

# The Complexity of Theory Revision

Russell Greiner\*

Siemens Corporate Research Princeton, NJ 08540-6632

greiner@scr Siemens com

## Abstract

A knowledge-based system uses its database (aka its "theory") to produce answers to the queries it receives. Unfortunately, these answers may be incorrect if the underlying theory is faulty. Standard "theory revision" systems use a given set of "labeled queries" (each a query paired with its correct answer) to transform the given theory, by adding and/or deleting either rules and/or antecedents, into a related theory that is as accurate as possible. After formally defining the theory revision task and bounding its sample complexity, this paper addresses the task's computational complexity. It first proves that, unless  $P = NP$ , no polynomial time algorithm can identify the optimal theory, even given the exact distribution of queries, except in the most trivial of situations. It also shows that, except in such trivial situations, no polynomial-time algorithm can produce a theory whose inaccuracy is even close (i.e., within a particular polynomial factor) to optimal. These results justify the standard practice of hill-climbing to a locally-optimal theory, based on a given set of labeled samples

## 1 Introduction

There are many fielded knowledge-based systems, ranging from expert systems and logic programs to production systems and database management systems [Lev84]. Each such system uses its database of general information (aka its "theory") to produce an answer to each given query, this can correspond to retrieving information from a database (eg, finding  $X$  such that "(make(X) & (color X red))" or to providing a diagnosis or repair, based on a given set of symptoms. Unfortunately, these responses may be incorrect if the underlying theory includes erroneous information. If we observe that some answers are incorrect (eg, if the patient does

"I gratefully acknowledge receiving helpful comments from George Drastal, R. Bharat Rao, Narendra Gupta, Tom Hancock, Sheila McIlraith, Edoardo Amaldi, Dan Roth and Rom Khardon

not get better, or the proposed repair does not correct the device's faults), we can then ask a human expert to supply the correct answer. *Theory revision* is the process of using such correctly-answered queries to modify the initial theory, to produce a new theory that is more accurate, i.e., which will not make those mistakes again, while remaining correct on the other queries.

Most theory revision algorithms use a set of transformations to hill-climb through successive theories, until reaching a theory whose accuracy is (locally) optimal, based on a set of correctly-answered queries, *c/*, [Pol85, MB88, Coh90, OM94, WP93, CS90, LDRG94]. This report addresses the obvious question: Is there a better approach, which will directly yield the *globally* optimal theory?<sup>7</sup>

Section 2 first states the theory revision objective more precisely: as finding the theory with the highest accuracy from the space of theories formed by applying a sequence of transformations to a given initial theory, here each transform involves either adding or deleting either a rule or an antecedent. It also proves that a polynomial number of training "labeled queries" (each a specific query paired with its correct answer) is sufficient, i.e., they provide the information needed to identify a transformation-sequence that will transform the given theory into a new theory whose accuracy is arbitrarily close to optimal, with arbitrarily high probability. Section 3 then addresses the *computational* complexity of the task of finding the optimal (or even near-optimal) revised theory. It first proves that the task of computing the optimal theory within this space of theories is intractable, even in trivial contexts — eg, even when dealing with propositional Horn theories, or when considering with only atomic queries, or when considering only a bounded number of transformations, etc.<sup>1</sup> We then show that this task cannot even be approximated, i.e., that no efficient algorithm can find a theory whose inaccuracy is even close to (i.e., within a particular small polynomial of) optimum'. We also prove that these negative results apply even when we are only generalizing, or only specializing, the initial theory. We also discuss the efficiency of other restricted variants of theory revision,

throughout, we will assume that  $P \neq NP$  [GJ79], which implies that any NP-hard problem is intractable. This also implies certain approximation claims, presented below

providing sharp boundaries that describe exactly when this task is, versus is not, tractable

We view these results as sanctioning the standard approach of using a set of transformations to hill-climb to a local optimum, based on a set of samples. The labeled training samples are required to obtain the needed distribution information, and the realization that no tractable algorithm will be able to find the *global* optimum justifies hill-climbing to a local optimum, within the space formed using specified transformations

We close this section by describing related research<sup>2</sup>

**Related Research** While most learning systems begin with an "empty theory" and attempt to learn a target function (perhaps a decision tree, or a logic program), theory revision processes work by modifying a given initial theory. There are several implemented theory revision systems. Most use essentially the same set of transformations we describe — e.g., AUDREY [WP93], FONTE [MB88], EITHER [OM94] and A [LDRG94] all consider adding or deleting antecedents or rules. Our analysis, and results, can easily be applied to many other types of modifications — e.g., specializing or generalizing antecedents [OM94], using "n-of-m rules" [BM93], or merging rules and removing chains paths of rules that produced incorrect results [Coh90, Coh92]<sup>3</sup>. While these projects provide empirical evidence of the effectiveness of their specific algorithms, and deal with classification (i.e., determining whether a given element is a member of some target class) rather than general derivation, our work formally addresses the complexities inherent in finding the best theory, for handling arbitrary queries.

Finally, note that, in some cases, our task can require extracting the best consistent sub-theory from a given inconsistent theory. From this perspective, our work is related to "Knowledge Representation" form of theory revision, a la Gardenfors [Gar88, AGM85], Katsuno and Mendelzon [KM91] and many others. Our work differs by using the notion of expected accuracy to dictate which of the "revisions" is best.

