

Autómatas de pilas y gramáticas libres de contexto

Pedro O. Pérez M., PhD.

Implementación de métodos computacionales
Tecnológico de Monterrey

pperezm@tec.mx

02-2022

① Autómatas de pila

Definición de autómata de pila (AP)

Definición formal de autómata de pila

② Análisis sintáctico

- Introducción

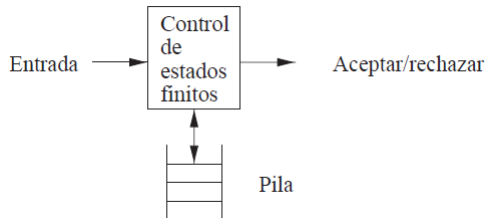
- Gramáticas libres de contexto

- Escribiendo una gramática

- Análisis Top-Down

Definición de autómatas de pila (AP)

- Existe un tipo de autómata que define los lenguajes independientes de contexto. Dicho autómata, conocido como “autómata de pila” (PDA, pushdown automata), es una expresión del autómata finito no determinista con transiciones- ϵ , el cual constituye una forma de definir los lenguajes regulares. El autómata de pila es, fundamentalmente, un AFN- ϵ con la adición de una pila. La pila se puede leer, se pueden introducir elementos en ella y extraer sólo el elemento que está en la parte superior de la misma, exactamente igual que la estructura de datos “pila”.
- La presencia de una pila significa que, a diferencia del autómata finito, el autómata de pila puede “recordar” una cantidad infinita de información.



El AP lee las entradas, un símbolo cada vez. El autómata puede observar el símbolo colocado en la parte superior de la pila y llevar a cabo su transición basándose en el estado actual, el símbolo de entrada y el símbolo que hay en la parte superior de la pila.

Alternativamente, puede hacer una transición “espontánea”, utilizando ϵ como entrada en lugar de un símbolo de entrada.

En una transición, el AP:

- ① Consuma de la entrada el símbolo que se usa en la transición. Si como entrada se utiliza ϵ , entonces no se consume ningún símbolo de entrada.
- ② Pasa a un nuevo estado, que puede o no ser el mismo que el estado anterior.
- ③ Reemplaza el símbolo de la parte superior de la pila por cualquier cadena.
 - Puede ser el mismo símbolo del tope de la pila (no hace nada).
 - Saca el tope de la pila (pop, corresponde a ϵ).
 - Cambia el primer elemento por otro (reemplazo sin *push* o *pop* o haciendo ambos).
 - Coloca un elemento en el tope de la pila posiblemente cambiando el primer elemento.

Definición formal de autómatas de pila (AP)

La notación formal de un AP incluye siete componentes. Escribiremos la especificación de un autómata de pila P de la siguiente forma:

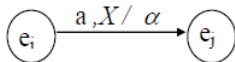
$$P = (Q, \Sigma, \Gamma, \delta, q_0, Z_0, F).$$

El significado de cada uno de los componentes es el siguiente:

- 1 Q : Un conjunto finito de estados, como los estados de un autómata finito.
- 2 Σ : Un conjunto finito de *símbolos de entrada*, también análogo al componente correspondiente de un autómata finito.
- 3 Γ : Un *alfabeto de pila* finito. Este componente, que no tiene análogo en los autómatas finitos, es el conjunto de símbolos que pueden introducirse en la pila.

- 4 δ : La *función de transición*. Como en el autómata finito, δ controla el comportamiento del autómata. Formalmente, δ toma como argumento $\delta(q, a, X)$, donde q es un estado de Q , a es cualquier símbolo de entrada de Σ o ϵ y X es un símbolo de la pila que pertenece a Γ . La salida de δ es un conjunto finito de pares (p, γ) , donde P es el nuevo estado y γ es la cadena de símbolos de la pila que reemplaza a X en la parte superior de la pila.
- 5 q_0 : El *Estado inicial*. El autómata de pila se encuentra en este estado antes de realizar ninguna transición.
- 6 Z_0 : El *Símbolo inicial*. Inicialmente, la pila del autómata consta de una instancia de este símbolo.
- 7 F : El conjunto de *estados de aceptación o estados finales*.

La función de transición de estados de un AP puede ser representada por un grafo donde los nodos representan los estados y los arcos transiciones.



