

LL(k) LANGUAGES ARE CLOSED UNDER UNION
WITH FINITE LANGUAGES

Ileana Streinu
Faculty of Mathematics
University of Bucharest
Str. Academiei 14
Bucuresti 7000 / Romania

1. Introduction

LL(k) grammars were first defined by Lewis and Stearns [3]. Algebraic properties, as closure under usual operations (union, intersection, complement, product, etc.) have been first studied by Rosenkrantz and Stearns [4]. They have proved that the class of LL languages is not closed under any of these operations.

Nevertheless, closure properties have been investigated. Wood [5] has proved that the LL languages are not closed under pre-product with finite or regular languages. Now we shall prove that they are closed under union and difference with finite languages. Moreover, the union (or difference) of an LL(k) language with a finite one does not modify the initial "k" : the resulting language still remains LL(k).

2. Previous definitions and results

In order to prove the results announced in the introduction, we shall use the following notations, definitions and theorems from Aho and Ullman [1] :

(2.1) Let $G = (N, \Sigma, P, S)$ be a context-free grammar. For $\alpha \in (N \cup \Sigma)^*$ and an integer $k > 0$, we define :

$FIRST_k^G(\alpha) = \{ w \in \Sigma^* / \text{either } |w| < k \text{ and } \alpha \xrightarrow{G} w, \text{ or } |w| = k \text{ and } \alpha \xrightarrow{G} wx, \text{ for}$

some $x \in \Sigma^* \}$. If $\alpha \in \Sigma^*$, the definition is independent of G .

For $L \subset \Sigma^*$, $FIRST_k(L) = \{ FIRST_k(w) / w \in L \}$

(2.2) A context-free grammar, $G = (N, \Sigma, P, S)$, is LL(k), for some inte-

ger $k > 0$, if whenever there are two leftmost derivations:

$$S \xrightarrow{\text{lm}}^* wA\alpha \xrightarrow{\text{lm}} w\beta\alpha \xrightarrow{*} wx$$

$$S \xrightarrow{\text{lm}}^* wA\alpha \xrightarrow{\text{lm}} w\gamma\alpha \xrightarrow{*} wy$$

such that $\text{FIRST}_k(x) = \text{FIRST}_k(y)$, it follows that $\beta = \gamma$.

We say that: a grammar is LL, if it is LL(k) for some integer k; a language is LL(k), if there is an LL(k) grammar G, such that $L = L(G)$.

(2.3) Let $G = (N, \Sigma, P, S)$ be a context-free grammar. Then G is LL(k) if and only if the following condition holds: if $A \rightarrow \beta$ and $A \rightarrow \gamma$ are distinct productions in P, then $\text{FIRST}_k(\beta\alpha) \cap \text{FIRST}_k(\gamma\alpha) = \emptyset$, for all $wA\alpha$ such that $S \xrightarrow{\text{lm}}^* wA\alpha$.

(2.4) Every LL(k) grammar has an equivalent LL(k+1) grammar in Greibach normal form.

(2.5) If a language L has an LL(k) grammar without e-productions, $k \geq 2$, then L has an LL(k-1) grammar.

3. Main results

(3.1) Lemma. Let L_1 and L_2 be LL languages, such that $L_1 \cap L_2 = \emptyset$. A sufficient condition for $L_1 \cup L_2$ to be LL is the existence of an integer $k > 0$ such that:

$$(3.1.1) \text{FIRST}_k(L_1) \cap \text{FIRST}_k(L_2) = \emptyset.$$

The condition is not necessary.

Proof. Let $G_i = (N_i, \Sigma, P_i, S_i)$ be LL(k_i) grammars, such that $L_i = L(G_i)$ ($i=1,2$). Without loss of generality, we can assume that $N_1 \cap N_2 = \emptyset$. We shall define a new grammar $G = (N, \Sigma, P, S)$ such that: $N = N_1 \cup N_2 \cup \{S\}$, where S is a new symbol, $S \notin N_1 \cup N_2$; $P = P_1 \cup P_2 \cup \{S \rightarrow S_1, S \rightarrow S_2\}$.

