

TM-0604

Inference of Reversible Context-Free
Grammars

by
Y. Sakakibara (Fujitsu)

August, 1988

©1988, ICOT

ICOT

Mita Kokusai Bldg. 21F
4-28 Mita 1-Chome
Minato-ku Tokyo 108 Japan

(03) 456-3191~5
Telex ICOT J32964

Institute for New Generation Computer Technology

Inference of Reversible Context-Free Grammars

リバーシブル文脈自由文法の推論

Yasubumi SAKAKIBARA

榊原 康文

*International Institute for Advanced Study of Social Information Science,
FUJITSU LIMITED,*

140 Miyamoto, Numazu, Shizuoka 410-03, JAPAN

沼津市宮本140 富士通(株) 国際情報社会科学研究所

Abstract

We consider the problem of learning a context-free grammar from examples. We present an efficient algorithm for learning a context-free grammar from positive examples of structural descriptions. Structural descriptions of a context-free grammar are unlabelled parse trees of the grammar, the *shapes* of parse trees. Thus the input to the learning algorithm is a finite set of shapes of parse trees. Our learning algorithm has some desirable features that the output grammar has the intended structure and the algorithm learns a grammar from positive-only examples efficiently. We show that the learning algorithm learns a grammar which is structurally equivalent to the unknown grammar and achieves the polynomial time bound.

構造的な正の例だけからリバーシブル文脈自由文法と呼ばれる文脈自由文法の部分クラスを、過剰一般化することなく帰納的に効率良く学習するアルゴリズムを提案する。このアルゴリズムは未知の文法と構造的に等価な文法を極限において同定し、また新しい正の例が入力されてから次の推測を出力するまでに多項式時間で実行する。

1 Introduction

We consider the problem of learning a context-free grammar from examples. The problem of learning a “correct” grammar for the unknown language from finite examples of the language is known as the grammatical inference problem. In the grammatical inference problem, a “correct” grammar only means a grammar which correctly generates the language. In this paper, the problem is slightly different from the usual grammatical inference problem. We consider the problem of learning a context-free grammar from positive examples of structural descriptions. The learning algorithm that we present in this problem setting outputs a grammar with the following properties.

(1) *The learned grammar has the intended structure.* The traditional grammatical inference problem is defined to identify a grammar G from examples of the unknown language L such that G correctly generates the language L , i.e., $L = L(G)$. However for any context-free language L there exist infinitely many grammars G such that $L = L(G)$. Furthermore, those grammars may have different structures. Consider the following example.

The grammar G_1 below describes the set of all valid arithmetic expressions involving a variable “ v ” and the operations of multiplication “ \times ” and addition “ $+$ ”.

$$\begin{aligned} S &\rightarrow v \mid Av \\ A &\rightarrow v+ \mid v\times \mid v+A \mid v\times A \\ &\text{(the grammar } G_1) \end{aligned}$$

However the structure assigned by the grammar G_1 to sentences is semantically meaningless. The same language can be specified by the grammar G_2 below which has a different structure from G_1 .

$$\begin{aligned} S &\rightarrow E \\ E &\rightarrow F \mid F+E \\ F &\rightarrow v \mid v\times F \\ &\text{(the grammar } G_2) \end{aligned}$$

Here the phrases are all significant in terms of the rules of arithmetic. Although G_1 and G_2 are equivalent (i.e. $L(G_1) = L(G_2)$), this fact is not very relevant from a practical point of view since it would be unusual to consider such a grammar as G_1 which assigns the structures to the sentences in a nonsignificant manner. Thus if the learned grammar must be used in a practical situation entailing the translation or interpretation of sentences like in a compiler, the structure of the learned grammar is more significant. However in the framework of the usual grammatical inference problem, it is impossible to learn such a grammar (e.g. not the grammar G_1 but G_2) which has the correct (intended) structure. To do so, it is necessary for us to assume that information on the structure of the grammar is available to the learning algorithm. In the case of context-free grammars, the structure of a grammar is usually described by the *shapes* of the parse trees, called *structural descriptions*. A structural description is a kind of tree whose internal nodes have no label. The algorithm that we present learns a context-free grammar which has the intended structure from structural descriptions.