Si el estado actual es e_i y la cabeza lectora apunta al símbolo a , y el tope de la pila es X , entonces cambia al estado e_j , avanza la cabeza lectura y sustituye el símbolo del tope, X , de la pila por la cadena α . Por ejemplo:

- Si $\alpha = ZYX$, deja X , agrega Y y Z (tope = Z), donde $X, Y, Z \in P$.
- Si $\alpha = XX$, deja X , agrega X (tope = X).
- Si $\alpha = X$, deja X como el mismo tope (no altera la pila).
- Si $\alpha = \epsilon$, elimina X , y el nuevo tope es el símbolo por debajo de la pila.

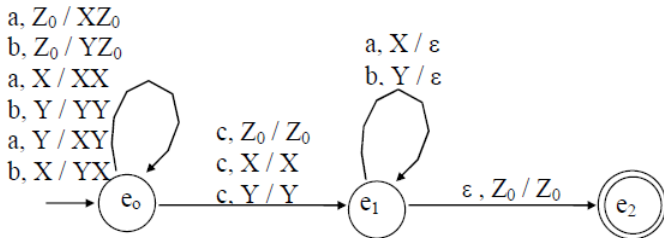
$$L_1 = \{wcw^R \mid w \in \{a, b\}^*\}$$

$$APD_1 = \langle \{e_0, e_1, e_2\}, \{a, b, c\}, \{X, Y, Z_0\}, \delta, e_0, Z_0, \{e_2\} \rangle$$

δ

¿Cuál cadena es
aceptada por APD_1 ?

- abcba
- c
- abcab
- a



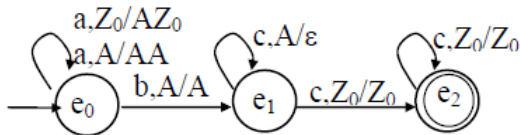
$$L_2 = \{a^i bc^k \mid i, k \geq 1 \text{ y } i < k\}$$

$$APD_2 = \langle \{e_0, e_1, e_2\}, \{a, b, c\}, \{A, Z_0\}, \delta, e_0, Z_0, \{e_2\} \rangle$$

 δ

¿Cuál cadena es aceptada por APD_2 ?

- aacbbb
- aacbb



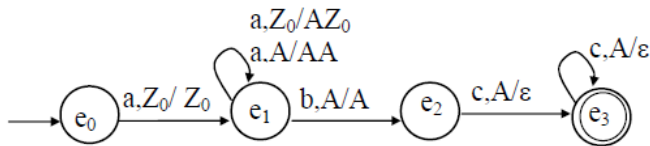
$$L_3 = \{a^i bc^k \mid i, k \geq 1 \text{ y } i > k\}$$

$$APD_3 = \langle \{e_0, e_1, e_2, e_3\}, \{a, b, c\}, \{A, Z_0\}, \delta, e_0, Z_0, \{e_3\} \rangle$$

 δ

¿Cuál cadena es aceptada por APD_1 ?

- aacbbb
- aacbb



Un ejemplo: analizador sintáctico

La sintaxis de las construcciones de un lenguaje de programa deben ser especificados a través de una gramática libre de contexto o en la notación BNF (Backus-Naur Form). Pero, ¿cuáles son sus beneficios?

- Una gramática da una especificación precisa y fácil de entender de un lenguaje de programación.
- Para ciertas gramáticas, nosotros podemos construir automáticamente un analizador sintáctico eficiente; con el beneficio adicional de que el proceso de construcción nos puede revelar ambigüedades sintácticas.
- La estructura de una gramática es útil para una correcta traducción de código fuente a código objeto y para la detección de errores.
- La gramática permite que un lenguaje se desarrolle o evolucione iterativamente.

Existen tres tipos generales de analizadores sintácticos:

- Universal (Algoritmos Cocke-Younger-Kasami y Earley).
- Top-down (LL).
- Bottom-up (LR).

En cualquier caso, la entrada del analizador sintáctico es revisada de izquierda a derecha, un símbolo a la vez.

