

MASSACHUSETTS INSTITUTE OF TECHNOLOGY
ARTIFICIAL INTELLIGENCE LABORATORY
and
CENTER FOR BIOLOGICAL COMPUTATIONAL LEARNING
WHLAKER COLLEGE

A.I. Memo No. 1449
C.B.C.L. Paper No. 86

November, 1993

Formalizing Triggers: A Learning Model for Finite Spaces

Partha Niyogi and Robert C. Berwick

Abstract

In a recent seminal paper, Gibson and Wexler ([1], GW) take important steps to formalizing the notion of language learning in a (finite) space whose grammars are characterized by a finite number of *parameters*. One of their aims is to characterize the complexity of learning in such spaces. For example, they demonstrate that even in finite spaces, convergence may be a problem since it is possible under some single-step gradient ascent methods to remain at a local maximum. From the standpoint of learning theory, however, GW leave open several questions that can be addressed by a more precise formalization in terms of Markov structures (a possible formalization suggested but left unpursued in a footnote of GW). In this paper we explicitly formalize learning in a finite parameter space as a Markov structure whose states are parameter settings. Several important results that follow directly from this characterization, include (1) a corrected version of GW's central convergence proof; (2) an explicit formula for calculating the transition probabilities between hypotheses and the existence of "problem states" in addition to local maxima; (3) an explicit calculation of the time needed to converge, in terms of number of (positive) examples; (4) the convergence and comparison of several variants of the GW learning procedure, e.g., random walk; (5) batch- and PAC-style learning bounds for the model.

Copyright © Massachusetts Institute of Technology, 1993

This report describes research done within the Center for Biological and Computational Learning in the Department of Brain and Cognitive Sciences, and at the Artificial Intelligence Laboratory. This research is supported by NSF grant 9217041-ASC and ARPA under the HRC program. Correspondence by e-mail could be directed to pn@ai.mit.edu or berwick@ai.mit.edu.

1 Introduction: The Triggering Model as a Markov structure

Recently, Gibson and Wexler ([1], GW) have begun to formalize the notion of language learning in a (finite) space whose grammars (and languages) are characterized by a finite number of parameters or 1-dimensional Boolean-valued arrays, n long. A grammar in this space is simply a particular n -length array of 0's and 1's; hence there are 2^n possible grammars (languages). One of Gibson and Wexler's aims is to establish that under some simple hill-climbing learning regimes, namely, single-step gradient ascent, some linguistically natural, finite, spaces are unlearnable, in the sense that positive-only examples lead to *local maxima*—incorrect hypotheses from which a learner can never escape. More broadly, they wish to show that learnability in such spaces is still an interesting problem in that there is a substantive learning theory concerning feasibility, convergence time, and the like, that must be addressed beyond traditional linguistic theory and that might even choose between otherwise adequate linguistic theories.

In this paper, we choose as a convenient starting point their Triggering Learning Algorithm (TLA) to focus our investigation of parameter learning. Our central result is that the performance of this algorithm is completely modeled by a Markov chain. The remainder of the current paper is devoted to exploring the basic consequences of this fact.

Let us first review the GW model and the TLA. Following Gold [2] the basic framework is that of identification in the limit. The learner (child) starts out in an arbitrary state—some setting of the n parameter values. The learner (child) receives a (countably infinite) sequence of positive example sentences drawn from some target language. After each presentation, the learner can either (i) stay in the same state; or (ii) move to a new hypothesis state, using the algorithm given below. If after some finite number of examples the learner converges to the correct target language (= parameter settings) and never changes state, then it has correctly identified the target language; otherwise, it does not converge.

In addition, in the GW model the language learner obeys two fundamental constraints: (1) the *single-value constraint*—the learner can change only 1 parameter value at a time; and (2) the *greediness constraint*—the learner is given a positive example it cannot recognize (accept), and if the learner changes one parameter value and finds that it can accept the example, then the learner retains that new parameter value. Finally, we also recall GW's definition of a *local trigger* (minor variations aside): given values for all parameter values, a *local trigger* for value v of parameter p is a sentence s from the target grammar G such that s is grammatical iff $f_i(pv) = v$. GW then state their TLA as follows:

- [Initialize] Step 1. Start at some random point in the (finite) space of possible parameter settings, specifying a single hypothesized grammar with the resulting extension as a language;
- [Process input sentence] Step 2. Receive a positive

example sentence e_i at time t (examples drawn from the language of a single target grammar, $L(G_t)$), from a uniform distribution on the language (we shall be able to relax this distributional constraint later on);

[Learnability on error detection] Step 3. If the current grammar parses (generates) e_i then go to Step 2; otherwise, continue.

[Single-step gradient-ascent] Select a single parameter at random uniformly with probability $1/n$, step from its current setting, and change it (0 to 1, 1 to 0) iff that change allows the current sentence to be analyzed; otherwise go to Step 2;

Of course, this algorithm never halts in the usual sense. GW aim to show under what conditions this algorithm converges “in the limit”—that is, after some number, n , of steps, where n is unknown, the correct target parameter settings will be selected and never be changed. Their central claim is stated as their Theorem 7 in their manuscript).

Theorem 1 *As long as the probability is always greater than a lower bound b ($b > 0$) that the learner will 1) encounter a local trigger for some incorrectly-set parameter P , and 2) then reset P accordingly to the target value, it turns out that the target grammar can always be learned using the Triggering Learning Algorithm*

1.1 The Markov formulation

From the standpoint of learning theory, however, GW have opened several questions that can be addressed by a more precise formalization of this model in terms of Markov chains (a possible formalization suggested but unpursued in footnote 9 of GW). We can picture the hypothesis space, of size 2^n , as a set of points, each corresponding to one particular array of parameter settings (languages, grammars). Call each point a *hypothesis state* or simply *state* of this space. As is conventional, we define these languages over some alphabet Σ as a subset of Σ^n . One of them is the target language (grammar). We arbitrarily place the (single) target grammar at the center of this space. Since by the TLA the learner is restricted to moving at most 1 binary value in a single step, the theoretically possible transitions between states can be drawn as (directed) lines connecting parameter arrays (the hypotheses) that differ by at most 1 binary digit (a 0 or 1) in some corresponding position in their arrays). Recall that this is the so-called *Hamming distance*.

Now, if we further place *weights* on the transitions from one state i to state j corresponding to the nonzero b 's mentioned in the theorem above; these correspond to the probabilities that the learner will move from hypothesis state i to state j . In fact, as we shall show below, given a distribution over $L(G)$, we can further carry out the calculation of the actual b 's themselves. Thus, we

Note that the notion of “trigger” does not enter into the statement of the TLA or the constraints the TLA employs,

but only into the statement of the theorem

can picture the TLA learning space as a directed, therefore C is not learnable, a contradiction. In the labeled graph V with 2^k vertices. More precisely, we can second case, without loss of generality, assume there are make the following remarks about the TLA system GW exactly two absorbing states, the first S corresponding describe. to the target parameter setting, and the second S' corresponding to some other setting. By the definition of an absorbing state, in the limit C will with some nonzero probability enter S and never exit. S' Then C is not learnable, a contradiction. Hence our assumption that there is not exactly 1 AS must be false.