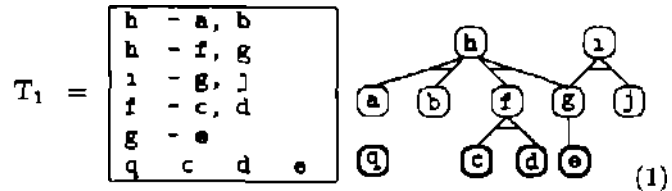
## 2 Framework

We define a "theory" as a collection of (propositional or first order) Horn clauses, where each clause is a disjunction of literals, at most one of which is positive. Borrowing from [Lev84, DP91], we also view a theory T as a function that maps each query to its proposed answer, hence,  $T: Q \rightarrow A$ , where Q is a (possibly infinite) set of Horn queries, and  $A = \{ \text{No}, \text{Yes} \}$  is the set of possible

<sup>2</sup>The technical report [Gre95b] provides a more extensive literature survey, as well as proofs of the theorems.

<sup>3</sup>The companion paper [Gre95a] considers yet other ways of modifying a theory, by rearranging its component rules or antecedents.

answers<sup>4</sup>. Hence, given



$T_1(h) = \text{Yes}$ ,  $T_1(i) = \text{No}$  and  $T_1(i - e, j) = \text{Yes}$ . Of course, different theories can return different answers to a given query. For example, let  $T_2$  be a theory that differs from  $T_1$  only by excluding the "g ← e" rule, then  $T_2(h) = \text{No}$ .

While the non atomic queries may seem unusual at first, they are actually quite common. For example, a medical expert system typically collects relevant data  $\{f_1(p), \dots, f_n(p)\}$  about an individual patient p, then determines whether p has some specific disease  $\text{disease}_i$ , i.e., if  $T \cup \{f_1(p), \dots, f_n(p)\} \models \text{disease}_i(p)$ , where T is the expert system's initial theory that contains general information about diseases, etc. Notice this entailment condition holds iff  $T \models \neg f_1(p) \vee \dots \vee \neg f_n(p) \vee \text{disease}_i(p)$ , i.e., iff the Horn query "disease<sub>i</sub>(p) ← f<sub>1</sub>(p), ..., f<sub>n</sub>(p)" follows from the initial theory (See also "entailment queries" [FP93]).

We assume there is a single correct answer to each question, and represent it using the real-world oracle  $O: Q \rightarrow A$ . Here, perhaps,  $O(h) = \text{No}$ , meaning that "h" should not hold. We say an oracle is consistent if it produces the same answers as a Horn theory, over the set of queries Q. Note, we will sometimes deal with inconsistent oracles, e.g., where  $O(a) = \text{Yes}$ ,  $O(b - a) = \text{Yes}$ , and  $O(b) = \text{No}$ .

Our goal is to find a theory that is as close to  $O(\cdot)$  as possible. To quantify this, we first define the "accuracy function"  $a(\cdot, \cdot)$  where  $a(T, q)$  is the accuracy of the answer the theory T returned for the query q.

$$a(T, q) \stackrel{\text{def}}{=} \begin{cases} 1 & \text{if } T(q) = O(q) \\ 0 & \text{otherwise} \end{cases}$$

(Notice  $a(T, \cdot)$  implicitly depends on the oracle  $O(\cdot)$ ). Hence, as  $O(h) = \text{No}$ ,  $a(T_2, "h") = 1$  as  $T_2$  provides the correct answer while  $a(T_1, "h") = 0$  as  $T_1$  returns the wrong answer.

This  $a(T, \cdot)$  function measures T's accuracy for a single query. In general, our theories must deal with a range of queries. We model this using a stationary, but unknown, probability function  $Pr: Q \rightarrow [0, 1]$ , "where  $Pr(Q)$  is the probability that the query q will be posed. Given this distribution, we can compute the "expected accuracy" of a theory, T.

$$A(T) = E[a(T, q)] = \sum_{q \in Q} Pr(q) \times a(T, q)$$

We will consider various sets of possible theories,  $\mathcal{T} = \{T_i\}$ , where each such T contains the set of theories

<sup>4</sup>To simplify the presentation, the bulk of this paper will deal only with propositional logic, Section 2.3 below describes the extensions needed to deal with predicate calculus.

formed by applying various sequences of transformations to a given initial theory, see Section 2.1 below. Our challenge is to identify the theory  $T_{opt} \in \mathcal{T}$  whose expected accuracy is optimal, i.e.,

$$\forall T \in \mathcal{T} \quad A(T_{opt}) \geq A(T) \quad (2)$$

There are two challenges to finding such optimal theories. The first is based on the observation that the expected accuracy of a theory depends on the distribution of queries, which means different theories will be optimal for different distributions. While this distribution is not known initially, it can be estimated by observing a set of samples (each a query/answer pair), drawn from that distribution. Section 2.2 below discusses the number of samples required to be confident of obtaining the information needed to identify a good  $T^* \in \mathcal{T}$ , with high probability.

We are then left with the challenge of computing the best theory, once given this distributional estimate. Section 3 addresses the computational complexity of this process, showing that the task is not just intractable,<sup>5</sup> but it is also not approximatable — i.e., no efficient algorithm can even find a theory whose expected accuracy is even close (in a sense defined below) to the optimal value.

We first close this section by describing the transformations we will use to define the various spaces of theories, then discussing the sample complexity of the implied learning process and finally providing the extensions needed to deal with predicate calculus.

