



Complexity in the case against accuracy estimation

Richard Nock*

*GRIMAAG-Département Scientifique Interfacultaire, Université des Antilles-Guyane,
Campus de Schoelcher, BP 7209, 97233 Schoelcher, Martinique, France*

Received 28 February 2001; received in revised form 4 February 2002; accepted 30 May 2002

Communicated by G. Ausiello

Abstract

Some authors have repeatedly pointed out that the use of the accuracy, in particular for comparing classifiers, is not adequate. The main argument concerns some assumptions of seldom validity or correctness underlying the use of this criterion. In this paper, we study the computational burden of the accuracy's replacement for building and comparing classifiers, using the framework of Inductive Logic Programming. Replacement is investigated in three ways: completion of the accuracy with an additional requirement, replacement of the accuracy with a bi-criterion recently introduced from statistical decision theory: the Receiver Operating Characteristic analysis, and replacement of the accuracy by a single criterion. We prove very hard results for most of the possible replacements. A first result shows that allowing the arbitrary multiplication of clauses appears to be totally useless. "Arbitrary" is to be taken in its broadest meaning, in particular exponential. The second point is the sudden appearance of the negative result, which is not a function of the criteria's demands. The third point is the equivalence in difficulty of all these different criteria. In contrast, the single accuracy's optimization appears to be tractable in this framework.

© 2002 Published by Elsevier Science B.V.

1. Introduction

An essential task of Machine Learning (ML) and Data Mining (DM) systems is related to classification. This basically consists in giving the most accurate answer

* Tel.: +33-596-72-73-64; fax: +33-596-72-73-62.

E-mail address: mock@martinique.univ-ag.fr (R. Nock).

URL: <http://www.martinique.univ-ag.fr/~mock>

1 to the affectation of some observations (or patterns) to a finite number of classes.
2 Historically, measuring the quality of the system's answer has mostly been a matter of
3 computing its accuracy, i.e., the frequency (or probability) of correct predictions made
4 [3,17]. The advent of new technological media authorizing the storing of databases of
5 huge sizes, together with the increasing diversity of the problems (and goals) addressed,
6 have favored the emergence of a research trend in ML and DM. This trend discusses,
7 mostly experimentally, this standardized use of the accuracy to assess the merit of a
8 system (see e.g. [17]), as well as to compare different classifier learning algorithms to
9 decide which one should be used preferentially (see e.g. [22]). Following this trend, we
10 question again the replacement of the accuracy itself by other performance measures.
11 However, apart from our purely theoretical standpoint on this question, we deem our
12 approach original and distinguished from the others in that it is, to our knowledge, the
13 first to evaluate the computational burden of the accuracy's replacement/completion.

14 The primary inadequacy of the accuracy stems from a tacit assumption that the
15 overall accuracy controls by-class accuracies, or similarly that class distributions among
16 examples are constant and relatively balanced, see for example [20]. This is obviously
17 not true: skewed distributions are frequent in agronomy, or more generally in life or
18 earth sciences. For example, no more than 6% of the human DNA represents coding
19 genes [22]. Another example is the oil spill detection problem of [17], in which roughly
20 4% of the data represent oil slicks, the remaining being lookalikes. Even more extreme
21 cases exist, in information retrieval, in which the minority class can scarcely represent
22 0.2% of the data [17]. In all these cases, the interesting, unusual class is often the rare
23 one, and the well-balanced hypothesis may simply lead to the elusion of its elements
24 when building a classifier. In [17], a simple classifier labeling all patterns as lookalikes
25 (this is the so-called majority rule) would achieve an accuracy of 96%. As pointed out
26 by Kubat et al. [17], this looks like a high accuracy, but the classifier is totally useless
27 since it completely fails to achieve the goal of oil spill detection. On the other hand,
28 a system achieving only 94% detection on oil spills, and 94% detection on lookalikes,
29 would have a worse accuracy, and yet would be deemed highly successful [17].

30 This last example shows two important and typical phenomena in real-world
31 problems. First, the balanced distributions assumption is actually false. Second, the
32 misclassification of some examples may be of heavy consequences, a cost which is not
33 integrated in the accuracy. Fraud detection is another good example of such a cost-
34 sensitive situation [22], but there are many others. In database marketing, a prominent
35 application consists in targeting the people likely to respond to a mailing. In that case,
36 the cost of mailing to a non-respondent is small, but the cost of not mailing to someone
37 who would respond is the entire profit lost [3]. Solving the cost problem by the inte-
38 gration of the costs in the accuracy, to shift its behavior towards the crucial examples,
39 is also far from being obvious, as it involves "multiple considerations whose units are
40 incommensurable" [17].

41 Furthermore, the accuracy may be inadequate in some cases because other parameters
42 are to be taken into account. Some works [16] report the need to add an information
43 measure to the accuracy, to eliminate the influence of prior probabilities. Constraints
44 on size parameters (see [18,19]) are sometimes to be used because we want to obtain
45 small formulas, to ease their interpretability by the system end-user.

1 Finally, some works also report the experimental convenience that reducing the size
2 of the data itself can have when simply optimizing the accuracy [24]. Indeed, it is
3 well known in ML and DM that removing some parts of the data, such as features
4 (or variables), is a good experimental solution to reduce the size of the models built
5 afterwards, while avoiding to damage their accuracy too much. Sometimes, it can
6 even provide a way to improve their accuracy on hard problems. Whereas experiments
7 show that feature reduction can be a good criterion to optimize in conjunction with
8 the accuracy, one may wonder how these two constraints computationally interact.

9 To examine the possible influence of all these completion/replacement criteria, we
10 have chosen as our framework a field particularly sensitive to the computational com-
11 plexity factor, Inductive Logic Programming (ILP). ILP is a rapidly growing research
12 field, concerned with the use of variously restricted subclasses of Horn clauses to
13 build ML algorithms. According to a census of [26], in 1998, almost 70 applications
14 were using ILP formalism, 20 of which were science applications, which can be par-
15 titioned into biological (four) and drug design (16) applications. With the increasing
16 popularity of ILP, in particular to address complex domains, this number has certainly
17 increased since then. ILP–ML algorithms have been applied with some success in
18 areas of biochemistry and molecular biology [26]. Using ILP formalism, we argue that
19 the replacement of the accuracy raises computational complexity issues. This is all the
20 more important in ILP studies, as ILP is a field which can be naturally concerned with
21 intractability or even undecidability issues [13], and keeping tractability is of primary
22 importance to keep the ILP formalism power affordable to practical learning systems.
23 For this reason, ILP is certainly a domain of choice for a computational study of
24 the accuracy’s replacement/completion. More precisely, here is the structuring of the
25 argument.

26 First, we explain that the single accuracy requirement can be completed by an addi-
27 tional requirement to provide more adequate criteria. We integrate various constraints
28 over two important kind of parameters: by-class error functions, and representation
29 parameters such as feature selection ratios, size constraints. These criteria are inspired
30 by the works of [17–19,24]. We do not integrate in our criteria the information measure
31 of [16], as it is mainly designed to handle classifiers with probabilistic answers, and
32 is therefore not suited to ordinary Horn clauses.

33 Then, we study the replacement of the accuracy criterion using a general method
34 derived from statistical decision theory, based on a specific bi-criteria optimization (see
35 e.g. [20–22]).

36 Finally, we investigate the replacement of the error by a single replacement criterion.
37 Two candidates we study are criteria proposed in [20], and used in [17].

38 In this paper, we show that any of such integration leads to a very negative structural
39 complexity result, which is not faced by the accuracy optimization alone. The result
40 has a side effect which can be presented as a “loss” in the formalism’s expressiveness,
41 a seldom property in classical ML complexity issues. Indeed, it authorizes the construc-
42 tion of Horn clauses sets of unbounded size (even exponential), but, which we prove,
43 having no more expressive power than a single Horn clause. We prove a threshold
44 in intractability since it appears immediately with the additional requirement, and is
45 not a function of its tightness. Furthermore, the effects of the constraints on optimal

1 accuracies vanish as the number of predicates increases, since optimal accuracies with
2 or without the additional constraints are asymptotically equal. This phenomenon tends
3 to strengthen the threshold effect in intractability. Finally, for some criteria, their blend-
4 ing with the accuracy brings the most negative result: not only does the intractability
5 appear immediately with the criterion, but also the error cannot be dropped down
6 under that of the unbiased coin.

7 The reductions are presented for a subclass of Horn formalism simple enough to be
8 an element of the intersection of all classically encountered in theoretical ILP studies or
9 practical ILP learning systems. As a consequence, our results also hold in all these other
10 settings. The following section details the bases of learnability and ILP. It is followed
11 by a section introducing the possible criteria to address the accuracy's drawbacks, and
12 the tools used in our proofs. Then, all the results are presented in the last section,
13 along with some possible extensions to other formalisms, or to learning models. For
14 the reader's convenience, the proofs and technical aspects, not necessary to understand
15 the results, have been bulked in the two appendices. The first presents a synthetic view
16 of all proofs, the second presents in-depth reductions. In order not to laden the paper
17 with a collection of extensive proofs, some cases have been voluntarily omitted.