Remark. The TLA system is *memoryless*, that is, given a sequence s of sentences up to time t , the selection of hypothesis h depends only on sentence s_t and not (directly) on previous sentences, i. e.,

$$p\{h(s) \leq x_i | x(t), t \leq t_1\} = P\{x(t_i) \leq x_i | x(t_{i-1})\}$$

In other words, the TLA system is a classical discrete stochastic process, in particular, a discrete Markov process or Markov chain. We can now use the theory of Markov chains to describe TLA parameter spaces [3]. For example, as is well known, we can convert the graphical representation of an n -dimensional Markov chain M to an $n \times n$ matrix T , where each matrix entry (i, j) represents the transition probability from state i to state j . A single step of the Markov process is computed via the matrix multiplication $T \times T$; n steps is given by T^n . A "1" entry in any cell (i, j) means that the system will converge with probability 1 to state j , given that it starts in state i .

As mentioned, not all these transitions will be possible in general. For example, by the single value hypothesis, the system can only move 1 Hamming bit at a time. Also, by assumption, only differences in surface strings can force the learner from one hypothesis state to another. For instance, if state i corresponds to a grammar that generates a language that is a proper subset of another grammar hypothesis j , there can never be a transition (nonzero b) from j to i , and there must be one from i to j . Further, by assumption and the TLA, it is clear that once we reach the target grammar, there is nothing that can move the learner from this state, since all remaining positive evidence will not cause the learner to change its hypothesis. Thus, there must be a loop from the target state to itself, with some positive label b and no exit arcs. In the Markov chain literature, this is known as an *Absorbing State* (AS). Obviously, a state that only leads to an AS will also drive the learner to that AS. Finally, if a state corresponds to a grammar that generates some sentences of the target language, it is always a loop from many state to itself, that has some nonzero probability. Clearly, one can conclude at once the following learnability result:

Theorem 2 Given a Markov chain C corresponding to a GW TLA learner, \exists exactly 1 AS (corresponding to the target grammar/language) iff C is learnable.

Proof. \Leftarrow . By assumption, C is learnable. Now assume for sake of contradiction that there is not exactly one AS. Then there must be either 0 AS or > 1 AS. In the first case, by the definition of an absorbing state, there is no hypothesis in which the learner will remain forever. In brief, GW attempt to show that the probability of the learner avoiding the target forever is zero by showing that the fact that some minimal cycle occurs infinitely often makes the probability of the infinite sequence zero. In other words every way in which the learner avoids the target has probability zero. Thus they conclude that probability of the event

$$\text{Event} = \text{Learner avoids target forever}$$

is zero, more precisely, they claim

$$Pr[\cup W_\alpha] = 0$$

where each W_α is a path avoiding the target and $\cup W$ is set of all such paths. However, as is well known, this union computation is true iff it is taken over a countable number of elements. In the example at hand, the crucial omission in the argument is that there are an uncountable number of ways in which the learner can avoid the target. This is because there are an uncountable number of sequences of numbers between 1 and $M-1$. The base $M-1$ expansion of any real number in the

²GW construct an identical transition diagram in the description of their computer program for calculating local maxima. However, this diagram is not explicitly presented as a Markov structure; it does not include transition probabilities. Of course, topologically both structures must be identical.

interval $[0, 1)$ would yield such a sequence (e.g., $\cos(\sqrt{2})$). Suppose SOV (setting #5 = $[0\ 1\ 0]$) is the target grammar (language). With the GW3-parameter system

Since there are an uncountable number of ways there are $2^3 = 8$ possible hypotheses, so we can draw which the event of avoiding the target forever can be as an 8-point Markov configuration space, as shown realized, the fact that each such way has probability zero. The shaded rings represent increasing *does not* imply that the total event has probability zero. The shaded rings represent increasing Hamming distances from the target. Each labeled as well. To see this consider a random variable X while is a Markov state, a possible array of parameter a uniform distribution on $[0, 1]$. Now consider the settings or grammar, hence extensionally specifies a possible target language. Each state is exactly 1 binary digit away from its possible transition neighbors. Each directed arc between the points is a possible (nonzero) transition from state i to state j ; we shall show how to

Event: $X < 1/2$

There are many ways in which this event could occur e.g. $X = 1/4$, $X = 1/3$, $X = 0.234$ etc. Each of these ways has probability zero i.e., $P[X = 1/4] = 0$, $P[X = 1/3] = 0$, $P[X = 0.234] = 0$, . . . and so on. However we know that the probability of the event $X < 1/2$ is $1/2$ not zero. This is because there are an *uncountable* number of ways in which the event $X < 1/2$ could take place. Thus the proof as given in [1] is incorrect. One correct way to formulate the proof is by resorting to an explicit Markov formulation, as suggested but not executed in GWs footnote 9, and as we established above. A similar conceptual difficulty seemingly leads to their failure to note that there are other states *besides* local maxima, for which convergence may not occur.

Corollary 1 *Given a Markov chain corresponding to a (finite) family of grammars in a GW learning system if there exist 2 or more AS, then that family is not learnable.*

Example.

Consider the GW3-parameter system. Its binary parameters are: (1) Spec(ifier) first (0) or last (1); (2) Comp(lement) first (0) or last (1); and Verb Second (V2) does not exist (0) or does exist (1). By *Specifier* GW means the standard linguistic convention of whether a phrase is part of a phrase that “specifies” that phrase, roughly like *the old in the old book*; by *Complement* GW means a phrase’s arguments, like *a nice-cream in John ate an ice-cream with envy in green with envy*. There are also 7 possible “words” in this language: S, V, O, Q, Adv, and Aux, corresponding to Subject, Verb, Object, Direct Object, Indirect Object, Adverb, and Adjective. There are 12 possible surface strings for each

(-V2) grammar and 18 possible surface strings for each (+V2) grammar if we restrict ourselves to unembedded or “degree-0” examples for reasons of psychological plausibility (see GW for discussion). Note that the “surface strings” of these languages are actually *phrases* such as Subject, Verb, and Object. Figure (3) of GW summarizes the possible binary parameter settings in this system. For instance, parameter setting (5) corresponds to the array $[0\ 1\ 0]$ = Specifier first, Comp last, and -V2, which works out to the possible basic English surface phrase order of Subject-Verb-Object (SVO). As shown in GWs figure (3), the other possible arrangements of surface strings corresponding to this parameter setting include SV; SVOI O2 (two objects, as in *give John an ice-cream*); S Aux V (as in *John is a nice guy*); S Aux V O; S Aux VOI O2; Adv S V (where Adv is an Adverb, like *quickly*); Adv S VO; Adv S VOI O2; Adv S Aux V; Adv S Aux VO; and Adv S Aux VOI O2.

