

Procesadores de Lenguajes I:

Una introducción a la fase de análisis

*J. M. Benedí Ruíz
V. Gisbert Giner
L. Moreno Boronat
E. Vivancos Rubio*

Capítulo 3

Análisis sintáctico descendente

En este capítulo se realiza la presentación de las técnicas de análisis sintáctico empleadas en la construcción de un compilador. Tras una breve introducción al análisis sintáctico descendente y ascendente, se presentan los conceptos fundamentales que permiten la utilización del formalismo de los autómatas de pila para la especificación del analizador sintáctico. El resto del capítulo está dedicado a mostrar la técnica de análisis sintáctico descendente predictiva LL y los métodos de recuperación de errores.

3.1. Conceptos fundamentales

3.1.1. Especificación sintáctica de los lenguajes de programación

Mientras que para la especificación de las unidades léxicas de un lenguaje se requiere un formalismo como las Expresiones Regulares, en el caso de la especificación de las estructuras sintácticas de un lenguaje de programación se necesita también un formalismo. El formalismo que utilizaremos para este último propósito es el de las Gramáticas Independientes del Contexto (GIC).

Podemos apuntar diferentes propiedades de este formalismo que justifican su utilización como herramienta de especificación sintáctica de los lenguajes de programación:

1. Es precisa y fácil de entender.
2. Para algunos lenguajes hay algoritmos de análisis eficientes basados en este formalismo.
3. La ampliación de nuevas estructuras del lenguaje es sencilla.

Una **Gramática Independiente del Contexto (o Gramática Incontextual)** se define como: $G = (N, \Sigma, P, S)$, donde N representa el conjunto de símbolos no-terminales, Σ representa el conjunto de símbolos terminales, P el conjunto de reglas de producción donde cada regla $(A \rightarrow \alpha)$

$\in P$ es tal que $A \in N$ y $\alpha \in (N \cup \Sigma)^*$ (siendo $N \cap \Sigma = \emptyset$); y S es el símbolo inicial de la gramática tal que $S \in N$.

Ejemplo 3.1.- $G = (\{E, O\}, \{+, -, *, /, (,), id\}, P, E)$ donde P es:

$$E \rightarrow E O E \mid (E) \mid - E \mid id$$

$$O \rightarrow + \mid - \mid * \mid /$$

Estas reglas de producción pueden verse como reglas de reescritura. Así, la regla de producción $O \rightarrow +$, representa que se reemplaza una instancia de O por $+$.

Respecto a la aplicación de las reglas de producción aparecen los siguientes conceptos:

Derivación Directa. Representa la aplicación de una regla de producción.

$$\delta A \gamma \rightarrow \delta \beta \gamma \quad \text{si} \quad \exists (A \rightarrow \beta) \in P; \quad \delta, \gamma \in (N \cup \Sigma)^*$$

Derivación. Representa la aplicación a una cadena de 0 o más derivaciones directas.

$$\alpha \rightarrow^* \beta \quad \text{si} \quad \exists \alpha_0, \dots, \alpha_m \in (N \cup \Sigma)^*; \alpha_0 = \alpha \rightarrow \alpha_1 \rightarrow \dots \alpha_{m-1} \rightarrow \alpha_m = \beta$$

Forma Sentencial. Es una cadena que puede derivarse a partir del símbolo inicial de la gramática:

$$\alpha \in (N \cup \Sigma)^* \quad \text{si} \quad \exists S \rightarrow^* \alpha$$

Sentencia. Es una cadena formada por terminales que puede derivarse a partir del símbolo inicial de la gramática:

$$\alpha \in \Sigma^* \quad \text{si} \quad \exists S \rightarrow^* \alpha$$

Lenguaje Generado por una gramática G . Lenguaje formado por el conjunto de cadenas de símbolos terminales (o sentencias) que pueden derivarse a partir del símbolo inicial de la gramática.

$$L(G) = \{x / x \in \Sigma^*: S \rightarrow^* x\}$$

Ejemplo 3.2.- Dadas las siguientes reglas de producción de una gramática:

$$S \rightarrow a b \mid a S b,$$

podemos obtener distintas formas sentenciales por sucesivas derivaciones como

$$S \rightarrow a S b \rightarrow a a S b b \rightarrow a a a b b b$$

donde la última cadena ($a a b b b$) es una *sentencia* por estar formada únicamente por símbolos terminales.

3.1.2. El papel del analizador sintáctico

El objetivo fundamental del analizador sintáctico es construir el árbol sintáctico conocidas las hojas (terminales) y el símbolo inicial de la gramática.

Dada una gramática $G = (N, \Sigma, P, S)$, un árbol sintáctico, también denominado *árbol de derivación*, se define de la siguiente forma:

1. La raíz está etiquetada con el símbolo inicial S .
2. Cada hoja está etiquetada con un símbolo de $\Sigma \cup \{ \varepsilon \}$.
3. Cada nodo interior está etiquetado con un símbolo de N .
4. Sea $A \in N$ la etiqueta de un nodo, y x_1, \dots, x_n las etiquetas de sus nodos hijos (de izquierda a derecha). Entonces: $A \rightarrow x_1 \dots x_n \in P$

A la cadena obtenida por concatenación de las hojas del árbol de derivación se le denomina *frontera*.

Si identificamos cada regla de la gramática por un número, el *análisis* o “*parse*” de una cadena consiste en una secuencia de números resultante de la numeración sucesiva de cada una de las reglas utilizadas en el proceso de derivación, representando internamente el árbol de derivación,

Si en el proceso de derivación se sustituye siempre el símbolo no-terminal que se encuentra más a la izquierda de la forma sentencial decimos que se ha realizado una *derivación a izquierdas*. Si se sustituye siempre el no-terminal más a la derecha, decimos que se ha realizado una *derivación a derechas*.

