



4º Ingeniería Informática

II26 Procesadores de lenguaje

Analizador sintáctico

Esquema del tema

1. Introducción
2. Gramáticas incontextuales
3. Algunas construcciones de los lenguajes de programación
4. Extensiones de las gramáticas incontextuales
5. Análisis descendente
6. Implementación del analizador sintáctico
7. Resumen del tema

1. Introducción

El analizador sintáctico tiene como objetivo encontrar las estructuras presentes en su entrada. Estas estructuras se pueden representar mediante el árbol de análisis sintáctico, que explica cómo se puede derivar la cadena de entrada en la gramática que especifica el lenguaje. Aunque en la práctica es habitual que el árbol de análisis no llegue a construirse, se trata de una abstracción que nos permite entender mejor todo el proceso.

Para construir la especificación sintáctica de los lenguajes de programación, se suelen emplear gramáticas incontextuales, generalmente restringidas para que el análisis se pueda realizar de manera eficiente.

Para que sea posible construir el árbol de análisis, es necesario que la entrada no presente errores sintácticos. En caso de que los haya, el analizador debe informar de su presencia adecuadamente y, si es posible, intentar continuar el análisis.

2. Gramáticas incontextuales

Ya hemos comentado que las gramáticas incontextuales se emplean para especificar la sintaxis de los lenguajes de programación. Sabemos que son un paso intermedio entre las gramáticas regulares y las sensibles al contexto. Es natural preguntarse por qué nos detenemos en ellas. Hay diversas razones:

- Tienen suficiente capacidad expresiva para representar la mayor parte de las construcciones presentes en el compilador. Además, es un formalismo que puede aumentarse con sencillez para incorporar las restricciones que no son puramente incontextuales.
- Son razonablemente sencillas de diseñar. Pese a que el poder expresivo de las gramáticas sensibles al contexto es mucho mayor, su diseño es excesivamente complejo para ser empleado en la práctica.
- Permiten un análisis eficiente en el caso general y muy eficiente si se introducen una serie de restricciones que no suponen una limitación excesiva de su poder expresivo.

Ahora recordaremos los conceptos más importantes de las gramáticas incontextuales.

2.1. Definición

Como seguramente recordarás, una gramática incontextual es una cuádrupla $G = (N, \Sigma, P, \langle S \rangle)$ donde

- N es un alfabeto de *no terminales*.
- Σ es un alfabeto de *terminales*, disjunto con N .
- $P \subseteq N \times (N \cup \Sigma)^*$ es un conjunto de *producciones* o *reglas*.
- $\langle S \rangle \in N$ es el *símbolo inicial*.

Como ves, las producciones de una gramática incontextual tienen únicamente un no terminal en la parte izquierda.

2.2. Derivaciones

El funcionamiento de las gramáticas se basa en el concepto de *derivación*. Comenzando por una cadena formada únicamente por el símbolo inicial, se aplican repetidamente las reglas de P hasta llegar a una cadena formada únicamente por terminales. Aplicar una regla consiste simplemente en encontrar, dentro de la cadena actual, la parte izquierda de la regla y sustituirla por la parte derecha correspondiente.

Por ejemplo, sea la gramática: $G = (\{\langle S \rangle, \langle E \rangle\}, \{\text{id}, +, *\}, P, \langle S \rangle)$, con

$$\begin{aligned}
 P &= \{ \langle S \rangle \rightarrow \langle E \rangle, \\
 &\quad \langle E \rangle \rightarrow \langle E \rangle + \langle E \rangle, \\
 &\quad \langle E \rangle \rightarrow \langle E \rangle * \langle E \rangle, \\
 &\quad \langle E \rangle \rightarrow \text{id} \};
 \end{aligned}$$

para derivar **id+id*id** podemos hacer:

$$\begin{aligned}
 \langle S \rangle &\Rightarrow \langle E \rangle \\
 &\Rightarrow \langle E \rangle + \langle E \rangle \\
 &\Rightarrow \langle E \rangle + \langle E \rangle * \langle E \rangle \\
 &\Rightarrow \text{id} + \langle E \rangle * \langle E \rangle \\
 &\Rightarrow \text{id} + \text{id} * \langle E \rangle \\
 &\Rightarrow \text{id} + \text{id} * \text{id}
 \end{aligned}$$

EJERCICIO 1

Encuentra otra derivación de **id+id*id**.

EJERCICIO* 2

¿Cuántas derivaciones distintas tiene la cadena **id+id*id**?

De una manera más formal, tenemos las siguientes definiciones:

Derivación directa: η deriva directamente γ en G si $\eta = \omega_1 \langle A \rangle \omega_2$, $\gamma = \omega_1 \beta \omega_2$ y $\langle A \rangle \rightarrow \beta \in P$.

Lo denotamos con $\eta \Rightarrow \gamma$.

Derivación: η deriva γ (en G) si hay una secuencia de cadenas $\{\xi_1, \xi_2, \dots, \xi_n\}$ de $(N \cup \Sigma)^*$ tal que $\eta = \xi_1$, $\gamma = \xi_n$ y $\xi_i \Rightarrow \xi_{i+1}$ para $1 \leq i < n$. Lo denotamos con $\eta \xRightarrow{*} \gamma$.

Forma sentencial: Cadena α de $(N \cup \Sigma)^*$ tal que $\langle S \rangle \xRightarrow{*} \alpha$.

Sentencia: forma sentencial perteneciente a Σ^* .

Lenguaje definido (o generado) por la gramática G : $L(G) = \{x \in \Sigma^* \mid \langle S \rangle \xRightarrow{*} x\}$.

2.3. Convenios de notación

En lo sucesivo y con objeto de no tener que indicar explícitamente a qué alfabeto (N o Σ) pertenece cada símbolo o cadena, asumiremos unos convenios que permitirán deducirlo inmediatamente:

- Representamos los **terminales** con un tipo de letra mecanográfico (**a**, **c**, **+**, **:=**), en negrita (**identificador**, **entero**) o entrecomillado (“”, “”).
- Los **no terminales** se representan encerrados entre ángulos ($\langle A \rangle$, $\langle \text{Programa} \rangle$).
- Un **símbolo que puede ser terminal o no terminal**, indistintamente, se representará con alguna de las últimas letras mayúsculas del alfabeto latino (X , Y , ...).
- Las **cadena de símbolos terminales** (elementos de Σ^*) se representan con la últimas letras minúsculas del alfabeto latino (u , v , w , ...).
- Las **cadena de símbolos terminales o no terminales** (elementos de $(N \cup \Sigma)^*$) se representan con letras minúsculas del alfabeto griego (α , β , γ , ...).

Además, cuando escribamos una gramática, sólo mostraremos sus reglas y asumiremos que el símbolo inicial es el que está en la parte izquierda de la primera regla.

También asumiremos que la gramática no tiene ni *símbolos inútiles* ni *reglas inútiles*. Esto es, para cada símbolo o regla, habrá al menos una derivación de una sentencia del lenguaje que contenga ese símbolo o regla.

2.4. Algunas definiciones adicionales

Si una producción tiene el símbolo $\langle A \rangle$ en su parte izquierda, la llamamos $\langle A \rangle$ -producción. Si su parte derecha es únicamente la cadena vacía, decimos que es una λ -producción.

Si α deriva β en uno o más pasos, escribiremos $\alpha \xRightarrow{+} \beta$.

Las producciones del tipo $\langle A \rangle \rightarrow \langle B \rangle$ se llaman *producciones simples*. Si tenemos que para algún no terminal $\langle A \rangle \xRightarrow{+} \langle A \rangle$, decimos que la gramática tiene *ciclos*. Si $\langle A \rangle \xRightarrow{+} \langle A \rangle \beta$ diremos que la gramática tiene *recursividad por la izquierda*. Análogamente, la gramática tiene *recursividad por la derecha* si $\langle A \rangle \xRightarrow{+} \beta \langle A \rangle$.

2.5. Árboles de análisis y ambigüedad

Supongamos que tenemos la gramática $G = (\{ \langle S \rangle \}, \{ \text{id}, *, + \}, P, \langle S \rangle)$ con las siguientes producciones:

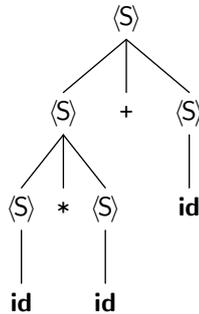
$$\begin{aligned} \langle S \rangle &\rightarrow \langle S \rangle + \langle S \rangle, \\ \langle S \rangle &\rightarrow \langle S \rangle * \langle S \rangle, \\ \langle S \rangle &\rightarrow \text{id}. \end{aligned}$$

La sentencia **id*id+id**, tiene distintas derivaciones:

$$\begin{aligned} \langle S \rangle &\Rightarrow \langle S \rangle + \langle S \rangle \Rightarrow \langle S \rangle * \langle S \rangle + \langle S \rangle \Rightarrow \langle S \rangle * \text{id} + \langle S \rangle \Rightarrow \text{id} * \text{id} + \langle S \rangle \Rightarrow \text{id} * \text{id} + \text{id}, \\ \langle S \rangle &\Rightarrow \langle S \rangle + \langle S \rangle \Rightarrow \langle S \rangle + \text{id} \Rightarrow \langle S \rangle * \langle S \rangle + \text{id} \Rightarrow \langle S \rangle * \text{id} + \text{id} \Rightarrow \text{id} * \text{id} + \text{id}, \\ \langle S \rangle &\Rightarrow \langle S \rangle + \langle S \rangle \Rightarrow \langle S \rangle * \langle S \rangle + \langle S \rangle \Rightarrow \text{id} * \langle S \rangle + \langle S \rangle \Rightarrow \text{id} * \text{id} + \langle S \rangle \Rightarrow \text{id} * \text{id} + \text{id}, \\ \langle S \rangle &\Rightarrow \langle S \rangle + \langle S \rangle \Rightarrow \langle S \rangle * \langle S \rangle + \langle S \rangle \Rightarrow \langle S \rangle * \text{id} + \langle S \rangle \Rightarrow \langle S \rangle * \text{id} + \text{id} \Rightarrow \text{id} * \text{id} + \text{id}, \\ &\dots \end{aligned}$$

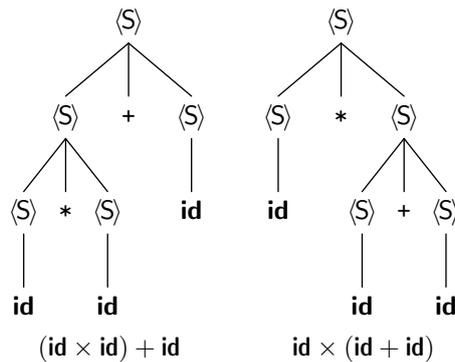
Sin embargo, en cierto sentido, todas estas derivaciones representan una misma “estructura”. De alguna manera nos transmiten la idea de que se está sumando el producto de los dos primeros **id**

con el tercero. Podemos representar esto de forma compacta mediante un árbol como el siguiente:



Esto es lo que se conoce como *árbol de análisis* o *de derivación*. Intuitivamente, lo que hacemos es representar “en paralelo” las reglas necesarias para derivar la cadena de entrada. La raíz es el símbolo inicial de la gramática y cada nodo interior representa una producción: está etiquetado con un no terminal $\langle A \rangle$ y sus hijos de izquierda a derecha están etiquetados con X_1, X_2, \dots, X_n de modo que $\langle A \rangle \rightarrow X_1 X_2 \dots X_n$ es una regla de la gramática.

