



# Inductive-deductive systems: a mathematical logic and statistical learning perspective

Nicolas Baskiotis, Michèle Sebag, Olivier Teytaud

## ► To cite this version:

Nicolas Baskiotis, Michèle Sebag, Olivier Teytaud. Inductive-deductive systems: a mathematical logic and statistical learning perspective. CAP, 2007, Grenoble, France. inria-00173259v1

**HAL Id: inria-00173259**

**<https://inria.hal.science/inria-00173259v1>**

Submitted on 19 Sep 2007 (v1), last revised 1 Nov 2010 (v2)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Inductive-Deductive Systems: A mathematical logic and statistical learning perspective

Nicolas Baskiotis, Michele Sebag, Olivier Teytaud

Equipe TAO (Inria), LRI, UMR 8623 (CNRS - Université Paris-Sud),  
bât 490 Université Paris-Sud 91405 Orsay Cedex France,  
olivier.teytaud@inria.fr, michele.sebag@lri.fr, nicolas.baskiotis@lri.fr

**Résumé** : The theorems about incompleteness of arithmetic have often been cited as an argument against automatic theorem proving and expert systems. However, these theorems rely on a worst-case analysis, which might happen to be overly pessimistic with respect to real-world domain applications.

For this reason, a new framework for a probabilistic analysis of logical complexity is presented in this paper. Specifically, the rate of non-decidable clauses and the convergence of a set of axioms toward the target one when the latter exists in the language are studied, by combining results from mathematical logic and from statistical learning.

Two theoretical settings are considered, where learning relies respectively on Turing oracles guessing the provability of a statement from a set of statements, and computable approximations thereof. Interestingly, both settings lead to similar results regarding the convergence rate towards completeness.

Key words : Deduction, Induction, Algorithmic Learning Theory

## 1 Introduction

Inspired by the “Learning to Reason” framework Khardon & Roth (1997), this paper investigates the conditions for a hybrid inductive-deductive system (IDS). This system is provided with a set of axioms or statements (e.g. examples), and its goal is to determine the truth value of any further statement  $e$ . We consider a framework for dealing with undecidable theories as well ; this is a main difference with many previous works (Shapiro (1981)). We will often refer to arithmetic or set theory, but many other essentially undecidable theories could be considered instead of this.

From a mathematical logic perspective, the question is whether i) the available set of statements is complete, and ii) the logical setting is complete. Under these assumptions, the truth value of  $e$  is determined using mathematical deduction ; the algorithmic challenge is to provide an efficient search engine for constructing a proof of  $e$  or  $\neg e$ .

When the set of statements is not complete, by definition there exists statements  $e$  which can neither be proved nor refuted ; the famous Gödel’s theorem (1931) states that sufficiently powerful logical settings (e.g. including arithmetic) are incomplete.

When the set of statements is not sufficient for deciding relevant statements, inductive reasoning is needed to find additional axioms, consistent with the available ones and sufficient for determining the truth value of  $e$ . The challenge here is to compare the different natural methods available for adding new axioms.

From a hybrid inductive-deductive perspective, the logical setting considered must thus be examined with respect to both its completeness (deduction-oriented performances), and its VC-dimension or PAC learnability (induction-oriented performances, Appendix A.2).

Typically, statements  $C(1)$ ,  $C(3)$ ,  $C(5)$ ,  $C(7)$ ,  $\neg C(4)$ ,  $\neg C(6)$ ,  $\neg C(2)$ ,  $C(217)$ ,  $\neg C(200)$  do not allow deduction of  $\forall n C(2n+1) \wedge \neg C(2n)$ . In the meanwhile, inductive logic programming might learn the hypothesis  $\forall n C(2n+1) \wedge \neg C(2n)$ , which could in turn allow for many other deductions.

This paper examines the convergence properties of an inductive-deductive system, i.e. the probability that the  $n+1$ -st example can be proved from the axioms learned from the previous  $n$  examples. The originality of the work is to propose a probabilistic analysis of logical decidability and completeness, contrasting with the worst-case analysis and undecidability results used in the literature. Indeed, a worst-case perspective does not account for the fact that many relevant statements can yet be proved in an undecidable setting.

The rest of this section describes the proposed framework, discusses the relevant work and introduces the results reported in the paper.

**Formalisation.** This paper considers a first order logic language, where the initial set of axioms  $\mathfrak{Z}$  is an essentially undecidable ((Kleene, 1967 Dover 2002, p277), Tarski (1949))<sup>1</sup> set of axioms with finite description length<sup>2</sup> such as the Zermelo-Fraenkel set of axioms.

Let us consider a sequence of examples or statements  $e_i$ , independently and identically distributed from a probability distribution  $M$ . We further assume that  $M$  is consistent with  $\mathfrak{Z}$ , in the sense that  $\mathfrak{Z} \cup \{e \text{ s.t. } M(e) > 0\}$  is consistent.

From each set of examples  $\mathfrak{E}_n = \{e_1, \dots, e_n\}$ , the system extracts a recursive set of axioms noted  $A_n$ , which together with  $\mathfrak{Z}$  allows for proving every example in  $\mathfrak{E}_n$ .

The main subject of this paper is the analysis of three types of induction :

- in *deduction*,  $A_n$  includes all examples in  $\mathfrak{E}_n$ , except those examples  $e_i$  which could be proved from the theory learned from the previous examples (i.e. all  $e_i$  except those such that  $A_{i-1}, \mathfrak{Z} \vdash e_i$ ).  $A_n$  thus is an independent axiom set.
- in *pruned-deduction*,  $A_n$  is a minimal subset of  $\mathfrak{E}_n$ , sufficient to prove every example in  $\mathfrak{E}_n$  ( $A_n, \mathfrak{Z} \vdash \mathfrak{E}_n$ ).
- in *induction-deduction*,  $A_n$  is a set of axioms, minimal wrt its description length, such that every example in  $\mathfrak{E}_n$  can be proved from  $A_n$  and  $\mathfrak{Z}$ . Contrasting with deduction and pruned-deduction,  $A_n$  is no longer necessarily included in  $\mathfrak{E}_n$ .

The theoretical case where (deductive, pruned-deductive and inductive-deductive) learning is based on Turing oracles will first be considered in sections 3 and 4, respectively devoted to the cases where the target set of axioms is finite and infinite.

<sup>1</sup> A set of axioms is essentially undecidable if any recursive extension of this set is undecidable.

<sup>2</sup> In all the paper, the description length refers to any classical mathematical notation of statements or proofs. Note that a set of axioms with finite description length can include an infinite axiom schema.

Section 5 extends the analysis, considering Turing-computable-approximations of Turing oracles, based on finite-length proofs.

**Goals of the study.** The behavior of an inductive-deductive system is examined with respect to three stochastic variables, modeling respectively the completeness, the accuracy and the compactness of the current axiom set (noted  $A_n$  in the following instead of  $A_n, \mathfrak{Z}$  for simplicity of notation) :

- the incompleteness of  $A_n$  refers to the probability  $L_n$  that  $A_n$  does not allow for deciding on further examples ( $L_n = M(\{e \text{ s.t. } A_n \not\vdash e \wedge A_n \not\vdash \neg e\})$ ) ; in order to distinguish this incompleteness from the standard logical one, it will be referred to as the *relative* incompleteness ;
- the error or falsity of  $A_n$  is the probability that  $A_n$  decides wrongly on further examples<sup>3</sup> ;
- the compactness is measured as the description size  $DL(A_n)$  of the current axiom set (this is not related to other definitions of compactness, but we keep this notion as no ambiguity arises).

