

# Teoría de Lenguajes

## Segunda parte

## Algoritmos de Parsing de Earley

Verónica Becher

Segundo cuatrimestre 2020

Bibliografía para esta clase:

A. V. Aho, J. D. Ullman, The Theory of Parsing, Translation, and Compiling, Vol. 1 , Parsing. Prentice Hall, 1972.

<https://www-2.dc.uba.ar/staff/becher/Aho-Ullman-Parsing-V1.pdf> Capítulo 4.2

### **Algoritmo de parsing de Earley**

Dada una gramática libre de contexto el algoritmo de parsing de Earley con una entrada de longitud  $n$  hace una cantidad de operaciones en el orden de  $n^3$  y usa espacio  $n^2$ .

En caso de que la gramática sea no-ambigua, la cantidad de operaciones es del orden  $n^2$ .

Para la mayoría de las gramáticas que definen lenguajes de programación el algoritmo de Earley se puede modificar para obtener cantidad de operaciones y espacio lineales respecto a la longitud de la entrada.

Es un método bottom-up que encuentra la derivación más a la derecha.

### Algoritmo de Parsing de Earley

Sea  $G = (N, T, P, S)$  una GLC y sea  $w = a_1 a_2 \dots a_n$  una cadena de input en  $T^*$ .

**Definición** (Item de Earley). *Decimos que  $[A \rightarrow X_1 X_2 \dots X_k \cdot X_{k+1} \dots X_m, i]$  es un item para  $w$  si  $A \rightarrow X_1 X_2 \dots X_m$  es una producción en  $P$  y  $0 \leq i \leq n$ .*

Si la producción es  $A \rightarrow \lambda$ , el item es  $[A \rightarrow \cdot, i]$ .

### Algoritmo de Parsing de Earley

Dada  $G = (N, T, P, S)$  y una cadena de entrada  $a_1 \dots a_n$  en  $T^*$  construimos las listas de items  $\ell_1, \ell_2, \dots, \ell_n$  tal que

$$[A \rightarrow \alpha \cdot \beta, i] \in \ell_j$$

si y solo si para algún  $\gamma$  y  $\delta$

$$S \xRightarrow{*} \gamma A \delta, \quad \gamma \xRightarrow{*} a_1 \dots a_i \quad \text{y} \quad \alpha \xRightarrow{*} a_{i+1} \dots a_j,$$

Luego,  $w \in L(G)$  si y solo si hay  $\alpha$  tal que  $[S \rightarrow \alpha \cdot, 0] \in \ell_n$ .

### Algoritmo de Earley

Algorithm 4.5 Aho Ullman vol 1

*Input.* Gramática GLC  $G = (N, T, P, S)$  y cadena  $w = a_1 \dots a_n$  en  $T^*$ .

*Output.* Las listas  $\ell_0, \ell_1, \dots, \ell_n$ .

Paso 1 Si  $S \rightarrow \alpha$  es una producción en  $P$  agregar a  $\ell_0$ ,  $[S \rightarrow \cdot \alpha, 0]$ .

Paso 2 Si  $[A \rightarrow \alpha \cdot B \beta, 0] \in \ell_0$  y  $[B \rightarrow \gamma \cdot, 0] \in \ell_0$  (en particular  $\gamma$  puede ser  $\lambda$ ) entonces agregar a  $\ell_0$ ,  $[A \rightarrow \alpha B \cdot \beta, 0]$ .

Paso 3 Si  $[A \rightarrow \alpha \cdot B \beta, 0] \in \ell_0$ , agregar a  $\ell_0$  para toda  $B \rightarrow \gamma$  en  $P$ , (si es que aún no está),  $[B \rightarrow \cdot \gamma, 0]$ .

