is either $x_1$ or $x_3$ or $x_4$. A literal is *simple* iff it contains a single state. Hence, $x_3$ is a simple literal but $x_{134}$ is not. Since a simple literal corresponds to a state, these two notions are interchangeable.

A *formula* is either a constant $\top$, $\bot$, literal $\ell$, negation $\overline{\alpha}$, conjunction $\alpha \cdot \beta$ or disjunction $\alpha + \beta$ where $\alpha$, $\beta$ are formulas. The set of variables appearing in a formula $\Delta$ are denoted by $vars(\Delta)$. A *term* is a conjunction of literals for distinct variables. A *clause* is a disjunction of literals for distinct variables. A *DNF* is a disjunction of terms. A *CNF* is a conjunction of clauses. An *NNF* is a formula without negations. These definitions imply that terms cannot be inconsistent, clauses cannot be valid, and negations are not allowed in DNFs, CNFs, or NNFs. Finally, we say a term/clause is *simple* iff it contains only simple literals.

A *world* maps each variable in $\Sigma$ to one of its states and is typically denoted by $\omega$. A world $\omega$ is called a *model* of formula $\alpha$, written $\omega \models \alpha$, iff $\alpha$ is satisfied by $\omega$ (that is, $\alpha$ is true at $\omega$). The constant $\top$ denotes a valid formula (satisfied by every world) and the constant $\bot$ denotes an unsatisfiable formula (has no

models). Formula $\alpha$ implies formula $\beta$, written $\alpha \models \beta$, iff every model of $\alpha$ is also a model of $\beta$. A term $\tau_1$ subsumes another term $\tau_2$ iff $\tau_2 \models \tau_1$. A clause $\sigma_1$ subsumes another clause $\sigma_2$ iff $\sigma_1 \models \sigma_2$. Formula $\alpha$ is weaker than formula $\beta$ iff $\beta \models \alpha$ (hence $\beta$ is stronger than $\alpha$).

The *conditioning* of formula $\Delta$ on simple term $\tau$ is denoted $\Delta|\tau$ and obtained as follows. For each state $x$ of variable $X$ that appears in term $\tau$, replace each $X$-literal $\ell$ in $\Delta$ with $\top$ if $x \in \ell$ and with $\bot$ otherwise. Note that $\Delta|\tau$ does not mention any variable that appears in term $\tau$. A *prime implicant* for a formula $\Delta$ is a term $\alpha$ such that $\alpha \models \Delta$, and there does not exist a distinct term $\beta$ such that $\alpha \models \beta \models \Delta$. A *prime implicate* for a formula $\Delta$ is a clause $\alpha$ such that $\Delta \models \alpha$, and there does not exist a distinct clause $\beta$ such that $\Delta \models \beta \models \alpha$.

An *instance* of a classifier will be represented by a simple term which contains exactly one literal for each variable in $\Sigma$. A classifier with $n$ classes will be represented by a set of mutually exclusive and exhaustive formulas $\Delta^1, \ldots, \Delta^n$, where the models of formula $\Delta^i$ capture the instances in the $i^{th}$ class. That is, instance $\mathcal{I}$ is in the $i^{th}$ class iff $\mathcal{I} \models \Delta^i$. We refer to each $\Delta^i$ as a *class formula,* or simply a *class,* and say that instance $\mathcal{I}$ is in class $\Delta^i$ when $\mathcal{I} \models \Delta^i$.

Consider the decision diagram on the right which represents a classifier with three ternary features $(X, Y, Z)$ and three classes $c_1$, $c_2$, and $c_3$. This classifier can be represented by the class formulas $\Delta^1 = x_{12} + x_3 \cdot y_1 \cdot z_{13}$, $\Delta^2 = x_3 \cdot z_2$ and $\Delta^3 = x_3 \cdot y_{23} \cdot z_{13}$. This classifier has 27 instances, partitioned as follows: 20 instances in class $c_1$, 3 in class $c_2$ and 4 in class $c_3$. For example, instance $\mathcal{I} = x_3 \cdot y_2 \cdot z_2$ belongs to class $c_2$ since $\mathcal{I} \models \Delta^2$.



## 3   The General Reason for a Decision

An operator $\forall x$ which eliminates the state $x$ of a Boolean variable $X$ from a formula was introduced and studied in [19]. This operator, called universal literal quantification, was also generalized in [19] to the states of discrete variables but without further study. Later, [18] studied this discrete generalization, given next.

**Definition 1.** *For variable $X$ with states $x_1, \ldots, x_n$, the universal literal quantification of state $x_i$ from formula $\Delta$ is defined as $\forall x_i \cdot \Delta = \Delta|x_i \cdot \prod_{j \neq i} (x_i + \Delta|x_j)$.*

The operator $\forall$ is commutative so we can equivalently write $\forall x \cdot (\forall y \cdot \Delta)$, $\forall y \cdot (\forall x \cdot \Delta)$, $\forall x, y \cdot \Delta$ or $\forall \{x, y\} \cdot \Delta$. It is meaningful then to quantify an instance $\mathcal{I}$ from its class formula $\Delta$ since $\mathcal{I}$ is a set of states. As shown in [19], the quantified formula $\forall \mathcal{I} \cdot \Delta$ corresponds to the complete reason for the decision on instance $\mathcal{I}$. Hence, the prime implicants of $\forall \mathcal{I} \cdot \Delta$ are the sufficient reasons for the decision [17] and its prime implicates are the necessary reasons [18].

We next define a new operator $\overline{\forall}$ that we call a *selection operator* for reasons that will become apparent later. This operator will lead to the notion of a general reason for a decision which subsumes the decision's complete reason, and provides the basis for defining general necessary and sufficient reasons.

**Definition 2.** *For variable $X$ with states $x_1, \ldots, x_n$ and formula $\Delta$, we define $\overline{\forall} x_i \cdot \Delta$ to be $\Delta | x_i \cdot \Delta$.*

The selection operator $\overline{\forall}$ is also commutative, like $\forall$.

**Proposition 1.** $\overline{\forall} x \cdot (\overline{\forall} y \cdot \Delta) = \overline{\forall} y \cdot (\overline{\forall} x \cdot \Delta)$ *for states $x, y$.*

Since a term $\tau$ corresponds to a set of states, the expression $\overline{\forall} \tau \cdot \Delta$ is well-defined just like $\forall \tau \cdot \Delta$. We can now define our first major notion.

**Definition 3.** *Let $\mathcal{I}$ be an instance in class $\Delta$. The general reason for the decision on instance $\mathcal{I}$ is defined as $\overline{\forall} \mathcal{I} \cdot \Delta$.*

The complete reason $\forall \mathcal{I} \cdot \Delta$ can be thought of as a property/abstraction of instance $\mathcal{I}$ that justifies (i.e., can trigger) the decision. In fact, it is equivalent to the weakest NNF $\Gamma$ whose literals appear in the instance and that satisfies $\mathcal{I} \models \Gamma \models \Delta$ [19,18]. The next result shows that the general reason is a weaker property and, hence, a further abstraction that triggers the decision.

**Proposition 2.** *For instance $\mathcal{I}$ and formula $\Delta$ where $\mathcal{I} \models \Delta$, we have $\mathcal{I} \models \forall \mathcal{I} \cdot \Delta \models \overline{\forall} \mathcal{I} \cdot \Delta \models \Delta$. ($\mathcal{I} \not\models \Delta$ only if $\forall \mathcal{I} \cdot \Delta = \overline{\forall} \mathcal{I} \cdot \Delta = \bot$)*

The next result provides further semantics for the general reason and highlights the key difference with the complete reason.

**Proposition 3.** *The general reason $\overline{\forall} \mathcal{I} \cdot \Delta$ is equivalent to the weakest NNF $\Gamma$ whose literals are implied by instance $\mathcal{I}$ and that satisfies $\mathcal{I} \models \Gamma \models \Delta$.*

The complete and general reasons are abstractions of the instance that explain why it belongs to its class. The former can only reference simple literals in the instance but the latter can reference any literal that is implied by the instance. The complete reason can be recovered from the general reason and the underlying instance. Moreover, the two types of reasons are equivalent when all variables are binary since $\forall x \cdot \Delta = \overline{\forall} x \cdot \Delta$ when $x$ is the state of a binary variable.

We next provide a number of results that further our understanding of general reasons, particularly their semantics and how to compute them. We start with the following alternative definition of the operator $\overline{\forall} x_i$.

**Proposition 4.** *For formula $\Delta$ and variable $X$ with states $x_1, \ldots, x_n$, $\overline{\forall} x_i \cdot \Delta$ is equivalent to $(\Delta | x_i) \cdot \prod_{j \neq i} (\ell_j + (\Delta | x_j))$, where $\ell_j$ is the literal $\{x_1, \ldots, x_n\} \setminus \{x_j\}$.*

According to this definition, we can always express $\overline{\forall} x_i \cdot \Delta$ as an NNF in which every $X$-literal includes state $x_i$ (recall that $\Delta | x_i$ and $\Delta | x_j$ do not mention variable $X$). This property is used in the proofs and has a number of implications.[5]

---

[5] For example, we can use it to provide *forgetting* semantics for the dual operator $\overline{\exists} x_i \cdot \Delta = \overline{\overline{\forall} x_i \cdot \overline{\Delta}}$. Using Definition 2, we get $\overline{\exists} x_i \cdot \Delta = \Delta + \Delta | x_i$. Using Proposition 4, we get $\overline{\exists} x_i \cdot \Delta = \Delta | x_i + \sum_{j \neq i} (x_j \cdot \Delta | x_j)$. We can now easily show that (1) $\Delta \models \overline{\exists} x_i \cdot \Delta$ and (2) $\overline{\exists} x_i \cdot \Delta$ is equivalent to an NNF whose $X$-literals do not mention state $x_i$. That is, $\overline{\exists} x_i$ can be understood as forgetting the information about state $x_i$ from $\Delta$. This is similar to the dual operator $\exists x_i \cdot \Delta = \overline{\forall x_i \cdot \overline{\Delta}}$ studied in [31,19] except that $\overline{\exists} x_i$ erases less information from $\Delta$ since one can show that $\Delta \models \overline{\exists} x_i \cdot \Delta \models \exists x_i \cdot \Delta$.

When $\Delta$ is a class formula, [19] showed that the application of $\forall x$ to $\Delta$ can be understood as *selecting* a specific set of instances from the corresponding class. This was shown for states $x$ of Boolean variables. We next generalize this to discrete variables and provide a selection semantics for the new operator $\overline{\forall}$.

**Proposition 5.** *Let $\tau$ be a simple term, $\Delta$ be a formula and $\omega$ be a world. Then $\omega \models \forall \tau \cdot \Delta$ iff $\omega \models \Delta$ and $\omega' \models \Delta$ for any world $\omega'$ obtained from $\omega$ by changing the states of some variables that are set differently in $\tau$. Moreover, $\omega \models \overline{\forall} \tau \cdot \Delta$ iff $\omega \models \Delta$ and $\omega' \models \Delta$ for any world $\omega'$ obtained from $\omega$ by setting some variables in $\omega$ to their states in $\tau$.*

That is, $\forall \tau \cdot \Delta$ selects all instances in class $\Delta$ whose membership in the class does not depend on characteristics that are inconsistent with $\tau$. These instances are also selected by $\overline{\forall} \tau \cdot \Delta$ which further selects instances that remain in class $\Delta$ when any of their characteristics are changed to agree with $\tau$.

The complete reason is monotone which has key computational implications as shown in [17,19,18]. The general reason satisfies a weaker property called *fixation* which has also key computational implications as we show in Section 6.

**Definition 4.** *An NNF is locally fixated on instance $\mathcal{I}$ iff its literals are consistent with $\mathcal{I}$. A formula is fixated on instance $\mathcal{I}$ iff it is equivalent to an NNF that is locally fixated on $\mathcal{I}$.*

We also say in this case that the formula is $\mathcal{I}$-fixated. For example, if $\mathcal{I} = x_1 \cdot y_1 \cdot z_2$ then the formula $x_{12} \cdot y_1 + z_2$ is (locally) $\mathcal{I}$-fixated but $x_{12} \cdot z_1$ is not. By the selection semantic we discussed earlier, a formula $\Delta$ is $\mathcal{I}$-fixated only if for every model $\omega$ of $\Delta$, changing the states of some variables in $\omega$ to their states in $\mathcal{I}$ guarantees that the result remains a model of $\Delta$. Moreover, if $\Delta$ is $\mathcal{I}$-fixated, then $\mathcal{I} \models \Delta$ but the opposite does not hold (e.g., $\Delta = x_1 + y_1$ and $\mathcal{I} = x_1 \cdot y_2$). We now have the following corollary of Proposition 3.

**Corollary 1.** *The general reason $\overline{\forall} \mathcal{I} \cdot \Delta$ is $\mathcal{I}$-fixated.*

The next propositions show that the new operator $\overline{\forall}$ has similar computational properties to $\forall$ which we use in Section 5 to compute general reasons.

**Proposition 6.** *For state $x$ and literal $\ell$ of variable $X$, $\overline{\forall} x \cdot \ell = \ell$ if $x \in \ell$ ($x \models \ell$); else $\overline{\forall} x \cdot \ell = \bot$. Moreover, $\overline{\forall} x \cdot \Delta = \Delta$ if $X$ does not appear in $\Delta$.*

**Proposition 7.** *For formulas $\alpha$, $\beta$ and state $x_i$ of variable $X$, we have $\overline{\forall} x_i \cdot (\alpha \cdot \beta) = (\overline{\forall} x_i \cdot \alpha) \cdot (\overline{\forall} x_i \cdot \beta)$. Moreover, if variable $X$ does not occur in both $\alpha$ and $\beta$, then $\overline{\forall} x_i \cdot (\alpha + \beta) = (\overline{\forall} x_i \cdot \alpha) + (\overline{\forall} x_i \cdot \beta)$.*

An NNF is $\vee$-decomposable if its disjuncts do not share variables. According to these propositions, we can apply $\overline{\forall} \mathcal{I}$ to an $\vee$-decomposable NNF in linear time, by simply applying $\overline{\forall} \mathcal{I}$ to each literal in the NNF (the result is $\vee$-decomposable).

## 4    General Necessary and Sufficient Reasons

We next introduce generalizations of necessary and sufficient reasons and show
that they are prime implicates and implicants of the general reason for a deci-
sion. These new notions have more explanatory power and subsume their classi-
cal counterparts, particularly when explaining the behavior of a classifier beyond
a specific instance/decision. For example, when considering the classifier in Fig-
ure 1b, which is a variant of the one in Figure 1a, we will see that the two
classifiers will make identical decisions on some instances, leading to identical
simple necessary and sufficient reasons for these decisions but distinct general
necessary and sufficient reasons. Moreover, we will see that general necessary
and sufficient reasons are particularly critical when explaining the behavior of
classifiers with (discretized) numeric features.

### 4.1    General Sufficient Reasons (GSRs)

We start by defining the classical notion of a (simple) sufficient reason but using
a different formulation than [45] which was the first to introduce this notion
under the name of a PI-explanation. Our formulation is meant to highlight a
symmetry with the proposed generalization.

**Definition 5 (SR).** *A sufficient reason for the decision on instance $\mathcal{I}$ in class
$\Delta$ is a weakest simple term $\tau$ s.t. $\mathcal{I} \models \tau \models \Delta$.*

This definition implies that each literal in $\tau$ is a variable setting (i.e., character-
istic) that appears in instance $\mathcal{I}$. That is, the (simple) literals of sufficient reason
$\tau$ are a subset of the literals in instance $\mathcal{I}$. We now define our generalization.

**Definition 6 (GSR).** *A general sufficient reason for the decision on instance
$\mathcal{I}$ in class $\Delta$ is a term $\tau$ which satisfies (1) $\tau$ is a weakest term s.t. $\mathcal{I} \models \tau \models \Delta$
and (2) no term $\tau'$ satisfies the previous condition if $vars(\tau') \subset vars(\tau)$.*

This definition does not require the GSR $\tau$ to be a simple term, but it requires
that it has a minimal set of variables. Without this minimality condition, a GSR
will be redundant in the sense of the upcoming Proposition 8. For a term $\tau$ and
instance $\mathcal{I}$ s.t. $\mathcal{I} \models \tau$, we will use $\mathcal{I} \dot{\cap} \tau$ to denote the smallest subterm in $\mathcal{I}$ that
implies $\tau$. For example, if $\mathcal{I} = x_2 \cdot y_1 \cdot z_3$ and $\tau = x_{12} \cdot y_{13}$, then $\mathcal{I} \dot{\cap} \tau = x_2 \cdot y_1$.

**Proposition 8.** *Let $\mathcal{I}$ be an instance in class $\Delta$ and $\tau$ be a weakest term s.t.
$\mathcal{I} \models \tau \models \Delta$. If $\tau'$ is a weakest term s.t. $\mathcal{I} \models \tau' \models \Delta$ and $vars(\tau') \subset vars(\tau)$,
then $\mathcal{I} \dot{\cap} \tau \models \mathcal{I} \dot{\cap} \tau' \models \Delta$. Also, $\mathcal{I} \dot{\cap} \tau$ is a SR iff such a term $\tau'$ does not exist.*

According to this proposition, the term $\tau$ is redundant as an explanation in that
the subset of instance $\mathcal{I}$ which it identifies as being a culprit for the decision
($\mathcal{I} \dot{\cap} \tau$) is dominated by a smaller subset that is identified by the term $\tau'$ ($\mathcal{I} \dot{\cap} \tau'$).
    Consider the classifiers in Figures 1a and 1b and the patient Susan: AGE $\geq$55,
BTYPE=A and WEIGHT=OVER. Both classifiers will make the same decision YES on

Susan with the same SRs: ($\textsc{Age} \geq 55 \cdot \textsc{Btype}=\text{A}$) and ($\textsc{Age} \geq 55 \cdot \textsc{Weight}=\textsc{over}$). The GSRs are different for these two (equal) decisions. For the first classifier, they are ($\textsc{Age} \geq 55 \cdot \textsc{Btype} \in \{\text{A}, \text{B}\}$) and ($\textsc{Age} \geq 55 \cdot \textsc{Weight}=\textsc{over}$). For the second, they are ($\textsc{Age} \geq 55 \cdot \textsc{Btype} \in \{\text{A}, \text{O}\}$) and ($\textsc{Age} \geq 55 \cdot \textsc{Weight} \in \{\textsc{over}, \textsc{norm}\}$). GSRs encode all SRs and contain more information.[6]

**Proposition 9.** *Let $\tau$ be a simple term. Then $\tau$ is a SR for the decision on instance $\mathcal{I}$ iff $\tau = \mathcal{I} \mathbin{\dot{\cap}} \tau'$ for some GSR $\tau'$.*

Consider the instance Susan again, $\mathcal{I} = (\textsc{Age} \geq 55) \cdot (\textsc{Btype}=\text{A}) \cdot (\textsc{Weight}=\textsc{over})$ and the classifier in Figure 1b. As mentioned, the GSRs for the decision on Susan are $\tau'_1 = (\textsc{Age} \geq 55 \cdot \textsc{Btype} \in \{\text{A}, \text{O}\})$ and $\tau'_2 = (\textsc{Age} \geq 55 \cdot \textsc{Weight} \in \{\textsc{over}, \textsc{norm}\})$ so $\tau_1 = \mathcal{I} \mathbin{\dot{\cap}} \tau'_1 = (\textsc{Age} \geq 55 \cdot \textsc{Btype}=\text{A})$ and $\tau_2 = \mathcal{I} \mathbin{\dot{\cap}} \tau'_2 = (\textsc{Age} \geq 55 \cdot \textsc{Weight}=\textsc{over})$, which are the two SRs for the decision on Susan.

