

A recursive ascent Earley parser

René Leermakers

Instituut voor Perceptie Onderzoek,
P.O. Box 513, 5600 MB Eindhoven, The Netherlands
E-mail:leermake@prl.philips.nl

1 Introduction

Recently, the theory of LR-parsing gained new impetus by the discovery of the recursive ascent implementation technique for deterministic [2] and nondeterministic [3,4] LR-parsers. In short, the novelty is that LR-parsers can be implemented purely functionally and that this implementation has very simple correctness proofs. In its primary form, a recursive ascent parser consists of two functions for each state.

One of the major application areas of the Earley algorithm [1] is the field of computational linguistics. A few years ago, Tomita [5] has proposed a nondeterministic LR-parser to parse natural languages. This parser has stirred some research to establish the differences between the Earley and Tomita parsers. One of the results of this research was the insight that the Earley algorithm can be seen as a shift-reduce parser and as such fits in a large family of parsers to which also nondeterministic LR-parsers belong. A consequence of this insight is that implementation techniques for LR-parsers can also be applied to the Earley algorithm. In particular, this can be done with the recursive ascent technique, and this is the subject of this note.

2 The algorithm

Consider CF grammar G , with terminals V_T and nonterminals V_N . Let $V = V_N \cup V_T$. Recursive ascent parsers consist of a number of functions. Let us start with associating a function to each item. Items, grammar rules with a dot somewhere in the right hand side, are used ubiquitously in parsing theory to denote a partially recognized rule. A typical item is

written as $A \rightarrow \alpha.\beta$, with greek letters for arbitrary elements of V^* . We use the unorthodox embracing operator $[\cdot]$ to map each item to its function:

$$[A \rightarrow \alpha.\beta] : N \mapsto 2^N$$

where N is the set of integers, or a subset $0 \dots n_{max}$, with n_{max} the maximum sentence length, and 2^N is the powerset of N . The functions are to meet the following specification:

$$[A \rightarrow \alpha.\beta](i) = \{j | \beta \xrightarrow{*} x_{i+1} \dots x_j\},$$

with $x_1 \dots x_n$ the string to be parsed. So, the function reports which parts of the string can be derived from β starting from position i . Below we will find a constructive definition for this function that can be viewed as a functional implementation. If we add a grammar rule $S' \rightarrow S$ to G , with $S' \notin V$ then $S \xrightarrow{*} x_1 \dots x_n$ is equivalent to $n \in [S' \rightarrow .S](0)$, so that the algorithm is to be invoked by calling $[S' \rightarrow .S](0)$.

To be able to construct an implementation of $[A \rightarrow \alpha.\beta]$ that has no problems with left-recursive grammars, we need so-called predict sets. Let $predict(A \rightarrow \alpha.\beta)$ be the set of initial items, that are derived from $A \rightarrow \alpha.\beta$ by the closure operation:

$$predict(A \rightarrow \alpha.\beta) = \{B \rightarrow .\mu | B \rightarrow \mu \wedge \beta \xrightarrow{*} B\gamma\}.$$

The double arrow \Rightarrow denotes a left-most-symbol rewriting with a non- ϵ grammar rule, i.e.

$$\alpha \Rightarrow \beta \equiv \exists B\gamma\delta (\alpha = B\gamma \wedge \beta = \delta\gamma \wedge B \rightarrow \delta \wedge \delta \neq \epsilon).$$

A recursive ascent recognizer may be obtained by relating to each item $A \rightarrow \alpha.\beta$ not only the above $[A \rightarrow \alpha.\beta]$, but also a function that we take to be the result of applying operator $\overline{[\cdot]}$ to the item:

$$\overline{[A \rightarrow \alpha.\beta]} : V \times N \mapsto 2^N$$

It has the specification ($X \in V$)

$$\overline{[A \rightarrow \alpha.\beta]}(X, i) = \{j | \exists \gamma : \beta \xrightarrow{*} X\gamma \wedge \gamma \xrightarrow{*} x_{i+1} \dots x_j\}.$$

