# Explaining Naive Bayes and Other Linear Classifiers with Polynomial Time and Delay

Joao Marques-Silva[1], Thomas Gerspacher[1], Martin C. Cooper[2,1],
Alexey Ignatiev[3], and Nina Narodytska[4]

[1] ANITI, Université de Toulouse, France
joao.marques-silva@irit.fr, thomas.gerspacher@irit.fr
[2] IRIT, Université de Toulouse III, France, cooper@irit.fr
[3] Monash University, Australia, alexey.ignatiev@monash.edu
[4] VMware Research, CA, USA, nnarodytska@vmware.com

**Abstract.** Recent work proposed the computation of so-called PI-explanations of Naive Bayes Classifiers (NBCs) [29]. PI-explanations are subset-minimal sets of feature-value pairs that are sufficient for the prediction, and have been computed with state-of-the-art exact algorithms that are worst-case exponential in time and space. In contrast, we show that the computation of one PI-explanation for an NBC can be achieved in log-linear time, and that the same result also applies to the more general class of linear classifiers. Furthermore, we show that the enumeration of PI-explanations can be obtained with polynomial delay. Experimental results demonstrate the performance gains of the new algorithms when compared with earlier work. The experimental results also investigate ways to measure the quality of heuristic explanations.

## 1 Introduction

Approaches proposed in recent years for computing explanations of Machine Learning (ML) models can be broadly characterized as *heuristic* or *non-heuristic*[5]. Heuristic approaches denote those providing *no* formal guarantees on their results. In contrast, non-heuristic approaches *do* provide some sort of formal guarantee(s) on their results, usually at the cost of increased computational complexity. Among the heuristic approaches for finding explanations, two have been studied in greater detail. One line of work focuses on devising model-agnostic linear approximations of the underlying model [24,15]. Another line of work is exemplified by Anchor [25], and targets the computation of a set of feature-value pairs associated with a given instance as a way of explaining the prediction. To date, all non-heuristic methods have focused on computing sets of feature-value pairs that are sufficient for the prediction [29,9,30,5][6]. Moreover, in terms of formal guarantees, [29] studies two

---

[5] There is a large body of recent work on explaining ML models. Example recent overviews include [8,26,27,18,17,1,19,33,20].

[6] Earlier work imposed the additional restriction of considering boolean-valued features. Clearly, non-boolean features can be binarized, e.g. with the one hot encoding, at the cost of adding additional features.

distinct definitions of explanations. A PI-explanation represents a subset-minimal set of feature values that entails the outcome of the decision function for the predicted class whatever the values of the other features (i.e. it represents a *prime implicant* of the outcome of the decision function). PI-explanations have also been studied under the name of abductive explanations [9]. In contrast, and assuming binary features, an MC-explanation is a cardinality-minimal set of equal-valued features that entails the outcome of the decision function. Non-heuristic approaches are model-based, and so earlier work specifically considered Naive-Bayes Classifiers (NBCs) and Latent-Tree Classifiers (LTCs) [29,5], Bayesian Network Classifiers [30,5], and Neural Networks [9].

In the concrete case of computing (non-heuristic) PI-explanations for NBCs, earlier work [29] proposed algorithms that are worst-case exponential in both time and space. In contrast, in this paper we propose a novel non-heuristic solution for computing PI-explanations of NBCs and other linear classifiers [7], which exhibits two fundamental advantages over earlier work. First, the paper shows that computing PI-explanations for NBCs (but also for any linear classifier) is in P, by proposing a log-linear algorithm for computing one smallest size PI-explanation. Second, the paper proposes a polynomial (log-linear) delay algorithm for enumerating the PI-explanations of NBCs (and also of any linear classifier). Furthermore, the paper presents an experimental evaluation of different approaches for explaining NBCs with PI-explanations, including the heuristic solutions computed by Anchor [25] and SHAP [15][8]. Moreover, although (real-valued) linear classifiers can be viewed as interpretable [24], this does not equate with computing PI-explanations, particularly when features are categorical. To the best of our knowledge, proving the (polynomial) complexity of computing PI-explanations for linear classifiers (including NBCs) closes an open problem.

The paper is organized as follows. Section 2 introduces the concepts and notation used throughout the paper. Section 3 introduces XLCs (a simple extension of linear classifiers (LCs)), and develops a new approach for computing, in polynomial time, one PI-explanation for XLCs. Section 3 also proposes a polynomial delay algorithm for the enumeration of PI-explanations of XLCs. Section 4 compares dedicated approaches for explaining NBCs [29] with the algorithms proposed in this paper, but also with the explanations produced by heuristic approaches. The paper concludes in Section 5.

---

[7] In fact, the paper considers a generalization of linear classifiers, that accommodates both real-valued and categorical features, which serves to streamline the presentation. This generalization will be referred to as an *eXtended Linear Classifier* (XLC).

[8] It should be noted that for linear classifiers (including NBCs), heuristic explanation approaches based on linear approximations, such as those provided by LIME [24] or SHAP [15], can be regarded as uninteresting, since the model is itself linear. Nevertheless, aiming for coverage, we opt to include also results for SHAP.

## 2   Preliminaries

*Explanations of ML models.* We consider a classification problem with two classes $\mathcal{K} = \{\oplus, \ominus\}$, defined on a set of features (or attributes) $e_1, \ldots, e_n$, which will be represented by their indices $\mathcal{E} = \{1, \ldots, n\}$. The features can either be real-valued or categorical. For real-valued features, we have $\lambda_i \leq e_i \leq \mu_i$, where $\lambda_i$, $\mu_i$ are given lower and upper bounds. For categorical features, we have $e_i \in \{1, \ldots, d_i\}$. A concrete assignment to the features referenced by $\mathcal{E}$ is represented by an $n$-dimensional vector $\mathbf{a} = (a_1, \ldots, a_n)$, where $a_j$ denotes the value assigned to feature $j$, represented by variable $e_j$, such that $a_j$ is taken from the domain of $e_j$. The set of all $n$-dimensional vectors denotes the *feature space* $\mathbb{E}$. Given a classifier with features $\mathcal{E}$, a *decision function* [29] is a mapping from the feature space to the set of classes, i.e. $\tau : \mathbb{E} \to \mathcal{K}$. For example, for a linear classifier, the decision function picks $\oplus$ if $\sum_i w_i e_i > 0$, and $\ominus$ if $\sum_i w_i e_i \leq 0$. Given $\mathbf{a} \in \mathbb{E}$, with $\tau(\mathbf{a}) = \oplus$, we consider the set of feature literals of the form $(e_i = a_i)$, where $e_i$ denotes a variable and $a_i$ a constant. A PI-explanation [29] is a subset-minimal set $\mathcal{P} \subseteq \mathcal{E}$, denoting feature literals, such that,

$$\forall (\mathbf{e} \in \mathbb{E}). \bigwedge_{j \in \mathcal{P}} (e_j = a_j) \ \to \ \tau(\mathbf{e}) = \oplus \tag{1}$$

is true. Alternatively, we can represent (1) as a rule:

$$\textbf{IF } \bigwedge_{j \in \mathcal{P}} (e_j = a_j) \textbf{ THEN } \tau(\mathbf{e}) = \oplus \tag{2}$$

(The same definitions apply in the case of class $\ominus$ (given $\mathbf{a} \in \mathbb{E}$, with $\tau(\mathbf{a}) = \ominus$).)