## 2.1 Standard Transformations

Standard theory revision algorithms implicitly explore the space of possible theories  $\Sigma^\infty(T_0) = \{\sigma(T_0) \mid \sigma \in \Sigma^\infty\}$ , which contains the theories formed by applying some sequence of theory-to-theory transformations  $\sigma \in \Sigma^\infty$  to the given initial theory  $T_0$ . Each  $\sigma = \tau_1 \circ \tau_2 \circ \dots \circ \tau_k \in \Sigma^\infty$  sequence is formed from  $\Sigma = \Sigma_{DR} \cup \Sigma_{AR} \cup \Sigma_{DA} \cup \Sigma_{AA}$ , where each  $\tau_{DR} \in \Sigma_{DR}$  deletes a rule from the theory, each  $\tau_{AR} \in \Sigma_{AR}$  adds a new rule to the theory, each  $\tau_{DA} \in \Sigma_{DA}$  deletes an existing antecedent from an existing rule, and each  $\tau_{AA} \in \Sigma_{AA}$  adds a new antecedent to an existing rule. In some situations, we will consider “ $K$ -bounded sequences”

$$\Sigma^K = \{\sigma = \tau_1 \circ \tau_2 \circ \dots \circ \tau_k \mid \tau_i \in \Sigma \ \& \ c(\sigma) \leq K\}$$

whose members  $\sigma = \tau_1 \circ \tau_2 \circ \dots \circ \tau_k \in \Sigma^K$  are sequences of transformations whose total cost  $c(\sigma) = c(\tau_1) + c(\tau_2) + \dots + c(\tau_k)$  are at most  $K$ , where the cost  $c(\tau)$  of the transformation  $\tau$  is the number of symbols added to, or deleted from,  $T$  to form  $\tau(T)$ . In the propositional case,  $c(\tau_{AA}) = c(\tau_{DA}) = 1$  for each transformation that either adds or deletes an antecedent, and

<sup>5</sup>All  $a(T, q)$  requires computing  $T(q)$ , which can require proving an arbitrary theorem, this computation alone can be computationally intractable, if not undecidable. Our results show that the task of finding the optimal theory is intractable even given a *poly time oracle for these arbitrary derivations*. Of course, as we are considering only Horn theories, these computations are guaranteed to be poly time in the propositional case [BCH90].

$c(\tau_{AR}) = c(\tau_{DR}) = |\rho|$  for each add-rule (resp., delete-rule) transformation that adds (resp., deletes) the rule  $\rho$ , which has 1 conclusion and  $|\rho| - 1$  antecedent literals.

As an example, applying the 3-element sequence  $\sigma = \tau_{g-e+q}^{AA} \circ \tau_{b-f}^{AR} \circ \tau_{f-c,d,-c}^{DA}$  with total cost  $c(\sigma) = c(\tau_{g-e+q}^{AA}) + c(\tau_{b-f}^{AR}) + c(\tau_{f-c,d,-c}^{DA}) = 1 + 2 + 1 = 4$ , will transform  $T_1$  into  $\sigma(T_1) = \tau_{g-e+q}^{AA}(\tau_{b-f}^{AR}(\tau_{f-c,d,-c}^{DA}(T_1)))$  which is a theory with 8 clauses that differs from  $T_1$  by including the clause “ $g - e, q$ ” rather than “ $g - a$ ”, including the clause “ $f - d$ ” rather than “ $f - c, d$ ”, and by including an extra clause “ $b - f$ ”.

Finally, we will also consider various other restricted spaces of transformation-sequences, which are formed from specified types of transformations, e.g.,

$$\begin{aligned} \Sigma^{-R+A} &= \{\sigma = \tau_1 \circ \tau_2 \circ \dots \circ \tau_k \mid \tau_i \in \Sigma_{DR} \cup \Sigma_{AA}\} \\ \Sigma^{+R-A} &= \{\sigma = \tau_1 \circ \tau_2 \circ \dots \circ \tau_k \mid \tau_i \in \Sigma_{AR} \cup \Sigma_{DA}\} \end{aligned}$$

correspond to (unbounded) transformation sequences that produce more specific (resp., more general) theories, as well as the bounded variants, e.g.,  $\Sigma^{-R+A}(K) = \{\sigma = \tau_1 \circ \dots \circ \tau_k \mid \tau_i \in \Sigma_{DR} \cup \Sigma_{AA} \ \& \ c(\sigma) \leq K\}$ . Note that the earlier  $\Sigma^\infty = \Sigma^{+R-R+A-A}$  and  $\Sigma^K = \Sigma^{+R,-R+A-A}(K)$ .

## 2.2 Sample Complexity

We can use following standard Computational Learning Theory theorem to bound the number of samples required to obtain the information needed to identify a good  $T^* \in \mathcal{T}$  with high probability, showing in particular how this depends on the space of theories  $\mathcal{T}$  being considered.

**Theorem 1 (from [Vap82, Theorem 6.2])** *Given a class of theories  $\mathcal{T}$  and  $\epsilon, \delta > 0$ , let  $T_* \in \mathcal{T}$  be the theory with the largest empirical accuracy after*

$$M_{upper}(\mathcal{T}, \epsilon, \delta) = \left\lceil \frac{2}{\epsilon^2} \ln \left( \frac{|\mathcal{T}|}{\delta} \right) \right\rceil$$

*samples (each a labeled query), drawn from the stationary distribution,  $Pr(\cdot)$ . Then, with probability at least  $1 - \delta$ , the expected accuracy of  $T_*$  will be within  $\epsilon$  of the optimal theory in  $\mathcal{T}$ , i.e.,  $Pr[A(T_*) \geq A(T_{opt}) - \epsilon] \geq 1 - \delta$ , using the  $T_{opt}$  from Equation 2.*

This means a polynomial number of samples is sufficient to identify an  $\epsilon$ -good theory from  $\mathcal{T}$  with probability at least  $1 - \delta$ , whenever  $\ln(|\mathcal{T}|)$  is polynomial in the relevant parameters. Notice this is true for most of the classes of theories being considered, e.g., as  $\Sigma^{-R}(T)$  is the power-set of the rules in  $T$ ,  $|\Sigma^{-R}(T)| = 2^{|\text{Rules}(T)|}$ , which means  $|\ln(\Sigma^{-R}(T))| = |\text{Rules}(T)| = O(|T|)$ , which clearly is polynomial in the size of the initial theory. This “ $\ln(|\mathcal{T}|) = \text{poly}(|T|)$ ” claim is slightly problematic for transformations that can add symbols, notably for  $\Sigma^{+R}$  and  $\Sigma^{+A}$ . But even here, the sample complexity remains polynomial in the size of the revised theory, which effectively means again that sample-efficient learning remains possible, cf., “nonuniform” pac-learning [B188].