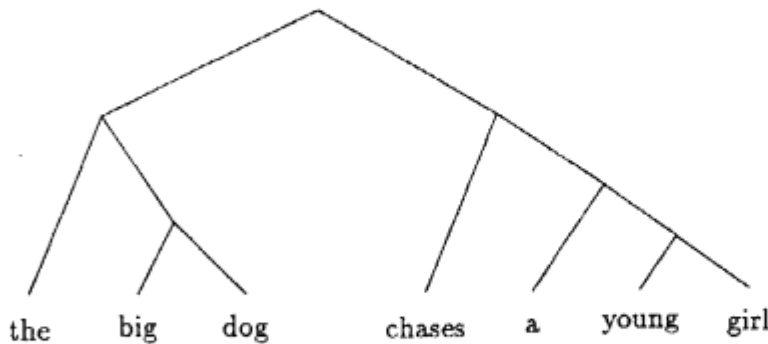


Figure 1: A structural description for “the big dog chases a young girl”

(2) *The grammar is learned from positive-only examples.* In the case of learning an unknown language L , there is a fundamental, important distinction between giving only positive information (members of L) and giving both positive and negative information (both members and nonmembers of L). A *positive presentation* of L is an infinite sequence giving all and only the elements of L . A *complete presentation* of L is a sequence of ordered pairs $\langle w, d \rangle$ from $\Sigma^* \times \{0, 1\}$ such that $d = 1$ iff w is a member of L , and such that every element w of Σ^* appears as the first component of some pair in the sequence, where Σ is the alphabet which the language L is defined over. A positive presentation eventually includes every member of L , whereas a complete presentation eventually classifies every element of Σ^* as to its membership in L . Intuitively, an added difficulty in trying to learn from positive rather than complete presentation is the problem of “overgeneralization”. Learning from positive presentation is strictly less powerful than learning from complete presentation. Gold [Gol67] shows that any set of languages containing all the finite languages and at least one infinite language cannot be identified in the limit from positive presentations. This result applies to many important classes of languages (e.g., the regular languages and the context-free languages). However Angluin [Ang80] gives a characterization of the sets of recursive languages that can be identified in the limit from positive presentation. In this paper, we consider the problem of learning a context-free grammar from positive presentation because assuming the teacher giving positive information of the grammar is acceptable in a practical use, whereas assuming the teacher giving complete information of it is not so easy for users. Since, in our problem setting, information of the grammar is the structural descriptions of it, it is assumed that positive presentation of structural descriptions is given to the learning algorithm. As we said before, the class of context-free grammars cannot be identified from positive presentation. We define a subclass of context-free grammars, called *reversible context-free grammars*, that is still powerful to define usual languages and invertible, and show that the class of reversible context-free grammars can be identified from positive presentation of structural descriptions.

(3) *The grammar is learned efficiently.* In practical use of the grammatical inference, the crucial point is the time efficiency of the learning algorithm. One of criteria for evaluating the time efficiency of the learning algorithm is the polynomial time bound. Several learning algorithms for different domains [Ang87, Sak88] have been studied to achieve the polynomial time bound. We investigate an algorithm for learning a reversible context-free grammar in polynomial time. In this paper, we extend Angluin’s efficient algorithm [Ang82] for learning a finite automaton from positive presentation and present an efficient algorithm for learning a reversible context-free grammar from positive presentation of structural descriptions.