Ejemplo 3.3.- Dadas las siguientes reglas de producción de una gramática:

$$S \rightarrow S + T \quad (1)$$

$$S \rightarrow T \quad (2)$$

$$T \rightarrow T * F \quad (3)$$

$$T \rightarrow F \quad (4)$$

$$F \rightarrow (S) \quad (5)$$

$$F \rightarrow a \quad (6)$$

$$F \rightarrow b \quad (7)$$

y siendo S el símbolo inicial de ésta, analizaremos la cadena $a * (a + b)$ siguiendo las dos estrategias antes mencionadas, derivación a izquierdas (DI) y derivación a derechas (DD):

DI: $S \rightarrow^2 T \rightarrow^3 T * F \rightarrow^4 F * F \rightarrow^6 a * F \rightarrow^5 a * (S) \rightarrow^1 a * (S + T)$

$\rightarrow^2 a * (T + T) \rightarrow^4 a * (F + T) \rightarrow^6 a * (a + T) \rightarrow^4 a * (a + F) \rightarrow^7 a * (a + b)$

siendo por tanto el “parse” izquierdo = 2-3-4-6-5-1-2-4-6-4-7

DD: $S \rightarrow^2 T \rightarrow^3 T * F \rightarrow^5 T * (S) \rightarrow^1 T * (S + T) \rightarrow^4 T * (S + F) \rightarrow^7 T * (S + b)$

$\rightarrow^2 T * (T + b) \rightarrow^4 T * (F + b) \rightarrow^6 T * (a + b) \rightarrow^4 F * (a + b) \rightarrow^6 a * (a + b)$

siendo el “parse” derecho = 6-4-6-4-2-7-4-1-5-3-2. En ambos casos, como puede verse en la Fig. 3.1 el árbol de derivación es el mismo aunque se ha obtenido por métodos de análisis distintos.

Teorema.

Un árbol tiene un único “parse” izquierdo y un único “parse” derecho.

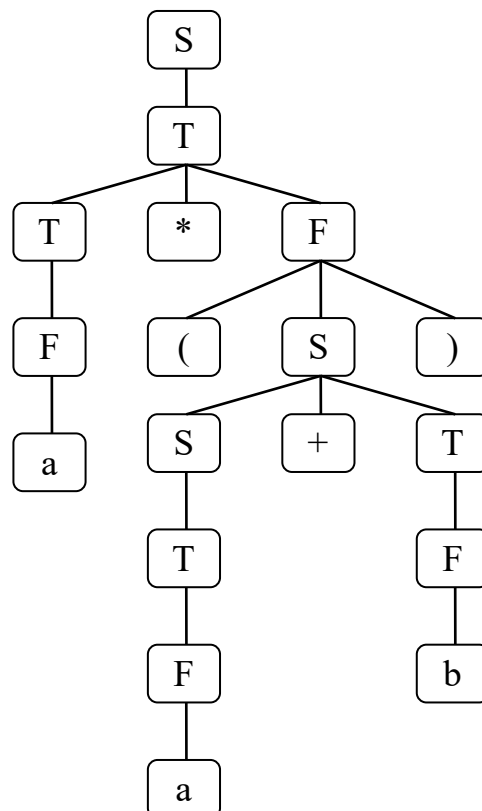


Figura 3.1. Árbol sintáctico correspondiente a $a*(a+b)$.

Teorema.

Dada $G(S)$, $S \rightarrow^* \alpha$ (con $\alpha \in \Sigma^*$) si y solo si existe un árbol de derivación que produce α .

Gramática Ambigua. Una gramática es ambigua si y solo si existe una cadena $x \in \Sigma^*$ / $x \in L(G)$ para la cual hay más de un árbol de derivación que produce x . En la Fig. 3.2 puede verse un ejemplo de gramática ambigua para una instrucción condicional if-then-else.

