

# Algebraic View of Grammatical Inference

Alexander S. Saidi<sup>1,2</sup>, Souad Tayeb-bey<sup>1</sup>

1- Laboratoire de Reconnaissance de Formes et Vision INSA de Lyon- Bât. 403  
20 Ave. Albert Einstein 69621 Villeurbanne 2- Dépt. Mathématiques, Informatique et Systèmes  
Ecole Centrale de Lyon B.P. 163. 69131 Ecully  
saidi@cc.ec-lyon.fr, tayebbey@rfv.insa-lyon.fr

**.Abstract.** We consider the problem of grammatical inference (GI) for classes of structured documents like summaries, dictionaries, bibliographic data basis, encyclopaedias and so on. The inference is based on examples of individual sample of these documents. In this paper, we present an algebraic framework of the GI in which rewrite rules will define the process of generalisation. The implementation algorithm discussed here is used in a document handling project in which paper documents are typographically tagged and then recognised. One of the current applications in this project is to translate paper documents into machine readable form

## I- Introduction

We consider the problem of grammatical inference for classes of structured documents like summaries, dictionaries, scientific reports, bibliographic basis, encyclopaedias and so on from examples. We propose an algorithm that can identify the class of regular languages from positive examples. This is the usual case in the grammatical inference for classes of documents whose structures are learned from examples. In the proposed method, we infer a regular grammar which generalise the Prefix-tree of the given positive examples. Although negative examples can be considered, the usual situation in this kind of document recognition is to deal only with the positive representation. Theses examples are structured representation of documents which are collected by optical typographic recognition techniques ([1]). It is known that any algorithm that would construct a DFA (deterministic finite automaton) with a minimum number of states compatible with all the data already processed can identify any regular language in the limit ([2]). The search space being a lattice, we propose in this paper an algebraic framework for the grammatical inference and show a function on the lattice that characterise the construction of partitions over the Prefix-tree of the representation. In this framework, an initial algebra  $A_G$  is assigned to the Prefix-tree of the strings in the sample. Then the main result is that a quotient automaton of the Prefix-tree denotes a quotient algebra  $A'_G$  whose terms are obtained by the application of a uniquely defined homomorphism from  $A_G$  to  $A'_G$ . The corresponding function induced by this homomorphism is then defined by an algorithm. We also discuss an alternative view of GI problem base on the construction of a table of successor and predecessors of any element of the alphabet. Then the relation between this table is discussed. We show how this table lead to a automaton of the lattice. The proposed algorithms are used for GI in a project on paper document processing whose one application is the translation of the document into HTML text.

## II- The GI problem

The Inductive Inference paradigm is the basis of the automatic learning problem ([1]). We are interested in the Syntactical Pattern Recognition framework where many grammatical inference algorithms exist that can be used in the learning step of pattern recognition tasks ([3],[4],[5],[6], [7]).

Most of the algorithms used to infer a regular grammar deal only with positive data. However, it is well known that the class of regular languages can not be correctly identified from only positive examples. Hence, any recursively enumerable class of language is identifiable using a complete representation with positive and negative data ([14]). However, the usual situation in (paper) document recognition is to propose a generalisation of the language presented by a set on positive examples. For example, from scientific reports, we may have :

report → abstract, acknowledgement, outline, chapter, chapter, references.

report → abstract, outline, chapter, subchapter, chapter, index, references.

report → abstract, acknowledgement, outline, figure-table, chapter, chapter, index.

Although one can propose some negative examples (for instance, no report without any chapter or no report without outline), the general case is the unsupervised inference and the negative information is not available. In this paper, first we give some basic definitions for reference. Then, in section IV, an algebraic view of GI problem is given in details. In section V, some practical issues and the implementation of the proposed algorithm are reported. An alternative view of GI based on tables are presented in section VI. Then, some relationships with other works in the field is recalled in the section VII.

## III- Basic definitions

We assume the reader is familiar with context-free grammars, regular grammars and regular expressions ([15]). We give some basic definitions used in the rest of the paper.

Let  $\Sigma$  be a finite alphabet. The set of all finite strings of symbols from  $\Sigma$  is denoted  $\Sigma^*$ . The empty string is denoted by  $\epsilon$ . The concatenation of two strings  $u$  and  $v$  is denoted  $uv$ . If a string is  $u=vw$  then  $v$  is a *prefix* of  $u$  and  $w$  a *suffix* of  $u$ . A *language* denoted by  $L$  is any subset of  $\Sigma^*$ .

A *finite automaton* (FA)  $A$  is a quintuplet  $(Q, \Sigma, \delta, q_0, F)$  where  $Q$  is the set of states,  $\Sigma$  is the set of input symbols,  $\delta : Q \times \Sigma^* \rightarrow 2^Q$  is the transition function,  $q_0 \in Q$  is the start state and  $F \subseteq Q$  is the set of final states. For an automaton  $A$ , the (regular) language accepted by  $A$  is denoted by  $L(A)$ .

An automaton  $A$  is *deterministic* (DFA) if for all  $q \in Q$  and for all  $a \in \Sigma$ ,  $\delta(q,a)$  has at most one (state) element. A language is regular iff it is accepted by a FA. If  $A=(Q, \Sigma, \delta, q_0, F)$  is a FA and  $\pi$  a *partition* of  $Q$ ,  $B(q,\pi)$  is the only block that contains  $q$  and we denote the quotient set by  $Q/\pi$  as the set of all partitions  $\{B(q,\pi) \mid q \in Q\}$ . Given a FA  $A$  and a partition  $\pi$  over  $Q$ , the *quotient* (or derived) *automaton*  $A/\pi = \{Q/\pi, \Sigma, \delta', B(q,\pi), \{B \in Q/\pi \mid \exists q \in B, q \in F\}\}$  where  $\delta'$  is defined by  $\forall B, B' \in Q/\pi, \forall a \in \Sigma, B \in \delta(B,a)$  iff  $\exists q, q' \in Q, q \in B, q' \in B' : q \in \delta(q,a)$ .