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

Ejemplo de gramática LR

$$\begin{aligned}
 E &\rightarrow T E' \\
 E' &\rightarrow + T E' \mid \epsilon \\
 T &\rightarrow F T' \\
 T' &\rightarrow * F T' \mid \epsilon \\
 F &\rightarrow (E) \mid \mathbf{id}
 \end{aligned}$$

Ejemplo de gramática LL

Definición formal de una gramática libre de contexto

Un gramática libre de contexto consiste de:

- **Terminales.** Son los símbolos básicos de cuales están formados las cadenas (strings).
- **No terminales.** Son las variables sintácticas que hacen referencia a un conjunto de cadenas.
- **Símbolo inicial.**
- **Producciones.** Una producción consiste en;
 - Un no terminal llamado *cabeza* o *lado izquierdo*.
 - El símbolo \rightarrow .
 - Un *cuerpo* o *lado derecho* que consiste en cero o más terminales y no terminales.

<i>expression</i>	\rightarrow	<i>expression</i> + <i>term</i>
<i>expression</i>	\rightarrow	<i>expression</i> - <i>term</i>
<i>expression</i>	\rightarrow	<i>term</i>
<i>term</i>	\rightarrow	<i>term</i> * <i>factor</i>
<i>term</i>	\rightarrow	<i>term</i> / <i>factor</i>
<i>term</i>	\rightarrow	<i>factor</i>
<i>factor</i>	\rightarrow	(<i>expression</i>)
<i>factor</i>	\rightarrow	id

Figure 4.2: Grammar for simple arithmetic expressions

1. These symbols are terminals:
 - (a) Lowercase letters early in the alphabet, such as *a*, *b*, *c*.
 - (b) Operator symbols such as $+$, $*$, and so on.
 - (c) Punctuation symbols such as parentheses, comma, and so on.
 - (d) The digits $0, 1, \dots, 9$.
 - (e) Boldface strings such as **id** or **if**, each of which represents a single terminal symbol.

2. These symbols are nonterminals:

- (a) Uppercase letters early in the alphabet, such as A , B , C .
- (b) The letter S , which, when it appears, is usually the start symbol.
- (c) Lowercase, italic names such as *expr* or *stmt*.
- (d) When discussing programming constructs, uppercase letters may be used to represent nonterminals for the constructs. For example, nonterminals for expressions, terms, and factors are often represented by E , T , and F , respectively.

3. Uppercase letters late in the alphabet, such as X , Y , Z , represent *grammar symbols*; that is, either nonterminals or terminals.
4. Lowercase letters late in the alphabet, chiefly u, v, \dots, z , represent (possibly empty) strings of terminals.
5. Lowercase Greek letters, α, β, γ for example, represent (possibly empty) strings of grammar symbols. Thus, a generic production can be written as $A \rightarrow \alpha$, where A is the head and α the body.
6. A set of productions $A \rightarrow \alpha_1, A \rightarrow \alpha_2, \dots, A \rightarrow \alpha_k$ with a common head A (call them *A-productions*), may be written $A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_k$. Call $\alpha_1, \alpha_2, \dots, \alpha_k$ the *alternatives* for A .
7. Unless stated otherwise, the head of the first production is the start symbol.

To understand how parsers work, we shall consider derivations in which the nonterminal to be replaced at each step is chosen as follows:

1. In *leftmost* derivations, the leftmost nonterminal in each sentential is always chosen. If $\alpha \Rightarrow \beta$ is a step in which the leftmost nonterminal in α is replaced, we write $\alpha \Rightarrow_{lm} \beta$.
2. In *rightmost* derivations, the rightmost nonterminal is always chosen; we write $\alpha \Rightarrow_{rm} \beta$ in this case.

$$E \rightarrow E + E \mid E * E \mid - E \mid (E) \mid \mathbf{id} \\ \quad \quad \quad -(\mathbf{id} + \mathbf{id})$$