2. Learnability and ILP

19 Denote as \mathcal{C} and \mathcal{H} two classes of concept representations, respectively called *target*
20 class and *hypothesis* concept class. Informally, our objective is to build a concept from
21 the hypothesis class, approximating as best as possible an unknown concept c , called
22 the target concept, element of \mathcal{C} . In real-world domains, we do not know the target
23 concept's class, that is why we have to make ad hoc choices for \mathcal{H} with a powerful
24 enough formalism, yet ensuring tractability. Even if some benchmark problems appear
25 to be easily solvable [10], ML applications, and particularly ILP, face more difficult
26 problems [26], for which the choice of \mathcal{H} is crucial.

27 After the choice of \mathcal{H} , approximating the target concept can only be achieved by
28 catching a glimpse of it, through its extensional representation, i.e. by drawing exam-
29 ples, classified according to c . Generally, the data collected can only account for a
30 small part of this very large set, and the objective is then to build the intensional rep-
31 resentation of some hypothesis, whose extensional representation shall hopefully match
32 as best as possible the target concept's. Most of the studies dealing with the accuracy's
33 replacement, as well as computational complexity results in ML, have been investigated
34 with two classes [22,21]. We also consider a two-class setting. It is not really important
35 for us, as results already become hard in that setting.

36 We shall see later in this section how examples and concepts are described in the
37 context of ILP. Before, it is important to clarify the way we "collect" the examples,
38 and then use it to obtain either positive or negative results in ML or DM. Theoretically
39 speaking, a large part of the modern approaches to obtain positive results for ML/DM
40 algorithms draws its roots in two fundamental bodies, the so-called Probably Approx-
41 imately Correct (PAC) learning model of [27], and the Statistical Learning Theory,
42 fathered by Vapnik [29]. The principle is that the examples are drawn from some

1 unknown, but fixed distribution D , and labeled according to an unknown $c \in \mathcal{C}$. If we
 2 suppose that the representation space is discretized, then we can denote the accuracy
 3 of some $h \in \mathcal{H}$ with respect to (w.r.t.) c by $P_D(h=c) = \sum_{h(x)=c(x)} D(x)$ (here x is an
 4 observation and $h(x)$, $c(x)$ are, respectively, the classes given to x by h and c). Note
 5 that this quantity is measured over the whole set of possible examples, a domain to
 6 which we do not have access, as specified before. We have only access to its estimator
 7 over the sample collected. This raises statistical issues to evaluate the quality of this
 8 estimator (and h), issues discussed in many papers (see e.g. [29]). The objective of our
 9 paper is not to discuss the statistical burden of the theory, but its computational issues.

10 Proving most negative complexity-theoretic results for ML/DM follows a quite stan-
 11 dardized approach. It consists in building a particular set of examples, supposed to be
 12 the set collected, and giving a frequency distribution over these examples “mimick-
 13 ing” D (a seminal paper to this approach is [12]), and then proving the result on the
 14 basis of this particular instance of the problem. Our results also rely on this scheme.
 15 It is important to note that the negative results are, in that case, complexity-theoretic,
 16 i.e. they raise the hardness of finding efficient (e.g. polynomial) algorithms to address
 17 the problem. As briefly exposed before, they do not address the statistical hardness of
 18 building h , since the instance built boils down to having access to the whole domain
 19 knowledge (all examples that are not present in the set are supposed to have zero prob-
 20 ability of occurrence). An interesting fact in negative computational results in ML/DM
 21 is that they may have two consequences. The first is what motivates this paper, i.e. the
 22 inexistence of affordable practical algorithms to solve these problems. The second is
 23 the extension of these results to negative results for learning in models derived from
 24 the PAC model of [27]. Some of our results can be extended to negative results on
 25 the PAC-derived robust learning model of [9,11]. This is described later.

26 We now introduce our formalism for the examples and the hypotheses, ILP. The
 27 field of ILP is concerned with the induction of first-order Horn clauses from examples
 28 and background knowledge. A Horn clause has the following form:

$$29 \quad q(\dots) \leftarrow a_1(\dots) \wedge a_2(\dots) \wedge \dots \wedge a_n(\dots)$$

30 $q(\dots)$ is called the head of the clause and the conjunction $a_1(\dots) \wedge a_2(\dots) \wedge \dots \wedge a_n(\dots)$
 31 is called the body of the clause. A Horn clause with no body is unit. A clause with no
 32 variable is ground. Given a Horn clause language \mathcal{L} and a correct inference relation
 33 on \mathcal{L} , the problem can be formalized in a general way as follows [11]:

Definition 1. Given:

- 34 • A background knowledge \mathcal{BK} expressed in a language $\mathcal{LB} \subseteq \mathcal{L}$,
- 35 • A set of examples \mathcal{S} in a language $\mathcal{LS} \subseteq \mathcal{L}$, consisting of positive examples,
 36 \mathcal{S}^+ , and negative examples, \mathcal{S}^- , such that $\mathcal{B} \not\models \mathcal{S}^+$ (\mathcal{B} does not entail the positive
 37 examples) and $\mathcal{B}, \mathcal{S}^- \not\models \square$ (\mathcal{S} is consistent with \mathcal{B}).
- 38 • A hypothesis class \mathcal{H} described over a language $\mathcal{LH} \subseteq \mathcal{L}$,
 39 find a hypothesis $h \in \mathcal{H}$ such that

$$40 \quad \mathcal{B} \wedge h \models \mathcal{S}^+, \tag{1}$$

$$41 \quad \mathcal{B} \wedge h \not\models \mathcal{S}^- \tag{2}$$

1 i.e., \mathcal{B} and h explain the positive examples whereas they do not explain the negative
 2 examples.

3 We now give some precisions on this definition. The background knowledge in ILP
 4 is usually restricted in order to avoid undecidability problems about the deduction
 5 process [13,4]. A usual restriction makes use of ground background knowledge, i.e.,
 6 consisting of ground unit clauses. A clause is ground if it does not contain any variables.
 7 Therefore, to ensure tractability, we suppose that the background knowledge consists
 8 of ground unit clauses, and examples are ground unit clauses too. Another restriction
 9 commonly encountered consists in preventing the use of function symbols of arity > 0 :

11 **Definition 2.** A clause is called function-free iff all its arguments are either variables
 12 or constants (function symbols of arity 0).

13 As in [14], we use θ -subsumption as the inference relation. θ -subsumption is a
 14 correct and complete inference procedure between function-free Horn clauses ($h \models h'$
 15 iff $h\theta \subseteq h'$). This however leads to a modification of the learning problem, as stated in
 16 the following lemma:

17 **Lemma 3** (Kietz, [13]). *The learning problem is equivalent to learning the same pro-
 18 gram with θ -subsumption, and empty background knowledge and examples defined
 19 as ground Horn clauses of the form $e \leftarrow b$, where $e \in \mathcal{S}$ and $b \in \mathcal{BK}$.*

20 This lemma allows us to incorporate the background knowledge in the new exam-
 21 ples (and is thus empty). Our results make use of a simple subclass of Horn clause
 22 formalism. Its property is that it is an element of common subclasses of Horn clause
 23 formalisms usually encountered in practice or in theoretical learning studies. Therefore,
 24 since our results are essentially negative, they hold also for all these other subclasses.
 25 The most important property of our subclass is that the predicates arity is one. There-
 26 fore,

- 27 • we can suppose that the Horn clauses contain the same variable, say X . In other
 28 words, the clauses are constrained.
- 29 • the clauses are 01-determinate as defined in [13,4]. In other words, the maximum
 30 predicate arity is 1 and the depth of each term is that of the head, 0. This represents
 31 the easiest case of determinacy.
- 32 • \mathcal{BK} does not contain the predicate to be inferred, and the Horn clauses are non-
 33 recursive.

34 The principle of our negative results, from an ILP point of view, is quite simple:
 35 we create a formalism so simple such that, given the constraints, there cannot always
 36 exist a set of Horn clauses solution of the learning problem. From that, the goal of the
 37 learning problem is relaxed to that of an approximation problem well known in robust
 38 learning [11]: find a hypothesis $h \in \mathcal{H}$ such that

$$39 \quad \mathcal{B} \wedge h \models \mathcal{S}^+, \tag{3}$$

$$40 \quad \mathcal{B} \wedge h \not\models \mathcal{S}^- \tag{4}$$

41 for the largest part of the examples in \mathcal{S} .