It is easy to see that for a partition  $\pi$  over  $Q$ ,  $L(A) \subseteq L(A/\pi)$ . The set of all automata derived from  $A$  is a (language inclusion) lattice  $\text{Lat}(A)$ . A *Regular expression* is defined by :

$\emptyset$  is a regular expression     $\epsilon$  is a regular expression    For each  $a \in \Sigma$ ,  $a$  is a regular expression

If  $r$  and  $s$  are regular expressions, then  $(r \mid s)$ ,  $(rs)$ ,  $(r^*)$  are regular expressions. An optional regular expression  $r$  is denoted by  $[r]=(r \mid \epsilon)$  and the regular expression  $(r+)=r^*r$ .

A *context-free grammar* (CFG) is denoted by  $G=(N, T, P, S)$  where  $N$  and  $T$  are finite sets. Elements of  $N$  are called nonterminals and those of  $T$  are called terminals.  $P$  is a finite set of productions. Each production rule is of the form  $A \rightarrow \omega$  with  $A \in N, \omega \in (N \cup T)^*, |\omega| \geq 1$ . The special nonterminal  $S$  is called the start symbol. A (right) *regular grammar* is a context-free grammar whose the productions rules are of the form  $A \rightarrow \alpha$  or  $A \rightarrow \alpha B$ , where  $\alpha \in T, A, B \in N$ . The *language*  $L(G)$  is any string  $\omega \in T^*$  such that there is a derivation from  $S$  to  $\omega$  (denoted by  $S \Rightarrow^* \omega$ ). By extension, the language of any nonterminal  $A \in N$  is any string  $\omega \in T^*$  such that for  $\tau, \sigma \in T^*, S \Rightarrow^* \tau A \sigma \Rightarrow^* \tau \omega \sigma$ .

Given  $I$ , the positive representation from a regular language  $L$ ,  $I$ , is said to be *structurally complete* if all transitions of (the unknown) automaton  $A(L)$  are used in the acceptance of strings in  $I$ . We call the  $A(L)$  the *canonical automaton* of a language  $L$  which is a DFA accepting  $L$  and has the minimal number of states. The *maximal canonical automaton*  $MAC$  with respect to  $I$ , (where  $I$  is structurally complete) is the automaton whose language is  $L$  and has the largest number of states. One can define the *prefix tree* acceptor of  $I$ , denoted by  $PT(I)$  from the MCA by merging states sharing the same prefix.  $PT(I)$  accepts only the strings of  $I$ . It is well known that if  $I$  is a structurally complete sample of a regular language  $L$ , then there exists a partition  $\pi$  over the states of  $PT(I)$  such that  $PT(I)/\pi$  is isomorphic to  $A(L)$ . The aim of this paper is to give an algebraic specification of the state partitions of  $PT(I)$  in order to formally characterise a function on  $PT(I)$ . This is done by the definition of a function over the terms of an algebra associated to  $PT(I)$  that produces a quotient-algebra whose terms are (by construction) isomorphic to those of  $PT(I)$ .

The following theorems delimit the search space of  $L(A)$  ([8]): If  $I$  is a positive sample, the set  $\Gamma$  of automata such that  $I$  is structurally complete with any automaton in  $\Gamma$  is  $\text{Lat}(MCA(I))$ . If  $I$  is a positive sample,  $L$  the target language and  $I$  is structurally complete with respect to  $A(L)$ , then  $A(L)$  is an element of  $\text{Lat}(PT(I))$ .

The automaton  $A(L)$  is learnable as a partition  $PT(I)/\pi$  isomorphic to  $A(L)$  in the lattice of all possible partitions of  $PT(I)/\pi$ . The main aim of grammatical (regular) inference is to characterise this lattice and to guide a search in it. Note that this search space grows exponentially with the size of the state set in  $PT(I)$  and therefore with the size of  $I$ . The automata we will construct in the sequel are constrained automata [9]: some kind of conditional automata where transitions are constrained to fulfil some conditions. For example, a transition  $\delta(p, a/b)$  means that the transition takes place if the symbol  $\alpha \in \Sigma^*$  is followed by  $\beta \in \Sigma^*$ . By this means, one can constrain the context (left and right) of a symbol in the automaton.

## IV- The Algebraic View

We use here the relation between an initial many sorted algebra and grammars [goguen-77-78]. This relation over context-free grammars is easily extended to regular grammars. Such a grammar is used as a specification (rewrite) system to define an algebra. The major property is that the defined algebra is initial in a given category. To construct the algebra associated to a context-free grammar  $G$ , each nonterminal of  $G$  is assigned to a class of derivation tree. Consequently, the nonterminals of  $G$  are sorts of a many sorted algebra whose operations are defined by the production of  $G$ . The derivation tree (and hence the language) of any nonterminal  $X$  denote the carriers of the sort  $X$ .