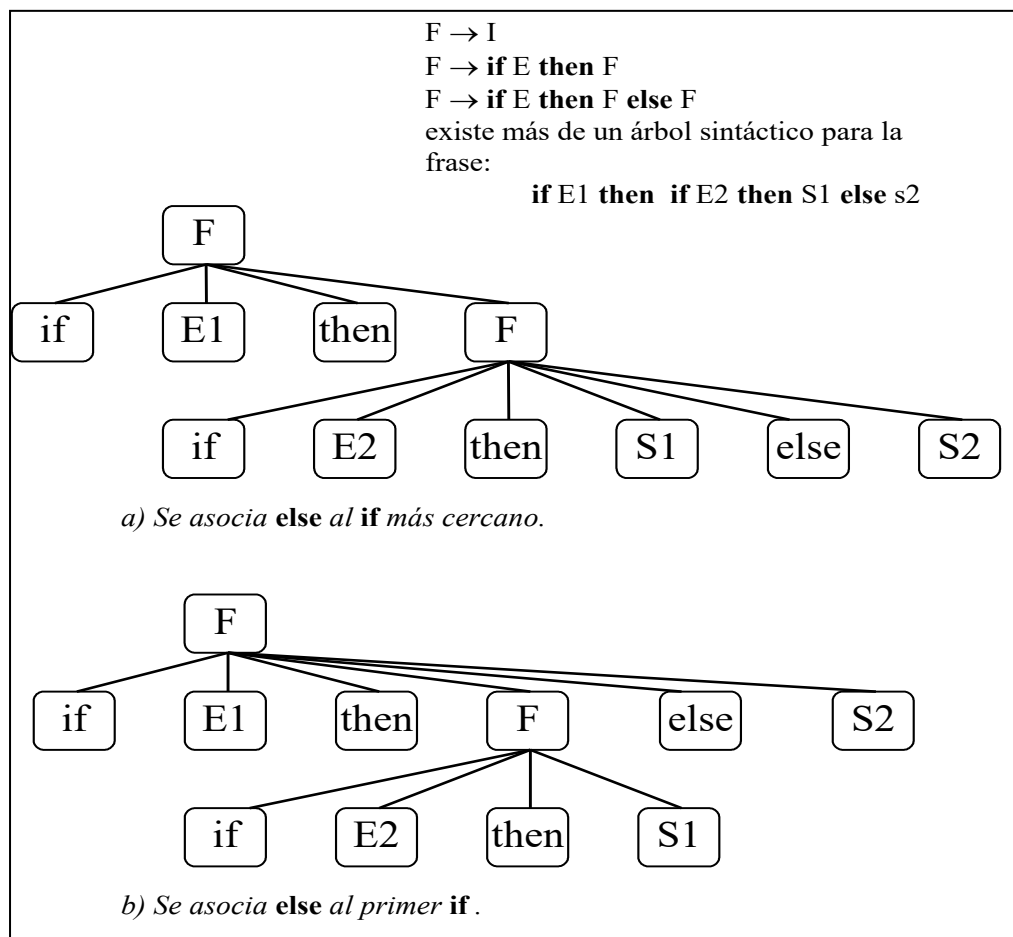


Figura 3.2. Gramática Ambigua.

Proposición.

Determinar si una gramática es ambigua, es un problema indecidible.

Es decir, no existe un algoritmo que determine en un tiempo finito si una gramática es ambigua o no. De la misma manera, tampoco existe un algoritmo general que produzca una gramática no ambigua a partir de una gramática ambigua, pero en algunos casos se puede transformar la gramática para eliminar la ambigüedad¹.

Una gramática ambigua crea problemas para la traducción de los programas ya que distintos árboles sintácticos pueden tener distintas interpretaciones semánticas.

¹ De hecho hay lenguajes intrínsecamente ambiguos, es decir, lenguajes para los que no existe una gramática no ambigua.

Ejemplo 3.4.- La gramática de la Fig. 3.2 se puede transformar de tal manera que produciendo el mismo lenguaje deja de ser ambigua.

$$\begin{aligned} F &\rightarrow F1 \mid F2 \\ F2 &\rightarrow \text{if } E \text{ then } F2 \text{ else } F2 \mid I \\ F1 &\rightarrow \text{if } E \text{ then } F \mid \text{if } E \text{ then } F2 \text{ else } F1 \end{aligned}$$

3.1.3. Aceptores: Reconocedores y Traductores

Una máquina reconocedora o autómata permite definir un lenguaje por medio de la vía de aceptación de las cadenas que lo componen.

Las máquinas reconocedoras deciden si una cadena pertenece o no a un lenguaje, en contraposición con las gramáticas formales que generan las cadenas a partir de las reglas de producción.

En la Fig. 1.1 del capítulo 1 se presentaron los distintos tipos de reconocedores junto a la fase de análisis en la que se utilizan cada uno de ellos. En el análisis sintáctico se utilizan los autómatas a pila como reconocedores de lenguajes independientes del contexto.

Autómata a Pila

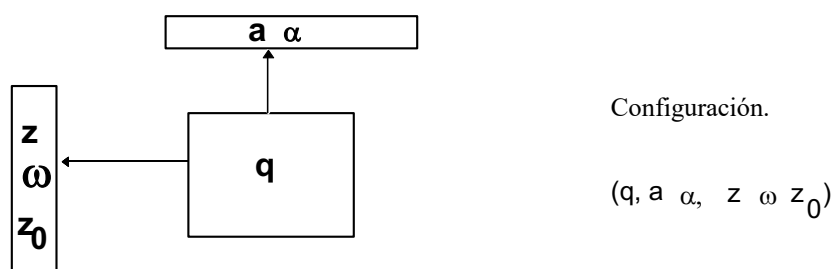
Un autómata a pila se define como:

$$A_p = (Q, \Sigma, \Gamma, \delta, q_0, z_0, F)$$

Donde Q representa un conjunto finito de estados, Σ representa el alfabeto de entrada, Γ representa el alfabeto de la pila, δ es la función de transición tal que

$\delta : Q \times (\Sigma \cup \{\varepsilon\}) \times \Gamma \rightarrow P(Q \times \Gamma^*)$, q_0 es el estado inicial tal que $q_0 \in Q$, z_0 representa el símbolo inicial de la pila, y F es el conjunto de estados finales: $F \subset Q$.

• Estructura:



En esta configuración podemos ver que q representa el estado actual del autómata, α la entrada de la cadena que queda por analizar, y z es el símbolo que se encuentra en la cima de la pila.

- **Movimiento.**

$(q, a \alpha, z \omega z_0) \mid \longrightarrow (q', \alpha, \beta \omega z_0)$ es un movimiento válido sii $\exists \delta(q, a, z) = (q', \beta)$
con $q, q' \in Q; a \in \Sigma \cup \{\epsilon\}; \alpha \in \Sigma^*; z \in \Gamma; \beta, \omega \in \Gamma^*$

- **Lenguaje Aceptado a Pila Vacía**

$$L_v(A_p) = \{x / x \in \Sigma^* : (q_0, x, z_0) \mid^* \longrightarrow (q, \epsilon, \epsilon)\}$$

- **Lenguaje Aceptado a Estado Final**

$$L_F(A_p) = \{x / x \in \Sigma^* : (q_0, x, z_0) \mid^* \longrightarrow (q, \epsilon, a); \quad q \in F\}$$

Si una máquina reconocedora es además capaz de emitir cadenas de algún alfabeto determinado, entonces a esta máquina se le denomina máquina traductora. Así, si un Autómata a Pila es capaz de emitir cadenas de un determinado alfabeto, éste se denominará *Traductor a pila*.

3.1.4. Introducción general al análisis sintáctico

En la construcción del árbol de derivación se puede seguir una estrategia de análisis *ascendente* o *descendente*.