### 2.3 Dealing with Predicate Calculus

To handle predicate calculus expressions, we have to consider answers of the form  $\{\text{Yes}[\{X_i/v_i\}]\}$ , where the expression within each  $\text{Yes}[\ ]$  is a binding list of the free variables, corresponding to a single answer to the query. For example, given the theory<sup>6</sup>

$$T_{pc} = \left\{ \begin{array}{ll} \text{tall}(\text{john}) & \text{short}(\text{fred}) \\ \text{rich}(\text{john}) & \text{rich}(\text{fred}) \\ \text{eligible}(X) & - \text{tall}(X), \text{rich}(X) \end{array} \right\}$$

the query  $\text{short}(Y)$  will return  $T_{pc}(\text{short}(Y)) = \{\text{Yes}[\{Y/\text{fred}\}]\}$ , the query  $\text{rich}(Z)$  will return the pair of answers  $T_{pc}(\text{rich}(Z)) = \{\text{Yes}[\{Z/\text{john}\}], \text{Yes}[\{Z/\text{fred}\}]\}$ , and  $T_{pc}(\text{eligible}(A)) = \{\text{Yes}[\{A/\text{john}\}]\}$ . As  $\mathcal{O}(\ )$  and  $\mathcal{T}(\ )$  each returns a set of answers to each query, we therefore define  $T$ 's accuracy score as  $a(T, \sigma) = \frac{|\mathcal{O}(\sigma) \cap \mathcal{T}(\sigma)|}{|\mathcal{O}(\sigma) \cup \mathcal{T}(\sigma)|} \in [0, 1]$ . All of the theorems in this paper hold even when considering only non-recursive Datalog (i.e., "function free") theories.

### 3 Computational Complexity

Our basic challenge is to produce a theory  $T_{opt}$  whose accuracy is as large as possible. The previous section supplied the number of samples needed to guarantee, with high probability, that the expected accuracy of the theory whose empirical accuracy is largest,  $T_{\epsilon}$ , will be within  $\epsilon$  of the expected accuracy of this  $T_{opt}$ . This section discusses the computational challenge of determining this  $T_{\epsilon}$ , given this distributional estimate. We show first that this task is tractable in degenerate trivial situations when considering (1) only atomic queries posed to a (2) propositional theory and being allowed (3) an arbitrarily large number of modifications to the initial theory, to produce (4) a perfect theory (i.e. one that returns the correct answer to every query). This task becomes intractable, however, if we remove (essentially) any of these restrictions — e.g., if we seek optimal (rather than only seeking "perfect") propositions! theories and are allowed to pose Horn queries, or if we consider predicate calculus theories. It also remains intractable even if we restrict the number of modifications allowed, which implies that the task of determining the smallest number of modifications required to find a perfect theory is intractable. We next show that these tasks are not just intractable but worse, they are not even approximatable, except in the most trivial of situations.

We also consider two special subtasks by restricting the allowed types of transformations, to consider revision processes that only specialize (respectively, only generalize) the initial theory. We show that these tasks, also, are intractable and non-approximatable in essentially all situations, i.e., except when all four of the above conditions hold<sup>7</sup>. Figures 1 and 2 summarize the various cases.

<sup>6</sup>Following PROLOG'S conventions, we will capitalize each variable, as in the "X" above.

<sup>7</sup>Actually, there is one other tractable case in the generalization situation, see Figure 1.

### 3.1 Basic Complexity Results

To formally state the problem

**Definition 1 (TR[ $\Sigma^{\dagger}$ ] Decision Problem)**

- INSTANCE**
- Initial theory  $T$ ,
  - Labeled training sample  $S = \{(q_i, \mathcal{O}(q_i))\}$  containing a set of Horn queries and the correct answers, and
  - Probability value  $p \in [0, 1]$

**QUESTION** Is there a theory  $T' \in \Sigma^{\dagger}[T]$  such that  $A_S(T') = \frac{1}{|S|} \sum_{(q, \mathcal{O}(q)) \in S} a(T', q) \geq p$ ?

The  $\Sigma^{\dagger}[\ ]$  function maps a theory to a set of candidate revised theories, here, we will consider various  $\Sigma^{\pm R \pm A}$  transformation sets. To simplify our notation, we will write  $A(T)$  for  $A_S(T)$ .

We will also consider the following special cases.  $\text{TR}_{\text{Prrf}}[\Sigma^{\dagger}]$  requires that  $p = 1$  (i.e., seeking perfect theories, rather than "optimal" theories  $\text{TR}_{(\text{Opt})}[\Sigma^{\dagger}]$ ),  $\text{TR}_{\text{Prop}}[\Sigma^{\dagger}]$  deals with propositional logic (rather than predicate calculus,  $\text{TR}_{(\text{PredCal})}[\Sigma^{\dagger}]$ ), and  $\text{TR}_{\text{Atom}}[\Sigma^{\dagger}]$  deals with only atomic queries (as opposed to Horn queries,  $\text{TR}_{(\text{Horn})}[\Sigma^{\dagger}]$ ). We will also use  $\text{TR}_{\text{Disj}}[\Sigma^{\dagger}]$  to refer to the task when the queries can be arbitrary disjunctions, which need not be Horn. (While the other subscripts are restrictions on  $\text{TR}[\Sigma^{\dagger}]$ , this *Disj* case is more permissive.) We will also combine subscripts, with the obvious meanings. When  $\text{TR}_X[\Sigma^{\dagger}]$  is a special case of  $\text{TR}_Y[\Sigma^{\dagger}]$ , finding that  $\text{TR}_X[\Sigma^{\dagger}]$  is hard immediately implies that  $\text{TR}_Y[\Sigma^{\dagger}]$  is hard. Similarly, seeing that  $\text{TR}_Y[\Sigma^{\dagger}]$  is easy immediately implies that each special case of  $\text{TR}_Y[\Sigma^{\dagger}]$  is easy. As a final note: all of the hardness results presented in this paper hold even if we only consider "3-CNF Horn theories" — i.e., rules whose antecedents contain at most 2 literals.