Assuming x_{n+1} to be some end of sentence marker that is not in V , it can never be that $\beta \xrightarrow{*} x_{n+1}\gamma$, hence $\overline{[A \rightarrow \alpha.\beta]}(x_{n+1}, n+1) = \emptyset$. For $i \leq n$ the above functions are recursively implemented by

$$\begin{aligned} [A \rightarrow \alpha.\beta](i) &= \overline{[A \rightarrow \alpha.\beta]}(x_{i+1}, i+1) \cup \\ &\quad \{j | \exists B : B \rightarrow .\epsilon \in \text{predict}(A \rightarrow \alpha.\beta) \wedge j \in \overline{[A \rightarrow \alpha.\beta]}(B, i)\} \cup \\ &\quad \{i | \beta = \epsilon\} \end{aligned}$$

$$\begin{aligned} \overline{[A \rightarrow \alpha.\beta]}(X, i) &= \{j | \exists \gamma : \beta = X\gamma \wedge j \in [A \rightarrow \alpha X.\gamma](i)\} \cup \\ &\quad \{j | \exists C\delta k : j \in \overline{[A \rightarrow \alpha.\beta]}(C, k) \wedge C \rightarrow .X\delta \in \text{predict}(A \rightarrow \alpha.\beta) \wedge \\ &\quad k \in [C \rightarrow X.\delta](i)\} \end{aligned}$$

The correctness of this implementation is shown by the following proof:

First we notice that

$$\begin{aligned} \beta \xrightarrow{*} x_{i+1} \dots x_j &\equiv \exists \gamma (\beta \xrightarrow{*} x_{i+1} \gamma \wedge \gamma \xrightarrow{*} x_{i+2} \dots x_j) \vee \\ &\quad \exists B \gamma (\beta \xrightarrow{*} B \gamma \wedge B \rightarrow \epsilon \wedge \gamma \xrightarrow{*} x_{i+1} \dots x_j) \vee \\ &\quad (\beta = \epsilon \wedge i = j). \end{aligned}$$

Substituting this in the specification of $[A \rightarrow \alpha.\beta]$ one gets

$$\begin{aligned} [A \rightarrow \alpha.\beta](i) &= \{j | \exists \gamma : \beta \xrightarrow{*} x_{i+1} \gamma \wedge \gamma \xrightarrow{*} x_{i+2} \dots x_j\} \cup \\ &\quad \{j | \exists B \gamma : B \rightarrow \epsilon \wedge \beta \xrightarrow{*} B \gamma \wedge \gamma \xrightarrow{*} x_{i+1} \dots x_j\} \cup \\ &\quad \{j | \beta = \epsilon \wedge i = j\}. \end{aligned}$$

This directly leads to the implementation given above, using the specifications of $\overline{[A \rightarrow \alpha.\beta]}(x_{i+1}, i+1)$, $\overline{[A \rightarrow \alpha.\beta]}(B, i)$ and $\text{predict}(A \rightarrow \alpha.\beta)$.

For establishing the correctness of $\overline{[A \rightarrow \alpha.\beta]}$ notice that $\beta \xrightarrow{*} X\gamma$ either consists of zero steps, in which case $\beta = X\gamma$, or it contains at least one step:

$$\begin{aligned} \exists \gamma (\beta \xrightarrow{*} X\gamma \wedge \gamma \xrightarrow{*} x_{i+1} \dots x_j) &\equiv \exists \gamma (\beta = X\gamma \wedge \gamma \xrightarrow{*} x_{i+1} \dots x_j) \vee \\ &\quad \exists C\delta \gamma k (\beta \xrightarrow{*} C\gamma \wedge C \rightarrow X\delta \wedge \delta \xrightarrow{*} x_{i+1} \dots x_k \wedge \gamma \xrightarrow{*} x_{k+1} \dots x_j) \end{aligned}$$

Hence, $\overline{[A \rightarrow \alpha.\beta]}(X, i)$ may be written as the union of two sets, S_0 and S_1 :

$$S_0 = \{j | \exists \gamma : \beta = X\gamma \wedge \gamma \xrightarrow{*} x_{i+1} \dots x_j\}$$

$$S_1 = \{j | \exists C\delta \gamma k : \beta \xrightarrow{*} C\gamma \wedge C \rightarrow X\delta \wedge \delta \xrightarrow{*} x_{i+1} \dots x_k \wedge \gamma \xrightarrow{*} x_{k+1} \dots x_j\}.$$