Paso 4 Supongamos que ya hemos construido  $\ell_0, \ell_1, \dots, \ell_{j-1}$ . Si  $[B \rightarrow \alpha \cdot a B \beta, i] \in \ell_{j-1}$  tal que  $a = a_j$ , agregar a  $\ell_j$   $[B \rightarrow \alpha a \cdot \beta, i]$

Paso 5 Si  $[B \rightarrow \alpha \cdot A \beta, k] \in \ell_i$  y  $[A \rightarrow \gamma \cdot, i] \in \ell_j$ , agregar a  $\ell_j$   $[B \rightarrow \alpha A \cdot \beta, k]$ .

Paso 6 Si  $[A \rightarrow \alpha \cdot B \beta, i] \in \ell_j$ , agregar a  $\ell_j$ , para todo  $B \rightarrow \gamma$  en  $P$ ,  $[B \rightarrow \cdot \gamma, j]$ .

Notemos que la consideración de un item con un terminal a la derecha del  $\cdot$  no produce nuevos items en los pasos 2, 3, 5 y 6.

### Ejemplo

	$I_0$ $\begin{aligned} [E \rightarrow \cdot T + E, 0] \\ [E \rightarrow \cdot T, 0] \\ [T \rightarrow \cdot F * T, 0] \\ [T \rightarrow \cdot F, 0] \\ [F \rightarrow \cdot (E), 0] \\ [F \rightarrow \cdot a, 0] \end{aligned}$	$I_1$ $\begin{aligned} [F \rightarrow (\cdot E), 0] \\ [E \rightarrow \cdot T + E, 1] \\ [E \rightarrow \cdot T, 1] \\ [T \rightarrow \cdot F * T, 1] \\ [T \rightarrow \cdot F, 1] \\ [F \rightarrow \cdot (E), 1] \\ [F \rightarrow \cdot a, 1] \end{aligned}$	$I_2$ $\begin{aligned} [F \rightarrow a \cdot, 1] \\ [T \rightarrow F \cdot * T, 1] \\ [T \rightarrow F \cdot, 1] \\ [E \rightarrow T \cdot + E, 1] \\ [E \rightarrow T \cdot, 1] \\ [F \rightarrow (E \cdot), 0] \end{aligned}$
Sea $G$ con las producciones	$I_3$ $\begin{aligned} [E \rightarrow T + \cdot E, 1] \\ [E \rightarrow \cdot T + E, 3] \\ [E \rightarrow \cdot T, 3] \\ [T \rightarrow \cdot F * T, 3] \\ [T \rightarrow \cdot F, 3] \\ [F \rightarrow \cdot (E), 3] \\ [F \rightarrow \cdot a, 3] \end{aligned}$	$I_4$ $\begin{aligned} [F \rightarrow a \cdot, 3] \\ [T \rightarrow F \cdot * T, 3] \\ [T \rightarrow F \cdot, 3] \\ [E \rightarrow T \cdot + E, 3] \\ [E \rightarrow T \cdot, 3] \\ [E \rightarrow T + E \cdot, 1] \\ [F \rightarrow (E \cdot), 0] \end{aligned}$	$I_5$ $\begin{aligned} [F \rightarrow (E) \cdot, 0] \\ [T \rightarrow F \cdot * T, 0] \\ [T \rightarrow F \cdot, 0] \\ [E \rightarrow T \cdot + E, 0] \\ [E \rightarrow T \cdot, 0] \end{aligned}$
1. $E \Rightarrow T + E$			
2. $E \Rightarrow T$			
3. $T \Rightarrow F * T$			
4. $T \Rightarrow F$			
5. $F \Rightarrow (E)$			
6. $F \Rightarrow a$			
Input $(a + a) * a$ , de longitud 7	$I_6$ $\begin{aligned} [T \rightarrow F * \cdot T, 0] \\ [T \rightarrow \cdot F * T, 6] \\ [T \rightarrow \cdot F, 6] \\ [F \rightarrow \cdot (E), 6] \\ [F \rightarrow \cdot a, 6] \end{aligned}$	$I_7$ $\begin{aligned} [F \rightarrow a \cdot, 6] \\ [T \rightarrow F \cdot * T, 6] \\ [T \rightarrow F \cdot, 6] \\ [T \rightarrow F * T \cdot, 0] \\ [E \rightarrow T \cdot + E, 0] \\ [E \rightarrow T \cdot, 0] \end{aligned}$	
Damos $\ell_0, \dots, \ell_7$ .			
Dado que $E \rightarrow T \cdot, 0] \in \ell_7$ ,			
$(a + a) * a \in L(G)$			