Árboles de análisis y derivaciones

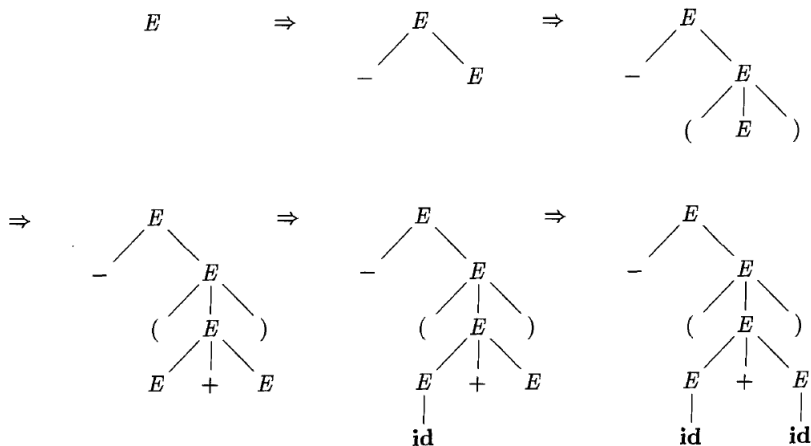


Figure 4.4: Sequence of parse trees for derivation (4.8)

Si una gramática produce más de un árbol de análisis para una misma entrada, entonces es una gramática ambigua.

Gramáticas libres de contexto vs. expresiones regulares.

Cada construcción que pueda ser descrita por una expresión regular puede ser descrita con una gramática libre de contexto, pero no lo contrario.
Alternativamente, cada lenguaje regular es un lenguaje libre de contexto, pero no lo contrario.

- Para cada estado i de un AFN, crea un no terminal A_i .
- Si el estado i tiene una transición al estado j con la entrada a , agregar la producción $A_i \rightarrow aA_j$. Si el estado i va al estado j con una transición ϵ , agrega la producción $A_i \rightarrow A_j$.
- Si i es un estado de aceptación, agrega $A_i \rightarrow \epsilon$.
- Si i es el estado inicial, haz A_i el símbolo inicial de la gramática.

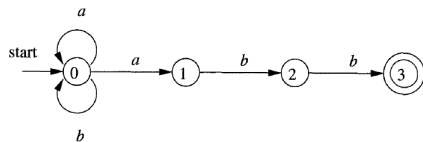


Figure 3.24: A nondeterministic finite automaton

Exercise 4.2.1: Consider the context-free grammar:

$$S \rightarrow S S + \mid S S * \mid a$$

and the string $aa + a*$.

- a) Give a leftmost derivation for the string.
- b) Give a rightmost derivation for the string.
- c) Give a parse tree for the string.
- ! d) Is the grammar ambiguous or unambiguous? Justify your answer.
- ! e) Describe the language generated by this grammar.

Si todo lo que puede ser descrito por una expresión regular también puede ser descrito por una gramática, ¿porqué usamos expresiones regulares para definir el análisis léxico de un lenguaje?

- Separando la estructura sintáctica de un lenguaje en un componente léxico y otro no léxico provee de una conveniente manera de modularizar el “Front End” en dos componentes manejables.
- Las reglas léxicas de un lenguaje son frecuentemente muy simples, no se requiere una notación tan poderosa como las gramáticas.
- Las expresiones regulares generalmente proveen de una notación más concisa y simple para los tokens.
- Analizadores léxico eficientes pueden ser contruidos automáticamente a partir de expresiones regulares que de gramáticas.

$$\begin{array}{lcl} stmt & \rightarrow & \text{if } expr \text{ then } stmt \\ & | & \text{if } expr \text{ then } stmt \text{ else } stmt \\ & | & \text{other} \end{array}$$

if E_1 then S_1 else if E_2 then S_2 else S_3

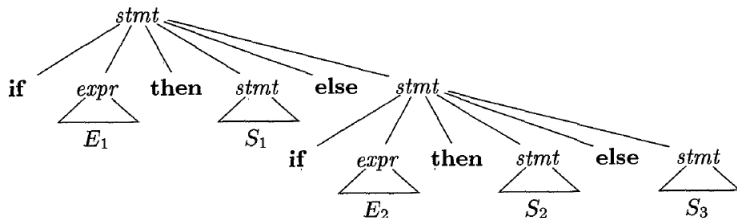


Figure 4.8: Parse tree for a conditional statement

$$\text{if } E_1 \text{ then if } E_2 \text{ then } S_1 \text{ else } S_2 \quad (4.15)$$

has the two parse trees shown in Fig. 4.9.