Let  $G=(N,T,P,S)$  be a non ambiguous context-free grammar and  $L$  be its language (we do not use variables in the sequel; the algebraic representation of grammatical programming is detailed in [11], [12]). Let  $L_G$  be the terms (strings represented by their derivation tree) of  $G$ . Let the algebra  $A_G$ -algebra be associated to  $G$  whose signature is  $((N \cup T), P)$  where  $P$  is the set of names given to productions in  $P$ . The terms of this algebra are derivation trees starting from any nonterminal of  $G$ . The terms of the  $A_G$ -algebra are constructed by using the names of the productions of  $G$ . Constants (0-ary operators of  $A_G$ -algebra) are elements of  $T$  in the grammar  $G$ . Two strings will be *equivalent* if their derivation tree (terms of  $A_G$ -algebra) are equivalent. By equivalence, we mean an equivalence relation over the sorts of  $A_G$ . Since the nonterminals  $N$  are sorts of  $A_G$ -algebra, any string  $\omega$  such that  $X \in N$ ,  $X \Longrightarrow^* \omega$  is typed by  $X$ , particularly when  $\omega \in L(G)$  where for the start non terminal  $S \Longrightarrow^* \omega$ . Note that the terms of  $L_G$  are strings whose derivation tree are the terms of  $A_G$ . An  $A_G$ -algebra is initial in a category  $C$  if for all algebra  $B$  of  $C$ , there exist a unique homomorphism  $h_G : A_G \rightarrow B$ . By the initiality of  $A_G$ , we have a homomorphism  $h_G : A_G \rightarrow sem$  where  $sem$  denote any algebra of  $C$ . One can consider the language generated by a context-free (or regular) grammar as an example of  $sem$  and hence to associate a *string semantics* to the terms of  $A_G$ . One can also define another algebra assigned to another grammar (context-free or regular)  $G'$  and to define a homomorphism from  $G$  to  $G'$ . This function can be a (syntax-

directed) translation function. Here, we are interested by  $h : A_0 \rightarrow L_0$  (see [11] for the details of the construction of an  $A_0$ -algebra and  $h$ ). We may note that by initiality,  $h$  always exists and is defined uniquely.  $h$  is surjective and if the grammar  $G$  is not ambiguous, then  $h$  is also injective. In this case,  $h$  will be an isomorphism from  $A_0$  to  $L_0$ . For any partition of  $PT(I_0)$  where some states are considered *equivalent* (merged by some function), the generated language is a super set of  $L_0$ . We will show that there exists a uniquely defined homomorphism  $h$  from the derivations in  $PT(I_0)$  to the derivations in any partition of  $PT(I_0)$ . We also show that for some function induced by  $h$ , terms of  $PT(I_0)/\pi$  are homomorphic to those of  $PT(I_0)$ . Thus, we will characterise different DFA of the lattice  $Lat(PT(I_0)/\pi)$  by the application of this function to the terms of the  $A_0$ -algebra. Let  $PT(I_0)$  and its automaton (a regular grammar)  $G$ , a partition  $PT(I_0)/\pi$  from  $PT(I_0)$  and its automaton  $G'$ . Let  $A_0$  and  $A'_0$  be the algebra assigned to  $G$  and  $G'$ . Note that  $A_0$  and  $A'_0$  have the same signature and  $A'_0$  is a quotient algebra of  $A_0$ . Suppose further that the sorts  $p$  and  $p'$  of  $A_0$  are equivalent ( $p \equiv p'$  are in the same block of states of  $PT(I_0)/\pi$ ). Let  $[p]$  denote the equivalence class of  $p$ .

*Theorem-1* : the value of a derivation tree  $Z$  in  $A'_0$  is the equivalence class of derivation tree of  $Z$  in  $A_0$ . That is :  $(A'_0)(Z) = [A_0(Z)]$ . *proof* : this is an immediate consequence of the quotient algebra  $A'_0$  of  $A_0$ . In other words, if  $p \equiv p' \in [p]$ , then the derivation trees of  $A_0$  whose root are  $p$  and  $p'$  are equivalent terms in  $A'_0$ . One can define  $i : A_0 \rightarrow \Sigma^*$  a homomorphism from  $PT(I_0)$  to  $L$  (or  $j : A'_0 \rightarrow \Sigma^*$ ) and show that all strings of  $L$  whose derivation tree root is  $[p]$  are equivalents. Let  $V=(N \cup T)$  be the set sort of the algebra  $A_0$  and  $A'_0$ , and  $h : A_0 \rightarrow A'_0$ . Let  $\equiv_h = (\equiv_{hv})_{v \in V}$ , be a family of relations on  $A_0$  defined by :  $a \equiv_h a'$  iff  $h(a)=h(a')$ ,  $a, a' \in V$ .  $\equiv_h$  is called the congruence relation induced by the homomorphism  $h$ . The main problem of grammatical inference is to define  $h : A_0 \rightarrow A'_0$  and to assign it a function. This function is defined by  $f : PT(I_0) \rightarrow PT(I_0)/\pi$ . Let us define the properties of  $h$  (and  $f$ ).

*Theorem-2* : Let  $A_0$  and  $A'_0$  be the algebra assigned to  $PT(I_0)$  and  $PT(I_0)/\pi$ . Furthermore, let  $h : A_0 \rightarrow A'_0$  be a homomorphism and  $\equiv_h$  be the congruence relation induced by  $h$ . If the homomorphism  $h$  is surjective, then  $A_0/\equiv_h \approx A'_0$ . *Proof* : according to the theorem of homomorphism [13]. Let  $f : A_0/\equiv_h \rightarrow A'_0$  be the family of functions defined by  $f([a])=h(a)$ ,  $a$  is of a sort of  $A_0$ . We have  $[a]=[a']$  implies that  $h(a)=h(a')$ . To show that  $f$  is an isomorphism : -  $f$  is bijective : as  $h$  is surjective by assumption, then  $f$  is surjective. Further, if  $f([a])=f([a'])$ , then  $h(a)=h(a')$  and consequently  $[a]=[a']$ . Hence  $f$  is injective.