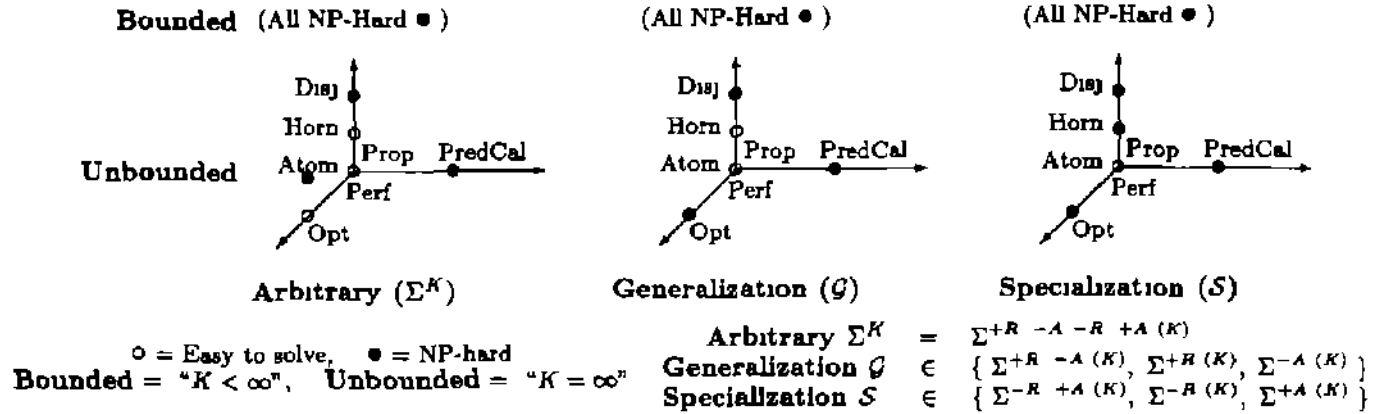
Here, it is easy to find the optimal theory in certain degenerate cases, where either the individual queries can be decoupled (e.g., when using atomic propositional queries) or when our actions are forced (e.g., when seeking perfect propositional theories) — just throw away the original theory, then add in propositions corresponding to the "Yes-labeled queries". In every other case, however, the task is intractable.

**Theorem 2** The  $\text{TR}_{\text{Prop Atom (Opt)}}[\Sigma^{\infty}]$  and  $\text{TR}_{\text{Prop (Horn) Perf}}[\Sigma^{\infty}]$  decision problems are easy, each other problem is NP-hard.

(This information is summarized in lower left "Unbounded, Arbitrary" graph of Figure 1.)

The above theorem describes the complexity of computing the best theory when we are allowed to use an arbitrarily expensive sequence of transformations. We get an even stronger negative results if we restrict the "expense" of the transformation sequence.

**Theorem 3** For some  $K \in \mathcal{Z}^+$ , the  $\text{TR}_{\text{Prop Atom Perf}}[\Sigma^K]$  decision problem is NP-hard. This is true even if we consider only labeled queries produced by a consistent oracle (i.e.,



Any task that "projects" down to an NP-hard task, along any axis, is NP-hard. Here, this means all of the "cross terms" are NP-hard. (For example  $\text{ThRev}_{\text{PredCal Horn Perf}}[\Sigma^\infty]$  is NP-hard, as its projection to the "Prop-PredCal  $\times$  Perf-Opt" plane,  $\text{ThRev}_{\text{PredCal Atom Perf}}[\Sigma^\infty]$  is NP-hard.) The  $\text{ThRev}_{\text{Prop Horn Opt}}[\Sigma^\infty]$  case is shown explicitly as each of its projections is easy, the figures omit all other cross-terms.

Figure 1 Tractability of Theory Revision Tasks

even when there is a Horn theory that correctly labels all of the queries)

The observation that determining such "it-step perfect theories" is NP-hard leads immediately to

**Corollary 3.1** *It is NP-hard to compute the minimal-cost transformation sequence required to produce a perfect theory (i.e., to compute the smallest  $K$  for which there is a  $T_{\text{perfect}} \in \Sigma^K[T]$  such that  $A(T_{\text{perfect}}) = 1$ ), even in the propositional case when considering only atomic queries. It is also NP-hard to compute the "minimal-length" transformation, where the length of the transformation sequence  $\tau_1 \circ \tau_2 \circ \dots \circ \tau_k$  is simply  $k$  — i.e., when each transformation has "unit cost"*

This negative result shows the intractability of the obvious proposal of using a breath-first transversal of the space of all possible theory revisions. First test the initial theory  $T_0$  against the labeled queries, and return  $T_0$  if it is 100% correct. If not, then consider all theories formed by applying a single (unit-cost) transformation, and return any perfect  $T_1 \in \Sigma^1[T_0]$ , and if not, consider all theories in  $\Sigma^2[T_0]$  (formed by applying sequences of transformations with cost at most two), and return any perfect  $T_2 \in \Sigma^2[T_0]$ , and so forth.

