

Inductive Inference of Context-free Languages

— Context-free Expression Method

Takashi YOKOMORI

International Institute for Advanced Study of
Social Information Science(IAS-SIS), FUJITSU LIMITED
140 Miyamoto, Numazu, Shizuoka 410-03 JAPAN

ABSTRACT

An inductive inference problem of context-free languages is investigated. There have been many attempts to this problem, and most of them are based on a problem setting in which a representation space for hypotheses is a class of context-free grammars. An inference algorithm given in this paper, on the contrary, employs a kind of extensions of regular expressions called context-free expressions as a representation space for context-free languages. The algorithm, based on the notion of an identification in the limit, is significantly concise when compared with existing algorithms.

1. Introduction

We consider the following model of inductive inference problem: Given an object L of inference, an inductive inference device (IID) tries to infer a representation H for the object from examples. It is assumed that IID has an enumeration mechanism by which any possible hypothesis from the representation space can be eventually enumerated at least once. It is also assumed that we can utilize an oracle for presenting examples concerning the object. IID asks the oracle for an example, and computes hypothesis and outputs it, and again asks another example for the next step, and this process is cycled. In a sequence of hypotheses H_1, H_2, \dots IID is said to identify L in the limit if there exists a positive integer n such that H_n represents L and H_{n+i} equals to H_n for all $i > 0$.

A simple algorithm for identification in the limit is the one based on the notion of identification by enumeration. Let H_1, H_2, \dots be an effective enumeration of the possible hypotheses, and suppose a set of examples e_1, e_2, \dots, e_k are presented. Then, ED provides as its next output the first hypothesis which is compatible with all these examples. Under the assumption of a perfect oracle, the sequence of hypotheses converges in the limit. ([Gold 1967])

In this paper, we deal with the inductive inference problem for context-free languages, and employ a representation space for hypotheses different from the ones in the existing methods. This enables us to make an elegant discussion on the problem and to obtain a simple algorithm for solving the problem.

2. Context-free Expressions — Extended Regular Expressions

The reader is assumed to be familiar with the rudiments in the formal language theory. (See, e.g., [Salomaa 1973] for definitions not mentioned here.)

For a given finite alphabet Σ , the set of all strings with finite length (including zero) is denoted by Σ^* . (An empty string is denoted by ϵ .) A language L over Σ is a subset of Σ^* . For an infinite alphabet Γ , L is a language over Γ if L is a language over some finite subset Σ of Γ .

The following operation plays a crucial role in this paper.

Definition 2.1 ([Gruska 1971])

(i) Let σ be a symbol and L_1, L_2 be languages. Then, σ -substitution of L_2 into L_1 , denoted by $L_1 \uparrow \sigma L_2$, is defined as follows:

$L_1 \uparrow \sigma L_2 = \{x_1 y_1 \dots x_k y_k x_{k+1} | x_1 \sigma \dots x_k \sigma x_{k+1} \in L_1, \sigma \text{ does not occur in the word } x_1 \dots x_{k+1} \text{ and } y_i \in L_2 \text{ for } 1 \leq i \leq k\}$

(ii) Let σ be a symbol and L be a language. Then, σ -iteration of L , denoted by L^σ , is defined by

$L^\sigma = \{x | x \in L \cup L \uparrow \sigma L \cup L \uparrow \sigma L \uparrow \sigma L \cup \dots, \text{ and } x \text{ has no occurrence of } \sigma\}$.

Remarks.

(1) If L_1 does not contain σ , then $L_1 \uparrow \sigma L_2 = L_1$.

(2) Let L be a language over T and suppose that T does not contain σ . Then, $L^\sigma = (L \cup \{\epsilon\})^\sigma$ and $L^{+\sigma} = (L \cup \{\epsilon\})^{\sigma+}$.

We introduce the notion of a context-free expression which provides the representation space for context-free languages.

Definition 2.2 ([Gruska 1971])

Let Γ be a (possibly infinite) alphabet, Γ' be the boldface version of Γ , i.e. $\Gamma' = \{\sigma | \sigma \in \Gamma\}$.