El **análisis descendente** puede verse como una derivación a izquierdas. En términos generales, podemos decir que las operaciones en este análisis consisten en realizar sucesivas *derivaciones* hasta que el símbolo de la cima de la pila corresponda con el primer símbolo que hay en la cadena a analizar, en cuyo caso se *consume* este símbolo de la cadena y se *extrae* el correspondiente de la *cima de la pila*. Equivale a construir el árbol de derivación de la cadena a analizar desde la raíz a las hojas.

El **análisis ascendente** puede verse como una derivación a derechas recorrida al revés. En la derivación a derechas del Ejemplo 3.3 vimos que pasamos de la forma sentencial $T * F$ a T aplicando al revés la regla de producción 3. En este caso se dice que $T * F$ se ha *reducido*. En términos generales, en el análisis ascendente las operaciones que se realizan son: *reducir* por la parte izquierda de una determinada regla de producción, si la cadena que se encuentra en la cima de la pila coincide con la parte derecha de dicha regla, o bien, si no es posible la reducción entonces *desplazar* a la pila un símbolo de la cadena a analizar. Esta técnica es equivalente a construir el árbol de derivación de la cadena a analizar desde las hojas a la raíz.

Durante el proceso de análisis sintáctico, puede presentarse el problema del indeterminismo. Por ejemplo, el **indeterminismo** se nos plantea cuando en la GIC se presentan varias reglas para un mismo símbolo auxiliar como $A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$. Podemos enfrentarnos a este problema aplicando distintos tipos de estrategias:

- Algoritmos basados en backtracking (o vuelta atrás).

Es una estrategia muy general que sirve para cualquier gramática independiente del contexto, pero no son los más adecuados para analizar los lenguajes de programación por su complejidad exponencial. La sección 4.1 de [Aho 72] está dedicada a discutir este tipo de técnicas.

- Algoritmos generales o tabulares, como el algoritmo de Early o el de Cocke-Younger-Kasami.

Tienen una complejidad más baja $O(n^3)$, que puede llegar a ser de $O(n^2)$ en el algoritmo de Early si la gramática no es ambigua. Sin embargo este coste sigue siendo muy elevado si tenemos en cuenta que el tamaño del problema viene determinado por el número de símbolos que aparecen en el programa fuente a analizar. En la sección 4.2 de [Aho 72] se pueden consultar estos dos algoritmos.

- Algoritmos deterministas.

Algoritmos descendentes y ascendentes deterministas que se caracterizan por encontrar una única alternativa a aplicar dada una determinada cadena de entrada y un no-terminal. Consecuentemente tendrán un coste $O(n)$. En este último tipo de algoritmos es en el que nos vamos a centrar en los próximos apartados y en el siguiente capítulo.

Lamentablemente, estos últimos métodos solo son apropiados para un subconjunto de GIC. Por esta razón, para poder realizar un análisis determinista de coste lineal nos centraremos en un subconjunto de gramáticas incontextuales: aquellas en las que desde cada configuración sólo hay un posible movimiento adecuado.

3.2. Análisis sintáctico descendente

Un analizador LL(1) analiza la cadena de izquierda a derecha (Left-to-right), siguiendo un proceso de derivación a izquierdas (Leftmost-derivation). Para decidir la regla de producción a utilizar es suficiente con examinar sólo un símbolo (cadena de símbolos terminales de longitud 1) de la cadena de entrada.

Ejemplo 3.5.- En la siguiente gramática podemos saber qué regla aplicar para el auxiliar S (la número 2) en el análisis de la cadena “bd” observando el primer símbolo de la cadena (“b”),

$$S \rightarrow a B B \quad (1)$$

$$S \rightarrow b E \quad (2)$$

$$B \rightarrow c \quad (3)$$

$$E \rightarrow d \quad (4)$$

Para construir un analizador sintáctico descendente predictivo LL(1), nos apoyaremos en dos funciones: Primeros y Siguientes.

El dominio de la función Primeros es: $(N \cup \Sigma)^* \rightarrow P(\Sigma \cup \{\varepsilon\})$ donde primeros de una cadena α de auxiliares y terminales se define como:

$$\text{Primeros}(\alpha) = \{ a \in \Sigma \mid \alpha \xrightarrow{*} a \beta \} \cup \{ \varepsilon \mid \alpha \xrightarrow{*} \varepsilon \}$$

con $\alpha, \beta \in (N \cup \Sigma)^*$

Se puede observar que se puede calcular Primeros de cadenas de terminales y no-terminales de cualquier tamaño, y el resultado será un conjunto de símbolos terminales, (incluyendo en algunos casos la cadena vacía). El algoritmo 3.1 muestra una forma de calcular el conjunto de elementos pertenecientes a Primeros de una cadena x .

Función Primeros ($x \in (N \cup \Sigma)^*$): Conjunto de $(\Sigma \cup \{\varepsilon\})$

Dada $G = (N, \Sigma, P, S)$; con C : Conjunto de $(\Sigma \cup \{\varepsilon\})$

Comienzo

$C := \emptyset$;

Si $x \in (\Sigma \cup \varepsilon)$ ent $C := \{x\}$

Si $x \in N$ ent

Para toda $(x \rightarrow \alpha) \in P$ hacer $C := C \cup \text{Primeros}(\alpha)$

Si $x = x_1 x_2 \dots x_m \wedge m > 1$ ent

$i := 1$;

Mientras $(i < m) \wedge (\varepsilon \in \text{Primeros}(x_i))$ hacer

$C := C \cup (\text{Primeros}(x_i) - \{\varepsilon\})$;

$i := i + 1$;

$C := C \cup (\text{Primeros}(x_i) - \{\varepsilon\})$;

Si $(i = m) \wedge (\varepsilon \in \text{Primeros}(x_m))$ entonces $C := C \cup \{\varepsilon\}$

Devolver C

Fin

Algoritmo 3.1. Función Primeros.

El algoritmo 3.1 calcula primeros de una cadena x en un conjunto C .