It must be noted that the behavior of the relative incompleteness rate  $L_n$  is known in some specific cases :

- In case  $\mathfrak{Z}$  is a complete set of axioms, no learning is required ;  $A_n = \mathfrak{Z}$  leads to  $L_n = 0$  as for any  $e$ ,  $\mathfrak{Z} \vdash e$  or  $\mathfrak{Z} \vdash \neg e$ .
- Otherwise, if  $M$  is modified at each time step  $n$  by a malign adversary with unrestricted computational power, then after Gödel's first theorem,  $L_n = 1$  for all  $n$ , as the adversary can choose  $M$  concentrated on some undecided  $e$ .

Thus, our framework lies between the (too simple, unrealistic) complete case, and the (too difficult, pessimistic) straightforward application of essential undecidability.

The goal is to examine the practical limitations of learning in a powerful language (e.g. including the axioms of set theory or arithmetic). The limitations of such languages regarding completeness and decidability issues are well known ; these limitations have significant impact on inductive learning as well. For example, in languages including arithmetic there are always infinitely many theories proving a finite consistent set of statements ; discussions around this fact are referred to as Quine's underdetermination thesis List (1999); J.D. (2003); Shook (2002).

However, it might be the case that the problems entailed by incompleteness and undecidability, though certain, are actually *not* frequent. And if there are an infinite number of solutions to a finite learning problem (Quine's underdetermination thesis), then it might be interesting to assess the average quality of these solutions.

Therefore, our goal is to provide a statistical study of the relative incompleteness, compactness and falsity of learned theories, applying the statistical learning methodology and body of results to other learning criteria, namely the probability of facing an undecided example or introducing inconsistencies in the theory.

**Related work.** As far as we know, the simultaneous use of deduction and induction has not been studied yet in a statistical perspective though the three domains involved (automatic deduction, inductive and statistical learning, mathematical logic) have some

---

<sup>3</sup>Note that in the deduction or pruned-deduction cases,  $A_n$  cannot be inconsistent with  $e_{n+1}$  since  $A_n \subseteq \mathfrak{E}_n$  and distribution  $M$  is assumed to be consistent with  $\mathfrak{Z}$ .

intersections<sup>4</sup>.

Along an inductive-deductive setting, the works related to Quine's underdetermination thesis List (1999); J.D. (2003); Shook (2002) focused on a worst case analysis ; they do not integrate the statistical learning and generalization aspects.

Other studies deal with some kinds of incomplete frameworks, e.g. involving recursion theory and referring to Turing machines with oracles or infinite-time Turing machines Hamkins (2002). However, this work focuses on extending the set of decidable statements and the role of induction is not considered.

**Overview of the paper.** As an alternative to the worst-case analysis, the framework proposed in this paper is based on a logically consistent probability distribution  $M$  over the set of statements. In each step  $n$ , the system outputs a set of axioms  $A_n$  from the first  $n$  statements, and one examines whether this set allows for proving further statements.

As noted earlier on, if these further statements are selected in a worst-case manner,  $A_n$  does not allow for deciding their truth value even with unbounded computational resources. However, a worst-case perspective often leads to overly pessimistic conclusions Cheeseman *et al.* (1991).

The probabilistic setting proposed is inspired by the standard Probably Approximately Correct (PAC) framework Valiant (1984), and the study borrows the standard statistical learning tools (VC-dimension Vapnik & Chervonenkis (1974)) in order to bound the relative incompleteness expectation  $L_n = M(\{e \text{ s.t. } A_n \not\vdash e \wedge A_n \not\vdash \neg e\})$ .

The paper is organized as follows. Section 2 introduces general definitions and lemmas used in the rest of the paper. Section 3 presents results about the induction of a target theory with bounded description length, comparing the *deductive*, *pruned-deductive* and *inductive-deductive* learning settings. It is shown that (corollaries 1-4) : i) in all cases, non-asymptotic performance depends on the underlying distribution  $M$  and it might be arbitrarily bad (as in the worst-case setting); ii) *induction-deduction*, and more generally restrictions on the description length entails faster convergence rates than *deduction*; iii) for any algorithm with a faster completeness convergence rate than *deduction*, there exists a distribution such that the error or falsity is not almost surely zero ( $\exists M, e \text{ s.t. } \forall n, P(A_n \vdash \neg e) > 0$  and  $M(e) > 0$ ); iv) *pruned* learning can behave arbitrarily badly in the sense of an infinite asymptotic description length.

Section 4 considers the case of a target theory with infinite description length, and presents negative results (corollaries 5-8) : i) arbitrarily slow convergence rates can occur; ii) the length of the axiom set can increase fast. However, the completeness rate goes to 1 as the number of examples goes to infinity.

While results presented in sections 3 and 4 are based on an oracle (axiomatic optimization or theorem proving with unbounded computational power), section 5 considers the case of Turing-computable approximations of such an oracle. Results similar to those of the oracle case are presented (with, unfortunately, a huge computational complexity). The paper ends with a discussion of the presented results ; a short introduction to the terminology and state of the art is given in Appendix A.

---

<sup>4</sup>Some advances in mathematical logic have been exploited for automatic theorem proving, for instance Craig's interpolation theorem is used to design a "partition-based" logic Amir & McIlraith (2003).

## 2 Formal background

Let  $\mathfrak{J}$  denote a consistent essentially undecidable set of axioms (e.g. Zermelo-Fraenkel). We note  $DL(A)$  the description length of an axiom set  $A$ , and by abuse we use also  $DL(e) = DL(\{e\})$ . Let  $T'$  denote the set of consistent theories including  $\mathfrak{J}$ , and let  $T \subset T'$  be the set of consistent theories defined from an axiom set with finite description length ( $T = \{t \in T', \exists A \text{ s.t. } A \vdash t, DL(A) < \infty\}$ ).  $T$  and  $T'$  can be viewed as boolean mappings from the set of well-formed statements ( $t(e) = 1$  iff  $t \vdash e$ ).

This section examines the VC-dimensions and shattering properties of both  $T$  and  $T'$  spaces. These results can be viewed as a partial statistical interpretation of Quine's under-determination thesis List (1999), J.D. (2003), Shook (2002).

**Theorem 1 :**  *$T$  has infinite VC-dimension (see Appendix B.1).*

The infinite VC-dimension also holds for propositional logic with infinitely many propositional variables. This theorem points out that this property is preserved in the framework defined above, i.e. the VC-dimension of the extensions of an essentially undecidable theory  $\mathfrak{J}$ .

Although  $T$  and therefore  $T'$  both have infinite VC-dimensions, they differ by their shattering properties :

**Theorem 2 :**  *$T'$  shatters an infinite set (see Appendix B.2).* The above theorem implies significant differences about learning in the search spaces  $T$  and  $T'$ . Specifically, in the case of  $T'$  there exists distributions leading to arbitrarily slow convergence rates, such as  $C/\log(\log(\log(n)))$  where  $n$  is the number of examples. In contrast, we shall see that a reasonable convergence rate is obtained within  $T$ , although the convergence can be delayed due to adverse distributions.