1 3. Replacement criteria and the hardness technique

3.1. Extending the accuracy

3 For any fixed positive rational v , we use the following adequate notion of distance
 [1] between two reals u, v : $d_v(u, v) = |u - v| / (u + v + v)$. We also use eight rates on the
 5 examples (definitions differ slightly from [22]):

$$TP = \sum_{h(x)=1=c(x)} D(x); \quad TPR = \frac{TP}{\sum_{c(x)=1} D(x)},$$

7 $FP = \sum_{h(x)=1 \neq c(x)} D(x); \quad FPR = \frac{FP}{\sum_{c(x)=0} D(x)},$

$$TN = \sum_{h(x)=0=c(x)} D(x); \quad TNR = \frac{TN}{\sum_{c(x)=0} D(x)},$$

9 $FN = \sum_{h(x)=0 \neq c(x)} D(x); \quad FNR = \frac{FN}{\sum_{c(x)=1} D(x)}.$

In many DM/ML domains, the user's desiderata are often the optimization of more than
 11 one basic criterion (accuracy, precision, recall, sizes, etc.). Various composed criteria
 exist, combining some of these, but it is hard to obtain a suitable combination into one
 13 criterion, so as to optimize in one step more than one of these basic demands. The
 accuracy is typical, but others are well known, such as the geometrical mean of Kubat
 15 et al. [17], which ignores precision. Some authors, such as [16], have proposed to take
 into account more than one criterion, such as information measures for probabilistic
 17 classifiers. In order to complete the accuracy requirements, we imagine seven types of
 additional constraints aiming at controlling the well-balanced drawback of the accuracy
 19 alone [20], or precision or recall measures [20,17], or size parameters [18,24]. Each of
 them is parameterized by a number ζ (between 0 and 1), and defines a subset of \mathcal{H} ,
 21 which shall be parameterized by D if the distribution controls the subset through the
 constraint. The first three subsets of \mathcal{H} contain hypotheses for which the FP and FN
 23 are not far from each other, or a one-side error is upper bounded:

$$\mathcal{H}_{D,1}(\zeta) = \{h \in H \mid d_v(FP, FN) \leq \zeta\}, \quad (5)$$

25 $\mathcal{H}_{D,2}(\zeta) = \{h \in H \mid FN \leq \zeta\}, \quad (6)$

$$\mathcal{H}_{D,3}(\zeta) = \left\{ h \in H \mid FN \leq \frac{1}{\zeta} FP \right\}. \quad (7)$$

27 The two following subsets are parameterized by constraints equivalent to some fre-
 quently encountered in the information retrieval community [25], respectively (1 minus)
 29 the precision and (1 minus) the recall criteria.

$$\mathcal{H}_{D,4}(\zeta) = \left\{ h \in H \mid \frac{FP}{TP + FP} \leq \zeta \right\}, \quad (8)$$

$$1 \quad \mathcal{H}_{D,5}(\zeta) = \left\{ h \in H \mid \frac{FN}{TP + FN} \leq \zeta \right\}. \quad (9)$$

Now, we give two more constraints specific to Horn clauses. Horn clauses shall be extensively defined in a section devoted to ILP formalism. We give some preliminary and necessary definitions for the two constraints we define. A Horn clause (a definite program clause) [4] has the following form:

$$q(\dots) \leftarrow a_1(\dots) \wedge a_2(\dots) \wedge \dots \wedge a_n(\dots).$$

7 Here, q, a_1, a_2, \dots, a_n are predicate symbols. Define $\#Predicates(h)$ as the total number of different predicates of h , $\#Whole_predicates(h)$ as the overall number of predicates of h (if one predicate is present k times, it is counted k times), and $\#Total_predicates$ as the total number of different available predicates to build a Horn clause for our specific problem. The two last subsets of \mathcal{H} are parameterized by formulas, respectively, having a sufficiently small fraction of the available predicates, or having a sufficiently small overall size:

$$\mathcal{H}_6(\zeta) = \left\{ h \in H \mid \frac{\#Predicates(h)}{\#Total_predicates} \leq \zeta \right\}, \quad (10)$$

$$15 \quad \mathcal{H}_7(\zeta) = \left\{ h \in H \mid \frac{\#Whole_predicates(h)}{\#Total_predicates} \leq \zeta \right\}. \quad (11)$$

The division by the total number of different predicates in $\mathcal{H}_7(\zeta)$ is made only for technical reasons: to obtain hardness results for small values of ζ . The first problem we address can be summarized as follows:

19 **Problem 1.** *Given ζ and $i \in \{1, 2, \dots, 7\}$, can we find an algorithm returning a set of Horn clauses from $\mathcal{H}_{(D),i}(\zeta)$ whose error is no more than a given γ , if such a hypothesis exists?*

3.2. Replacing the accuracy: the ROC analysis

23 Receiver Operating Characteristic (ROC) analysis is a traditional methodology from signal detection theory [5]. It has been used in machine learning recently [20–22] in order to correct the main drawbacks of the accuracy. In ROC space (this is the coordinate system), we visualize the performance of a classifier by plotting TPR on the Y -axis, and FPR on the X -axis. Fig. 1 presents the ROC analysis, along with three possible outputs which we present and analyze now. If a classifier produces a continuous output (such as an estimate of posterior probability of an instance's class membership [22], or a real-valued confidence such as in AdaBoost [23], for any possible value of FPR , we can get a value for TPR , by thresholding the output between its extreme bounds. If a classifier produces a discrete output (such as Horn clauses), then the classifier gives rise to a single point. If the classifier is the random choice of the class, either (if it is continuous) the curve is the line $y=x$, or (if it is discrete) there is a single dot, on the line $y=x$. One important thing to note is that the ROC representation gives the

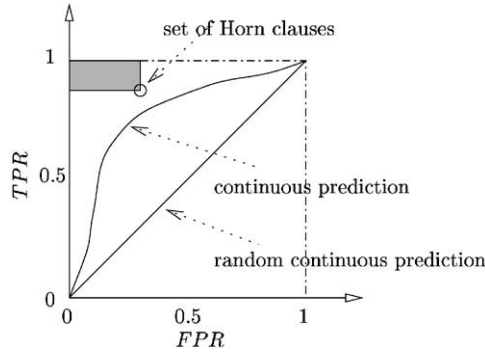


Fig. 1. The ROC analysis of a learning algorithm, with three general classifications: discrete, continuous, and default.

- 1 behavior of an algorithm without regarding the class distribution or the error cost [20].
 2 And it allows to choose the best of some classifiers, by the following procedure. Fix
 3 as K^+ the cost of misclassifying a positive example, and K^- the cost of misclassifying
 4 a negative example (these two costs depend on the problem). Then the *expected* cost
 5 of some classifier represented by point (FPR, TPR) is given by the following formula:

$$\sum_{c(x)=1} D(x)(1 - TPR)K^+ + \sum_{c(x)=0} D(x)FPRK^-. \quad (12)$$

- 7 Two algorithms, whose corresponding point are, respectively (FPR_1, TPR_1) and $(FPR_2,$
 $TPR_2)$, have the same expected cost iff

$$\frac{TPR_2 - TPR_1}{FPR_2 - FPR_1} = \frac{\sum_{c(x)=1} D(x)K^+}{\sum_{c(x)=0} D(x)K^-}. \quad (13)$$

- 11 This gives the slope of an *isoperformance* line, which only depends on the rela-
 12 tive weights of the examples, and the respective misclassification costs. Given one
 13 point on the ROC, the classifiers performing better are those on the “northwest” of
 14 the isoperformance line with the preceding slope, and to which the point belongs. If
 15 we want to find an algorithm A performing *surely* better than an algorithm B , we therefore
 16 should strive to find A such that its point lies into the rectangle whose opposite
 17 vertices are the $(0,1)$ point (the perfect classification) and B 's point (a grey rectangle is shown
 18 on the top left of Fig. 1). From that, the second problem we address is the following:

- 19 **Problem 2.** *Given one point (TPR_x, FPR_x) on the ROC, can we find an algorithm*
 20 *returning a set of Horn clauses whose point falls into the rectangle with opposite*
 21 *vertices $(0, 1)$ and (TPR_x, FPR_x) , if such a hypothesis exists?*

- 23 Note that the problem we address is based on weak constraints: indeed, we only
 require the algorithm to work on a *single* point (TPR_x, FPR_x) .

1 3.3. Replacing the accuracy by a single criterion

3 The ROC analysis is based on two criteria, controlling *FPR* and *TPR*. The question
 4 of whether the accuracy can be replaced by a single criterion instead of two has
 5 been raised in [20]. Some researchers [20] propose the use of the following criterion:
 6 $(1 - FPR) \times TPR$. A geometric interpretation of the criterion is the following [20]:
 7 it corresponds to the area of a rectangle whose opposite vertices are (FPR, TPR) and
 8 $(1, 0)$. The typical isoperformance curve is now a hyperbola. The third problem we
 9 address is therefore:

9 **Problem 3.** *Given γ , can we find an algorithm returning a set of Horn clauses such
 10 that $(1 - FPR) \times TPR \geq \gamma$, if such a hypothesis exists?*

11 In [17], a criterion is maximized which is the square-root of our criterion. Because
 12 of the monotonicity properties of this function, our negative results on problem 3 shall
 13 also hold for the criterion of [17].