-  $f$  is a homomorphism : let  $(r : p \leftarrow \alpha q)$  be an operation of  $A_0$ . The homomorphism condition may be satisfied by :  $f((A_0/\equiv_h) r([\alpha], [q])) = f((A_0/\equiv_h) (\alpha, q)) = h((A_0) (\alpha, q)) = (A'_0)(h(\alpha), h(q))$  and by homomorphism condition  $= (A'_0)(f([\alpha]), f([q]))$

Informally, this means that given  $h$  and  $f$  the function induced by  $h$ , if the states  $p \equiv p'$  under some conditions, then for every derivation in  $PT(I_0)/\pi$ , there is a derivation in  $PT(I_0)$ . Note that we have seen that the inverse is true in the lattice  $L(PT(I_0))$ .

## V- Practical issues : definition of the homomorphism $h$

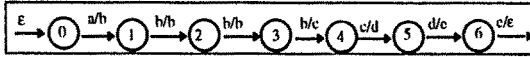
Below, we will give an algorithm (*infer*) that corresponds to the function  $f$ . The function *infer* induced by the homomorphism  $h$  is parametered by an integer  $k > 0$ . A derivation tree in  $A_0$  is of the form  $p_0(a_1/a_2, p_1(a_2/a_3, p_3(a_3/a_4, p_4(\dots))))$  where  $a_i/a_j$  means that  $a_i$  is followed by  $a_j$ . Informally speaking, if one can find two derivations.  $p_0(a_1/a_2/\dots/ak, p_1(a_2/a_3/\dots/ak+1, p_3(a_3/a_4/\dots/ak+2, p_4(\dots))))$  and  $q_0(b_1/b_2/\dots/bk, q_1(b_2/b_3/\dots/bk+1, q_3(b_3/b_4/\dots/bk+2, q_4(\dots))))$  where for  $i=1..k-1, j=i+1, a_i/a_j = b_i/b_j$ , then for  $l=0..k-1, p_l \equiv q_l$ .

Note that  $k > 0$  defines the extent of the function *infer* below where for  $k=0$ , the language of the generated partition is  $\Sigma^*$ . The more  $k$  is great, the "lower" is  $PT(I_0)/\pi$  in the lattice whose bottom element is  $PT(I_0)$  and whose top element is  $\Sigma^*$ . It is straightforward to see that for  $0 < i^j, L(PT(I_0)/\pi, i) \subseteq L(PT(I_0)/\pi, j)$ . Hence, the value of the parameter  $k$  defines the depth of the search in the lattice.

The function *infer* induced by the homomorphism  $h$  is defined as follows :

Given an integer value  $k^3$ , consider every string  $I$ . The inference process starts with  $infer(\langle q_0, \omega \rangle, \{ \}, k)$ ,  $I$ ,  $q_0$ , the start symbol (initial state). The function  $infer$  will generate a set of conditional rules for  $I$ , if  $k > 1$ . Let the call  $infer(\langle q, \omega \rangle, G, k)$  be the general case. A state variable is a state like  $q_i$ , where  $i$  is an unknown natural number to be defined at the end of each step.

Function  $infer(\langle q, \omega \rangle, G, k)$  returns  $G_{i+1} =$   
 if  $\langle q, \omega \rangle = \langle \_, \epsilon \rangle$  return  $G$ ; if  $! \text{tok} < k$  then let  $k = \text{tok}$  -- for this last turn !  
 if there exists a rule  $\langle p, v \rangle \leftarrow \langle p', \gamma \rangle$  in  $G$ , where  $v = \alpha_1 \dots \alpha_k \rho$  and  $\omega = \alpha_1 \dots \alpha_k \phi$  and unify( $q, p, G$ ) (1)  
 then return  $infer(\langle p', \alpha_1 \dots \alpha_k \phi \rangle, G)$  else create a rule  $r : \langle p, \omega \rangle \leftarrow \langle p', \alpha_1 \dots \alpha_k \phi \rangle \rightarrow p'$  of the form  $q, i$  will  
 return  $infer(\langle p', \alpha_1 \dots \alpha_k \phi \rangle, G \cup \{r\})$ ; -- be defined later. end if; end infer;



When returned from  $infer$ , different variables states of the grammar  $G_n$  of each step are assigned natural numbers beginning from the least natural not still used. Note that merges are done in  $infer$  function at (1).

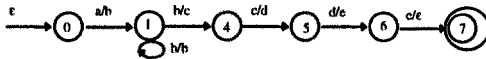
The unify function sets the relations between states.

Function  $unify(q, p, G)$  returns Boolean =  
 if  $q = p$  then return true; else if one of (or both)  $p$  and  $q$  is a state variable (say  $p$ )  
 then set  $(p=q)$  in  $G$ ; return true; else return false; end if; end if;  
 end unify;

### V-1- A simple example

We consider the following simple positive set for  $I$ :  $I = \{abbcbde, abcbde, abcbde\}$

Let  $K=2$ . Then consider the first example "abbcbde". The  $A_0$  term corresponding to this example is  $(x/y$  means "x" followed by "y") :  $p_0(a/b, p_1(b/b, p_2(b/b, p_3(b/c, p_4(c/d, p_5(d/e, p_6(e/\epsilon))))))$  depicted by the following path in  $PT(I)$ : Following  $infer$  function, the states with the same prefix of length  $KE$  are merged from 1..k-



1. Here the states 1, 2 and 3 are merged and give the following automata: In other words, we have  $\delta(p_1, b/b) = p_2$  and  $\delta(p_2, b/b) = p_3$ . Then  $p_1, p_2$  and  $p_3$  are merged. Nota bene : merging prefixes is done in other situations in an automaton : when designing an automaton, every derivation begins from  $q_0$ , where we want that every example of  $I$ , begins with  $q_0$ .