### 3.2 Approximability

Many decision problems correspond immediately to optimization problems, for example, the MIN GRAPH COLOR decision problem (given a graph  $G = (N, E)$  and a positive integer  $A$ , can each node be labeled by one of  $K$  colors in such a way that no edge connects two nodes of the same color, see [GJ79, pl91(Chromatic Number)]) corresponds to the obvious minimization problem. Find the minimal coloring of the given graph  $G$ . We can similarly view the  $\text{TR}_\chi[\Sigma^1]$  decision problem as either the maximization problem "Find the  $T' \in \Sigma^1[T]$  whose accuracy is maximal" or the minimization problem "Find the  $T' \in \Sigma^1[T]$  whose inaccuracy is minimal", where a theory's inaccuracy is obviously  $\text{INA}(T) = 1 - A(T)$ .

(While the maximally accurate theory is also minimally inaccurate, these two formulations can lead to different approximability results.) For notation, let " $\text{MAXTR}_\chi[\Sigma^1]$ " (resp., " $\text{MINTR}_\chi[\Sigma^1]$ ") refer to the maximization (resp., minimization) problem.

Now consider any algorithm  $B$  that, given any  $\text{MINTR}_\chi[\Sigma^1]$  instance  $x = (T, S)$  with initial theory  $T$  and labeled training sample  $S$ , computes a syntactically legal, but not necessarily optimal, revision  $B(x) \in \Sigma^1[T]$ . Then  $B$ 's "performance ratio for the instance  $x$ " is defined as

$$\text{MinPerf}_{\text{MINTR}_\chi[\Sigma^1]}(B, x) = \begin{cases} \frac{\text{INA}(B(x))}{\text{INA}(\text{opt}(x))} & \text{if } \text{INA}(\text{opt}(x)) \neq 0 \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

where  $\text{opt}(x) = \text{opt}_{\text{MINTR}_\chi[\Sigma^1]}(x)$  is the optimal solution for this instance, i.e.,  $\text{opt}(x)$  is the theory  $T_{\text{opt}} \in \Sigma^1[T]$  with minimal inaccuracy over  $S$ .

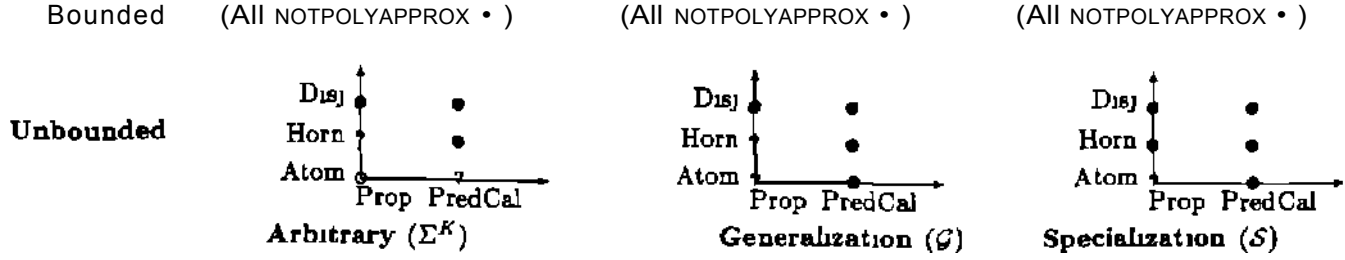
We say a function  $g(\cdot)$  "bounds  $B$ 's performance ratio (over  $\text{MINTR}_\chi[\Sigma^1]$ )" iff

$$\forall x \in \text{MINTR}_\chi[\Sigma^1], \text{MinPerf}_{\text{MINTR}_\chi[\Sigma^1]}(B, x) \leq g(|x|)$$

where  $|x|$  is the size of the instance  $x = (T, S)$ , which we define to be the number of symbols in  $T$  plus the number of symbols used in  $S$ . Intuitively, this  $g(\cdot)$  function indicates how closely the  $B$  algorithm comes to returning the best answer for  $x$ , over all  $\text{MINTR}_\chi[\Sigma^1]$  instances  $x$ .

Now let  $\text{Poly}(\text{MINTR}_\chi[\Sigma^1])$  be the collection of all polytime algorithms that return legal answers to  $\text{MINTR}_\chi[\Sigma^1]$  instances. It is natural to ask for the algorithm in  $\text{Poly}(\text{MINTR}_\chi[\Sigma^1])$  with the best performance ratio, this would indicate how close we can come to the optimal solution, using only a feasible computational time. For example, if this function was the constant 1 for  $\text{MINTR}_{\text{Prop}}[\Sigma^\infty]$ , then a poly-time algorithm could produce the optimal solution to any  $\text{MINTR}_{\text{Prop}}[\Sigma^\infty]$  instance, as  $\text{TR}_{\text{Prop}}[\Sigma^\infty]$  is NP-complete,<sup>8</sup> this would

<sup>8</sup>While Theorem 2 only proves  $\text{TR}_{\text{Prop}}[\Sigma^\infty]$  to be NP-



(● = NOTPOLYAPPROX, ○ = Easy (as poly time decision), ? = unknown approximability class)

Figure 2 Approximability of Theory Revision Tasks

mean  $P = NP$ , which is why we do not expect to obtain this result. Or if this bound was some constant  $c(x) = c \in \mathbb{R}$ , then we could efficiently obtain a solution within a factor of  $c$  of optimal, which may be good enough for some applications<sup>9</sup>.