14 3.4. Basic tools for the hardness results

15 Concerning problem 1, fix $a \in \{1, 2, 3, 4, 5, 6, 7\}$. We want to approximate the best
 16 concept in $\mathcal{H}_{(D),a}(\zeta)$ by one still in $\mathcal{H}_{(D),a}(\zeta)$. However, the best concept in $\mathcal{H}_{(D),a}(\zeta)$
 17 generally does not have an error equal to the optimal one over \mathcal{H} given D , $opt_{\mathcal{H}_D}(c)$.
 18 In fact, it has an error that we can denote

$$\begin{aligned} opt_{\mathcal{H}_{(D),a}(\zeta)}(c) &= \min_{h' \in \mathcal{H}_{(D),a}(\zeta)} \sum_{h(x) \neq c(x)} D(x). \\ &\geq opt_{\mathcal{H}_D}(c) \end{aligned}$$

19 The goodness of the accuracy of a concept taken from $\mathcal{H}_{(D),a}(\zeta)$ should be appreciated
 20 with respect to this “constrained” optimum. Our results on problem 1 are all obtained
 21 by showing the hardness of solving the following decision problem:

Definition 4 (*Approx-Constrained*($\mathcal{H}, (a, \zeta)$)). • *Name:* Approx-Constrained($\mathcal{H}, (a, \zeta)$).

- 23 • *Instance:* A set of negative examples \mathcal{S}^- , a set of positive examples \mathcal{S}^+ , a rational
 24 weight $0 < w(x_i) = n_i/d_i < 1$ for each example x_i , a rational $0 \leq \gamma < 1$. We assume that
 25 $\sum_{x \in \mathcal{S}^+ \cup \mathcal{S}^-} w(x_i) = 1$.
 26 • *Question:* Does there exist a hypothesis $h \in \mathcal{H}_{(D),a}(\zeta)$ that satisfies $\sum_{h(x) \neq c(x)} w(x)$
 27 $\leq \gamma$?

28 Define as n_e the size of the largest example we dispose of. Note that when the
 29 constraint is too tight, it can be the case that no hypothesis can actually satisfy it, and
 30 therefore

$$31 \quad \mathcal{H}_{(D),a}(\zeta) = \emptyset. \quad (14)$$

1 Define as $|h|$ the size of some $h \in \mathcal{H}$ (in our case, it is the number of Horn clauses
 of h). In the non-empty subset of \mathcal{H} where formulas are the most constrained (i.e.
 3 strengthening further the constraint gives an empty subset), define $n_{\mathcal{H}_{(D,a)}^*}^*$ as the size
 of the *smallest* hypothesis in $\mathcal{H}_{(D,a)}(\zeta)$ (therefore, it is the smallest hypothesis which
 5 satisfies the constraint). Then, our reductions all satisfy

$$n_{\mathcal{H}_{(D,a)}^*}^* \leq (n_e)^3. \quad (15)$$

7 Note that the constraint makes generally

$$opt_{\mathcal{H}_{(D,a)}(\zeta)}(c) > opt_{\mathcal{H}_D}(c), \quad (16)$$

9 which might seem to be a negative effect of the constraints. However, the reductions
 all satisfy

$$d_v(opt_{\mathcal{H}_D}(c), opt_{\mathcal{H}_{(D,a)}(\zeta)}(c)) = o(1) \quad (17)$$

i.e. asymptotic optimal accuracies coincide with or without the constraints; here, the
 13 limit is taken as the number of distinct predicates of the problem grows towards
 infinity ($\#Total_predicates \rightarrow \infty$). In addition, a principal corollary to all our
 15 results is that we can suppose that the whole time used to *write* the total set of
 Horn clauses is assimilated to $\mathcal{O}(n_e)$, for *any* set. By writing time, we mean time
 17 of a y procedure consisting only in writing down clauses. Examples of such a pro-
 cedure are “write down all clauses having k literals”, or even “write down *all* Horn
 19 clauses”. Such procedures can be viewed as *for-to*, or *repeat* algorithms. This prop-
 erty authorizes the construction of Horn clause sets having arbitrary sizes, even
 21 exponential.

Problem 2 is addressed by studying the complexity of the following decision problem.

23

Definition 5 (*Approx-Constrained-ROC*($\mathcal{H}, \gamma_{FPR}, \gamma_{TPR}$)). • *Name*: Approx-Constrained-
 25 ROC($\mathcal{H}, \gamma_{FPR}, \gamma_{TPR}$).

• *Instance*: A set of negative examples \mathcal{S}^- , a set of positive examples \mathcal{S}^+ , a rational
 27 weight $0 < w(x_i) = n_i/d_i < 1$ for each example x_i . We assume that $\sum_{x \in \mathcal{S}^+ \cup \mathcal{S}^-} w(x_i) = 1$.

• *Question*: Does there exist a hypothesis $h \in \mathcal{H}$ satisfying $1 - FPR \geq 1 - \gamma_{FPR}$ and
 29 $TPR \geq \gamma_{TPR}$?

31 Concerning problem 3, the reductions study a single replacement criterion Γ , and
 the following decision problem.

33 **Definition 6** (*Approx-Constrained-Single*($\mathcal{H}, \Gamma, \gamma$)). • *Name*: Approx-Constrained-
 Single($\mathcal{H}, \Gamma, \gamma$).

35 • *Instance*: A set of negative examples \mathcal{S}^- , a set of positive examples \mathcal{S}^+ , a ra-
 tional weight $0 < w(x_i) = \frac{n_i}{d_i} < 1$ for each example x_i . We assume that $\sum_{x \in \mathcal{S}^+ \cup \mathcal{S}^-}$
 37 $w(x_i) = 1$.

• *Question*: Does there exist a hypothesis $h \in \mathcal{H}$ satisfying $\Gamma(h) \leq \gamma$?

4. Results

For the sake of simplicity in stating our results, we abbreviate “Function free Horn Clauses” by the acronym “FfHC”.

4.1. Extending the accuracy

Theorem 5. *We have:*

- (i) [1] $\forall 0 < \zeta < 1$, *Approx-Constrained(FfHC, (1, ζ)) is Hard, when $v < (1 - \zeta)/\zeta$*
- (ii) [2] $\forall 0 < \zeta < \frac{1}{2}$, *Approx-Constrained(FfHC, (2, ζ)) is Hard.*
- (iii) [3] $\forall 0 < \zeta < 1$, *Approx-Constrained(FfHC, (3, ζ)) is Hard.*
- (iv) [4] $\forall 0 < \zeta < 1$, *Approx-Constrained(FfHC, (4, ζ)) is Hard.*
- (v) [5] $\forall 0 < \zeta < 1$, *Approx-Constrained(FfHC, (5, ζ)) is Hard.*
- (vi) [6] $\forall 0 < \zeta < 1$, *Approx-Constrained(FfHC, (6, ζ)) is Hard.*
- (vii) [7] $\forall 0 < \zeta < 1$, *Approx-Constrained(FfHC, (7, ζ)) is Hard.*

At that point, the notion of “hardness” needs to be clarified. By “Hard” we mean “cannot be solved in polynomial time under some particular complexity assumption”. The hypothesis we use is the same as [8] ($NP \not\subseteq ZPP$), which involves randomized complexity classes. We refer the reader to the paper of [8] for further details, not needed here.

Due to the fact that all proofs are essentially based on the same properties, only proof of point [1] is presented in details in Appendix B; the other results presented strictly use the same type of reduction, and are eventually sketched [6,7]. Also, in Appendix A, we give the proof that all distributions under which our negative results are proven lead to trivial positive results for the same problem when we remove the additional constraint, and optimize the accuracy alone.

Beyond the range of constraints that our negative results cover, note that any other additional constraint aside from the accuracy is a natural candidate to test the existence of negative results, unless pathological situations are created, such as when the constraint is so tight and removes so many hypotheses that the set of constrained hypotheses has small size (e.g. polynomial), and can be explored in polynomial time. Therefore, another incidence of our results is that in between the two extreme situations (no/over constrained requirement), optimizing the accuracy under constraint is a strictly more difficult problem, with non-trivial additional drawbacks. Furthermore, the upperbound error value (γ in Definition 4) in constraints 4–6 can be fixed arbitrarily in $]0, 1/2[$, which shows that almost removing the accuracy’s constraint does not make the problem easier: requiring the Horn clauses to perform slightly better than the unbiased coin leads also to intractability.

4.2. Replacing the accuracy: the ROC analysis

In this section, we show that the classical ROC components as described by Provost et al. [22] and Provost and Fawcett [21] lead to the same results as those we claimed for the preceding bi-criteria optimizations. The problem is all the more difficult as the

1 difficulty appears as soon as we choose to use ROC analysis, and is not a function of
the ROC bounds.