- a) Si se trata de un símbolo terminal o de la cadena vacía, C será igual al mismo símbolo terminal o la cadena vacía.
- b) Si x es un símbolo no-terminal, C será el conjunto de Primeros de los lados derechos de las producciones del símbolo no-terminal x .

- c) Finalmente, si x es una cadena, $\text{Primeros}(x)$ incluirá a Primeros del primer símbolo de la cadena x . Si la cadena vacía pertenece a Primeros de este primer símbolo de la cadena x_1 , no se incluirá en C , pero se añadirá a C el conjunto de Primeros del siguiente símbolo de la cadena x_2 . Se puede observar en el algoritmo que debe hacerse lo mismo con el resto de símbolos de la cadena x hasta que Primeros del símbolo que se está considerando (x_i) no incluya a la cadena vacía, o se haya llegado al final de la cadena. De esta manera, la cadena vacía pertenecerá a Primeros de una cadena x si y solo si la cadena vacía pertenece a Primeros de todos los símbolos que forman la cadena x .

Ejemplo 3.6.- Calcular el conjunto Primeros de los símbolos no-terminales de la gramática

$$E \rightarrow T E'$$

$$E' \rightarrow + T E' \mid \varepsilon$$

$$T \rightarrow F T'$$

$$T' \rightarrow * F T' \mid \varepsilon$$

$$F \rightarrow (E) \mid a$$

$$\text{Primeros}(E) = \text{Primeros}(T) - \{\varepsilon\} = \text{Primeros}(F) - \{\varepsilon\} = \{ (, a \}$$

$$\text{Primeros}(E') = \{ +, \varepsilon \} \quad \text{Primeros}(T') = \{ *, \varepsilon \}$$

La otra función que se usa para construir un analizador LL(1) es la función *Siguientes*. El dominio de la función *Siguientes* es: $N \longrightarrow P(\Sigma \cup \{\$ \})$ donde *Siguientes* de un auxiliar A se define como:

$$\text{Siguientes}(A) = \{ a \in \Sigma \mid S \xrightarrow{*} \alpha A a \beta \} \cup \{ \$ \mid S \xrightarrow{*} \alpha A \}$$

con $\alpha, \beta \in (N \cup \Sigma)^*$; $A \in N$

El algoritmo 3.2 muestra una forma de calcular el conjunto de elementos pertenecientes a *Siguientes* de los símbolos no-terminales de una gramática.

Función *Siguientes*: Conjunto de $(\Sigma \cup \{\$ \})$

Dada $G = (N, \Sigma, P, S)$; con $C = \text{vector}[N]$ de conjunto de $(\Sigma \cup \{\$ \})$

Comienzo

Para todo $A \in N$ hacer $C[A] = \phi$
 $C[S] = \{ \$ \}$

Para toda $(B \rightarrow \gamma) \in P$ hacer
Para todo $A \in N$ $\wedge \gamma = \alpha A \beta$ hacer
 $C[A] = C[A] \cup (\text{Primeros}(\beta) - \{ \varepsilon \})$

Repetir

Para toda $(B \rightarrow \gamma) \in P$ hacer
Para todo $A \in N$ $\wedge \gamma = \alpha A \beta$ $\wedge \beta \rightarrow^* \varepsilon$ hacer
 $C[A] = C[A] \cup C[B]$

hasta que $C[A]$ no se modifique para ningún $A \in N$

Para toda $A \in N$ hacer Siguietes (A) := C [A]
Devolver Siguietes

Fin

Algoritmo 3.2. Función Siguietes.

El algoritmo 3.2 va calculando Siguietes de los no-terminales de la gramática en un array indexado por los símbolos no-terminales: C. Inicialmente los conjuntos de Siguietes para todos los símbolos de la gramática están vacíos (excepto para el símbolo inicial de la gramática, ya que se le incluye un símbolo especial \$). Si en el lado derecho de una producción ($B \rightarrow \alpha A \beta$) aparece un no-terminal A seguido de una cadena β , se añaden los Primeros de β a los siguientes del no-terminal A². Si Primeros de dicha cadena β incluye a la cadena vacía, se añaden al conjunto C[A] los Siguietes del símbolo no-terminal del lado izquierdo de la producción (B).

El conjunto Siguietes de los símbolos no-terminales de la gramática del Ejemplo 3.6 sería:

$$\text{Siguietes}(E') = \text{Siguietes}(E) = \{ \$,) \}$$

$$\text{Siguietes}(T') = \text{Siguietes}(T) = \{ +, \$,) \}$$

$$\text{Siguietes}(F) = \{ *, +, \$,) \}$$

3.2.1. Condición LL(1)

Es posible construir un analizador sintáctico descendente predictivo usando un símbolo de anticipación para cualquier GIC que cumpla la siguiente condición, denominada Condición LL(1):

Definición

Una gramática independiente del contexto es LL(1), si para cualquier par de producciones de un mismo auxiliar A ($A \rightarrow \alpha \wedge A \rightarrow \beta$) se cumple la condición:

$$\text{Primeros}(\alpha \text{ Siguietes}(A)) \cap \text{Primeros}(\beta \text{ Siguietes}(A)) = \emptyset$$

Para comprobar si una gramática cumple la condición LL(1) se comprueba si la cumple para todos sus no-terminales. Dado un no-terminal A, se calcula Primeros del lado derecho de cada una de sus producciones. Si alguno de estos conjuntos incluye a la cadena vacía, se eliminará ésta, y en su lugar se añadirá Siguietes del no-terminal A. Si ninguno de los conjuntos calculados de esta forma para los lados derechos de las producciones de A incluye ningún elemento incluido

² Notad que no se incluye la cadena vacía porque ésta nunca pertenece al conjunto de Siguietes.

en otro de los conjuntos (es decir, si los conjuntos son disjuntos tomados dos a dos), la gramática cumple la condición LL(1) para ese no-terminal (A). La gramática cumplirá la condición LL(1) (también decimos que la gramática es LL(1)) si cumple la condición para todos sus no-terminales.