It is easy to verify that $L(G) = L_1 \cup L_2$.

Now we shall prove that G is an LL(n) grammar, where $n = \max(k, k_1, k_2)$. First of all, we shall remark that the new initial symbol S does not belong to the right side of any production from G (the only productions where S appears are $S \rightarrow S_1$ and $S \rightarrow S_2$). But, because $N_1 \cap N_2 = \emptyset$, it is true that:

$$\text{FIRST}_k(L_i) = \text{FIRST}_k^{G_i}(S_i) = \text{FIRST}_k^G(S_i), \quad i=1,2.$$

Then, from the hypothesis (3.1.1) it follows that:

(3.1.2) $\text{FIRST}_k^G(S_1) \cap \text{FIRST}_k^G(S_2) = \emptyset$, and this is also true for n, because $k \leq n$.

Moreover, from (2.3) it follows that for any $a \in N_i$ ($i=1,2$), the following conditions hold, for $i=1,2$: if $A \rightarrow \beta$ and $A \rightarrow \gamma$ are distinct

productions in P_i (then, also in P), then:

$$\text{FIRST}_{k_i}^{G_i}(\beta\alpha) \cap \text{FIRST}_{k_i}^{G_i}(\gamma\alpha) = \phi, \text{ for all } w\alpha \text{ such that } S_i \xrightarrow{*} w\alpha.$$

But, because $k_i \leq n$, $P_i \subseteq P$, $N_1 \cap N_2 = \phi$ and $S \rightarrow S_1$, $S \rightarrow S_2$ are the only S -productions in P , it follows that :

$$(3.1.3) \text{ FIRST}_n^G(\beta\alpha) \cap \text{FIRST}_n^G(\gamma\alpha) = \phi, \text{ for all } w\alpha \text{ such that } S \xrightarrow{*} w\alpha.$$

From (3.1.2) and (3.1.3), using (2.3), we conclude that G is $LL(n)$.

The condition (3.1.1) is not necessary, because $L_1 = \{a^n b^n / n \text{ odd}\}$ and $L_2 = \{a^n b^n / n \text{ even}\}$ are LL languages, generated by the $LL(1)$ grammars: $G_1: S \rightarrow aA$, $A \rightarrow aSbb / b$; $G_2: S \rightarrow aaSbb / e$. The condition (3.1.1) is not satisfied, but $L_1 \cup L_2 = \{a^n b^n / n \geq 0\}$ is an $LL(1)$ language.

(3.2) Corollary. If L is an LL language and w is a word, then $L\{w\}$ is also LL .

Proof. If $w \in L$, the statement is trivial. If $w \notin L$, the result follows from the preceding lemma, taking $k = |w| + 1$ (where $|w| =$ the length of w).

Remark. The proof of Lemma (3.1) shows that, if L is $LL(k)$ and $|w| = n$, then $L\{w\}$ has an $LL(m)$ grammar, where $m = \max(k, n+1)$. The following theorem will prove a stronger result, namely that $L\{w\}$ is still an $LL(k)$ language, if L is so.

(3.3) Theorem. Let L be an $LL(k)$ language, $w \notin L$. Then there exists an $LL(k)$ grammar which generates $L\{w\}$.

Proof. From (2.4) it follows that we can find an $LL(k+1)$ grammar in Greibach normal form, $G = (N, \Sigma, P, S)$, such that $L = L(G)$. We shall prove that $L\{w\}$ can be generated by an e -free $LL(k+1)$ grammar, and then, using (2.5), we may conclude that $L\{w\}$ can be generated by an $LL(k)$ grammar.