With the specification of $[A \rightarrow \alpha X.\gamma]$, S_0 may be rewritten as

$$S_0 = \{j | \exists \gamma : \beta = X\gamma \wedge j \in [A \rightarrow \alpha X.\gamma](i)\}.$$

The set S_1 may be rewritten using the specifications of $\overline{[A \rightarrow \alpha.\beta]}(C, k)$ and $\text{predict}(A \rightarrow \alpha.\beta)$:

$$\begin{aligned} S_1 &= \{j | \exists C\delta k : j \in \overline{[A \rightarrow \alpha.\beta]}(C, k) \wedge \\ &\quad C \rightarrow .X\delta \in \text{predict}(A \rightarrow \alpha.\beta) \wedge \delta \xrightarrow{*} x_{i+1} \dots x_k\}. \end{aligned}$$

With the definition of $[C \rightarrow X.\delta]$ one finally gets:

$$S_1 = \{j | \exists C\delta k : j \in \overline{[A \rightarrow \alpha.\beta]}(C, k) \wedge \\ C \rightarrow .X\delta \in \text{predict}(A \rightarrow \alpha.\beta) \wedge k \in [C \rightarrow X.\delta](i)\}.$$

□

3 Complexity and Variants

The above recognizer needs exponential time resources unless the functions are implemented as memo-functions. Memo-functions memorize for which arguments they have been called. If a function is called with the same arguments as before, the function returns the previous result without recomputing it. In conventional programming languages memo-functions are not available, but they can easily be implemented. The use of memo-functions obsoletes the introduction of devices like parse matrices [1]. The worst-case complexity analysis of the memo-ized recognizer is quite simple. Let n be the sentence length, $|G|$ the number of items of the grammar, p the maximum number of different left hand sides in a predict set (bounded by the number of nonterminals), and q the maximum number of items in a predict set with the same symbol after the dot. Then, there are $O(|G|pn)$ different invocations of recognizer functions. Each invocation of a function $\overline{[I]}$ invokes $O(qn)$ other functions, that all result in a set with $O(n)$ elements. The merge of these sets into one set with no duplicates can be accomplished in $O(qn^2)$ time on a random access machine. Hence, the total time-complexity is $O(|G|pqn^3)$. The space needed for storing function results is $O(n)$ per invocation, i.e. $O(|G|pn^2)$ for the whole recognizer. These complexity results are almost identical to the usual ones for Earley parsing. Only the dependence on the grammar variables $|G|$, p and q slightly differs. Worst-case complexities need not be relevant in practice. We claim that for many practical grammars the present algorithm is more efficient than existing implementations, for the following reason. The above formulae can be interpreted as to define two functions $[\cdot]$ and $\overline{[\cdot]}$, that will work for variable grammars and strings. This view is convenient when building prototypes. If efficiency is an issue, however, one should precompute as much as possible and actually create, for a fixed grammar, the functions $[I]$ and $\overline{[I]}$ for every item I . In the terminology of functional programming the functions $[\cdot]$ and $\overline{[\cdot]}$ are to be evaluated partially for each item. In this way, the grammar is compiled into a collection of functions, just like conventional parser generators compile a grammar into LR-tables or a recursive descent parser. Quite

some work that is done at parse time by the standard Earley parser, such as the creation of predict sets and the processing of item sets, is transferred to compile time when transforming the grammar into a functional parser. As a consequence, the compiled parser is more efficient than the standard implementations of the Earley parser.

The above considerations only hold if our algorithm terminates. If the grammar has a cyclic derivation $A \xrightarrow{+} A$, the execution of $\overline{[I]}(A, i)$ leads to a call of itself, and the algorithm does not terminate. Also, there may be a cycle of transitions labeled by nonterminals that derive ϵ . This occurs if for some k one has that for $i = 1 \dots k$

$$A_{i+1} \rightarrow \cdot \alpha_{i+1} \beta_{i+1} \in \text{predict}(A_i \rightarrow \alpha_i \cdot \beta_i) \wedge \alpha_i \xrightarrow{+} \epsilon$$

$$\text{while } A_1 = A_{k+1} \wedge \alpha_1 = \alpha_{k+1} \wedge \beta_1 = \beta_{k+1}.$$