(i) A context-free expressions over Γ are strings over $\Gamma \cup \Gamma' \cup \{\phi, +, (\cdot,)\}$ defined as follows:

(1) ϕ is a context-free expression,

(2) if a is in $\Gamma \cup \{\epsilon\}$, then a is a context-free expression,

(3) if E_1, E_2 are context-free expressions and $\sigma \in \Gamma'$, then $(E_1 + E_2), E_1 E_2, (E_1) \sigma$ are context-free expressions,

(4) nothing else is a context-free expression.

(ii) A mapping $||$ from context-free expressions to a class of languages is defined by:

(1) $|\phi| = \Phi$ (empty set)

(2) $|a| = \{a\}$ (for $\forall a \in \Gamma \cup \{\epsilon\}$)

(3) $|E_1 + E_2| = |E_1| \cup |E_2|, |E_1 E_2| = |E_1| |E_2|, |(E_1) \sigma| = |E_1| \sigma$.

(iii) A class of languages Ω_Γ is defined as follows:

- (1) $\Phi, \{\epsilon\} \in \Omega_\Gamma$,
- (2) if $a \in \Gamma$, then $\{a\} \in \Omega_\Gamma$,
- (3) if $L_1, L_2 \in \Omega_\Gamma$ and $\sigma \in \Gamma$, then $L_1 \cup L_2, L_1 L_2$, and $L_1^\sigma \in \Omega_\Gamma$.

Example 2.1

Let a, b, σ be in Γ , then $E = (aob + ab)\sigma$ is a context-free expression. Further, $|E| = \{(aob + ab)\sigma\} = \{aob, ab\}^\sigma = \{a^i b^i \mid i \geq 1\}$. \square

Now, the next result plays an important role in this paper.

Theorem 2.1 ([Gruska 1971])

Let Σ be a finite alphabet, and let L be a language over Σ . Then, L is a context-free language if and only if there exists a finite alphabet T such that $\Sigma \subseteq T$ and $L \in \Omega_T$.

Therefore, we have the following characterization of context-free languages in terms of context-free expressions.

Theorem 2.2

Let Σ be a finite alphabet, and let L be a language over Σ . Then, L is a context-free language if and only if there exists a finite alphabet T and a context-free expression E over T such that $\Sigma \subseteq T$ and $|E| = L$.

Proof. Obvious from Theorem 2.1 and definitions. \square

Thus, context-free expressions provide a way of representing context-free languages, and as is shown later, an inductive inference algorithm for context-free languages is obtained by naturally extending the one for regular sets in terms of regular expressions.

3. Inductive Inference of Context-free Languages

3.1 Inductive Inference Algorithm

In this paper, we formalize the inductive inference problem for context-free languages as follows:

<Inductive Inference Problem for Context-free Languages>

- (1) a semantic domain D is the class of context-free languages,
- (2) a target d_0 is a given context-free language,
- (3) a representation space Ω is the class of context-free expressions,
- (4) an oracle $EX()$ gives a complete presentation of d_0 , that is, for every $e \in \Omega$ such that $|e| \subseteq d_0$, $EX()$ eventually returns "+e" at least once, and for every $e \in \Omega$ such that $|e| \not\subseteq d_0$, $EX()$ eventually returns "-e" at least once.

Now, let Σ be a finite alphabet over which a target context-free language is defined. We fix an infinite alphabet Γ over which context-free expressions are defined, where $\Sigma \subseteq \Gamma$. (It is assumed that auxiliary symbols "+", "(", and ")") are not contained in Γ .)