The following is the trace of the construction of the grammar by the call

$infer(\langle p_0, "abbcbde" \rangle, \{ \}, 2)$ . For the sake of simplicity, we note "a"/"b". $\omega$  by "ab". $\omega$  (that is,  $\omega \in \Sigma^*$  is concatenated to the string "ab"):

$\langle p_0, "abbcbde" \rangle$	Current goal	There is no candidate rule	Conclusion : create the rule $\langle p_0, "ab" \cdot \omega \rightarrow \langle p_1, "b" \cdot \omega \rangle$
$\langle p_1, "bbcbde" \rangle$	Current goal	create the rule $\langle p_1, "bb" \cdot \omega \rightarrow \langle p_1, "b" \cdot \omega \rangle$	
$\langle p_1, "bbcbde" \rangle$	Candidate rule	$\langle p_1, "bb" \cdot \omega \rightarrow \langle p_1, "b" \cdot \omega \rangle$	rule found $i=j$ . The states $i$ and $j$ will be merged
$\langle p_1, "bcde" \rangle$	create the rule	$\langle p_1, "bc" \cdot \omega \rightarrow \langle p_1, "c" \cdot \omega \rangle$	
$\langle p_1, "cde" \rangle$	create the rule	$\langle p_1, "cd" \cdot \omega \rightarrow \langle p_1, "d" \cdot \omega \rangle$	
$\langle p_1, "de" \rangle$	create the rule	$\langle p_1, "de" \cdot \omega \rightarrow \langle p_m, "e" \cdot \omega \rangle$	create the rule $\langle p_m, "e" \cdot \omega \rightarrow \langle \rangle$

Here  $j=i$  means that these states are merged. When the derivation is completed, integer values are assigned to the variables beginning by 1. That would set  $i=1, k=2, l=3, m=4$ . In order to have a set of rules corresponding to the above DFA (see the figure above) and for the sake of clarity, let's say that  $k=4, l=5, m=6$ .

Note that during the construction of this automaton, the rule  $\langle p_1, "bb" \cdot \omega \rightarrow \langle p_1, "b" \cdot \omega \rangle$  is reused and handles the recursive part of the derivation. The set of rules of  $G_1$  obtained from this example analysis is :

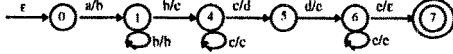
$\langle p_0, "ab" \cdot \omega \rightarrow \langle p_1, "b" \cdot \omega \rangle$     $\langle p_1, "bb" \cdot \omega \rightarrow \langle p_1, "b" \cdot \omega \rangle$     $\langle p_1, "bc" \cdot \omega \rightarrow \langle p_1, "c" \cdot \omega \rangle$     $\langle p_1, "cd" \cdot \omega \rightarrow \langle p_1, "d" \cdot \omega \rangle$   
 $\langle p_1, "de" \cdot \omega \rightarrow \langle p_m, "e" \cdot \omega \rangle$     $\langle p_m, "e" \cdot \omega \rightarrow \langle \rangle$

The regular expression : The extraction of the regular expression associated to  $G_1$  is straightforward. We have  $L(G_1) = a b^* c d e$

Note that we apply a usual generalisation rule applied in the Grammatical Inference where every  $v+$ ,  $v \in \Sigma^*$  is generalised to  $v^*$ . The call  $\text{infer}(\langle p0, "abcde" \rangle, G1, 2)$  for the next example adds two more rules to  $G1$  :  $\langle p_1, "cc".\omega \rangle \rightarrow \langle p_1, "c".\omega \rangle$  and  $\langle p_2, "ee".\omega \rangle \rightarrow \langle p_2, "e".\omega \rangle$ . Then  $G2$  is :

$\langle p0, "ab".\omega \rangle \rightarrow \langle p_1, "b".\omega \rangle$     $\langle p_1, "hb".\omega \rangle \rightarrow \langle p_1, "b".\omega \rangle$     $\langle p_1, "bc".\omega \rangle \rightarrow \langle p_2, "c".\omega \rangle$   
 $\langle p_1, "cd".\omega \rangle \rightarrow \langle p_2, "d".\omega \rangle$     $\langle p_2, "cc".\omega \rangle \rightarrow \langle p_2, "c".\omega \rangle$     $\langle p_2, "de".\omega \rangle \rightarrow \langle p_2, "e".\omega \rangle$   
 $\langle p_2, "ee".\omega \rangle \rightarrow \langle p_2, "e".\omega \rangle$     $\langle p_2, "e" \rangle \rightarrow \langle \epsilon \rangle$

The third example "abcde" is accepted by  $G2$  and does not change the grammar. The final DFA is :



The regular expression associated to  $G2$  is  $L(G2) = ab^*c^*de^*$

### V-2- The Constraint Logic Program associated to $G1$

The grammar  $G2$  can be immediately prototyped by a (constraint) logic program. We implemented the grammatical inference by a logic program which is parametered by the value of  $K$ .

The (constraint) logic program  $\text{Prg}^{\text{inf}}$  associated to  $G2$  where  $K=2$  is given by the predicate  $\text{aut}^{\text{inf}}$  :

$\text{aut}^{\text{inf}}(['a', 'b' | L], p0) \leftarrow \text{aut}^{\text{inf}}(['b' | L], p1).$     $\text{aut}^{\text{inf}}(['b', 'b' | L], p1) \leftarrow \text{aut}^{\text{inf}}(['b' | L], p1).$   
 $\text{aut}^{\text{inf}}(['b', 'c' | L], p1) \leftarrow \text{aut}^{\text{inf}}(['c' | L], p4).$     $\text{aut}^{\text{inf}}(['c', 'd' | L], p4) \leftarrow \text{aut}^{\text{inf}}(['d' | L], p5).$   
 $\text{aut}^{\text{inf}}(['c', 'c' | L], p4) \leftarrow \text{aut}^{\text{inf}}(['c' | L], p4).$     $\text{aut}^{\text{inf}}(['d', 'e' | L], p5) \leftarrow \text{aut}^{\text{inf}}(['e' | L], p6).$   
 $\text{aut}^{\text{inf}}(['e', 'e' | L], p6) \leftarrow \text{aut}^{\text{inf}}(['e' | L], p6).$     $\text{aut}^{\text{inf}}(['e', \epsilon | L], p6) \leftarrow \text{aut}^{\text{inf}}([\epsilon | L], p7).$   
 $\text{aut}^{\text{inf}}([\epsilon], p7).$