Derivation of Transition Probabilities for the Markov TLA Structure

The computation of the transition probabilities from the language family can be computed by a direct extension of the procedure given in GW. Let the target language consist of the strings s_1, s_2, s_3, \dots , i.e.,

$$L_t = \{s_1, s_2, s_3, \dots\}$$

Let there be a probability distribution P on these strings. Suppose the learner is in a state corresponding to the language L . Suppose it now receives the strings s_1, s_2, s_3, \dots with probability $P(s_i)$. There are two cases to consider depending upon whether or not the string s_i is parameterizable by the grammar corresponding to the current parameter setting.

Case I. Suppose the learner can syntactically analyze the received string s_i . By the TLA, it will not change its

parameter values. In the Markov chain formulation, can now be given as,

learner remains in the same state. Remember that this state corresponds to the language M also note that $P[s \rightarrow s] = 1 - \sum_{k \text{ is a neighboring state of } s} P[s \rightarrow k]$

this situation arises only when in the language L . Therefore the probability of the learner remaining in the state s is $P(s)$

Case II. Suppose the learner cannot syntactically analyze the string. Then $s \notin L_s$. By the TLA, the learner chooses a parameter at random flips it, and if the new parameter setting makes analyzable, it retains the value and moves to the corresponding state; otherwise it remains in its original state s . Let us examine this situation using the Markov chain formulation. The learner is in state s . It has n neighboring states each at a Hamming distance of 1 from itself. The learner picks one of these uniformly at random. Imagine that these neighboring states correspond to languages which contain s . If the learner picks any one of these states (which of course it does with probability $1/n$) it would stay in that state. If the learner picks any of the other states (with probability $(n-1)/n$) then it remains in state s . Note that p of course could be 0 which means that none of the neighboring states would allow the string to be analyzed. The maximum value p could take is n . Thus we see that the probability that the learner remains in state s is $P(s) = ((n-1)/n)$. The probability that it moves to each of the other states is $1/n$.

Clearly this allows us to compute the probability that the learner will remain in its original state s as the sum of the probabilities of the above two cases, namely the following expression:

$$\sum_{s_j \in L_s} P(s_j) + \sum_{s_j \notin L_s} (1 - n_j/n) P(s_j)$$

The above expression is still a little untidy because of the n_j 's in it. We would like to clean it up a little. To do this consider the way we would compute the transition probability of state s to some other neighboring state, say k in the chain. From the above analysis, we see that such a transition will occur with probability $1/n$ for all the strings that are in the language L but not in the language s . The strings themselves occur with probability $P(s)$ each and so the transition probability is:

$$P[s \rightarrow k] = \sum_{s_j \in L_t, s_j \notin L_s} (1/n) P(s_j)$$

Note that the above summation is done over all strings $s_j \in (L_t \cap L_k) \setminus L_s$ where \setminus is the set difference symbol. It is easy to see that

$$s_j \in (L_t \cap L_k) \setminus L_s \Leftrightarrow s_j \in (L_t \cap L_k) \setminus (L_t \cap L_s).$$

Thus we can rewrite the transition probability as

$$P[s \rightarrow k] = \sum_{s_j \in (L_t \cap L_k) \setminus (L_t \cap L_s)} (1/n) P(s_j)$$

Since we have shown this in generality where for any given target, we can compute the transition probability between any two states in the Markov chain formulation of the parameter space, the self-transition probability

Finally, given any parameter space with n parameters, we have 2^n languages. Fixing one of them as the target language we obtain the following procedure for constructing the corresponding Markov chain. Note that this is the GW procedure for finding local maxima, with the addition of a probability measure on the language family.
 (Assign distribution) First fix a probability measure P on the strings of the target language L
 (Enumerate states) Assign a state to each language i.e., each L
 (Normalize by the target language.) Intersect all languages with the target language to obtain for each i , the language $L_i \cap L_t$. Thus with state i associated with language L_i now associate the language L_i
 (Take set differences.) Now for any two states i and k , if they are more than 1 Hamming distance apart, then the transition $P[i \rightarrow k] = 0$. If they are 1 Hamming distance apart then $P[i \rightarrow k] = P(L_k \setminus L_i)$.

This model captures the dynamics of the TLA comb the sum.
Example.

Consider again the 3-parameter system in the previous figure with target language 5. We can calculate the following set differences to build the Markov figure straight forwardly.

- 1. $L_1 \cap L_5 = \emptyset$ (no strings in common between L_1 and target L_5).
- 2. $L_2 \cap L_5 = \{S V, S V O, S V O I O_2, S Aux V, S Aux V O, S Aux V O I O_2\}$.
- 3. $L_3 \cap L_5 = \emptyset$.
- 4. $L_4 \cap L_5 = \{S V, S V O, S Aux V\}$.
- 5. $L_5 \cap L_5 = L_5$.
- 6. $L_6 \cap L_5 = \{S V, S V O, S V O I O_2, S Aux V, S Aux V O, S Aux V O I O_2\}$
- 7. $L_7 \cap L_5 = \{S V, Adv S V\}$.
- 8. $L_8 \cap L_5 = \{S V, S V O, S Aux V\}$.

From these values alone, we can draw the figure illustrated, and find the local maxima. For example, since the normalized state set for state 1 is the empty set, the set difference between states 1 and 5 gives all of the target language; so there is a (high) transition probability from state 1 to state 5. Similarly, since states 7 and 8 share some target language strings in common, such as $S V$, and do not share others, such as $Adv S$ and $S V O$, they learner can move from state 7 to 8 and back again. Many additional properties of the triggering learning system now become evident once the mathematical formulation has been given. It is easy to imagine other

alternatives to the TLA that will avoid the local maxima problem. One could also look at whether any local maxima exist. For example, as it stands the learner only considers one parameter setting if that change allows the learner to analyze the sentence it could not analyze before. If we relax this condition so that in this case the learner picks a parameter at random to change then the problem with local maxima disappears, because there can be only 1 Absorbing State, namely the target grammar. All other states have exit arcs. Thus, by formulation is that it allows us to also analyze convergence times. Given the transition matrix of a Markov chain, the problem of how long it takes to converge has is, occasionally the learner gets strings that are not the target language. GW state (fn. 4, p. 5) that this is not a problem the learner need only pay attention to frequent data. But this is of course a serious problem for the model. Unless some kind of memory limit is added, the learner cannot know whether the example it receives is noise or not. This being so, then there is always some finite parameter spaces where consistency might not be, however small, of escaping a local maximum. It appears that the identification in the limit framework given is simply incompatible with the notion of noise unless a memory window of some kind is added.