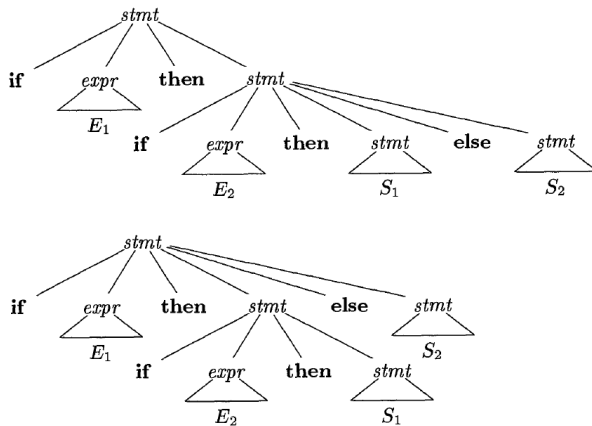


Figure 4.9: Two parse trees for an ambiguous sentence

<i>stmt</i>	→	<i>matched_stmt</i>
		<i>open_stmt</i>
<i>matched_stmt</i>	→	if <i>expr</i> then <i>matched_stmt</i> else <i>matched_stmt</i>
		other
<i>open_stmt</i>	→	if <i>expr</i> then <i>stmt</i>
		if <i>expr</i> then <i>matched_stmt</i> else <i>open_stmt</i>

Figure 4.10: Unambiguous grammar for if-then-else statements

Eliminando la recursión por la izquierda

$$\begin{array}{l} A \rightarrow \beta A' \\ A' \rightarrow \alpha A' \mid \epsilon \end{array}$$

$$\begin{array}{l} E \rightarrow E + T \mid E - T \mid T \\ T \rightarrow T * F \mid T / F \mid F \\ F \rightarrow (E) \mid \mathbf{id} \end{array}$$

Algorithm 4.19: Eliminating left recursion.

INPUT: Grammar G with no cycles or ϵ -productions.

OUTPUT: An equivalent grammar with no left recursion.

METHOD: Apply the algorithm in Fig. 4.11 to G . Note that the resulting non-left-recursive grammar may have ϵ -productions. \square

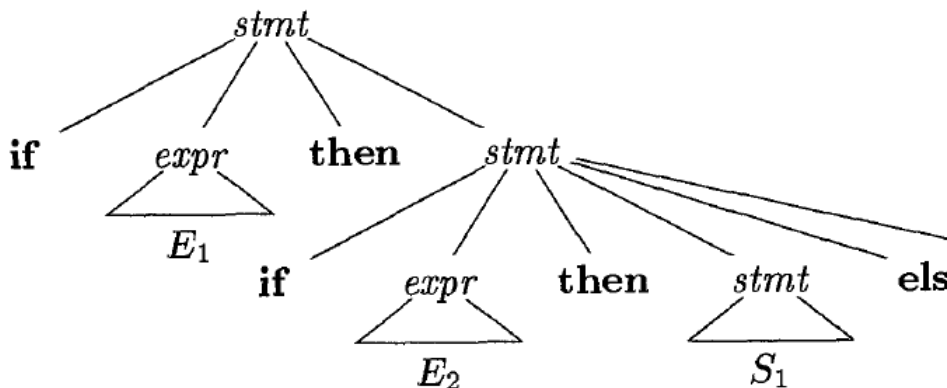
$$\begin{array}{l} S \rightarrow A a \mid b \\ A \rightarrow A c \mid S d \mid \epsilon \end{array}$$

- 1) arrange the nonterminals in some order A_1, A_2, \dots, A_n .
- 2) **for** (each i from 1 to n) {
- 3) **for** (each j from 1 to $i - 1$) {
- 4) replace each production of the form $A_i \rightarrow A_j \gamma$ by the
 productions $A_i \rightarrow \delta_1 \gamma \mid \delta_2 \gamma \mid \dots \mid \delta_k \gamma$, where
 $A_j \rightarrow \delta_1 \mid \delta_2 \mid \dots \mid \delta_k$ are all current A_j -productions
- 5) }
- 6) eliminate the immediate left recursion among the A_i -productions
- 7) }

Figure 4.11: Algorithm to eliminate left recursion from a grammar

if E_1 then if E_2 then S_1 else S_2

has the two parse trees shown in Fig. 4.9.



$$\begin{array}{l}
 A \rightarrow \beta A' \\
 A' \rightarrow \alpha A' \quad | \quad \epsilon
 \end{array}$$

Exercise 4.3.1: The following is a grammar for regular expressions over symbols a and b only, using $+$ in place of $|$ for union, to avoid conflict with the use of vertical bar as a metasyMBOL in grammars:

$$S \rightarrow S S + \mid S S * \mid a$$

- a) Left factor this grammar.
- b) Does left factoring make the grammar suitable for top-down parsing?
- c) In addition to left factoring, eliminate left recursion from the original grammar.
- d) Is the resulting grammar suitable for top-down parsing?

