

Apéndice B

Algoritmos de análisis sintáctico para CFG

En este capítulo se describen los sistemas de análisis sintáctico de aquellos algoritmos de análisis de gramáticas independientes del contexto que se han utilizado como base para la obtención de algoritmos de análisis sintáctico para LIG y TAG. En consecuencia, se abordan los algoritmos CYK, Earley ascendente y Earley.

B.1 El algoritmo CYK

El algoritmo de Cocke-Younger-Kasami (CYK) para análisis sintáctico de gramáticas independientes del contexto en forma normal de Chomsky [74] fue descubierto por J. Cocke, pero fue publicado independientemente por Younger [217] y Kasami [85], de ahí su nombre.

Es este un algoritmo ascendente punto basado en programación dinámica. A este respecto, la versión original del algoritmo hace uso de una matriz bidimensional indexada por posiciones de la cadena de entrada para almacenar los resultados parciales obtenidos, de tal modo que el elemento A se encuentra en la posición $[i, j]$ de dicha matriz si y sólo si $A \xrightarrow{*} a_{i+1} \dots a_j$. Sin embargo, nosotros seguiremos aquí el enfoque de Sikkel [174] de abstraer las estructuras de datos utilizadas puesto que estas forman parte de los detalles de implementación y no son esenciales al algoritmo. En este contexto, consideraremos que el algoritmo CYK construye un conjunto de ítems

$$\left\{ [A, i, j] \mid A \xrightarrow{*} a_{i+1} \dots a_n \right\}$$

Una cadena de entrada $a_1 \dots a_n$ pertenece al lenguaje generado por la gramática si y sólo si el ítem $[S, 0, n]$ se encuentra en dicho conjunto.

Esquema de análisis sintáctico B.1 El sistema de análisis sintáctico \mathbb{P}_{CYK} correspondiente al algoritmo CYK para una CFG $\mathcal{G} = (V_N, V_T, P, S)$ en forma normal de Chomsky, esto es, con producciones de la forma $A \rightarrow BC$ o $A \rightarrow a$, donde $A, B, C \in V_N$ y $a \in V_T$, y una cadena de entrada $a_1 \dots a_n$ es el que se muestra a continuación:

$$\mathcal{I}_{\text{CYK}} = \left\{ [A, i, j] \mid A \in V_N, 0 \leq i \leq j \right\}$$

$$\mathcal{H}_{\text{CYK}} = \left\{ [a, i-1, i] \mid a = a_i \right\}$$

$$\mathcal{D}_{\text{CYK}}^{\text{Scan}} = \left\{ [a, i, i+1] \vdash [A, i, i+1] \mid A \rightarrow a \in P \right\}$$

$$\mathcal{D}_{\text{CYK}}^{\text{Comp}} = \{ [B, i, k], [C, k, j] \vdash [A, i, j] \mid A \rightarrow BC \in P \}$$

$$\mathcal{D}_{\text{CYK}} = \mathcal{D}_{\text{CYK}}^{\text{Scan}} \cup \mathcal{D}_{\text{CYK}}^{\text{Comp}}$$

$$\mathcal{F}_{\text{CYK}} = \{ [S, 0, n] \}$$

§

Básicamente, podemos decir que el algoritmo CYK comienza por la creación de todos los ítems que se corresponden con reglas terminales (pasos $\mathcal{D}_{\text{CYK}}^{\text{Scan}}$) para posteriormente tratar de combinar en parejas los ítems generados a fin de reconocer las reglas binarias (pasos $\mathcal{D}_{\text{CYK}}^{\text{Comp}}$). El algoritmo termina cuando no se pueden combinar más ítems. En tal caso, si se ha generado algún ítem final la cadena de entrada ha sido reconocida.

B.2 Una versión ascendente del algoritmo de Earley

El algoritmo CYK presenta una importante limitación, ya que sólo es aplicable a gramáticas en forma normal de Chomsky. Nuestro objetivo ahora es extender el esquema CYK a la clase general de CFG, obteniendo un esquema Earley ascendente [176] en el que se construye un conjunto de ítems

$$\left\{ [A \rightarrow \alpha \bullet \beta, i, j] \mid \alpha \xrightarrow{*} a_{i+1} \dots j \right\}$$

que nos permitirán representar el reconocimiento parcial de producciones mediante la utilización de producciones con punto, al contrario de lo que ocurría en el caso del algoritmo CYK, el cual sólo permitía representar producciones binarias completas. En este sentido, un ítem CYK $[A, i, j]$ puede ser equivalentemente representado por $[A \rightarrow \alpha \bullet, i, j]$, donde $\alpha = BC$ o bien $\alpha = a_i$. El punto en las producciones indica que los elementos gramaticales situado a su izquierda han sido reconocidos. Las producciones son por tanto reconocidas de izquierda a derecha de tal modo que cuando el punto está a la derecha del último elemento del lado derecho de una producción si y sólo si esta ha sido completamente reconocida.