## 3 The finite description length case

Let  $M$  denote a probability distribution on the well-formed statements, such that the mass of  $M$  is restricted to a consistent theory  $t$  in  $T$  ( $\exists t \in T \text{ s.t. } M(t) = 1$ ). A first negative result concerns the incompleteness convergence rate. We show that there exists a distribution  $M$  such that the incompleteness rate is bounded from below.

In the corollary below, as in corollary 5, we will use a link between supervised learning (i.e. learning statements with their truth values) and unsupervised learning (i.e. finding theories covering true statements). This link is based on the fact that stating  $e$  is exactly equivalent to stating  $\neg e$ . A distribution  $M$  on couples  $(x, y) = (e, \text{true})$  or  $(x, y) = (e, \text{false})$ , where  $e$  is a statement, can be replaced by a distribution  $M$  such that  $M(e)$  is the probability of  $(e, \text{true})$  plus the probability of  $(e, \text{false})$ , as well as we can identify sets of axioms with classifiers (the associated classifier separates theorems and non-theorems). This allows the use of counter-examples from learning in the framework of this paper. A family of classifiers (for supervised learning) such that for any learning algorithm  $\mathbb{E}L_n \geq c$  for some distribution on these classifiers, is identified with a family of sets such that for any algorithm, for some distribution  $\mathbb{E}L_n \geq c$ .

**Corollary 1 :** *for any  $n > 0$ , for any  $\delta > 0$ , for any method generating  $A_n$ , there exists a generator of examples (distribution  $M$ ) such that  $\mathbb{E}[L_n] \geq \frac{1}{2\exp(1)} - \delta$ .*

**Remark :** This result can be reformulated as : for any  $n > 0$ , for any  $\delta > 0$ , there exists  $M$  such that after  $n$  examples the learned theory is at distance at least  $\frac{1}{2\exp(1)} - \delta$  from the target one. Note that the above distribution  $M$  depends on  $n$ .

**Proof of corollary 1 :** Follows from theorem 1 and the lower bound cited in Appendix A.2.  $\square$

Let us first consider the *fine* learning case, restricting the description length of the induced axiom set. By abuse of notation, in the following the description length of a theory derived from an axiom set  $A$  is set to  $DL(A)$ .

**Theorem 3 :** Let  $T_s$  denote the set of theories in  $T$  which can be generated from a set of axioms with description length less than  $s$ , and let  $V_s$  denote the VC-dimension of  $T_s$ .  $V_s$  is finite as  $T_s$  is finite. For each theory  $t$ , let  $V(t)$  be defined by  $V(t) = \inf\{V_s | t \in T_s\}$ .

Let  $V$  denote  $V(t^*)$  where  $t^*$  is the target theory, assuming it is finite. Let  $t_n$  denote the theory extracted along induction-deduction after  $n$  examples.

Then the following results hold (see Appendix B.3) :

1. **Convergence rate :** if  $n > V$ ,  $P(L_n > \epsilon) \leq 2(2\exp(1)n/V)^V 2^{-n\epsilon/2}$ .
2. **Asymptotic behavior :** almost surely, there exists  $n_0$  such that  $n \geq n_0 \Rightarrow L_n = 0$ .
3.  $T$  does not shatter any infinite set.
4.  $V(t_n) \leq V$ .

Let us now consider the *deduction* and *pruned-deduction* learning settings.

**Theorem 4 :** Consider the deduction or pruned-deduction settings. For any decreasing sequence  $a_n$  bounded by  $1/2$  and converging to 0, there exists a probability distribution  $M$  such that i)  $\forall n, L_n \geq a_n$  ; ii) there exists  $t \in T$  (i.e.  $DL(t) < \infty$ ) such that  $M(t) = 1$  (see Appendix B.4).

**Corollary 2 :** In the deduction or pruned-deduction framework, for any decreasing sequence  $a_n$  upper bounded by  $1/2$ , there exists  $M$  such that for any  $n$  the relative incompleteness of  $A_n$  is bounded from below by a sum of  $n$  independent binary random variables  $X_i$ , where  $X_i$  takes value  $1/0$  with probability  $(a_i, 1 - a_i)$ . In particular, under distribution  $M$ , the relative incompleteness of  $A_n$  is greater than  $\sum_{i \leq n} a_i$ .

**Corollary 3 :** Theorem 4 and corollary 2 can be extended to any learning method producing a minimal (wrt set inclusion) theory such that it covers examples  $e_1, \dots, e_n$ .

**Proof :** This is a direct corollary of the proof of theorem 4.  $\square$

**Corollary 4 :** Any method which does not incur the limitations stated by theorem 4 or corollary 2, can with non-zero probability select a theory  $A_n$  which is strictly larger (wrt set inclusion) than the minimal theory generated from  $\{e_1, \dots, e_n\}$ . In particular, for some distribution, there is a positive probability of generating a theory inconsistent with some statement of non-zero measure.

**Proof :** Reformulation of corollary 3.  $\square$

## 4 The infinite description length case

Removing the finite length assumption has significant impact on the convergence results obtained in the previous section, notably in relation to the lower bounds on the

convergence rate (see Appendix A.2). In the proof of the following corollary, we use the same correspondence as explained before corollary 1.

**Corollary 5 :** *for any decreasing sequence  $(a_n)$  running to 0 and upper-bounded by  $1/16$ , there exists a distribution of examples such that  $\mathbb{E}(L_n) \geq a_n$ .*

**Proof :** Direct consequence of the lower bound on the convergence rate in Appendix A.2.  $\square$

**Corollary 6 :** *for the deduction or pruned-deduction settings,  $\mathbb{E}[DL(A_n)] \geq \sum_{i \leq n} a_i$ . (see Appendix B.5)*

**Corollary 7 :** *In all learning cases such that the empirical error rate is null,  $L_n$  goes to 0 (in the sense that for any  $\epsilon > 0$ ,  $P(L_n > \epsilon) \rightarrow 0$ ). (see Appendix B.5)*

**Corollary 8 :** *For all learning methods such that  $A_n$  is consistent and proves all statements  $e_1, \dots, e_n$ , there exists a distribution  $M$  such that the compactness  $DL(A_n)$  goes to infinity.*

**Proof :** Consider  $M$  a distribution which support is a consistent non-recursively axiomatizable set of statements. For any axiom set  $A$  let  $M(A)$  be the measure for  $M$  of all statements proved from  $A$  ( $M(A) = \sum_{e \text{ s.t. } A \vdash e} M(e)$ ). Consider  $K(\epsilon)$ , the minimal description length over all axiom sets  $A$  such that  $M(A)$  greater than  $1 - \epsilon$ .  $K(\epsilon)$  is non-decreasing, and  $\lim_{\epsilon \rightarrow 0} K(\epsilon) = \infty$ . It is sufficient to see that  $L_n < \epsilon$  implies that  $DL(A_n) > K(\epsilon)$ .  $\square$

These results show that although the error rate goes to 0 as the number of examples increases, the convergence rate can be arbitrarily low. Moreover, the description length of the induced theory cannot be bounded, as showed above.

## 5 Turing-computable algorithms : proofs of bounded length

By definition, deduction, pruned-deduction and induction-deduction all rely on Turing machines with oracles. As a first step toward a practical analysis, this section considers instead approximate learning, based on Turing machines without oracles and bounded length reasoning<sup>5</sup>. The approximation is considered from an algorithmic complexity perspective.