We define an operator δ on Ω , the set of all context-free expressions over Γ , as follows:

Suppose E, E_1, E_2 are context-free expressions over Γ . As a notation, by $E_1 \rightarrow E_2$ we denote $E_2 \in \delta(E_1)$:

- (1) $\phi \rightarrow \phi\phi$
- (2) $\phi \rightarrow a$ ($\forall a \in \Gamma \cup \{\epsilon\}$)
- (3) $\phi \rightarrow (\phi) \sigma$ ($\forall \sigma \in \Gamma$)
- (4) $\phi \rightarrow (\phi + \phi)$
- (5) if $E_1 \rightarrow E$, then $E_1 + E_2 \rightarrow E + E_2$ and $E_2 + E_1 \rightarrow E_2 + E$
- (6) if $E_1 \rightarrow E$, then $E_1 \sigma \rightarrow E \sigma$ ($\forall \sigma \in \Gamma$)
- (7) if $E_1 \rightarrow E$, then $E_1 E_2 \rightarrow E E_2$ and $E_2 E_1 \rightarrow E_2 E$.

Lemma 3.1

The operator δ defined above has the following properties:

(i) δ is complete for the most specific expression ϕ in the sense that the set $\delta^*(\phi)$ of all expressions obtainable from ϕ in a finite number of applications of δ includes at least one expression for every context-free language.

(ii) For arbitrary expressions E_1, E_2 in Ω , whenever $E_1 \in \delta(E_2)$, $|E_2| \subseteq |E_1|$ holds.

Proof. We prove by induction on the way of constructing expressions in Definition 2.2. (i) It suffices to show that $\Omega \subseteq \delta^*(\phi)$ holds. First, by the rewriting rule (2) above, we have that $a \in \delta(\phi)$ for $\forall a \in \Gamma \cup \{\epsilon\}$. Now, suppose that E_1 and E_2 are in $\delta^*(\phi)$. (In what follows, by \rightarrow^* we denote a finite number of applications of \rightarrow . Then,

$\phi \rightarrow \phi + \phi$ (by (4)) $\rightarrow^* E_1 + \phi$ (by applying (5) together with the induction hypothesis) $\rightarrow^* E_1 + E_2$ (by applying (5) together with the induction hypothesis),

thus, we have $E_1 + E_2 \in \delta^*(\phi)$. Similarly,

$\phi \rightarrow \phi\phi$ (by (1)) $\rightarrow^* E_1\phi$ (by applying (7) together with the induction hypothesis) $\rightarrow^* E_1 E_2$ (by applying (7) together with the induction hypothesis),

hence we have $E_1 E_2 \in \delta^*(\phi)$. Further, $\phi \rightarrow \phi \sigma$ (by (3)) and since $\phi \rightarrow^* E_1$, by applying (6) repeatedly, we have $\phi \sigma \rightarrow^* E_1 \sigma$, hence $\phi \rightarrow^* E_1 \sigma$, i.e. $E_1 \sigma \in \delta^*(\phi)$.

(ii) Since, from the rules (1)-(4), $\delta(\phi) = \{\phi\phi, \epsilon, \phi + \phi, \sigma, \phi\sigma \mid \forall \sigma \in \Gamma\}$, we have that for each $E \in \delta(\phi)$, $\Phi = |\phi| \subseteq |E|$.

Now, let E' be in $\delta(E)$. Then, there are only three cases concerning E . (Note that if E is in $\Gamma \cup \{\epsilon\}$, then no rule in δ is applicable to E .)

① $E = E_1 + E_2$; Let $E_i' \in \delta^*(E_i)$ ($i=1,2$). Then, by the induction hypothesis, $|E_i| \subseteq |E_i'|$ holds. Hence, by (5) $|E| = |E_1 + E_2| = |E_1| \cup |E_2| \subseteq |E_1'| \cup |E_2'| = |E'|$ or $|E| = |E_1 + E_2| = |E_1| \cup |E_2| \subseteq |E_1| \cup |E_2'| = |E'|$ is obtained.

② $E = E_1 E_2$; By the same hypothesis, $|E| = |E_1 E_2| = |E_1| |E_2| \subseteq |E_1'| |E_2'| = |E'|$ or $|E| = |E_1 E_2| = |E_1| |E_2| \subseteq |E_1| |E_2'| = |E'|$ is obtained.