Puede suceder que una misma cadena tenga asociado más de un árbol de derivación. Por ejemplo, la cadena anterior tiene asociados dos árboles, cada uno con un significado distinto:



Decimos que una sentencia es *ambigua* (con respecto a una gramática) si tiene más de un árbol de análisis. Si al menos una de las sentencias que se derivan de una gramática es ambigua, decimos que la gramática es ambigua. En algunos casos, es posible “resolver” la ambigüedad modificando la gramática adecuadamente, como veremos para el caso de las expresiones aritméticas. Sin embargo, existen lenguajes que tienen la propiedad de ser *inherentemente* ambiguos; no existe ninguna gramática no ambigua que los genere.

Está claro que la ambigüedad es perjudicial para un lenguaje de programación y debe ser evitada. Desgraciadamente, averiguar si una gramática es ambigua es un problema indecidible. Lo que sí que podremos garantizar es que las gramáticas pertenecientes a las familias que estudiaremos son no ambiguas.

2.6. Derivaciones canónicas

Hemos visto que para una sentencia de la gramática existen en general bastantes derivaciones posibles. De entre todas ellas, se suelen destacar dos, a las que se les llama derivaciones canónicas. Estas se caracterizan por aplicar las reglas de manera sistemática sobre el no terminal más a la izquierda (*derivación canónica por la izquierda*) o más a la derecha (*derivación canónica por la derecha*).

Así una derivación canónica por la izquierda de **id*id+id** es:

$$\langle S \rangle \Rightarrow \langle S \rangle + \langle S \rangle \Rightarrow \langle S \rangle * \langle S \rangle + \langle S \rangle \Rightarrow \text{id} * \langle S \rangle + \langle S \rangle \Rightarrow \text{id} * \text{id} + \langle S \rangle \Rightarrow \text{id} * \text{id} + \text{id}.$$

Cada árbol de derivación tiene asociada una derivación canónica por la izquierda y otra por la derecha. Por eso, si la sentencia es ambigua, habrá, al menos, dos derivaciones canónicas por la izquierda y otras tantas por la derecha. Dado que la gramática del ejemplo es ambigua, tenemos otra posible derivación por la izquierda para la expresión:

$$\langle S \rangle \Rightarrow \langle S \rangle * \langle S \rangle \Rightarrow \text{id} * \langle S \rangle \Rightarrow \text{id} * \langle S \rangle + \langle S \rangle \Rightarrow \text{id} * \text{id} + \langle S \rangle \Rightarrow \text{id} * \text{id} + \text{id}.$$

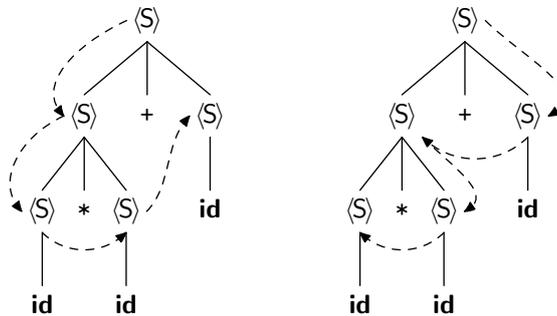
Las derivaciones canónicas por la derecha de **id*id+id** son:

$$\begin{aligned} \langle S \rangle &\Rightarrow \langle S \rangle + \langle S \rangle \Rightarrow \langle S \rangle + \text{id} \Rightarrow \langle S \rangle * \langle S \rangle + \text{id} \Rightarrow \langle S \rangle * \text{id} + \text{id} \Rightarrow \text{id} * \text{id} + \text{id}, \\ \langle S \rangle &\Rightarrow \langle S \rangle * \langle S \rangle \Rightarrow \langle S \rangle * \langle S \rangle + \langle S \rangle \Rightarrow \langle S \rangle * \langle S \rangle + \text{id} \Rightarrow \langle S \rangle * \text{id} + \text{id} \Rightarrow \text{id} * \text{id} + \text{id}. \end{aligned}$$

Es sencillo obtener, a partir del árbol de análisis, las correspondientes derivaciones canónicas por la derecha y por la izquierda:

- La derivación canónica por la izquierda se obtiene recorriendo el árbol de modo que se visita cada nodo, se escribe la regla correspondiente, y después, recursivamente, se visitan sus hijos de izquierda a derecha.
- La derivación canónica por la derecha se obtiene recorriendo el árbol de modo que se visita cada nodo, se escribe la regla correspondiente, y después, recursivamente, se visitan sus hijos de derecha a izquierda.

En nuestro ejemplo, los órdenes de visita para obtener las derivaciones canónicas por la izquierda y por la derecha, respectivamente, son:



EJERCICIO 3

Dada la gramática con las reglas:

$$\begin{aligned} \langle E \rangle &\rightarrow \langle E \rangle + \langle T \rangle | \langle T \rangle \\ \langle T \rangle &\rightarrow \langle T \rangle * \langle F \rangle | \langle F \rangle \\ \langle F \rangle &\rightarrow \text{id} \end{aligned}$$

Encuentra el árbol de análisis de **id*id+id**. A partir de este árbol, escribe la derivación canónica por la derecha y por la izquierda.

¿Existe más de un árbol de análisis para esa sentencia?

EJERCICIO* 4

¿Puede una sentencia tener infinitos árboles de análisis con una gramática?

Pista: piensa en una gramática con ciclos.

3. Algunas construcciones de los lenguajes de programación

Ahora veremos cómo se expresan algunas de las construcciones más habituales de los lenguajes de programación mediante gramáticas. Lógicamente, la lista de construcciones presentadas no será completa. Debes tomarla más como una lista de ejemplos que como una referencia. Es más, las menciones que se hacen a C y Pascal no implican que ésta sea la manera de escribir las estructuras en estos lenguajes; simplemente sirven para que tengas una idea de a qué nos referimos.

3.1. Estructura del programa

Los lenguajes tipo Pascal tienen una estructura similar a:

$$\langle \text{Programa} \rangle \rightarrow \langle \text{Cabecera} \rangle \langle \text{Bloque} \rangle .$$

con:

$$\begin{aligned} \langle \text{Cabecera} \rangle &\rightarrow \text{program id;} \\ \langle \text{Bloque} \rangle &\rightarrow \langle \text{Declaraciones} \rangle \text{ begin } \langle \text{ListaSentencias} \rangle \text{ end} \\ \langle \text{Declaraciones} \rangle &\rightarrow \langle \text{DConstantes} \rangle \langle \text{DTipos} \rangle \langle \text{DVariables} \rangle \langle \text{DSubrutinas} \rangle \end{aligned}$$

Los lenguajes tipo C tienen una estructura similar a:

$$\langle \text{Programa} \rangle \rightarrow \langle \text{ListaDeclaraciones} \rangle$$

con:

$$\begin{aligned} \langle \text{ListaDeclaraciones} \rangle &\rightarrow \langle \text{Declaración} \rangle \langle \text{ListaDeclaraciones} \rangle | \lambda \\ \langle \text{Declaración} \rangle &\rightarrow \langle \text{DeclVariable} \rangle | \langle \text{DeclTipo} \rangle | \langle \text{DeclFunción} \rangle \end{aligned}$$

EJERCICIO 5

Construye una gramática para un lenguaje tipo Pascal en que las declaraciones de constantes, tipos, variables y subprogramas puedan aparecer en cualquier orden y cualquier número de veces.

3.2. Secuencias de instrucciones

Habitualmente se distinguen dos tipos de secuencias de instrucciones:

- Sentencias *separadas* por punto y coma (Pascal):

$$\begin{aligned} \langle \text{ListaSentencias} \rangle &\rightarrow \langle \text{Sentencia} \rangle ; \langle \text{ListaSentencias} \rangle \\ \langle \text{ListaSentencias} \rangle &\rightarrow \langle \text{Sentencia} \rangle \end{aligned}$$

- Sentencias *terminadas* por punto y coma (C):

$$\begin{aligned} \langle \text{ListaSentencias} \rangle &\rightarrow \langle \text{Sentencia} \rangle \langle \text{ListaSentencias} \rangle \\ \langle \text{ListaSentencias} \rangle &\rightarrow \langle \text{Sentencia} \rangle \end{aligned}$$

En el segundo caso, las producciones para $\langle \text{Sentencia} \rangle$ deben incorporar el punto y coma, ya que, aunque informalmente se dice que C es un lenguaje en que las sentencias terminan en punto y coma, no todas lo hacen. Por ejemplo, la sentencia compuesta no termina en punto y coma, lo que hace que las sentencias de control de flujo tampoco terminen en punto y coma si afectan a una sentencia compuesta.

EJERCICIO 6

¿Qué se hace en Pascal para que una lista de sentencias separadas por punto y coma pueda, en determinadas circunstancias, terminar en un punto y coma?

3.3. Sentencia compuesta

La sentencia compuesta es simplemente una lista de sentencias encerrada entre dos marcadores adecuados. Por ejemplo, en Pascal:

$$\langle \text{Sentencia} \rangle \rightarrow \mathbf{begin} \langle \text{ListaSentencias} \rangle \mathbf{end}$$

En C es posible incluir declaraciones dentro de las sentencias compuestas:

$$\langle \text{Sentencia} \rangle \rightarrow \{ \langle \text{Declaraciones} \rangle \langle \text{ListaSentencias} \rangle \}$$

En C++ se pueden intercalar las declaraciones con las sentencias:

$$\langle \text{Sentencia} \rangle \rightarrow \{ \langle \text{ListaSentenciasYDeclaraciones} \rangle \}$$

3.4. Declaraciones de variables

Podemos distinguir dos tipos de declaraciones de variables:

- Un tipo precediendo a una lista de identificadores (C):

$$\begin{aligned} \langle \text{DeclVariables} \rangle &\rightarrow \langle \text{Tipo} \rangle \langle \text{ListaIds} \rangle ; \\ \langle \text{ListaIds} \rangle &\rightarrow \langle \text{Id} \rangle , \langle \text{ListaIds} \rangle | \langle \text{Id} \rangle \end{aligned}$$

donde $\langle \text{Id} \rangle$ es un identificador, posiblemente afectado por modificadores como * o corchetes.