However, not all problems can be so approximated. Following [CP91, Kan92], we define

**Definition 2** A minimization problem MINP is NOTPOLYAPPROX if there is a  $\gamma \in \mathbb{R}^+$  such that

$$\forall B \in \text{Poly}(\text{MINP}), \exists x \in \text{MINP}, \text{MinPerf}[\text{MINP}](B, x) \geq |x|^\gamma$$

Lund and Yannakakis [LY93] prove that the “MIN-GRAPH-COLOR minimization problem” is NOTPOLYAPPROX. We can use that result to prove

**Theorem 4** Unless  $P = NP$ , each of  $\text{MINTR}_{\text{Prop Disj}}[\Sigma^\infty]$ ,  $\text{MINTR}_{(\text{PredCal})(\text{Horn})}[\Sigma^\infty]$  and  $\text{MINTR}_{\text{Prop Atom}}[\Sigma^K]$  is NOTPOLYAPPROX.

While these results may at first seem trivial, given that it is NP-hard to determine if a perfect theory exists, notice from Equation 3 that  $\text{MinPerf}[\text{MINTR}[\Sigma^\infty]](\cdot)$  essentially ignores such perfect theories. Note also that this result holds in the context based on an “inconsistent” oracle, in such situations, no theory can be perfect.

As  $|x|$  can get arbitrary large, this result means that these  $\text{MINTR}_x[\Sigma^\dagger]$  tasks cannot be approximated by any constant, nor even by any logarithmic factor nor any sufficiently small polynomial, etc.

### 3.3 Special Cases

If the theory is too general (i.e., returns Yes too often), then we may want to consider “specializing” it by applying only the “delete rule” and “add antecedent” transformations. In particular, recall that  $\Sigma^{+A-R}[T]$  is the set of theories obtained using an arbitrary number of such transformations, and  $\Sigma^{-R}[T]$  (resp.,  $\Sigma^{+A}[T]$ ), is the set of theories obtained by applying an arbitrary number of “delete rule” (respectively, “add antecedent”) transformations. Similarly, if the theory is too specific (i.e.,

hard, this problem is clearly in NP.

<sup>9</sup>There are such constants for some other NP-hard minimization problems. For example, there is a polynomial-time algorithm that computes a solution whose cost is within a factor of 1.5 for any TRAVELINGSALESMAN WITH-TRIANGLE\_EQUALITY problem, see [GJ79, Theorem 6.5].

returns No too often), then we may want to consider “generalizing” it by applying only the “add rule” and “delete antecedent” transformations, here, we consider  $\Sigma^{+R-A}[T]$ ,  $\Sigma^{+R}[T]$  and  $\Sigma^{-A}[T]$ , which are the set of theories obtained by applying an arbitrary number of such transformations.

Even using only these transformations, almost all of these tasks remain intractable.

**Theorem 5** It is easy to solve

$$\text{TR}_{\text{Prop Perf}}[\mathcal{G}] \quad \text{for } \mathcal{G} \in \{\Sigma^{+R-A}, \Sigma^{+R}, \Sigma^{-A}\}$$

$$\text{TR}_{\text{Prop Atom Perf}}[\mathcal{S}] \quad \text{for } \mathcal{S} \in \{\Sigma^{-R+A}, \Sigma^{-R}, \Sigma^{+A}\}$$

However, every other situation, formed by any other combination of restrictions (read “subscripts”) is NP hard (See middle and right of Figure 1).

Worse.

**Theorem 6** Unless  $P = NP$ , each of the following is NOTPOLYAPPROX

- $\text{MINTR}_{(\text{PredCal}) \text{Atom}}[\mathcal{S}]$ ,  $\text{MINTR}_{\text{Prop}(\text{Horn})}[\mathcal{S}]$   
for  $\mathcal{S} \in \{\Sigma^{-R+A}, \Sigma^{-R}, \Sigma^{+A}\}$
- $\text{MINTR}_{(\text{PredCal}) \text{Atom}}[\mathcal{G}]$ ,  $\text{MINTR}_{\text{Prop Disj}}[\mathcal{G}]$   
for  $\mathcal{G} \in \{\Sigma^{+R-A}, \Sigma^{+R}, \Sigma^{-A}\}$
- $\text{MINTR}_{\text{Prop Atom}}[\Sigma^\dagger]$   
for  $\Sigma^\dagger \in \left\{ \begin{array}{l} \Sigma^{+A-R(K)}, \Sigma^{-R(K)}, \Sigma^{+A(K)} \\ \Sigma^{-A+R(K)}, \Sigma^{+R(K)}, \Sigma^{-A(K)} \end{array} \right\}$

(See middle and right of Figure 2.)

In each of these cases, however, there is a trivial polynomial-time algorithm that can produce a theory whose accuracy (not inaccuracy) is within a factor of 2 of optimal. That is, using the ratio of an algorithm’s accuracy to the optimal value,

$$\text{MaxPerf}[\text{MAXTR}_x[\Sigma^\dagger]](B, x) = \frac{A(\text{opt}(x))}{A(B(x))}$$

**Theorem 7**

For  $\Psi \in \{\Sigma^{-R+A}, \Sigma^{-R}, \Sigma^{+A}, \Sigma^{+R-A}, \Sigma^{+R}, \Sigma^{-A}\}$ ,

$$\exists B \in \text{Poly}(\text{MAXTR}[\Psi]), \text{MaxPerf}[\text{MAXTR}[\Psi]](B, x) \leq 2$$

The companion paper [Gre95a] considers other related cases, including the above special cases in the context where our underlying theories can use the not ( $\cdot$ ) operator to return Yes if the specified goal cannot be proven, i.e., using Negation-as-Failure [Cla78]. It also considers the effect of re-ordering the rules and the antecedents, in the context where such shufflings can affect the answers returned. In most of these cases, we show that the corresponding maximization problem is not approximatable within a particular polynomial.