Perhaps the significant advantage of the Markov chain formulation is that it allows us to also analyze convergence times. Given the transition matrix of a Markov chain, the problem of how long it takes to converge has is, occasionally the learner gets strings that are not the target language. GW state (fn. 4, p. 5) that this is not a problem the learner need only pay attention to frequent data. But this is of course a serious problem for the model. Unless some kind of memory limit is added, the learner cannot know whether the example it receives is noise or not. This being so, then there is always some finite parameter spaces where consistency might not be, however small, of escaping a local maximum. It appears that the identification in the limit framework given is simply incompatible with the notion of noise unless a memory window of some kind is added.

We may now proceed to ask the following questions about the TLA more precisely:

1. Does it converge?
2. How fast does it converge? How does this vary with distributional assumptions on the input examples?
3. Can we now compute the dynamics for other “natural” parameter systems, like the 10-parameter system for the acquisition of stress in languages opened by [4]?
4. Variants of TLA would correspond to other Markov structures. Do they converge? If so, how fast?
5. How does the convergence time scale up with the number of parameters?
6. What is the computational complexity of learning parametrized language families?
7. What happens if we move from on-line to batch learning? Can we get PAC-style bounds [6]?
8. What does it mean to have non-stationary (non-ergodic) Markov structures? How does this relate to assumptions about parameter ordering and naturalization?
9. What other parametrizations can we consider?

In the remainder of this paper we shall consider these and other questions. We turn first to the question of convergence and convergence times.

3.1 Some Transition Matrices and Their Convergence Curves

Let us begin by following the procedure detailed in the previous section to actually obtain a few transition matrices. Consider the example which we looked at informally in the previous section. Here the target grammar was grammar 5 and the L languages have already been obtained. For simplicity, let us first assume a uniform distribution on the strings in L , the probability the learner sees a particular string is $1/12$ because there are 12 (degree-0) strings in L . We can now compute the transition matrix entries if not otherwise specified:

	L_1	L_2	L_3	L_4	L_5	L_6	L_7	L_8
L_1	$\frac{1}{2}$	$\frac{1}{6}$			$\frac{1}{3}$			
L_2		1						
L_3			$\frac{3}{4}$	$\frac{1}{12}$			$\frac{1}{6}$	
L_4		$\frac{1}{12}$		$\frac{11}{12}$				
L_5					1			
L_6					$\frac{1}{6}$	$\frac{5}{6}$		
L_7					$\frac{5}{18}$		$\frac{2}{3}$	$\frac{1}{18}$
L_8						$\frac{1}{12}$	$\frac{1}{36}$	$\frac{1}{9}$

3 Convergence Times for the Markov Chain Model

The Markov chain formulation gives us some distinct advantages in theoretically characterizing the language acquisition problem. First, we have already seen that given a Markov Chain one could investigate whether or not it has exactly one absorbing state corresponding to the target grammar. This is equivalent to the question of whether any local maxima exist. Note also (following the previous figure as well) that state 4 only exits to either itself or to state 5. More precisely, if T is the transition probability matrix of a chain, then t_{ij} is the probability that the learner moves from state i to state j in one step. It is a well-known fact that if one

considers the corresponding i, j element of T . This is not clear, presumably the issue of learnability even in the 3-parameter case deserves re-examination in light of this possibility.

Of course, one can examine other details of this part of the paper. However, let us now look at a case where there is no local maximum. This is the case when the matrix should contain 0's everywhere in the 3-parameter case. Consider the transition matrix obtained when the target language is again we assume a uniform distribution on strings of the target.

	L_1	L_2	L_3	L_4	L_5	L_6	L_7	L_8
L_1		$\frac{1}{3}$			$\frac{2}{3}$			
L_2		1						
L_3		$\frac{1}{3}$			$\frac{2}{3}$			
L_4		1						
L_5					1			
L_6					1			
L_7					1			
L_8					1			

	L_1	L_2	L_3	L_4	L_5	L_6	L_7	L_8
L_1	1							
L_2	$\frac{1}{6}$	$\frac{5}{6}$						
L_3	$\frac{5}{18}$		$\frac{2}{3}$	$\frac{1}{18}$				
L_4		$\frac{3}{36}$	$\frac{1}{36}$	$\frac{8}{9}$				
L_5	$\frac{1}{3}$				$\frac{23}{36}$	$\frac{1}{36}$		
L_6		$\frac{5}{36}$				$\frac{31}{36}$		
L_7			$\frac{1}{18}$				$\frac{11}{12}$	$\frac{1}{36}$
L_8				$\frac{1}{18}$				$\frac{17}{18}$

Examining this matrix we see that if the learner starts out in states 2 or 4, it will certainly end up in state 2 in the limit. These two states correspond to local maxima in the GW framework. If the learner starts in either of these two states, it will never reach the target. From the matrix we also see that if the learner starts in states 5 through 8, it will certainly converge in the limit to the target grammar.

The situation regarding states 1 and 3 is more interesting. If the learner starts in either of these states, it will reach the target grammar with probability 2/3 and reach state 2, the other absorbing state with probability 1/3. Thus we see that local maxima are *not* the only problem for learnability. GW (p. 26 in manuscript) focuses exclusively on local maxima, and indirectly implies that these are the only difficult states: “rest of the source grammars have local triggers that enable the learner to get to the target. . . however, there exist pairs of source and target grammars from the parameter space given in the table in Figure 3, such that no data point for the target grammar will ever shift the learner out of the source grammar. . . There are six such pairs of source local maximum and target grammars” They then go on to list in their figure 4, *two* such local maxima for the target grammar 5, corresponding to states 2 and 4. The quantity $p(m)$ is easy to interpret. Thus $p(m) =$

$$\begin{bmatrix} p_1(m) \\ p_2(m) \\ p_3(m) \\ p_4(m) \\ p_5(m) \\ p_6(m) \\ p_7(m) \\ p_8(m) \end{bmatrix}$$

$$\lim_{m \rightarrow \infty} p_i(m) = 1$$

While this statement is strictly true, it does not mean that for every initial state of the learner the set of source states that never lead to the target grammar is at least 0.95. Further there is one initial state (states 2 and 4) which do not converge to the target grammar, it is also true that states 1 and 3 will not converge to the target. Thus we find on looking at the curve that the learner converges with high probability within 100 to 200 (degree-0) than that presented in Figure 4 of GW. This difference is again due to the new probabilistic framework introduced in the current paper, and in fact is related to the difficulty found earlier with the central convergence “partial” distributions of examples, and we are currently looking just at minimal paths and cycles in fact (some possible learning paths. In the appendix of this paper, we provide a complete list of all starting states which might result in non-learnability. While the implication is that the existence of additional non-learnable starting points in the 3-parameter space,

and then, if an input sentence cannot be analyzed, ~~we~~ ~~revert~~ ~~to~~ ~~the~~ ~~initial~~ ~~state~~. This matrix has non-zero elements (transition probabilities) exactly where the earlier matrix had non-zero elements. However, the value of each transition is always some finite probability of exiting a non-target state. In particular if we choose $a = 1/12$, $b = 2/12$, $c = 3/12$, $d = 1/12$

To satisfy the reader's curiosity, we provide the convergence curves for a random walk algorithm (RWA) on the 8 state space. We find that the convergence times are actually faster than for the TLA; see figure 2. Since the RWA is also superior in that it does not suffer from the same local maxima problems as TLA, the conceptual support for the TLA is by no means clear. Of course it may be that the TLA has empirical support, in the sense of independent evidence that children do use this procedure (given by the pattern of their errors, et al.) but this evidence is lacking, as far as we know.