Section 5.1 is devoted to a complexity analysis of axiomatic optimization. This result is used in section 5.2 to provide a bound on the convergence of pruned-deductive and deductive-inductive learning toward the target theory, in the case where the latter is finite.

### 5.1 Algorithmic complexity

Let us consider as hypothesis space the sets of axioms with finite description length (possibly including axiom schemas ; proofs using axiom schemas are allowed as well).

---

<sup>5</sup>Since recursion theory provides negative results in the case of proofs with arbitrary length (ie, the set of statements that can be proved, in many cases, is not recursive but only recursively enumerable), we restricted this study to proofs with bounded length.



Let us define the  $k$ -deduction as follows : statement  $e$  is  $k$ -proved from the set of axioms  $A$ , noted  $A \vdash_k e$ , if there exists a proof of  $e$  from  $A$  with description length less than  $k$ .

The algorithmic complexity of  $k$ -deduction (i.e. the complexity required to decide the fact that a statement  $e$  is  $k$ -proved from  $A$ ) is upper bounded by a function noted  $Complexity(DL(A), k, DL(e))$ . In the arithmetic setting considered,  $Complexity(DL(A), k, DL(e))$  is dominated by a  $2^k$  term (considering all  $2^k$  strings of length  $k$  and determining whether they are proofs of  $A \vdash e$ ).

Along the same lines, the  $k$ -consistency of a set of axioms is defined as follows :  $A$  is  $k$ -consistent, noted  $A \nvdash_k \perp$  if there is no proof with length smaller than  $k$  that  $A$  is inconsistent.

Similarly, the algorithmic complexity of  $k$ -consistency is upper bounded by  $Complexity(DL(A), k, DL(\perp))$ . Therefore, the generality and consistency tests (respectively,  $B \vdash_k A$  and  $B \nvdash_k \perp$ ) can be performed with complexity  $\sum_{e \in A} Complexity(DL(B), k, DL(e))$ , over all statements or axiom schemas  $e$  in  $A$ . The following proposition is then straightforward.

**Proposition : Complexity of axiomatic optimization**

*Assume that there are at most  $2^n$  sets of axioms with description length less than  $n$ . Given a set  $A$  of statements, axiomatic optimization aims at a set of axioms  $B$  with minimal description length such that it entails all statements in  $A$  : Find  $\text{Arg min}_B \{DL(B) | B \vdash_k A, B \nvdash_k \perp\}$ . Its complexity is upper bounded by  $O(2^{DL(A)} \times DL(A) \times Complexity(DL(A), k, DL(A)))$ .*

## 5.2 Axiomatic optimization

Given a set  $\mathfrak{E}_n$  of statements, the point here is to find a minimal set  $A_n$  of axioms, using a finite number  $\delta_n$  of computation steps to check  $A_n$  consistency and completeness wrt  $\mathfrak{E}_n$ . We show that if  $\delta_n$  increases sufficiently fast, there exists an algorithm with essentially same convergence results as in the oracle-based analysis (section 3).

Practically, let  $\delta_1, \dots, \delta_n$  denote a sequence of integers. Then :

- Let  $T_{n+1}$  be the set of statements that can be proved with proofs of length at most  $\delta_n$  from  $A_n$  ;
- $A_n$  is a<sup>6</sup> minimal description length set of axioms such that : i)  $A_n$  proves all examples in  $\mathfrak{E}_n$  with proof of length at most  $\delta_n$  ; ii)  $A_n$  does not prove  $\perp$  with proof of length at most  $\delta_n$ .

Note that  $A_n$  is not necessarily a minimal set of axioms in the usual sense, e.g. one of the axioms could be proved from the others (but its presence makes it feasible to prove  $k$ -completeness).

We note  $A^*$  the shortest axiom set capable of proving any  $e$  in the target theory (i.e. such that  $M(e) > 0$ ). Let  $L_n$  here denote  $L_n = M(\{e; e \notin T_{n+1} \vee (\neg e) \in T_{n+1}\})$ .

**Theorem 5 :** *Consider  $e$  a random variable on statements with probability law  $M$  and assume that the mean and the variance of the shortest proof of  $e$  from  $A^*$  are finite, and assuming further that  $\delta_n = \Omega(n^3)$ , then (Appendix B.6) ,*

---

<sup>6</sup>In case of equality, the first axiom set in lexicographic order is retained.

1.  $P(DL(A_n) \geq DL(A^*) + \epsilon) \leq p_{n,\epsilon}$  with  $p_{n,\epsilon}$  is  $O(1/n^2)$ , as soon as  $n = \Omega(1/\sqrt{\epsilon})$ .
2.  $A_n$  reaches  $A^*$  almost surely, and thus is consistent for  $n$  sufficiently large.
3.  $L_n \leq \epsilon'$  with probability at least  $1 - p_{n,\epsilon} - 2S2^{-n\epsilon'/2}$  for any  $\epsilon, \epsilon' > 0$ , where  $S$  is the number of axioms sets with description length bounded by  $DL(A^*) + \epsilon$ .

This result shows that the theory extracted by induction-deduction (using Turing machines with no oracle and bounded-deduction) is consistent, for sufficiently large number of examples ; that its description length converges toward the optimal one, and finally, that its relative incompleteness goes to 0 as  $O(1/\sqrt{n})$ .

In summary, Theorem 5 shows that a Turing-machine algorithm with no oracle can implement an inductive-deductive system, with essentially the same performances and limitations regarding consistency and completeness as in the theoretical case. Indeed the complexity of this algorithm is exponential in  $n$ .

## 6 Discussion

A probabilistic relational setting has been proposed in this paper to study inductive-deductive systems (IDS). Precisely, from a random generator providing statements and their truth values, the IDS extracts a set of axioms via one among three settings : the *deduction* one corresponds to a purely deductive algorithm ; the *pruned-deduction* one extracts a minimal excerpt of the statements, sufficient to prove all seen statements ; and the *inductive-deductive* setting selects the set of axioms with minimal description length such that it proves all seen statements.

Two cases are distinguished : the “finitely describable” (FDR) and “non-finitely describable” (NFDR) realities respectively correspond to the case where the target set of axioms has a finite (resp. infinite) description length.

FDR and NFDR cases are confronted to *deduction*, *pruned-deduction* and *induction-deduction* settings, considering two criteria : relative incompleteness (proportion of statements which cannot be proved from the current theory) and compactness (description length of the current theory).

Though relative incompleteness always goes to 0 as the number of examples goes to infinity, its convergence rate can be arbitrarily low in all cases, except when reality is finitely describable and in a *inductive-deductive* setting, in other words, when the system actually performs induction. In this favorable case, the target concept is reached almost surely in finite time.

Along the same lines, the description length of the extracted theory is unbounded (for adverse distributions) in all cases, except again when reality is finitely describable and in a *inductive-deductive* setting.

This result provides additional precisions related to Quine’s under-determination thesis. Despite the multiplicity of theories consistent with a finite set of statements, if the IDS system extracts the theory consistent with the statements already seen, that is minimal wrt its description length (as opposed to, wrt its set inclusion), then a fast convergence in terms of both incompleteness and length can occur granted that the reality is “finitely describable”.