In fact we shall no more complicate the notation and we'll prove that : if $G = (N, \Sigma, P, S)$ is an $LL(k)$ grammar in Greibach normal form ($k \geq 2$), then there exists an e -free $LL(k)$ grammar which generates $L(G)\{w\}$.

Let $w = a_1 a_2 \dots a_m$. Let $n \geq 0$ be an integer such that for at least one word in L , $a_1 a_2 \dots a_n$ is a prefix, and for no word in L $a_1 a_2 \dots a_{n+1}$ is a prefix. The case $n = m$ is when w is a prefix for some word in L (in this case we suppose $a_{n+1} = e$). Note that here we use the fact that $w \notin L$.

If $n < k$, the statement follows from the proof of Lemma (3.1). So we shall suppose that $n \geq k$. Let us consider in G a leftmost derivation of the form:

$$S \xrightarrow{p_1} a_1 \alpha_1 \xrightarrow{p_2} a_1 a_2 \alpha_2 \xrightarrow{p_3} \dots \xrightarrow{p_{n-k+1}} a_1 a_2 \dots a_{n-k+1} \alpha_{n-k+1}$$

such that: $a_{i+1} \dots a_{i+k} \in \text{FIRST}_k(\alpha_i)$, $1 \leq i \leq n-k$, $a_{n-k+2} \dots a_{n+1} \notin \text{FIRST}_k(\alpha_{n-k+1})$.

Such a leftmost derivation can be found because of the choice of n , and it is even unique, because G is an LL(k) grammar. Additionally, if we set $\alpha_0 = S$ and if for a value of i , $0 \leq i \leq n-k$, we write the sequence of nonterminals α_i as $\alpha_i = A\beta_i$ ($A \in N$, $\beta_i \in N^*$), then for any rule $A \rightarrow \beta$ which is not the rule P_{i+1} , it follows that:

$$(3.3.1) \text{FIRST}_k^G(\beta\beta_i) \cap \text{FIRST}_k^G(a_{i+1}\alpha_{i+1}) = \emptyset, \quad 0 \leq i \leq n-k.$$

Let us now modify the initial grammar $G = (N, \Sigma, P, S)$ in order to obtain a new LL(k) grammar, $G' = (N', \Sigma, P', X_0)$ such that $L(G') = L \cup \{w\}$. Let $X_0, X_1, \dots, X_{n-k+1}$ be new symbols not in $N \cup \Sigma$, all different. Let $N' = N \cup \{X_0, X_1, \dots, X_{n-k+1}\}$, where X_0 is the new initial symbol. P' is the set obtained by joining to P the following new productions:

$$(3.3.2) X_i \rightarrow \beta\beta_i, \text{ if } \alpha_i = A\beta_i \text{ and } A \rightarrow \beta \text{ is a rule in } P, \text{ but not the rule } P_{i+1}.$$

$$(3.3.3) X_i \rightarrow a_{i+1}X_{i+1}$$

for all $i=0, 1, \dots, n-k$, and

$$(3.3.4) X_{n-k+1} \rightarrow \alpha_{n-k+1}$$

$$(3.3.5) X_{n-k+1} \rightarrow a_{n-k+2} \dots a_m.$$

Note that for $k \geq 2$ this is an ϵ -free grammar.

Let us now prove that G' is LL(k). We shall use theorem (2.3). Let $A \in N'$, $A \rightarrow \beta$, $A \rightarrow \gamma$ be two different productions from P' , and $\alpha \in N'^*$ such that $X_0 \xrightarrow{\text{IM}} xA\alpha$, $x \in \Sigma^*$. We must prove that $\text{FIRST}_k^{G'}(\beta\alpha) \cap \text{FIRST}_k^{G'}(\gamma\alpha) = \emptyset$. If $A \in N$, it is easy to remark that $\alpha \in N^*$, and $A \rightarrow \beta$, $A \rightarrow \gamma$ are from P . Then $\text{FIRST}_k^{G'}(\beta\alpha) = \text{FIRST}_k^G(\beta\alpha)$ and $\text{FIRST}_k^{G'}(\gamma\alpha) = \text{FIRST}_k^G(\gamma\alpha)$, which are disjoint, because G is LL(k). If $A = X_i$ ($0 \leq i \leq n-k+1$), then $\alpha = \epsilon$. Then $\text{FIRST}_k^{G'}(\beta) \cap \text{FIRST}_k^{G'}(\gamma) = \emptyset$, because we have the statement (3.3.1) true, if $A \rightarrow \beta$ or $A \rightarrow \gamma$ are (3.3.3) or (3.3.5) rules, and because G is LL(k), if both are (3.3.2) rules. So it follows that G' is LL(k).