Now that we have made a first attempt to quantify convergence time, several other questions can be raised. How does convergence time depend upon the distribution of the data? How does it compare with other kinds of Markov structures with the same number of states? How will the convergence time be affected if the number of states increases, i.e. the number of parameters increases? How does it depend upon the way in which the parameters relate to the surface strings? Are there other ways to characterize convergence times? We proceed to answer some of these questions.

3.2 Distributional Assumptions

In the earlier section we assumed that the data was uniformly distributed. We computed the transition matrix for a particular target language and showed that convergence times were of the order of 100-200 samples. In this section we show that the convergence times depend especially upon the distribution. In particular we can choose a distribution which will make the convergence time as large as we want. Thus the distribution-free convergence time for the 3-parameter system is infinite.

As before, we consider the situation where the target language is L . There are no local maxima problems for this choice. We begin by letting the distribution be parametrized by the variables a, b, c, d where

$$\begin{aligned} a &= P(A = \{\text{Adv VS}\}) \\ b &= P(B = \{\text{Adv VOS, Adv Aux VS}\}) \\ c &= P(C = \{\text{Adv VOI O2 S, Adv Aux VOS, Adv Aux VOI O2 S}\}) \\ d &= P(D = \{\text{VS}\}) \end{aligned}$$

Thus each of the sets A, B, C and D contain different degree-0 sentences of L . Clearly the probability of the set $L \setminus \{A \cup B \cup C \cup D\}$ is $1 - (a + b + c + d)$. The elements of each defined subset are equally likely with respect to each other. Setting positive values for a, b, c, d such that $a + b + c + d < 1$ now defines a unique probability for each degree(0) sentence. For example, the probability of *Adv VOS* is $b/2$, the probability of *Adv Aux VOS* is $c/3$, that of *VOS* is $(1 - (a + b + c + d)) / 6$ and so on.

We can now obtain the transition matrix corresponding to this distribution. This is shown in Table 1.

Compare this matrix with that obtained with a uniform distribution on the sentences of L (the earlier curves.)

3.3 Absorption Times

In the previous sections, we computed the transition matrix for a variety of distributions and showed the rate of convergence. In particular we plotted $p(m)$, (the probability of converging from the most unfavorable initial state) against m (the number of samples). However, this is not the only way to characterize convergence times. Given an initial state, the time taken to reach the absorption state (known as the absorption time) is a random variable. One can compute the mean and variance of this random variable. For the case when the target language is L we have seen that the transition matrix has the form

$$T = \begin{pmatrix} 1 & 0 \\ R & Q \end{pmatrix}$$

Here Q is a 7-dimensional square matrix. The mean absorption times from states 2 through 8 is given by the vector (see Isaacson and Madsen [3])

$$\mu = (I - Q)^{-1} \mathbf{1}$$

where $\mathbf{1}$ is a 7-dimensional column vector of ones. The vector of second moments is given by

$$\mu' = (I - Q)^{-1} (2\mu - \mathbf{1}).$$

Using this result, we can now compute the mean and standard deviation of the absorption time from the most unfavorable initial state of the learner. (We note that the second moment is fairly skewed in such cases and so

Learning scenario	Mean abs. time	St. Dev. of abs. time
TLA (uniform)	34.8	22.3
TLA ($a = 0.99$)	45000	33000
TLA ($a = 0.9999$)	4.5×10^6	3.3×10^6
RW	9.6	10.1

3.4 Eigenvalue Rates of Convergence

In classical Markov chain theory, there are also well-known convergence theorems derived from a consideration of the eigenvalues of the transition matrix. We state these matrices in terms of its eigenvalues.

Theorem 3 Let T be an $n \times n$ transition matrix with n linearly independent left eigenvectors $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ corresponding to eigenvalues $\lambda_1, \dots, \lambda_n$. Let \mathbf{x}_0 (an n -dimensional vector) represent the starting probability of being in each state of the chain and π be the limiting probability of being in each state. Then after k transitions, the probability of being in each state is described by

$$\| \mathbf{x}_0 T^k - \pi \| = \left\| \sum_{i=1}^n \lambda_i^k \mathbf{x}_0 y_i \mathbf{x}_i \right\| \leq \max_{2 \leq i \leq n} |\lambda_i|^k \sum_{i=2}^n \| \mathbf{x}_0 y_i \mathbf{x}_i \|$$

where the \mathbf{x}_i 's are the right eigenvectors of T .

This theorem thus bounds the rate of convergence to the limiting distribution π (in cases where there is one absorption state, π will have a 1 corresponding to that state and 0 everywhere else). Using this result we can now bound the rates of convergence (in terms of number k of samples) by:

Learning scenario	Rate of Convergence
TLA (uniform)	$O(0.94^k)$
TLA ($a = 0.99$)	$O((1 - 10^{-4})^k)$
TLA ($a = 0.9999$)	$O((1 - 10^{-6})^k)$
RW	$O(0.89^k)$

This theorem also helps us to see the connection between the number of examples and the number of parameters since a chain with n states (corresponding to an $n \times n$ transition matrix) represents a language family with $\log(n)$ parameters.

4 Batch Learning Upper and Lower Bounds: An Aside

So far we have discussed a memoryless learner moving from state to state in parameter space and hopefully converging to the correct target in finite time. As we saw, this was well-modeled by our Markov formulation. In this section however we step back and consider upper and lower bounds for learning finite language families. The learner was allowed to remember all the strings encountered and optimize over them. Needless to say, this might not be a psychologically plausible assumption, but it can shed light on the information-theoretic complexity of the learning problem.

Consider a situation where there are n languages L_1, L_2, \dots, L_n over an alphabet Σ . Each language can

be represented as a subset of Σ^* .

$$L_i = \{\omega_{i1}, \omega_{i2}, \dots\} \subseteq \Sigma^*$$

The learner is provided with positive data (strings that belong to the language) drawn according to distribution P on the strings of a particular target language. The learner is to identify the target. It is quite possible that the learner receives strings that are in more than one language. In such a case the learner will not be able to uniquely identify the target. However, as more and more data becomes available, the probability of having received only ambiguous strings becomes smaller and smaller and eventually the learner will be able to identify the target uniquely. An interesting question to ask then is how many samples does the learner need to see so that with high confidence it is able to identify the target, i.e. the probability that after seeing that many samples, the learner is still ambiguous about the target is less than δ . The following theorem provides a lower bound.