Esquema de análisis sintáctico B.2 El sistema de análisis \mathbb{P}_{buE} correspondiente al algoritmo Earley ascendente para una gramática independiente del contexto \mathcal{G} y una cadena de entrada $a_1 \dots a_n$ es el siguiente:

$$\mathcal{I}_{\text{buE}} = \{ [A \rightarrow \alpha \bullet \beta, i, j] \mid A \rightarrow \alpha\beta \in P, 0 \leq i \leq j \}$$

$$\mathcal{H}_{\text{buE}} = \mathcal{H}_{\text{CYK}}$$

$$\mathcal{D}_{\text{buE}}^{\text{Init}} = \{ \vdash [A \rightarrow \bullet \alpha, i, i] \}$$

$$\mathcal{D}_{\text{buE}}^{\text{Scan}} = \{ [A \rightarrow \alpha \bullet a\beta, i, j], [a, j, j+1] \vdash [A \rightarrow \alpha a \bullet \beta, i, j+1] \}$$

$$\mathcal{D}_{\text{buE}}^{\text{Comp}} = \{ [A \rightarrow \alpha \bullet B\beta, i, k], [B \rightarrow \gamma \bullet, k, j] \vdash [A \rightarrow \alpha B \bullet \beta, i, j] \}$$

$$\mathcal{D}_{\text{buE}} = \mathcal{D}_{\text{buE}}^{\text{Init}} \cup \mathcal{D}_{\text{buE}}^{\text{Scan}} \cup \mathcal{D}_{\text{buE}}^{\text{Comp}}$$

$$\mathcal{F}_{\text{buE}} = \{ [S \rightarrow \alpha \bullet, 0, n] \}$$

§

Proposición B.1 $\text{CYK} \xrightarrow{\text{ir}} \xrightarrow{\text{sr}} \xrightarrow{\text{ext}} \text{buE}$.

Efectivamente, el sistema de análisis \mathbb{P}_{buE} se deriva del sistema de análisis \mathbb{P}_{CYK} mediante la aplicación de un refinamiento de los ítems y de un refinamiento de pasos deductivos, puesto que se ha descompuesto el paso $\mathcal{D}_{\text{CYK}}^{\text{Comp}}$ en dos pasos $\mathcal{D}_{\text{buE}}^{\text{Init}}$ y $\mathcal{D}_{\text{buE}}^{\text{Comp}}$. Finalmente se ha realizado una extensión del sistema de análisis al considerar, no ya gramáticas en forma normal de Chomsky, sino en forma arbitraria. La prueba formal puede encontrarse en [174].

Con respecto al comportamiento del algoritmo, podemos resumirlo indicando que el análisis comienza con la creación por parte de los pasos $\mathcal{D}_{\text{buE}}^{\text{Init}}$ de los ítems $[A \rightarrow \bullet \alpha, i, i]$ para toda regla de la gramática y para toda posición en la cadena de entrada. A continuación se aplican los pasos $\mathcal{D}_{\text{buE}}^{\text{Scan}}$ y $\mathcal{D}_{\text{buE}}^{\text{Comp}}$ con el fin de ir desplazando el punto de las producciones hacia la derecha.

Un paso deductivo $\mathcal{D}_{\text{buE}}^{\text{Scan}}$ es aplicable a ítems de la forma $[A \rightarrow \alpha \bullet a\beta, i, j]$ cuando $a_{j+1} = a$, obteniéndose un nuevo ítem $[A \rightarrow \alpha a \bullet \beta, i, j+1]$. Es decir, se ha reconocido el símbolo terminal que estaba justo a la derecha del punto.