- Un tipo siguiendo a una lista de identificadores (Pascal):

$$\begin{aligned} \langle \text{DeclVariables} \rangle &\rightarrow \langle \text{ListaIds} \rangle : \langle \text{Tipo} \rangle \\ \langle \text{ListaIds} \rangle &\rightarrow \mathbf{id} , \langle \text{ListaIds} \rangle | \mathbf{id} \end{aligned}$$

3.5. Declaración de subrutinas

Las subrutinas suelen declararse como una cabecera seguida de un bloque (Pascal):

$$\langle \text{DeclSubrutinas} \rangle \rightarrow \langle \text{Cabecera} \rangle \langle \text{Bloque} \rangle$$

Si no hay anidamiento, en lugar de un bloque se puede emplear una lista de declaraciones y sentencias delimitadas (C):

$$\langle \text{DeclSubrutinas} \rangle \rightarrow \langle \text{Cabecera} \rangle \{ \langle \text{Declaraciones} \rangle \langle \text{ListaSentencias} \rangle \}$$

3.6. Estructuras de control condicional

La estructura de control condicional puede llevar un fin explícito como en:

$$\begin{aligned} \langle \text{Sentencia} \rangle &\rightarrow \mathbf{if} \langle \text{Condición} \rangle \mathbf{then} \langle \text{ListaSentencias} \rangle \mathbf{endif} \\ \langle \text{Sentencia} \rangle &\rightarrow \mathbf{if} \langle \text{Condición} \rangle \mathbf{then} \langle \text{ListaSentencias} \rangle \mathbf{else} \langle \text{ListaSentencias} \rangle \mathbf{endif} \end{aligned}$$

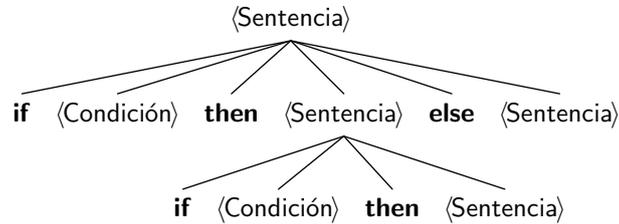
Si no lo lleva, podemos utilizar las reglas siguientes:

$$\begin{aligned} \langle \text{Sentencia} \rangle &\rightarrow \mathbf{if} \langle \text{Condición} \rangle \mathbf{then} \langle \text{Sentencia} \rangle \\ \langle \text{Sentencia} \rangle &\rightarrow \mathbf{if} \langle \text{Condición} \rangle \mathbf{then} \langle \text{Sentencia} \rangle \mathbf{else} \langle \text{Sentencia} \rangle \end{aligned}$$

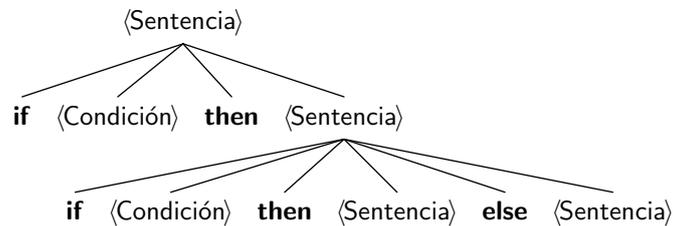
El problema de esta construcción es que es ambigua. Sea la siguiente forma sentencial:

if ⟨Condición⟩ **then if** ⟨Condición⟩ **then** ⟨Sentencia⟩ **else** ⟨Sentencia⟩

Un árbol de derivación es:



En este caso, el **else** pertenece al primer **if**. Pero también tenemos el árbol:



Este caso se corresponde a la regla habitual de asociar el **else** al **if** más cercano. Para incorporar esta regla a la gramática, tendremos dos no terminales que reflejarán si los **if-then** tienen que llevar obligatoriamente un **else** o no. El no terminal ⟨Sentencia⟩ podrá derivar condicionales con y sin parte **else** y el resto de sentencias. Por otro lado, las sentencias que se deriven del no terminal ⟨ConElse⟩ serán condicionales completos o sentencias de otro tipo. La gramática queda:

```

<Sentencia> → if <Condición> then <Sentencia>
            | if <Condición> then <ConElse> else <Sentencia>
            | <OtrasSentencias>
<ConElse>  → if <Condición> then <ConElse> else <ConElse>
            | <OtrasSentencias>
  
```

En la práctica, esta modificación complica la construcción del analizador por la presencia de prefijos comunes (lo veremos más adelante en este mismo tema) y lo que se hace es manipular el analizador sintáctico para que cumpla la regla.

EJERCICIO 7

Encuentra el árbol de derivación de **if id then if id then print else return** con la gramática:

```

<Sentencia> → if <Condición> then <Sentencia>
            | if <Condición> then <ConElse> else <Sentencia>
            | <OtrasSentencias>
<ConElse>  → if <Condición> then <ConElse> else <ConElse>
            | <OtrasSentencias>
<OtrasSentencias> → print | return
<Condición>      → id
  
```

EJERCICIO* 8

Nos han propuesto la siguiente gramática para resolver el problema de la ambigüedad del **if-then**:

$$\begin{aligned} \langle \text{Sentencia} \rangle &\rightarrow \text{if } \langle \text{Condición} \rangle \text{ then } \langle \text{Sentencia} \rangle \\ \langle \text{Sentencia} \rangle &\rightarrow \langle \text{ConElse} \rangle \\ \langle \text{ConElse} \rangle &\rightarrow \text{if } \langle \text{Condición} \rangle \text{ then } \langle \text{ConElse} \rangle \text{ else } \langle \text{Sentencia} \rangle \\ \langle \text{ConElse} \rangle &\rightarrow \langle \text{OtrasSentencias} \rangle \end{aligned}$$

Demuestra que esta gramática es ambigua.

3.7. Estructuras de control repetitivas

Las estructuras de control repetitivas suelen tener una condición al principio:

$$\langle \text{Sentencia} \rangle \rightarrow \text{while } \langle \text{Condición} \rangle \text{ do } \langle \text{Sentencia} \rangle$$

o al final:

$$\langle \text{Sentencia} \rangle \rightarrow \text{repeat } \langle \text{ListaSentencias} \rangle \text{ until } \langle \text{Condición} \rangle$$

Si queremos que en un bucle del primer tipo se admitan varias sentencias, sin utilizar la sentencia compuesta, necesitamos algún tipo de terminador:

$$\langle \text{Sentencia} \rangle \rightarrow \text{while } \langle \text{Condición} \rangle \text{ do } \langle \text{ListaSentencias} \rangle \text{ endwhile}$$

Si dentro de la estructura de control se permite la salida prematura del bucle, por ejemplo con **break**, tenemos dos alternativas para especificar que sólo pueden aparecer dentro del bucle:

- Reflejarlo directamente en la gramática. El problema es que la solución inmediata es más complicada de lo que parece a primera vista. Si probamos con:

$$\begin{aligned} \langle \text{Sentencia} \rangle &\rightarrow \text{repeat } \langle \text{SentenciasConBreak} \rangle \text{ until } \langle \text{Condición} \rangle \\ \langle \text{SentenciasConBreak} \rangle &\rightarrow \langle \text{Sentencia} \rangle \langle \text{SentenciasConBreak} \rangle \\ \langle \text{SentenciasConBreak} \rangle &\rightarrow \text{break}; \langle \text{SentenciasConBreak} \rangle \\ \langle \text{SentenciasConBreak} \rangle &\rightarrow \lambda \end{aligned}$$

Tenemos el problema de que **break** puede estar, por ejemplo, dentro de un condicional o una sentencia compuesta. Esto obliga a que esencialmente haya que replicar las reglas de $\langle \text{Sentencia} \rangle$ en un nuevo no terminal, por ejemplo $\langle \text{SentenciaOBreak} \rangle$ para incluir esta posibilidad.

- Tratar el **break** como una sentencia más y añadir una comprobación durante el análisis semántico.

Esto es un ejemplo de lo que suele suceder al diseñar el compilador: determinadas fases se pueden simplificar según cómo se hayan diseñado las otras. Aquí simplificamos el análisis sintáctico a costa de complicar ligeramente el semántico.

EJERCICIO* 9

Reescribe la gramática siguiente para excluir la posibilidad de que **break** aparezca fuera del bucle.

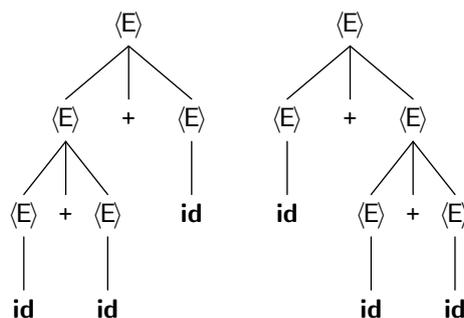
$$\begin{aligned}
 \langle \text{Programa} \rangle &\rightarrow \langle \text{Declaraciones} \rangle \langle \text{Compuesta} \rangle \\
 \langle \text{Declaraciones} \rangle &\rightarrow \text{variable } \langle \text{Declaraciones} \rangle \\
 \langle \text{Declaraciones} \rangle &\rightarrow \text{constant } \langle \text{Declaraciones} \rangle \\
 \langle \text{Declaraciones} \rangle &\rightarrow \lambda \\
 \langle \text{Compuesta} \rangle &\rightarrow \text{begin } \langle \text{ListaSentencias} \rangle \text{ end} \\
 \langle \text{ListaSentencias} \rangle &\rightarrow \langle \text{Sentencia} \rangle \langle \text{ListaSentencias} \rangle | \lambda \\
 \langle \text{Sentencia} \rangle &\rightarrow \langle \text{Compuesta} \rangle \\
 \langle \text{Sentencia} \rangle &\rightarrow \text{repeat } \langle \text{ListaSentencias} \rangle \text{ until } \langle \text{Expresión} \rangle \\
 \langle \text{Sentencia} \rangle &\rightarrow \text{if } \langle \text{Expresión} \rangle \text{ then } \langle \text{Sentencia} \rangle \\
 \langle \text{Sentencia} \rangle &\rightarrow \text{if } \langle \text{Expresión} \rangle \text{ then } \langle \text{Sentencia} \rangle \text{ else } \langle \text{Sentencia} \rangle \\
 \langle \text{Sentencia} \rangle &\rightarrow \text{break} \\
 \langle \text{Sentencia} \rangle &\rightarrow \text{other} \\
 \langle \text{Expresión} \rangle &\rightarrow \text{id} \mid \text{cte}
 \end{aligned}$$

3.8. Expresiones aritméticas

Si queremos permitir operaciones en notación *infija* con los operadores de suma, resta, producto, división y exponenciación, podemos utilizar una gramática como la siguiente:

$$\begin{aligned}
 \langle E \rangle &\rightarrow \langle E \rangle + \langle E \rangle \\
 &| \langle E \rangle - \langle E \rangle \\
 &| \langle E \rangle * \langle E \rangle \\
 &| \langle E \rangle / \langle E \rangle \\
 &| \langle E \rangle \uparrow \langle E \rangle \\
 &| \text{"("} \langle E \rangle \text{"}")} \\
 &| \text{id}
 \end{aligned}$$

Como ya te habrás dado cuenta, por su similitud con el ejemplo de la sección 2.5, esta gramática es ambigua. Por ejemplo, la expresión **id+id+id**, tiene dos árboles de análisis:

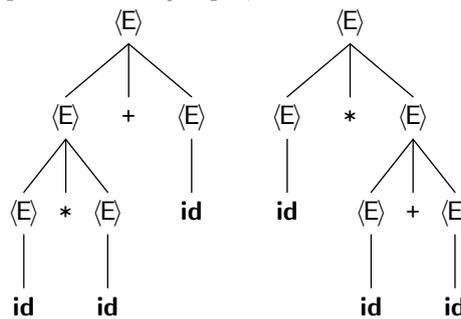


Lo habitual es elegir la primera interpretación, siguiendo las *reglas de asociatividad*. Normalmente, se sigue el siguiente convenio:

- Los operadores de suma, resta, producto y división son asociativos por la izquierda.
- El operador de exponenciación es asociativo por la derecha.