The use of general terms to explain the decision on an instance $\mathcal{I}$ in class $\Delta$ was first suggested in [12]. This work proposed the notion of a general PI-explanation as a prime implicant of $\Delta$ that is consistent with instance $\mathcal{I}$. This definition is equivalent to Condition (1) in our Definition 6 which has a second condition relating to variable minimality. Hence, the definition proposed by [12] does not satisfy the desirable properties stated in Propositions 8 and 9 which require this minimality condition. The merits of using general terms were also discussed when explaining decision trees in [27], which introduced the notion of an *abductive path explanation (APXp)*. In a nutshell, each path in a decision tree corresponds to a general term $\tau$ that implies the formula $\Delta$ of the path's class. Such a term is usually used to explain the decisions made on instances that follow that path. As observed in [27], such a term can often be shortened, leading to an APXp that still implies the class formula $\Delta$ and hence provides a better explanation. An APXp is an implicant of the class formula $\Delta$ but not necessarily a prime implicant (or a variable-minimal prime implicant). Moreover, an APXp is a property of the specific decision tree (syntax) instead of its underlying classifier (semantics). See Appendix B for further discussion of these limitations.[7]

## 4.2 General Necessary Reasons (GNRs)

We now turn to simple necessary reasons and their generalizations. A necessary reason is a property of the instance that will flip the decision if violated in a certain way (by changing the instance). As mentioned earlier, the difference between the classical necessary reason and the generalized one is that the latter comes with stronger guarantees. Again, we start with a definition of classical necessary reasons using a different phrasing than [24] which formalized them under the name of contrastive explanations [32]. Our phrasing, based on [18], highlights a symmetry with the generalization and requires the following notation.

---

[6] Unlike SRs, two GSRs may mention the same set of variables. Consider the class formula $\Delta = (x_1 \cdot y_{12}) + (x_{12} \cdot y_1)$ and instance $\mathcal{I} = x_1 \cdot y_1$. There are two GSRs for the decision on $\mathcal{I}$, $x_1 \cdot y_{12}$ and $x_{12} \cdot y_1$, and both mention the same variables $X, Y$.

[7] A dual notion, contrastive path explanation (CPXp), was also proposed in [27].

For a clause $\sigma$ and instance $\mathcal{I}$ s.t. $\mathcal{I} \models \sigma$, we will use $\mathcal{I} \backslash\!\backslash \sigma$ to denote the largest subterm of $\mathcal{I}$ that does not imply $\sigma$. For example, if $\mathcal{I} = x_2 \cdot y_1 \cdot z_3$ and $\sigma = x_{12} + y_{13}$ then $\mathcal{I} \backslash\!\backslash \sigma = z_3$. We will also write $\mathcal{I} \models \sigma$ to mean that instance $\mathcal{I}$ implies every literal in clause $\sigma$. For instance $\mathcal{I} = x_2 \cdot y_1 \cdot z_3$, we have $\mathcal{I} \models x_{12} + y_{13}$ but $\mathcal{I} \not\models x_{12} + y_{23}$ even though $\mathcal{I} \models x_{12} + y_{23}$.

**Definition 7 (NR).** *A necessary reason for the decision on instance $\mathcal{I}$ in class $\Delta$ is a strongest simple clause $\sigma$ s.t. $\mathcal{I} \models \sigma$ and $(\mathcal{I} \backslash\!\backslash \sigma) \cdot \overline{\sigma} \not\models \Delta$ (if we minimally change the instance to violate $\sigma$, it is no longer guaranteed to stay in class $\Delta$).*

A necessary reason guarantees that *some* minimal change to the instance which violates the reason will flip the decision. But it does not guarantee that *all* such changes will. A general necessary reason comes with a stronger guarantee.

**Definition 8 (GNR).** *A general necessary reason for the decision on instance $\mathcal{I}$ in class $\Delta$ is a strongest clause $\sigma$ s.t. $\mathcal{I} \models \sigma$, $(\mathcal{I} \backslash\!\backslash \sigma) \cdot \overline{\sigma} \models \overline{\Delta}$, and no clause $\sigma'$ satisfies the previous conditions if $vars(\sigma') \subset vars(\sigma)$.*

The key difference between Definitions 7 and 8 are the conditions $(\mathcal{I} \backslash\!\backslash \sigma) \cdot \overline{\sigma} \not\models \Delta$ and $(\mathcal{I} \backslash\!\backslash \sigma) \cdot \overline{\sigma} \models \overline{\Delta}$. The first condition guarantees that *some* violation of a NR will flip the decision (by placing the modified instance outside class $\Delta$) while the second condition guarantees that *all* violations of a GNR will flip the decision.

The next proposition explains why we require GNRs to be variable-minimal. Without this condition, the changes identified by a GNR to flip the decision may not be minimal (we can flip the decision by changing a strict subset of variables).

For instance $\mathcal{I}$ and clause $\sigma$ s.t. $\mathcal{I} \models \sigma$, we will use $\mathcal{I} \dot{\cap} \sigma$ to denote the disjunction of states that appear in both $\mathcal{I}$ and $\sigma$ (hence, $\mathcal{I} \dot{\cap} \sigma \models \sigma$). For example, if $\mathcal{I} = x_1 \cdot y_1 \cdot z_1$ and $\sigma = x_{12} + y_{23} + z_1$, then $\mathcal{I} \dot{\cap} \sigma = x_1 + z_1$.

**Proposition 10.** *Let $\mathcal{I}$ be an instance in class $\Delta$ and let $\sigma$ be a strongest clause s.t. $\mathcal{I} \models \sigma$ and $(\mathcal{I} \backslash\!\backslash \sigma) \cdot \overline{\sigma} \models \overline{\Delta}$. If $\sigma'$ is another strongest clause satisfying these conditions and $vars(\sigma') \subset vars(\sigma)$, then $\mathcal{I} \backslash\!\backslash \sigma' \models \mathcal{I} \backslash\!\backslash \sigma$. Moreover, $\mathcal{I} \dot{\cap} \sigma$ is a NR iff such a clause $\sigma'$ does not exist.*

That is, if violating $\sigma$ requires changing some characteristics $C$ of instance $\mathcal{I}$, then $\sigma'$ can be violated by changing a strict subset of these characteristics $C$.

Consider the classifiers in Figures 1a and 1b which make the same decision, YES, on Susan (AGE $\geq$ 55, BTYPE=A, WEIGHT=OVER). The NRs for these equal decisions are the same: (AGE $\geq$ 55) and (WEIGHT=OVER+BTYPE=A). The GNRs for the classifier in Figure 1a are (AGE $\geq$ 55), (BTYPE $\in \{$A, B, AB$\}$+ WEIGHT=OVER$\}$) and (BTYPE $\in \{$A, B$\}$+WEIGHT $\in \{$UNDER, OVER$\}$). If the instance is changed to violate any of them, the decision will change. For example, if we set BTYPE to AB and WEIGHT to NORM, the third GNR will be violated and the decision on Susan becomes NO. For the classifier in Figure 1b, the GNRs for the decision are different: (AGE $\geq$ 55) and (BTYPE $\in \{$A, O$\}$ + WEIGHT $\in \{$NORM, OVER$\}$). However, both sets of GNRs contain more information than the NRs since the minimal changes they identify to flip the decision include those identified by the NRs.

**Proposition 11.** *Let $\sigma$ be a simple clause. Then $\sigma$ is a NR for the decision on instance $\mathcal{I}$ iff $\sigma = \mathcal{I} \dot{\cap} \sigma'$ for some GNR $\sigma'$.*

Consider the instance Susan again, $\mathcal{I} = (\text{AGE} \geq 55) \cdot (\text{BTYPE=A}) \cdot (\text{WEIGHT=OVER})$ and the classifier in Figure 1b. As mentioned earlier, the GNRs for the decision on Susan are $\sigma_1' = (\text{AGE} \geq 55)$ and $\sigma_2' = (\text{BTYPE} \in \{\text{A, O}\} + \text{WEIGHT} \in \{\text{NORM, OVER}\})$. Then $\sigma_1 = \mathcal{I} \dot{\cap} \sigma_1' = (\text{AGE} \geq 55)$ and $\sigma_2 = \mathcal{I} \dot{\cap} \sigma_2' = (\text{WEIGHT=OVER} + \text{BTYPE=A})$, which are the two NRs for the decision on Susan.

GSRs and GNRs are particularly significant when explaining the decisions of classifiers with numeric features, a topic which we discuss in Appendix C.

We next present a fundamental result which allows us to compute GSRs and GNRs using the general reason for a decision (we use this result in Section 6).

**Definition 9.** *A prime implicant/implicate $c$ of formula $\Delta$ is variable-minimal iff there is no prime implicant/implicate $c'$ of $\Delta$ s.t. $vars(c') \subset vars(c)$.*

**Proposition 12.** *Let $\mathcal{I}$ by an instance in class $\Delta$. The GSRs/GNRs for the decision on instance $\mathcal{I}$ are the variable-minimal prime implicants/implicates of the general reason $\overline{\forall}\mathcal{I} \cdot \Delta$.*

The disjunction of SRs is equivalent to the complete reason which is equivalent to the conjunction of NRs. However, the disjunction of GSRs implies the general reason but is not equivalent to it, and the conjunction of GNRs is implied by the general reason but is not equivalent to it; see Appendix D. This suggests that more information can potentially be extracted from the general reason beyond the information provided by GSRs and GNRs.

## 5    The General Reasons of Decision Graphs

Decision graphs are DAGs which include decision trees [7,9], OBDDs [10], and can have discrete or numeric features. They received significant attention in the work on explainable AI since they can be compiled from other types of classifiers such as Bayesian networks [46], random forests [12] and some types of neural networks [44]. Hence, the ability to explain decision graphs has a direct application to explaining the decisions of a broad class of classifiers. Moreover, the decisions undertaken by decision graphs have closed-form complete reasons as shown in [18]. We provide similar closed forms for the general reasons in this section. We first review decision graphs to formally state our results.

Each leaf node in a decision graph is labeled with some class $c$. An internal node $T$ that *tests* variable $X$ has outgoing edges $\xrightarrow{X, S_1} T_1, \ldots, \xrightarrow{X, S_n} T_n$, $n \geq 2$. The children of node $T$ are $T_1, \ldots, T_n$ and $S_1, \ldots, S_n$ is a partition of *some* states of variable $X$. A decision graph will be represented by its root node. Hence, each node in the graph represents a smaller decision graph. Variables can be tested more than once on a path if they satisfy the *weak test-once property* discussed next [18,22]. Consider a path $\ldots, T \xrightarrow{X, S_j} T_j, \ldots, T' \xrightarrow{X, R_k} T_k, \ldots$ from the root to a leaf (nodes $T$ and $T'$ test $X$). If no nodes between $T$ and $T'$ on the path

test variable $X$, then $\{R_k\}_k$ must be a partition of states $S_j$. Moreover, if $T$ is the first node that tests $X$ on the path, then $\{S_j\}_j$ must be a partition of *all* states for $X$. Discretized numeric variables are normally tested more than once while satisfying the weak test-once property; see Appendix C for an illustration.

**Proposition 13.** *Let $T$ be a decision graph, $\mathcal{I}$ be an instance in class $c$, and $\mathcal{I}[X]$ be the state of variable $X$ in instance $\mathcal{I}$. Suppose $\Delta^c[T]$ is the class formula of $T$ and class $c$. The general reason $\overline{\forall}\, \mathcal{I} \cdot \Delta^c[T]$ is given by the NNF circuit:[8]*

$$\Gamma^c[T] = \begin{cases} \top & \text{if } T \text{ is a leaf with class } c \\ \bot & \text{if } T \text{ is a leaf with class } c' \neq c \\ \prod_j (\Gamma^c[T_j] + \ell) & \text{if } T \text{ has outgoing edges } \xrightarrow{X, S_j} T_j \end{cases}$$

*Here, $\ell$ is the $X$-literal $\{x_i \mid x_i \notin S_j\}$ if $\mathcal{I}[X] \notin S_j$, else $\ell = \bot$.*

The following proposition identifies some properties of the above closed form, which have key computational implications that we exploit in the next section.

**Proposition 14.** *The NNF circuit in Proposition 13 is locally fixated on instance $\mathcal{I}$. Moreover, every disjunction in this circuit has the form $\ell + \Delta$ where $\ell$ is an $X$-literal, and for every $X$-literal $\ell'$ in $\Delta$ we have $\ell' \neq \ell$ and $\ell \models \ell'$.*

## 6   Computing Prime Implicants & Implicates

Computing the prime implicants/implicates of Boolean formulas was studied extensively for decades; see, e.g., [47,29,30]. The classical methods are based on *resolution* when computing the prime implicates of CNFs, and *consensus* when computing the prime implicants of DNFs; see, e.g., [20,15]. More modern approaches are based on passing encodings to SAT-solvers; see, e.g., [40,34,28]. In contrast, the computation of prime implicants/implicates of discrete formulas has received very little attention in the literature. One recent exception is [12] which showed how an algorithm for computing prime implicants of Boolean formulas can be used to compute simple prime implicants of discrete formulas given an appropriate encoding. Computing prime implicants/implicates of NNFs also received relatively little attention; see [41,18,14] for some exceptions. We next provide methods for computing variable-minimal prime implicants/implicates of some classes of discrete formulas that are relevant to GSRs and GNRs.

A set of terms $S$ will be interpreted as a DNF $\sum_{\tau \in S} \tau$ and a set of clauses $S$ will be interpreted as a CNF $\prod_{\sigma \in S} \sigma$. If $S_1$ and $S_2$ are two sets of terms, then $S_1 \times S_2 = \{\tau_1 \cdot \tau_2 \mid \tau_1 \in S_1, \tau_2 \in S_2\}$. For a set of terms/clauses $S$, $\ominus(S)$ denotes the result of removing subsumed terms/clauses from $S$.

---

[8] An NNF circuit is a DAG whose leaves are labeled with $\bot, \top$, or literals; and whose internal nodes are labelled with $\cdot$ or $+$.

---

**Algorithm 1** GSR($\Delta$) — without Line 10, this is **Algorithm 2** PI($\Delta$)

---

**Input:** NNF circuit $\Delta$ which satisfies the properties in Proposition 14
1: **if** CACHE($\Delta$) $\neq$ NIL **then return** CACHE($\Delta$)
2: **else if** $\Delta = \top$ **then return** $\{\top\}$
3: **else if** $\Delta = \bot$ **then return** $\emptyset$
4: **else if** $\Delta$ is a literal **then return** $\{\Delta\}$
5: **else if** $\Delta = \alpha \cdot \beta$ **then**
6:      $S \leftarrow \ominus(\text{GSR}(\alpha) \times \text{GSR}(\beta))$
7: **else if** $\Delta = \alpha + \beta$ **then**
8:      $S \leftarrow \ominus(\text{GSR}(\alpha) \cup \text{GSR}(\beta))$
9: **end if**
10: $S \leftarrow \boxtimes(S, ivars(\Delta))$
11: CACHE($\Delta$) $\leftarrow S$
12: **return** $S$

---

### 6.1   Computing General Sufficient Reasons

Our first result is Algorithm 1 which computes the variable-minimal prime implicants of an NNF circuit that satisfies the properties in Proposition 14 and, hence, is applicable to the general reasons of Proposition 13. If we remove Line 10 from Algorithm 1, it becomes Algorithm 2 which computes all prime implicants instead of only the variable-minimal ones. Algorithm 2 is the same algorithm used to convert an NNF into a DNF (i.e., no consensus is invoked), yet the resulting DNF is guaranteed to be in prime-implicant form. Algorithm 2 is justified by the following two results, where the first result generalizes Proposition 40 in [37].

In the next propositions, pi($\Delta$) denotes the prime implicants of formula $\Delta$.

**Proposition 15.** $pi(\alpha \cdot \beta) = \ominus(pi(\alpha) \times pi(\beta))$.

**Proposition 16.** *For any disjunction $\alpha + \beta$ that satisfies the property of Proposition 14, $pi(\alpha + \beta) = \ominus(pi(\alpha) \cup pi(\beta))$.*

We will next explain Line 10 of Algorithm 1, $S \leftarrow \boxtimes(S, ivars(\Delta))$, which is responsible for pruning prime implicants that are not variable-minimal (hence, computing GSRs). Here, $\Delta$ is a node in the NNF circuit passed in the first call to Algorithm 1, and $ivars(\Delta)$ denotes variables that appear only in the sub-circuit rooted at node $\Delta$. Moreover, $\boxtimes(S, V)$ is the set of terms obtained from terms $S$ by removing every term $\tau \in S$ that satisfies $vars(\tau) \supset vars(\tau')$ and $V \cap (vars(\tau) \setminus vars(\tau')) \neq \emptyset$ for some other term $\tau' \in S$.[9] That is, term $\tau$ will be removed only if some variable $X$ in $vars(\tau) \setminus vars(\tau')$ appears only in the sub-circuit rooted at node $\Delta$ (this ensures that term $\tau$ will not participate in constructing any variable-minimal prime implicant). This incremental pruning technique is enabled by the local fixation property (Definition 4).

**Proposition 17.** *Algorithm 1, GSR($\Delta$), returns the variable-minimal prime implicants of NNF circuit $\Delta$.*

---

[9] The condition $V \cap (vars(\tau) \setminus vars(\tau')) \neq \emptyset$ is trivially satisfied when $\Delta$ is the root of the NNF circuit since $V$ will include all circuit variables in this case.

### 6.2   Computing General Necessary Reasons

We can convert an NNF circuit into a CNF using a dual of Algorithm 2 but the result will not be in prime-implicate form, even for ciruits that satisfy the properties Proposition 14.[10] Hence, we next propose a generalization of the Boolean resolution inference rule to discrete variables, which can be used to convert a CNF into its prime-implicate form. Recall first that Boolean resolution derives the clause $\alpha + \beta$ from the clauses $x + \alpha$ and $\overline{x} + \beta$ where $X$ is a Boolean variable.

**Definition 10.** *Let $\alpha = \ell_1 + \sigma_1$, $\beta = \ell_2 + \sigma_2$ be two clauses where $\ell_1$ and $\ell_2$ are $X$-literals s.t. $\ell_1 \not\models \ell_2$ and $\ell_2 \not\models \ell_1$. If $\sigma = (\ell_1 \cdot \ell_2) + \sigma_1 + \sigma_2 \neq \top$, then the $X$-resolvent of clauses $\alpha$ and $\beta$ is defined as the clause equivalent to $\sigma$.*

We exclude the cases $\ell_1 \models \ell_2$ and $\ell_2 \models \ell_1$ to ensure that the resolvent is not subsumed by clauses $\alpha$ and $\beta$. If $\sigma = \top$, it cannot be represented by clause since a clause is a disjunction of literals over distinct variables so it cannot be trivial.

**Proposition 18.** *Closing a (discrete) CNF under resolution and removing subsumed clauses yields the CNF's prime implicates.*

The following proposition shows that we can incrementally prune clauses that are not variable-minimal after each resolution step. This is significant computationally and is enabled by the property of local fixation (Definition 4) which is satisfied by the general reasons in Proposition 13 and their CNFs.

**Proposition 19.** *Let $S$ be a set of clauses (i.e., CNF) that is locally fixated. For any clauses $\sigma$ and $\sigma'$ in $S$, if $vars(\sigma') \subset vars(\sigma)$, then the variable-minimal prime implicates of $S$ are the variable-minimal prime implicates of $S \setminus \{\sigma\}$.*

In summary, to compute GNRs, we first convert the general reason in Proposition 13 into a CNF, then close the CNF under resolution while removing subsumed clauses and ones that are not variable-minimal after each resolution step.

## 7   Conclusion

We considered the notions of sufficient, necessary and complete reasons which have been playing a fundamental role in explainable AI recently. We provided generalizations of these notions for classifiers with non-binary features (discrete or discretized). We argued that these generalized notions have more explanatory power and reveal more information about the underlying classifier. We further provided results on the properties and computation of these new notions.

### Acknowledgments

---

[10] The number of clauses in this CNF will be no more than the number of NNF nodes if the NNF is the general reason of a decision tree (i.e., the NNF has a tree structure).