Exercise 4.3.1: The following is a grammar for regular expressions over symbols a and b only, using $+$ in place of $|$ for union, to avoid conflict with the use of vertical bar as a metasyMBOL in grammars:

a) $S \rightarrow 0 S 1 \mid 0 1$ with string 000111.

- a) Left factor this grammar.
- b) Does left factoring make the grammar suitable for top-down parsing?
- c) In addition to left factoring, eliminate left recursion from the original grammar.
- d) Is the resulting grammar suitable for top-down parsing?

Exercise 4.3.1: The following is a grammar for regular expressions over symbols a and b only, using $+$ in place of $|$ for union, to avoid conflict with the use of vertical bar as a metasympol in grammars:

$$! \text{ c) } S \rightarrow S (S) S \mid \epsilon \text{ with string } (()()).$$

- a) Left factor this grammar.
- b) Does left factoring make the grammar suitable for top-down parsing?
- c) In addition to left factoring, eliminate left recursion from the original grammar.
- d) Is the resulting grammar suitable for top-down parsing?

Revisar código fuente.

- $\text{FIRST}(\alpha)$, donde α es cualquier cadena de símbolos de la gramática, es el conjunto de terminales que empieza una secuencia de cadenas que derivan de α .
- $\text{FOLLOW}(A)$, para un no terminal A , es el conjunto de terminales a que pueden aparecer inmediatamente a la derecha de A en alguna sentencia.

$$\begin{array}{lcl}
E & \rightarrow & T E' \\
E' & \rightarrow & + T E' \mid \epsilon \\
T & \rightarrow & F T' \\
T' & \rightarrow & * F T' \mid \epsilon \\
F & \rightarrow & (E) \mid \text{id}
\end{array}$$

1. If X is a terminal, then $\text{FIRST}(X) = \{X\}$.
2. If X is a nonterminal and $X \rightarrow Y_1 Y_2 \cdots Y_k$ is a production for some $k \geq 1$, then place a in $\text{FIRST}(X)$ if for some i , a is in $\text{FIRST}(Y_i)$, and ϵ is in all of $\text{FIRST}(Y_1), \dots, \text{FIRST}(Y_{i-1})$; that is, $Y_1 \cdots Y_{i-1} \xRightarrow{*} \epsilon$. If ϵ is in $\text{FIRST}(Y_j)$ for all $j = 1, 2, \dots, k$, then add ϵ to $\text{FIRST}(X)$. For example, everything in $\text{FIRST}(Y_1)$ is surely in $\text{FIRST}(X)$. If Y_1 does not derive ϵ , then we add nothing more to $\text{FIRST}(X)$, but if $Y_1 \xRightarrow{*} \epsilon$, then we add $\text{FIRST}(Y_2)$, and so on.
3. If $X \rightarrow \epsilon$ is a production, then add ϵ to $\text{FIRST}(X)$.

$$\begin{array}{lcl}
E & \rightarrow & T E' \\
E' & \rightarrow & + T E' \mid \epsilon \\
T & \rightarrow & F T' \\
T' & \rightarrow & * F T' \mid \epsilon \\
F & \rightarrow & (E) \mid \text{id}
\end{array}$$

1. Place \$ in FOLLOW(S), where S is the start symbol, and \$ is the input right endmarker.
2. If there is a production $A \rightarrow \alpha B \beta$, then everything in FIRST(β) except ϵ is in FOLLOW(B).
3. If there is a production $A \rightarrow \alpha B$, or a production $A \rightarrow \alpha B \beta$, where FIRST(β) contains ϵ , then everything in FOLLOW(A) is in FOLLOW(B).