Ejemplo 3.7- La gramática del Ejemplo 3.6 es LL(1) ya que cumple la condición LL(1):

$$\text{Primeros}(+TE' \text{ Siguietes}(E')) \cap \text{Primeros}(\varepsilon \text{ Siguietes}(E')) = \{+\} \cap \{\$, \, \} = \emptyset$$

$$\text{Primeros}(*FT' \text{ Siguietes}(T')) \cap \text{Primeros}(\varepsilon \text{ Siguietes}(T')) = \{*\} \cap \{+, \$, \, \} = \emptyset$$

$$\text{Primeros}((E) \text{ Siguietes}(F)) \cap \text{Primeros}(a \text{ Siguietes}(F)) = \{(\} \cap \{a \} = \emptyset$$

3.2.2. Transformaciones de las gramáticas

Factorización:

Cuando en una GIC, cualquier par de producciones de un mismo no-terminal comienzan por el mismo prefijo α la gramática no puede ser LL(1):

$$A \rightarrow \alpha \beta_1 \mid \alpha \beta_2 \mid \dots \mid \alpha \beta_n \mid \gamma$$

Una solución podría ser factorizar, transformando la gramática en

$$A \rightarrow \alpha A' \mid \gamma$$

$$A' \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n \quad \text{donde } A' \text{ es un nuevo no-terminal.}$$

Ejemplo 3.8.- Dadas las siguientes reglas de producción de una gramática,

$$I \rightarrow aB \mid aBc \mid d$$

$$B \rightarrow \varepsilon$$

se puede observar que no es LL(1),

$$\begin{array}{l} \text{Primeros}(aB \text{ Siguietes}(I)) = \{a\} \\ \text{Primeros}(aBc \text{ Siguietes}(I)) = \{a\} \end{array} \quad \left| \quad \{a\} \cap \{a\} \neq \emptyset \right.$$

Si factorizamos por la izquierda esta gramática se convierte en,

$$I \rightarrow aBI' \mid d$$

$$I' \rightarrow c \mid \varepsilon$$

$$B \rightarrow \varepsilon$$

que si que es LL(1):

$$\text{Primeros}(aBI' \text{ Siguietes}(I)) = \{a\} \quad \left| \quad \right.$$

$$\begin{array}{l|l}
 \text{Primeros}(d \text{ Siguientes}(I)) = \{d\} & \{a\} \cap \{d\} = \emptyset \\
 \text{Primeros}(c \text{ Siguientes}(I')) = \{c\} & \{c\} \cap \{\$ \} = \emptyset \\
 \text{Primeros}(\varepsilon \text{ Siguientes}(I')) = \{\$ \} &
 \end{array}$$

Eliminación de la recursión a izquierdas:

Definición

Una GIC es recursiva a izquierdas sii $\exists A \in N: A \rightarrow^* A \alpha$.

Teorema

Una GIC LL(1) no es ambigua ni recursiva a izquierdas.

Según este teorema, una gramática que sea recursiva a izquierdas ($\exists A \in N: A \rightarrow^* A \alpha$) no puede ser LL(1). No obstante, frecuentemente se puede realizar un análisis descendente predictivo, transformando la gramática en otra equivalente que genere el mismo lenguaje.

Un caso particular de gramática con recursión a la izquierda es el de las gramáticas con recursión a izquierdas inmediata. Una GIC tiene recursión a izquierdas inmediata si y solo si \exists una producción de la forma $A \rightarrow A \alpha \in P$.

En general, las producciones de un auxiliar con recursión a izquierdas inmediata son de la forma:

$$A \rightarrow A \alpha_1 \mid A \alpha_2 \mid \dots \mid A \alpha_n \mid \beta_1 \mid \beta_2 \mid \dots \mid \beta_p$$

donde los β_i no comienzan por A , $\forall i$.

Una solución para eliminar la recursión a izquierdas inmediata es transformar las reglas genéricas anteriores en,

$$\begin{array}{l}
 A \rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_p A' \\
 A' \rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_n A' \mid \varepsilon
 \end{array}$$

donde A' es un nuevo no-terminal.

En la sección 4.4 de [Aho 90] aparece un algoritmo (algoritmo 4.1) para transformar una gramática con recursión a izquierdas no inmediata en otra sin recursión.

Ejemplo 3.9.- Dadas las siguientes reglas de producción de una gramática recursiva a izquierdas,

$$E \rightarrow E + T \mid T$$

$$\begin{aligned} T &\rightarrow T * F \mid F \\ F &\rightarrow (E) \mid a \end{aligned}$$

Se puede eliminar la recursividad a izquierdas inmediata aplicando la transformación vista anteriormente, obteniendo una gramática equivalente:

$$\begin{aligned} E &\rightarrow T E' \\ E' &\rightarrow + T E' \mid \varepsilon \\ T &\rightarrow F T' \\ T' &\rightarrow * F T' \mid \varepsilon \\ F &\rightarrow (E) \mid a \end{aligned}$$

Ejemplo 3.10.- Dadas las siguientes reglas de producción de una gramática en la que se observa el problema de la ambigüedad que introduce el *else*,

$$\begin{aligned} I &\rightarrow \text{if } E \text{ then } I \mid \text{if } E \text{ then } I \text{ else } I \mid a \\ E &\rightarrow b \end{aligned}$$

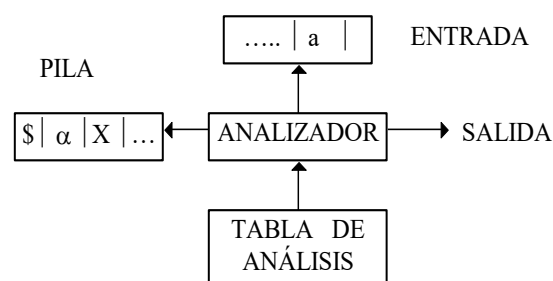
Si factorizamos por la izquierda esta gramática se convierte en,

$$\begin{aligned} I &\rightarrow \text{if } E \text{ then } I I' \mid a \\ I' &\rightarrow \text{else } I \mid \varepsilon \\ E &\rightarrow b \end{aligned}$$

Es evidente, que a pesar de la transformación la gramática, no es LL(1), ya que sigue siendo ambigua: para I' no se sabrá que alternativa elegir cuando la entrada sea *else*. Una solución a este problema se vio en el Ejemplo 3.4.