# References

1. Albini, E., Rago, A., Baroni, P., Toni, F.: Relation-based counterfactual explanations for bayesian network classifiers. In: IJCAI. pp. 451–457. ijcai.org (2020)
2. Amgoud, L.: Explaining black-box classifiers: Properties and functions. Int. J. Approx. Reason. **155**, 40–65 (2023)
3. Amgoud, L., Ben-Naim, J.: Axiomatic foundations of explainability. In: IJCAI. pp. 636–642. ijcai.org (2022)
4. Audemard, G., Bellart, S., Bounia, L., Koriche, F., Lagniez, J., Marquis, P.: On the explanatory power of Boolean decision trees. Data Knowl. Eng. **142**, 102088 (2022)
5. Audemard, G., Koriche, F., Marquis, P.: On tractable XAI queries based on compiled representations. In: KR. pp. 838–849 (2020)
6. Audemard, G., Lagniez, J., Marquis, P., Szczepanski, N.: Computing abductive explanations for boosted trees. CoRR **abs/2209.07740** (2022)
7. Belson, W.A.: Matching and prediction on the principle of biological classification. Journal of the Royal Statistical Society. Series C (Applied Statistics) **8**(2), 65–75 (1959), http://www.jstor.org/stable/2985543
8. Boumazouza, R., Alili, F.C., Mazure, B., Tabia, K.: ASTERYX: A model-agnostic sat-based approach for symbolic and score-based explanations. In: CIKM. pp. 120–129. ACM (2021)
9. Breiman, L., Friedman, J.H., Olshen, R.A., Stone, C.J.: Classification and Regression Trees. Wadsworth (1984)
10. Bryant, R.E.: Graph-based algorithms for Boolean function manipulation. IEEE Trans. Computers **35**(8), 677–691 (1986)
11. Chan, H., Darwiche, A.: Reasoning about bayesian network classifiers. In: UAI. pp. 107–115. Morgan Kaufmann (2003)
12. Choi, A., Shih, A., Goyanka, A., Darwiche, A.: On symbolically encoding the behavior of random forests. CoRR **abs/2007.01493** (2020)
13. Choi, A., Xue, Y., Darwiche, A.: Same-decision probability: A confidence measure for threshold-based decisions. Int. J. Approx. Reason. **53**(9), 1415–1428 (2012)
14. de Colnet, A., Marquis, P.: On the complexity of enumerating prime implicants from decision-DNNF circuits. In: IJCAI. pp. 2583–2590. ijcai.org (2022)
15. Crama, Y., Hammer, P.L.: Boolean functions - theory, algorithms, and applications. In: Encyclopedia of mathematics and its applications (2011)
16. Darwiche, A.: Logic for explainable AI. In: 38th Annual ACM/IEEE Symposium on Logic in Computer Science, LICS. pp. 1–11. IEEE (2023), CoRR abs/2305.05172
17. Darwiche, A., Hirth, A.: On the reasons behind decisions. In: ECAI. Frontiers in Artificial Intelligence and Applications, vol. 325, pp. 712–720. IOS Press (2020)
18. Darwiche, A., Ji, C.: On the computation of necessary and sufficient explanations. In: AAAI. pp. 5582–5591. AAAI Press (2022)
19. Darwiche, A., Marquis, P.: On quantifying literals in Boolean logic and its applications to explainable AI. J. Artif. Intell. Res. **72**, 285–328 (2021)
20. Gurvich, V., Khachiyan, L.: On generating the irredundant conjunctive and disjunctive normal forms of monotone Boolean functions. Discrete Applied Mathematics **96**, 363–373 (1999)
21. Huang, X., Izza, Y., Ignatiev, A., Cooper, M.C., Asher, N., Marques-Silva, J.: Efficient explanations for knowledge compilation languages. CoRR **abs/2107.01654** (2021)

22. Huang, X., Izza, Y., Ignatiev, A., Marques-Silva, J.: On efficiently explaining graph-based classifiers. In: KR. pp. 356–367 (2021)
23. Ignatiev, A., Izza, Y., Stuckey, P.J., Marques-Silva, J.: Using maxsat for efficient explanations of tree ensembles. In: AAAI. pp. 3776–3785. AAAI Press (2022)
24. Ignatiev, A., Narodytska, N., Asher, N., Marques-Silva, J.: From contrastive to abductive explanations and back again. In: AI*IA. Lecture Notes in Computer Science, vol. 12414, pp. 335–355. Springer (2020)
25. Ignatiev, A., Narodytska, N., Marques-Silva, J.: Abduction-based explanations for machine learning models. In: Proceedings of the Thirty-Third Conference on Artificial Intelligence (AAAI). pp. 1511–1519 (2019)
26. Ignatiev, A., Narodytska, N., Marques-Silva, J.: On validating, repairing and refining heuristic ML explanations. CoRR **abs/1907.02509** (2019)
27. Izza, Y., Ignatiev, A., Marques-Silva, J.: On tackling explanation redundancy in decision trees. J. Artif. Intell. Res. **75**, 261–321 (2022)
28. Izza, Y., Marques-Silva, J.: On explaining random forests with SAT. In: IJCAI. pp. 2584–2591. ijcai.org (2021)
29. Jackson, P.: Computing prime implicates. In: Proceedings of the 1992 ACM Annual Conference on Communications. p. 65–72. CSC '92, Association for Computing Machinery, New York, NY, USA (1992). https://doi.org/10.1145/131214.131223, https://doi.org/10.1145/131214.131223
30. Kean, A., Tsiknis, G.: An incremental method for generating prime implicants/implicates. Journal of Symbolic Computation **9**(2), 185–206 (1990)
31. Lang, J., Liberatore, P., Marquis, P.: Propositional independence: Formula-variable independence and forgetting. J. Artif. Intell. Res. **18**, 391–443 (2003)
32. Lipton, P.: Contrastive explanation. Royal Institute of Philosophy Supplements **27**, 247–266 (1990). https://doi.org/10.1017/S1358246100005130
33. Liu, X., Lorini, E.: A unified logical framework for explanations in classifier systems. J. Log. Comput. **33**(2), 485–515 (2023)
34. Luo, W., Want, H., Zhong, H., Wei, O., Fang, B., Song, X.: An efficient two-phase method for prime compilation of non-clausal boolean formulae. In: 2021 IEEE/ACM International Conference On Computer Aided Design (ICCAD). pp. 1–9 (2021). https://doi.org/10.1109/ICCAD51958.2021.9643520
35. Marques-Silva, J., Gerspacher, T., Cooper, M.C., Ignatiev, A., Narodytska, N.: Explanations for monotonic classifiers. In: ICML. Proceedings of Machine Learning Research, vol. 139, pp. 7469–7479. PMLR (2021)
36. Marques-Silva, J., Ignatiev, A.: Delivering trustworthy AI through formal XAI. In: AAAI. pp. 12342–12350. AAAI Press (2022)
37. Marquis, P.: Consequence finding algorithms. In: Handbook of defeasible reasoning and uncertainty management systems, pp. 41–145. Springer (2000)
38. Miller, T.: Explanation in artificial intelligence: Insights from the social sciences. Artif. Intell. **267**, 1–38 (2019)
39. Narodytska, N., Kasiviswanathan, S.P., Ryzhyk, L., Sagiv, M., Walsh, T.: Verifying properties of binarized deep neural networks. In: Proc. of AAAI'18. pp. 6615–6624 (2018)
40. Previti, A., Ignatiev, A., Morgado, A., Marques-Silva, J.: Prime compilation of non-clausal formulae. In: IJCAI. pp. 1980–1988. AAAI Press (2015)
41. Ramesh, A., Becker, G., Murray, N.V.: CNF and DNF considered harmful for computing prime implicants/implicates. Journal of Automated Reasoning **18**(3), 337–356 (1997)
42. Ribeiro, M.T., Singh, S., Guestrin, C.: "Why should I trust you?": Explaining the predictions of any classifier. In: KDD. pp. 1135–1144. ACM (2016)

43. Ribeiro, M.T., Singh, S., Guestrin, C.: Anchors: High-precision model-agnostic explanations. In: AAAI. pp. 1527–1535. AAAI Press (2018)
44. Shi, W., Shih, A., Darwiche, A., Choi, A.: On tractable representations of binary neural networks. In: KR. pp. 882–892 (2020)
45. Shih, A., Choi, A., Darwiche, A.: A symbolic approach to explaining bayesian network classifiers. In: IJCAI. pp. 5103–5111. ijcai.org (2018)
46. Shih, A., Choi, A., Darwiche, A.: Compiling bayesian network classifiers into decision graphs. In: AAAI. pp. 7966–7974. AAAI Press (2019)
47. Slagle, J., Chang, C.L., Lee, R.: A new algorithm for generating prime implicants. IEEE Transactions on Computers **C-19**(4), 304–310 (1970). https://doi.org/10.1109/T-C.1970.222917
48. Wang, E., Khosravi, P., den Broeck, G.V.: Probabilistic sufficient explanations. In: IJCAI. pp. 3082–3088. ijcai.org (2021)

## A    Proofs

### Proposition 4

We prove Proposition 4 and 5 first, which do not depend on Propositions 1, 2, and 3.

*Proof (of Proposition 4).* The proof will use the following observations:

$$\text{(A)} \quad \omega \models x_i \cdot \Delta \text{ only if } \omega \models \Delta|x_i$$
$$\text{(B)} \quad \omega \models x_i \cdot (\Delta|x_i) \text{ only if } \omega \models \Delta$$

(A) is justified as follows: $\omega \models \Delta$ iff $\Delta|\omega = \top$, and $(\Delta|x_i)|\omega = \Delta|\omega$ since $\omega \models x_i$; thus, $(\Delta|x_i)|\omega = \top$ and $\omega \models \Delta|x_i$. (B) is justified as follows: $\omega \models \Delta|x_i$ implies $(\Delta|x_i)|\omega = \top$ and $\omega \models x_i$ implies $\Delta|\omega = (\Delta|x_i)|\omega$; hence, $\Delta|\omega = \top$ and $\omega \models \Delta$.

We next prove both directions of the equivalence while noting that $\ell_j$ in the proposition statement is equivalent to $\overline{x_j}$.

$\overline{\forall}\, x_i \cdot \Delta \models (\Delta|x_i) \cdot \prod_{j\neq i}(\overline{x_j} + (\Delta|x_j))$. Suppose $\omega \models \overline{\forall}\, x_i \cdot \Delta$. Then $\omega \models \Delta \cdot (\Delta|x_i)$ by Definition 2. If $\omega \models x_k$ for some $k$, then (1) $\omega \models (\overline{x_j} + \Delta|x_j)$ for all $j \neq k$ since $x_k \models \overline{x_j}$ and (2) $\omega \models (\overline{x_k} + \Delta|x_k)$ since $\omega \models \Delta|x_k$ which follows from $\omega \models \Delta$ and $\omega \models x_k$ by (A). Hence, $\omega \models (\overline{x_k} + \Delta|x_k)$ for all $k$ and, therefore, $\omega \models \prod_{j\neq i}(\overline{x_j}+(\Delta|x_j))$ and $\omega \models (\Delta|x_i)\cdot\prod_{j\neq i}(\overline{x_j}+(\Delta|x_j))$. Hence, $\overline{\forall}\, x_i \cdot \Delta \models (\Delta|x_i) \cdot \prod_{j\neq i}(\overline{x_j} + (\Delta|x_j))$.

$(\Delta|x_i) \cdot \prod_{j\neq i}(\overline{x_j} + (\Delta|x_j)) \models \overline{\forall}\, x_i \cdot \Delta$. Suppose $\omega \models (\Delta|x_i) \cdot \prod_{j\neq i}(\overline{x_j} + (\Delta|x_j))$. If $\omega \models x_i$, then $\omega \models \Delta$ since $\omega \models \Delta|x_i$ and given (B). If $\omega \not\models x_i$, then $\omega \models x_k$ for some $k \neq i$, and $\omega \models \Delta|x_k$ since $\omega \models (\overline{x_k}+(\Delta|x_k))$, which implies $\omega \models \Delta$ given (B). Hence, $\omega \models \Delta$ in either case and also $\omega \models \Delta \cdot \Delta|x_i = \overline{\forall}\, x_i \cdot \Delta$. Therefore, $(\Delta|x_i) \cdot \prod_{j\neq i}(\overline{x_j} + (\Delta|x_j)) \models \Delta \cdot \Delta|x_i$.

### Proposition 5

*Proof (of Proposition 5).* We first prove the semantics of $\forall\tau \cdot \Delta$ and then $\overline{\forall}\,\tau \cdot \Delta$ by induction on the length of simple term $\tau$.

*Semantics of $\forall\tau \cdot \Delta$.*

Base case: $\tau = x_i$.

By definition of $\forall$, a world $\omega \models \forall x_i \cdot \Delta$ iff $\omega \models (\Delta|x_i) \cdot \prod_{j\neq i}(x_i + \Delta|x_j)$. If $\omega \models x_i$, then $\omega \models \forall x_i \cdot \Delta$ iff $\omega \models \Delta$ by observations (A) and (B) in the proof of Proposition 4. If $\omega \not\models x_i$, then $\omega \models \forall x_i \cdot \Delta$ iff $\omega \models \Delta|x_j$ for all $j$. Hence, $\omega \models \forall x_i \cdot \Delta$ iff $\omega \models \Delta$ and ($\omega \not\models x_i$ only if $\omega \models \Delta|x_j$ for all $j$). The condition "$\omega \models \Delta|x_j$ for all $j$" is equivalent to "$\omega' \models \Delta$ for $\omega'$ obtained from $\omega$ by changing its state $x_j$ if $x_j \neq x_i$." If $\omega \models x_i$, the previous property holds trivially as there is no such $\omega'$. Hence, $\omega \models \forall x_i \cdot \Delta$ iff $\omega \models \Delta$ and $\omega' \models \Delta$ for $\omega'$ obtained from $\omega$ by changing its state $x_j$ to any other state if $x_j \neq x_i$. The semantics of $\forall x_i \cdot \Delta$ holds.

Inductive step: $\tau = x_i \cdot \tau'$.

Suppose the proposition holds for $\forall \tau' \cdot \Delta$. We next show that it holds for $\forall \tau \cdot \Delta$. Let $\Gamma = \forall \tau' \cdot \Delta$. By the base case, $\omega \models \forall x_i \cdot \Gamma$ iff (1) $\omega \models \Gamma$ and (2) $\omega' \models \Gamma$ for $\omega'$ obtained from $\omega$ by changing its state $x_j$ if $x_j \neq x_i$. By the induction hypothesis, (1) can be replaced by "$\omega \models \Delta$ and $\omega' \models \Delta$ for $\omega'$ obtained from $\omega$ by changing the states of variables set differently in $\tau'$." Moreover, (2) can be replaced by "$\omega' \models \Delta$ and $\omega'' \models \Delta$ for $\omega'$ obtained from $\omega$ by changing its state $x_j$ if $x_j \neq x_i$ and for $\omega''$ obtained from $\omega'$ by changing the states of variables set differently in $\tau'$." Replacing (1), (2) as suggested above gives: $\omega \models \forall \tau \cdot \Delta$ iff $\omega \models \Delta$ and $\omega' \models \Delta$ for $\omega'$ obtained from $\omega$ by changing the states of variables set differently in $\tau$. The semantics of $\forall \tau \cdot \Delta$ holds.

*Semantics of $\overline{\forall} \tau \cdot \Delta$.*

Base case $\tau = x_i$.

By definition of $\overline{\forall}$, $\omega \models \overline{\forall} x_i \cdot \Delta$ iff $\omega \models \Delta \cdot (\Delta | x_i)$. We next prove: if $\omega \models \Delta$, then $\omega \models \Delta | x_i$ is equivalent to "$\omega' \models \Delta$ for $\omega'$ obtained by setting $X$ to $x_i$ in $\omega$," which proves the semantics of $\overline{\forall} x_i \cdot \Delta$. Suppose $\omega \models \Delta$. We next show both directions of the equivalence.

Suppose $\omega \models \Delta | x_i$ and let $\omega'$ be a world obtained by setting variable $X$ to $x_i$ in world $\omega$. Then $\omega' \models \Delta | x_i$ given $\omega \models \Delta | x_i$ and since $\Delta | x_i$ does not mention variable $X$. Hence, $\omega' \models \Delta$ by observation (B) in the proof of Proposition 4.

Suppose $\omega' \models \Delta$ for $\omega'$ obtained by setting $X$ to $x_i$ in $\omega$. Then $\omega' \models \Delta | x_i$ by observation (A) in the proof of Proposition 4. Moreover, $\omega \models \Delta | x_i$ given $\omega' \models \Delta | x_i$ and since $\Delta | x_i$ does not mention variable $X$.

This proves the semantics of $\overline{\forall} x_i \cdot \Delta$.

Inductive step: $\tau = x_i \cdot \tau'$.

Suppose the proposition holds for $\overline{\forall} \tau' \cdot \Delta$. We next show that it holds for $\overline{\forall} \tau \cdot \Delta$. Let $\Gamma = \overline{\forall} \tau' \cdot \Delta$. By the base case, $\omega \models \overline{\forall} x_i \cdot \Gamma$ iff (1) $\omega \models \Gamma$ and (2) $\omega' \models \Gamma$ for $\omega'$ obtained from $\omega$ by setting $X$ to $x_i$. By the induction hypothesis, (1) can be replaced by "$\omega \models \Delta$ and $\omega' \models \Delta$ for $\omega'$ obtained from $\omega$ by setting some variables to their states in $\tau'$." Moreover, (2) can be replaced by "$\omega' \models \Delta$ and $\omega'' \models \Delta$ for $\omega'$ obtained from $\omega$ by setting $X$ to $x_i$ and $\omega''$ obtained from $\omega'$ by setting some variables to their states in $\tau'$." Replacing (1), (2) as suggested above gives: $\overline{\forall} \tau \cdot \Delta$ iff $\omega \models \Delta$ and $\omega' \models \Delta$ for $\omega'$ obtained from $\omega$ by setting some variables to their states in $\tau$, which proves the semantics of $\overline{\forall} \tau \cdot \Delta$.

## Proposition 1

*Proof (of Proposition 1).* We have:

$$
\begin{aligned}
\overline{\forall}\, y \cdot (\overline{\forall}\, x \cdot \Delta) &= \overline{\forall}\, y \cdot (\Delta \cdot \Delta|x) \\
&= (\Delta \cdot \Delta|x) \cdot ((\Delta \cdot \Delta|x)|y) \\
&= (\Delta \cdot \Delta|x) \cdot (\Delta|y \cdot \Delta|x, y) \\
&= (\Delta \cdot \Delta|y) \cdot (\Delta|x \cdot \Delta|y, x) \\
&= (\Delta \cdot \Delta|y) \cdot ((\Delta \cdot \Delta|y)|x) \\
&= (\overline{\forall}\, y \cdot \Delta) \cdot ((\overline{\forall}\, y \cdot \Delta)|x) \\
&= \overline{\forall}\, x \cdot (\overline{\forall}\, y \cdot \Delta).
\end{aligned}
$$

## Proposition 2

*Proof (of Proposition 2).* It suffices to prove that $\forall \mathcal{I} \cdot \Delta \models \overline{\forall} \mathcal{I} \cdot \Delta \models \Delta$ since $\mathcal{I} \models \forall \mathcal{I} \cdot \Delta$ [19]. By multiple applications of Definition 2, $\overline{\forall} \mathcal{I} \cdot \Delta$ is equivalent to $\Delta \cdot \Gamma$ for some formula $\Gamma$. Thus, $\overline{\forall} \mathcal{I} \cdot \Delta \models \Delta$. Moreover, $\forall \mathcal{I} \cdot \Delta \models \overline{\forall} \mathcal{I} \cdot \Delta$ by Proposition 5 (already proven).