### Correctitud de algoritmo de Earley

**Teorema 1** (Theorem 4.9 Aho Ullman vol 1). Dada  $G = (N, T, P, S)$  libre de contexto,  $[S \rightarrow \alpha, 0] \in \ell_n$  si y solo si  $S \rightarrow \alpha$  está en  $P$  y  $\alpha \xrightarrow{*}_R a_1 \dots a_n$ .

### Complejidad de algoritmo de Earley

**Teorema 2** (Theorem 4.10 en Aho Ullman vol 1). Si la gramática es no-ambigua, para una entrada de longitud  $n$  el algoritmo de Earley construye las listas  $\ell_1, \dots, \ell_n$  en una cantidad de operaciones del orden  $n^2$

Si la gramática es ambigua esta cantidad es del orden  $n^3$ .

Cada lista  $\ell_j$  tiene a lo sumo  $(j+1)k$  items, donde  $k$  es una constante. Mantenemos una lista encadenada entre todos ellos y una table indicando si una producción ya está o no.

Mantenemos además otra lista encadenada entre los items que tiene el mismo símbolo a la derecha del  $\cdot$ . Esto sirve para los pasos donde se requiere buscar los items que tienen un determinado símbolo a derecha del  $\cdot$ .

Cuando la gramática es no ambigua cada item se considera una sola vez. Y hay una cantidad constante  $c$  de operaciones a realizar para cada item. Los puntos 4,5 y 6 del algoritmo aseguran que la cantidad de operaciones es a lo sumo

$$c \sum_{j=0}^n (j+1)k = ck(n+1)(n+2)/2 = c'n^2$$

Los pasos 1,2,3 construyen  $\ell_0$  y la cantidad de operaciones está acotada por la cantidad de producciones. Concluimos que algoritmo realiza en el orden de  $n^2$  operaciones. En caso de que la gramática sea ambigua son  $n^3$ .

### Derivacion a derecha de Earley

Derivación a derecha del Algoritmo de Earley, Algorithm 4.6 Aho Ullman vol 1

*Input.* Una GLC  $G = (N, T, P, S)$  sin ciclos, con las producciones numeradas  $1, \dots, p$ , una cadena de input  $w = a_1 \dots a_n$  y las listas  $\ell_0, \dots, \ell_n$

*Output.* Derivación a derecha de  $w$ , o “error”.

*Método.* Si  $[S \rightarrow \alpha, 0]$  en  $\ell_n$  entonces  $w$  está en  $L(G)$ . Fijar variable global  $\pi = \lambda$ . Ejecutar  $R([S \rightarrow \alpha, 0], n)$   
Emitir secuencia de producciones  $\pi$ .

**Rutina**  $R([A \rightarrow X_1 \dots X_m, i], j)$

$\pi = h\pi$ , donde  $h$  es el número de producción  $A \rightarrow X_1 \dots X_m$   $k = m$   $j' = j$

Mientras ( $k > 0$ )

Si  $X_k \in T$  entonces

$k = k - 1$

$j' = j' - 1$

Sino (es decir,  $X_k \in N$ )

Encontrar un item  $[X_k \rightarrow \gamma, i']$  en  $\ell_{j'}$  para algún  $i'$  tal que  $[A \rightarrow X_1 X_2 \dots X_{k-1} \cdot X_k \dots X_m, i] \in \ell_{i'}$ .

Ejecuta  $R([X_k \rightarrow \gamma, i'], j')$ .