EJERCICIO 10
¿Qué árbol de análisis asociamos a **id*id*id**?

Sin embargo, la asociatividad no nos permite resolver la ambigüedad que aparece cuando hay operadores distintos en la expresión. Por ejemplo, **id*id+id** tiene dos árboles asociados:



Para estos casos utilizamos *reglas de prioridad*. Habitualmente:

- El operador de exponenciación es el más prioritario.
- Los operadores de multiplicación y división tienen la misma prioridad, que es menor que la de la exponenciación y mayor que la de la suma y resta.
- Los operadores de suma y resta tienen la misma prioridad, y ésta es la más baja.

Con estas reglas, el árbol que utilizaríamos sería el de la izquierda.

EJERCICIO 11
¿Qué árbol de análisis corresponde a la cadena **id+id*id*id*id+id**?

Para diseñar una gramática que recoja tanto la asociatividad como las prioridades, utilizaremos las siguientes ideas:

- Cada nivel de prioridad se asociará con un no terminal.
- La asociatividad por la izquierda se reflejará en recursividad por la izquierda; la asociatividad por la derecha en recursividad por la derecha.

En nuestro caso, tendremos que:

- Una expresión será bien la suma de una expresión con un término, bien la diferencia entre una expresión y un término, bien un único término.
- Un término será bien el producto de un término por un factor, bien el cociente entre un término y un factor, bien un único factor.
- Un factor será bien una base elevada a un factor, bien una única base.
- Una base será un identificador o una expresión encerrada entre paréntesis.

De esta manera, nuestra gramática es:

$$\begin{aligned}
 \langle E \rangle &\rightarrow \langle E \rangle + \langle T \rangle \mid \langle E \rangle - \langle T \rangle \mid \langle T \rangle \\
 \langle T \rangle &\rightarrow \langle T \rangle * \langle F \rangle \mid \langle T \rangle / \langle F \rangle \mid \langle F \rangle \\
 \langle F \rangle &\rightarrow \langle B \rangle \uparrow \langle F \rangle \mid \langle B \rangle \\
 \langle B \rangle &\rightarrow \text{id} \mid \langle \langle E \rangle \rangle
 \end{aligned}$$

EJERCICIO 12

Halla el árbol de análisis de $\text{id} * \text{id} \uparrow (\text{id} + \text{id})$.

El añadir operadores unarios no supone mayor complicación. Si queremos completar la gramática anterior para que acepte el operador de cambio de signo con una prioridad mayor que la exponenciación, haríamos:

$$\begin{aligned} \langle E \rangle &\rightarrow \langle E \rangle + \langle T \rangle \mid \langle E \rangle - \langle T \rangle \mid \langle T \rangle \\ \langle T \rangle &\rightarrow \langle T \rangle * \langle F \rangle \mid \langle T \rangle / \langle F \rangle \mid \langle F \rangle \\ \langle F \rangle &\rightarrow \langle B \rangle \uparrow \langle F \rangle \mid \langle B \rangle \\ \langle B \rangle &\rightarrow \text{id} \mid \text{"("} \langle E \rangle \text{"} \mid - \langle B \rangle \end{aligned}$$

EJERCICIO* 13

Analiza $-\text{id} \uparrow -\text{id}$ con la gramática anterior y con esta:

$$\begin{aligned} \langle E \rangle &\rightarrow \langle E \rangle + \langle T \rangle \mid \langle E \rangle - \langle T \rangle \mid \langle T \rangle \\ \langle T \rangle &\rightarrow \langle T \rangle * \langle F \rangle \mid \langle T \rangle / \langle F \rangle \mid \langle F \rangle \\ \langle F \rangle &\rightarrow \langle B \rangle \uparrow \langle F \rangle \mid - \langle F \rangle \mid \langle B \rangle \\ \langle B \rangle &\rightarrow \text{id} \mid \text{"("} \langle E \rangle \text{"} \end{aligned}$$

¿Por qué son distintas?

Es habitual que los operadores no aparezcan totalmente separados sino en categorías que asocien operadores con características similares. Por ejemplo, en casi todos los contextos donde puede aparecer un $+$ también puede aparecer un $-$. Esto nos lleva a agruparlos en la categoría **opad** (por operador aditivo). Si también creamos la categoría **opmul** para los operadores multiplicativos, nuestra gramática queda:

$$\begin{aligned} \langle E \rangle &\rightarrow \langle E \rangle \text{ opad } \langle T \rangle \mid \langle T \rangle \\ \langle T \rangle &\rightarrow \langle T \rangle \text{ opmul } \langle F \rangle \mid \langle F \rangle \\ \langle F \rangle &\rightarrow \langle B \rangle \uparrow \langle F \rangle \mid \langle B \rangle \\ \langle B \rangle &\rightarrow \text{id} \mid \text{"("} \langle E \rangle \text{"} \mid \text{opad } \langle B \rangle \end{aligned}$$

Fíjate en que ahora se acepta como expresión válida $+\text{id}$. Para evitarlo, lo que se hace es que el analizador semántico la rechace. Eso tiene además la ventaja de que los mensajes de error del semántico suelen ser más informativos. Este es otro ejemplo de simplificación del análisis sintáctico a costa de una ligera complicación del semántico.

EJERCICIO 14

Define conjuntos de producciones adecuados para modelar las siguientes construcciones de un lenguaje de programación:

- Una declaración de variables consiste en un identificador de tipo (**id**, **int** o **float**) seguido de una lista, separada por comas, de identificadores, cada uno de los cuales puede estar seguido del signo igual y una expresión de inicialización.
- Una declaración de función consiste en un identificador de tipo o en la palabra reservada **void**, un identificador que da nombre a la función, una lista de parámetros formales entre paréntesis y un cuerpo de función. Una lista de parámetros formales es una sucesión (posiblemente vacía) de parámetros formales separados por comas. Un argumento formal empieza opcionalmente por la palabra **var** seguida de un identificador, dos puntos y un identificador de tipo.
- Un elemento de un vector se referencia mediante un identificador seguido de uno o más índices. Cada índice es una expresión aritmética encerrada entre corchetes.

4. Extensiones de las gramáticas incontextuales

Ahora veremos cómo podemos extender la definición de gramática incontextual de modo que resulte más fácil escribirlas. La extensión que utilizaremos tiene el mismo poder expresivo que la definición habitual, por lo que conservará sus propiedades.

4.1. Definiciones

Una *gramática con partes derechas regulares* (GPDR) es una cuádrupla $G = (N, \Sigma, P, \langle S \rangle)$ donde

- N es un alfabeto de *no terminales*.
- Σ es un alfabeto de *terminales*, disjunto con N .
- $P \subseteq N \times \mathcal{R}(N \cup \Sigma)$ es un conjunto de *producciones* o *reglas*.
- $\langle S \rangle \in N$ es el *símbolo inicial*.

Usamos $\mathcal{R}(N \cup \Sigma)$ para representar el conjunto de las expresiones regulares básicas sobre $N \cup \Sigma$. Recuerda que las expresiones regulares básicas se construían a partir del conjunto vacío, la cadena vacía, los símbolos del alfabeto y las operaciones de disyunción, concatenación y clausura. Para simplificar la formulación, en esta parte, excluirémos el conjunto vacío del conjunto de expresiones regulares.

Así pues, una gramática con partes derechas regulares es similar a una gramática incontextual con la diferencia de que la parte derecha de las reglas contiene expresiones regulares formadas a partir de terminales y no terminales.

Dado que hemos modificado la definición de gramática, tendremos que cambiar la definición de derivación. Ahora las formas sentenciales son expresiones regulares sobre $N \cup \Sigma$. Las reglas que se aplican son:

1. $\alpha \langle A \rangle \beta \Rightarrow \alpha \gamma \beta$, si $\langle A \rangle \rightarrow \gamma$ pertenece a P . Esta es la definición habitual.
2. $\alpha(\rho_1 | \rho_2 | \dots | \rho_n)\beta \Rightarrow \alpha \rho_i \beta$, si $1 \leq i \leq n$. Es decir, si tenemos una disyunción, podemos escoger una cualquiera de las alternativas.
3. $\alpha(\rho)^* \beta \Rightarrow \alpha \beta$. Si tenemos una clausura, podemos sustituirla por la cadena vacía.
4. $\alpha(\rho)^* \beta \Rightarrow \alpha \rho(\rho)^* \beta$. Si tenemos una clausura, podemos sustituirla por una copia de su base y una copia de ella misma.

Como puedes imaginar, las reglas referentes a la clausura son las que dan su utilidad a las gramáticas con partes derechas regulares.

El lenguaje definido por una gramática con partes derechas regulares es:

$$L(G) = \{x \in \Sigma^* \mid \langle S \rangle \xRightarrow{*} x\}.$$

Esta es la misma definición que para las gramáticas incontextuales, como era de esperar.