2 Preliminaries

A *ranked alphabet* V is a finite set of symbols associated with a finite relation called the *rank relation* $r_V \subseteq V \times \{0, 1, 2, \dots, m\}$. V_n denotes the subset $\{f \in V \mid (f, n) \in r_V\}$ of V . Especially, we call V_0 , denoted Σ (i.e. $\Sigma = V_0$), the *terminal alphabet*. In many

cases the symbols in V_n are considered as *function symbols*. The rank of a function symbol is called its *arity* and a symbol of arity 0 is called a *constant symbol*. A *tree* over V is a mapping $t : \text{Dom}_t \mapsto V$, which labels the nodes of the tree domain Dom_t . V^T denotes the set of all trees over V . A *tree language* is any subset of V^T . A *terminal node* in Dom_t is one which has no descendant. For a set of trees T , the set of subtrees of elements of T is denoted by $\text{Sub}(T)$.

A *context-free grammar* is denoted $G = (N, \Sigma, P, S)$, where N and Σ are alphabets of *nonterminals* and *terminals* respectively such that $N \cap \Sigma = \emptyset$. P is a finite set of productions; each production is of the form $A \rightarrow \alpha$, where A is a nonterminal and α is a string of symbols from $(N \cup \Sigma)^*$. Finally, S is a special nonterminal called the *start symbol*. If $A \rightarrow \beta$ is a production of P and α and γ are any strings in $(N \cup \Sigma)^*$, then $\alpha A \gamma \Rightarrow \alpha \beta \gamma$. \Rightarrow is the reflexive and transitive closure of \Rightarrow . The *language generated* by G , denoted $L(G)$, is $\{w \mid w \text{ is in } \Sigma^* \text{ and } S \Rightarrow w\}$. Two context-free grammars G_1 and G_2 are said to be *equivalent* if $L(G_1) = L(G_2)$. A *parenthesis grammar* is a context-free grammar $G = (N, \Sigma, P, S)$ such that the productions in P are restricted to the form $A \rightarrow \langle \alpha \rangle$, where \langle and \rangle are special symbols not in Σ and α contains neither \langle nor \rangle . Without loss of generality, we restrict our consideration to only ϵ -free context-free grammars.

Let $G = (N, \Sigma, P, S)$ and $G' = (N', \Sigma, P', S')$ be context-free grammars. G is *isomorphic* to G' iff there exists a bijection φ of N onto N' such that $\varphi(S) = S'$, and for every $A, B_1, \dots, B_k \in N \cup \Sigma$, $A \rightarrow B_1 \dots B_k \in P$ iff $\varphi(A) \rightarrow B'_1 \dots B'_k \in P'$ where $B'_i = \varphi(B_i)$ if $B_i \in N$ and $B'_i = B_i$ if $B_i \in \Sigma$ for $1 \leq i \leq k$.

Let $G = (N, \Sigma, P, S)$ be a context-free grammar. For A in $N \cup \Sigma$, the set $D_A(G)$ of trees over $N \cup \Sigma$ is recursively defined as :

$$D_A(G) = \begin{cases} \{a\} & \text{if } A = a \in \Sigma, \\ \{A(t_1, \dots, t_k) \mid A \rightarrow B_1 \dots B_k, t_i \in D_{B_i}(G) (1 \leq i \leq k)\} & \text{if } A \in N. \end{cases}$$

A tree in $D_A(G)$ is called a *parse tree* of G from A . For the set $D_S(G)$ of parse trees of G from the start symbol S , the S -subscript will be deleted.

A *skeletal alphabet* Sk is a ranked alphabet consisting of only the special symbol σ with the rank relation $r_{Sk} \subseteq \{\sigma\} \times \{1, 2, 3, \dots, m\}$. A tree defined over $Sk \cup \Sigma$ is called a *skeleton*. Let $t \in V^T$. The *skeletal* (or *structural*) *description* of t , denoted $s(t)$, is a skeleton with $\text{Dom}_{s(t)} = \text{Dom}_t$ such that

$$s(t)(x) = \begin{cases} t(x) & \text{if } x \text{ is a terminal node,} \\ \sigma & \text{otherwise.} \end{cases}$$

Let T be a set of trees. The *corresponding skeletal set*, denoted $K(T)$, is $\{s(t) \mid t \in T\}$.