*Naive Bayes Classifier (NBC).* NBCs [6] can be viewed as special cases of Bayesian Network Classifiers (BNCs) [7], that make strong conditional independence assumptions among the features. Graphically, NBCs are represented as depicted in Figure 1 for a concrete example. Given some evidence $\mathbf{e}$ (in our case, this is an assignment to the features), the predicted class is given by:

$$\tau(\mathbf{e}) = \mathrm{argmax}_{c \in \mathcal{K}} \left( \Pr(c | \mathbf{e}) \right) \tag{3}$$

It is well known that $\Pr(c|\mathbf{e})$ can be computed as follows: $\Pr(c|\mathbf{e}) = \frac{\Pr(c, \mathbf{e})}{\Pr(\mathbf{e})}$. However, $\Pr(\mathbf{e})$ is constant for every $c \in \mathcal{K}$. Hence, (3) can be rewritten as follows:

$$\tau(\mathbf{e}) = \mathrm{argmax}_{c \in \mathcal{K}} \left( \Pr(c, \mathbf{e}) \right) \tag{4}$$

Finally, assuming features to be mutually conditional independent, (4) can be rewritten as follows:

$$\tau(\mathbf{e}) = \mathrm{argmax}_{c \in \mathcal{K}} \left( \Pr(c) \times \prod_i \Pr(e_i | c) \right) \tag{5}$$

A standard transformation is to apply logarithms, thus getting:

$$\tau(\mathbf{e}) = \mathrm{argmax}_{c \in \mathcal{K}} \left( \log \Pr(c) + \sum_i \log \Pr(e_i | c) \right) \tag{6}$$
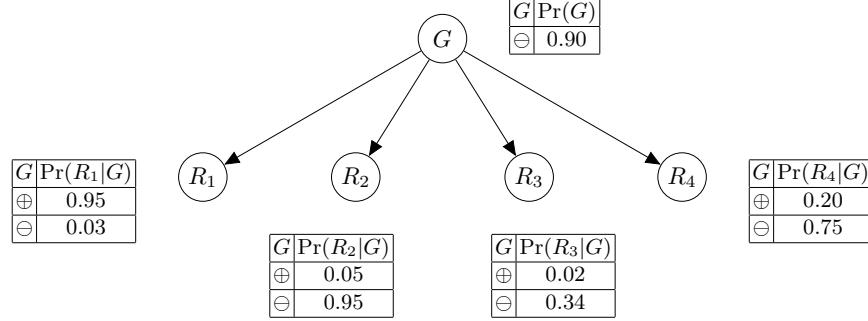
Fig. 1: Running example.

Also, if $\Pr(e_i|c) = 0$, then we use instead a sufficiently large negative value $\mathbb{M}$ [22] [9], i.e. we pick $\max(\mathbb{M}, \log(\Pr(e_i|c))) \in [\mathbb{M}, 0]$. (A simple solution is to use the sum of the logarithms of all the non-zero probabilities plus some $\epsilon < 0$.) For simplicity, i.e. to work with positive values, we can add a sufficiently large positive threshold $\mathbb{T}$ to each probability, to serve as a reference, thus obtaining:

$$\tau(\mathbf{e}) = \text{argmax}_{c \in \mathcal{K}} \left( (\mathbb{T} + \log \Pr(c)) + \sum_i (\mathbb{T} + \log \Pr(e_i|c)) \right) \qquad (7)$$

(For example, we can set $\mathbb{T}$ to the complement of the negative value with the largest absolute value.) Also for simplicity, we use the notation $\text{lPr}(\alpha) \triangleq \mathbb{T} + \max(\mathbb{M}, \log(\Pr(\alpha)))$.

*Running Example.* Consider the NBC shown in Figure 1 [10]. The features are the random variables $R_1$, $R_2$, $R_3$ and $R_4$. Each $R_i$ can take values $\mathbf{t}$ or $\mathbf{f}$ denoting, respectively, whether a listener likes or not that radio station. Random variable $G$ denotes an age class, which can take values $\mathsf{Y}$ and $\mathsf{O}$, denoting young and older listeners, respectively. Using the notation proposed earlier, we will use $\oplus$ for $\mathsf{Y}$ and $\ominus$ for $\mathsf{O}$. We also associate $\oplus$ with 1 or $\mathbf{t}$ and $\ominus$ with 0 or $\mathbf{f}$. In general we have,

$$\Pr(G, R_1, R_2, R_3, R_4) =$$
$$\Pr(G) \times \Pr(R_1|G) \times \Pr(R_2|G) \times \Pr(R_3|G) \times \Pr(R_4|G) \qquad (8)$$

Considering the assignment $(G, R_1, R_2, R_2, R_3) = (\oplus, \mathbf{t}, \mathbf{f}, \mathbf{t}, \mathbf{f})$, and using $g$ to denote $G = \oplus$, $r_i$ to denote $R_i = \mathbf{t}$ and $\neg r_i$ to denote $R_i = \mathbf{f}$, (8) can be written as follows:

$$\Pr(g, r_1, \neg r_2, r_3, \neg r_4) = \Pr(g) \times \Pr(r_1|g) \times \Pr(\neg r_2|g) \times \Pr(r_3|g) \times \Pr(\neg r_4|g)$$

---

[9] This section follows [22] throughout. An alternative would be to use Laplace smoothing [16].