Una GPDR para las declaraciones de variables consistentes en listas de identificadores seguidas por un tipo entero o real podría ser:

$$\begin{aligned} \langle \text{DeclVariables} \rangle &\rightarrow \langle \text{ListaVariables} \rangle : (\text{integer} \mid \text{float}); \\ \langle \text{ListaVariables} \rangle &\rightarrow \text{id} (, \text{id})^* \end{aligned}$$

Con esta gramática, la derivación por la izquierda de `id, id: integer;` es:

$$\begin{aligned}
 \langle \text{DeclVariables} \rangle &\xRightarrow{1} \langle \text{ListaVariables} \rangle : (\text{integer} \mid \text{float}); \\
 &\xRightarrow{1} \text{id} (, \text{id})^* : (\text{integer} \mid \text{float}); \\
 &\xRightarrow{4} \text{id} , \text{id} (, \text{id})^* : (\text{integer} \mid \text{float}); \\
 &\xRightarrow{3} \text{id} , \text{id} : (\text{integer} \mid \text{float}); \\
 &\xRightarrow{2} \text{id} , \text{id} : \text{integer};
 \end{aligned}$$

(El número encima de las flechas indica cuál de las reglas para la derivación se ha aplicado.)

4.2. Equivalencia con las gramáticas incontextuales

Hemos afirmado que las GPDR tienen el mismo poder expresivo que las incontextuales. Esto quiere decir que si un lenguaje tiene una gramática incontextual que lo represente, también tendrá una GPDR. Análogamente, si un lenguaje tiene una GPDR, también tendrá una gramática incontextual. Ahora justificaremos estas afirmaciones.

En primer lugar, parece claro que toda gramática incontextual es un caso particular de GPDR. Así la primera parte está trivialmente demostrada.

En cuanto a la segunda parte, utilizaremos un procedimiento constructivo. Empezamos con las reglas de nuestra gramática. Iremos sustituyendo aquellas reglas que utilicen extensiones por otras que tengan menos extensiones. Finalmente, todas las reglas tendrán en la parte derecha únicamente secuencias de terminales y no terminales, con lo que habremos obtenido la gramática incontextual equivalente.

Si la regla tiene la forma $\langle A \rangle \rightarrow \alpha(\rho_1|\rho_2|\dots|\rho_n)\beta$, creamos un nuevo no terminal $\langle N \rangle$ y sustituimos la regla por el conjunto de reglas:

$$\begin{aligned}
 \langle A \rangle &\rightarrow \alpha\langle N \rangle\beta \\
 \langle N \rangle &\rightarrow \rho_1 \\
 \langle N \rangle &\rightarrow \rho_2 \\
 &\dots \\
 \langle N \rangle &\rightarrow \rho_n
 \end{aligned}$$

Si la regla tiene la forma $\langle A \rangle \rightarrow \alpha(\rho)^*\beta$, creamos un nuevo no terminal $\langle N \rangle$ y sustituimos la regla por el conjunto de reglas:

$$\begin{aligned}
 \langle A \rangle &\rightarrow \alpha\langle N \rangle\beta \\
 \langle N \rangle &\rightarrow \rho\langle N \rangle \\
 \langle N \rangle &\rightarrow \lambda
 \end{aligned}$$

Es fácil demostrar que estas transformaciones no cambian el lenguaje generado por la gramática. También está claro que cada vez que aplicamos una de las transformaciones se reduce en uno el número de extensiones de la gramática. De esta manera, terminaremos con una gramática incontextual sin partes derechas regulares.

Así, la gramática que hemos empleado para las declaraciones de variables es equivalente a:

$$\begin{aligned}
 \langle \text{DeclVariables} \rangle &\rightarrow \langle \text{ListaVariables} \rangle : \langle \text{DeclVariables2} \rangle ; \\
 \langle \text{DeclVariables2} \rangle &\rightarrow \text{integer} \mid \text{float} \\
 \langle \text{ListaVariables} \rangle &\rightarrow \text{id} \langle \text{ListaVariables2} \rangle \\
 \langle \text{ListaVariables2} \rangle &\rightarrow , \text{id} \langle \text{ListaVariables2} \rangle \\
 \langle \text{ListaVariables2} \rangle &\rightarrow \lambda
 \end{aligned}$$

EJERCICIO 15

Transforma la siguiente GPDR en una gramática incontextual:

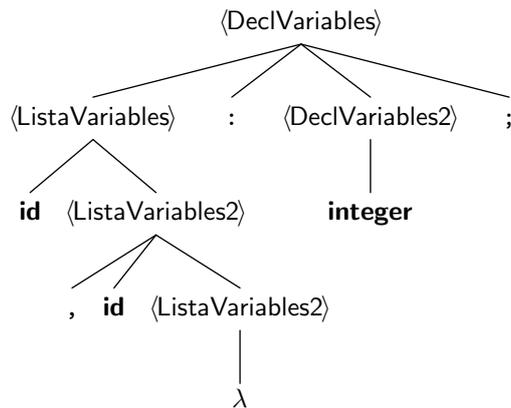
$$\begin{aligned} \langle S \rangle &\rightarrow ((\text{print } \langle \text{Expr} \rangle | \text{int } \langle \text{Listald} \rangle);)^* \\ \langle \text{Expr} \rangle &\rightarrow \text{id} (+ \text{id})^* \\ \langle \text{Listald} \rangle &\rightarrow \text{id} (, \text{id})^* \end{aligned}$$

Encuentra una derivación de **int id, id; print id;** con la gramática original y con la nueva.

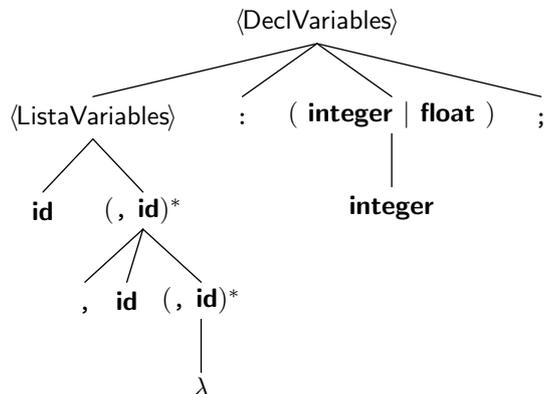
4.3. Árboles de análisis

A la hora de dibujar el árbol de derivación de una sentencia del lenguaje, existen dos filosofías, que reflejan dos formas de interpretar las GPDR. Una opción es dibujar el árbol como si se hubiera transformado la GPDR en una gramática incontextual equivalente. La otra opción es considerar que las expresiones regulares en las partes derechas suponen de hecho condensar un número potencialmente infinito de reglas.

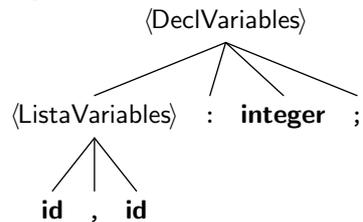
Según la primera interpretación, el árbol de análisis de **id, id: integer;** con la GPDR para las declaraciones sería:



Si no queremos introducir los nuevos no terminales y mantenemos las expresiones regulares, podemos escribir:



Sin embargo, si utilizamos la segunda interpretación, obtenemos:



Esta última interpretación es la que seguiremos: la primera manera de construir los árboles introduce una estructura que realmente no estaba presente en la gramática.

EJERCICIO 16

Construye el árbol de análisis de `int id, id; print id;` con la gramática del ejercicio 15.

4.4. Reescritura de algunas construcciones

Ahora repasaremos algunas de las construcciones que hemos visto antes, adaptándolas a las GPDR.

4.4.1. Secuencias de instrucciones

Las secuencias de instrucciones separadas por punto y coma, las podemos expresar mediante la regla:

$$\langle \text{ListaSentencias} \rangle \rightarrow \langle \text{Sentencia} \rangle (; \langle \text{Sentencia} \rangle)^*$$

Si queremos sentencias terminadas en punto y coma, con el punto y coma explícito:

$$\langle \text{ListaSentencias} \rangle \rightarrow (\langle \text{Sentencia} \rangle ;)^*$$

EJERCICIO 17

La regla anterior permite que la secuencia de sentencias sea vacía, ¿cómo puedes hacer que haya al menos una sentencia?

4.4.2. Declaraciones de variables

Si tenemos un tipo precediendo a una lista de identificadores (C):

$$\langle \text{DeclVariables} \rangle \rightarrow \langle \text{Tipo} \rangle \langle \text{Id} \rangle (, \langle \text{Id} \rangle)^* ;$$

El caso del tipo siguiendo a la lista de identificadores (Pascal):

$$\langle \text{DeclVariables} \rangle \rightarrow \text{id} (, \text{id})^* : \langle \text{Tipo} \rangle$$

4.4.3. Expresiones aritméticas

Podemos simplificar las producciones correspondientes a las expresiones aritméticas. Sin embargo, esta simplificación tiene un precio; no podremos expresar directamente en la gramática la asociatividad de los operadores. Tendremos que relegar en fases posteriores la responsabilidad de definir la asociatividad. En la práctica esto no supondrá ningún problema. Sin embargo, seguiremos teniendo un no terminal por cada nivel de prioridad. La gramática que teníamos para las expresiones la podemos reescribir así:

$$\begin{array}{l}
 \langle E \rangle \rightarrow \langle T \rangle ((+ | -) \langle T \rangle)^* \\
 \langle T \rangle \rightarrow \langle F \rangle ((* | /) \langle F \rangle)^* \\
 \langle F \rangle \rightarrow \langle B \rangle (\uparrow \langle B \rangle)^* \\
 \langle B \rangle \rightarrow \text{id} | \text{"("} \langle E \rangle \text{"}
 \end{array}$$

Si agrupamos + junto con - en la categoría léxica **opad** (operador aditivo) y * junto con / en **opmult** (operador multiplicativo), nos queda una gramática ligeramente más clara:

$$\begin{aligned}\langle E \rangle &\rightarrow \langle T \rangle (\text{opad } \langle T \rangle)^* \\ \langle T \rangle &\rightarrow \langle F \rangle (\text{opmult } \langle F \rangle)^* \\ \langle F \rangle &\rightarrow \langle B \rangle (\uparrow \langle B \rangle)^* \\ \langle B \rangle &\rightarrow \text{id} | \langle \text{E} \rangle\end{aligned}$$

EJERCICIO 18

Halla el árbol de análisis de **id*id↑(id+id)**.

EJERCICIO 19

Aumenta la gramática para permitir operadores unarios. Dales la máxima prioridad.

EJERCICIO 20

Reescribe las construcciones del ejercicio 14 utilizando GPDR.

5. Análisis descendente

Una vez definido nuestro lenguaje mediante una gramática G , estamos interesados en resolver el siguiente problema:

Dada una cadena x y una gramática G , ¿pertenece x a $L(G)$?

Esto es lo que se conoce como el *problema del análisis*.

Existen diversos algoritmos que resuelven este problema de manera general. Dos de los más conocidos son:

Algoritmo de Cocke-Younger-Kasami o CYK. Este algoritmo utiliza el esquema algorítmico de Programación Dinámica. Su coste es $O(|x|^3)$ y necesita una gramática en forma normal de Chomsky.