The query  $\leftarrow \text{aut}^{\text{inf}}(\omega, p0)$  submitted to  $\text{Prg}^{\text{inf}}$  will succeed if  $\omega \in L(G2)$ . In order to apply an ascending evaluation procedure, the initial state  $p_0$  should appear in a basic clause (not in the query). To do this, the above program is transformed to a dual one (using e.g. magic set techniques [16]) in which any rule of the form

$\text{aut}^{\text{inf}}([\alpha, \beta | L], p_i) \leftarrow \text{aut}^{\text{inf}}([\beta | L], p_j)$  is transformed to  $\text{aut}^{\text{inf}}([\alpha | L], \beta, p_j) \leftarrow \text{aut}^{\text{inf}}(L, \alpha, p_i)$ . From  $\text{Prg}^{\text{inf}}$ , we obtain  $\text{Prg}^{\text{inf}*}$ :

$\text{aut}^{\text{inf}*}([\epsilon], \epsilon, p0).$     $\text{aut}^{\text{inf}*}(['a' | L], 'b', p1) \leftarrow \text{aut}^{\text{inf}*}(L, 'a', p0).$     $\text{aut}^{\text{inf}*}(['b' | L], 'b', p1) \leftarrow \text{aut}^{\text{inf}*}(L, 'b', p1).$   
 $\text{aut}^{\text{inf}*}(['b' | L], 'c', p4) \leftarrow \text{aut}^{\text{inf}*}(L, 'b', p1).$     $\text{aut}^{\text{inf}*}(['c' | L], 'c', p4) \leftarrow \text{aut}^{\text{inf}*}(L, 'c', p4).$   
 $\text{aut}^{\text{inf}*}(['c' | L], 'd', p5) \leftarrow \text{aut}^{\text{inf}*}(L, 'c', p4).$     $\text{aut}^{\text{inf}*}(['d' | L], 'e', p6) \leftarrow \text{aut}^{\text{inf}*}(L, 'd', p5).$   
 $\text{aut}^{\text{inf}*}(['e' | L], 'e', p6) \leftarrow \text{aut}^{\text{inf}*}(L, 'e', p6).$     $\text{aut}^{\text{inf}*}(['e' | L], \epsilon, p7) \leftarrow \text{aut}^{\text{inf}*}(L, 'e', p6).$

where a rule  $\text{aut}^{\text{inf}*}([\alpha | L], \beta, q) \leftarrow \text{aut}^{\text{inf}*}(L, \alpha, p)$  states that in the transition  $\delta(p, \alpha) = q$ ,  $\beta$  has been the last symbol (the left context of  $\alpha$ ) and  $\alpha$  will then be the left context for the next transition. To verify if a string  $\omega \in L(G2)$ , one should reverse  $\omega$  (denoted by  $\bar{\omega}$ ) and then submit the query  $\leftarrow \text{aut}^{\text{inf}*}(\bar{\omega}, \epsilon, X)$ . This query will succeed (and will assign a value to  $X$ ) in  $\text{Prg}^{\text{inf}*}$  if the query  $\leftarrow \text{aut}^{\text{inf}}(\omega, p_0)$  succeeds in  $\text{Prg}^{\text{inf}}$  (See e.g. [9] for details).

In order to prove the properties of such a program, we would simplify  $\text{Prg}^{\text{inf}}$  by omitting the transition symbols and all other variables. We obtain a logic program on which we would apply the immediate consequence operator ([18]). The least fix-point of this operator is defined by  $T^{\uparrow} = \cup_n T^{\uparrow n}$  denoted by  $T^{\uparrow \infty}$  where  $T^{\uparrow 0} = \{\text{aut}^{\text{inf}}(p0)\}$  and  $T^{\uparrow n+1} = T^{\uparrow n}$ .  $T^{\uparrow \infty}$  characterise the set of available states of the automaton. Furthermore, for an automaton  $A = (Q, \Sigma, \delta, p_0, F)$ ,  $\text{aut}(f) \in T^{\uparrow \infty}$ ,  $f \in F$ , if there exists a valid string (starting from  $q0$ ) which ends to the final state  $f$ .

### V-3- A realistic example

Let us consider the following example of final research reports. We apply the function *infer* for  $K=1$  and  $K=2$ .

Content  $\rightarrow$  Title, Intr, Abs, Ack, Oln, Chap, Chap, Ref, Ind. Content  $\rightarrow$  Title, Abs, Oln, Chap, Chap, Ref, Ann, Anx.

Content  $\rightarrow$  Title, Intr, Abs, Oln, Chap, SChap, Chap, SChap, Ref, Ind, Anx.

Content  $\rightarrow$  Title, Intr, Oln, Chap, SChap, Chap, SChap, Chap, SChap, Ref. Content  $\rightarrow$  Title, Oln, Chap, SChap, Chap, SChap, Chap, SChap, Ref.