## Proposition 3

*Proof (of Proposition 3).* By Proposition 4 (already proven), $\overline{\forall} x_i \cdot \Delta$ can be written as an NNF over formulas that either do not mention variable $X$ or are $X$-literals implied by $x_i$. Hence, $\overline{\forall} \mathcal{I} \cdot \Delta$ can always be written as an NNF whose literals are implied by instance $\mathcal{I}$ (by repeated application of the previous observation). We also have $\mathcal{I} \models \overline{\forall} \mathcal{I} \cdot \Delta \models \Delta$ by Proposition 2. Hence, it suffices to show that if $\Gamma$ is an NNF such that (1) $\mathcal{I}$ satisfies the literals of $\Gamma$ and (2) $\mathcal{I} \models \Gamma \models \Delta$, then $\Gamma \models \overline{\forall} \mathcal{I} \cdot \Delta$. We next prove this by contradiction. Suppose $\Gamma$ is an NNF that satisfies properties (1) and (2), and $\Gamma \not\models \overline{\forall} \mathcal{I} \cdot \Delta$. Then $\omega \models \Gamma$ and $\omega \not\models \overline{\forall} \mathcal{I} \cdot \Delta$ for some world $\omega$. Let $\omega'$ be a world obtained from $\omega$ by setting some variables in $\omega$ to their states in $\mathcal{I}$. Then $\omega' \models \Gamma$ since $\mathcal{I}$ satisfies all literals of $\Gamma$ by (1), and $\mathcal{I} \models \Gamma$ by (2). We now have $\omega \models \Gamma \models \Delta$ by (2), and $\omega' \models \Gamma \models \Delta$ for all such worlds $\omega'$, which implies $\omega \models \overline{\forall} \mathcal{I} \cdot \Delta$ by Proposition 5 (already proven). This is a contradiction so $\Gamma \models \overline{\forall} \mathcal{I} \cdot \Delta$.

## Proposition 6

*Proof (of Proposition 6).* If $x \models \ell$, then $\overline{\forall} x \cdot \ell = \ell \cdot \ell|x = \ell \cdot \top = \ell$. If $x \not\models \ell$, then $\overline{\forall} x \cdot \ell = \ell \cdot \ell|x = \ell \cdot \bot = \bot$. If $X$ does not appear in $\Delta$, then $\overline{\forall} x \cdot \Delta = \Delta \cdot \Delta|x = \Delta \cdot \Delta = \Delta$.

## Proposition 7

*Proof (of Proposition 7).* For the distribution over conjuncts,

$$\overline{\forall}\, x_i \cdot (\alpha \cdot \beta) = (\alpha \cdot \beta) \cdot (\alpha \cdot \beta)|x_i$$
$$= \alpha \cdot \beta \cdot \alpha|x_i \cdot \beta|x_i$$
$$= (\alpha \cdot \alpha|x_i) \cdot (\beta \cdot \beta|x_i)$$
$$= (\overline{\forall}\, x_i \cdot \alpha) \cdot (\overline{\forall}\, x_i \cdot \beta).$$

For the distribution over disjuncts, suppose variable $X$ does not occur in $\alpha$. Then $\overline{\forall}\, x_i \cdot \alpha = \alpha$ by Proposition 6. Moreover,

$$\overline{\forall}\, x_i \cdot (\alpha + \beta) = (\alpha + \beta) \cdot (\alpha + \beta)|x_i$$
$$= (\alpha + \beta) \cdot (\alpha|x_i + \beta|x_i)$$
$$= (\alpha + \beta) \cdot (\alpha + \beta|x_i)$$
$$= \alpha + (\alpha \cdot (\beta|x_i)) + (\alpha \cdot \beta) + (\beta \cdot (\beta|x_i))$$
$$= \alpha + (\beta \cdot (\beta|x_i))$$
$$= \overline{\forall}\, x_i \cdot \alpha + \overline{\forall}\, x_i \cdot \beta.$$

The proof is symmetric when $X$ does not occur in $\beta$.

## Proposition 9

We prove Proposition 9 first, which does not depend on Proposition 8.

**Lemma 1.** *Let $\mathcal{I}$ be an instance, $\tau$ be a simple term and $\tau'$ be a GSR for the decision on $\mathcal{I}$. Then $\tau = \mathcal{I} \,\dot{\cap}\, \tau'$ iff $\mathcal{I} \models \tau \models \tau'$ and $vars(\tau) = vars(\tau')$.*

*Proof.* Let $\mathcal{I}$ be an instance, $\tau$ be a simple term and $\tau'$ be a GSR for the decision on $\mathcal{I}$. We next prove both directions of the equivalence.

$\tau = \mathcal{I} \,\dot{\cap}\, \tau'$ only if $\mathcal{I} \models \tau \models \tau'$ and $vars(\tau) = vars(\tau')$.
  Suppose $\tau = \mathcal{I} \,\dot{\cap}\, \tau'$. Recall that $\mathcal{I} \,\dot{\cap}\, \tau'$ denotes the smallest subterm in $\mathcal{I}$ that implies $\tau'$. Hence, $\mathcal{I} \models \mathcal{I} \,\dot{\cap}\, \tau' \models \tau'$ and $\mathcal{I} \models \tau \models \tau'$. Moreover, $vars(\mathcal{I} \,\dot{\cap}\, \tau') = vars(\tau')$ by definition of $\dot{\cap}$ so $vars(\tau) = vars(\tau')$.
$\mathcal{I} \models \tau \models \tau'$ and $vars(\tau) = vars(\tau')$ only if $\tau = \mathcal{I} \,\dot{\cap}\, \tau'$.
  Suppose $\mathcal{I} \models \tau \models \tau'$ and $vars(\tau) = vars(\tau')$. Since $\mathcal{I} \models \tau \models \tau'$ and term $\tau$ is simple, then (1) $\tau$ is a subterm in $\mathcal{I}$ and (2) $\tau$ implies $\tau'$. It then suffices to show that no strict subset of $\tau$ satisfies (1) and (2). Since $vars(\tau) = vars(\tau')$, and $vars(\mathcal{I}\dot{\cap}\tau') = vars(\tau')$ by definition of $\dot{\cap}$, we get $vars(\tau) = vars(\mathcal{I}\dot{\cap}\tau')$. Hence, no strict subset of $\tau$ satisfies (1) and (2), so $\tau = \mathcal{I} \,\dot{\cap}\, \tau'$.

*Proof (of Proposition 9).* Let $\mathcal{I}$ be an instance and $\tau$ be a simple term. Given Lemma 1, it suffices to show that $\tau$ is a SR iff $\mathcal{I} \models \tau \models \tau'$ and $vars(\tau) = vars(\tau')$ for some GSR $\tau'$. Recall that $\tau$ is a SR iff (1) $\mathcal{I} \models \tau \models \Delta$ and (2) $\tau \models \tau'' \models \Delta$ for simple term $\tau''$ only if $\tau = \tau''$. We next prove both directions of equivalence.

$\tau$ is a SR only if $\mathcal{I} \models \tau \models \tau'$ and $vars(\tau) = vars(\tau')$ for some GSR $\tau'$.

Suppose $\tau$ is a SR. Then $\mathcal{I} \models \tau \models \Delta$ and $\tau \models \tau'' \models \Delta$ for simple term $\tau''$ only if $\tau = \tau''$. If $\tau$ is a GSR, then $\mathcal{I} \models \tau \models \tau$ and $vars(\tau) = vars(\tau)$ so the result holds trivially. Suppose $\tau$ is not a GSR. By definition of a GSR and $\mathcal{I} \models \tau \models \Delta$, there must exist a GSR $\tau'$ such that $\tau \models \tau' \models \Delta$ and $\tau' \neq \tau$. Since $\tau \models \tau'$, $vars(\tau') \subseteq vars(\tau)$. Moreoever, by Proposition 8, $\mathcal{I} \dot{\cap} \tau' \models \Delta$ and $\mathcal{I} \dot{\cap} \tau'$ is a simple term. Therefore, $vars(\mathcal{I} \dot{\cap} \tau') = vars(\tau') \subseteq vars(\tau)$. Since $vars(\mathcal{I} \dot{\cap} \tau') \subseteq vars(\tau)$ and both $\mathcal{I} \dot{\cap} \tau'$ and $\tau$ are simple terms implied by $\mathcal{I}$, we get $\tau \models \mathcal{I} \dot{\cap} \tau' \models \Delta$. Since $\tau$ is a SR, we now have $\tau = \mathcal{I} \dot{\cap} \tau'$, so $vars(\tau') = vars(\tau)$. Hence, $\mathcal{I} \models \tau \models \tau'$ and $vars(\tau) = vars(\tau')$ for GSR $\tau'$.

$\mathcal{I} \models \tau \models \tau'$ and $vars(\tau) = vars(\tau')$ for some GSR $\tau'$ only if $\tau$ is a SR.

Suppose $\mathcal{I} \models \tau \models \tau'$ and $vars(\tau) = vars(\tau')$ for some GSR $\tau'$. By definition of a GSR, $\tau' \models \Delta$ and, hence, (1) $\mathcal{I} \models \tau \models \tau' \models \Delta$. Suppose now that $\tau \models \tau'' \models \Delta$ and $\tau \neq \tau''$ for some simple term $\tau''$. We will next show a contradiction which implies (2) $\tau \models \tau'' \models \Delta$ only if $\tau = \tau''$ for any simple term $\tau''$. Let $\tau''$ be the weakest simple term satisfying our supposition. We then have $vars(\tau'') \subset vars(\tau)$. Moreover, $\tau''$ must be a SR. By the first direction, there exists a GSR $\tau'''$ where $vars(\tau''') = vars(\tau'') \subset vars(\tau) = vars(\tau')$. Hence $\tau'$ is not variable-minimal (compared to $\tau'''$) so it cannot be a GSR, a contradiction. Hence, (2) holds. Given (1) and (2), $\tau$ is a SR.

## Proposition 8

*Proof (of Proposition 8).* Suppose $\mathcal{I}$ is an instance in class $\Delta$ and $\tau$ is a weakest term s.t. $\mathcal{I} \models \tau \models \Delta$. We next prove both parts of the proposition.

*Part* 1. Suppose $\tau'$ is a weakest term s.t. $\mathcal{I} \models \tau' \models \Delta$ and $vars(\tau') \subset vars(\tau)$. We will next show $\mathcal{I} \dot{\cap} \tau \models \mathcal{I} \dot{\cap} \tau' \models \Delta$. Since $\mathcal{I} \models \tau$ and $\mathcal{I} \models \tau'$, then $\mathcal{I} \models \ell$ for every literal $\ell$ in $\tau$ or $\tau'$. Hence, $\mathcal{I} \dot{\cap} \tau$ is the subset $\mathcal{J}$ of $\mathcal{I}$ such that $vars(\mathcal{J}) = vars(\tau)$ and $\mathcal{I} \dot{\cap} \tau'$ is the subset $\mathcal{J}'$ of $\mathcal{I}$ such that $vars(\mathcal{J}') = vars(\tau')$. Since $vars(\tau') \subset vars(\tau)$, $vars(\mathcal{J}') \subset vars(\mathcal{J})$ and, hence, $\mathcal{I} \dot{\cap} \tau = \mathcal{J} \models \mathcal{J}' = \mathcal{I} \dot{\cap} \tau'$. Moreover, since $\mathcal{I} \dot{\cap} \tau' \models \tau'$ and $\tau' \models \Delta$, we get $\mathcal{I} \dot{\cap} \tau \models \mathcal{I} \dot{\cap} \tau' \models \Delta$.

*Part* 2(*a*). Suppose $\tau'$ is a weakest term s.t. $\mathcal{I} \models \tau' \models \Delta$ and $vars(\tau') \subset vars(\tau)$. By Part 1, $\mathcal{I} \models \mathcal{I} \dot{\cap} \tau \models \mathcal{I} \dot{\cap} \tau' \models \Delta$. Hence, $\mathcal{I} \dot{\cap} \tau$ is not a SR since $\mathcal{I} \dot{\cap} \tau'$ is weaker than $\mathcal{I} \dot{\cap} \tau$, yet $\mathcal{I} \models \mathcal{I} \dot{\cap} \tau' \models \Delta$.

*Part* 2(*b*). Suppose there is no weakest term $\tau'$ s.t. $\mathcal{I} \models \tau' \models \Delta$ and $vars(\tau') \subset vars(\tau)$. Then $\tau$ is a GSR. By Proposition 9, $\mathcal{I} \dot{\cap} \tau$ is a SR.

## Proposition 11

We prove Proposition 11 first, which does not depend on Proposition 10.

**Lemma 2.** *Let $\mathcal{I}$ be an instance, $\sigma$ be a simple clause, and $\sigma'$ be a GNR for the decision on $\mathcal{I}$. Then $\sigma = \mathcal{I} \dot{\cap} \sigma'$ iff $\mathcal{I} \dot{\models} \sigma \models \sigma'$ and $vars(\sigma) = vars(\sigma')$.*

*Proof.* Suppose $\mathcal{I}$ is an instance, $\sigma$ is a simple clause, and $\sigma'$ is a GNR for the decision on $\mathcal{I}$. We next prove both directions of the equivalence.

$\sigma = \mathcal{I} \mathbin{\dot\cap} \sigma'$ only if $\mathcal{I} \models \sigma \models \sigma'$ and $vars(\sigma) = vars(\sigma')$.

Suppose $\sigma = \mathcal{I} \mathbin{\dot\cap} \sigma'$. Recall that $\mathcal{I} \mathbin{\dot\cap} \sigma'$ denotes the disjunction of states that appear in both $\mathcal{I}$ and $\sigma'$. Therefore, we have $\mathcal{I} \models \mathcal{I} \mathbin{\dot\cap} \sigma'$ and $\mathcal{I} \mathbin{\dot\cap} \sigma' \models \sigma'$, which implies $\mathcal{I} \models \sigma \models \sigma'$. Since $\sigma'$ is a GNR, we have $\mathcal{I} \models \sigma'$. Therefore, $vars(\mathcal{I} \mathbin{\dot\cap} \sigma') = vars(\sigma')$, so $vars(\sigma) = vars(\sigma')$.

$\mathcal{I} \models \sigma \models \sigma'$ and $vars(\sigma) = vars(\sigma')$ only if $\sigma = \mathcal{I} \mathbin{\dot\cap} \sigma'$.

Suppose $\mathcal{I} \models \sigma \models \sigma'$ and $vars(\sigma) = vars(\sigma')$. Since $\sigma$ is a simple clause and $\mathcal{I} \models \sigma$, $\sigma$ is a disjunction of some states $S$ in $\mathcal{I}$. By definition of $\dot\cap$, $\mathcal{I} \mathbin{\dot\cap} \sigma'$ is a disjunction of some states $S'$ in $\mathcal{I}$. Since $\sigma \models \sigma'$, $S \subseteq S'$. Since $vars(\sigma) = vars(\sigma')$, $S = S'$. Hence, $\sigma = \mathcal{I} \mathbin{\dot\cap} \sigma'$.

*Proof (of Proposition 11).* Let instance $\mathcal{I}$ be in class $\Delta$ and $\sigma$ be a simple clause. By Lemma 2, it suffices to show that $\sigma$ is a NR iff $\mathcal{I} \models \sigma \models \sigma^\star$ and $vars(\sigma) = vars(\sigma^\star)$ for some GNR $\sigma^\star$. Recall that $\sigma$ is a NR for the decision on $\mathcal{I}$ iff (1) $\mathcal{I} \models \sigma$ and $(\mathcal{I} \backslash\!\backslash \sigma) \cdot \overline{\sigma} \not\models \Delta$ and (2) a simple clause $\sigma'$ satisfies the previous condition and $\sigma' \models \sigma$ only if $\sigma = \sigma'$. We will reference (1) and (2) next as we prove both directions of the equivalence in Proposition 11.

$\sigma$ is a NR only if $\mathcal{I} \models \sigma \models \sigma^\star$ and $vars(\sigma) = vars(\sigma^\star)$ for some GNR $\sigma^\star$.

Suppose $\sigma$ is a NR. We prove this direction by finding a GNR $\sigma^\star$ that satisfies the properties above. Given (1) and (2), there is no simple clause $\sigma'$ s.t. $\mathcal{I} \models \sigma'$, $(\mathcal{I} \backslash\!\backslash \sigma') \cdot \overline{\sigma'} \not\models \Delta$, $\sigma' \models \sigma$ and $\sigma' \neq \sigma$ (equivalent to $vars(\sigma') \subset vars(\sigma)$). Hence, there is no GNR $\sigma''$ such that $vars(\sigma'') \subset vars(\sigma)$. Next, since $(\mathcal{I} \backslash\!\backslash \sigma) \cdot \overline{\sigma} \not\models \Delta$, there is a world $\omega$ such that $\omega \models (\mathcal{I} \backslash\!\backslash \sigma) \cdot \overline{\sigma}$ and $\omega \not\models \Delta$. Our goal is to construct a clause $\sigma'$ such that the only world that satisfies $(\mathcal{I} \backslash\!\backslash \sigma') \cdot \overline{\sigma'}$ is $\omega$. This gives us $(\mathcal{I} \backslash\!\backslash \sigma') \cdot \overline{\sigma'} \models \overline{\Delta}$. Then, either $\sigma'$ is a GNR, or $\sigma'$ is subsumed by some GNR. If we can find such a clause $\sigma'$, we can also find the sought GNR $\sigma^\star$ that finishes this direction of the proof. This is shown next. Consider the clause $\sigma'$ equivalent to $\overline{\omega \setminus (\omega \backslash\!\backslash \sigma)}$. Note that $\omega \models (\mathcal{I} \backslash\!\backslash \sigma)$ and $\omega \models \overline{\sigma}$ where $\mathcal{I} \models \sigma$. It follows that $\sigma'$ is equivalent to the negation of a conjunction of literals in $\omega$ whose variables are mentioned by $\sigma$. Every $X$-literal $\ell'$ in $\sigma'$ is entailed by an $X$-literal $\ell$ in $\sigma$ because $\ell'$ contains all states of $X$ but the one from $\omega$ and $\omega \models \overline{\sigma}$. Thus, $\mathcal{I} \models \sigma'$. Then $(\mathcal{I} \backslash\!\backslash \sigma') \cdot \overline{\sigma'} = (\mathcal{I} \backslash\!\backslash \sigma) \cdot \overline{\sigma'}$, and the only model of $(\mathcal{I} \backslash\!\backslash \sigma) \cdot \overline{\sigma'}$ is $\omega$. Therefore, $(\mathcal{I} \backslash\!\backslash \sigma') \cdot \overline{\sigma'} \models \overline{\Delta}$. Note that $vars(\sigma') = vars(\sigma)$, and we already showed that there is no GNR $\sigma''$ such that $vars(\sigma'') \subset vars(\sigma)$. Thus, either $\sigma'$ is a GNR, in which case we let $\sigma^\star = \sigma'$, or $\mathcal{I} \models \sigma^\star \models \sigma'$ and $vars(\sigma^\star) = vars(\sigma') = vars(\sigma)$ for some GNR $\sigma^\star$. Either way, $\sigma \models \sigma^\star$ follows from $vars(\sigma) = vars(\sigma^\star)$, $\mathcal{I} \models \sigma$ and $\mathcal{I} \models \sigma^\star$.

$\mathcal{I} \models \sigma \models \sigma'$ and $vars(\sigma) = vars(\sigma')$ for some GNR $\sigma'$ only if $\sigma$ is a NR.

Suppose $\mathcal{I} \models \sigma \models \sigma'$ and $vars(\sigma) = vars(\sigma')$ for some GNR $\sigma'$. We next prove (1) and then prove (2). Since $\sigma'$ is a GNR, $(\mathcal{I} \backslash\!\backslash \sigma') \cdot \overline{\sigma'} \models \overline{\Delta}$. Moreover, $\mathcal{I} \backslash\!\backslash \sigma = \mathcal{I} \backslash\!\backslash \sigma'$ since $vars(\sigma) = vars(\sigma')$, $\mathcal{I} \models \sigma$ and $\mathcal{I} \models \sigma'$. We also have $\overline{\sigma'} \models \overline{\sigma}$, given $\sigma \models \sigma'$, which implies $(\mathcal{I} \backslash\!\backslash \sigma') \cdot \overline{\sigma'} \models (\mathcal{I} \backslash\!\backslash \sigma) \cdot \overline{\sigma}$. Therefore, $\omega \models (\mathcal{I} \backslash\!\backslash \sigma') \cdot \overline{\sigma'} \models \overline{\Delta}$ only if $\omega \models (\mathcal{I} \backslash\!\backslash \sigma) \cdot \overline{\sigma}$ and $\omega \models \overline{\Delta}$, which implies $(\mathcal{I} \backslash\!\backslash \sigma) \cdot \overline{\sigma} \not\models \Delta$. Together with $\mathcal{I} \models \sigma$, this gives (1). We prove (2) by contradiction. Suppose (2) does not hold. Then there exists a NR $\sigma''$ such that $\sigma'' \models \sigma$ and $\sigma'' \neq \sigma$,

so $vars(\sigma'') \subset vars(\sigma)$. By the first direction, there exists a GNR $\sigma'''$ such that $vars(\sigma''') = vars(\sigma'') \subset vars(\sigma) = vars(\sigma')$. This is a contradiction since $\sigma'$ is a GNR. Thus, (2) holds.