Algoritmo de Early este algoritmo tiene un coste igual al CYK si la gramática es ambigua. Si no lo es, el coste se reduce a $O(|x|^2)$ aunque, para muchas gramáticas, el coste es $O(|x|)$.

Nuestro objetivo será buscar algoritmos que garanticen un coste lineal con la entrada. Para ello, tendremos que poner algunas restricciones a las gramáticas que se podrán utilizar.

Dependiendo del modo en que se construye el árbol de análisis, distinguimos dos familias de algoritmos de análisis:

Análisis descendente se construye el árbol partiendo de la raíz hasta llegar a las hojas. En este proceso se va “prediciendo” lo que vendrá después, por lo que también se le denomina *análisis predictivo*.

Análisis ascendente se construye el árbol desde las hojas. La detección de las construcciones se realiza una vez que se han visto completamente.

En este tema estudiaremos un tipo de análisis descendente: el análisis LL(1), que es fácil de implementar manualmente y es razonablemente potente; a final de curso, estudiaremos métodos de análisis ascendente. Comenzaremos por ver cómo analizar las gramáticas incontextuales y después pasaremos a las GPDR.

5.1. Un ejemplo

Supongamos que tenemos la siguiente gramática:

$$\begin{aligned} \langle S \rangle &\rightarrow \text{if } \langle E \rangle \text{ then } \langle S \rangle \\ \langle S \rangle &\rightarrow \text{while } \langle E \rangle \text{ do } \langle S \rangle \\ \langle S \rangle &\rightarrow \text{id} := \langle E \rangle \\ \langle E \rangle &\rightarrow \text{id} | \text{cte} \end{aligned}$$

Queremos analizar la sentencia **if cte then id:=cte**. Inicialmente sabemos que tendremos un árbol que en la raíz tendrá el no terminal $\langle S \rangle$ que derivará el resto. Además, sabemos que lo que derive $\langle S \rangle$ estará seguido por el fin de la entrada, que representaremos mediante **\$**. En este momento, estaremos mirando el primer símbolo de la entrada, el componente **if**.

Podemos resumir esta información de la siguiente manera:

$$\begin{array}{c|c} \text{if} & \text{cte then id:=cte\$} \\ \hline \langle S \rangle & \$ \end{array} \quad \langle S \rangle$$

En la parte derecha tenemos el árbol que estamos construyendo (de momento sólo conocemos la raíz). El cuadro de la izquierda tiene en la parte superior la entrada que tenemos que analizar y en la inferior la predicción del analizador. Hemos separado el primer símbolo de cada parte para que quede claro qué se está comparando en cada momento.

Mirando la gramática, vemos que la única regla que podemos aplicar es la primera. La nueva situación es:

$$\begin{array}{c|c} \text{if} & \text{cte then id:=cte\$} \\ \hline \text{if} & \langle E \rangle \text{ then } \langle S \rangle \$ \end{array} \quad \begin{array}{c} \langle S \rangle \\ / \quad | \quad \backslash \\ \text{if} \quad \langle E \rangle \quad \text{then} \quad \langle S \rangle \end{array}$$

Este cuadro lo podemos interpretar de la siguiente manera: el analizador predice que va a ver un **if** seguido de algo que expanda $\langle E \rangle$, un **then**, después algo que expanda $\langle S \rangle$ y finalmente el final de la entrada. Dado que estamos viendo el **if** que se ha predicho, podemos avanzar en la entrada:

$$\begin{array}{c|c} \text{cte} & \text{then id:=cte\$} \\ \hline \langle E \rangle & \text{then } \langle S \rangle \$ \end{array}$$

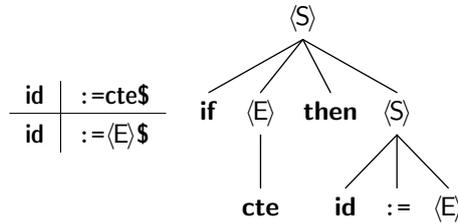
Nuevamente, analizando la gramática, vemos que sólo podemos aplicar una regla: $\langle E \rangle \rightarrow \text{cte}$. Cambiamos nuestra predicción, y aumentamos ligeramente el árbol:

$$\begin{array}{c|c} \text{cte} & \text{then id:=cte\$} \\ \hline \text{cte} & \text{then } \langle S \rangle \$ \end{array} \quad \begin{array}{c} \langle S \rangle \\ / \quad | \quad \backslash \\ \text{if} \quad \langle E \rangle \quad \text{then} \quad \langle S \rangle \\ | \\ \text{cte} \end{array}$$

Como coinciden los terminales, avanzamos. Vuelven a coincidir y avanzamos nuevamente. Estos dos pasos son:

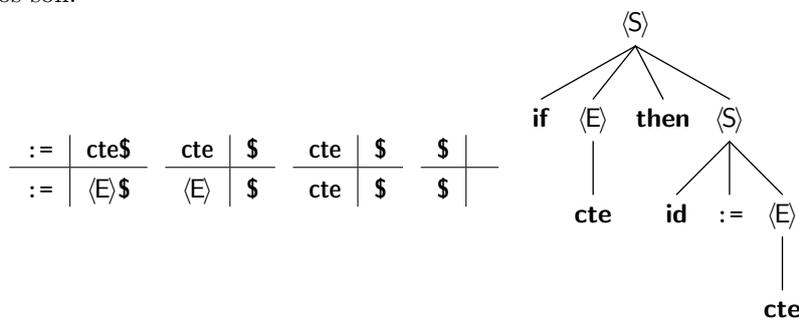
$$\begin{array}{c|c} \text{then} & \text{id:=cte\$} \\ \hline \text{then} & \langle S \rangle \$ \end{array} \quad \begin{array}{c|c} \text{id} & :=\text{cte\$} \\ \hline \langle S \rangle & \$ \end{array}$$

Ahora tenemos que encontrar una expansión de $\langle S \rangle$ que comience por **id**, lo que supone aplicar la regla $\langle S \rangle \rightarrow \text{id} := \langle E \rangle$. Aumentamos el árbol y llegamos a:



Vuelve a haber coincidencias de terminales, hasta que finalmente llegamos a tener que expandir $\langle E \rangle$ por **cte**.

Estos pasos son:



Hemos llegado a una situación en la que vemos que hemos predicho el fin de la entrada y lo hemos encontrado. Esto quiere decir que el análisis ha concluido con éxito y ya tenemos el árbol de análisis.

5.2. El algoritmo de análisis

Vamos a intentar transformar el ejemplo anterior en un algoritmo. En primer lugar, observemos que la predicción que hacemos (la fila de abajo en los cuadros) se comporta como una pila. Además podemos darnos cuenta de que hacemos dos tipos de acción según sea el símbolo en la entrada y la predicción:

- Si hemos predicho un no terminal, vemos qué regla de ese no terminal comienza por el símbolo en la entrada y utilizamos la parte derecha de esa regla para completar nuestra predicción.
- Si ambos son terminales, comprobamos si son iguales. Si lo son, avanzamos.

Observando el primer paso, nos damos cuenta de que la información acerca de qué regla aplicar en cada caso, la podemos guardar en una tabla. Esta tabla tendrá, para nuestro ejemplo, el aspecto siguiente:

	cte	do	id	if	then	while	:=	\$
$\langle S \rangle$			id := ⟨E⟩	if ⟨E⟩ then ⟨S⟩		while ⟨E⟩ do ⟨S⟩		
$\langle E \rangle$	cte		id					

Las entradas vacías corresponden a errores sintácticos.

Si llamamos M a nuestra tabla, el algoritmo de análisis es el siguiente:

```

Algoritmo Análisis predictivo
  push( $\Pi$ ,  $\$$ ); push( $\Pi$ ,  $\langle S \rangle$ ); // Predicción inicial:  $\langle S \rangle \$$ 
   $a := siguiente()$ ;
  mientras top( $\Pi$ )  $\neq$   $\$$  hacer
     $X := pop(\Pi)$ ;
    si  $X \in N$ : //  $X$  es un no terminal
      si  $M[X, a] \neq \text{ERROR}$  entonces
        push( $\Pi$ ,  $M[X, a]^{-1}$ )
      si no
        devolver falso
      fin si
    si no //  $X$  es un terminal
      si  $a = X$  entonces
         $a := siguiente()$ 
      si no
        devolver falso
      fin si
    fin si
  fin mientras
  devolver  $a = \$$ ;
fin Análisis predictivo
  
```

Este algoritmo devolverá cierto si la cadena pertenece al lenguaje definido por la gramática y falso en otro caso. Más adelante veremos que nos interesará será tratar los errores adecuadamente, pero de momento, esto nos vale. Utilizamos la pila Π para guardar la predicción. Fíjate que cuando introducimos la nueva parte derecha en la pila, lo hacemos “al revés”, para que sea analizada en el orden correcto.

Si además estamos interesados en obtener el árbol, bastará con tomar nota cada vez que se haga una predicción de la regla utilizada. Como puedes comprobar, el coste del algoritmo es lineal con la longitud de la entrada. Lo que vamos a hacer ahora es buscar la manera de construir las tablas de manera automática a partir de la gramática.

EJERCICIO 21

Utiliza el algoritmo para analizar las cadenas:

- **while cte do id :=id**
- **while cte if id :=id**

5.3. Cálculo de primeros

En el ejemplo anterior hemos podido decidir en cada caso qué producción aplicar porque los primeros símbolos de cada alternativa eran terminales diferentes. Sin embargo, exigir esta condición para que podamos analizar cadenas con una gramática es excesivo. Vamos a ver cómo podemos relajarla sin perder la capacidad de escribir la tabla de análisis.

Supongamos que modificamos la gramática anterior para permitir que se emplee el operador unario $*$ de modo que sea posible utilizar punteros. La nueva gramática tiene este aspecto:

$$\begin{aligned}
 \langle S \rangle &\rightarrow \text{if } \langle E \rangle \text{ then } \langle S \rangle | \text{while } \langle E \rangle \text{ do } \langle S \rangle | \langle l \rangle := \langle E \rangle \\
 \langle E \rangle &\rightarrow \langle l \rangle | \text{cte} \\
 \langle l \rangle &\rightarrow * \langle l \rangle | \text{id}
 \end{aligned}$$

Ahora tenemos un problema para rellenar en la tabla las filas de $\langle S \rangle$ y $\langle E \rangle$: tenemos que saber cuándo utilizar las reglas que comienzan por $\langle l \rangle$. Sin embargo, un análisis rápido nos permite ver

que las cadenas producidas por $\langle l \rangle$ sólo pueden comenzar por $*$ o por **id**. Así pues, la tabla de análisis es la siguiente:

	cte	do	id	if	then	while	*	:=	\$
$\langle S \rangle$			$\langle l \rangle := \langle E \rangle$	if $\langle E \rangle$ then $\langle S \rangle$		while $\langle E \rangle$ do $\langle S \rangle$		$\langle l \rangle := \langle E \rangle$	
$\langle E \rangle$	cte		$\langle l \rangle$					$\langle l \rangle$	
$\langle l \rangle$			id				$*\langle l \rangle$		

Vamos a formalizar el proceso de construcción de la tabla. En primer lugar, para decidir entre las alternativas de una producción, debemos saber cómo empiezan las cadenas que puede generar. Si tenemos una forma sentencial α , podemos definir $\text{primeros}(\alpha)$ como:

$$\text{primeros}(\alpha) = \{a \in \Sigma \mid \alpha \xRightarrow{*} a\beta\}.$$

Es decir, $\text{primeros}(\alpha)$ es el conjunto de terminales que pueden ser el primer símbolo de las cadenas generadas por α^1 .