[10] This example of an NBC is adapted from [2, Ch.10], with some of the conditional probabilities changed.

| | $\Pr(g)$ | $\Pr(r_1\|g)$ | $\Pr(\neg r_2\|g)$ | $\Pr(r_3\|g)$ | $\Pr(\neg r_4\|g)$ | $\mathrm{lPr}(\oplus\|\mathbf{a})$ |
|---|---|---|---|---|---|---|
| $\Pr(\cdot)$ | 0.10 | 0.95 | 0.95 | 0.02 | 0.80 | |
| $\mathrm{lPr}(\cdot)$ | 1.70 | 3.95 | 3.95 | 0.09 | 3.78 | 13.47 |

(a) Computing $\mathrm{lPr}(\oplus\|\mathbf{a})$

| | $\Pr(\neg g)$ | $\Pr(r_1\|\neg g)$ | $\Pr(\neg r_2\|\neg g)$ | $\Pr(r_3\|\neg g)$ | $\Pr(\neg r_4\|\neg g)$ | $\mathrm{lPr}(\ominus\|\mathbf{a})$ |
|---|---|---|---|---|---|---|
| $\Pr(\cdot)$ | 0.90 | 0.03 | 0.05 | 0.34 | 0.25 | |
| $\mathrm{lPr}(\cdot)$ | 3.89 | 0.49 | 1.00 | 2.92 | 2.61 | 10.91 |

(b) Computing $\mathrm{lPr}(\ominus\|\mathbf{a})$

Fig. 2: Deciding prediction for $\mathbf{a} = (\mathbf{t}, \mathbf{f}, \mathbf{t}, \mathbf{f})$

Let us consider $\mathbf{a} = (R_1, R_2, R_3, R_4) = (\mathbf{t}, \mathbf{f}, \mathbf{t}, \mathbf{f})$. Since all probabilities are strictly positive, we set $\mathbb{M}$ to a very large negative (irrelevant) value. In addition, we set $\mathbb{T}$ to a value above the complement of the logarithm of the smallest probability (i.e. 0.02), e.g we can set $\mathbb{T} = 4 > -\log(0.02)$. Using (7), we get the values shown in Figure 2. As can be concluded, the prediction will be $\oplus$. Observe that neither the value of $\mathbb{M}$ nor of $\mathbb{T}$ affect the prediction.

## 3   Explaining Extended Linear Classifiers

This section first introduces Extended Linear Classifiers (XLCs) and then details how PI-explanations can be computed for predictions of XLCs.

### 3.1   Extended Linear Classifiers

Let $\mathcal{E}$ be partitioned into $\mathcal{R}$ and $\mathcal{C}$, denoting respectively the real-valued and the categorical features. Each real-valued feature with index $i \in \mathcal{R}$ takes bounded values $\lambda_i \le e_i \le \mu_i$. For each categorical feature $j \in \mathcal{C}$, $e_j \in \{1, \ldots, d_j\}$.
We consider an XLC, that encompasses real-valued and categorical features. Let,

$$\nu(\mathbf{e}) \triangleq w_0 + \sum\nolimits_{i \in \mathcal{R}} w_i e_i + \sum\nolimits_{j \in \mathcal{C}} \sigma(e_j, v_j^1, v_j^2, \ldots, v_j^{d_j}) \qquad (9)$$

$\sigma$ is a selector function that picks the value $v_j^r$ iff $e_j$ takes value $r$. Moreover, let us define the decision function, $\tau(\mathbf{e}) = \oplus$ if $\nu(\mathbf{e}) > 0$ and $\tau(\mathbf{e}) = \ominus$ if $\nu(\mathbf{e}) \le 0$.

*Reducing linear classifiers to XLCs.* For a linear classifier, with only real-valued features, simply set $\mathcal{C} = \emptyset$. For an NBC with boolean features[11], we consider a different reduction with $\mathcal{R} = \emptyset$, starting from (7). Moreover, the argmax operator in (7) can be replaced by an inequality, from which we get,

$\mathrm{lPr}(\oplus) - \mathrm{lPr}(\ominus) +$
$\quad \sum_{i=1}^{n} (\mathrm{lPr}(e_i|\oplus) - \mathrm{lPr}(e_i|\ominus))e_i + \sum_{i=1}^{n} (\mathrm{lPr}(\neg e_i|\oplus) - \mathrm{lPr}(\neg e_i|\ominus))\neg e_i > 0 \qquad (10)$

---

[11] Given the proposed reductions, it is immediate to represent an NBC with categorical features as an XLC.

| $w_0$ | $v_1^1$ | $v_1^2$ | $v_2^1$ | $v_2^2$ | $v_3^1$ | $v_3^2$ | $v_4^1$ | $v_4^2$ |
|---|---|---|---|---|---|---|---|---|
| -2.19 | -2.97 | 3.46 | 2.95 | -2.95 | 0.4 | -2.83 | 1.17 | -1.32 |

(a) Example reduction of NBC to XLC (Example 1)

| $\Gamma$ | $\delta_1$ | $\delta_2$ | $\delta_3$ | $\delta_4$ | $\Phi$ |
|---|---|---|---|---|---|
| 2.56 | 6.43 | 5.90 | 0.00 | 2.49 | 12.26 |

(b) Computing $\delta_j$'s for the XLC (Example 2)

Fig. 3: Values used in the running example (Example 1 and Example 2)

The reduction is completed by setting: $w_0 \triangleq \mathrm{lPr}(\oplus) - \mathrm{lPr}(\ominus)$, $v_j^1 \triangleq \mathrm{lPr}(\neg e_j | \oplus) - \mathrm{lPr}(\neg e_j | \ominus)$, $v_j^2 \triangleq \mathrm{lPr}(e_j | \oplus) - \mathrm{lPr}(e_j | \ominus)$, and $d_j \triangleq 2$.

*Example 1.* Figure 3a shows the resulting XLC formulation for the example in Figure 2. We also let $\mathbf{f}$ be associated with value 1 and $\mathbf{t}$ be associated with value 2, and $d_j = 2$.

### 3.2   Explaining XLCs

We now investigate how (smallest or cardinality-minimal) PI-explanations can be computed for XLCs, and also how (minimal) PI-explanations can be enumerated. For this, we need to assess how *free* some of the features are. For a given instance $\mathbf{e} = \mathbf{a}$, define a *constant* slack (or gap) value $\Gamma$ given by,

$$\Gamma \triangleq \nu(\mathbf{a}) = w_0 + \sum_{i \in \mathcal{R}} w_i a_i + \sum_{j \in \mathcal{C}} \sigma(a_j, v_j^1, v_j^2, \ldots, v_i^{d_j}) \qquad (11)$$

i.e. this is the value obtained when deciding $\oplus$ to be the picked class, given the assignment $\mathbf{e} = \mathbf{a}$.