## Proposition 10

*Proof (of Proposition 10).* Let $\mathcal{I}$ be an instance in class $\Delta$ and $\sigma$ be a strongest clause s.t. $\mathcal{I} \models \sigma$ and $(\mathcal{I} \backslash\!\backslash \sigma) \cdot \overline{\sigma} \models \overline{\Delta}$. We next prove both parts of the proposition.

*Part 1.* Suppose $\sigma'$ is a strongest clause s.t. $\mathcal{I} \models \sigma'$ and $(\mathcal{I} \backslash\!\backslash \sigma') \cdot \overline{\sigma'} \models \overline{\Delta}$ and $vars(\sigma') \subset vars(\sigma)$. We next show that $\mathcal{I} \backslash\!\backslash \sigma' \models \mathcal{I} \backslash\!\backslash \sigma$. Since $\mathcal{I} \models \sigma$ and $\mathcal{I} \models \sigma'$, every literal $\ell$ in $\sigma$ or $\sigma'$ satisfies $\mathcal{I} \models \ell$. Therefore, $\mathcal{I} \backslash\!\backslash \sigma$ is the subset $\mathcal{J}$ of $\mathcal{I}$ such that $vars(\mathcal{J}) = vars(\mathcal{I}) \setminus vars(\sigma)$ and $\mathcal{I} \dot{\cap} \sigma'$ is the subset $\mathcal{J}'$ of $\mathcal{I}$ such that $vars(\mathcal{J}') = vars(\mathcal{I}) \setminus vars(\sigma')$. Since $vars(\sigma') \subset vars(\sigma)$, we have $vars(\mathcal{J}) \subset vars(\mathcal{J}')$ and, hence, $\mathcal{I} \backslash\!\backslash \sigma' = \mathcal{J}' \models \mathcal{J} = \mathcal{I} \backslash\!\backslash \sigma$.

*Part 2(a).* Suppose there is no strongest clause $\sigma'$ s.t. $\mathcal{I} \models \sigma'$ and $(\mathcal{I} \backslash\!\backslash \sigma') \cdot \overline{\sigma'} \models \overline{\Delta}$ and $vars(\sigma') \subset vars(\sigma)$. Then $\sigma$ is a GNR. By Proposition 11, $\mathcal{I} \dot{\cap} \sigma$ is a NR.

*Part 2(b).* Suppose $\sigma'$ is a strongest clause s.t. $\mathcal{I} \models \sigma'$ and $(\mathcal{I} \backslash\!\backslash \sigma') \cdot \overline{\sigma'} \models \overline{\Delta}$ and $vars(\sigma') \subset vars(\sigma)$. Since $\mathcal{I} \models \sigma$, $\mathcal{I} \models \sigma'$, and $vars(\sigma') \subset vars(\sigma)$, we have $\mathcal{I} \dot{\cap} \sigma' \models \mathcal{I} \dot{\cap} \sigma$ and $\mathcal{I} \backslash\!\backslash \sigma' = \mathcal{I} \backslash\!\backslash (\mathcal{I} \dot{\cap} \sigma')$. Since $\overline{\sigma'} \models \overline{\mathcal{I} \dot{\cap} \sigma'}$ and $(\mathcal{I} \backslash\!\backslash \sigma') \cdot \overline{\sigma'} \models \overline{\Delta}$, we have $(\mathcal{I} \backslash\!\backslash (\mathcal{I} \dot{\cap} \sigma')) \cdot \overline{\mathcal{I} \dot{\cap} \sigma'} \not\models \Delta$. Hence, $\mathcal{I} \dot{\cap} \sigma$ is not a NR because $\mathcal{I} \dot{\cap} \sigma'$ is stronger than $\mathcal{I} \dot{\cap} \sigma$, yet $\mathcal{I} \models \mathcal{I} \dot{\cap} \sigma'$, and $(\mathcal{I} \backslash\!\backslash (\mathcal{I} \dot{\cap} \sigma')) \cdot \overline{\mathcal{I} \dot{\cap} \sigma'} \not\models \Delta$.

## Proposition 12

**Lemma 3.** *Let $\Delta$ be a formula with discrete variables and $\sigma_1, \ldots, \sigma_n$ be the prime implicates of $\Delta$. Then $\Delta$ is equivalent to $\prod_{i=1}^{n} \sigma_i$. That is, $\Delta$ is equivalent to the conjunction of its prime implicates.*

*Proof.* We prove both directions of the equivalence. $\Delta \models \prod_{i=1}^{n} \sigma_i$ since $\Delta \models \sigma_i$ for all $i$. We next prove that $\prod_{i=1}^{n} \sigma_i \models \Delta$ by contradiction. Let $\omega$ be a world s.t. $\omega \models \prod_{i=1}^{n} \sigma_i$ but $\omega \not\models \Delta$. Then $\omega \models \overline{\Delta}$ and, hence, $\Delta \models \overline{\omega}$. Since $\overline{\omega}$ is a clause, it must be subsumed by some prime implicate $\sigma_j$. Hence, $\omega \models \sigma_j \models \overline{\omega}$ which is a contradiction.

*Proof (of Proposition 12).*

We first prove the part about GSRs, then the one for GNRs.

**GSRs are the variable-minimal prime implicants of $\overline{\forall} \mathcal{I} \cdot \Delta$.** We prove two directions next.

All variable-minimal prime implicants of $\overline{\forall} \mathcal{I} \cdot \Delta$ are GSRs.

Let $\tau$ be a variable-minimal prime implicant of $\overline{\forall} \mathcal{I} \cdot \Delta$. By Proposition 3, it suffices to prove **(1)** $\tau$ is a weakest term such that $\mathcal{I} \models \tau \models \Delta$ and **(2)** no term $\tau'$ satisfies the previous condition if $vars(\tau') \subset vars(\tau)$. We prove these next.

**(1)** We have $\tau \models \overline{\forall} \mathcal{I} \cdot \Delta \models \Delta$ by Proposition 2 and given that $\tau$ is a prime implicant of $\overline{\forall} \mathcal{I} \cdot \Delta$. We next prove $\mathcal{I} \models \tau$ by contradiction.

Suppose $\mathcal{I} \not\models \tau$. Then $\mathcal{I} \not\models \ell$ for some $X$-literal $\ell$ in $\tau$. Let $\ell' = \ell \cup \{\mathcal{I}[X]\}$ where $\mathcal{I}[X]$ is the state of variable $X$ in $\mathcal{I}$. Consider the term $\tau'$ obtained from $\tau$ by replacing literal $\ell$ by $\ell'$. Then the models of $\tau'$ are the models of $\tau$ plus the worlds $\omega'$ obtained from a model $\omega$ of $\tau$ by setting the value of variable $X$ to $\mathcal{I}[X]$. If we can prove $\omega' \models \overline{\forall}\mathcal{I}\cdot\Delta$, then it follows that $\tau'$ is an implicant of $\overline{\forall}\mathcal{I}\cdot\Delta$, which gives us a contradiction since $\tau$ is a prime implicant and $\tau \models \tau'$. We have $\omega' \models \Delta$ by Proposition 5 and since $\omega \models \tau \models \overline{\forall}\mathcal{I} \cdot \Delta$. Moreover, for any world $\omega''$ obtained from $\omega'$ by setting some variables to their states in $\mathcal{I}$, $\omega''$ can also be obtained from some model of $\tau$ by setting some variables to their states in $\mathcal{I}$. Hence, by Proposition 5, $\omega'' \models \Delta$. Proposition 5 further tells us that $\omega' \models \overline{\forall}\mathcal{I} \cdot \Delta$. Thus, $\tau \models \tau' \models \overline{\forall}\mathcal{I} \cdot \Delta$ and $\tau \neq \tau'$, which is a contradiction since $\tau'$ is a prime implicant. Hence, $\mathcal{I} \models \tau$.

We now have $\mathcal{I} \models \tau \models \Delta$. To prove (1), we need to prove that $\tau$ is the weakest term satisfying the previous property. We prove this by contradiction. Suppose $\mathcal{I} \models \tau \models \tau' \models \Delta$ and $\tau' \neq \tau$ for some term $\tau'$. Let $\omega$ be a world such that $\omega \models \tau' \models \Delta$. Since $\mathcal{I} \models \tau'$, all literals in $\tau'$ are consistent with $\mathcal{I}$. Therefore, $\omega' \models \tau' \models \Delta$ for any world $\omega'$ obtained from $\omega$ by setting some variables in $\omega$ to their states in $\mathcal{I}$. By Proposition 5, $\omega \models \overline{\forall}\mathcal{I} \cdot \Delta$, which means all models of $\tau'$ are models of $\overline{\forall}\mathcal{I} \cdot \Delta$, so $\tau'$ is an implicant of $\overline{\forall}\mathcal{I} \cdot \Delta$. This is a contradiction since $\tau$ is a prime implicant of $\overline{\forall}\mathcal{I} \cdot \Delta$. Hence, $\tau$ must be a weakest term satisfying $\mathcal{I} \models \tau \models \Delta$, so (1) holds.

When proving (1), we did not use the variable-minimality of $\tau$. Hence, we make the following observation which we use later in the proof: **(A)** every prime implicant $\tau$ of $\overline{\forall}\mathcal{I} \cdot \Delta$ satisfies $\mathcal{I} \models \tau \models \Delta$.

**(2)** We prove this by contradiction. Suppose there exists a weakest term $\tau'$ satisfying $\mathcal{I} \models \tau' \models \Delta$ and $vars(\tau') \subset vars(\tau)$. Since $\mathcal{I} \models \tau'$, all literals in $\tau'$ are consistent with $\mathcal{I}$. Therefore, for a world $\omega \models \tau' \models \Delta$, every $\omega'$ obtained from $\omega$ by setting some variables in $\omega$ to their state in $\mathcal{I}$ satisfies $\omega' \models \tau' \models \Delta$. By Proposition 5, $\tau'$ is an implicant of $\overline{\forall}\mathcal{I}\cdot\Delta$. Since $vars(\tau') \subset vars(\tau)$, there is a prime implicant $\tau''$ of $\overline{\forall}\mathcal{I} \cdot \Delta$ satisfying $vars(\tau'') \subseteq vars(\tau') \subset vars(\tau)$. This is a contradiction since $\tau$ is variable-minimal. Hence, (2) holds.

All GSRs are variable-minimal prime implicants of $\overline{\forall}\mathcal{I} \cdot \Delta$.

Let $\tau$ be a GSR. We will prove **(1)** $\tau$ is a prime implicant of $\overline{\forall}\mathcal{I} \cdot \Delta$ and **(2)** there is no prime implicant $\tau'$ of $\overline{\forall}\mathcal{I} \cdot \Delta$ satisfying $vars(\tau') \subset vars(\tau)$.

**(1)** Since $\tau$ is a GSR, $\tau$ is a weakest term such that $\mathcal{I} \models \tau \models \Delta$. Therefore, all literals in $\tau$ are consistent with $\mathcal{I}$. Thus, for a world $\omega \models \tau \models \Delta$, every world $\omega'$ obtained from $\omega$ by setting some variables in $\omega$ to their states in $\mathcal{I}$ satisfies $\omega' \models \tau \models \Delta$. By Proposition 5, $\omega \models \overline{\forall}\mathcal{I} \cdot \Delta$, so $\tau$ is an implicant of $\overline{\forall}\mathcal{I}\cdot\Delta$. To prove (1), we next show that $\tau$ is prime by contradiction. Suppose $\tau$ is not prime. Then there must be a prime implicant $\tau'$ of $\overline{\forall}\mathcal{I} \cdot \Delta$ satisfying $\tau \models \tau' \models \overline{\forall}\mathcal{I} \cdot \Delta$ and $\tau \neq \tau'$. By observation (A) in the first direction, $\tau'$ satisfies $\mathcal{I} \models \tau' \models \Delta$. This is a contradiction since $\tau$ is a GSR but not the weakest given $\mathcal{I} \models \tau \models \tau' \models \Delta$. Hence, $\tau$ is prime, so (1) holds.

**(2)** We prove this by contradiction. Assume (2) does not hold. Then there are some prime implicants $\tau'$ of $\overline{\forall}\mathcal{I}\cdot\Delta$ satisfying $vars(\tau') \subset vars(\tau)$. Let $\tau''$ be a variable-minimal prime implicant among all such prime implicants $\tau'$.

By the first direction, $\tau''$ is a GSR, which implies $\tau$ is not a GSR because $\tau$ is not variable minimal. This is a contradiction, so (2) holds.

**GNRs are the variable-minimal prime implicates of $\overline{\forall}\mathcal{I}\cdot\Delta$.** We prove both directions next.

All variable-minimal prime implicates of $\overline{\forall}\mathcal{I}\cdot\Delta$ are GNRs.

Let $\sigma$ be a variable-minimal prime implicate of $\overline{\forall}\mathcal{I}\cdot\Delta$. Our goal is to prove **(1)** $\sigma$ is a strongest clause satisfying $\mathcal{I}\models\sigma$ and $(\mathcal{I}\backslash\!\backslash\sigma)\cdot\overline{\sigma}\models\overline{\Delta}$ and **(2)** no clause $\sigma'$ satisfies the previous condition if $vars(\sigma')\subset vars(\sigma)$.

**(1)** We first prove $\mathcal{I}\stackrel{.}{\models}\sigma$. Assume the opposite: there is an $X$-literal $\ell$ in $\sigma$ such that $\mathcal{I}\not\models\ell$. Let $\sigma'$ be a clause obtained from $\sigma$ by removing $\ell$. If we show $\overline{\forall}\mathcal{I}\cdot\Delta\models\sigma'$, we get a contradiction because $\sigma$ is a prime implicate of $\overline{\forall}\mathcal{I}\cdot\Delta$ and $\sigma'\models\sigma$. This would prove $\mathcal{I}\stackrel{.}{\models}\sigma$.

We show $\overline{\forall}\mathcal{I}\cdot\Delta\models\sigma'$ by showing $\overline{\forall}\mathcal{I}\cdot\Delta\not\models\sigma'$ is impossible. $\overline{\forall}\mathcal{I}\cdot\Delta\not\models\sigma'$ only if $\omega\models\overline{\forall}\mathcal{I}\cdot\Delta\models\sigma$ and $\omega\not\models\sigma'$ for some world $\omega$ (recall, $\sigma$ is a prime implicate of $\overline{\forall}\mathcal{I}\cdot\Delta$). Such a world $\omega$ satisfies $\omega\models\ell$ (since $\ell$ is the only distinction from $\sigma$ and $\sigma'$) and, hence, $\omega\models\overline{\sigma'}$. If such a world $\omega$ exists, let $\omega'$ be a world obtained from such a $\omega$ by setting variable $X$ to its state in $\mathcal{I}$. Note that $\sigma'$ does not mention variable $X$, so $\omega'\models\overline{\sigma'}$. We have $\omega'\models\Delta$ by Proposition 5 since $\omega\models\overline{\forall}\mathcal{I}\cdot\Delta$. Moreover, for any world $\omega''$ obtained from $\omega'$ by setting some variables to their states in $\mathcal{I}$, $\omega''$ can also be obtained from $\omega$ by setting some variables to their state in $\mathcal{I}$. Hence, by Proposition 5, $\omega''\models\Delta$ for all such $\omega''$, which means $\omega'\models\overline{\forall}\mathcal{I}\cdot\Delta$. Finally, $\omega'\not\models\sigma$ since $\omega'\models\overline{\sigma'}$ and $\omega'\not\models\ell$. This is a contradiction since $\overline{\forall}\mathcal{I}\cdot\Delta\models\sigma$. Hence, no such world $\omega$ exists, which shows $\overline{\forall}\mathcal{I}\cdot\Delta\not\models\sigma'$ is impossible. Thus, $\overline{\forall}\mathcal{I}\cdot\Delta\models\sigma'$, which gives us another contradiction. Hence, $\mathcal{I}\stackrel{.}{\models}\sigma$.

We did not use the fact that $\sigma$ is variable-minimal when proving $\mathcal{I}\stackrel{.}{\models}\sigma$. Therefore, we have the following observation which we use later in the proof: **(B)** $\mathcal{I}\stackrel{.}{\models}\sigma$ holds for any prime implicate $\sigma$ of $\overline{\forall}\mathcal{I}\cdot\Delta$.

We next prove $(\mathcal{I}\backslash\!\backslash\sigma)\cdot\overline{\sigma}\models\overline{\Delta}$. Since $\mathcal{I}\stackrel{.}{\models}\sigma$, $vars(\mathcal{I}\backslash\!\backslash\sigma)=vars(\mathcal{I})\setminus vars(\sigma)$. Let $\omega$ be a world such that $\omega\models(\mathcal{I}\backslash\!\backslash\sigma)\cdot\overline{\sigma}$. We next prove $\omega\models\overline{\Delta}$ by contradiction. Assume $\omega\models\Delta$. Then for every world $\omega'$ obtained from $\omega$ by setting some variables to their states in $\mathcal{I}$, if $\omega'\neq\omega$, then $\omega'\models\sigma$ because $\mathcal{I}\stackrel{.}{\models}\sigma$ and $\omega\models(\mathcal{I}\backslash\!\backslash\sigma)$. By observation (B) above, all literals in every prime implicate of $\overline{\forall}\mathcal{I}\cdot\Delta$ are consistent with $\mathcal{I}$. Then, consider any $\omega'\neq\omega$, which immediately implies $\omega'\models\sigma$. $\omega'$ must satisfy all prime implicate of $\overline{\forall}\mathcal{I}\cdot\Delta$. Otherwise, since the subset of $\omega'$ that disagrees with $\mathcal{I}$ mentions fewer variables than $\sigma$, $\sigma$ cannot be a variable-minimal prime implicate. By Lemma 3, since $\omega'$ satisfies all prime implicate of $\overline{\forall}\mathcal{I}\cdot\Delta$, $\omega'\models\overline{\forall}\mathcal{I}\cdot\Delta$ for all such $\omega'\neq\omega$. Since $\omega'\models\Delta$ and $\omega\models\Delta$, by Proposition 5, $\omega\models\overline{\forall}\mathcal{I}\cdot\Delta$, which means $\omega\models\sigma$. This is a contradiction, since $\omega\models\overline{\sigma}$. Hence, $\omega\models\overline{\Delta}$, which implies $(\mathcal{I}\backslash\!\backslash\sigma)\cdot\overline{\sigma}\models\overline{\Delta}$.

We now prove that $\sigma$ is a strongest clause satisfying $\mathcal{I}\stackrel{.}{\models}\sigma$ and $(\mathcal{I}\backslash\!\backslash\sigma)\cdot\overline{\sigma}\models\overline{\Delta}$, which finishes the proof of (1). We prove this by contradiction. Assume there is a clause $\sigma'$ such that $\mathcal{I}\models\sigma'$, $(\mathcal{I}\backslash\!\backslash\sigma')\cdot\overline{\sigma'}\models\overline{\Delta}$, $\sigma'\models\sigma$, and $\sigma'\neq\sigma$.