Para escribir el algoritmo de cálculo de los primeros de una forma sentencial, necesitamos el concepto de *anulable*. Decimos que una forma sentencial es anulable (lo representamos como $\text{anulable}(\alpha)$) si de ella se puede derivar la cadena vacía.

Si la gramática no tiene recursividad por la izquierda, podemos calcular el conjunto de primeros de α mediante el siguiente algoritmo recursivo:

```

Algoritmo primeros ( $\alpha : (N \cup \Sigma)^*$ )
   $C = \emptyset$ 
  si  $\alpha = a\beta$  entonces
     $C := \{a\}$ ;
  si no si  $\alpha = \langle A \rangle\beta$  entonces
    para todo  $\langle A \rangle \rightarrow \gamma \in P$  hacer
       $C := C \cup \text{primeros}(\gamma)$ ;
    fin para
    si  $\text{anulable}(\langle A \rangle)$  entonces
       $C := C \cup \text{primeros}(\beta)$ ;
    fin si
  fin si
  devolver  $C$ ;
fin primeros
    
```

El funcionamiento del algoritmo es sencillo:

- Si α comienza por un terminal, éste es el único primero.
- Si α comienza por el no terminal $\langle A \rangle$, tendremos que obtener los primeros de $\langle A \rangle$. Además, si $\langle A \rangle$ es anulable también habrá que mirar los primeros de lo que le sigue.

¹En algunos textos, cuando la cadena vacía puede ser derivada de α , se la considera parte de los primeros de α . Esto supone una leve modificación de la formalización.

Podemos averiguar si α es anulable mediante el siguiente algoritmo:

```

Algoritmo anulable ( $\alpha : (N \cup \Sigma)^*$ )
  si  $\alpha = \lambda$  entonces
    devolver cierto;
  si no si  $\alpha = a\beta$  entonces
    devolver falso;
  si no si  $\alpha = \langle A \rangle \beta$  entonces
    si existe  $\gamma$  tal que  $\langle A \rangle \rightarrow \gamma \in P$  y anulable( $\gamma$ ) entonces
      devolver anulable( $\beta$ );
    si no
      devolver falso;
    fin si
  fin si
fin anulable
    
```

EJERCICIO 22

Calcula los primeros de los no terminales de la siguiente gramática:

$$\begin{aligned}
 \langle S \rangle &\rightarrow \langle A \rangle \langle B \rangle | \langle B \rangle c \\
 \langle A \rangle &\rightarrow a b | c d \\
 \langle B \rangle &\rightarrow b a | \lambda
 \end{aligned}$$

Estos algoritmos sirven si la gramática no tiene recursividad a izquierdas. En caso contrario, los cálculos de primeros y anulables se complican ligeramente (mira el ejercicio 24). Sin embargo, dado que las gramáticas con recursividad a izquierdas no son LL(1), no estudiaremos ese caso aquí y veremos más adelante cómo eliminar la recursividad a izquierdas.

Ahora podemos construir la tabla de análisis de la siguiente manera:

- En primer lugar calculamos los primeros de las partes derechas de las reglas de la gramática.
- Después rellenamos la celda $M[\langle A \rangle, a]$ de la tabla con α , si $\langle A \rangle \rightarrow \alpha \in P$ y $a \in \text{primeros}(\alpha)$.

Este método de construcción de la tabla funcionará siempre y cuando ningún no terminal sea anulable.

5.4. Cálculo de siguientes

¿Qué sucede cuando algún no terminal es anulable? Vamos a estudiarlo aumentando la gramática de nuestro ejemplo para permitir accesos a vectores:

$$\begin{aligned}
 \langle S \rangle &\rightarrow \text{if } \langle E \rangle \text{ then } \langle S \rangle | \text{while } \langle E \rangle \text{ do } \langle S \rangle | \langle l \rangle := \langle E \rangle \\
 \langle E \rangle &\rightarrow \langle l \rangle | \text{cte} \\
 \langle l \rangle &\rightarrow * \langle l \rangle | \text{id } \langle A \rangle \\
 \langle A \rangle &\rightarrow \lambda | [\langle E \rangle]
 \end{aligned}$$

Intentemos analizar la cadena **id:=id[cte]**. La situación inicial es:

id	:=id[cte]\$
(S)	\$

Ahora podemos aplicar la regla $\langle S \rangle \rightarrow \langle l \rangle := \langle E \rangle$ y después $\langle l \rangle \rightarrow \text{id } \langle A \rangle$:

id	:=id[cte]\$	id	:=id[cte]\$
(l)	:=(E)\$	id	(A):=(E)\$

Como coinciden los terminales, avanzamos:

$$\frac{:= \mid \text{id}[\text{cte}]\$}{\langle A \rangle \mid :=\langle E \rangle\$}$$

Y ahora tenemos un problema. Evidentemente, $:=$ no está entre los primeros de $\langle A \rangle$. Por tanto, nos preguntamos si en esta situación podemos utilizar la regla $\langle A \rangle \rightarrow \lambda$. Para encontrar la respuesta, debemos plantearnos cuándo esta regla es válida. La podremos emplear cuando nos encontremos en la entrada algo que “pueda ir detrás” de $\langle A \rangle$. Si examinamos la predicción, vemos que, en este caso, la asignación efectivamente puede seguir a $\langle A \rangle$, por lo que concluimos que es correcto sustituir $\langle A \rangle$ por la cadena vacía y continuar el análisis. Así, los pasos que siguen son

$$\begin{array}{c} := \mid \text{id}[\text{cte}]\$ \quad \text{id} \mid [\text{cte}]\$ \quad \text{id} \mid [\text{cte}]\$ \quad \text{id} \mid [\text{cte}]\$ \quad [\mid \text{cte}]\$ \\ := \mid \langle E \rangle\$ \quad \langle E \rangle \mid \$ \quad \langle \rangle \mid \$ \quad \text{id} \mid \langle A \rangle\$ \quad \langle A \rangle \mid \$ \end{array}$$

Ahora nos encontramos analizando uno de los primeros de $\langle A \rangle$; seguimos el análisis como hasta ahora:

$$\begin{array}{c} [\mid \text{cte}]\$ \quad \text{cte} \mid]\$ \quad \text{cte} \mid]\$ \quad] \mid \$ \quad \$ \mid \$ \\ [\mid \langle E \rangle\$ \quad \langle E \rangle \mid]\$ \quad \text{cte} \mid]\$ \quad] \mid \$ \quad \$ \mid \$ \end{array}$$

Y concluimos el análisis con éxito.

Hemos visto que podemos aplicar una λ -producción siempre que nos encontremos en la entrada alguno de los símbolos que puedan aparecer tras el no terminal correspondiente. Vamos a ver cómo calcular estos símbolos. Llamaremos *siguientes* de un no terminal $\langle A \rangle$ a los miembros del conjunto:

$$\text{siguientes}(\langle A \rangle) = \begin{cases} \{a \in \Sigma \mid \langle S \rangle \xrightarrow{*} \alpha \langle A \rangle \beta, a \in \text{primeros}(\beta)\} \cup \{\$\} & \text{si } \langle S \rangle \xrightarrow{*} \alpha \langle A \rangle, \\ \{a \in \Sigma \mid \langle S \rangle \xrightarrow{*} \alpha \langle A \rangle \beta, a \in \text{primeros}(\beta)\} & \text{si no.} \end{cases}$$

La idea es que los siguientes de un no terminal $\langle A \rangle$ son aquellos terminales que pueden llegar a aparecer detrás de $\langle A \rangle$ en alguna forma sentencial. En el caso del símbolo inicial y de aquellos que pueden aparecer al final de una forma sentencial, además incluimos el fin de la entrada.

Para encontrar los siguientes de cada no terminal, podemos emplear las siguientes ecuaciones:

- En primer lugar:

$$\{\$\} \subseteq \text{siguientes}(\langle S \rangle). \tag{1}$$

- Si $\langle A \rangle \rightarrow \alpha \langle B \rangle \beta \in P$, entonces:

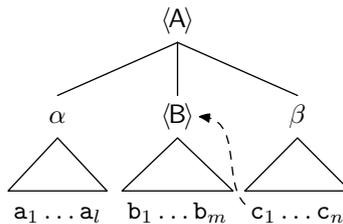
$$\text{primeros}(\beta) \subseteq \text{siguientes}(\langle B \rangle); \tag{2}$$

- si, además, anulable(β) entonces:

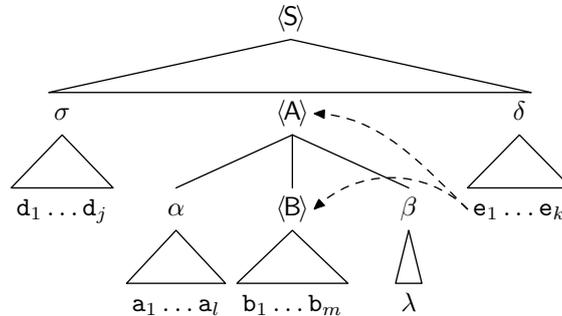
$$\text{siguientes}(\langle A \rangle) \subseteq \text{siguientes}(\langle B \rangle). \tag{3}$$

La ecuación (1) es fácil de entender: cualquier cadena generada a partir del símbolo inicial estará seguida del fin de la entrada.