Interestingly, there exists some distributions in the latter case which entail errors and not only undecidabilities ; i.e. there are cases such that the event  $\exists n; A_n \vdash \neg e$  has strictly positive probability. This result can be interpreted in the light of Popper's notion of falsifiability, central to the history of science ; as shown by the very general corollary 4, if one abstains from producing hypotheses which can be falsified by examples, the convergence rate of the IDS is not better than that of rote learning.

The last part of this paper has shown that the above theoretical results, obtained for Turing machines with oracles, essentially hold for Turing machines *without* oracles — although the considered algorithms are indeed of limited use due to their huge computational complexity.

In summary, the main ambition of this paper is to contribute to a less pessimistic view of inductive-deductive systems in relational logic, than allowed by a worst-case analysis and based on undecidability results. In particular, we show that better rates than purely deductive systems are possible. Also we point out the similarity with mathematical reasoning : replacing a long axiom by a shorter (non-equivalent and more general) one is possibly a good idea. For example, arguments in favor or against the continuum hypothesis (CH) or the axiom of choice (AC) are "shorter" axioms : see e.g. Woodin (2001) and Freiling (1986) for CH, Solovay's axiom for AC, or links between AC and generalized-CH shown by Sierpinski.

## Acknowledgements

This work was partially supported by the Pascal Network. We thank reviewers for their help.

## A State of the art and definitions

This appendix briefly summarizes the notations and learnability results used in the paper.

Notation  $\sup X$ , where  $X$  is a real-valued random variable, denotes the (possibly infinite) supremum of the  $x$  such that  $P(X > x) > 0$  ;

### A.1 Logical notations

A theorem is a statement which can be proved, which depends on both the logical setting and the axiom set considered. The paper only considers classical logic. Each axiom set  $A$  includes  $\exists$  and has finite description length. After Gödel's theorem, there exists thus  $e$  such that neither  $e$  nor  $\neg e$  can be proved from  $A$ .

For the feasibility of the study, it is assumed that  $\exists$  is consistent, although in many cases of interest, this has not been proved (and cannot be, e.g. for Zermelo Fraenkel, after Gödel's theorem).

Notation  $A \vdash e$  (respectively  $A \vdash_k e$ ) denotes the fact that  $e$  can be proved from  $A$  (resp. with proof of description length less than  $k$ ).

In the whole paper, the description length  $DL(\cdot)$  (of sets of axioms or proofs) refers to a standard logic coding (with no compression).

### A.2 Statistical learning theory

The interested reader is referred to Devroye *et al.* (1997), Vidyasagar (1997) for an exhaustive presentation.

Let  $Z$  denote the example space, and let  $F$  denote the hypothesis space, in which each hypothesis is viewed as a subset of  $Z$ . A set  $X$  of examples is said to be shattered by  $F$  if for any subset  $X'$  of  $X$  there exists  $f \in F$  such that  $f \cap X = X'$ . The **VC-dimension** of  $F$  is the cardinal of the largest finite set that is shattered by  $F$ . If arbitrarily large such sets exists, then the VC-dimension is said infinite. Hypothesis  $h$  is consistent with a set of examples  $X$  iff  $h \cap X = h^* \cap X$ , where  $h^*$  denotes the target concept. A learning algorithm associates a consistent hypothesis  $h_n$  to each training set  $X_n$  made of  $n$  iid examples drawn according to some probability distribution  $M$ . Accordingly, the loss variable  $L_n$  stands for the error expectation of  $h_n$  ( $M\{z, z \in Z, h_n(z) \neq h^*(z)\}$ ).

Fundamental results in the statistical learning theory can be summarized as : if the VC-dimension is finite, then the error expectation goes to 0 reasonably fast as the number of examples goes to infinity.

#### Upper-bounds depending upon the VC-dimension :

**Theorem, case of null empirical error** (see Vapnik & Chervonenkis (1974), (Devroye et al., 1997, Th. 12.7, p202)) :

Define  $\hat{L}(P) = \frac{1}{s} \sum_{i=1}^s \mathbf{1}_{P(x_i) \neq y_i}$  and  $L(P) = \mathbb{E} \mathbf{1}_{P(X) \neq Y}$ , with the  $(x_i, y_i)$  a sample of size  $s$  iid according to the law of the random variable  $(X, Y)$ .

Consider  $\mathfrak{F}$  a family of boolean functions on a domain  $X$  and let  $V$  be its VC-dimension. Then, for any  $\epsilon > 0$  if  $s > V$ ,

$$P\left(\sup_{P \in \mathfrak{F}; \hat{L}(P)=0} |L(P) - \hat{L}(P)| \geq \epsilon\right) \leq 2(2\exp(1)s/V)^V 2^{-s\epsilon/2}$$

where  $(2\exp(1)s/V)^V$  can be replaced by the  $2s$ -shattering coefficient of  $\mathfrak{F}$ .

**Lower bound :** ((Devroye et al., 1997, p239), theorem 14.3). Assume that the VC-dimension of  $F$  is infinite. Then for any  $n > 0$ , for any  $\delta > 0$ , for any classification rule, there exists at least one distribution such that  $\mathbb{E}L_n \geq \frac{1}{2\exp(1)} - \delta$  and  $F$  contains at least a function  $f$  such that  $L(f) = 0$ .

**Lower bound on the convergence rate :** Assume that  $F$  shatters an infinite set. Then for any sequence  $(a_n)$  decreasing to 0 and upper bounded by  $\frac{1}{16}$ , for any classification rule, there exists at least one distribution such that  $\forall n \mathbb{E}L_n \geq a_n$  whereas  $F$  contains  $f$  such that  $L(f) = 0$ .

Mainly, the difference with the previous result is that the distribution does not depend upon  $n$ .

**Lemma : learning on countable domains** Consider learning on a countable domain with a distribution and an algorithm ensuring that the empirical error  $\hat{L}$  is zero. Then, the generalization error almost surely converges toward 0.

## B Proof details

### B.1 Proof of Theorem 1

**Theorem 1 :**  $T$  has infinite VC-dimension.

**Proof :** By definition, there exists at least one statement that cannot be decided in  $\mathfrak{J}$ . Let  $e_1$  denote this statement. By definition, the theory  $e_1^*$  (respectively  $\neg e_1^*$ ) generated from  $\mathfrak{J} \cup \{e_1\}$  (resp.  $\mathfrak{J} \cup \{\neg e_1\}$ ) is consistent.

Therefore, there exists another statement  $e_2$  that cannot be decided in  $e_1^*$ ; similarly, one defines  $(e_1 \wedge e_2)^*$  (resp.  $(e_1 \wedge \neg e_2)^*$ ) as the theory generated from  $e_1^*$  and  $e_2$  (resp.  $e_1^*$  and  $\neg e_2$ ).

Symmetrically, another statement  $e_3$  that cannot be decided in  $(\neg e_1)^*$  is selected and theories  $(\neg e_1 \wedge e_3)^*$  and  $(\neg e_1 \wedge \neg e_3)^*$  are defined.

All four theories are consistent ; any pair of them is inconsistent. By induction, a tree of theories is constructed ; at depth  $d$ ,  $2^d$  consistent theories generated from a conjunction of  $e_i$  or  $\neg e_i$  are obtained.

For a VC-dimension analysis, it is desirable that the  $2^d$  theories be based on the same  $d$  examples (which is not the case in the above, as  $e_2 \neq e_3$ ). The tree is thus rewritten as follows.