If we prove $\sigma'$ is an implicate of $\overline{\forall}\mathcal{I}\cdot\Delta$, we get a contradiction. Consider $\omega\models\overline{\forall}\mathcal{I}\cdot\Delta$. If we prove $\omega\models\sigma'$, then $\overline{\forall}\mathcal{I}\cdot\Delta\models\sigma'$ follows. Assume $\omega\not\models\sigma'$, i.e., $\omega\models\overline{\sigma'}$. Let $\omega'$ be a world obtained from $\omega$ by setting all variables mentioned by $\omega$ but not by $\sigma'$ in $\omega$ to their states in $\mathcal{I}$. Then, $\omega'\models\overline{\sigma'}$ since the variables mentioned by $\sigma'$ are unchanged. Moreover, $\omega'\models(\mathcal{I}\backslash\!\backslash\sigma')$ since $\mathcal{I}\dot{\models}\sigma'$, so $\omega'\models(\mathcal{I}\backslash\!\backslash\sigma')\cdot\overline{\sigma'}\models\overline{\Delta}$. However, by Proposition 5, $\omega'\models\Delta$ because $\omega\models\overline{\forall}\mathcal{I}\cdot\Delta$. This is a contradiction, so $\omega\models\sigma'$, which shows $\sigma'$ is an implicate. Hence, (1) holds.

In the previous paragraph, we proved a property which we use later in the proof: **(C)** Every clause $\sigma'$ satisfying $\mathcal{I}\dot{\models}\sigma'$ and $(\mathcal{I}\backslash\!\backslash\sigma')\cdot\overline{\sigma'}\models\overline{\Delta}$ is an implicate of $\overline{\forall}\mathcal{I}\cdot\Delta$.

**(2)** We prove this by contradiction. Assume there is a strongest clause $\sigma'$ that satisfies $\mathcal{I}\dot{\models}\sigma'$, $(\mathcal{I}\backslash\!\backslash\sigma')\cdot\overline{\sigma'}\models\overline{\Delta}$, and $vars(\sigma')\subset vars(\sigma)$. By property (C) above, $\sigma'$ is an implicate of $\overline{\forall}\mathcal{I}\cdot\Delta$. Since $vars(\sigma')\subset vars(\sigma)$, $\sigma$ cannot be variable-minimal, which is a contradiction. Thus, (2) holds.

All GNRs are variable-minimal prime implicates of $\overline{\forall}\mathcal{I}\cdot\Delta$.

Let $\sigma$ be a GNR. We next prove **(1)** $\sigma$ is a strongest clause satisfying $\overline{\forall}\mathcal{I}\cdot\Delta\models\sigma$ and **(2)** there is no clause $\sigma'$ such that $vars(\sigma')\subset vars(\sigma)$ and $\sigma'$ satisfies the previous condition, i.e., $\sigma'$ is a prime implicate of $\overline{\forall}\mathcal{I}\cdot\Delta$.

**(1)** We first prove that $\overline{\forall}\mathcal{I}\cdot\Delta\models\sigma$, then prove $\sigma$ is the strongest such clause. Consider $\omega\models\overline{\forall}\mathcal{I}\cdot\Delta$. If we prove $\omega\models\sigma$, then $\overline{\forall}\mathcal{I}\cdot\Delta\models\sigma$ follows. Assume $\omega\not\models\sigma$, i.e., $\omega\models\overline{\sigma}$. Let $\omega'$ be a world obtained from $\omega$ by setting all variables mentioned by $\omega$ but not by $\sigma$ in $\omega$ to their states in $\mathcal{I}$. Then, $\omega'\models\overline{\sigma}$ since the variables mentioned by $\sigma$ are unchanged. Moreover, $\omega'\models(\mathcal{I}\backslash\!\backslash\sigma)$ since $\mathcal{I}\dot{\models}\sigma$, so $\omega'\models(\mathcal{I}\backslash\!\backslash\sigma)\cdot\overline{\sigma}\models\overline{\Delta}$ since $\sigma$ is a GNR. However, by Proposition 5, $\omega'\models\Delta$ because $\omega\models\overline{\forall}\mathcal{I}\cdot\Delta$. This is a contradiction, so $\omega\models\sigma$, which implies $\overline{\forall}\mathcal{I}\cdot\Delta\models\sigma$.

We next prove $\sigma$ is the strongest clause satisfying $\overline{\forall}\mathcal{I}\cdot\Delta\models\sigma$, which finishes the proof of (1). Assume $\sigma$ is not the strongest, which means there is a clause $\sigma'$ satisfying $\overline{\forall}\mathcal{I}\cdot\Delta\models\sigma'\models\sigma$. Let $\sigma'$ be the weakest such clause, i.e., $\sigma'$ is a prime implicate of $\overline{\forall}\mathcal{I}\cdot\Delta$. If $\sigma'$ is not variable-minimal, or if $vars(\sigma')\subset vars(\sigma)$, then there must exist a variable-minimal prime implicate of $\overline{\forall}\mathcal{I}\cdot\Delta$ that mentions no more variables than $\sigma'$. By the first direction, this variable-minimal prime implicate is a GNR, which means there is a GNR that mentions fewer variables than $\sigma$. Therefore, $\sigma$ cannot be a GNR, which is a contradiction. Hence, $\sigma$ must be variable-minimal and $vars(\sigma')\supseteq vars(\sigma)$. Since $\sigma'\models\sigma$, we have $vars(\sigma')=vars(\sigma)$. By the first direction, since $\sigma'$ is a variable-minimal prime implicate of $\overline{\forall}\mathcal{I}\cdot\Delta$, $\sigma'$ is a GNR. This is contradiction since $\sigma$ is a GNR and $\sigma'\models\sigma$, so (1) holds.

**(2)** We prove this by contradiction. Assume there are some clauses $\sigma'$ that are prime implicates of $\overline{\forall}\mathcal{I}\cdot\Delta$ and $vars(\sigma')\subset vars(\sigma)$. Let $\sigma''$ be a variable-minimal prime implicate among those prime implicates. By the first direction, $\sigma''$ is a GNR that mentions fewer variables than $\sigma$. This is a contradiction since $\sigma$ is a GNR. Hence, (2) holds.

**Proposition 13**

The proof of Proposition 13 uses the next two lemmas which we state and prove first.

**Lemma 4.** *Let $\alpha$ be an NNF and $\ell$ by an $X$-literal. If $\ell \models \ell'$ for every $X$-literal $\ell'$ that occurs in $\alpha$ then*

$$\overline{\forall}\, x_i \cdot (\alpha + \ell) = \overline{\forall}\, x_i \cdot \alpha + \overline{\forall}\, x_i \cdot \ell.$$

*Proof.* We consider two cases.

Case $x_i \models \ell$.
   Then $\overline{\forall}\, x_i \cdot (\alpha + \ell) = (\alpha + \ell) \cdot (\alpha|x_i + \ell|x_i) = (\alpha + \ell) \cdot (\alpha|x_i + \top) = \alpha + \ell = \alpha + \overline{\forall}\, x_i \cdot \ell$. We next show $\overline{\forall}\, x_i \cdot \alpha = \alpha$. If $\ell \models \ell'$ then $x_i \models \ell'$. Hence, $\alpha \models \alpha|x_i$ since $\alpha|x_i$ is obtained by replacing every $X$-literal $\ell'$ in $\alpha$ with $\top$. We now have $\alpha \models \alpha \cdot \alpha|x_i$ and, hence, $\alpha = \alpha \cdot \alpha|x_i = \overline{\forall}\, x_i \cdot \alpha$. Thus, $\overline{\forall}\, x_i \cdot (\alpha + \ell) = \overline{\forall}\, x_i \cdot \alpha + \overline{\forall}\, x_i \cdot \ell$.
Case $x_i \not\models \ell$.
   Since $\overline{\forall}\, x_i \cdot \ell = \ell \cdot \ell|x_i = \ell \cdot \bot = \bot$, it suffices to show $\overline{\forall}\, x_i \cdot (\alpha + \ell) = \overline{\forall}\, x_i \cdot \alpha$. We have $\overline{\forall}\, x_i \cdot (\alpha + \ell) = (\alpha + \ell) \cdot (\alpha|x_i + \ell|x_i) = (\alpha + \ell) \cdot (\alpha|x_i + \bot) = (\alpha + \ell) \cdot \alpha|x_i = (\alpha \cdot \alpha|x_i) + (\ell \cdot \alpha|x_i) = \overline{\forall}\, x_i \cdot \alpha + (\ell \cdot \alpha|x_i)$. We next show that $\ell \cdot \alpha|x_i \models \overline{\forall}\, x_i \cdot \alpha$ which finishes the proof. We have $\alpha|x_i = \alpha|x_j$ for all $x_j \models \ell$ since $\ell \models \ell'$ for every $X$-literal $\ell'$ in $\alpha$. Hence, $\ell \cdot \alpha|x_i = \sum_{x_j \models \ell}(x_j \cdot \alpha|x_i) = \sum_{x_j \models \ell}(x_j \cdot \alpha|x_j) = \sum_{x_j \models \ell}(x_j \cdot \alpha) = \ell \cdot \alpha$ which implies $\ell \cdot \alpha|x_i \models \alpha \cdot \alpha|x_i = \overline{\forall}\, x_i \cdot \alpha$.

**Lemma 5.** *The formula of class $c$ in decision graph $T$ is equivalent to an NNF $\Delta^c[T]$ defined as follows:*

$$\Delta^c[T] = \begin{cases} \top & \text{if } T \text{ has class } c \\ \bot & \text{if } T \text{ has a class } c' \neq c \\ \prod_j (\Delta^c[T_j] + \ell) & \text{if } T \text{ has edges } \xrightarrow{X, S_j} T_j \end{cases}$$

*where $\ell$ is the $X$-literal $\{x_i | x_i \notin S_j\}$.*

*Proof.* This result is proven in [18].

*Proof (of Proposition 13).* We will show how to compute the general reason $\overline{\forall}\, \mathcal{I} \cdot \Delta^c[T]$ by using the definition of class formula $\Delta^c[T]$ as given by Lemma 5. By Proposition 7, $\overline{\forall}$ distributes over the and-nodes of $\Delta^c[T]$. Every disjunction in this NNF has the form $\Delta^c[T_j] + \ell$ where $\ell = \{x_i | x_i \notin S_j\}$ is an $X$-literal. Every $X$-literal in the NNF $\Delta^c[T_j]$ has the form $\ell' = \{x_i | x_i \notin S'_k\}$, where $S'_k \subseteq S_j$ by the weak test-once property. Hence, $\ell \models \ell'$. Thus, by Lemma 4, $\overline{\forall}$ distributes over the or-nodes of $\Delta^c[T]$. Hence, we can compute $\overline{\forall}\, \mathcal{I} \cdot \Delta^c[T]$ by simply applying $\overline{\forall}\, \mathcal{I}$ to the literals of $\Delta^c[T]$. If we do this using Proposition 6, we get the closed-from of $\overline{\forall}\, \mathcal{I} \cdot \Delta^c[T]$ as shown in Proposition 13.

## Proposition 14

*Proof (of Proposition 14).* The given formula for the general reason has no negations so it is an NNF. Every $X$-literal in this NNF has the form $\ell = \{x_i | x_i \notin S_j\}$ where $\mathcal{I}[X] \notin S_j$. Hence, $\mathcal{I}[X] \models \ell$ which implies $\mathcal{I} \models \ell$. Thus, every literal in the NNF is consistent with instance $\mathcal{I}$.

Every disjunction in this NNF has the form $\Gamma^c[T_j] + \ell$ where $\ell = \{x_i | x_i \notin S_j\}$ is an $X$-literal. Every $X$-literal in the NNF $\Gamma^c[T_j]$ has the form $\ell' = \{x_i | x_i \notin S'_k\}$, where $S'_k \subset S_j$ by the weak test-once property. Hence, $\ell \models \ell'$ and $\ell \neq \ell'$.

## Proposition 15

*Proof (of Proposition 15).* We prove the two directions.

- $\tau \in \mathrm{pi}(\alpha \cdot \beta)$ only if $\tau \in \ominus(\mathrm{pi}(\alpha) \times \mathrm{pi}(\beta))$. Suppose $\tau \in \mathrm{pi}(\alpha \cdot \beta)$; that is, $\tau \models \alpha \cdot \beta$ and $\tau \models \tau' \models \alpha \cdot \beta$ for term $\tau'$ only if $\tau = \tau'$. It suffices to show (1) $\tau \in \mathrm{pi}(\alpha) \times \mathrm{pi}(\beta)$ and (2) $\tau \models \tau' \in \mathrm{pi}(\alpha) \times \mathrm{pi}(\beta)$ only if $\tau = \tau'$.
  Let $\tau_\alpha$ be the weakest term such that $\tau \models \tau_\alpha \models \alpha$ and define $\tau_\beta$ analogously. We next show that $\tau_\alpha \in \mathrm{pi}(\alpha)$, $\tau_\beta \in \mathrm{pi}(\beta)$ and $\tau = \tau_\alpha \cdot \tau_\beta$ which implies (1). Suppose $\tau_\alpha \notin \mathrm{pi}(\alpha)$: $\tau_\alpha \models \tau'_\alpha \models \alpha$ and $\tau_\alpha \neq \tau'_\alpha$ for some term $\tau'_\alpha$. This contradicts the definition of $\tau_\alpha$ so $\tau_\alpha \in \mathrm{pi}(\alpha)$. We can similarly show $\tau_\beta \in \mathrm{pi}(\beta)$. Finally, if $\tau \neq \tau_\alpha \cdot \tau_\beta$, then $\tau \notin \mathrm{pi}(\alpha \cdot \beta)$ since $\tau \models \tau_\alpha \cdot \tau_\beta \models \alpha \cdot \beta$, a contradiction, so $\tau = \tau_\alpha \cdot \tau_\beta$. Hence, (1) holds. Suppose now (2) does not hold: $\tau \models \tau' \in \mathrm{pi}(\alpha) \times \mathrm{pi}(\beta)$ and $\tau \neq \tau'$ for some term $\tau'$. Let $\tau'_\alpha \in \mathrm{pi}(\alpha)$ and $\tau'_\beta \in \mathrm{pi}(\beta)$ such that $\tau' = \tau'_\alpha \cdot \tau'_\beta$. Then $\tau \models \tau' = \tau'_\alpha \cdot \tau'_\beta \models \alpha \cdot \beta$ and $\tau \neq \tau'$ which is a contradiction with $\tau \in \mathrm{pi}(\alpha \cdot \beta)$. Hence, (2) holds and we have $\tau \in \ominus(\mathrm{pi}(\alpha) \times \mathrm{pi}(\beta))$.
- $\tau \in \ominus(\mathrm{pi}(\alpha) \times \mathrm{pi}(\beta))$ only if $\tau \in \mathrm{pi}(\alpha \cdot \beta)$. Suppose $\tau \in \ominus(\mathrm{pi}(\alpha) \times \mathrm{pi}(\beta))$: $\tau \in \mathrm{pi}(\alpha) \times \mathrm{pi}(\beta)$ and $\tau \models \tau' \in \mathrm{pi}(\alpha) \times \mathrm{pi}(\beta)$ only if $\tau = \tau'$. We show (1) $\tau \models \alpha \cdot \beta$ and (2) $\tau \models \tau' \models \alpha \cdot \beta$ for term $\tau'$ only if $\tau = \tau'$.
  Let $\tau = \tau_\alpha \cdot \tau_\beta$ where $\tau_\alpha \in \mathrm{pi}(\alpha)$ and $\tau_\beta \in \mathrm{pi}(\beta)$. Then $\tau \models \alpha \cdot \beta$ which establishes (1). Suppose (2) does not hold: $\tau \models \tau' \models \alpha \cdot \beta$ and $\tau \neq \tau'$ for some term $\tau'$. Let $\tau'$ be the weakest term satisfying the previous property. Then $\tau' \in \mathrm{pi}(\alpha \cdot \beta)$. Let $\tau'_\alpha$ be the weakest term such that $\tau' \models \tau'_\alpha \models \alpha$ and define $\tau'_\beta$ analogously. Then $\tau'_\alpha \in \mathrm{pi}(\alpha)$, $\tau'_\beta \in \mathrm{pi}(\beta)$ and $\tau' = \tau'_\alpha \cdot \tau'_\beta$ as shown in the first direction. Hence, $\tau \models \tau' \in \mathrm{pi}(\alpha) \times \mathrm{pi}(\beta)$. Since $\tau \neq \tau'$, we get a contradiction with $\tau \in \ominus(\mathrm{pi}(\alpha) \times \mathrm{pi}(\beta))$. Hence, (2) holds and we have $\tau \in \mathrm{pi}(\alpha \cdot \beta)$.

## Proposition 16

*Proof (of Proposition 16).* For literal $\ell$, $\mathrm{pi}(\ell) = \{\ell\}$. Hence, what we need to show is $\mathrm{pi}(\ell + \beta) = \ominus(\{\ell\} \cup \mathrm{pi}(\beta))$. We next prove both directions.

- $\tau \in \mathrm{pi}(\ell + \beta)$ only if $\tau \in \ominus(\{\ell\} \cup \mathrm{pi}(\beta))$. Suppose $\tau \in \mathrm{pi}(\ell + \beta)$: $\tau \models \ell + \beta$ and $\tau \models \tau' \models \ell + \beta$ for term $\tau'$ only if $\tau = \tau'$. We need to show (1) $\tau \in \{\ell\} \cup \mathrm{pi}(\beta)$ and (2) $\tau \models \tau' \in \{\ell\} \cup \mathrm{pi}(\beta)$ only if $\tau = \tau'$.

Our goal is to first prove either $\tau \models \ell$ or $\tau \models \beta$, and then prove $\tau \in \{\ell\} \cup$ pi$(\beta)$. To prove either $\tau \models \ell$ or $\tau \models \beta$, showing $\tau \not\models \ell$ only if $\tau \models \beta$ suffices. Suppose $\tau \not\models \ell$. Assume there exists a world $\omega$ such that $\omega \models \tau \models \ell + \beta$ but $\omega \not\models \beta$. If we find a contradiction, then all models of $\tau$ are models of $\beta$, i.e. $\tau \models \beta$, which is exactly what we want. Since $\tau \not\models \ell$, either $\tau$ does not mention the variable of $\ell$ or $\tau \models \ell'$ for $\ell \models \ell'$ and $\ell' \neq \ell$. Both cases suggest there exists a world $\omega'$ obtained from $\omega$ by setting the state of the variable of $\ell$ to some state not in $\ell$ such that $\omega' \models \tau$, $\omega' \not\models \beta$ by the second property in Proposition 14, and $\omega' \not\models \ell$. That is, $\omega' \models \tau$ but $\omega' \not\models \beta$ and $\omega' \not\models \ell$, so $\tau \not\models \ell + \beta$. This is a contradiction. Thus, $\omega \models \tau$ only if $\omega \models \beta$, so $\tau \models \beta$. Therefore, $\tau \not\models \ell$ only if $\tau \models \beta$. Equivalently, $\tau \models \ell$ or $\tau \models \beta$. We can now prove (1) by considering two cases: $\tau \models \ell$ and $\tau \not\models \ell$. If $\tau \models \ell$, then $\tau \models \ell \models \ell + \beta$ where $\ell$ is a term, so $\tau = \ell$, which means $\tau \in \{\ell\} \cup$ pi$(\beta)$. If $\tau \not\models \ell$, then $\tau \models \beta$. For any term $\tau' \models \beta$, if $\tau \models \tau'$, then $\tau \models \tau' \models \beta \models \ell + \beta$, which implies $\tau = \tau'$. Thus, $\tau \in$ pi$(\beta)$. Hence, (1) must hold.
Suppose (2) does not hold: $\tau \models \tau' \in \{\ell\} \cup$ pi$(\beta)$ for a term $\tau' \neq \tau$. Then $\tau' \models \ell + \beta$ since $\tau' \in \{\ell\} \cup$ pi$(\beta)$. We now have $\tau \models \tau' \models \ell + \beta$ and $\tau \neq \tau'$, which is a contradiction with $\tau \in$ pi$(\ell + \beta)$. Hence, (2) must hold.