- Los analizadores predictivos, o analizadores de descenso recursivo sin “backtracking” pueden ser contruidos para una clase de gramáticas conocidas como LL(1).
- Una gramática G es LL(1), sí y solo si siendo $A \rightarrow \alpha | \beta$ son distintas producciones de G , se dan las siguientes condiciones:
 - Para ningún terminal a , ambas cadenas empiezan con el mismo símbolo.
 - A lo mas un no terminal, α o β , pueden derivar en ϵ .
 - Si $\beta \Rightarrow \epsilon$, entonces α no debe derivar ninguna cadena que empiece con una terminal en FOLLOW(A). De igual forma, $\alpha \Rightarrow \epsilon$, entonces β no puede derivar ninguna cadena que empiece con un termina en FOLLOW(A).

Para la siguiente gramática, genera:

- Los conjuntos de FIRST y FOLLOW.
- ¿Es una gramática LL(1)?

$$S \rightarrow A a$$

$$A \rightarrow B D$$

$$B \rightarrow b \mid \epsilon$$

$$D \rightarrow d \mid \epsilon$$

Algorithm 4.31: Construction of a predictive parsing table.

INPUT: Grammar G .

OUTPUT: Parsing table M .

METHOD: For each production $A \rightarrow \alpha$ of the grammar, do the following:

1. For each terminal a in $\text{FIRST}(A)$, add $A \rightarrow \alpha$ to $M[A, a]$.
2. If ϵ is in $\text{FIRST}(\alpha)$, then for each terminal b in $\text{FOLLOW}(A)$, add $A \rightarrow \alpha$ to $M[A, b]$. If ϵ is in $\text{FIRST}(\alpha)$ and $\$$ is in $\text{FOLLOW}(A)$, add $A \rightarrow \alpha$ to $M[A, \$]$ as well.

Para la siguiente gramática:

- Genera los conjuntos de FIRST y FOLLOW.

$$\begin{aligned} S &\rightarrow iEtSS' \mid a \\ S' &\rightarrow eS \mid \epsilon \\ E &\rightarrow b \end{aligned}$$

NON - TERMINAL	INPUT SYMBOL					
	a	b	e	i	t	$\$$
S	$S \rightarrow a$			$S \rightarrow iEtSS'$		
S'			$S' \rightarrow \epsilon$ $S' \rightarrow eS$			$S' \rightarrow \epsilon$
E		$E \rightarrow b$				

Figure 4.18: Parsing table M for Example 4.33

METHOD: Initially, the parser is in a configuration with $w\$$ in the input buffer and the start symbol S of G on top of the stack, above $\$$. The program in Fig. 4.20 uses the predictive parsing table M to produce a predictive parse for the input. \square

```
set  $ip$  to point to the first symbol of  $w$ ;  
set  $X$  to the top stack symbol;  
while (  $X \neq \$$  ) { /* stack is not empty */  
    if (  $X$  is  $a$  ) pop the stack and advance  $ip$ ;  
    else if (  $X$  is a terminal )  $error()$ ;  
    else if (  $M[X, a]$  is an error entry )  $error()$ ;  
    else if (  $M[X, a] = X \rightarrow Y_1 Y_2 \cdots Y_k$  ) {  
        output the production  $X \rightarrow Y_1 Y_2 \cdots Y_k$ ;  
        pop the stack;  
        push  $Y_k, Y_{k-1}, \dots, Y_1$  onto the stack, with  $Y_1$  on top;  
    }  
    set  $X$  to the top stack symbol;  
}
```

Figure 4.20: Predictive parsing algorithm

NON - TERMINAL	INPUT SYMBOL					
	id	+	*	()	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$		
E'		$E' \rightarrow +TE'$			$E' \rightarrow \epsilon$	$E' \rightarrow \epsilon$
T	$T \rightarrow FT'$			$T \rightarrow FT'$		
T'		$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
F	$F \rightarrow \text{id}$			$F \rightarrow (E)$		

Figure 4.17: Parsing table M for Example 4.32

id + id * id

Para la siguiente gramática, genera:

- La tabla de análisis predictivo. Emplea el algoritmo 4.31.
- ¿Se reconoce la entrada **bda**? Emplea el algoritmo 4.34.
- ¿Se reconoce la entrada **dba**? Emplea el algoritmo 4.34.

$$S \rightarrow A a$$

$$A \rightarrow B D$$

$$B \rightarrow b \mid \epsilon$$

$$D \rightarrow d \mid \epsilon$$