3 **Theorem 6.** $\forall 0 < \gamma_{FPR}, \gamma_{TPR} < 1$, *Approx-Constrained-ROC(FfHC, $\gamma_{FPR}, \gamma_{TPR}$) is hard.*

The distribution under which the negative result is proven is an easy distribution for
5 the accuracy's optimization alone, similarly to those of Section 4.1.

4.3. Replacing the accuracy by a single criterion

7 The negative result stated in the following theorem is to be read with all additional
drawbacks mentioned for the seven constraints. Again, the distribution under which the
9 theorem is proven is an easy distribution when optimizing the accuracy alone.

11 **Theorem 7.** $\exists \gamma_{\max} > 0$ such that $\forall 0 < \gamma < \gamma_{\max}$, *Approx-Constrained-Single(FfHC, $(1 - FPR) \times TPR, \gamma$) is Hard.*

(Proof included in Appendix B). As far as we know, $\gamma_{\max} \geq \frac{175}{41,616}$ (roughly $4.2 \times$
13 10^{-3}), but we think that this bound can be much improved. The accuracy can some-
times be conveniently replaced by the F_β statistics [2], which is an accurate composition
15 of precision and recall (see Section 3.1 for their definition), useful for text categoriza-
tion problems [2]. So far, we have not been able to conclude to the hardness of using
17 this criterion in our framework.

4.4. Beyond computational complexity and ILP

19 It is well known since [12] that negative results on such problems can sometimes be
extended to negative results for PAC-type learning models [27]. Such a model typically
21 brings a statistical and a computational constraint for an algorithm to be qualified as a
learning algorithm. Consider for example Definition 4, and the following learning model
23 arising in exactly the same setting, but in which we replace the set of examples by a so-
called oracle [12], drawing examples on demand, following a probability distribution D
25 unknown, but fixed. Suppose that the requirement on the constraint defining $\mathcal{H}_{(D),a}(\zeta)$
remains exactly the same, but the one limiting the accuracy on the “learning” sample is
27 replaced by a condition which states that, with sufficiently high probability ($> 1 - \delta$),
the accuracy over the whole domain is lower than some threshold ($< \text{opt}_{\mathcal{H}_{(D),a}(\zeta)}(c) + \varepsilon$),
29 for some parameters $\varepsilon, \delta > 0$. If we require that the computational time be a polynomial
in $1/\varepsilon, 1/\delta$, as well as in n_e and the (smallest) size of the optimal constrained hypothesis,
31 then the learning model we obtain corresponds to the robust learning model of [9,11],
to which add the requirement that the outputs satisfy a constraint (among our seven
33 first constraints). In that case, following a standardized approach [9,11,12], it is easy
to show that a negative result regarding Definition 4 can be translated to a negative
35 constrained robust learning result.

37 Apart from the extension of the results to learning models, a natural question is
their extension to other formalisms, outside the ILP field. So far, as ILP is a complex

1 formalism, the results can be extended to simpler formalisms such as some subclasses
of Boolean formulas. One example is the subclass of DNF (disjunctive normal form
3 formulas [28]) containing all monotonous formulas (without negative literals). Note
that we do not put any restriction on the size of the formulas, a very seldom result in
5 the huge quantity of theoretical ML results on DNF. Indeed, DNF is one of the most
central classes to the PAC learning model of [27], studied early by Valiant himself
7 [28], and still raising some of the most important problems in computational learning
theory [15], in particular for its learnability or approximability properties. In that setting,
9 removing the monotonicity constraint in our results is certainly a problem which would
deserve further investigations.

11 5. Conclusion

In this paper, we have presented a new approach to the problem of the accuracy's
13 replacement in ML and DM, a problem recently addressed in a growing number of
papers. We have argued that the usual criticisms, against the use of the accuracy for
15 comparing the reliability of classifiers as well as for being optimized to build classi-
fiers, face complexity issues. The case against accuracy, as initially brought in [22], is
17 therefore more complicated than usually presented in ML or DM papers. This justifies
the title of our paper, which can be read in two ways, either presenting new aspects of
19 the difficulty (complexity) of the task to find new criteria to replace the accuracy, or
presenting (structural) complexity issues about the possible replacements/completions.
21 One important thing about our results is that the complexity results go beyond the
usual intractability results related to ML (or DM). In our case indeed, there are some
23 side effects, rather surprising, proving that the difficulty of the learning task, when
the accuracy is replaced, is accompanied by severe drawbacks on the formalism's
25 expressiveness. In deep contrast, the optimization of the accuracy alone in our setting
is trivial, since the optimal solutions can be found directly without any algorithmic
27 effort.

Recently, a new approach to building classifiers has been proposed, arguing against
29 the use of the accuracy as the optimization criterion for the induction of classifiers
[23]. This approach, called boosting, has been plebiscited as one of the best currently
31 available in classification [6]. However, it raises conjectures about the tractability of
the optimization of these new criteria [23] in some cases. This shall certainly be the
33 subject of future studies.

Appendix A. the global reduction

35 Reductions are achieved from the *NP*-Complete problem "Clique" [7]:

Definition 8 (Clique). • *Name*: Clique.

37 • *Instance*: A graph $G = (X, E)$, an integer k .

• *Question*: Does there exists a clique of size $\geq k$ in G ?

1 Of course, “Clique” is not hard to solve for any value of k . The following lemma
 2 establishes values of k for which we can suppose that the problem is hard to solve
 3 ($\binom{n}{k} = n!/((n-k)!k!)$ is the binomial coefficient):

Theorem 9. (i) *We can suppose that $\binom{k}{2} \leq |E|$, and k is not a constant, otherwise
 5 “Clique” is polynomial.* (ii) *For any $\alpha \in]0, 1[$, “Clique” is hard for the value $k = \alpha|X|$
 or $k = |X|^\alpha$.*

7 **Proof.** (i) is immediate; (ii) follows from [8]: it is proven that the largest clique
 8 size is not approximable to within $|X|^\beta$, for any constant $0 < \beta < 1$. Therefore, the
 9 graphs generated have a clique number which is either l , or greater than $l \times |X|^\beta$, with
 $l < |X|^{1-\beta}$. The idea is then to make k fall somewhere in between l and $l \times |X|^\beta$. For
 11 $k = |X|^\alpha$ ($\alpha \in]0, 1[$), this is immediate (if $\alpha > \frac{1}{2}$, we pick $1 > \beta > \alpha$ and if $\alpha \leq 1/2$, we
 pick $1 > \beta > 1 - \alpha$); for $k = \alpha|X|$, whenever the graph is large enough and satisfies
 13 $l \times |X|^\beta < \alpha|X|$, then we simply add u new vertices, each linked to all other vertices.
 Picking $u = (\alpha/(1 - \alpha))|X| - l(1 + |X|^\beta)/(2(1 - \alpha))$ is enough to make k fall in the
 15 interval $]l + u, l|X|^\beta + u[$. This ends the proof of the theorem. \square

The structure of the examples is the same for any of our reductions.

- 17 • Define a set of $|X|$ unary literals $a_1(\cdot), \dots, a_{|X|}(\cdot)$, in bijection with the vertices of G .
 To this set of literals, we add two unary literals, $s(\cdot)$ and $t(\cdot)$. The inferred predicate
 19 is denoted q . The choice of unary predicates is made only for a simplicity purpose.
 We could have replaced each of them by l -ary predicates without changing our
 21 proof.
- Define a set of constant symbols useful for the description of the examples:

$$23 \quad \{l_{i,j}, \forall (i,j) \in E\} \cup \{l_1, l_2, l_3, l_4\} \cup \{m_i, \forall i \in \{1, \dots, |X|\}\}.$$

Examples are described in the following way:

- 25 • Positive examples from \mathcal{S}^+ :

$$\forall (i,j) \in E, p_{i,j} = q(l_{i,j}) \leftarrow \bigwedge_{k \in \{1, \dots, |X|\} \setminus \{i,j\}} a_k(l_{i,j}) \wedge t(l_{i,j}), \quad (\text{A.1})$$

$$27 \quad p_1 = q(l_1) \leftarrow \bigwedge_{k \in \{1, \dots, |X|\}} a_k(l_1) \wedge t(l_1), \quad (\text{A.2})$$

$$p_2 = q(l_2) \leftarrow a_1(l_2). \quad (\text{A.3})$$

- 29 • Negative examples from \mathcal{S}^- :

$$\forall i \in \{1, \dots, |X|\}, n_i = q(m_i) \leftarrow \bigwedge_{k \in \{1, \dots, |X|\} \setminus \{i\}} a_k(m_i) \wedge t(m_i) \quad (\text{A.4})$$

$$31 \quad n'_1 = q(l_3) \leftarrow \bigwedge_{k \in \{1, \dots, |X|\}} a_k(l_3) \wedge s(l_3) \wedge t(l_3) \quad (\text{A.5})$$

$$n'_2 = q(l_4) \leftarrow \bigwedge_{k \in \{1, \dots, |X|\}} a_k(l_4) \wedge s(l_4). \quad (\text{A.6})$$

33 It comes that we always have $n_{\mathcal{H}(D, \omega(\zeta))} = \mathcal{O}(|X|^3)$ (this is the coding size of the
 positive examples) and $n_e = \mathcal{O}(|X|)$. Non-uniform weights are given to each example,
 35 depending on the constraint to be tackled with. The common-point to all reductions is

- 1 that the weights of all examples n_j (resp. all $p_{i,j}$) are equal (resp. to w^- and w^+). In each reduction, examples and clauses satisfy:
- 3 **H₁** p_2 is forced to be badly classified.
H₂ n'_1 is always badly classified.
5 **H₃** $w(n'_2)$ ensures that n'_2 is always given the right class, forcing any clause to contain literal $t(\cdot)$. \square
- 7 When we remove n'_2 , we also ensure that p_2 is removed too.