(The extended [Gre95b] explain the asymmetry between  $TR_{prop}$   $Perf[E^R]$  and  $TR_{PROP}$   $Perf[E+^R]$ , and discusses how these results relate to both inductive logic programming, and to default theories )

#### 4 Conclusion

A knowledge-based system can produce incorrect answers to queries if its underlying theory is faulty. A "theory revision" system attempts to transform a given theory into a related one that is as accurate as possible, using a given set of correctly answered "training queries". This report describes both the sample and computational complexity of this task. It first provides the number of samples required to obtain the statistics needed to identify a theory (from within a class of theories defined by applying various standard transformations to a given initial theory) whose accuracy will be within  $\epsilon$  of the optimal theory in this class, with probability at least  $1 - \delta$ . It then shows that, in general, the task of computing this globally optimal theory is intractable — and worse, that no polynomial time algorithm can be guaranteed to find a solution that is even close to optimal (given the standard  $P \neq NP$  assumption). We also present special cases of these tasks, which pin-point exactly when the task becomes tractable.

#### References

- [AGM85] Carlos E Alchourron, Peter Gardenfors, and David Makinson. On the logic of theory change. Partial meet contraction and revision functions. *Journal of Symbolic Logic*, 50:510-30, 1985.
- [BCH90] E Boros, Y Crama, and P L Hammer. Polynomial time inference of all valid implications for horn and related formulae. *Annals of Mathematics and Artificial Intelligence*, 1:21-32, 1990.
- [BI88] G Benedek and A Itai. Nonuniform learnability. In *Proceedings ICALP 88*, pages 82-92, 1988.
- [BM93] Paul T Baffes and Raymond J Mooney. Symbolic revision of theories with M-of-N rules. In *Proceedings of IJCAI 93*, August 1993.
- [Cla78] K Clark. Negation as failure. In H Gallaire and J Minker, editors, *Logic and Data Bases*, pages 293-322. Plenum Press, New York, 1978.
- [Coh90] William W Cohen. Learning from textbook knowledge: A case study. In *Proceeding of AAAI 90*, 1990.
- [Coh92] William W Cohen. Abductive explanation based learning: A solution to the multiple inconsistent explanation problems. *Machine Learning*, 8(2):167-219, March 1992.
- [CP91] P Crescenzi and A Panconesi. Completeness in approximation classes. *Information and Computation*, 93(2):241-62, 1991.
- [CS90] Susan C raw and Derek Sleeman. Automating the refinement of knowledge-based systems. In L C Aiello, editor, *Proceedings of ECAI 90*. Pitman, 1990.
- [DP91] Jon Doyle and Ramesh Patil. Two theses of knowledge representation: Language restrictions, taxonomic classification, and the utility of representation service\*. *Artificial Intelligence*, 48(3):1991.
- [FP93] Michael Frazier and Leonard Pitt. Learning from entailment: An application to prepositional horn sentences. In *Proceedings of IML 93*, pages 120-27. Morgan Kaufmann, 1993.
- [Gar88] Peter Gardenfore. *Knowledge in Flux: Modeling the Dynamics of the Epistemic States*. Bradford Book, MIT Press, Cambridge MA, 1988.
- [GJ79] Michael R Garey and David S Johnson. *Computers and Intractability: A Guide to the Theory of NP Completeness*. W H Freeman and Company, New York, 1979.
- [Gre95a] Russell Greiner. The challenge of revising impure theories. In *Proceedings of the Twelfth International Machine Learning Conference*, 1995.
- [Gre95b] Russell Greiner. The complexity of theory revision. Technical Report SCR 95-TR-539. Siemens Corporate Research, 1995. ftp://scr.Eiamani.com/pub/learning/P»peral/greinar/comp-tr.pE
- [Kan92] Viggo Kann. *On the Approximability of NP Complete Optimization Problems*. PhD thesis, Royal Institute of Technology, Stockholm, 1992.
- [KM91] Hirofumi Katsuno and Alberto Mendelzon. On the difference between updating a knowledge base and revising it. In *Proceedings of KR 91*, pages 387-94, Boston, April 1991.
- [LDRG94] Pat Langley, George Drastal, R Bharat Rao, and Russell Greiner. Theory revision in fault hierarchies. In *Proceedings of The Fifth International Workshop on Principles of Diagnosis (DX 94)*. New Paltz, NY, 1994.
- [Lev84] Hector J Levesque. Foundations of a functional approach to knowledge representation. *Artificial Intelligence*, 23:155-212, 1984.
- [LY93] Carsten Lund and Mihalis Yannakakis. On the hardness of approximating minimization problems. In *Proceeding of Twenty fifth Annual ACM Symposium on Theory of Computation (STOC 93)*, pages 286-93, 1993.
- [MB88] S Muggleton and W Buntine. Machine invention of first order predicates by inverting resolution. In *Proceedings of IML 88*, pages 339-51. Morgan Kaufmann, 1988.
- [OM94] Dirk Ourston and Raymond J Mooney. Theory refinement combining analytical and empirical methods. *Artificial Intelligence*, 66(2):273-310, 1994.
- [Pol85] PG Politakia. *Empirical Analysis for Expert Systems*. Pitman Research Notes in Artificial Intelligence, 1985.
- [Vap82] V N Vapnik. *Estimation of Dependences Based on Empirical Data*. Springer Verlag, New York, 1982.
- [WP93] James Wogulis and Michael J Pazzani. A methodology for evaluating theory revision systems: Results with Audrey II. In *Proceedings of IJCAI 93*, pages 1128-1134, 1993.