Theorem 4 The learner needs to draw at least $M = \max_{j \neq i} \frac{1}{\ln(1/p_j)} \ln(1/\delta)$ samples (where $p_j = P(L_i \cap L_j)$) in order to be able to identify the target with confidence greater than $1 - \delta$.

Proof. Suppose the learner draws m (less than M) samples. Let $k = \arg \max_{j \neq i} p_j$. This means 1) $M = \frac{1}{\ln(1/p_k)} \ln(1/\delta)$ and 2) that with probability p the learner receives a string which is in both L_i and L_k . Hence it will be unable to discriminate between the target and the k th language. After drawing m samples, the probability that all of them belong to the set $L_i \cap L_k$ is $(p)^m$. In such a case even after seeing m samples, the learner will be in an ambiguous state. Now $(p_k)^m > (p_k)^M$ since $m < M$ and $p_k < 1$. Finally since $M \ln(1/p_k) = \ln(1/\delta)$, we see that $(p_k)^m > \delta$. Thus the probability of being ambiguous after m examples is greater than δ which means that the confidence of being able to identify the target is less than $1 - \delta$. ■

This simple result allows us to assess the number of samples we need to draw in order to be confident of correctly identifying the target. Note that if the distribution of the data is very unfavorable, that is, the probability of receiving ambiguous strings is quite high, then the number of samples needed can actually be quite large. While the previous theorem provides the number of samples necessary to identify the target, the following theorem provides an upper bound for the number of samples that are sufficient to guarantee identification with high confidence.

Theorem 5 If the learner draws more than $M = \frac{1}{\ln(1/(1-b))} \ln(1/\delta)$ samples, then it will identify the target with confidence greater than $1 - \delta$. (Here $b = P(L_i \setminus \cup_{j \neq i} L_j)$).

Proof. Consider the set $L = \cup_{j \neq i} L_j$. Any element of this set is present in the target language L and in many other language. Consequently upon receiving such a string, the learner will be able to instantly identify the target. After $m > M$ samples, the probability that the learner has not received any member of this set

is $(1 - P(L))^M = (1 - b_t)^M < (1 - b_t)^M = \delta$. Hence the probability of seeing some member of L in those samples is greater than $1 - \delta$. But seeing such a member enables the learner to identify the target so the probability that the learner is able to identify the target is greater than $1 - \delta$ if it draws more than M samples.

To summarize, this section provides a simple upper and lower bound on the sample complexity of exact identification of the target language from positive data. They do not suffer from local maxima problems. It should be pointed out, however, that the differences from PAC [6] formulation. However there is a crucial difficulty for L as the target language. Ideally the convergence rates have to be computed for each target language approximation to the target language with at least a certain confidence. In our case, this is not so. Since we are allowed to approximate the target, the sample complexity shoots up with choice of unfavorable distributions.

There are some interesting directions one could follow within this batch learning framework. One could try to get true PAC-style distribution-free bounds for various kinds of language families. Alternatively one could use the exact identification results here for linguistically plausible language families with “reasonable” probability distributions on the data. It might be an interesting exercise to recompute the bounds for cases where the learner receives both positive and negative data. Finally the bounds obtained here could be sharpened further. We intend to look into some of these questions in the future.

5 Variants of the Learning Model

We have so far focused on the TLA scheme for learning. TLA observes the single value and greediness constraints. There could be several variants of this algorithm and many of these are captured completely by our Markov formulation. We consider the following three simple variants by dropping either or both of Single Value and Greediness constraints:

Random walk with neither greediness nor single value constraints: We have already seen this example before. The learner is in a particular state receiving a new sentence, it remains in that state if the sentence is analyzable. If not, the learner moves uniformly at random to any of the other states and stays there waiting for the next sentence. This is done without regard to whether the new state allows the sentence to be analyzed.

Random walk with no greediness but with single value constraint: The learner remains in its original state if the new sentence is analyzable. Otherwise the learner chooses one of the parameters uniformly at random and flips it thereby moving to an adjacent state in the Markov structure. Again this is done without regard to whether the new state allows the sentence to be analyzed. However since only one parameter is changed at a time, the learner can only move to neighboring states at any given time.

Random walk with no single value constraint but with greediness: The learner remains in its original

state if the new sentence is analyzable. Otherwise the learner moves uniformly at random to any of the other states and stays there iff the sentence can be analyzed. The target remains in its original state.

Fig. 4 shows the convergence times for these three algorithms when L is the target language. Interestingly, all three perform better than the TLA for this task. Further, they do not suffer from local maxima problems. It should be pointed out, however, that the differences from TLA are marginal and this convergence has been shown for L as the target language. Ideally the convergence rates have to be computed for each target language approximation to the target language with at least a certain confidence. In our case, this is not so. Since we are allowed to approximate the target, the sample complexity for the algorithm on the language family as a whole.

6 Conclusion, Open Questions, and Future Directions

As the number of parameters n increases, the size of the corresponding Markov matrix grows as 2^n . Thus in the case of a 10 parameter system as found in models of English stress ([4]) the corresponding Markov structure will be a 1024×1024 matrix. We are currently conducting an analysis of this larger system to find its local maxima, analyze its convergence times, and see if its convergence times correspond to what one might find in practice with real stress systems.

Additional questions remain to be answered. One issue has to do with the “smoothness” relation between the parameter settings and the resulting surface strings. In principles-and-parameters theory, it has often been suggested that a small parameter change could lead to a large deductive change in the grammar, hence a large change in the surface language generated. In all the examples considered so far there is a smooth relation between surface sentences and parameters, in that switching from a $V2$ to a non- $V2$ system for instance, leads us to a Markov state that is not too far away from the previous one. If this is not so, it is not so clear that the TLA will work as before. In fact, the whole question of how to formulate the notion of “smoothness” in a language-grammar framework is unclear. We know that in the case of continuous functions, for example, that if the learner is allowed to choose examples (which can be simulated by selective attention), then such an “active” learner can approximate such functions much more quickly than a “passive” learner, like the one presented in GW. Is there an analog to this in the discrete, digital domain of language? How can one approximate a language? Here too mathematics may play a helpful role. Recall that there is an analog to a functional analysis of languages—namely, the algebraic approach advanced by Chomsky and Schutzenberger ([5]). In this model, a language is described by an (infinite) polynomial generating function, where the coefficients on the polynomial term x gives the number of ways of deriving the string x . A (weak, string) approximation to a language can then be defined in terms of an approximation to the generating function. If this method can be deployed,

then one might be able to carry over the results of func-grammar is (VCS-V2). For cases when the target is learnable, the learner converges to the target in 100-200 samples with high (greater than 0.99) probability. Further, the variants of the TLA all previously underutilized mathematical tools to analyze and outperform the TLA in terms of convergence times. language learnability.