Lemma 10. *Any clause containing literal $s(\cdot)$ can be removed.*

- 9 **Proof.** Suppose that one clause contains $s(\cdot)$. Then it can be θ -subsumed by n'_1 and by no other example (even if n'_2 exists, because of **H₃**); but n'_1 θ -subsumes any clauses and also the empty clause. Therefore, removing the clause does not modify the value of any criteria based on the examples weights. Concerning the sixth (resp. seventh) constraint, the fraction of predicates used after removing the clause is at most the one before, thus, if the clause is an element of $\mathcal{H}_6(\zeta)$ (resp. $\mathcal{H}_7(\zeta)$) before, it is still an element after.

- As a consequence, p_1 is always given the positive class (even by the empty clause!).
17 We now give a general outline of the proof for Problem 1; reductions are similar for the other problems. Given $h = \{h_1, \dots, h_l\}$ a set of Horn clauses, we define the set

$$19 \quad \mathcal{I} = \{i \in \{1, \dots, |X|\} : \exists j \in \{1, \dots, l\}, a_i(\cdot) \notin h_j\}$$

- and we fix $|\mathcal{I}| = k'$. In our proofs, we define two functions taking rational values,
21 $E(k')$ and $F_a(k')$ ($k' \in \{1, \dots, |X|\}$, $a = 1, 2, 3, 4, 5, 6, 7$). They are chosen such that:

- $E(k')$ is strictly increasing, $\sum_{(x \in \mathcal{S}^+ \wedge h(x)=0) \vee (x \in \mathcal{S}^- \wedge h(x)=1)} w(x) \geq E(k')$ and $E(k) = \gamma$.
 - $F_a(k')$ is strictly decreasing, it is a lowerbound of the function inside $\mathcal{H}_{(D),a}(\zeta)$, and $F_a(k) = \zeta$ (excepted for $a = 3$, $F_3(k) = 1/\zeta$)
- 25 $\forall a \in \{1, 2, 3, 4, 5, 6, 7\}$, if there exists an unbounded set of Horn clauses $h \in \mathcal{H}_{(D),a}(\zeta)$ satisfying $\sum_{(x \in \mathcal{S}^+ \wedge h(x)=0) \vee (x \in \mathcal{S}^- \wedge h(x)=1)} w(x) \leq \gamma$, its error rate implies $k' \leq k$ and constraint implies $k' \geq k$. So $|\mathcal{I}| = k' = k$. The interest of the weights is then to force $\binom{k}{2}$ positive examples from the set $\{p_{i,j}\}_{(i,j) \in E}$ to be well classified, while we ensure the misclassification of at most k negative examples of the set $\{n_i\}_{i \in \{1, \dots, |X|\}}$. It comes that the $\binom{k}{2}$ correspond to the $\binom{k}{2}$ edges linking the $|\mathcal{I}| = k$ vertices corresponding to negative examples badly classified. We therefore dispose of a clique of size $\geq k$.

- 31 Conversely, $\forall a \in \{1, 2, 3, 4, 5, 6, 7\}$, given some clique of size k whose set of vertices is denoted \mathcal{I} , we show that the singleton

$$33 \quad h = q(X) \leftarrow \bigwedge_{i \in \{1, \dots, |X|\} \setminus \mathcal{I}} a_i(X) \wedge t(X)$$

- 35 is an element from $\mathcal{H}_{(D),a}(\zeta)$ satisfying $\sum_{(x \in \mathcal{S}^+ \wedge h(x)=0) \vee (x \in \mathcal{S}^- \wedge h(x)=1)} w(x) \leq \gamma$. In this case, $n_{\mathcal{H}_{(D),a}(\zeta)}$ drops down to $\mathcal{O}(n_e)$.

- 37 All distributions used in Theorems 5 and 7 are such that $w^+ < w^-/|X|$, at least for graphs exceeding a fixed constant size. Also, due to the negative examples of weights w^- , if we remove the additional constraints and optimize the accuracy alone,

1 we can suppose that the optimal Horn clause is a singleton: merging all clauses by
 2 keeping among predicates $a_j(\cdot)$ only those present in all clauses does not decrease
 3 the accuracy. Under such a distribution, the optimal Horn clause necessarily contains
 4 all predicates $a_j(\cdot)$, and the problem becomes trivial. The distribution in Theorem 6
 5 satisfies $w^+ = w^-$. This is also a simple distribution for the accuracy's optimization
 6 alone: indeed, the optimal Horn clause over predicates $a_j(\cdot)$ is such that it contains
 7 no predicates $a_j(\cdot)$ that does not appear at least in one positive example. If the graph
 8 instance of “Clique” is connex (and we can suppose so, otherwise the problem boils
 9 down to find the largest clique in one of the connected components), then the optimal
 Horn clause does not contain any of the $a_j(\cdot)$. \square

11 Appendix B. proofs of negative results

B.1. Proof of point [1], Theorem 5

13 We fix the following weights for positive examples:

$$w(p_2) = \frac{1}{2(1-\zeta)}(\zeta v + |X|^2 w^-(1+\zeta)),$$

$$15 \quad \forall (i, j) \in E, w(p_{ij}) = w^+ = \frac{w^-}{(|X| + k)^2},$$

$$\begin{aligned} w(p_1) = & \frac{1}{2} \left(1 - \frac{\zeta v}{1-\zeta} \right) \\ & - \frac{1}{2} \left(w^- \left[|X|^2 \left(\frac{1+\zeta}{1-\zeta} + |X| - k \right) \right] \right) \\ & - \frac{1}{2} \left(w^+ \left[\frac{1-\zeta}{1+\zeta} \left(|X| - \binom{k}{2} \right) + |X| \right] \right). \end{aligned}$$

We fix the following weights for negative examples:

$$17 \quad w(n'_2) = \frac{1}{2}$$

$$\forall j \in \{1, \dots, |X|\}, w(n_j) = w^- = \frac{1}{|X|^2 |E|^2},$$

$$19 \quad w(n'_1) = \frac{1}{2} \left(\frac{1-\zeta}{1+\zeta} \left(|E| - \binom{k}{2} \right) w^+ + (|X|^2 - k) w^- \right).$$

21 Fix $\gamma = w(p_2) + w(n'_1) + k w^- + (|E| - \binom{k}{2}) w^+ / 2$ (note that $w(n'_2)$ ensures that n'_2 is
 22 given the right class), and $k_{\max} = 1 + \max_{2 \leq k'' \leq |X|; |E| - \binom{k''}{2} \geq 0} k''$. From the choice of
 23 weights, $\text{lcm}(\bigcup_{x_i \in \mathcal{S}^+ \cup \mathcal{S}^-} d_i) = \mathcal{O}(|X|^8)$ (“lcm” is the least common multiple), which

1 is polynomial. Define the functions:

$$\forall k' \in \{0; 1\}, E(k') = |E|w^+ + k'w^- + w(p_2) + w(n_1),$$

3 $\forall 2 \leq k' \leq k_{\max}, E(k') = \left(|E| - \binom{k'}{2}\right) w^+ + k'w^- + w(p_2) + w(n_1),$

$$\forall k_{\max} < k' \leq |X|, E(k') = k'w^- + w(p_2) + w(n_1)$$

5 (from the choice of weights, $E(k) = \gamma$),

$$\forall k' \in \{0; 1\}, F_1(k') = \frac{||E|w^+ - k'w^- + w(p_2) - w(n_1)|}{v + |E|w^+ + k'w^- + w(p_2) + w(n_1)},$$

$$\forall 2 \leq k' \leq k_{\max}, F_1(k') = \frac{|(|E| - \binom{k'}{2})w^+ - k'w^- + w(p_2) - w(n_1)|}{v + |E|w^+ + k'w^- + w(p_2) + w(n_1)},$$

$$\forall k_{\max} < k' \leq |X|, F_1(k') = \frac{|-k'w^- + w(p_2) - w(n_1)|}{v + |E|w^+ + k'w^- + w(p_2) + w(n_1)}$$

(from the choice of weights, $F_1(k) = \zeta$).