We are interested in computing one PI-explanation [29] of an XLC, but we are also interested in enumerating PI-explanations. As argued in Section 2, this corresponds to finding a subset-minimal set of literals $\mathcal{P} \subseteq \mathcal{E}$ such that (1) holds, or alternatively,

$$\forall (\mathbf{e} \in \mathbb{E}). \bigwedge_{j \in \mathcal{P}} (e_j = a_j) \rightarrow (\nu(\mathbf{e}) > 0) \qquad (12)$$

under the assumption that $\nu(\mathbf{a}) > 0$. In what follows, we partition $\mathcal{E}$ into $\mathcal{P}$ and $\mathcal{N}$, respectively the picked and the non-picked attributes from $\mathcal{E}$.

*Categorical case.* Let us first consider $\mathcal{R} = \emptyset$. Each feature $e_j$ is assigned value $a_j$, which results in selecting some value $v_j^{a_j}$, i.e. the value from the weights associated with $e_j$ which is picked when $e_j = a_j$. Thus, $\Gamma$ is computed as follows: $\Gamma = w_0 + \sum_{j \in \mathcal{C}} v_j^{a_j}$.

Moreover, let $v_j^{\omega}$ denote the *smallest* (or *worst-case*) value associated with $e_j$. Then, by letting every $e_j$ take *any* value, the *worst-case* value of $\nu(\mathbf{e})$ is,

$$\Gamma^{\omega} = w_0 + \sum_{j \in \mathcal{C}} v_j^{\omega} \qquad (13)$$

We are interested in cases where $\Gamma^\omega \leq 0$, corresponding to predicting $\ominus$ instead of $\oplus$. (Otherwise the prediction would not change from $\oplus$.) The expression above can be rewritten as follows,

$$\begin{aligned}
\Gamma^\omega &= w_0 + \sum_{j \in \mathcal{C}} v_j^{a_j} - \sum_{j \in \mathcal{C}} (v_j^{a_j} - v_j^\omega) \\
&= \Gamma - \sum_{j \in \mathcal{C}} \delta_j = -\Phi
\end{aligned} \tag{14}$$

where we use $\delta_j \triangleq v_j^{a_j} - v_j^\omega$, and $\Phi \triangleq \sum_{j \in \mathcal{C}} \delta_j - \Gamma = -\Gamma^\omega$. Our goal is to find a smallest (or subset-minimal) set $\mathcal{P}$ such that the prediction is still $\oplus$ (whatever the values of the other features):

$$w_0 + \sum_{j \in \mathcal{P}} v_j^{a_j} + \sum_{j \notin \mathcal{P}} v_j^\omega = -\Phi + \sum_{j \in \mathcal{P}} \delta_j > 0 \tag{15}$$

i.e. we want to pick a smallest (or subset-minimal) set of literals that ensures that the prediction will be $\oplus$. In turn, (15) can be represented as the following optimization problem:

$$\begin{aligned}
\min \quad & \sum_{i=1}^n p_i \\
\text{s.t.} \quad & \sum_{i=1}^n \delta_i p_i > \Phi \\
& p_i \in \{0, 1\}
\end{aligned} \tag{16}$$

where the variables $p_i$ assigned value 1 denote the indices included in $\mathcal{P}$. Although solving (16) seems to equate to solving an NP-hard optimization, concretely the minimization version of the knapsack problem [12], the fact that the coefficients in the cost function are all equal to 1 makes the problem solvable in log-linear time[12]. Concretely, we can now develop a greedy algorithm that computes a smallest PI-explanation, representing one optimal solution of (16). At each step, we simply pick the largest $\delta_i$ that has not yet been picked.

**Proposition 1.** *Let $\mathcal{S} = \langle l_1, \ldots, l_n \rangle$ represent indices of $\mathcal{E}$ sorted by non-increasing value of $\delta_j$. Pick $k$ such that $\sum_{j \in \{l_1, \ldots, l_k\}} \delta_j > \Phi$ and $\sum_{j \in \{l_1, \ldots, l_{k-1}\}} \delta_j \leq \Phi$. Then (12) holds for $\mathcal{P} = \{p_{l_r} | 1 \leq r \leq k\}$, and $\mathcal{P}$ represents an optimal solution of (16).*

Optimality of the computed solution is given by Proposition 1 (proof included in Section A.2).

*Example 2.* Figure 3b shows the values used for computing explanations for the example in Figure 2.
For this example, the sorted $\delta_j$'s become $\langle \delta_1, \delta_2, \delta_4, \delta_3 \rangle$. By picking $\delta_1$ and $\delta_2$, we ensure that the prediction is $\oplus$, independently of the evidence provided for features $e_3$ and $e_4$. Thus $(e_1) \wedge (\neg e_2)$ is a PI-explanation for the NBC shown in Figure 1, with evidence $(e_1, e_2, e_3, e_4) = (\mathbf{t}, \mathbf{f}, \mathbf{t}, \mathbf{f})$. (It is easy to observe that $\tau(\mathbf{t}, \mathbf{f}, \mathbf{f}, \mathbf{f}) = \tau(\mathbf{t}, \mathbf{f}, \mathbf{f}, \mathbf{t}) = \tau(\mathbf{t}, \mathbf{f}, \mathbf{t}, \mathbf{f}) = \tau(\mathbf{t}, \mathbf{f}, \mathbf{t}, \mathbf{t}) = \oplus$).

---

[12] Pseudo-polynomial time algorithms for the knapsack problem are well-known [4,21]. One concrete example [21] yields a polynomial (cubic) time algorithm in the setting of computing a smallest PI-explanation of an XLC. We show that it is possible to devise a more efficient solution.

**Function** OneExplanation(Vs,Flip,$\Delta$,$\Phi^R$,Idx,Xpl) ;

    **Input:** Vs: Values of instance being explained; Flip: Array reference of
                decision steps; $\Delta$: Sorted $\delta_j$'s; $\Phi^R$: Explanation threshold; Idx: Index
                for $\Delta$; Xpl: Set reference of explanation literals
    **Output:** $\Phi^R$: Updated threshold; Idx: Updated index for $\Delta$

**1**       **while** $\Phi^R \geq 0$ **do**
**2**            Idx $\leftarrow$ Idx $+ 1$ ;
**3**            Flip[Idx] $\leftarrow 0$ ;
**4**            $\Phi^R \leftarrow \Phi^R - \Delta$[Idx] ;
**5**            Xpl $\leftarrow$ Xpl $\cup \{(e_{\mathsf{Idx}}, \mathsf{Vs}[\mathsf{Idx}])\}$ ;
**6**        ReportExplanation (Xpl) ;
**7**       **return** $(\Phi^R, \mathsf{Idx})$ ;

**Algorithm 1:** Finding one explanation

In the concrete case of NBCs, if the goal is to compute a single explanation, then the algorithm detailed in this section is exponentially more efficient (in the worst case) than earlier work [29]. However, in some settings one wants to be able to analyze some or even all explanations for a given instance (this is further discussed in Section 4). We describe next a polynomial (log-linear) delay algorithm for enumeration of explanations for XLCs (and so for NBCs).