③ $E = E_1 \sigma$ ($\forall \sigma \in \Gamma$); Let $E_1' \in \delta^*(E_1)$. Then, by the induction hypothesis, $|E_1| \subseteq |E_1'|$ holds. Hence by (6), $|E| = |E_1 \sigma| = |E_1| \sigma \subseteq |E_1'| \sigma = |E'|$ is obtained. This completes the proof. \square

Now, using the operator δ defined above, we obtain an algorithm for the inductive inference problem of context-free languages, which is quite simple and based on the principle of "identification by enumeration".

Definition 3.1 (admissible presentation [Laird 1986])

(i) Let d_0 be a context-free language of target. An oracle EX is called *complete* and *sufficient* for d_0 if there exists a signed subset K of Ω satisfying the following:

- (1) the set $\{E \in \Omega \mid \text{for all } e \text{ in } K, \text{ if } e \text{ is positive, then } |e| \subseteq |E|, \text{ else } |e| \not\subseteq |E|\}$ is exactly the set $\{E \in \Omega \mid |E| = d_0\}$,
- (2) for every $e \in K$ such that $|e| \subseteq d_0$, EX() eventually returns "+e" at least once, and for every $e \in K$ such that $|e| \not\subseteq d_0$, EX() eventually returns "-e" at least once.

(ii) A presentation of examples of d_0 is *admissible* if it has an oracle EX which is complete and sufficient for d_0 .

Note. There exists such a K for our case, that is, if we take K as the set $\{E \in \Omega \mid |E| \text{ is a singleton, and if it is in } d_0, \text{ then } E \text{ has a sign } +, \text{ otherwise it has a sign } -\}$, then K satisfies the conditions mentioned above. Hence, we may assume the existence of the admissible presentation of d_0 .

Before presenting an algorithm, we need some preliminaries.

For a given target context-free language L , let $T = \{a_1, \dots, a_n\}$ be the alphabet over which L is defined. Further, for each $k \geq 1$, let $\Gamma_k = T \cup \Delta_k \cup \Delta'_k \cup \{\phi, +, (\cdot)\}$, where $\Delta_k = \{\sigma_1, \sigma_2, \dots, \sigma_k\}$, $\Delta'_k = \{\sigma_1, \sigma_2, \dots, \sigma_k\}$.

For a given expression E and $i \geq 0$, $k \geq 1$, let $\delta(E, i, k) = \{E^j | E \rightarrow^i E^j \text{ and } E^j \text{ is a string over } \Gamma_k\}$. Then, $\delta^*(E) = \cup_{i \geq 0} \cup_{k \geq 1} \delta(E, i, k)$. We abbreviate $\delta(\phi, i, k)$ as $\delta(i, k)$. The algorithm requires an enumeration procedure which, starting with ϕ , enumerates every expression in $\delta^*(E)$.

As a notation, we denote an element of $\delta(i, k)$ by $[E, (i, k)]$. Note that for each i, k , $\delta(i, k)$ is finite. Further, let $\delta(i+1, k) = \{E^j | E \rightarrow E^j \text{ and } E \in \delta(i, k), E^j \in \Gamma_k^*\}$, that is, an element of $\delta(i+1, k)$ is obtained from an element of $\delta(i, k)$ by applying δ_k once, where δ_k is a restriction of δ to Γ_k . Each $\delta(1, k)$ is obtained from ϕ by one application of δ_k . [The outline of the algorithm]

Now, the algorithm works as follows: By applying δ_k to ϕ , it enumerates a hypothesis (an expression) E and stores into a queue Q in the form of $[E, (i, k)]$. then refines (generalizes) each hypothesis by examples. The enumeration is performed in the order as in

$\delta(1,1), \delta(2,1), \delta(1,2), \delta(3,1), \delta(2,2), \delta(1,3), \delta(4,1), \dots$. That is, the algorithm *dovetails* the enumeration of $\delta(i,1), \delta(i,2), \delta(i,3), \dots$ with that of $\delta(1,k), \delta(2,k), \delta(3,k), \dots$.

For a hypothesis E , if it does not cover some positive example (it is called "too specific"), then the algorithm generalizes E in some way. Otherwise, if its language $|E|$ includes some negative example (it is called "too general"), then the algorithm simply discards it. This is repeated until a correct hypothesis is found.

The algorithm is given as Algorithm A_1 .