Note  $p_i(x)$  the  $i^{th}$  binary digit of  $x$  and  $f_i^d$  the conjunctions of statements  $e_j$  until depth  $d$  defined by

$$f_i^d = \bigwedge_{p_j(i)=0} e_j \bigwedge_{p_j(i)=1} \neg e_j, 1 \leq j \leq d$$

for  $i \in [0; 2^d - 1]$ . Intuitively each conjunction is a path from the root of the tree, where the left branch (respectively, right branch) is selected at depth  $i$  if  $p_i(x) = 0$  (resp.  $p_i(x) = 1$ ).

A new tree is built, in which each left branch at depth  $d$  is labelled with the statement  $l_d$ , defined as the disjunction of  $f_i^d$  over all  $i$  such that  $p_d(i) = 0$  ; symmetrically, each right branch at depth  $d$  is labelled with  $r_d$ , defined as the disjunction of  $f_i^d$  over all  $i$  such that  $p_d(i) = 1$ . For  $d = 3$ , the following holds

$$\begin{aligned} l_1 &= e_1, r_1 = \neg e_1 \\ l_2 &= (e_1 \wedge e_2) \vee (\neg e_1 \wedge e_3), r_2 = (e_1 \wedge \neg e_2) \vee (\neg e_1 \wedge \neg e_3) \\ l_3 &= (e_1 \wedge e_2 \wedge e_4) \vee (e_1 \wedge \neg e_2 \wedge e_5) \vee (\neg e_1 \wedge e_3 \wedge e_6) \vee (\neg e_1 \wedge \neg e_3 \wedge e_7) \\ r_3 &= (e_1 \wedge e_2 \wedge \neg e_4) \vee (e_1 \wedge \neg e_2 \wedge \neg e_5) \vee (\neg e_1 \wedge e_3 \wedge \neg e_6) \vee (\neg e_1 \wedge \neg e_3 \wedge \neg e_7) \end{aligned}$$

Consider  $v$  a binary vector of size  $n$  ; by abuse of notation, the same notation is used for  $v$  and the associated integer value. We use the notation  $v_{|p}$  for the vector of size  $p$  restricted to the  $p$  first coordinates of  $v$ . Let  $B_v$  denote the theory of  $T$  generated by statements  $\mathfrak{J} \cup \{l_i \text{ s.t. } v_i = 0\} \cup \{r_i \text{ s.t. } v_i = 1\}$ . Then, the following lemma holds :

**Lemma 1 :**  $l_i \in B_v$  if and only if  $v_i = 0$ .

**Proof :** A consistent theory containing  $f_i^d$  cannot contain  $f_j^d$  for  $j \neq i$ , by construction. The choice of  $v$  implies that  $B_v$  contains  $f_{v_{|p}}^p$  and none of the  $f_x^p$  for  $x \neq v_{|p}$ . If  $v_i = 0$ , by definition  $l_i$  is in  $B_v$ . If  $v_i = 1$ ,  $f = f_{v_{|i}}^i$  is in  $B_v$ ,  $l_i$  contains only some  $f_x^i$ 's different of  $f$ , therefore  $B_v$  does not contain  $l_i$ .

The set  $\{l_i, i \in [1; n]\}$  is therefore shattered by  $T$ .

This proves that arbitrarily large finite sets can be shattered by  $T$ . Therefore, the VC-dimension of  $T$  is infinite.  $\square$

## B.2 Proof of Theorem 2

**Theorem 2 :**  $T'$  shatters an infinite set.

**Proof of theorem 2 :** We note  $B$  the set of paths on the tree above and we identify it to corresponding elements of  $T'$  (theories generated by the set of statements along the path). We consider the family of statements  $\{l_n, n \geq 1\}$ . The following lemma achieves the proof :  $\square$

**Lemma 2 :** The set  $\{l_n, n \geq 1\}$  is shattered by  $B$ .

**Proof of lemma 2 :** Consider  $(V_n)$  an infinite boolean random sequence. The tree path that branches on  $l_i$  iff  $v_i = 0$  is associated to a theory  $t$  which belongs to  $T'$ . This theory satisfies  $l_n \in t$  iff  $V_n = 0$ .  $\square$

## B.3 Proof of Theorem 3

**Theorem 3 :** Let  $T_s$  denote the set of theories in  $T$  which can be generated from a set of axioms with description length less than  $s$ , and let  $V_s$  denote the VC-dimension of  $T_s$ .  $V_s$  is finite as  $T_s$  is finite. For each theory  $t$ , let  $V(t)$  be defined by  $V(t) = \inf\{V_s / t \in T_s\}$ .

Let  $V$  denote  $V(t^*)$  where  $t^*$  is the target theory, assuming it is finite. Let  $t_n$  denote the theory extracted along inductive-deductive learning after  $n$  examples.

Then the following results hold :

1. **Convergence rate** : if  $n > V$ ,

$$P(L_n > \epsilon) \leq 2(2 \exp(1)n/V)^V 2^{-n\epsilon/2}$$

2. **Asymptotic behavior** : almost surely, there exists  $n_0$  such that  $n \geq n_0 \Rightarrow L_n = 0$

3.  $T$  does not shatter any infinite set.

4.  $V(t_n) \leq V$ .

**Proof of theorem 3** : Let  $t_n$  be a<sup>7</sup> minimal length theory such that it covers all examples in  $\mathfrak{E}_n$  ;  $t_n$  is viewed as a random variable in  $T$ .

Since by construction  $t^*$  covers  $\mathfrak{E}_n$ ,  $V(t_n)$  is less than  $V(t^*)$ , which proves point 4.

Point 1 follows from the "null empirical error case" bound (Appendix A.2) :  $P(L_n > \epsilon) \leq 2(2 \exp(1)n/V)^V 2^{-n\epsilon/2}$ . Accordingly, one has :  $\forall \epsilon \sum_n P(L_n > \epsilon) < \infty$ .

The Borell-Cantelli lemma then shows that  $M(e \notin t_n) \rightarrow 0$  almost surely. However, this does not imply that a null error is reached in finite time.

On the other hand,  $t_n$  lies in a finite set as  $V(t_n)$  is upper bounded by  $V$ . Therefore  $\epsilon$ , the smallest of the  $M(e \notin t_n)$  is strictly positive. Thus, after a finite time, thanks to almost sure convergence,  $L_n < \epsilon$ , and therefore  $M(e \notin t_n) = 0$ .

The third point of the theorem can be shown in two different manners :

- if  $T$  shatters an infinite set, the lower bound theorem holds, which contradicts the convergence rate above. Therefore, there does not exist a shattered infinite set.
- $T$  is countable ; as a family of functions shattering an infinite countable set is at least in one-to-one mapping with  $P(\mathbb{N})$ , and thus cannot be countable, this shows that  $T$  cannot shatter an infinite set.  $\square$

## B.4 Proof of Theorem 4

### Theorem 4 :

Consider the deduction or pruned-deduction settings. For any decreasing sequence  $a_n$  bounded by  $1/2$  and converging to 0, there exists a probability distribution  $M$  such that i)  $\forall n, L_n \geq a_n$  ; ii) there exists  $t \in T$  (i.e.  $DL(t) < \infty$ ) such that  $M(t) = 1$ .