7 The equation obtained when $k' < k_{\max}$ takes its maximum for integer values when $k' = (|X| + k)^2 + 0.5 \pm 0.5 > |X|$. Furthermore,

9 $\forall 1 \leq k_{\max} \leq |X|, \left(|E| - \binom{k_{\max} - 1}{2}\right) w^+ < w^-,$

11 which leads to $E(k_{\max} - 1) < E(k_{\max})$. In a more general way, $E(k')$ is strictly increasing
12 over natural integers. Now remark that the numerator of $F_1(k')$ is strictly decreasing,
13 and its denominator strictly increasing. Therefore, $F_1(k')$ is strictly decreasing. Further-
more

$$d_v \left(\sum_{h(x) \neq 1=c(x)} w(x); \sum_{h(x) \neq 0=c(x)} w(x) \right) \geq F_1(k').$$

15 If $\exists h \in \mathcal{H}_{\{w_i\}, 1}(\zeta)$ satisfying $\sum_{h(x) \neq c(x)} w(x) \leq \gamma$, the error rate implies $k' \leq k$ and the
16 constraint implies $k' \geq k$. Thus $|\mathcal{S}| = k' = k$. As pointed out in the preceding appendix,
17 this leads to the existence of a clique of size $\geq k$.

18 Reciprocally, the Horn clause constructed in Appendix A satisfies both relations
19 $h \in \mathcal{H}_{\{w_i\}, 1}(\zeta)$, and $\sum_{h(x) \neq c(x)} w(x) \leq \gamma$. Indeed, we have

$$\sum_{h(x) \neq 1=c(x)} w(x) = \left(|E| - \binom{k}{2}\right) w^+ + w(p_2)$$

21 but also

$$\sum_{h(x) \neq 0=c(x)} w(x) = kw^- + w(n_1).$$

1 Therefore,

$$d_v \left(\sum_{h(x) \neq 1=c(x)} w(x); \sum_{h(x) \neq 0=c(x)} w(x) \right) = F_1(k) = \zeta$$

3 and $h \in \mathcal{H}_{\{w_i\},1}(\zeta)$. We also have $\sum_{h(x) \neq c(x)} w(x) = E(k) = \gamma$.

The reduction is achieved. We end by a remark on $d_v(\text{opt}_{\mathcal{H}_{\{w_i\}}}(c), \text{opt}_{\mathcal{H}_{\{w_i\},1}(\zeta)}(c))$.

5 We have

$$\begin{aligned} |\text{opt}_{\mathcal{H}_{\{w_i\}}}(c) - \text{opt}_{\mathcal{H}_{\{w_i\},1}(\zeta)}(c)| &\leq 1 - w(p_1) - w(p_2) - w(n'_1) - w(n'_2) \\ &\leq |E|w^+ + |X|w^- \end{aligned}$$

and

$$\begin{aligned} \text{opt}_{\mathcal{H}_{\{w_i\}}}(c) + \text{opt}_{\mathcal{H}_{\{w_i\},1}(\zeta)}(c) + v &\geq 2(w(p_2) + w(n'_1)) + v \\ &\geq \frac{\zeta v}{(1 - \zeta)} + v. \end{aligned}$$

7 Therefore, we get

$$d_v(\text{opt}_{\mathcal{H}_{\{w_i\}}}(c), \text{opt}_{\mathcal{H}_{\{w_i\},1}(\zeta)}(c)) = o(1).$$

9 *B.2. Sketch of proof of points [6] and [7], Theorem 5*

11 The proof of these two points is easier than the others. Let us consider the sixth constraint to illustrate it. The function F_6 is exactly a decreasing function of the “holes” k' , which we can write

$$13 \quad \forall k' \in \{0, 1, |X|\}, F_6(k') = \frac{|X| - k'}{|X|}.$$

15 Fix γ strictly between 0 and $\frac{1}{2}$ (thus, the error is only slightly better than that of the unbiased coin). Weights are as follows for positive examples (we do not use p_1):

$$\begin{aligned} \forall (i, j) \in E, w(p_{ij}) &= w^+ = \frac{w^-}{(|X| + k)^2}, \\ 17 \quad w(p_2) &= \gamma - kw^- - \left(|E| - \binom{k}{2} \right) w^+. \end{aligned}$$

Weights are as follows for negative examples (we do not use n'_1):

$$\begin{aligned} 19 \quad \forall j \in \{1, \dots, |X|\}, w(n_j) &= w^- = \frac{1}{|X|^2 |E|^2}, \\ w(n'_2) &= 1 - \gamma - (|X| - k)w^- - \binom{k}{2} w^+. \end{aligned}$$

1 Fix $k_{\max} = 1 + \max_{2 \leq k'' \leq |X|; |E| - \binom{k''}{2} \geq 0} k''$. From the weights, $\text{lcm}(\cup_{x_i \in \mathcal{G}^+ \cup \mathcal{G}^-} d_i) = \mathcal{O}(|X|^8)$, which is polynomial. Define the function:

3 $\forall k' \in \{0; 1\}, E(k') = |E|w^+ + k'w^- + w(p_2),$

$$\forall 2 \leq k' \leq k_{\max}, E(k') = \left(|E| - \binom{k'}{2} \right) w^+ + k'w^- + w(p_2),$$

5 $\forall k_{\max} < k' \leq |X|, E(k') = k'w^- + w(p_2),$

(We have $E(k) = \gamma$). From that, it comes that the predicates that are not used can form a clique.

7 There remains to check the constraint values ζ which we allowed to take any value in $]0, 1[$. From Lemma 10, we may use $k = \theta(|X|^\alpha)$, for any $0 < \alpha < 1$. The fraction of authorized predicates is therefore upperbounded by

11 $\frac{|X| - k}{|X|} \leq 1 - \frac{1}{|X|^\alpha} \rightarrow_{\infty} 1.$

13 By considering sufficiently large sized graphs, the right side is greater than any chosen constant $0 < \zeta < 1$. Point [7] is achieved in the same way.

B.3. Proof of Theorem 7

15 Remark that $TPR \times (1 - FPR) = TPR \times TNR$. Weights are as follows for positive examples (we do not use p_2):

$$\begin{aligned} \forall (i, j) \in E, w(p_{i,j}) &= w^+ \\ &= \frac{\gamma}{(|X| - k)w^- \times \left(\binom{k}{2} + \frac{(|X| + 1)^2 - \left(k - \frac{|X| + 1}{3}\right)^2 - 3|X|}{6} \right)}, \end{aligned}$$

17 $w(p_1) = w^+ \times \left(\frac{(|X| + 1)^2 - \left(k - \left(\frac{|X| + 1}{3}\right)\right)^2 - 3|X|}{6} \right).$

Weights are as follows for negative examples (we do not use n'_2):

19 $\forall j \in \{1, \dots, |X|\}, w(n_j) = w^- = \frac{1}{|X| + k},$

$$w(n'_1) = 1 - |E|w^+ - |X|w^- - w(p_1).$$

21 The choice of γ_{\max} comes from the necessity of a tight calculation of the weights, in order to keep them in correct limits. In order to illustrate this, we proceed through the

1 proof of the correct values for the weights. The positive values of all weights (except
 2 from n'_1 , whose correctness stems from the study of all the other weights) is easily
 3 checked. However, we need to prove that they all take values which do not give a
 4 negative weight value to n'_1 .

5 Fix $k = \alpha|X|$, where α takes the adequate value $\alpha = \frac{5}{12}$ (other ones are possible, also
 6 valid according to Theorem 9, but we concentrate on this one).

7 Remark that $|X|w^- = \frac{1}{(1+\alpha)} = \frac{12}{17} < 1$. Now, we study $|E|w^+$. We have

$$w^+ = \frac{\gamma(|X| + k)}{(|X| - k) \left(\binom{k}{2} + \frac{(|X| + 1)^2 - (k - \frac{|X|+1}{3})^2 - 3|X|}{6} \right)}.$$

9 Note that $\gamma(|X| + k)/(|X| - k) = 17\gamma/7$. Suppose we choose $\gamma \leq 1/17 \times 7 \times 25/17 \times 144$
 10 (details of γ are given for the sake of clarity in the proof). Then, with such values,
 11 we have for the denominator of w^+ :

$$\left(\binom{k}{2} + \frac{(|X| + 1)^2 - (k - \frac{|X|+1}{3})^2 - 3|X|}{6} \right) \leq \frac{25}{144 \times 2} |X|^2,$$

13 which leads to an upperbound for $|E|w^+$ which is (taking into account that $|E| < |X|^2/2$):

$$|E|w^+ < \gamma \frac{17 \cdot 144 \times 2}{7 \cdot 25 \times 2} = \frac{1}{17}.$$

15 An upperbound of $\gamma \times 17/7 \times 144 \times 2/25 \times |X|^2$ is also available for w^+ , which leads
 16 to (for $w(p_1)$):

$$17 \quad w(p_1) < \gamma \frac{17 \cdot 144 \times 2}{7 \cdot 25 \times |X|^2} \frac{|X|^2}{6} < \frac{1}{17}.$$

This shows that the weight of n'_1 is positive, as we claimed.