3.2.3. Modelo de un analizador sintáctico descendente predictivo.

Estructura:



Configuración:

$$(x, \beta, \pi) \text{ siendo } x \in \Sigma^*; \beta \in (N \cup \Sigma \cup \{\$\})^*; \pi \in \{1, 2, \dots, N\}^*$$

donde “x” representa la cadena de entrada que queda por analizar, “ β ” representa el contenido de la pila, y “ π ” la cadena formada por la concatenación de los números de las reglas que se han utilizado en las sucesivas derivaciones durante el proceso de análisis (parse izquierdo).

Algoritmo de Análisis:

El algoritmo de análisis LL(1) es el que se muestra en 3.4. Como se observa hace uso de una tabla de análisis que previamente debe construirse según el algoritmo 3.3.

ALGORITMO Construcción de la T.A. LL(1)

ENTRADA $G = (N, \Sigma, P, S);$

SALIDA $TA: (N \cup \Sigma \cup \{\$, \}) \times (\Sigma \cup \{\$, \}) \longrightarrow \{(r: A \rightarrow \beta), \text{sacar}, \text{aceptar}, \text{error}\}$

METODO

1) Inicializar TA con la acción “error”;

2) para toda $(r: A \rightarrow \beta) \in P$ hacer

para todo $a \in (\Sigma \cup \{\$, \})$: $a \in \text{PRIMEROS}(\beta \text{ SIGUIENTES}(A))$ hacer

$TA[A, a] := (r: A \rightarrow \beta);$

3) para todo $a \in \Sigma$ hacer $TA[a, a] := \text{sacar};$

4) $TA[\$, \$] := \text{aceptar};$

FIN

Algoritmo 3.3. Construcción de la Tabla de análisis.

Una tabla de análisis LL(1) tendrá una fila por cada símbolo de la gramática, y una columna por cada terminal y otra para el símbolo \$. Para completar la fila del símbolo no-terminal A se calcula, para cada una de las producciones de A ($A \rightarrow \beta$), $\text{Primeros}(\beta \text{ Siguientes}(A))$. Para cada uno de los terminales resultantes, se añadirá en su columna la acción correspondiente a aplicar la producción $A \rightarrow \beta$.

Teorema

Una gramática es LL(1) si y solo si el algoritmo de construcción de la tabla de análisis no genera ningún conflicto, es decir, si $\forall X, Y / X \in V \wedge Y \in V$ se cumple que $|TA[X, Y]| = 1$.

Ejemplo 3.11.- La Tabla de análisis correspondiente a la gramática del Ejemplo 3.9, siguiendo el algoritmo 3.3 es la siguiente:

	a	+	*	()	\$
E	(TE',1)			(TE',1)		
E'		(+TE',2)			(ε,3)	(ε,3)
T	(FT',4)			(FT',4)		
T'		(ε,6)	(*FT',5)		(ε,6)	(ε,6)

F	(a,8)			((E),7)		
a	sacar					
+		sacar				
*			sacar			
(sacar		
)					sacar	
\$						aceptar

ALGORITMO A.S.D: basado en la T.A. LL(1)

ENTRADA $\omega \in \Sigma^*$; TA, para una $G = (N, \Sigma, P, S)$;

SALIDA Si $\omega \in L(G)$ entonces π else error()

METODO

```

apilar (S$);    sim = obtsym;     $\pi := \epsilon$ ; fin = falso;
repetir
    caso    TA [cima,sim] sea
        "(r: A  $\rightarrow$   $\beta$ )": desapilar;    apilar( $\beta$ );     $\pi := \pi \cdot r$ ;
        "sacar":    desapilar;    sim := obtsym;
        "aceptar":    fin := verdad;
        "error":    error( );
    end;
hasta fin

```

FIN

Algoritmo 3.4. Algoritmo de análisis LL(1).

Como se puede apreciar en el Algoritmo 3.4, las entradas que pueden aparecer en una tabla de análisis LL(1) son:

$r: A \rightarrow \beta$ que consiste en aplicar la producción indicada. Para ello se desapila de la cima del analizador el no-terminal A, y se apilan los símbolos que forman el lado derecho de la producción (β).

Sacar: que consiste en quitar de la cima de la pila un terminal, y pasar a considerar el siguiente símbolo de la cadena de entrada.

Aceptar: Fin del análisis, habiendo reconocido la cadena de la entrada.

Error (normalmente indicado por una celda vacía): La cadena de la entrada no pertenece al lenguaje reconocido por el autómata LL(1), o lo que es lo mismo, no puede ser generada por la gramática usada para construir la tabla de análisis.