For the sake of helping one understand the process of enumerating hypotheses in the algorithm, Figure 1 illustrates how the contents of a queue Q changes during the enumeration. (The contents of Q is represented as a rectangle whose length may change as the time goes. The shadow portion of Q denoting $\delta(1, k)$ is created only

Algorithm A_1 (Inference Algorithm for Context-free Languages)

Input: A recursively enumerable set of context-free expressions Ω

An enumeration operator δ

An admissible presentation of a target language d_0

Output: A sequence of expressions E_1, E_2, \dots such that E_n is correct for the first n examples.

Procedure:

$Q \leftarrow \delta(1,1)$; (elements of $\delta(1,1)$ are stored in the queue Q)

EXAMPLES $\leftarrow \Phi$ (empty set)

$X \leftarrow \text{next}(Q)$; (*next* removes the top element of Q)

do (forever):

EXAMPLES \leftarrow EXAMPLES \cup EX() (get next example)

while (let $X = [E_j, (i, k)]$ be the j -th element of $\delta(i, k)$, then)

$\exists e \in$ EXAMPLES s.t. E_j is not correct for e

if E_j is "too specific"

then do

if $i \geq 2$

then if $k=1$ and $j=1$

then append $\delta(1, i) \delta(E_j, 1, 1)$ to the tail of Q ;

$X \leftarrow \text{next}(Q)$;

else append $\delta(E_j, 1, k)$ to the tail of Q ;

$X \leftarrow \text{next}(Q)$;

else append $\delta(E_j, 1, k)$ to the tail of Q ;

$X \leftarrow \text{next}(Q)$;

else (E_j is "too general")

discard E_j ;

$X \leftarrow \text{next}(Q)$;

Output E_j as the next hypothesis.

where E is "too specific" :=

if $\exists +e \in$ EXAMPLES s.t. $e \notin |E|$, then return *true*

else return *false*

E is "too general" :=

if $\exists -e \in$ EXAMPLES s.t. $e \in |E|$, then return *true*

else return *false*

when $\delta(k+1, 1)$ is produced from $\delta(k, 1)$, and it is placed before $\delta(k+1, 1)$.

Theorem 3.1

For any given context-free language d_0 , the algorithm A_1 identifies d_0 in the limit.

Proof. We need to show the following two: First, the algorithm A_1 converges some hypothesis E , secondly, the hypothesis E is correct for the target d_0 .

From the completeness property of δ ((i) of Lemma 3.1) and Theorem 2.2, there is a chain of generalization steps from $\phi : \phi = E_0 \rightarrow E_1 \rightarrow \dots \rightarrow E_n = E$ and $|E| = d_0$. Here we assume that n is as small as possible for the given d_0 . Property (ii) of Lemma 3.1 together with the minimality of n ensures that for each i , $|E_i| \subset |E_{i+1}|$. Then, there are strings w_i such that $w_i \in |E_i|$ and $w_i \notin |E_{i+1}|$. Let E_n be in $\delta(n, k)$.

[Proof for correctness] Assume the algorithm converges to some expression E : that is, there exist i, k

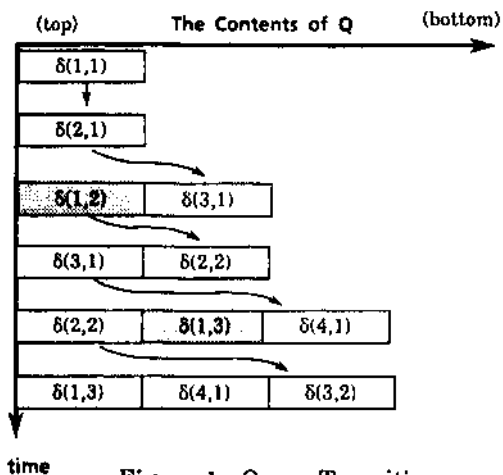


Figure 1. Queue Transition

≥ 0 and E in $\delta(i,k)$ such that E is correct for not only all examples in EXAMPLES but also every example given in the future. (That is the definition of "convergence".) Hence, E is correct for d_0 .