*Enumerating explanations with polynomial delay.* As shown above, a smallest PI-explanation can be computed in log-linear time by sorting the $\delta_i$ values and picking the first $k$ literals that ensure the prediction. We start by presenting a more elaborate description of the algorithm, which we then use for devising the enumeration of explanations with polynomial delay[13]. Algorithm 1 shows the pseudo-code for computing one smallest explanation. $\Delta$ denotes the array of sorted $\delta_j$'s. (The pseudo-code assumes that the order $1, 2, \ldots, n$ represents the literals in sorted order.) $\Phi^R$ is initialized with the value of $\Phi$, being updated as the algorithm(s) progress(es). Algorithm 1 corresponds to the direct application of Proposition 1. This algorithm can now be exploited for implementing a polynomial delay algorithm for enumerating PI-explanations. Algorithm 2 depicts the enumeration of PI-explanations. The algorithm implements a (restricted) backtrack search procedure, which in some circumstances can be shown to yield polynomial delay algorithms [3]. Idx denotes the depth of the search tree and Flip (if assigned 0) records which $\delta_j$'s are used for updating $\Phi^R$. (The entries of Flip take value -1 if unused, and value 1 if have been backtracked upon.) A key aspect of the algorithm is that it only branches when it is guaranteed that a PI-explanation can still be found, given the prefix (of picked or not picked $\delta_j$'s) defined by Flip and Idx. Otherwise, the algorithm must backtrack and enter a consistent state (with at most a linear backtracking effort). Algorithm 3 shows the backtrack step of the PI-enumeration algorithm. Algorithm 3 terminates if no more PI-explanations can be found, or with the guarantee that another PI-

---

[13] For a knapsack constraint, it is known that feasible solutions can be enumerated with quadratic delay [14,10]. Nevertheless, we exploit the problem's special structure to achieve a log-linear enumeration delay.

**Function** ALLEXPLANATIONS($\mathsf{Vs}, \Delta, \Phi^R$) ;
    **Input:** $\mathsf{Vs}$: Values of instance being explained; $\Delta$: Sorted $\delta_j$'s;
        $\Phi^R$: Explanation threshold

**1**     $(\mathsf{Xpl}, \mathsf{Flip}, \mathsf{Idx}) \leftarrow (\emptyset, [-1, \ldots, -1], 0)$ ;
**2**     **while** $\mathsf{Idx} \geq 0$ **do**
**3**         $(\Phi^R, \mathsf{Idx}) \leftarrow \text{ONEEXPLANATION}(\mathsf{Vs}, \mathsf{Flip}, \Delta, \Phi^R, \mathsf{Idx}, \mathsf{Xpl})$ ;
**4**         $(\Phi^R, \mathsf{Idx}) \leftarrow \text{ENTERVALIDSTATE}(\mathsf{Vs}, \mathsf{Flip}, \Delta, \Phi^R, \mathsf{Idx}, \mathsf{Xpl})$ ;

**Algorithm 2:** Finding all explanations

**Function** ENTERVALIDSTATE($\mathsf{Vs}, \mathsf{Flip}, \Delta, \Phi^R, \mathsf{Idx}, \mathsf{Xpl}$) ;
    **Input:** $\mathsf{Vs}$: Values of instance being explained; $\mathsf{Flip}$: Array reference of
        decision steps; $\Delta$: Sorted $\delta_j$'s; $\Phi^R$: Explanation threshold; $\mathsf{Idx}$: Index
        for $\Delta$; $\mathsf{Xpl}$: Set reference of explanation literals
    **Output:** $\Phi^R$: Updated threshold; $\mathsf{Idx}$: Updated index for $\Delta$

**1**     **while** $\Phi^R < 0$ **or** $\sum_{i=\mathsf{Idx}}^{n} \Delta[i] < \Phi^R$ **do**
**2**         **while** $\mathsf{Idx} \geq 0 \wedge \mathsf{Flip}[\mathsf{Idx}] = 1$ **do**
**3**             $\mathsf{Flip}[\mathsf{Idx}] \leftarrow -1$ ;
**4**             $\mathsf{Idx} \leftarrow \mathsf{Idx} - 1$ ;
**5**         **if** $\mathsf{Idx} < 0$ **then return** $(\Phi^R, \mathsf{Idx})$ ;
**6**         $\mathsf{Xpl} \leftarrow \mathsf{Xpl} \setminus \{(e_{\mathsf{Idx}}, \mathsf{Vs}[\mathsf{Idx}])\}$;
**7**         $\Phi^R \leftarrow \Phi^R + \Delta[\mathsf{Idx}]$;
**8**         $\mathsf{Flip}[\mathsf{Idx}] \leftarrow 1$;
**9**     **return** $(\Phi^R, \mathsf{Idx})$ ;

**Algorithm 3:** Entering a valid state

explanation can be extracted with Algorithm 1. It is straightforward to conclude that both Algorithm 1 and Algorithm 3 run in linear time on the size of the current depth of the search tree (which is linear on the number of features). Thus, we can list PI-explanations of XLC's with polynomial delay (proof included in Section A.2).

**Proposition 2.** *PI-explanations of an XLC can be enumerated with log-linear delay.*

*Real-valued & mixed case.* Let us now consider $\mathcal{R} \neq \emptyset$. As before, the prediction is assumed to be $\oplus$. For each feature, if $w_i > 0$, then we are interested in assessing the impact of reducing the value of $e_i$. Hence, the worst-case scenario is achieved when $e_i = \lambda_i$. In this case, we define $\delta_i = (a_i - \lambda_i)w_i$. A no-change constraint on the value of $e_i$ is formulated as $e_i \geq a_i$ (i.e. we *clamp* the value of $e_i$ by imposing a lower bound on its value). In contrast, if $w_i < 0$, then we are interested in assessing the impact of increasing the value of $e_i$. The worst-case scenario is now $e_i = \mu_i$. In this case, we define $\delta_i = (a_i - \mu_i)w_i$. Moreover, a no-change constraint on the value of $e_i$ is formulated as $e_i \leq a_i$ (i.e. in this case we *clamp* the value of $e_i$ by imposing an upper bound on its value). Given the definition of the $\delta_i$ constants for real-valued features, and associated literals in case of a no-change constraint, we can compute explanations using the restricted knapsack problem formulation as above. Thus, we can also compute one cardinality optimal solution in log-linear time, and enumerate subset-minimal solutions with polynomial delay.