Let us now prove that $L(G') = L \cup \{w\}$. In order to prove that $L \cup \{w\} \subset L(G')$, we shall remark that any word from L which has no prefix in common with w may be derived in G' by using at first a (3.3.2)-rule, and then rules from P . Any word in L having $a_1 a_2 \dots a_i$ as a prefix may be derived in G' in the same way as before, if $i < k$, and otherwise by using $(i-k+1)$ -times (3.3.3)-rules, then a (3.3.2) or a (3.3.4)-rule, and finally using rules from P . At last, w may be derived by using all the (3.3.3)-rules and the (3.3.5)-rule. In order to show the other inclusion, it may be proved, by induction on the number of steps used in

the derivation $X_i \xrightarrow{1m}^* x$, that $x \in \{y \in Z^* / a_1 a_2 \dots a_i y \in L \cup \{w\}\}$ ($0 \leq i \leq n-k+1$).

So, for $i=0$, $X_0 \xrightarrow{*} x$ implies that $x \in L \cup \{w\}$. The proof by induction is not difficult and it is left to the reader. So, we may conclude that the LL(k) e-free grammar G' generates $L \cup \{w\}$, and the theorem is thereby proved.

(3.4) Corollary. If L_1 is LL(k) and L_2 is finite, then $L_1 \cup L_2$ is LL(k).

(3.5) Corollary. Lemma (3.1) is valid even if $L_1 \cap L_2$ is finite.

Remark. Using an analogous proof to that of theorem (3.3), one might show that, if L is LL(k) and w is a word in L , then $L - \{w\}$ can be generated by an LL(k) grammar. An important consequence of these facts is the following:

(3.6) Theorem. If L is LL(k) but not LL(k-1), then $L \cup \{w\}$ and $L - \{w\}$ are also LL(k), but not LL(k-1).

(3.7) Corollary. The class of LL(k) languages is closed under union and difference with finite languages.

Remark. Examining the proof of Lemma (3.1), we conclude that $L_1 \cup L_2$ is LL(n), where $n = \max(k, k_1, k_2)$. This value can be improved, remarking that we need k only for expanding S : here we must make a choice between $S \rightarrow S_1$ and $S \rightarrow S_2$. Using a method of "linearization" of a derivation analogous to that from the proof of theorem (3.3), one might obtain an LL(n) grammar to generate $L_1 \cup L_2$, where $n = \max(k_1, k_2)$.

Bibliografy

- [1] A.V.Aho, J.D.Ullman: The Theory of Parsing, Translation and Compiling, vol.1(1972), vol.2(1973), Prentice-Hall, Englewood Cliffs.
- [2] J.Bordier, H.Saya: A necessary and sufficient condition for a power language to be LL(k), Computer Journal, 1973, vol.16, pp.451-456.
- [3] P.M.Lewis, R.E.Stearns: Syntax-directed transduction, Journal ACM, 1968, vol.15:3, pp.464-488.
- [4] D.J.Rosenkrantz, R.E.Stearns: Properties of deterministic top-down grammars, Information and Control, 1970, vol.17:3, pp.226-256.
- [5] D.Wood: A further note on TD deterministic languages, Computer Journal, 1971, vol.14, pp.396-403.