$\tau \in \ominus(\{\ell\} \cup$ pi$(\beta))$ only if $\tau \in$ pi$(\ell + \beta)$. Suppose $\tau \in \ominus(\{\ell\} \cup$ pi$(\beta))$: $\tau \in \{\ell\} \cup$ pi$(\beta)$, and $\tau \models \tau' \in \{\ell\} \cup$ pi$(\beta)$ only if $\tau = \tau'$. We next show (1) $\tau \models \ell + \beta$ and (2) $\tau \models \tau' \models \ell + \beta$ for term $\tau'$ only if $\tau = \tau'$.
If $\tau \in \{\ell\}$, then $\tau \models \ell$. If $\tau \in$ pi$(\beta)$, then $\tau \models \beta$. Thus, $\tau \models \ell + \beta$ follows from $\tau \in \{\ell\} \cup$ pi$(\beta)$, which establishes (1). Suppose (2) does not hold: $\tau \models \tau' \models \ell + \beta$ and $\tau \neq \tau'$ for some term $\tau'$. Let $\tau'$ be the weakest term satisfying the previous property. Then $\tau' \in$ pi$(\ell + \beta)$. By the first direction, $\tau' \in \ell \cup$ pi$(\beta)$, so $\tau \models \tau' \in \ell \cup$ pi$(\beta)$ and $\tau \neq \tau'$, which is a contradiction. Hence, (2) holds and we have $\tau \in$ pi$(\ell + \beta)$.

## Proposition 17

**Lemma 6.** *Let $\gamma$ be a node in the NNF passed to Algorithm 1, GSR$(.)$. The terms in GSR$(\gamma)$ are implicants of $\gamma$.*

*Proof.* The proof is by induction on the structure of the NNF passed to Algorithm 1.
*Base case:* $\gamma$ is a literal or constant. This case is immediate.
*Inductive step:* $\gamma = \alpha \cdot \beta$. By the induction assumption, GSR$(\alpha)$ are implicants of $\alpha$ and GSR$(\beta)$ are implicants of $\beta$. For any $\tau_1 \in$ GSR$(\alpha)$ and $\tau_2 \in$ GSR$(\beta)$, we have $\tau_1 \cdot \tau_2 \models \alpha \cdot \beta$ so $S =$ GSR$(\alpha) \times$ GSR$(\beta)$ are implicants of $\alpha \cdot \beta$. Algorithm 1 returns a subset of $S$ so the result holds.
*Inductive step:* $\gamma = \alpha + \ell$ where $\ell$ is a literal. By the induction assumption, GSR$(\alpha)$ are implicants of $\alpha$ and GSR$(\ell)$ is the implicant of $\ell$. Then GSR$(\alpha) \cup$ GSR$(\ell)$ are implicants of $\alpha + \ell$, since every implicant of $\alpha$ or $\ell$ implies $\alpha + \ell$. Algorithm 1 returns a subset of $S$ so the result holds.

*Proof (of Proposition 17).*

Let $\Delta$ be the NNF passed in the first call GSR($\Delta$) to Algorithm 1. We next prove two directions.

*First direction:* If $\tau$ is a variable-minimal prime implicant of $\Delta$, then $\tau \in$ GSR($\Delta$). We prove this by contradiction.

We first note that Algorithm 1, GSR($\Delta$), without Line 10 (variable minimization) corresponds to Algorithm 2, PI($\Delta$), which computes the prime implicants of $\Delta$. Hence, we will say Algorithm 2 to mean Algorithm 1 without Line 10.

Suppose now that $\tau$ is a variable-minimal prime implicant of $\Delta$ and $\tau \notin$ GSR($\Delta$). Since $\tau$ is a prime implicant of $\Delta$, it must be equivalent to the conjunction of some terms $S^*$ constructed by Algorithm 2, where at least one of these terms is dropped on Lines 6, 8 or 10 of Algorithm 1. By Lemma 6, for each node $\gamma$ of the NNF $\Delta$, GSR($\gamma$) are implicants of $\gamma$. Therefore, no prime implicant of $\gamma$ can be subsumed by any distinct term in GSR($\gamma$). Thus, one of the terms $\tau^*$ in $S^*$ must have been removed by variable minimization on Line 10 of Algorithm 1; that is, not by the subsumption checks on Lines 6 or 8 of the algorithm. Let $\Delta^*$ be the NNF node where the term $\tau^*$ is dropped by Algorithm 1. Then, there is a term $\tau^+$ generated by Algorithm 1 at node $\Delta^*$ such that $vars(\tau^+) \subset vars(\tau^*)$ and $ivars(\Delta^*) \cap (vars(\tau^*) \setminus vars(\tau^+)) \neq \emptyset$. It follows that term $\tau$ is equivalent to the conjunction of term $\tau^*$ and some other terms $S^o$ constructed by Algorithm 2 at nodes outside NNF $\Delta^*$. Consider term $\tau'$ that is equivalent to the conjunction of $\tau^+$ and the terms in $S^o$. Since the NNF $\Delta$ is locally fixated, the set of variables mentioned by $\tau'$ is equal to the union of the set of variables mentioned by the terms $S^o \cup \{\tau^+\}$. The same applies to term $\tau$ and terms $S^o \cup \{\tau^*\}$. Since $ivars(\Delta^*) \cap (vars(\tau^*) \setminus vars(\tau^+)) \neq \emptyset$, we have $vars(\tau') \subset vars(\tau)$. Note that $\tau'$ is an implicant of $\Delta$. Thus, $\tau'$ is either a prime implicant of $\Delta$ or is subsumed by some distinct prime implicant of $\Delta$. Hence, $\tau$ cannot be a variable-minimal prime implicant of $\Delta$, which is a contradiction.

*Second direction:* If $\tau \in$ GSR($\Delta$), then $\tau$ is a variable-minimal prime implicant of $\Delta$. Suppose $\tau \in$ GSR($\Delta$). It suffices to show (1) $\tau \models \Delta$, (2) there is no prime implicant $\tau'$ of $\Delta$ such that $vars(\tau') \subset vars(\tau)$, and (3) there is no distinct prime implicant $\tau'$ of $\Delta$ such that $\tau \models \tau'$.

Lemma 6 implies (1) immediately. We now show (2). By the first direction, GSR($\Delta$) contains all variable-minimal prime implicants of $\Delta$. Thus, to prove (2), it suffices to prove that there does not exist a term $\tau' \in$ GSR($\Delta$) such that $vars(\tau') \subset vars(\tau)$. By the definition of $ivars(.)$, $ivars(\Delta) = vars(\Delta)$ when $\Delta$ is the NNF passed to the first call to Algorithm 1. Thus, $\boxtimes(S, ivars(\Delta))$ on Line 10 of the algorithm removes all terms from $S$ that are not variable-minimal in this case. Therefore, (2) holds. We next prove (3) by contradiction. Assume there is a distinct prime implicant $\tau'$ of $\Delta$ such that $\tau \models \tau'$ (i.e., $\tau'$ subsumes $\tau$). Since, by the first direction, GSR($\Delta$) contains all variable-minimal prime implicants of $\Delta$, $\tau'$ cannot be a variable-minimal prime implicant of $\Delta$; otherwise $\tau'$ will be in GSR($\Delta$) so $\tau$ will not be in GSR($\Delta$) as it will be removed by the subsumption checks on Line 6 or Line 8, which is a contradiction. By (2), no variable-minimal prime implicant of $\Delta$ has a strict subset of the variables in $\tau$. Therefore, $vars(\tau') \not\subseteq vars(\tau)$; otherwise $\tau'$ must be a variable-minimal prime

implicant of $\Delta$. Note that $\tau'$ subsumes $\tau$ only if $vars(\tau') \subseteq vars(\tau)$. Therefore, $\tau'$ cannot subsume $\tau$, which is a contradiction. Hence, (3) holds.

## Proposition 18

We first prove a dual of Proposition 18 using several lemmas. The dual is for the *consensus* operation which can be used to compute the prime implicants of a DNF. The proof uses the same structure as the proof of Theorem 3.5 in [15] which treats the Boolean case of consensus.

**Definition 11.** *Let $\ell_1 \cdot \gamma_1$ and $\ell_2 \cdot \gamma_2$ be terms where $\ell_1$ and $\ell_2$ are $X$-literals such that $\ell_1 \not\models \ell_2$ and $\ell_2 \not\models \ell_1$. Then $\gamma = (\ell_1 + \ell_2) \cdot \gamma_1 \cdot \gamma_2$ is an $X$-consensus of the terms if $\gamma \neq \bot$.*

We use $\{\ell_1 \cdot \gamma_1, \ell_2 \cdot \gamma_2\} \blacktriangleleft X$ to denote the consensus of terms $\ell_1 \cdot \gamma_1$ and $\ell_2 \cdot \gamma_2$ on variable $X$. We also use Consensus$(\Delta)$ to denote the result of closing DNF $\Delta$ under consensus and then removing all subsumed terms. Our proofs will also use the following definition for consensus over multiple terms (can be emulated by Definition 11 over two terms if we skip subsumed consensus).

**Definition 12.** *Let $\ell_1 \cdot \gamma_1, \ldots, \ell_n \cdot \gamma_n$ be terms where $\ell_1, \ldots, \ell_n$ are $X$-literals. Then $\gamma = (\sum_{i=1}^{n} \ell_i) \cdot \prod_{i=1}^{n} \gamma_i$ is an $X$-consensus of the terms if $\gamma \neq \bot$.*

**Lemma 7.** *We have $(\ell_1 + \ell_2) \cdot \gamma_1 \cdot \gamma_2 \models \ell_1 \cdot \gamma_1 + \ell_2 \cdot \gamma_2$. Moreover, Consensus$(\Delta)$ is equivalent to $\Delta$.*

*Proof.* If $\omega \models (\ell_1 + \ell_2) \cdot \gamma_1 \cdot \gamma_2$, then $\omega \models \ell_1 \cdot \gamma_1 \cdot \gamma_2$ or $\omega \models \ell_2 \cdot \gamma_1 \cdot \gamma_2$. In either case, $\omega \models \ell_1 \cdot \gamma_1 + \ell_2 \cdot \gamma_2$. Hence, $(\ell_1 + \ell_2) \cdot \gamma_1 \cdot \gamma_2 \models \ell_1 \cdot \gamma_1 + \ell_2 \cdot \gamma_2$. This means that we can add to a DNF $\Delta$ the consensus of any of its terms without changing the models of $\Delta$. Hence, Consensus$(\Delta) = \Delta$.

**Lemma 8.** *Let $\tau$ be a simple term that mentions all variables in DNF $\Delta$. If $\tau \models \Delta$, then $\tau \models \tau'$ for some term $\tau'$ in $\Delta$ (that is, $\tau$ is subsumed by some term in $\Delta$).*

*Proof.* Since $\tau$ is simple and mentions all variables of $\Delta$, then $\tau'|\tau = \top$ or $\tau'|\tau = \bot$ for every term $\tau'$ in $\Delta$. Since $\tau \models \Delta$, $\Delta|\tau = \top$ so $\tau'|\tau = \top$ for at least one term $\tau'$ in $\Delta$. This term must satisfy $\tau \models \tau'$ and, hence, $\tau$ is subsumed by $\tau'$.

Lemma 8 does not hold if term $\tau$ is not simple. Counterexample: $\tau = x_{123}$ and $\Delta = x_{12} + x_{23}$.

**Lemma 9.** *A prime implicant of DNF $\Delta$ can mention only variables mentioned by Consensus$(\Delta)$.*

*Proof.* If Consensus$(\Delta)$ does not mention variable $X$, then $\Delta$ does not depend on $X$ since Consensus$(\Delta)$ is equivalent to $\Delta$ by Lemma 7. Hence, any implicant of $\Delta$ will remain an implicant of $\Delta$ if we drop any $X$-literal from it. Hence, a prime implicant of $X$ cannot mention variable $X$.

**Lemma 10.** *Consensus($\Delta$) is the set of prime implicants for DNF $\Delta$.*

*Proof.* We first show that every prime implicant of $\Delta$ is in Consensus($\Delta$), and then show the second direction: every term in Consensus($\Delta$) is a prime implicant of $\Delta$.

To show the first direction, suppose $\tau_0$ is a prime implicant of $\Delta$ and $\tau_0 \notin$ Consensus($\Delta$). We next show a contradiction. Let $S$ be the set of terms $\tau$ such that:

1. $\tau$ only mentions variables present in Consensus($\Delta$).
2. $\tau \models \tau_0$.
3. $\tau$ is not subsumed by any term in Consensus($\Delta$).

By Lemma 9, $\tau_0$ can only mention variables in Consensus($\Delta$). Thus, $S$ must be non-empty because $\tau_0 \in S$. Let $\tau_m$ be the term in $S$ that mentions the largest number of variables (i.e. with the maximal length).

Case: $\tau_m$ mentions all variables of Consensus($\Delta$).

Apply the following procedure which may change the value of $\tau_m$ but will keep the set $S$ intact:

While $\tau_m \in S$:
- Write $\tau_m$ as $x_{12\ldots n} \cdot \tau_m'$ for some variable $X$, term $\tau_m'$ and $n > 1$. This can be done since $\tau_m$ is not a simple term by definition of $S$ and Lemma 8.
- For $i = 1, \ldots, n$: If $x_i \cdot \tau_m' \in S$, set $\tau_m$ to $x_i \cdot \tau_m'$ and exit for-loop (variables of $\tau_m$ are invariant).
- Exit while-loop if the for-loop did not set $\tau_m$.

When the procedure terminates, $\tau_m$ will be such that $\tau_m \in S$ but for some variable $X$, $x_i \cdot \tau_m' \notin S$ for all $i$ in $1, \ldots, n$. The procedure will always terminate because, by Lemma 8, simple terms that mention all variables cannot be in $S$. Since $\tau_m \in S$ upon termination, we have $\tau_m \models \tau_0$. And since $x_i \cdot \tau_m' \models \tau_m$ for all $i$ in $1, \ldots, n$, we have $x_i \cdot \tau_m' \models \tau_0$ for all $i$ in $1, \ldots, n$. Since $x_i \cdot \tau_m' \notin S$, $x_i \cdot \tau_m'$ must be subsumed by some respective term in Consensus($\Delta$) for each $i$. Since $\tau_m = x_{1\ldots n} \cdot \tau_m'$ is not subsumed by these respective terms, each $x_i \cdot \tau_m'$ must be subsumed by some term $\alpha_i \in$ Consensus($\Delta$) that mentions state $x_i$.

Let $\beta_i$ be $\alpha_i$ but without its $X$-literal. Since $x_i \cdot \tau_m' \models \alpha_i$ for all $i$, we have $\tau_m' \models \beta_i$ for all $i$. Hence, $\tau_m' \models \prod_{i=1}^n \beta_i$. This means that the consensus of $\alpha_1, \ldots, \alpha_n$ on variable $X$ exists since $\alpha_i \cdot \ldots \cdot \alpha_n$ are consistent. Since $\tau_m' \models \prod_{i=1}^n \beta_i$, we have $x_{1\ldots n} \cdot \tau_m' \models x_{1\ldots n} \cdot \prod_{i=1}^n \beta_i$. And since each $\alpha_i$ mentions $x_i$, we have $x_{1\ldots n} \cdot \prod_{i=1}^n \beta_i \models \{\alpha_1, \ldots, \alpha_i\} \blacktriangleleft X$. Therefore, $\tau_m \models x_{1\ldots n} \cdot \prod_{i=1}^n \beta_i \models \{\alpha_1, \ldots, \alpha_i\} \blacktriangleleft X$ since $\tau_m = x_{1\ldots n} \cdot \tau_m'$. Hence, $\tau_m$ is subsumed by the consensus of $\alpha_1, \ldots, \alpha_i$ on variable $X$. Since $\alpha_i \in$ Consensus($\Delta$) for all $i$, their consensus must be subsumed by some term in Consensus($\Delta$). Therefore, $\tau_m$ is subsumed by some term in Consensus($\Delta$). This contradicts $\tau_m \in S$.

Case: $\tau_m$ does not mention all variables of Consensus($\Delta$).

Suppose $\tau_m$ does not mention variable $X$ which appears in Consensus($\Delta$). Consider the terms $x_1 \cdot \tau_m, \ldots, x_k \cdot \tau_m$ where $x_1, \ldots, x_k$ are the states of variable $X$. Since $\tau_m$ is a term in $S$ of maximal length, terms $x_1 \cdot \tau_m, \ldots, x_k \cdot \tau_m$ cannot be in set $S$. Because $x_i \cdot \tau_m$ satisfies the first two requirements of set $S$ for all $i$ between 1 and $k$, $x_i \cdot \tau_m$ must be subsumed by some term $\gamma_i \in$ Consensus($\Delta$). Since $\tau_m$ is not subsumed by $\gamma_i$ for any $i$, $\gamma_i$ must mention state $x_i$. Similarly, taking the consensus of $\gamma_1, \ldots, \gamma_k$ is allowed because $\tau_m \models \prod_{i=1}^{k}(\gamma_i | x_i)$. Note that $\{\gamma_1, \ldots, \gamma_k\} \blacktriangleleft X$ does not mention variable $X$. Since $x_i \cdot \tau_m \models \gamma_i$ for all $i$, we have $\tau_m \models \{\gamma_1, \ldots, \gamma_k\} \blacktriangleleft X$. Since $\gamma_i$ are all in Consensus($\Delta$), $\{\gamma_1, \ldots, \gamma_k\} \blacktriangleleft X$ must be subsumed by some term in Consensus($\Delta$). This implies that $\tau_m$ is subsumed by some term in Consensus($\Delta$), which contradicts the assumption that $\tau_m$ is in $S$.

Our assumption that $\tau_0$ is a prime implicant of $\Delta$ but $\tau_0 \notin$ Consensus($\Delta$) leads to a contradiction in both cases above. Thus, Consensus($\Delta$) includes all prime implicants of $\Delta$.

We next show the second direction: every term in Consensus($\Delta$) is a prime implicant of $\Delta$. Every term in Consensus($\Delta$) is an implicant of $\Delta$ by Lemma 7. Moreover, by definition of Consensus($\Delta$), no term in Consensus($\Delta$) can subsume another term in Consensus($\Delta$). Hence, given the first direction, every term in Consensus($\Delta$) is a prime implicants of $\Delta$.

**Lemma 11.** *The prime implicates of $\Delta$ are the negations of the prime implicates of $\overline{\Delta}$.*

*Proof.* This follows since $\tau \models \overline{\Delta}$ iff $\Delta \models \overline{\tau}$, and since $\tau$ is equivalent to a term iff $\overline{\tau}$ is equivalent to a clause.

**Lemma 12.** *For terms $\ell_1 \cdot \tau_1$ and $\ell_2 \cdot \tau_2$ where $\ell_1$, $\ell_2$ are $X$-literals, the negation of the consensus of $\ell_1 \cdot \tau_1$ and $\ell_2 \cdot \tau_2$ on $X$ is equivalent to the resolvent of $\overline{\ell_1} + \overline{\tau_1}$ and $\overline{\ell_2} + \overline{\tau_2}$ on $X$.*

*Proof.* The consensus of $\ell_1 \cdot \tau_1$ and $\ell_2 \cdot \tau_2$ is $(\ell_1 + \ell_2) \cdot \tau_1 \cdot \tau_2$. The resolvent of $\overline{\ell_1} + \overline{\tau_1}$ and $\overline{\ell_2} + \overline{\tau_2}$ is $(\overline{\ell_1 \cdot \ell_2}) + \overline{\tau_1} + \overline{\tau_2}$. Finally, $(\overline{\ell_1 \cdot \ell_2}) + \overline{\tau_1} + \overline{\tau_2} = \overline{(\ell_1 + \ell_2) \cdot \tau_1 \cdot \tau_2}$.

*Proof (of Proposition 18).* Let $\Delta$ be a CNF. By Lemma 10, closing the DNF $\overline{\Delta}$ under consensus and removing subsumed terms yields the prime implicants of $\overline{\Delta}$. By Lemma 12, the negations of consensus generated while closing DNF $\overline{\Delta}$ under consensus correspond to resolvents generated while closing CNF $\Delta$ under resolution. By Lemma 11, the prime implicates of $\Delta$ are the negations of the prime implicants of $\overline{\Delta}$. Hence, closing $\Delta$ under resolution generates all the negations of the prime implicants of $\overline{\Delta}$, which are the prime implicates of $\Delta$. Therefore, closing $\Delta$ under resolution and removing subsumed clauses yields exactly the prime implicates of $\Delta$.