**Proof** : Define  $h_i = l_{i+1} \vee e_1$  with  $l_i$  defined as in Theorem 1. Then, for any  $i \geq 1$ , the theory  $t$  learned from  $(\mathfrak{Z} \cup \{e_1\})$  contains all the  $h_i$ . Consider  $M$  a law of probability on the  $h_i$  such that  $\sum_{i>n} M(h_i) \geq a_i$  and  $M(h_i)$  decreases as a function of  $i$ . Such a probability law can be found for instance in (Devroye *et al.*, 1997, Lemma 7.1, p114).

As deduction and pruned-deduction learning produces at best a subset of the encountered statements, they can only cover the  $h_i$  which have been seen by definition. Therefore, the measure for  $M$  of the  $h_i$  which can be proved is upper bounded by the measure of the  $h_i$ 's which have been seen, itself upper-bounded by  $\sum_{i=1}^n M(h_i)$  as the sequence of the  $M(h_i)$  decreases. Therefore,  $L_n \geq a_n$ .  $\square$

## B.5 Proof of section 4

**Corollary 6** : for the deduction or pruned-deduction settings,  $\mathbb{E}[DL(A_n)] \geq \sum_{i \leq n} a_i$ .

---

<sup>7</sup>As  $t_n$  may not be uniquely determined.

**Proof :** The distribution used in the proof of corollary 5 does only use the statements  $l_n$  of lemma 2. As these statements are independent (since they are shattered), the pruned-deduction step can only remove already-seen statements; therefore, both the deduction and pruned-deduction algorithm add one axiom per example  $e_i$  such that  $\exists, e_1, \dots, e_{i-1} \not\models e_i$ , which gives the result.  $\square$

**Corollary 7 :** *In all learning cases such that the empirical error rate is null,  $L_n$  goes to 0 (in the sense that for any  $\epsilon > 0$ ,  $P(L_n > \epsilon) \rightarrow 0$ ).*

**Proof :**  $1 - L_n$  is lower-bounded by the sum of the probabilities of the  $e_i$ , for  $i \leq n$  (counting each statement only once). Let us sort the statements (they are countable) by decreasing order of probability for  $M$ . Note them  $f_1, \dots, f_n, \dots$ . For any  $\epsilon$ , there exists an integer  $N$  such that  $\sum_{i=1}^N M(f_i) > 1 - \epsilon$ .

With probability one as  $n$  goes to infinity, all statements  $f_i, i \leq N$  are seen by the system. Therefore,  $P(L_n > \epsilon)$  goes to 0 as  $n$  goes to infinity.

The above holds for all learning methods such that  $A_n$  entails all previously seen statements  $e_1, \dots, e_n$ .  $\square$

## B.6 Proof of Theorem 5

We note  $A^*$  the shortest axiom set capable of proving any  $e$  in the target theory (i.e. such that  $M(e) > 0$ ). Let  $L_n$  here denote  $L_n = M(\{e; e \notin T_{n+1} \vee (\neg e) \in T_{n+1}\})$ .

In the sequel we note  $\sigma(X)$  the standard deviation of a random variable  $X$ ,  $\mu_e$  the expectation of  $DL(e)$ ,  $\sigma_e = \sigma(DL(e))$ ,  $P_{(k)}$  the probability that the shortest proof of  $e$  from  $A^*$  is larger than  $k$ , and  $\mathbf{1}(e) = 1$  if  $A^* \not\models_{\delta_n} e$  and  $\mathbf{1}(e) = 0$  otherwise.

### Theorem 5 :

Consider  $e$  a random variable on statements with probability law  $M$  and assume that the mean and the variance of the shortest proof of  $e$  from  $A^*$  are finite, and assuming further that  $\delta_n = \Omega(n^3)$ , then,

1.  $P(DL(A_n) \geq DL(A^*) + \epsilon) \leq p_{n,\epsilon}$  with  $p_{n,\epsilon}$  is  $O(1/n^2)$ , as soon as  $n = \Omega(1/\sqrt{\epsilon})$ .
2.  $A_n$  reaches  $A^*$  almost surely, and thus is consistent for  $n$  sufficiently large.
3.  $L_n \leq \epsilon'$  with probability at least  $1 - p_{n,\epsilon} - 2S2^{-n\epsilon'/2}$  for any  $\epsilon, \epsilon' > 0$ , where  $S$  is the number of axioms sets with description length bounded by  $DL(A^*) + \epsilon$ .

**Proof :** Note  $A$  the axiom set including  $A^*$  and all the  $e_i$  such that  $A^* \not\models_{\delta_n} e_i$ . Define  $\mathbf{1}_e = 1 \iff A^* \not\models_{\delta_n} e$  and  $\mathbf{1}_e = 0$  otherwise.

**Step 0 :**  $DL(A) \leq DL(A^*) + \sum_{i \leq n; A^* \not\models_{\delta_n} e_i} DL(e_i)$  by definition.

**Step 1 :** upper-bounding  $\mu_A = \mathbb{E}[DL(A)]$ .

$$\mu_A \leq DL(A^*) + n\mathbb{E}[DL(e)\mathbf{1}(e)] \leq DL(A^*) + n\mu_e P_{(\delta_n)} + n\sigma_e \sqrt{P_{(\delta_n)}}$$

thanks to  $\mathbb{E}fg \leq \mathbb{E}f\mathbb{E}g + Cov(f, g)$ , with  $Cov(DL(e), \mathbf{1}(e)) \leq \sigma(DL(e))\sigma(\mathbf{1}(e))$ , and  $\sigma(\mathbf{1}(e)) \leq \sqrt{P_{(\delta_n)}}$ .

**Step 2 :** upper-bounding  $\sigma_A = \sigma(DL(A))$ .

$$\begin{aligned} \sigma_A^2 &\leq n\sigma(DL(e)\mathbf{1}(e))^2 \leq n\mathbb{E}[DL(e)^2\mathbf{1}(e)^2] \\ &\leq n\mathbb{E}[DL(e)^2]\mathbb{E}[\mathbf{1}(e)^2] + n\sigma(DL(e)^2)\sigma(\mathbf{1}(e)^2) \end{aligned}$$

thanks to  $\mathbb{E}fg \leq \mathbb{E}f\mathbb{E}g + Cov(f, g)$ , applied with  $f = DL(e)^2$  and  $g = \mathbf{1}(e)^2$ . Therefore as  $\mathbf{1}(e)^2 = \mathbf{1}(e)$ ,

$$\begin{aligned} \sigma_A^2 &\leq n\mathbb{E}(DL(e)^2)\mathbb{E}[\mathbf{1}(e)] + n\sigma(DL(e)^2)\sigma(\mathbf{1}(e)) \\ &\leq n\mathbb{E}(DL(e)^2)P_{(\delta_n)} + n\sigma(DL(e)^2)\sqrt{P_{(\delta_n)}} \end{aligned}$$

**Step 3 :** upper-bounding  $P_{(\delta_n)}$ .

Thanks to the inequality of Tchebychev-Cantelli ( $P(X > \mathbb{E}X + t) \leq \frac{\sigma(X)^2}{\sigma(X)^2 + t^2}$ ) applied to  $X$  the length of the shortest proof of  $e$ , if the two first moments of the proof length of  $e$  are finite, then  $P_{(\delta_n)} = O(1/\delta_n^2)$ . In particular,  $\delta_n = \Omega(n^3)$  implies  $P_{(\delta_n)} = O(1/n^6)$ .