Then the execution of $[A_1 \rightarrow \alpha_1 \cdot \beta_1](i)$ leads to a call of itself, and the algorithm does not terminate. A cycle of this form occurs *iff* there is a derivation $A \xrightarrow{+} \alpha A \beta$ such that $\alpha \xrightarrow{+} \epsilon$. It is easy, however, to define a variant of the recognizer that has no problems with these derivations. It is obtained from dropping the restriction that the left most symbol derivation \Rightarrow may not use ϵ -rules. Then this paper's analysis can be repeated. The redefinition of \Rightarrow affects the function *predict* and the way $\beta \xrightarrow{*} x_{i+1} \dots x_j$ is to be expressed in terms of \Rightarrow :

$$\beta \xrightarrow{*} x_{i+1} \dots x_j \equiv \exists \gamma (\beta \xrightarrow{*} x_{i+1} \gamma \wedge \gamma \xrightarrow{*} x_{i+2} \dots x_j) \vee (\beta \xrightarrow{*} \epsilon \wedge i = j).$$

Also the decomposition of $\beta \xrightarrow{*} X \gamma$ must be rephrased: X is introduced by a grammar rule $C \rightarrow \mu X \delta$ or it is not. The result is the following recognizer:

$$[A \rightarrow \alpha \cdot \beta](i) = \overline{[A \rightarrow \alpha \cdot \beta]}(x_{i+1}, i+1) \cup \{i | \beta \xrightarrow{*} \epsilon\}$$

$$\begin{aligned} \overline{[A \rightarrow \alpha \cdot \beta]}(X, i) = & \{j | \exists \gamma \mu : \beta = \mu X \gamma \wedge \mu \xrightarrow{*} \epsilon \wedge j \in [A \rightarrow \alpha \mu X \cdot \gamma](i)\} \cup \\ & \{j | \exists C \delta \mu k : j \in \overline{[A \rightarrow \alpha \cdot \beta]}(C, k) \wedge C \rightarrow \cdot \mu X \delta \in \text{predict}(A \rightarrow \alpha \cdot \beta) \wedge \\ & \mu \xrightarrow{*} \epsilon \wedge k \in [C \rightarrow \mu X \cdot \delta](i)\} \end{aligned}$$

To conclude, the above offers attractive alternatives to the standard Earley parser. Only for cyclic grammars there is a problem with termination. For other grammars one gains efficiency. The above recognition algorithms are not the most efficient ones, however. For instance, there is an additional improvement following from the observation that $[A \rightarrow \alpha \cdot \beta]$ and $\overline{[A \rightarrow \alpha \cdot \beta]}$ only depend on β . Therefore, functions for different items may be identical

and can be identified. A similar improvement is applicable to many implementations of the Earley algorithm. A recursive ascent recognizer with this improvement is

$$\begin{aligned}
 [\beta](i) &= \overline{[\beta]}(x_{i+1}, i+1) \cup \{i | \beta \xrightarrow{*} \epsilon\} \\
 \overline{[\beta]}(X, i) &= \{j | \exists \gamma \mu : \beta = \mu X \gamma \wedge \mu \xrightarrow{*} \epsilon \wedge j \in [\gamma](i)\} \cup \\
 &\quad \{j | \exists C \delta \mu k : j \in \overline{[\beta]}(C, k) \wedge C \rightarrow \cdot \mu X \delta \in \text{predict}(\beta) \wedge \mu \xrightarrow{*} \epsilon \wedge k \in [\delta](i)\}
 \end{aligned}$$

with $\text{predict}(\beta) = \{B \rightarrow \cdot \mu | B \rightarrow \mu \wedge \beta \xrightarrow{*} B \gamma\}$. The algorithm is to be invoked by calling $[S](0)$.

References

- 1 J.C. Earley, An efficient context-free parsing algorithm, *Commun. ACM* (1970) 13(2):94-102.
- 2 F.E.J. Kruseman Aretz, On a recursive ascent parser, *Information Processing Letters* (1988) 29:201-206.
- 3 R. Leermakers, Non-deterministic Recursive Ascent Parsing, Proceedings of the 5th Conference of the European Chapter of the Association of Computational Linguistics, 1991.
- 4 R. Leermakers, L. Augusteijn and F.E.J. Kruseman Aretz, A functional LR-parser, *Theoretical Computer Science* (accepted).
- 5 M. Tomita, Efficient Parsing for Natural Language, A Fast Algorithm for Practical Systems (Kluwer Academic Publishers, 1986).