7 Acknowledgements

We would like to thank Ken Wexler and Ted Gibson, for valuable discussions that led to this work; all residual errors are ours. This research is supported by NSF grant 9217041-ASC and ARPA under the HPCC program

Appendix

A Learnable Grammars: The Full Story

A.1 Problem States

We provide in Table 2 a complete list of problem states. In other words we list all the initial starting grammar-target grammar pairs for which the learner is not guaranteed to converge to the target with probability 1. In fact, assuming a uniform distribution on the strings for the target grammar, it is possible to compute the probability of not converging to the target for each of these pairs. Note that this probability is non-zero for the pairs listed.

A.2 Remarks

1. We have provided a complete list of initial starting grammars from which some target is not learnable (i.e. learnable with probability 1). We notice that there are three kinds of such problem starting states. Some states correspond to sinks in the Markov Structure with respect to some target grammar. Here the learner gets stuck, never leaves it and correspondingly never converges to the target. Then there are states which are not sinks (OVS+V2 when the target is SVO-V2) but which can only move to some non-target sink, and so never converge to the target. These two kinds of problem states (starred in our table) have been listed by Gibson and Wexler in Fig. 4 (pg. 27 of manuscript). Finally there are states which are not sinks, but which can with a non zero probability converge to some non-target sink. They can also with a non-zero probability converge to the target and in this respect are distinguished from problem states of type 2.
2. We would like to observe that of the 56 possible initial grammar-target grammar combinations possible, 12 result in non-learnable situations in the 3-parameter system investigated here. This is a fairly high density of unfavourable initial configurations. It would be interesting to see how this changes with other lingual subsystems with a larger number of parameters.
3. We also did an analysis of convergence times under uniform distribution for the each target grammar. We find that the results are similar to the results displayed in the paper for the case when the target

References

- [1] E. Gibson and K. Wexler, Triggers. Linguistic Inquiry, 1993, to appear.
- [2] E. Gold, Language Identification in the Limit. Information and Control 10 (1967) 447-474.
- [3] D. Isaacson and J. Masden, Markov Chains, John Wiley, New York, 1976.
- [4] B. E. Drescher and J. Kaye, A computational learning model for metrical phonology. Cognition, 1990, 137-195.
- [5] N. Chomsky and M. Schutzenberger, The Algebraic Theory of Context-free Languages. Computer Programming and Formal Systems, North Holland, Amsterdam 1963, 53-77.
- [6] L. G. Valiant, A theory of the Learnable. Proc. of the 1984 STOC, 1984, 436-445

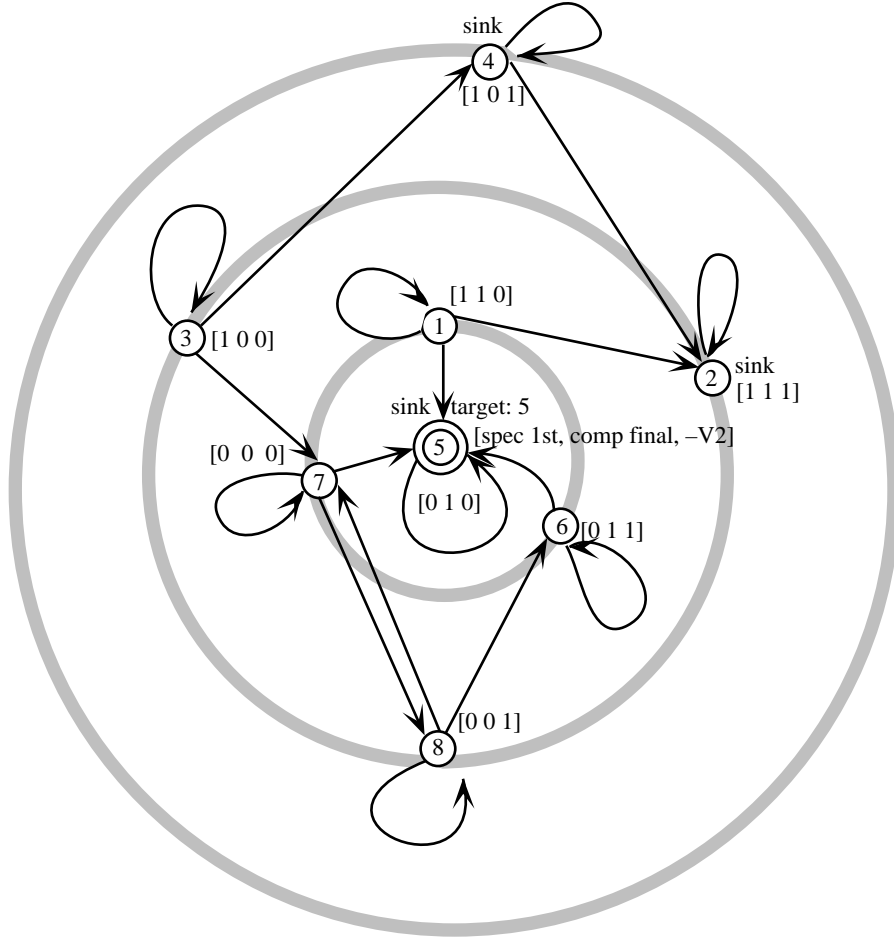


Figure 1: The 8 parameter settings in the GWexample, shown as a Markov structure, with transition probabilities omitted. (Without transition probabilities, this diagram corresponds exactly to that in GW’s appendix, as mentioned above.) Directed arrows between circles (states) represent possible nonzero (possible learner) transitions in the grammar (in this case, number 5, setting [0 1 0]), lies at dead center. Around it are the three settings that differ from the target by exactly one binary digit; surrounding those are the 3 hypotheses that differ from the target by two binary digits; the third ring out contains the single hypothesis that differs from the target by 3 binary digits. In this structure, the learner can either cycle or step in or out one ring (binary digit) at a time, according to the single-step hypothesis; but some transitions are not possible because there is no data to drive the learner from one state to another under the TLA.

	L_1	L_2	L_3	L_4	L_5	L_6	L_7	L_8
L_1	1							
L_2	$\frac{1-a-b-c}{3}$	$\frac{2+a+b+c}{3}$						
L_3	$\frac{1-a-d}{3}$		$\frac{2+a+d-b}{3}$					
L_4		$\frac{c}{3}$	$\frac{a}{3}$	$\frac{3-c-d}{3}$				
L_5	$\frac{1}{3}$				$\frac{2-a}{3}$	$\frac{a}{3}$		
L_6		$\frac{b+c}{3}$				$\frac{3-b-c}{3}$		
L_7			$\frac{a+d}{3}$				$\frac{3-2a-d}{3}$	$\frac{a}{3}$
L_8				$\frac{b}{3}$				$\frac{3-b}{3}$

Table 1: Transition matrix corresponding to a parametrized choice for the distribution on the target strings. In this case the target is 5 and the distribution is parametrized according to Section 3.2.