For the sake of simplicity, we abbreviated symbols as follows : Intr:Introduction Abs:Abstract SChap: subchapter Ack :

Acknowledgement Chap: Chapter Ref: Reference Ann: Annexes Oln:Outline Ind: Index. The set of rules generated for  $K=1$ :

$q0, "Title" \rightarrow q7, \epsilon$ $q1, "Aux" \rightarrow q2, \epsilon$ $q1, \epsilon \rightarrow q2, "Aux"$ $q2, "Ind" \rightarrow q1, \epsilon$ $q2, \epsilon \rightarrow q3, "SChap"$ $q3, "SChap" \rightarrow q3, \epsilon$ $q3, "Chap" \rightarrow q3, \epsilon$ $q3, "Ref" \rightarrow q2, \epsilon$ $q4, "Oln" \rightarrow q3, \epsilon$ $q5, "Oln" \rightarrow q3, \epsilon$ $q5, "Ack" \rightarrow q4, \epsilon$ $q6, "Oln" \rightarrow q3, \epsilon$ $q6, "Abs" \rightarrow q5, \epsilon$ $q7, "Oln" \rightarrow q3, \epsilon$ $q7, "Abs" \rightarrow q4, \epsilon$ $q7, "Intr" \rightarrow q6, \epsilon$
---

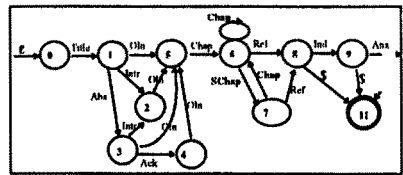
The language of this automaton is :  $L = Title [Intr [Abs [Ack]] | Abs] (SChap | Chap)^* Ref ([Ind] Anx)^* [Ind]$ . It is easy to note that for  $K=1$ , the language accepted by the automaton is larger set over  $\Sigma^*$  than the one for  $K=2$  (see also the next section). The set of rules generated for  $K=2$  is given in the next page. Note that the recursive part of the language for  $K=2$  is on the Chapter  $\times$  Sub\_Chapter part. The associated language to this automaton is :  $L = Title (Abs | Intr [Abs [Ack]]) L1$  where  $L1 = Oln (Chap SChap)^* (Chap SChap | Chap)^* L2$   $L2 = Ref [Anx Anx | Ind [Anx]]$ . The set of rules generated for  $K=2$  :

$q0, "Title"."Oln" \rightarrow q5, "Oln"$ $q0, "Title"."Abs" \rightarrow q11, "Abs"$ $q0, "Title"."Intr" \rightarrow q8, "Intr"$ $q1, "Ind"."Aux" \rightarrow q9, "Aux"$ $q1, "Ind" \rightarrow \epsilon$ $q2, "Ref"."Aux" \rightarrow q10, "Aux"$ $q2, "Ref"."Ind" \rightarrow q1, "Ind"$ $q2, "Ref" \rightarrow \epsilon$ $q3, "Chap"."Chap" \rightarrow q3, "Chap"$ $q3, "Chap"."Ref" \rightarrow q2, "Ref"$ $q4, "Chap"."SChap" \rightarrow q12, "SChap"$ $q4, "Chap"."Chap" \rightarrow q3, "Chap"$ $q5, "Oln"."Chap" \rightarrow q4, "Chap"$ $q6, "Ack"."Oln" \rightarrow q5, "Oln"$ $q7, "Abs"."Oln" \rightarrow q5, "Oln"$ $q7, "Abs"."Ack" \rightarrow q6, "Ack"$ $q8, "Intr"."Oln" \rightarrow q5, "Oln"$ $q8, "Intr"."Abs" \rightarrow q7, "Abs"$ $q9, "Aux" \rightarrow \epsilon$ $q10, "Aux"."Aux" \rightarrow q9, "Aux"$ $q11, "Abs"."Oln" \rightarrow q5, "Oln"$ $q12, "SChap"."Chap" \rightarrow q4, "Chap"$ $q12, "SChap"."Ref" \rightarrow q2, "Ref"$
---

### VI- An alternative view : The table of contexts

We give here an alternative construction of the result of *infer* and give the related automaton. We believe that this representation is a suitable alternative for the GI problem with a set of positive examples. This is done by the construction of a table where every symbol of  $\Sigma$  is given with its successor and predecessor. For the above example, we have :

Symbol $\alpha \in \Sigma$	Predecessor of $\alpha$	Successor of $\alpha$
Title	$\epsilon$	Intr, Abs, Oln
Intr	Title	Abs, Oln
Abs	Intr, Title	Ack, Oln
Ack	Abs	Oln
Oln	Ack, Abs, Intr, Title	Chap
Chap	Oln, Chap, SChap	Chap, Ref, SChap
Ref	Chap, SChap	Ind, Anx, \$
Anx	Ref, Ind, Anx	Anx, \$
SChap	Chap	Chap, Ref
Ind	Ref	Anx, \$



The special symbol '\$' denotes the end of a string (a final state of the automaton; a '\$' is added at the end of each string before analysing) while  $\epsilon$  denote the empty string at the beginning of a string. In fact, the table shows that some symbol of  $\Sigma$  can be followed (resp. preceded) by a set of other symbols in  $\Sigma$ . The canonical automaton associated to this table is (once again, for the sake of clarity, we distinguish between  $\epsilon$  - the empty string at start - and '\$' which is the very last symbol of each string) : Note that there are three recursive parts on the nodes 6, 6-7-6, and 10. In order to generalise the examples of  $L$ , the context in which the symbol  $\alpha \in \Sigma$  is present can be "forgotten" and let the symbol  $\alpha$  be followed by its successor anywhere in a string whatever be the predecessor of  $\alpha$ . Let a symbol  $\alpha \in \Sigma$ , P the set of its predecessors and S the set of its successors. Consider also the (language inclusion) lattice  $(PT(L))$  whose top element is  $\Sigma^*$  and whose bottom is  $PT(L)$ . An interesting use of this table is based on the well known but yet simple generalisation rule in Inductive Learning ([19]) which "forgets" the context of a sequence. Consequently, in the case of the above table, we can :