Thus a skeleton is a tree which has a special symbol σ for the internal nodes. The skeletal description of a tree preserves the structure of the tree, but not the label names describing that structure.

The *structural description* of a context-free grammar G is the skeletal set $K(D(G))$. Two context-free grammars G_1 and G_2 are said to be *structurally equivalent* if $K(D(G_1))$

$= K(D(G_2))$. Note that if G_1 and G_2 are structurally equivalent, they are equivalent, too.

3 Structural Identification

Gold's theoretical study of language learning introduces a fundamental concept that is very important in inductive inference : *identification in the limit*. In the Gold's traditional definition, for an inductive inference algorithm IA that is attempting to learn the unknown language L , an infinite sequence of examples of L is presented. Then after some finite number of example presentations, IA guesses the correct conjecture of the language and never changes (converges to) its guess after this. In the case that the conjectures are in the form of grammars, IA identifies in the limit a grammar G such that $L(G) = L$.

On the other hand, as in [Sak88], in order to identify a grammar which has the intended structure, it is necessary to assume that information on the structure of the grammar is available to the learning algorithm. In the case of context-free grammars, the structure of the grammar is the structural description of it. Suppose G is the unknown grammar (not the unknown language). This is the grammar that we assume has the intended structure, and that is to be learned (up to structural equivalence) by the learning algorithm. In this case, a sequence of examples of the language $L(G)$ is replaced by a sequence of examples of the structural description $K(D(G))$. Then a learning algorithm identifies in the limit a grammar G' such that $K(D(G')) = K(D(G))$ (i.e. structurally equivalent to G). This type of identification is called *structural identification in the limit*.

4 Condition for Positive Inference

In order to do correct identification in the limit from positive presentation, we must avoid the problem of "overgeneralization". Angluin has shown in [Ang80] various conditions for identification from positive presentation that avoids overgeneralization. In her framework, the domain is a family of languages $\mathcal{L} = \{L_1, L_2, L_3, \dots\}$. A *positive sample* of the language L is a finite subset of L . One of conditions for identification from positive presentation is following.

Condition-1

A family of language *satisfies Condition-1* iff there exists an effective procedure which on any input $i \geq 1$ enumerates a positive sample S_i of L_i such that for all $j \geq 1$, if $S_i \subseteq L_j$ then L_j is not a proper subset of L_i .

This condition requires that for every language L_i of the family \mathcal{L} , there exists a "telltale" finite subset S_i of L_i such that no language of the family \mathcal{L} that also contains S_i is a proper subset of L_i .

These discussions and formulations can be applied to the structural identification.

5 Reversible Context-Free Grammars

A context-free grammar $G = (N, \Sigma, P, S)$ is said to be *invertible* iff $A \rightarrow \alpha$ and $B \rightarrow \alpha$ in P implies $A = B$. Invertible grammar is one of normal forms for context-free grammars. Thus for any context-free language L , there is an invertible grammar G such that $L(G) = L$. A context-free grammar $G = (N, \Sigma, P, S)$ is *reset-free* iff for any two nonterminals B, C and $\alpha, \beta \in (N \cup \Sigma)^*$, $A \rightarrow \alpha B \beta$ and $A \rightarrow \alpha C \beta$ in P implies $B = C$. A context-free grammar G is said to be *reversible* iff G is invertible and reset-free. A context-free language L is defined to be *reversible* iff there exists a reversible context-free grammar G such that $L = L(G)$.

The idea of the reversible context-free grammars comes from the “reversible automata” and “reversible languages” in [Ang80].

We now consider characteristic structural samples for the reversible context-free grammars. A *positive structural sample* of a context-free grammar G is a finite subset of $K(D(G))$. A positive structural sample CS of a reversible context-free grammar G is a *characteristic structural sample* for G iff for any reversible context-free grammar G' , $K(D(G')) \supseteq CS$ implies $K(D(G)) \subseteq K(D(G'))$. The following result is necessary for the proof of correct structural identification in the limit of the reversible context-free grammars from positive presentation of structural descriptions.