**Step 4 :** upper-bounding  $P(DL(A) \geq DL(A^*) + \epsilon)$ .

Consider any  $\epsilon > 0$ . Thanks to the inequality of Tchebychev-Cantelli, for  $t > 0$ ,

$$P(DL(A) \geq \mu_A + t) \leq \frac{1}{1 + \frac{t^2}{\sigma_A^2}}$$

Define  $t = DL(A^*) + \epsilon - \mu_A$ . Thanks to step 1 and 3, if  $\delta_n = \Omega(n^3)$ , then  $\mu_A = DL(A^*) + O(1/n^2)$ . Therefore,  $t > 0$  if  $n$  is sufficiently large, and  $t = \epsilon - O(1/n^2)$ . Then,  $P(DL(A) \geq DL(A^*) + \epsilon) = P(DL(A) \geq \mu_A + t)$ . Thanks to Tchebychev-Cantelli,

$$P(DL(A) \geq \mu_A + t) \leq \frac{\sigma_A^2}{\sigma_A^2 + t^2}$$

$$\text{So, } P(DL(A) \geq DL(A^*) + \epsilon) \leq \frac{\sigma_A^2}{\sigma_A^2 + t^2}$$

Thanks to step 2 and step 3,  $\sigma_A^2 = O(1/n^2)$ . Then,  $P(DL(A) \geq DL(A^*) + \epsilon) \leq O(1/n^2)$ . More precisely, we have proved that  $P(DL(A) \geq DL(A^*) + \epsilon) \leq p_{n,\epsilon}$ , with  $p_{n,\epsilon} = \frac{\sigma_A^2}{\sigma_A^2 + (DL(A^*) + \epsilon - \mu_A)^2}$ .

**Step 5 :** applying the Borell-Cantelli lemma to prove  $\limsup DL(A) \leq DL(A^*)$ .

We have claimed  $P(DL(A) \geq DL(A^*) + \epsilon) \leq O(1/n^2)$ . This implies that for any  $\epsilon > 0$ , for  $n$  sufficiently large,  $\sum_n P(DL(A) \geq DL(A^*) + \epsilon) < \infty$  (as  $P(DL(A) \geq DL(A^*) + \epsilon) \leq 1/n^2$ ). The Borell-Cantelli lemma precisely states that this implies that  $\limsup DL(A) \leq DL(A^*)$ .

**Step 6 :** Concluding on the convergence of  $A_n$ .

For  $n$  sufficiently large,  $DL(A) \leq DL(A^*)$  ( $DL(\cdot)$  is discrete, so this is a direct consequence of step 5). By definition of  $A_n$ ,  $DL(A_n) \leq DL(A)$ .

Therefore,  $DL(A_n) \leq DL(A^*)$  for  $n$  sufficiently large. As  $A_n$  lives in a finite space (the space of axiom sets with description length bounded by  $DL(A^*)$ ), and as any axiom set shorter than  $A_n$  is excluded for  $n$  sufficiently large<sup>8</sup>  $A_n$  reaches  $A^*$  almost surely.

**Step 7 :** Concluding on the convergence of  $L_n$ .

By definition  $\mathfrak{E}_n$  is included in  $T_{n+1}$  and as  $A_n \not\models_{2\delta_n} e \wedge \neg e$ ,  $T_{n+1}$  does not contain any  $\neg e$  for  $e \in \mathfrak{E}_n$ . Thus, our algorithm is exact on all the examples in the sense that all examples in  $\mathfrak{E}_n$  are well classified. Moreover, we have proved in step 4 that  $P(DL(A_n) \geq DL(A^*) + \epsilon)$  is upper-bounded by  $p_{n,\epsilon}$ . So, with probability at least  $1 - p_{n,\epsilon}$ ,  $A_n$  lies in a family of cardinal bounded by  $S$ , where  $S$  is the number of axioms sets with description length bounded by  $DL(A^*) + \epsilon$ . Therefore (see e.g. (Devroye *et al.*, 1997, chap 12.7)), for any  $\epsilon$  and  $\epsilon'$ ,  $L_n \leq \epsilon'$  with probability at least  $1 - p_{n,\epsilon} - 2S2^{-n\epsilon'/2}$ .  $\square$

<sup>8</sup>If an axiom set  $C$  verifies  $DL(C) < DL(A^*)$ ,

- either it is not consistent and thus  $A_n \neq C$  if  $n$  is larger than some  $n(C)$  such that  $C \vdash_{\delta_{n(C)}} \perp$ ;
- or there is  $e$  such that  $M(e) > 0$  and  $C \not\models e$  and thus  $A_n \neq C$  if  $n \geq n(C)$  where  $n(C) = \min\{i; e_i = e\}$ .

So,  $DL(A_n) \leq DL(A^*)$  implies that  $A_n = A^*$  if  $n \geq \sup_{C; DL(C) < DL(A^*)} n(C)$ .



## Références

- AMIR E. & MCILRAITH S. (2003). Partition-based logical reasoning for first-order and propositional theories. *Artificial intelligence*, **162**(1,2), 49 :88.
- CHEESEMAN P., KANEFSKY B. & TAYLOR W. (1991). Where the really hard problems are. In *proceedings of IJCAI91*, p. 331–337.
- DEVROYE L., GYÖRFI L. & LUGOSI G. (1997). *A probabilistic Theory of Pattern Recognition*. Springer.
- FREILING C. (1986). Axioms of symmetry : Throwing darts at the real number line. *Journal of Symbolic Logic*, **51** (1), 190–200.
- HAMKINS J. (2002). Infinite time turing machines. *Minds and Machines (special issue on hypercomputation)*, **12**(4), 521–539.
- J.D. N. (2003). Must evidence under-determine theory ? In *First Notre Dame-Bielefeld Interdisciplinary Conference on Science and Values, Zentrum für Interdisziplinäre Forschung, Universität Bielefeld*.
- KHARDON R. & ROTH D. (1997). Learning to reason. *Journal of the ACM*, **44**(5), 697–725.
- KLEENE S. (1967 (Dover 2002)). *Mathematical Logic*. J. Wiley (reprint :Dover).
- LIST C. (1999). Craig’s theorem and the empirical underdetermination thesis reassessed. *Disputatio* 7, p. 28–39.
- SHAPIRO E. Y. (1981). Inductive inference of theories from facts. *Research Report* 192.
- SHOOK J. (2002). Dewey and quine on the logic of what there is. In D. M. H. TOM BURKE & R. TALISSE, Eds., *Dewey’s Logical Theory : New Studies and Interpretations* : Vanderbilt University Press.
- TARSKI A. (1949). On essential undecidability. *Journal of Symbolic Logic*, **14**, 75–76.
- VALIANT L. (1984). A theory of the learnable. *Communication of the ACM*, **27**, 1134–1142.
- VAPNIK V. & CHERVONENKIS A. (1974). *Theory of Pattern Recognition*. Nauka, Moskow. (in Russian).
- VIDYASAGAR M. (1997). *A Theory of Learning and Generalization, with Applications to Neural Networks and Control Systems*. Springer-Verlag.
- WOODIN W. H. (2001). The continuum hypothesis, i & ii. *Notices Amer. Math. Soc.*, p. 48–6, 567–576 & 8–7, 681–690.