## 4   Experimental Evaluation

This section evaluates the PI-explanation enumerator XPXLC, that implements the algorithms described in this paper[14]. XPXLC was tested in Debian Linux on an Intel Xeon CPU 5160 3.00 GHz with 64 GByte of memory. When testing scalability, XPXLC was run with 8GByte limit on RAM and two hours time limit. The experiment was divided into 3 parts: (1) evaluating the raw performance of XPXLC, (2) comparing it with the state-of-the-art compilation approach STEP [29,30], and (3) using complete enumeration of PI-explanations to assess the quality of explanations of the well-known heuristic explainers Anchor [25] and SHAP [15].

**Datasets.** We selected a set of widely-used, publicly available, datasets from [31,23,11]. The total number of datasets used is 37. For each dataset, we trained a Naive Bayes classifier[15] using 80% of the training data. The average test accuracy assessed for the 20% remaining instances is 77.7%. (All the datasets and the trained classifiers are available in the online repository.) The experiments targeted XPXLC's ability to enumerate a given number of explanations within a time limit.

**Raw performance.** Figure 4a shows the scalability of XPXLC. Here, XPXLC was set to compute $10^6$ distinct explanations for each instance of each dataset. For the cases having fewer than $10^6$ explanations, XPXLC terminates as soon as all explanations are computed. The smallest number of observed explanations per instance is 1, the maximum number is at least $10^6$, while on average 29207.5 PI-explanations are reported per each instance. The total number of instances to explain in this experiment is 94174. The line drawn through point $(x, y)$ in Figure 4a shows how many instances on the $X$-axis are solved by the time shown on the $Y$-axis. As can be observed, performance is not an issue for XPXLC – it never exceeds 12 seconds to enumerate $10^6$ explanations for each of the target instances. On average, XPXLC finishes complete enumeration (of at most $10^6$ explanations) in 0.23 seconds.

**Enumerative vs. compilation-based approaches.** The state of the art for finding PI-explanations for NBCs is the STEP compilation-based approach [29,30,32]. Concretely, STEP consists of (1) compilation of a BNC classifier into a *sentential decision diagram* (SDD) and (2) enumeration of PI-explanations using efficient algorithms for SDD-based prime implicant enumeration. The existing implementation of STEP can only handle binary features. Therefore, and in order to compare the relative performance of XPXLC and STEP, we apply a one-hot encoding (OHE) to categorical features, retrain the Naive Bayes classifiers and run both tools on the OHE instances[16], targeting the complete enumeration of explanations. Moreover, despite its worst-case exponential complexity in time

---

[14] The source code of XPXLC, as well as the datasets, documentation, and additional examples can be obtained from the authors.

[15] The CategoricalNB classifier of scikit-learn [28] was used for this purpose.

[16] This solution is not ideal, since the use of OHE impacts the assumption of feature independence of NBCs, and only serves to enable the comparison between STEP and XPXLC.

(a) Raw performance of XPXLC

(b) Performance of STEP (with MOs & TOs)
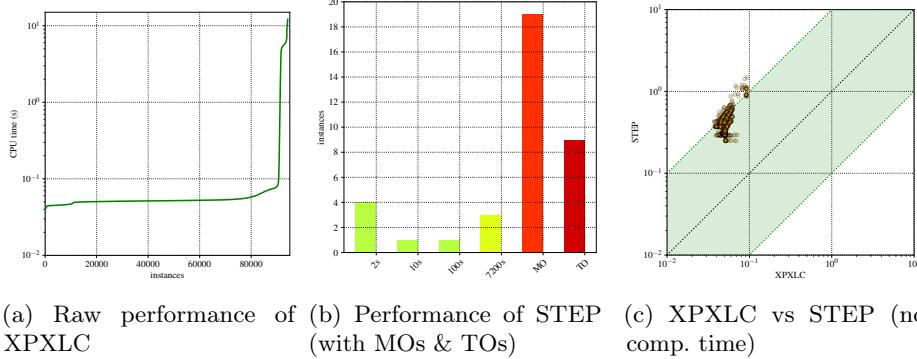
(c) XPXLC vs STEP (no comp. time)

Fig. 4: Scalability of XPXLC targeting $10^6$ PI-explanations, performance of STEP, and comparative performance of XPXLC and STEP.

and space, STEP can still compile into SDDs 9 (out of 37) NBC classifiers, i.e. close to 25% of the classifiers, within the 2 hours time limit and 32 GByte memory limit. Once an NBC classifier is compiled into an SDD, enumeration of all PI-explanations is relatively easy — concretely, it takes 0.39 seconds for the compilation-based approach to enumerate all explanations. However, the SDD compilation step itself takes between 1 and 4300 seconds for the classifiers that can be compiled. If the compilation time is amortized over all data instances of each dataset, its impact ranges from a fraction of a second to ≈50 seconds. Figure 4b shows a histogram summarizing the performance of STEP's compiler. The bars in the histogram represent the classifiers that STEP is able to compile within 2 seconds (there are 4 of them), 10 seconds (1), 100 seconds (1), 2 hours (3) and also classifiers that STEP fails to compile due to reaching the memory (MO) or time (TO) limits. The last two bars represent 19 and 9 classifiers, respectively. Finally, Figure 4c summarizes the performance comparison between XPXLC and STEP. In this comparison, the SDD compilation time is *ignored*, and the plot shows only instances for the classifiers that STEP is able to compile within the 2 hour time limit. Also note that both tools finish complete enumeration of PI-explanations for each of these instances. A point $(x, y)$ in the plot represents the time (in seconds) spent by XPXLC (shown on the $X$-axis) and by STEP (shown on the $Y$-axis) for a concrete data instance. Observe that, even if the compilation time is ignored, STEP's enumeration phase is still between 4 and 20 times slower than XPXLC.