$k = k - 1$

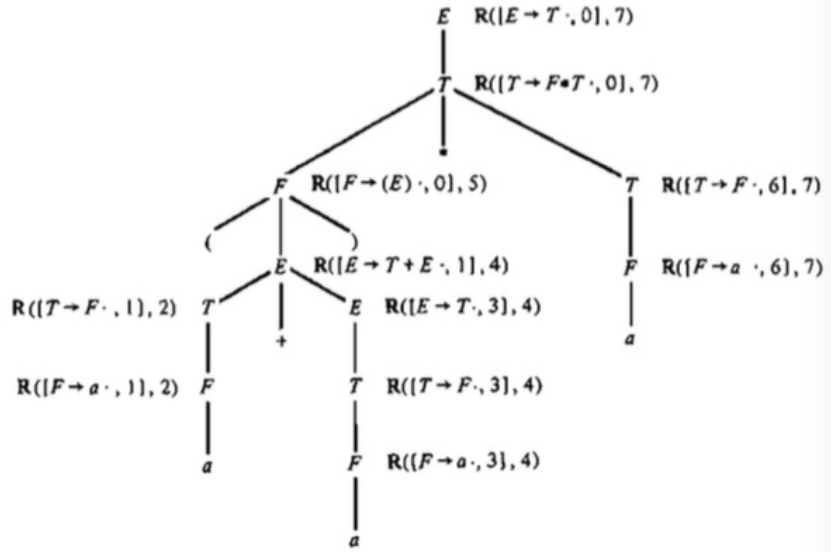
$j' = i'$ .

### Ejemplo

Sea  $G$  con las producciones

1.  $E \Rightarrow T + E$
2.  $E \Rightarrow T$
3.  $T \Rightarrow F * T$
4.  $T \Rightarrow F$
5.  $F \Rightarrow (E)$
6.  $F \Rightarrow a$

La figura muestra la rutina  $R$  sobre el árbol de derivación para  $(a + a) * a$ , de longitud 7, y  $\pi = 64642156432$ .



## Complejidad de algoritmo de Earley

**Teorema 3** (Theorem 4.12 en Aho Ullman vol 1). *Si existe una derivación más a la derecha de  $a_1 \dots a_n$  el algoritmo de Earley la encuentra en  $c n^2$  operaciones, para una constante  $c$ .*

## Complejidad de algoritmo de Earley

**Demostración.** Recordemos que  $[A \rightarrow \alpha \cdot \beta, i] \in \ell_j$  si y solo si  $\alpha \xrightarrow{*}_R a_{i+1} \dots a_j$  y existen  $\gamma$  y  $\delta$  tales que  $S \xrightarrow{*}_R \gamma A \delta \xrightarrow{*}_R a_1 \dots a_i$ .

En  $\ell_j$  unimos todos los items con la misma segunda componente.

Acceso y consulta en esta lista son operaciones elementales.

Este preprocesamiento se puede hacer en tiempo constante  $n^2$ .

Veamos que la ejecución de  $R([A \rightarrow \beta \cdot, i], j)$  lleva del orden de  $(j - i)^2$  operaciones. El ciclo realiza  $|\beta|$  iteraciones. En cada iteración lo costoso es buscar un valor  $i'$  examinando  $j - i + 1$  listas (cada exámen es elemental) y realizar una llamada recursiva  $R([A \rightarrow \beta \cdot, i'], j')$ , donde  $j' - i' < j - i$ .

La ejecución de  $R([A \rightarrow \beta \cdot, i], j)$  invoca a  $R([A \rightarrow \beta \cdot, i'], j')$  a lo sumo  $|P| \times M$  veces (cada vez con un valor distinto de  $--$ ) donde  $M$  es la máxima cantidad d esímbolso de una producción.

El total de invocaciones a  $R$  (una por cada iteración más las invocaciones recursivas anidadas) es a lo sumo  $(j - i) \times |P| \times M$ .  $\square$