- 1- Allow  $\alpha \in \Sigma$  to be followed by an element of S and preceded by an element of P according to I. This corresponds to  $K=3$  in function *infer*. For example, if  $\alpha = "Abs"$ , the only possible triplets of the form  $\beta."Abs".\delta$  are  $\langle Intr, Abs, Ack \rangle$ ,  $\langle Title, Abs, Oln \rangle$  and  $\langle Intr, Abs, Oln \rangle$ .
- 2- Allow  $\alpha \in \Sigma$  to be followed by any element of S and preceded by any element of P. This gives a good degree of generalisation. This is a special case of *infer* with  $K=2$ . For  $\alpha = "Abs"$ , the possible cases are :  $\{ "Intr", "Title" \} \times "Abs" \times \{ "Ack", "Oln" \}$  which makes 4 possibilities (generalisation).
- 3- Allow  $\alpha \in \Sigma$  to be preceded and followed by any element of  $\Sigma^*$ . This is a special case of *infer* with  $K=0$ . This gives the Top element of the lattice. For  $\alpha = "Abs"$ , the possible cases are  $\beta \times "Abs" \times \delta$  with  $\delta, \beta \in \Sigma$ . Other configurations are possible which correspond to  $K=1$  : allowing  $\alpha$  to be only followed (resp. preceded) by any symbol of S (resp. P) and

free for the rest. However, during our experimentation on paper documents, the case (2) above showed to be a reasonable configuration. While an automaton specifies rules to construct sequences of, say characters, the choices above (variations of  $K$ ) let generate strings with various degree of constraints. Note that if negative examples are taken into account, more constraints must be verified and the precedence table alone will not suffice.

## VII- Conclusion

A new algorithm for GI has been presented which is immediately prototyped by a constraint logic program. The algebraic specification allows to show that the homomorphism  $h$  exists and we gave an implementation of it by the *infer* function. An alternative view of GI based on the construction of the table of successors and predecessors of any symbol in  $\Sigma^*$  has been outlined. The relationship between this alternative and *infer* function and the formalisation of this approach are in hands. Most of the related works on Grammatical Inference deal with positives and negative examples. When only positive examples are available (which describe the characteristic cases), researchs concern rather the Structured Documents field and have led to several document standards like ODA and SGML. Among other works in the field, [20] and [21] proposed similar methods for document analysis. But the algebraic framework of the grammatical Inference, the logical aspects and the table manipulation for the direct grammar extraction have, as well as known, not yet been investigated. This work is developed inside a paper document processing project where GI results are used to classify and then translate documents into machine readable form. The generated logic program is augmented to handle some attributes of the logical structure of paper documents such as typographic attributes.

## VIII- References

- [1] S. Tayeb-Bey, A. S. Saïdi "Grammatical Formalism for Document Understanding System : From Document towards HTML Text". BSDIA97, November 1997, Brasilia.
- [2] E.M. Gold. "Language identification in the limit". Information and Control, 10(5)- 1967.
- [3] H.S. Fu and T. Booth: "Grammatical Inference: Introduction and Survey", parts 1 & 2. IEEE Trans. Sys. man and Cyber. SMC-5:95-11.
- [4] R. C. Gonzalez and M. G. Thomason. "Syntactic Pattern Recognition, an Introduction". Addison Wesley, Reading Mass. 1978.
- [5] H.S. Fu. "Syntactic Pattern Recognition and Applications". Prentice hall, N.Y. 1982.
- [6] L. Mictet. "Grammatical Inference". Syntactic and Structural Pattern Recognition. H. Bunk and San Feiui eds. World Scientific.
- [7] J. Onica, P. Garcia. "Inferring regular Languages in Polynomial Update time". Pattern Recognition and Image Analysis. 1992.
- [8] P. Dupont, L. Mictet & E. Vidal. "What is the search space of Regular Inference?". ICGI94, Grammatical Inference and Applications. Springer-Verlag-94.
- [9] L. Fribourg, M. V. Peixoto. "Automates concourants à Contraintes". TSI.13 (6). 1994.
- [10] J. A. Goguen, J.W. Thatcher, E.G. Wagner, J.B. Wright. "Initial Algebra Semantics and Continuous Algebra". JACM 24(1), 1977.
- [11] A. S. Saïdi : "Extensions Grammaticales de la Programmation Logique". PhD. 1992.
- [12] A. S. Saïdi. "On the unification of phrases". IFIP-94.
- [13] H. Ehrig, B. Mahr. "Fundamentals of Algebraic Specification", Vol-1 & 2. Springer-Verlag 1985.
- [14] E.M. Gold. "Complexity of automaton identification from given data". Information and Control, 37- 1978.
- [15] J.E. Hopcroft, J.D. Ullmann. "Formal Languages and their Relation to Automata". Addison-Wesley 1969.
- [16] F. Bancilhon & all. "Magic Sets and Other Strange Ways to Implement Logic Programs". Proc. ACM Symp. on principles of Databases Systems. Boston 1986.
- [17] F. Coste, J. Nicols : "Regular Inference as a graph coloring Problem". ICML97. 1997.
- [18] K.R. Apt, M.H. Van Emden : "Contribution to the Theory of Logic Programming". JACM. 29(3). 1982.
- [19] R.S. Michalski & all. "Machine Learning : An Artificial Intelligence Approach", vol. 1 & 2. Springer-Verlag 1984 and Morgan Kaufmann 1986.
- [20] H. Ahoen, H. Mannila. "Forming Grammars for structured documents". Research report. University of Helsinki. 1994.
- [21] P. Frankhauser, Y. Xu. "MarkitUp! an incremental approach to document structure recognition". Electronic Publishing-Organisation, Dissemination and Design, 6(4). 1994.