**Assessing heuristic approaches.** Exhaustive enumeration of PI-explanations can serve to assess heuristic explanations. Exhaustive enumeration provides a distribution of how many times feature-value pairs appear in explanations, and thus which are likely to be more *relevant* for the given prediction. As a result, one can evaluate how many features in a heuristic explanation "hit" the set of most relevant (commonly-occurring) features. This strategy may be beneficial in some practical settings where trustable explanations are of concern. While our "hit" metric is a heuristic evaluation measure to compare the quality of explanations, we demonstrate its usefulness experimentally. For example, our metric does show

a strong correlation between features of heuristic explanations and common features that we identify via enumeration. Figure 5 depicts the percentage of features in explanations of Anchor [25] and SHAP [15] "hitting" the set of common features. Here, we focus on 2 datasets *Adult* [13,25] and *Spambase* [31] and use the following methodology. For an explanation $E$ of Anchor, we keep the top $|E|$ features most commonly-occurring in all PI-explanations[17]; then we count the number of features in $E$ that hit the set of common features. As SHAP assigns numerical weights to *all* features, we take 5 features reported by SHAP as most relevant and count how many of them intersect the set of 5 most common features of PI-explanations. The rationale of this choice is that larger explanations are typically harder for a user to reason about and so 5 features is normally deemed enough to make a conclusion wrt. the cause of prediction. As can be observed, both Anchor and SHAP are successful at hitting the most common features. However, in some cases both tools' explanations do not overlap our important features, e.g. Anchor has zero overlap with the common features in more than 2000 instances. Given a significant overlap in the majority of cases, a zero hit suggests that Anchor's explanation might be using less influential features and is hence less trustworthy. This experiment illustrates another setting where PI-explanations can be useful, i.e. not only to output a provably correct explanation but also to provide the user with an alternative evaluation toolkit to measure confidence in heuristic explanations. Finally, we observe that both Anchor and SHAP are significantly slower than XPXLC: on average, Anchor takes 1.55 seconds to compute one explanation of an instance, whereas SHAP takes 99.58 seconds. In contrast, as highlighted above, XPXLC never exceeds a few tens of $\mu$sec for computing a single explanation.

## 5   Conclusions

This paper presents a log-linear algorithm for computing a smallest PI-explanation of linear classifiers. Moreover, the paper shows that PI-explanations for linear classifiers can be enumerated with polynomial delay. The results in the paper also apply to NBCs (among other classifiers), and so should be contrasted with earlier work [29], which proposes a worst-case exponential time and space solution for computing PI-explanations of NBCs. A natural line of research is to investigate extensions of XLCs that also admit polynomial time algorithms for computing PI-explanations.

## References

1. S. Anjomshoae, A. Najjar, D. Calvaresi, and K. Främling. Explainable agents and robots: Results from a systematic literature review. In *AAMAS*, pages 1078–1088, 2019.

---

[17] If $> |E|$ features are in the top due to having the same frequency, all of them are marked as common. Also, the experiment is performed only for instances for which complete PI-explanation enumeration finishes.

2. D. Barber. *Bayesian reasoning and machine learning*. Cambridge University Press, 2012.
3. D. A. Cohen. Tractable decision for a constraint language implies tractable search. *Constraints An Int. J.*, 9(3):219–229, 2004.
4. G. B. Dantzig. Discrete-variable extremum problems. *Operations research*, 5(2):266–288, 1957.
5. A. Darwiche. Three modern roles for logic in AI. *CoRR*, abs/2004.08599, 2020.
6. R. O. Duda, P. E. Hart, and D. G. Stork. *Pattern classification*. John Wiley & Sons, 1973.
7. N. Friedman, D. Geiger, and M. Goldszmidt. Bayesian network classifiers. *Mach. Learn.*, 29(2-3):131–163, 1997.
8. R. Guidotti, A. Monreale, S. Ruggieri, F. Turini, F. Giannotti, and D. Pedreschi. A survey of methods for explaining black box models. *ACM Comput. Surv.*, 51(5):93:1–93:42, 2019.
9. A. Ignatiev, N. Narodytska, and J. Marques-Silva. Abduction-based explanations for machine learning models. In *AAAI*, pages 1511–1519, 2019.
10. D. S. Johnson, C. H. Papadimitriou, and M. Yannakakis. On generating all maximal independent sets. *Inf. Process. Lett.*, 27(3):119–123, 1988.
11. Kaggle Machine Learning Community. `https://www.kaggle.com/`.
12. H. Kellerer, U. Pferschy, and D. Pisinger. *Knapsack problems*. Springer, 2004.
13. R. Kohavi. Scaling up the accuracy of naive-bayes classifiers: A decision-tree hybrid. In *KDD*, pages 202–207, 1996.
14. E. L. Lawler, J. K. Lenstra, and A. H. G. Rinnooy Kan. Generating all maximal independent sets: NP-hardness and polynomial-time algorithms. *SIAM J. Comput.*, 9(3):558–565, 1980.
15. S. M. Lundberg and S. Lee. A unified approach to interpreting model predictions. In *NeurIPS*, pages 4765–4774, 2017.
16. C. D. Manning, P. Raghavan, and H. Schütze. *Introduction to information retrieval*. Cambridge University Press, 2008.
17. T. Miller. "But why?" understanding explainable artificial intelligence. *ACM Crossroads*, 25(3):20–25, 2019.
18. T. Miller. Explanation in artificial intelligence: Insights from the social sciences. *Artif. Intell.*, 267:1–38, 2019.
19. B. D. Mittelstadt, C. Russell, and S. Wachter. Explaining explanations in AI. In *FAT*, pages 279–288, 2019.
20. S. T. Mueller, R. R. Hoffman, W. J. Clancey, A. Emrey, and G. Klein. Explanation in human-AI systems: A literature meta-review, synopsis of key ideas and publications, and bibliography for explainable AI. *CoRR*, abs/1902.01876, 2019.
21. C. H. Papadimitriou and K. Steiglitz. *Combinatorial Optimization: Algorithms and Complexity*. Prentice-Hall, 1982.
22. J. D. Park. Using weighted MAX-SAT engines to solve MPE. In *AAAI*, pages 682–687, 2002.
23. Penn Machine Learning Benchmarks. `https://github.com/EpistasisLab/penn-ml-benchmarks`.
24. M. T. Ribeiro, S. Singh, and C. Guestrin. "Why should I trust you?": Explaining the predictions of any classifier. In *KDD*, pages 1135–1144, 2016.
25. M. T. Ribeiro, S. Singh, and C. Guestrin. Anchors: High-precision model-agnostic explanations. In *AAAI*, pages 1527–1535, 2018.
26. W. Samek, G. Montavon, A. Vedaldi, L. K. Hansen, and K. Müller, editors. *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*. Springer, 2019.
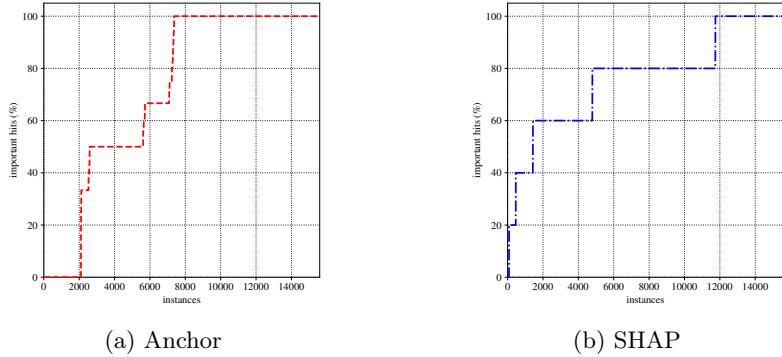
(a) Anchor

(b) SHAP

Fig. 5: Percentage of important "hits" of explanations produced by Anchor and SHAP.