Un paso deductivo $\mathcal{D}_{\text{buE}}^{\text{Comp}}$ se aplica cuando un ítem tiene una regla con el punto en el extremo derecho. Dado un ítem $[B \rightarrow \gamma \bullet, k, j]$ se buscan todos los posibles $[A \rightarrow \alpha \bullet B\beta, i, k]$ y se generan nuevos ítems $[A \rightarrow \alpha B \bullet \beta, i, j]$ que representan que la subcadena $a_{k+1} \dots a_j$ puede ser reducida a B y por consiguiente, como la subcadena $a_{i+1} \dots a_k$ se reduce a α , la subcadena $a_{i+1} \dots a_j$ se reduce a αB .

El proceso termina cuando no se pueden combinar más ítems. En tal caso, si se ha generado algún ítem final la cadena de entrada pertenece al lenguaje generado por la gramática.

B.3 El algoritmo de Earley

Se puede derivar un esquema de análisis correspondiente al algoritmo de Earley [62] a partir del algoritmo de Earley ascendente mediante la aplicación de un *filtrado dinámico* a este último, de modo que los ítems de la forma $[A \rightarrow \bullet \alpha, i, i]$ no sean generados por los pasos Init para todas las posibles producciones y posiciones en la cadena de entrada, sino que sean generados o no dependiendo de la validez de otros ítems mediante un paso deductivo predictivo, encargándose el paso Init únicamente de la creación de los ítems correspondientes a las producciones del axioma de la gramática. Los ítems válidos son por tanto de la forma

$$\left\{ [A \rightarrow \alpha \bullet \beta, i, j] \mid \alpha \xrightarrow{*} a_{i+1} \dots j, S \xrightarrow{*} a_1 \dots a_i A \delta \right\}$$

Esquema de análisis sintáctico B.3 El sistema de análisis $\mathbb{P}_{\text{Earley}}$ correspondiente al algoritmo de Earley para una gramática independiente del contexto \mathcal{G} y una cadena de entrada $A_1 \dots a_n$ es el que se muestra a continuación:

$$\mathcal{I}_{\text{Earley}} = \mathcal{I}_{\text{buE}}$$

$$\mathcal{H}_{\text{Earley}} = \mathcal{H}_{\text{CYK}}$$

$$\mathcal{D}_{\text{Earley}}^{\text{Init}} = \{ \vdash [S \rightarrow \bullet\alpha, 0, 0] \}$$

$$\mathcal{D}_{\text{Earley}}^{\text{Scan}} = \mathcal{D}_{\text{buE}}^{\text{Scan}}$$

$$\mathcal{D}_{\text{Earley}}^{\text{Pred}} = \{ [A \rightarrow \alpha \bullet B\beta, i, j] \vdash [B \rightarrow \bullet\gamma, j, j] \}$$

$$\mathcal{D}_{\text{Earley}}^{\text{Comp}} = \mathcal{D}_{\text{buE}}^{\text{Comp}}$$

$$\mathcal{D}_{\text{Earley}} = \mathcal{D}_{\text{Earley}}^{\text{Init}} \cup \mathcal{D}_{\text{Earley}}^{\text{Scan}} \cup \mathcal{D}_{\text{Earley}}^{\text{Pred}} \cup \mathcal{D}_{\text{Earley}}^{\text{Comp}}$$

$$\mathcal{F}_{\text{Earley}} = \mathcal{F}_{\text{buE}}$$

§

Proposición B.2 $\text{buE} \xrightarrow{\text{df}} \text{Earley}$.

La prueba de dicha transformación puede encontrarse en [174]. Con respecto al comportamiento del algoritmo, lo resumiremos indicando que el proceso de análisis comienza con la generación, por parte del paso $\mathcal{D}_{\text{Earley}}^{\text{Init}}$, de los ítems de la forma $[S \rightarrow \bullet\alpha, 0, 0]$, donde $S \rightarrow \alpha \in P$ es una regla del axioma de la gramática. Estos ítems indican que se está tratando reconocer el axioma de la gramática desde el inicio de la cadena de entrada. Los pasos deductivos $\mathcal{D}_{\text{Earley}}^{\text{Scan}}$ y $\mathcal{D}_{\text{Earley}}^{\text{Comp}}$ se comportan como en el caso del algoritmo de Earley ascendente, mientras que los pasos $\mathcal{D}_{\text{Earley}}^{\text{Pred}}$ se encarga de realizar la fase descendente o predictiva del algoritmo, ya que a partir de ítems de la forma $[A \rightarrow \alpha \bullet B\beta, i, j]$ genera los ítems $[B \rightarrow \bullet\gamma, j, j]$, esto es, se predicen todas las reglas que potencialmente pueden ser útiles en el reconocimiento de la cadena de entrada.