[Proof for convergence] Suppose that the algorithm diverges. Then, we show, by the induction on j , that every hypothesis E_j ($0 \leq j \leq n$) in the chain above appears on the top of Q and is generalized. When $j=0$, it is trivial. Suppose that E_j appears on the top of Q as a hypothesis in $\delta(i,p)$ and is generalized. Since $E_{j+1} \in \delta(E_j)$, $E_{j+1} \in \delta(i+1,q)$ for some q . The divergence of the algorithm implies that the finite number of expressions preceding E_{j+1} will all appear on the top and be generalized or simply discarded. Hence, E_{j+1} eventually appears on the top of Q and is generalized due to $w+i$. Thus E_n eventually appears on the top. However, because of the assumption of the divergence, E_n is not correct for some w , which contradicts the fact that $|E_n| = d_0$. Hence, the algorithm converges.

Thus, we conclude that the algorithm converges to a correct expression. \square

3.2 Inferring Semilinear Languages

If we restrict our attention to a subclass of context-free languages called semilinear languages ([Gruska 1971]), then we can easily get an operator for enumerating all expressions over T_2 , where $T_2 = T \cup \{\sigma_1, \sigma_2\}$, and obtain more efficient and simpler inference algorithm for the class. This is due to the fact that Any semilinear language over T is contained in nT_2^+ that is for any semilinear language L over T there exists a context-free expression E over T_2 such that $L = \setminus E$ holds. (Note that the rules (1)-(4) in 8 form a context-free grammar, provided that T is finite.)

As a natural conclusion, we have:

Theorem 3.2

There exists a simple algorithm for inferring semilinear languages in which the enumerator δ is achieved by a context-free grammar.

(Proof/Similar to that of Theorem 3.1 and omitted.)

4. Meta Inference

As we have mentioned, an expression enumerator δ for semilinear case is realized as a context-free grammar, which implies that an enumerator in the IID is identified with a context-free expression. Therefore, one can think of a meta inference problem in which for a given target, the meta-IID infers a representation (context-free expression, or enumerator) denoting the target from examples of expressions, where the target is a class of expressions (or, a class of semilinear languages denoted by the expressions). Note that since the input of meta-IID, which is an expression, can be regarded as a string over some alphabet, the inference schema of the meta-inference problem is structurally equivalent to that of the inference problem in the usual sense, which is illustrated by Figure 2.

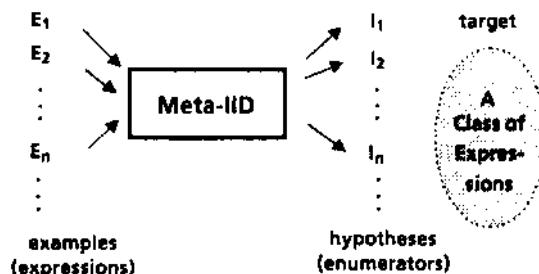


Figure 2. Meta Inductive Inference Schema

ACKNOWLEDGEMENTS

The author is grateful to Dr. T. Kitagawa, the president of IAS-SIS, for ceaseless encouragements. He is also indebted to Dr. H. Enomoto, the director of IAS-SIS, for providing useful reference papers as well as invaluable advice. Last but not least, discussion with the colleagues Y. Takada, Y. Sakakibara, and H. Ishizaka was very fruitful.

This work is a part of the major research and development of FGCS Project conducted under the program setup by MITI.

REFERENCES

- [1] Gold, E.M., Language Identification in the Limit, Information and Control 10, 447-474 (1967).
- [2] Gruska, J., A Characterization of Context-free Languages, Journal of Computer and System Sciences. 5, 353-364 (1971).
- [3] Laird, P.D., "Inductive Inference by Refinement", Technical Report TR-376, Dept. of Computer Science, Yale University, 1986.
- [4] Salomaa, A., "Formal Languages", Academic Press, 1973.
- [5] Yokomori, T., "Inductive Inference of Context-free Languages -Context-free Expression Method", Research Report 71, IAS-SIS FUJITSU LIMITED, 1986.