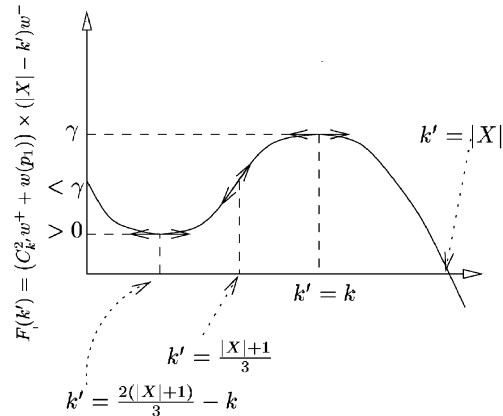
19 Now, we explain more in depth the proof scheme by describing a polynomial of order
 20 3, $F(k')$ which upperbounds $TPR \times TNR$, and of course has the desirable property of
 21 having its maximum for $k' = k$, with value γ , and with no other equal or greater values
 22 on the interval $[0, |X|]$. Similarly to the other proofs, the value γ can only be reached
 23 when $k' = k$ represents k “holes” among predicates $\{a_j(\cdot)\}$, and this induces a size- k
 24 clique in the graph.

25 Fix $k_{\max} = 1 + \max_{2 \leq k'' \leq |X|; |E| - \binom{k''}{2} \geq 0} k''$. From the weights, $\text{lcm}(\bigcup_{x_i \in \mathcal{G}^+ \cup \mathcal{G}^-} d_i) =$
 $\mathcal{O}(|X|^8)$, which is polynomial. Define the function:

27 $\forall k' \in \{0; 1\}, F(k') = w(p_1) \times (|X| - k')w^-$,

$$28 \quad \forall 2 \leq k' \leq k_{\max}, F(k') = \left(\binom{k'}{2} w^+ + w(p_1) \right) \times (|X| - k')w^-,$$

29 $\forall k_{\max} < k' \leq |X|, F(k') = (|E|w^+ + w(p_1)) \times (|X| - k')w^-$.

Fig. 2. Scheme of $F(k')$.

- 1 With our choice of weights, and inside the values of k' for which we described k
 2 (clearly, in the second curve), F describes a polynomial of degree 3, with a second-
 3 order derivative taking its zero for $k' = k'' = |X| + 1/3$. Its first-order derivative takes its
 4 zeroes, respectively, for $k' = k'_0 \in [0, k'']$ and $k' = k'_1 = k > k''$ (note that the choice of α
 5 respects this latter inequality). Outside $[k'_0, k]$, F is decreasing, and increasing inside.
 6 Since the choice of weights was also made so as to have $F(0) < \gamma$, and $F(k) = \gamma$, it
 7 is sufficient to prove that there is only one point for $k' = k$ where F takes a value of
 8 γ , with lower values elsewhere (Fig. 2 shows a simplified view of the function, for
 9 the sake of clarity). As we pointed out before, F upperbounds the product of TPR and
 10 TNR of any set of Horn clauses, which leads to a single favorable case: the “holes”
 11 inside the set of Horn clauses describe a clique of size k in the graph.

References

- 13 [1] S. Boucheron, Théorie de l'apprentissage, de l'approche formelle aux enjeux cognitifs. Hermes, 1992.
 14 [2] W.W. Cohen, Text categorization and relational learning, in: Proc. 12th Internat. Conf. on Machine
 15 Learning, 1995, pp. 124–132.
 16 [3] P. Domingos, MetaCost: a general method for making classifiers cost-sensitive, in: Proc. 5th Internat.
 17 Conf. on Knowledge Discovery in Databases, 1999, pp. 155–164.
 18 [4] S. Dzeroski, S. Muggleton, S. Russel, Pac-learning of determinate logic programs, in: Proc. 5th Internat.
 19 Conf. on Computational Learning Theory, 1992, pp. 128–137.
 20 [5] J.P. Egan, Signal detection theory and ROC analyses, Academic Press, New York, 1975.
 21 [6] J. Friedman, T. Hastie, R. Tibshirani, Additive logistic regression: a statistical view of boosting, Ann.
 22 Statist. 28 (2000) 337–374.
 23 [7] M. Garey, D. Johnson, Computers and Intractability, a guide to the theory of NP-Completeness, Bell
 24 Telephone Laboratories, 1979.
 25 [8] J. Håstad, Clique is hard to approximate within $n^{1-\epsilon}$, in: Proc. 37th IEEE Symp. on the Foundations
 26 of Computer Science, 1996, pp. 627–636.
 27 [9] K.-U. Höfgen, H. Simon, Robust trainability of single neurons, in: Proc. 5th Internat. Conf. on
 Computational Learning Theory, 1992.

- 1 [10] R. Holte, Very simple classification rules perform well on most commonly used datasets, *Mach. Learning*
11 (1993) 63–91.
- 3 [11] P. Jappy, R. Nock, O. Gascuel, Negative robust learning results for horn clause programs, in: *Proc.*
13th Internat. Conf. on Machine Learning, Morgan Kaufman, Los Altos, CA, 1996, pp. 258–265.
- 5 [12] M. Kearns, M. Li, L. Pitt, L. Valiant, On the learnability of boolean formulae, in: *Proc. 19th ACM*
7 Symp. on the Theory of Computing, 1987, pp. 285–295.
- 7 [13] J.-U. Kietz, Some lower bounds for the computational complexity of inductive logic programming, in:
9 *Proc. 8th European Conf. on Machine Learning*, 1993, pp. 115–123.
- 9 [14] J. Kietz, S. Dzeroski, Inductive logic programming and learnability, *Sigart Bull.* 5 (1994) 22–32.
- [15] A.R. Klivans, R. Servedio, Learning DNF in time $2^{\tilde{O}(n^{\frac{1}{3}})}$, in: *Proc. 33th ACM Symp. on the Theory*
11 *of Computing*, 2001, pp. 258–265.
- 13 [16] M. Kukar, N. Besic, I. Kononenko, M. Auersperg, M. Robnik-Sikonja, Prognosing the survival time
15 of the patients with anaplastic thyroid carcinoma using machine learning, in: N. Lavrac, E. Karavnou,
B. Zupan (Eds.), *Intelligent data analysis in medicine and pharmacology*, Kluwer Academic Publishers,
Dordrecht, 1997, pp. 116–129.
- 17 [17] M. Kubat, R.C. Holte, S. Matwin, Machine Learning for the detection of oil spills in satellite radar
19 images, *Mach. Learning* 30 (1998) 195–215.
- 19 [18] R. Nock, O. Gascuel, On learning decision committees, in: *Proc. 12th Internat. Conf. on Machine*
21 *Learning*, Morgan Kaufmann, Los Altos, CA, 1995, pp. 413–420.
- 21 [19] R. Nock, P. Jappy, On the power of decision lists, in: *Proc. 15th Internat. Conf. on Machine Learning*,
23 *Morgan Kaufmann*, Los Altos CA, 1998, pp. 413–420.
- 23 [20] F. Provost, T. Fawcett, Analysis and visualization of classifier performance : comparison under imprecise
25 class and cost distributions, in: *Proc. 3rd Internat. Conf. on Knowledge Discovery in Databases*, 1997,
pp. 43–48.
- 25 [21] F. Provost, T. Fawcett, Robust classification systems for imprecise environments, in: *Proc. 15th National*
27 *Conf. on Artificial Intelligence*, 1998, pp. 706–713.
- 27 [22] F. Provost, T. Fawcett, R. Kohavi, The case against accuracy estimation for comparing induction
29 algorithms, in: *Proc. 15th Internat. Conf. on Machine Learning*, 1998, pp. 445–453.
- 29 [23] R.E. Schapire, Y. Singer, Improved boosting algorithms using confidence-rated predictions, in: *Proc.*
31 *11th Internat. Conf. on Computational Learning Theory*, 1998, pp. 80–91.
- 31 [24] M. Sebban, R. Nock, A hybrid filter/wrapper approach of feature selection using information theory,
33 *Pattern Recognition* 35 (2002) 835–846.
- 33 [25] S. Slattery, M. Craven, Combining statistical and relational methods for learning in hypertext domains,
35 in: *Proc. 8th Internat. Conf. on Inductive Logic Programming*, 1998, pp. 38–52.
- 35 [26] A. Srinivasan, Applications of ILP to problems in chemistry and biology, in: *Invited talk at the 8th*
37 *Internat. Conf. on Inductive Logic Programming*, 1998.
- 37 [27] L.G. Valiant, A theory of the learnable, *Commun. ACM* 27 (1984) 1134–1142.
- 39 [28] L.G. Valiant, Learning disjunctions of conjunctions, in: *Proc. 9th Internat. Joint Conf. on Artificial*
Intelligence, 1985, pp. 560–566.
- [29] V. Vapnik, *Statistical learning theory*, John Wiley, New York, 1998.