**Proposition 19**

The proof of this proposition uses two lemmas which effectively say that when applying resolution to a locally fixated CNF, the variables of resolvents grow monotonically. That is, if a clause $\sigma^*$ was derived using a clause $\sigma$, then the variables of $\sigma^*$ are a superset of the variables of $\sigma$.

**Lemma 13.** *Let $\alpha = \ell_1 + \sigma_1$, $\beta = \ell_2 + \sigma_2$ be two clauses which are locally fixated on some instance $\mathcal{I}$. If $\ell_1$ and $\ell_2$ are $X$-literals, and if $\sigma$ is the $X$-resolvent of clauses $\alpha$ and $\beta$, then $vars(\sigma) = vars(\alpha) \cup vars(\beta)$.*

*Proof.* Recall that a clause is a disjunction of literals over distinct variables. Suppose that $\ell_1$ and $\ell_2$ are $X$-literals. If $\sigma$ is the $X$-resolvent of clauses $\alpha$ and $\beta$, then $\sigma$ is the clause equivalent to $(\ell_1 \cdot \ell_2) + \sigma_1 + \sigma_2$ and $\sigma \neq \top$. Since $\alpha$ and $\beta$ are locally fixated on $\mathcal{I}$, all literals in $\alpha$ and $\beta$ are consistent with $\mathcal{I}$, so $\ell_1 \cdot \ell_2 \neq \bot$ and $X \in vars(\sigma)$. Since $\sigma \neq \top$, then $\sigma_1 + \sigma_2 \neq \top$ so the variables of the clause equivalent to $\sigma_1 + \sigma_2$ are $vars(\sigma_1) \cup vars(\sigma_2)$. Hence, $vars(\sigma) = vars(\alpha) \cup vars(\beta)$.

We will say that clause $\sigma^*$ is a *descendant resolvent* of clause $\sigma$ if $\sigma^* = \sigma$ or if $\sigma^*$ was obtained by a sequence of resolutions that involved clause $\sigma$.

**Lemma 14.** *Let $S$ be a set of clauses which are locally fixated on some instance $\mathcal{I}$, and let $S^*$ be the result of closing $S$ under resolution. If $\sigma^* \in S^*$ is a descendant resolvent of some $\sigma \in S$, then $vars(\sigma) \subseteq vars(\sigma^*)$.*

*Proof.* This lemma follows directly from Lemma 13.

*Proof (of Proposition 19).* We prove both directions.

*First direction:* If $\sigma^*$ is a variable-minimal prime implicate of $S$, then $\sigma^*$ is a variable-minimal prime implicate of $S \setminus \{\sigma\}$. Let $\sigma^*$ be a variable-minimal prime implicate of $S$. Our goal is to show that (1) $\sigma^*$ is a prime implicate of $S \setminus \{\sigma\}$ and (2) there does not exist another prime implicate $\sigma^+$ of $S \setminus \{\sigma\}$ such that $vars(\sigma^+) \subset vars(\sigma^*)$.

To prove (1), it suffices to show that (1a) $\sigma^*$ is an implicate of $S \setminus \{\sigma\}$ and (1b) $\sigma^*$ is not subsumed by any other implicate of $S \setminus \{\sigma\}$. To prove (1a), we first recall that $vars(\sigma) \supset vars(\sigma')$ for some clause $\sigma' \in S$ by the conditions of Proposition 19. Since $\sigma^*$ is a prime implicate of $S$, it must be derivable from $S$ using resolution by Proposition 18. Suppose there is a resolution proof of $\sigma^*$ that involves clause $\sigma$. We will now show a contradiction, therefore establishing $\sigma^*$ as an implicate of $S \setminus \{\sigma\}$. First, $\sigma^*$ is descendant resolvent of $\sigma$ in this case so $vars(\sigma^*) \supseteq vars(\sigma)$ by Lemma 14. This implies that $vars(\sigma^*) \supseteq vars(\sigma) \supset vars(\sigma')$ for some clause $\sigma' \in S$. If $\sigma'$ is a prime implicate of $S$, then $\sigma^*$ cannot be a variable-minimal prime implicate of $S$ since $vars(\sigma^*) \supset vars(\sigma')$. If $\sigma'$ is not a prime implicate of $S$, then it must be subsumed by some prime implicate of $S$ which must mention a subset of the variables in $\sigma'$ so $\sigma^*$ cannot be a variable-minimal prime implicate of $S$. In either case, we have a contradiction. Hence, $\sigma^*$ can be derived from $S$ using resolution without involving clause $\sigma$.

This means that $\sigma^*$ is an implicate of $S \setminus \{\sigma\}$ so (1a) holds. We next show (1b) by contradiction. Suppose $\sigma^*$ is subsumed by some other implicate $\sigma^{**}$ of $S \setminus \{\sigma\}$. Then $\sigma^*$ cannot be a prime implicate of $S$ as it is subsumed by $\sigma^{**}$ which must also be an implicate of $S$. This is a contradiction so $\sigma^*$ is not subsumed by any other implicate of $S \setminus \{\sigma\}$ and (1b) holds. Hence, (1) holds.

We now prove (2) by contradiction. Suppose $\sigma^+$ is a prime implicate of $S \setminus \{\sigma\}$ such that $vars(\sigma^+) \subset vars(\sigma^*)$. Since $\sigma^+$ is an implicates of $S \setminus \{\sigma\}$, it is also an implicate of $S$. Hence, either $\sigma^+$ is a prime implicate of $S$ or a clause that subsumes $\sigma^+$ (mentions a subset of $\sigma^+$'s variables) is a prime implicate of $S$. Either way, $\sigma^*$ cannot be a variable-minimal prime implicate of $S$, which is a contradiction, so (2) holds.

*Second direction:* If $\sigma^*$ is a variable-minimal prime implicate of $S \setminus \{\sigma\}$, then $\sigma^*$ is a variable-minimal prime implicate of $S$. Let $\sigma^*$ be a variable-minimal prime implicate of $S \setminus \{\sigma\}$. Our goal is to show that (1) $\sigma^*$ is a prime implicate of $S$ and (2) there does not exist another prime implicate $\sigma^+$ of $S$ such that $vars(\sigma^+) \subset vars(\sigma^*)$.

To prove (1), it suffices to show that (1a) $\sigma^*$ is an implicate of $S$ and (1b) $\sigma^*$ is not subsumed by any other prime implicate of $S$. Since $\sigma^*$ is an implicate of $S \setminus \{\sigma\}$, it must be an implicate of $S$ so (1a) holds immediately. We next show (1b). Since $\sigma^*$ is a prime implicate of $S \setminus \{\sigma\}$, it cannot be subsumed by any other prime implicate of $S \setminus \{\sigma\}$. Suppose $\sigma^{**}$ is a prime implicate of $S$ but not a prime implicate of $S \setminus \{\sigma\}$. Then $\sigma^{**}$ can be derived from $S$ using a resolution proof that involves $\sigma$. Hence, $\sigma^{**}$ is a descendent resolvent of $\sigma$ so $vars(\sigma^{**}) \supseteq vars(\sigma)$ by Lemma 14. Moreover, $vars(\sigma) \supset vars(\sigma')$ for some clause $\sigma' \in S$ by the conditions of Proposition 19. Since $\sigma^{**}$ subsumes $\sigma^*$ only if $vars(\sigma^{**}) \subseteq vars(\sigma^*)$ and since $vars(\sigma') \subset vars(\sigma^{**})$, then $\sigma^{**}$ cannot subsume $\sigma^*$; otherwise, $vars(\sigma') \subset vars(\sigma^*)$, which implies there is a prime implicate of $S \setminus \{\sigma\}$ that subsumes $\sigma'$ and that mentions only a strict subset of the variables in $\sigma^*$, so $\sigma^*$ cannot be a variable-minimal prime implicate of $S \setminus \{\sigma\}$ which is a contradiction. As such, $\sigma^*$ cannot be subsumed by any other prime implicate of $S$ so (1b) and (1) hold.

We next prove (2) by contradiction. Suppose there is a prime implicate $\sigma^+$ of $S$ such that $vars(\sigma^+) \subset vars(\sigma^*)$. Then, $\sigma^+$ can be derived from $S$ using resolution by Proposition 18. We consider two cases. First case: the resolution proof does not involve clause $\sigma$. Then $\sigma^+$ is an implicant of $S \setminus \{\sigma\}$. Since some clause that subsumes $\sigma^+$ must be a prime implicate of $S \setminus \{\sigma\}$ and must mention only a subset of the variables in $\sigma^+$, $\sigma^*$ cannot be a variable-minimal prime implicate of $S \setminus \{\sigma\}$ which is a contradiction. Second case: the resolution proof involves clause $\sigma$. In this case, $\sigma^+$ is a descendant resolvent of $\sigma$ so $vars(\sigma^+) \supseteq vars(\sigma)$ by Lemma 14. This further implies $vars(\sigma^+) \supset vars(\sigma')$ for some clause $\sigma' \in S$ by the conditions of Proposition 19, and also $vars(\sigma^*) \supset vars(\sigma^+) \supset vars(\sigma')$. Since $\sigma' \in S \setminus \{\sigma\}$, some prime implicate of $S \setminus \{\sigma\}$ must subsume $\sigma'$ and must mention only a subset of its variables. Therefore, $\sigma^*$ cannot be a variable-minimal prime implicate of $S \setminus \{\sigma\}$, a contradiction. We get a contradiction in both cases, so (2) holds.

## B    Path Explanations

Consider the decision tree in Figure 2 which classifies the instance $(x_1=1, x_2=1, x_3=1, x_4=1)$ as Y using the red path. This path corresponds to the term $(x_1 \in \{1,2\}, x_2 \in \{1,2\}, x_3 \in \{1\}, x_4 \in \{1\})$ which implies the formula $\Delta_Y$ for class Y. This term is normally viewed as an explanation for the decisions on instances that follow this path. However, the shorter term $(x_1 \in \{1,2\}, x_2 \in \{1,2\}, x_3 \in \{1\})$ also implies the class formula $\Delta_Y$ and can therefore be viewed as a better explanation since feature $x_4$ is irrelevant to such decisions. This phenomena was observed in [27] which introduced the notion of an abductive path explanation (APXp): a minimal subset of the literals on a path that implies the corresponding class formula. The APXp is a syntactic notion as it depends on the specific decision tree. That is, two different decision trees that represent the same classifier may lead to different APXps. This is in contrast to the notion of a GSR that we propose which is a semantic notion that depends only on the underlying classifier (i.e., its class formulas). That is, two distinct decision trees that represent the same classifier always lead to the same GSRs for any instance.
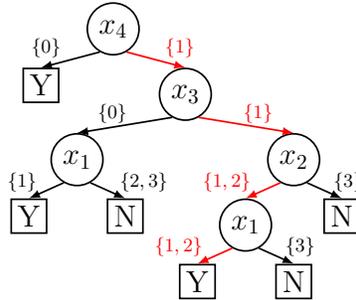


Fig. 2: A decision tree with two classes: Y and N. Variables $x_1$, $x_2$ are ternary. Variables $x_3$, $x_4$ are binary.

For the decision tree in Figure 2, the decision on instance $(x_1=1, x_2=1, x_3=1, x_4=1)$ has a GSR, $(x_1 \in \{1\}, x_2 \in \{1,2\})$, which does not correspond to any APXp of any path in the decision tree. Moreover, this GSR generates the SR $(x_1=1, x_2=1)$ as ensured by our Proposition 9.
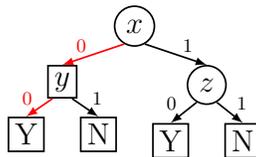


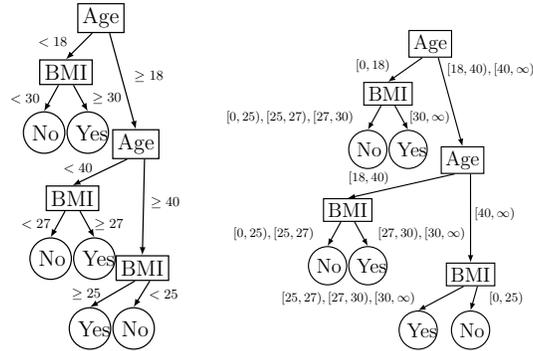Fig. 3: A decision tree with two classes: Y and N. All variables are binary.

Fig. 4: Numeric features (left) and their discretization (right).

For another example, consider Figure 3 in which all variables are binary so one does not need to go beyond simple explanations that are subsets of instances. The instance $(x=0, y=0, z=0)$ is classified as Y using the red path. This decision has two SRs, $(x=0, y=0)$ and $(y=0, z=0)$. The APXp for the red path is $(x=0, y=0)$. The APXps for the other paths are $(x=0, y=1)$, $(x=1, z=0)$, $(x=1, z=1)$. None correspond to the SR $(y=0, z=0)$. This shows that APXps cannot even generate all simple explanations (i.e., subsets of the instance). In contrast, Proposition 9 guarantees that every SR will be generated by some GSR. A similar argument applies to the notion of contrastive path explanation (CPXp) proposed in [27]. See also Example 6 in [27] for a related discussion of limitations.

## C   Numeric Features

GSRs and GNRs are particularly significant when explaining the decisions of classifiers with numeric features, such as decision trees and random forests. Consider the decision tree in Figure 4(left). One can discretize its numeric features to yield the decision tree in Figure 4(right) as is commonly practiced. For example, AGE is discretized into three intervals: $[0, 18), [18, 40)$ and $[40, \infty)$ so it can be treated as a ternary discrete variable. Similarly, BMI is discretized into four intervals: $[0, 25), [25, 27), [27, 30), [30, \infty)$. The numeric and discrete decision trees are equivalent as they make the same decision on every instance. This follows since two distinct instances will be classified equally by the numeric decision tree if the point values of their features fall into the same intervals.

The decision on instance $(\text{AGE}=42 \cdot \text{BMI}=28)$ is YES. To explain this decision, one usually works with the discrete decision tree which views this as the discrete instance $(\text{AGE}=[40, \infty) \cdot \text{BMI}=[27, 30))$, which can be notated equivalently as $(\text{AGE} \geq 40) \cdot (27 \leq \text{BMI} < 30)$. There is only one SR for the decision on this instance, which is $(\text{AGE} \geq 40) \cdot (27 \leq \text{BMI} < 30)$; that is, the instance itself. But there are two GSRs: $(\text{AGE} \geq 18 \cdot \text{BMI} \geq 27)$ and $(\text{AGE} \geq 40 \cdot \text{BMI} \geq 25)$ which are significantly more informative. SRs are quite limited in this case as they can only reference simple literals that appear in the instance: $\text{AGE}=[40, \infty)$ and $\text{BMI}=[27, 30)$. GSRs can

reference any literal implied by the instance, such as AGE $\in \{[18,40),[40,\infty)\}$, which allows them to provide more informative explanations.

The NRs for the above decision are AGE $\geq$ 40 and 27 $\leq$ BMI $<$ 30. All we can learn from the second one, as an example, is that it is possible to flip the decision by changing BMI to some value $\notin [27,30)$. If we change BMI to 32, keeping AGE the same, this NR is violated but the decision is not changed (we are only guaranteed that *some* change that violates the NR will flip the decision). In contrast, the GNRs are AGE $\geq$ 18 and BMI $\geq$ 25 which come with stronger guarantees as mentioned earlier. For example, the second GNR, BMI $\geq$ 25, tells us that changing BMI to $<$ 25, while keeping AGE the same, is guaranteed to flip the decision which is significantly more informative.

## D  More on General Reasons, GSRs and GNRs

Suppose $\Gamma$ is a general reason; $\tau_1,\ldots,\tau_n$ are the GSRs (variable-minimal prime implicants of $\Gamma$), and $\sigma_1,\ldots,\sigma_m$ are the GNRs (variable-minimal prime implicates of $\Gamma$). Then it is possible that $\Gamma \neq \sum_{i=1}^{n} \tau_i$, $\Gamma \neq \prod_{i=1}^{m} \sigma_i$ and/or $\sum_{i=1}^{n} \tau_i \neq \prod_{i=1}^{m} \sigma_i$. To illustrate this, consider the class formula $\Delta = x_1 \cdot y_1 + x_{12} \cdot y_{12} \cdot z_1$ and instance $\mathcal{I} = x_1 \cdot y_1 \cdot z_1$. The general reason is $\overline{\forall} \mathcal{I} \cdot \Delta = \Delta$. The only GSR is $x_1 \cdot y_1$ and the GNRs are $x_{12}$ and $y_{12}$. We have, $\Delta \neq x_1 \cdot y_1$; $\Delta \neq x_{12} \cdot y_{12}$; and $x_1 \cdot y_1 \neq x_{12} \cdot y_{12}$. This is different from the case for simple explanations where the disjunction of SRs, the conjunction of NRs, and the complete reason are all equivalent. Therefore, neither GSRs nor GNRs capture all the information contained in the general reason, which suggest that general reasons may have futher applications beyond GSRs and GNRs.

We now turn to another key observation. Suppose $\Gamma$ is a general reason for instance $\mathcal{I}$ and let $\sigma$ be one of its prime implicates ($\sigma$ is not necessarily variable-minimal and, hence, may not be a GNR). We can minimally change instance $\mathcal{I}$ to violate $\sigma$ yet without necessarily flipping the decision on $\mathcal{I}$. This can never happen though if $\sigma$ is variable-minimal (by Definition 8 and Proposition 12).

Consider the following example with ternary variables $X, Y, Z$, instance $\mathcal{I} = x_1 \cdot y_1 \cdot z_1$ and its class formula $\Delta = \overline{\Delta_n}$ where

$$
\begin{aligned}
\Delta_n = \ &(x_1 \cdot y_2 \cdot z_3) \ + (x_1 \cdot y_3 \cdot z_2) \ + (x_1 \cdot y_3 \cdot z_3) \ + \\
&(x_2 \cdot y_1 \cdot z_2) \ + (x_3 \cdot y_1 \cdot z_2) \ + (x_3 \cdot y_1 \cdot z_3) \ + \\
&(x_2 \cdot y_2 \cdot z_1) \ + (x_2 \cdot y_3 \cdot z_1) \ + (x_3 \cdot y_2 \cdot z_1).
\end{aligned}
$$

The general reason $\overline{\forall} \mathcal{I} \cdot \Delta$ for the decision on instance $\mathcal{I}$ is

$$
\begin{aligned}
&(\bot + y_{13} + z_{12}) \cdot (\bot + y_{12} + z_{13}) \cdot (\bot + y_{12} + z_{12}) \cdot \\
&(x_{13} + \bot + z_{13}) \cdot (x_{12} + \bot + z_{13}) \cdot (x_{12} + \bot + z_{12}) \cdot \\
&(x_{13} + y_{13} + \bot) \cdot (x_{13} + y_{12} + \bot) \cdot (x_{12} + y_{13} + \bot).
\end{aligned}
$$

which simplifies to

$$
\overline{\forall} \mathcal{I} \cdot \Delta = (y_{12} + z_1) \cdot (y_1 + z_{12}) \cdot (x_{12} + z_1) \cdot (x_1 + z_{13}) \cdot (x_{13} + y_1) \cdot (x_1 + y_{13}).
$$

Note that $\sigma = x_1 + y_1 + z_1$ is a prime implicate of $\overline{\forall} I \cdot \Delta$ which can be obtained by resolving $y_1 + z_{12}$ with $x_1 + z_{13}$ on variable $Z$. However, any instance $\mathcal{I}'$ that does not satisfy $\sigma$ is a model of $\Delta$. This follows since $\mathcal{I}' \models \overline{\sigma} = x_{23} \cdot y_{23} \cdot z_{23}$ and all models of $\overline{\Delta} = \Delta_n$ contain $x_1$, $y_1$ or $z_1$. Therefore, violating the prime implicate $\sigma$ of the general reason does not flip the decision. Note further that this prime implicate $\sigma$ is not variable-minimal (i.e., not a GNR) since $\sigma' = y_{12} + z_1$ is also a prime implicate of the general reason and $vars(\sigma') \subset vars(\sigma)$.