Aplicando el algoritmo de análisis a la cadena ‘a + a \$’, tomando la tabla de análisis del Ejemplo 3.11, obtenemos la traza del análisis LL(1):

$$\begin{aligned}
 (a + a \$, E \$, \epsilon) &\vdash (a + a \$, T E' \$, 1) \\
 &\vdash (a + a \$, F T' E' \$, 14) \\
 &\vdash (a + a \$, a T' E' \$, 148) \\
 &\vdash (+ a \$, T' E' \$, 148) \\
 &\vdash (+ a \$, E' \$, 1486) \\
 &\vdash (+ a \$, + T E' \$, 14862) \\
 &\vdash (a \$, T E' \$, 14862) \\
 &\vdash (a \$, F T' E' \$, 148624) \\
 &\vdash (a \$, a T' E' \$, 1486248) \\
 &\vdash (\$, T' E' \$, 1486248) \\
 &\vdash (\$, E' \$, 14862486) \\
 &\vdash (\$, \$, 148624863)
 \end{aligned}$$

3.2.4. Recuperación de errores en el análisis sintáctico predictivo

Durante el análisis sintáctico descendente predictivo se detecta un error cuando el primer símbolo a de la cadena que queda por analizar no concuerda con el símbolo terminal b que se encuentra en la cima de la pila, o bien, cuando en la cima de la pila tenemos un símbolo no-terminal A y en la tabla de análisis para la entrada $TA [A, a]$ tenemos *error* (o blanco).

Cuando un compilador que emplea un analizador LL(1) detecta un error, informa del error y, normalmente se “recupera” del error para continuar con el análisis sintáctico y mostrar el rsto de errores sintácticos del programa fuente analizado. A este proceso se le denomina “*recuperación del error sintáctico*”. Una posible técnica de recuperación de error es la denominada ***recuperación en modo pánico*** que consiste en saltarse símbolos de la cadena de entrada hasta que aparezca un símbolo que pertenezca a un conjunto especial de símbolos terminales denominada conjunto de ***símbolos de sincronización***.

Como conjunto de símbolos de sincronización de un no-terminal A deben seleccionarse aquellos que permitan al analizador sintáctico recuperarse de los errores y así continuar el análisis. Hay varias formas de definir el conjunto de símbolos de sincronización. Por ejemplo, podemos definir como conjunto de símbolos de sincronización para el no-terminal A al conjunto de Primeros(A). De esta manera, si durante el análisis está el no-terminal A en la cima de la pila, el terminal x en la cadena de entrada, y en la celda $[A, x]$ aparece la acción error (celda en blanco), se saltarán (quitarán) símbolos de la cadena de entrada hasta encontrar uno que pertenezca a

$\text{Sinc}(A)$, definido como $\text{Primeros}(A)$. En ese momento se continuará el análisis de la cadena con normalidad en busca de otros errores sintácticos.

De igual manera se puede definir como conjuntos símbolos de sincronización de A ($\text{Sinc}(A)$) al conjunto de $\text{Siguietes}(A)$. En este caso, cuando durante el proceso de sincronización se termine de saltar símbolos porque se ha alcanzado uno que pertenece a $\text{Siguietes}(A)$, se desapilará el símbolo A de la pila. Esto se hace porque al sincronizar con un símbolo perteneciente a $\text{Siguietes}(A)$, estamos supuestamente llegando a un estado equivalente a haber aplicado una producción para el no-terminal A .

Ejemplo 3.12.- Partiendo de la Tabla de análisis correspondiente al Ejemplo 3.11, podemos marcar los símbolos de sincronización con *sinc* (siempre que en la casilla correspondiente no exista ya una acción de derivación). Definimos el conjunto de símbolos de sincronización de cada no-terminal como la unión de los Primeros y Siguietes. Los símbolos de sincronización correspondientes a $\text{Primeros}(A)$ ya están marcados con una acción (aplicar la producción correspondiente), así que únicamente marcamos las casillas correspondientes a $\text{Siguietes}(A)$ que estuviesen en blanco todavía. Si el analizador sintáctico al buscar la entrada $TA[A, a]$ se encuentra la acción de error (o blanco) entonces informará del error y comenzará la sincronización en modo pánico. En este proceso de sincronización, se saltará símbolos de la cadena de entrada hasta alcanzar uno (x) cuya celda ($TA[A, x]$) esté marcada con una acción (habrá sincronizado con un símbolo perteneciente a $\text{Primeros}(A)$), o con *sinc* (habrá sincronizado con $\text{Siguietes}(A)$). Si es una sincronización con $\text{Primeros}(A)$, se aplicará la producción indicada en la celda, si es una sincronización con $\text{Siguietes}(A)$, deberá desapilarse el símbolo A de la pila y continuar el análisis.

	a	+	*	()	\$
E	(TE',1)			(TE',1)	<u>sinc</u>	<u>sinc</u>
E'		(+TE',2)			(ε,3)	(ε,3)
T	(FT',4)	<u>sinc</u>		(FT',4)	<u>sinc</u>	<u>sinc</u>
T'		(ε,6)	(*FT',5)		(ε,6)	(ε,6)
F	(a,8)	<u>sinc</u>	<u>sinc</u>	((E),7)	<u>sinc</u>	<u>sinc</u>
a	sacar					
+		sacar				
*			sacar			
(sacar		
)					sacar	
\$						aceptar

$$(a * + a \$, E \$, \epsilon) \mid\!\!\!-\! (a * + a \$, T E' \$, 1)$$

$$\mid\!\!\!-\! (a * + a \$, F T' E' \$, 14)$$

$\vdash (a * + a \$, a T' E' \$, 148)$
 $\vdash (* + a \$, T' E' \$, 148)$
 $\vdash (* + a \$, * F T' E' \$, 1485)$
 $\vdash (+ a \$, F T' E' \$, 1485) \text{ ERROR (sincronizado)}$
 $\vdash (+ a \$, T' E' \$, 1485)$
 $\vdash (+ a \$, E' \$, 14856)$
 $\vdash (+ a \$, + T E' \$, 148562)$
 $\vdash (a \$, T E' \$, 148562)$
 $\vdash (a \$, F T' E' \$, 148562)$
 $\vdash (a \$, a T' E' \$, 1485628)$
 $\vdash (\$, T' E' \$, 1485628)$
 $\vdash (\$, E' \$, 14856286)$
 $\vdash (\$, \$, 1485624863)$