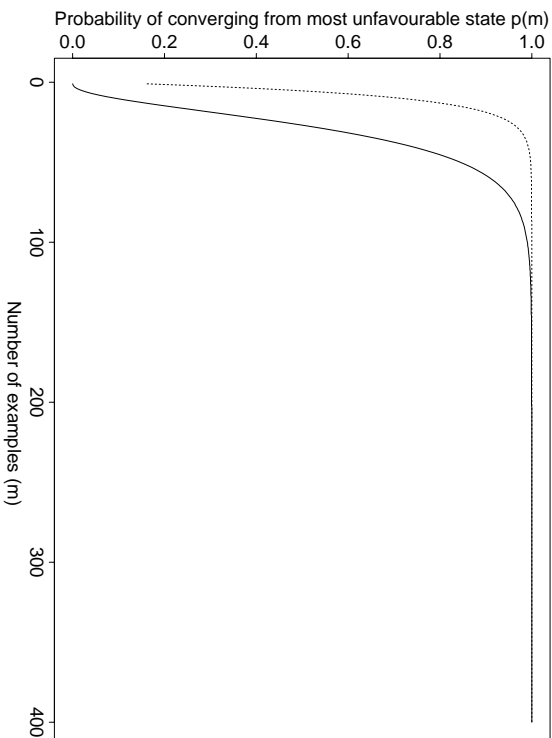


Figure 2: Convergence as function of number of examples. The horizontal axis denotes the number of examples received and the vertical axis represents the probability of converging to the target state. The data from this figure is assumed to be distributed uniformly over degree-0 sentences. The solid line represents TIA convergence and the dotted line is a random walk learning algorithm (RM). Note that random walk actually converges faster than the TIA in this case.

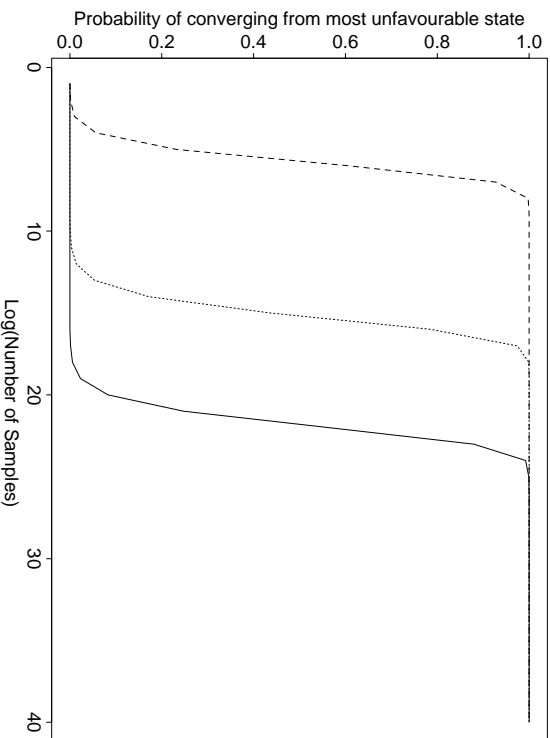


Figure 3: Rates of convergence for TIA with the target language for different distributions. The y -axis plots the probability of converging to the target after m samples and the x -axis is on a log scale, i.e., it shows $\log(m)$ and not m . The solid line denotes the choice of an “unfavorable” distribution characterized by $a = 0.9999$; $b = c = d = 0.0001$. The dotted line denotes the choice of $a = 0.99$; $b = c = d = 0.0001$ and the dashed line is the convergence curve for a uniform distribution, the same curve as plotted in the earlier figure.

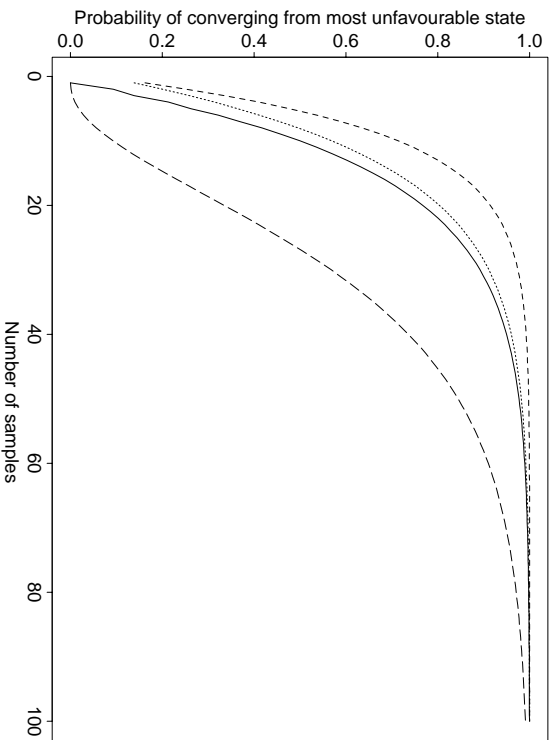


Figure 4: Convergence rates for different learning algorithms when target language. The curve with the slowest rate (large dashes) represents the TLA. The curve with the fastest rate (small dashes) is the Random (RM) with no greediness or single value constraints. Randomwalks with exactly one of the greediness and single value constraints have performances in between these two and are very close to each other.

Initial Grammar	Target Grammar	State of Initial Grammar (Markov Structure)	Probability of Not Converging to Target
(SVO-V2)	(OVS-V2)	Not Sink	0.5
(SVO+V2)*	(OVS-V2)	Sink	1.0
(SOV-V2)	(OVS-V2)	Not Sink	0.15
(SOV+V2)*	(OVS-V2)	Sink	1.0
(VOS-V2)	(SVO-V2)	Not Sink	0.33
(VOS+V2)*	(SVO-V2)	Sink	1.0
(OMS-V2)	(SVO-V2)	Not Sink	0.33
(OMS+V2)*	(SVO-V2)	Not Sink	1.0
(VOS-V2)	(SOV-V2)	Not Sink	0.33
(VOS+V2)*	(SOV-V2)	Sink	1.0
(OMS-V2)	(SOV-V2)	Not Sink	0.08
(OMS+V2)*	(SOV-V2)	Sink	1.0

Table 2: Complete list of problemstates, i.e., all combinations of starting grammar and target grammar which are in non-linearability of the target. The items marked with an asterisk are those listed in the original paper by Axelsson and Wexler [1].