Proposition 1 *For any reversible context-free grammar G , there exists a characteristic structural sample.*

6 Learning Algorithm

In this section we describe and analyze the algorithm RC to learn a reversible context-free grammar from positive structural samples.

The input to RC is a finite nonempty set of skeletons Sa . The output is a particular reversible context-free grammar $G = RC(Sa)$. The learning algorithm RC begins with the primitive context-free grammar for Sa and generalizes it by merging nonterminals.

A *partition* of some set X is a set of pairwise disjoint nonempty subsets of X whose union is X . If π is a partition of X , then for any element $x \in X$ there is a unique element of π containing x , which we call the *block* of π containing x . A partition π is *finer* than another partition π' iff every block of π' is a union of blocks of π . The *trivial partition* of a set X is the class of all sets $\{x\}$ such that $x \in X$.

Let $G = (N, \Sigma, P, S)$ be any context-free grammar. If π is any partition of N , we define the context-free grammar $G/\pi = (N', \Sigma, P', S')$ induced by π as follows. N' is the set of blocks of π (i.e. $N' = \pi$). S' is the block of π that contains S . The production $Bl \rightarrow Bl_1 \cdots Bl_k$ is in P' whenever there exist $A \in Bl$ and $A_i \in Bl_i \in \pi$ or $A_i = Bl_i \in \Sigma$ for $1 \leq i \leq k$ such that $A \rightarrow A_1 \cdots A_k$ is in P .

Let Sa be a finite set of skeletons. Define the *primitive context-free grammar* for Sa , $G(Sa) = (N, \Sigma, P, S)$, as follows :

$$\begin{aligned}
N &= (Sub(Sa) - \Sigma) \cup \{S\}, \\
P &= \{\sigma(A_1, \dots, A_k) \rightarrow A_1 \cdots A_k \mid \sigma(A_1, \dots, A_k) \in N\} \\
&\quad \cup \{S \rightarrow A_1 \cdots A_k \mid \sigma(A_1, \dots, A_k) \in Sa\}.
\end{aligned}$$

Then $G(Sa)$ is a context-free grammar such that $K(D(G(Sa))) = Sa$.

Algorithm *RC*

Input : a nonempty positive structural sample Sa ;

Output : a reversible context-free grammar G ;

Procedure :

On input Sa , *RC* first constructs $G_0 = G(Sa)$, the primitive context-free grammar for Sa . It then constructs the finest partition π_f of the set N_0 of nonterminals of G_0 with the property that G_0/π_f is reversible, and outputs G_0/π_f .

To construct π_f , *RC* begins with the trivial partition of N_0 and repeatedly merges any two distinct blocks Bl_1 and Bl_2 if either of the following conditions is satisfied.

1. There exist two productions of the forms $A \rightarrow A_1 \cdots A_k$ and $A' \rightarrow A'_1 \cdots A'_k$ in P_0 such that $A \in Bl_1$ and $A' \in Bl_2$, and for $1 \leq j \leq k$, A_j and A'_j both are in the same block or are the same terminal symbols.
2. There exist two productions of the forms $A \rightarrow A_1 \cdots A_k$ and $A' \rightarrow A'_1 \cdots A'_k$ in P_0 and an integer l ($1 \leq l \leq k$) such that $A_l \in Bl_1$ and $A'_l \in Bl_2$, A and A' are in the same block, and for $1 \leq j \leq k$, $j \neq l$, A_j and A'_j both are in the same block or are the same terminal symbols.

When there no longer remains any such pair of blocks, the resulting partition is π_f .

This completes the description of the algorithm *RC*, and we next analyze its correctness and time efficiency.