27. W. Samek and K. Müller. Towards explainable artificial intelligence. In Samek et al. [26], pages 5–22.
28. scikit-learn: Machine Learning in Python. `https://scikit-learn.org/`.
29. A. Shih, A. Choi, and A. Darwiche. A symbolic approach to explaining bayesian network classifiers. In *IJCAI*, pages 5103–5111, 2018.
30. A. Shih, A. Choi, and A. Darwiche. Compiling bayesian network classifiers into decision graphs. In *AAAI*, pages 7966–7974, 2019.
31. UCI Machine Learning Repository. `https://archive.ics.uci.edu/ml`.
32. Automated Reasoning Group UCLA. `http://reasoning.cs.ucla.edu/xai/`.
33. F. Xu, H. Uszkoreit, Y. Du, W. Fan, D. Zhao, and J. Zhu. Explainable AI: A brief survey on history, research areas, approaches and challenges. In *NLPCC*, pages 563–574, 2019.

# A   Appendix

## A.1   Additional Plots

Additional plots are shown in Figure 5.

## A.2   Proofs

**Proposition 1.** *Let $\langle l_1, \ldots, l_n \rangle$ represent indices $\mathcal{E}$ sorted by non-increasing value of $\delta_j$. Pick $k$ such that $\sum_{j \in \{l_1, \ldots, l_k\}} \delta_j > \Phi$ and $\sum_{j \in \{l_1, \ldots, l_{k-1}\}} \delta_j \leq \Phi$. Then* (12) *holds for $\mathcal{P} = \{l_r | 1 \leq r \leq k\}$, and $\mathcal{P}$ represents an optimal solution of* (16).

*Proof.* We prove that an optimal solution to (16) can be obtained with the greedy algorithm that picks features in non-increasing order of $\delta_j$'s. Let $\mathcal{P}^* = \langle i_1, \ldots, i_k \rangle$ denote the $k$ indices in some optimal solution, such that $\delta_{i_1} \geq \ldots \geq \delta_{i_k}$. Moreover, let $\mathbb{V}(\mathcal{P}^*) = \sum_{j \in \{i_1, \ldots i_k\}} \delta_j$. Clearly, $\mathbb{V}(\mathcal{P}^*) > \Phi$; otherwise $\mathcal{P}^*$ would not satisfy the constraint in (16).

We prove by induction that one can construct another optimal solution $\mathcal{P} = \langle l_1, \ldots, l_k \rangle$, where $l_1, \ldots, l_k$ denote the first $k$ features with highest $\delta_j$. For the base case, we consider the first pick, and suppose that $i_1 \neq l_1$ (and so $l_1$ does not occur in $\mathcal{P}^*$). We can construct another sequence $\mathcal{P}' = \langle l_1, i_2, \ldots, i_k \rangle$, such that $\mathbb{V}(\mathcal{P}') = \sum_{j \in \{l_1, i_2, \ldots, i_k\}} \delta_j \geq \mathbb{V}(\mathcal{P}^*) > \Phi$. Hence, $\mathcal{P}'$ is still an optimal solution, and starts with a greedy choice. For the general case, we assume that the first $r-1$ picks can be made to respect the greedy choice, and that the $r^{\text{th}}$ does not. The reasoning now can be mimicked again, and so we can construct another optimal solution such that the $r^{\text{th}}$ choice is also greedy. Thus, Proposition 1 yields a smallest PI-explanation.                                                              □

**Proposition 2.** *PI-explanations of an XLC can be enumerated with log-linear delay.*

*Proof.* For simplicity of presentation, we assume that the values $\delta_i$ are sorted in non-increasing order, i.e. $\delta_1 \geq \ldots \geq \delta_n$. This sorting operation can be achieved in log-linear time. Recall that $\delta_i \geq 0$ $(i = 1, \ldots, n)$ and that a PI-explanation represented by the bit vector $p$ must satisfy the two constraints: (C1) $\sum_{i=1}^{n} \delta_i p_i > \Phi$ and (C2) $\forall j \in \{1, \ldots, n\}$ such that $p_j = 1$, $(\sum_{i=1}^{n} \delta_i p_i) - \delta_j p_j \leq \Phi$ (subset-minimality).

Consider an exhaustive depth-first binary search (DFS) in which at depth $r$ the two branches correspond to $p_r = 1$ and $p_r = 0$. It is critical for the correctness of this search that on each branch, the $p_i$ variables are instantiated in non-increasing order of the corresponding values $\delta_i$. For a depth-$r$ node $\alpha$ of this search tree, let $S_\alpha$ be the sum $\sum_{i=1}^{r} \delta_i p_i$. A node $\alpha$ is declared a leaf (and is hence not expanded) if $S_\alpha > \Phi$. Assuming that, by default, the remaining values $\delta_{r+1}, \ldots, \delta_n$ are assigned 0, node $\alpha$ satisfies (C1). Clearly, any other descendant nodes (at which at least one of $\delta_{r+1}, \ldots, \delta_n$ is 1) would not satisfy (C2) and hence does not need to be considered. This means that all PI-explanations will be found. It remains to show that all leaves $\alpha$ satisfy subset-minimality and hence are PI-explanations. To see that $\alpha$ satisfies (C2), let $\beta$ be its parent node. Since $\beta$ is not a leaf, we must have $S_\beta = S_\alpha - \delta_r p_r \leq \Phi$. But then $S_\alpha - \delta_j p_j \leq \Phi$ for all $j$ such that $p_j = 1$ since $\delta_j \geq \delta_r$ $(j = 1, \ldots, r-1)$. Thus, all leaves correspond to PI-explanations.

We add to our DFS the pruning rule that a depth-$r$ node $\alpha$ is only created if $S_\alpha + \sum_{i=r+1}^{n} \delta_i > \Phi$. This sum is calculated incrementally, so only requires $O(1)$ time at each node. The reason behind this rule is that if it is not satisfied, then no descendant of $\alpha$ can satisfy (C1). On the other hand, if this rule is satisfied then we know that at least one descendant of $\alpha$ will be a leaf (and as explained above will correspond to a PI-explanation). It is well known that a depth-first search in a search tree with no dead-end nodes provides a polynomial delay algorithm [3]. In our DFS, the delay between visiting two leaves is linear in $n$. Since finding the first PI-explanation also requires a sorting step, with a log-linear complexity, we can conclude that the worst-case delay is log-linear.                             □