Theorem 2 *Let Sa be a nonempty positive structural sample of skeletons, and G_f be the output of the context-free grammar by the algorithm *RC* on input Sa . Then G_f is reversible and for any reversible context-free grammar G , $K(D(G)) \supseteq Sa$ implies $K(D(G_f)) \subseteq K(D(G))$.*

Theorem 3 *The algorithm *RC* may be implemented to run in time polynomial in the sum of the sizes of the input skeletons, where the size of a skeleton (or tree) is the number of symbols in its textual representation.*

Next we show that the algorithm *RC* may be used at the finite stages of an infinite learning process to identify the reversible context-free grammars in the limit from positive presentation of structural descriptions. The idea is simply to run *RC* on the sample at the n th stage and output the result as the n th guess. Define an operator RC_∞ from

infinite sequences of skeletons s_1, s_2, s_3, \dots to infinite sequences of context-free grammars G_1, G_2, G_3, \dots by

$$G_n = RC(\{s_1, s_2, \dots, s_n\}) \quad \text{for all } n \geq 1.$$

We need to show that this converges to a correct guess after a finite number of stages.

An infinite sequence of skeletons s_1, s_2, s_3, \dots is defined to a *positive structural presentation* of a context-free grammar G iff the set $\{s_1, s_2, s_3, \dots\}$ is precisely $K(D(G))$. An infinite sequence of context-free grammars G_1, G_2, G_3, \dots is said to *converge to* a context-free grammar G iff there exists an integer N such that for all $n \geq N$, G_n is isomorphic to G . By Proposition 1 and Theorem 2, we conclude the following result.

Theorem 4 *Let G be a reversible context-free grammar, s_1, s_2, s_3, \dots be a positive structural presentation of G , and G_1, G_2, G_3, \dots be the output of RC_∞ on this input. Then G_1, G_2, G_3, \dots converges to a reversible context-free grammar G' such that $K(D(G')) = K(D(G))$.*

We may modify RC by a simple updating scheme to have good incremental behavior so that G_{n+1} may be obtained from G_n and s_{n+1} .

7 Concluding Remarks

In this paper, we consider the problem of learning a context-free grammar from positive examples of structural descriptions. We make much more of the “operationality” of the grammar learned by the learning algorithm in contrast to traditional grammatical inference problems. We set up the new learning problem for context-free grammars that is slightly different from the usual grammatical inference problem. Then the grammar learned by our algorithm has some desirable properties for a practical use. Thus this problem setting makes our learning algorithm practicable.

Lastly we remark on related work. Crespi [Cre72] is most closely related, as it describes a constructive method for learning a context-free grammar from positive examples of structural descriptions. However his algorithm and our one use completely different methods and learn different classes of context-free grammars. Since our formalism is based on tree automata, one of merits of our way is the simplicity of the theoretical analysis and the easiness of understanding the algorithm, whereas the time efficiency of his algorithm [Cre72] is still not clear. Perhaps there may be a useful synthesis of these two approaches. The investigation that we must do but have not done yet is the characterization of the “reversible context-free languages”. Especially it is interesting to contrast them with noncounting context-free languages [CGM78].

This is part of the work in the major R&D of FGCP, conducted under program set up by MITI.

参考文献

- [Ang80] Dana Angluin. Inductive inference of formal languages from positive data. *Information and Control*, 45:117–135, 1980.

- [Ang82] Dana Angluin. Inference of reversible languages. *Journal of the ACM*, 29:741–765, 1982.
- [Ang87] Dana Angluin. Learning regular sets from queries and counter-examples. *Information and Computation*, 75:87–106, 1987.
- [CGM78] Stefano Crespi-Reghizzi, Giovanni Guida, and Dino Mandrioli. Noncounting context-free languages. *Journal of the ACM*, 25:571–580, 1978.
- [Cre72] Stefano Crespi-Reghizzi. An effective model for grammar inference. In B. Gilchrist, editor, *Information Processing 71*, pages 524–529, Elsevier North-Holland, 1972.
- [Gol67] E Mark Gold. Language identification in the limit. *Information and Control*, 10:447–474, 1967.
- [Sak88] Yasubumi Sakakibara. Learning context-free grammars from structural data in polynomial time. In *Proceedings of 1st Workshop on Computational Learning Theory*, pages 296–310, 1988.