Para entender (2), nos planteamos cómo se genera la cadena $a_1 \dots a_l b_1 \dots b_m c_1 \dots c_n$ a partir de $\langle A \rangle$:



Es fácil ver qué “va detrás” (los siguientes) de $\langle B \rangle$: los primeros de β . Por otro lado, si β puede anularse, tenemos una situación como la siguiente:



Con lo que $\langle B \rangle$ y $\langle A \rangle$ comparten siguientes, lo que explica la ecuación (3). Para resolver las ecuaciones podemos ir aplicándolas en tres pasos:

1. Como inicialización, asociamos a todos los no terminales distintos del inicial el conjunto vacío y al símbolo inicial le asociamos el conjunto $\{\$ \}$ para que se cumpla la ecuación (1).
2. Para cumplir con la ecuación (2), añadimos al conjunto de siguientes de cada no terminal $\langle A \rangle$ los primeros de los β tales que $\alpha \langle A \rangle \beta$ constituye la parte derecha de alguna regla.
3. Recorremos ordenadamente las reglas aplicando la ecuación (3). Para ello, por cada regla $\langle A \rangle \rightarrow \alpha \langle B \rangle \beta$ con β anulable, añadimos los siguientes de $\langle A \rangle$ a los de $\langle B \rangle$. Este paso tendremos que repetirlo hasta que no cambie ninguno de los conjuntos de siguientes en un recorrido completo de las reglas.

Vamos a aplicar el algoritmo a nuestro ejemplo. En primer lugar, inicializamos los siguientes al conjunto vacío salvo para el símbolo inicial:

No terminal	Siguientes
$\langle S \rangle$	$\{\$ \}$
$\langle E \rangle$	\emptyset
$\langle I \rangle$	\emptyset
$\langle A \rangle$	\emptyset

Después, añadimos los símbolos que aparecen inmediatamente después de cada no terminal:

No terminal	Siguientes
$\langle S \rangle$	$\{\$ \}$
$\langle E \rangle$	$\{], \text{do}, \text{then}\}$
$\langle I \rangle$	$\{:=\}$
$\langle A \rangle$	\emptyset

Ahora aplicamos ordenadamente la ecuación (3) hasta que no haya modificaciones. Así, a partir de $\langle S \rangle \rightarrow \langle I \rangle := \langle E \rangle$, aumentamos los siguientes de $\langle E \rangle$:

No terminal	Siguientes
$\langle S \rangle$	$\{\$ \}$
$\langle E \rangle$	$\{], \$, \text{do}, \text{then}\}$
$\langle I \rangle$	$\{:=\}$
$\langle A \rangle$	\emptyset

Con la regla $\langle E \rangle \rightarrow \langle l \rangle$, aumentamos los siguientes de $\langle l \rangle$:

No terminal	Siguientes
$\langle S \rangle$	$\{\$ \}$
$\langle E \rangle$	$\{], \$, \text{do}, \text{then}\}$
$\langle l \rangle$	$\{:=,], \$, \text{do}, \text{then}\}$
$\langle A \rangle$	\emptyset

A partir de $\langle l \rangle \rightarrow \text{id}\langle A \rangle$, aumentamos los siguientes de $\langle A \rangle$:

No terminal	Siguientes
$\langle S \rangle$	$\{\$ \}$
$\langle E \rangle$	$\{], \$, \text{do}, \text{then}\}$
$\langle l \rangle$	$\{:=,], \$, \text{do}, \text{then}\}$
$\langle A \rangle$	$\{:=,], \$, \text{do}, \text{then}\}$

Si volvemos a repasar las reglas, vemos que los siguientes no varían, así que podemos dar por terminado el proceso.

En algunos casos, hay que repetir el tercer paso más veces. Por ejemplo, vamos a calcular los siguientes para la gramática:

$$\begin{aligned} \langle A \rangle &\rightarrow \langle B \rangle c \mid \langle C \rangle d \\ \langle B \rangle &\rightarrow b \mid a \langle A \rangle \\ \langle C \rangle &\rightarrow a \langle B \rangle \end{aligned}$$

Inicialmente, tenemos vacíos los conjuntos de siguientes, salvo el de $\langle A \rangle$:

No terminal	Siguientes
$\langle A \rangle$	$\{\$ \}$
$\langle B \rangle$	\emptyset
$\langle C \rangle$	\emptyset

Ahora, por la ecuación (2), añadimos c a los siguientes de $\langle B \rangle$ y d a los de $\langle C \rangle$:

No terminal	Siguientes
$\langle A \rangle$	$\{\$ \}$
$\langle B \rangle$	$\{c\}$
$\langle C \rangle$	$\{d\}$

La ecuación (3) sólo se puede aplicar a $\langle B \rangle \rightarrow a \langle A \rangle$ y $\langle C \rangle \rightarrow a \langle B \rangle$. Así que, copiamos los siguientes de $\langle B \rangle$ en los de $\langle A \rangle$:

No terminal	Siguientes
$\langle A \rangle$	$\{c, \$ \}$
$\langle B \rangle$	$\{c\}$
$\langle C \rangle$	$\{d\}$

Y los de $\langle C \rangle$ en los de $\langle B \rangle$:

No terminal	Siguientes
$\langle A \rangle$	{c, \$}
$\langle B \rangle$	{c, d}
$\langle C \rangle$	{d}

Y todavía no hemos terminado. Como ha habido modificaciones en el último recorrido, tenemos que dar una nueva pasada. Copiamos los siguientes de $\langle B \rangle$ en los de $\langle A \rangle$:

No terminal	Siguientes
$\langle A \rangle$	{c, d, \$}
$\langle B \rangle$	{c, d}
$\langle C \rangle$	{d}

Al copiar los siguientes de $\langle C \rangle$ en los de $\langle B \rangle$ no hay cambios. Haciendo ahora un último recorrido, vemos que no hay más cambios y terminamos.

EJERCICIO 23

Calcula los primeros y siguientes de los no terminales de la siguiente gramática:

$$\begin{aligned} \langle S \rangle &\rightarrow \langle \text{Id} \rangle := \langle E \rangle ; \\ \langle S \rangle &\rightarrow \text{while } \langle E \rangle \text{ do } \langle \text{LS} \rangle \text{ endwhile} \\ \langle S \rangle &\rightarrow \langle \text{Cond} \rangle \text{ fi} \\ \langle S \rangle &\rightarrow \langle \text{Cond} \rangle \text{ else } \langle \text{LS} \rangle \text{ fi} \\ \langle \text{LS} \rangle &\rightarrow \langle S \rangle \langle \text{LS} \rangle | \lambda \\ \langle \text{Cond} \rangle &\rightarrow \text{if } \langle E \rangle \text{ then } \langle \text{LS} \rangle \\ \langle E \rangle &\rightarrow \text{cte} | \langle \text{Id} \rangle \\ \langle \text{Id} \rangle &\rightarrow * \langle \text{Id} \rangle | \text{id} \langle A \rangle \\ \langle A \rangle &\rightarrow [\langle E \rangle] | \lambda \end{aligned}$$

EJERCICIO* 24

Si la gramática tiene ciclos, el cálculo de primeros y de anulables se hace mediante ecuaciones similares a las que hemos utilizado para el cálculo de siguientes. Escribe esas ecuaciones y pruébalas con la siguiente gramática:

$$\begin{aligned} \langle S \rangle &\rightarrow \langle E \rangle ; \langle S \rangle | \lambda \\ \langle E \rangle &\rightarrow \langle E \rangle + \langle T \rangle | \langle T \rangle \\ \langle T \rangle &\rightarrow \langle T \rangle \langle M \rangle \langle F \rangle | \langle F \rangle \\ \langle M \rangle &\rightarrow * | \lambda \\ \langle F \rangle &\rightarrow \text{id} | \text{"("} \langle E \rangle \text{"} \end{aligned}$$

5.5. Análisis LL(1)

Ahora estamos en condiciones de rellenar la tabla de análisis. Los pasos que seguiremos son:

- Averiguar qué no terminales son anulables.
- Calcular los primeros de las partes derechas de las reglas.
- Calcular los siguientes de los no terminales anulables.
- Por cada regla $\langle A \rangle \rightarrow \alpha$, poner α en las celdas $M[\langle A \rangle, a]$ tales que:
 - $a \in \text{primeros}(\alpha)$
 - $\text{o anulable}(\alpha)$ y $a \in \text{siguientes}(\langle A \rangle)$.

Si en cada celda hay como máximo una parte derecha, decimos que la gramática es LL(1) y podemos aplicar nuestro algoritmo de análisis. El nombre LL(1) viene de:

- **Left to right**: el orden en que se lee la entrada.
- **Leftmost derivation**: se encuentra una derivación por la izquierda.
- **1 lookahead symbol**: se utiliza un símbolo de anticipación.

La tabla de nuestro ejemplo es:

	cte	do	id	if	then	while	*	:=	[]	\$
$\langle S \rangle$			$\langle l \rangle := \langle E \rangle$	$\text{if } \langle E \rangle \text{ then } \langle S \rangle$		$\text{while } \langle E \rangle \text{ do } \langle S \rangle$	$\langle l \rangle := \langle E \rangle$				
$\langle E \rangle$	cte		$\langle l \rangle$				$\langle l \rangle$				
$\langle l \rangle$			id $\langle A \rangle$				$* \langle l \rangle$				
$\langle A \rangle$		λ			λ				λ	$[\langle E \rangle]$	$\lambda \lambda$

Así pues, la gramática era LL(1).

Algunas propiedades de las gramáticas LL(1):

- No son ambiguas.
- No tienen recursividad por la izquierda.
- Los conjuntos de primeros de las partes derechas de un no terminal son disjuntos.
- Si $\langle A \rangle \xrightarrow{*} \lambda$, los siguientes de $\langle A \rangle$ son disjuntos con los primeros de la parte derecha de cualquier $\langle A \rangle$ -producción. Además, $\langle A \rangle$ tiene como mucho una parte derecha anulable.

Hasta ahora hemos hablado de gramáticas LL(1). Si un lenguaje tiene una gramática LL(1), decimos que ese lenguaje es LL(1). Existen lenguajes incontextuales que no son LL(1), es decir que no existe ninguna gramática LL(1) capaz de modelarlos. Un ejemplo sería:

$$\{a^n 0 b^n \mid n \geq 1\} \cup \{a^n 1 b^{2n} \mid n \geq 1\}.$$

Afortunadamente, gran parte de las construcciones habituales de los lenguajes de programación pueden modelarse mediante gramáticas LL(1).

5.5.1. Conflictos LL(1)

Si al rellenar la tabla de análisis aparece una celda con más de una entrada, tenemos lo que se llama un *conflicto LL(1)*. Esta situación indica que el analizador no sabe qué debe predecir en ese momento. Normalmente, basta modificar ligeramente la gramática para evitar estos conflictos. Veremos ahora algunas de las posibles modificaciones.

Eliminación de la recursividad por la izquierda Si la gramática tiene recursividad por la izquierda es seguro que aparecerá un conflicto (¿por qué?). Hay dos tipos de recursividad por la izquierda:

- Directa: existe al menos una *regla* de la forma $\langle A \rangle \rightarrow \langle A \rangle \alpha$.
- Indirecta: existe una *derivación* de la forma $\langle A \rangle \xRightarrow{\pm} \langle A \rangle \alpha$.

En primer lugar, veremos cómo podemos eliminar la recursividad directa. Supongamos que tenemos el siguiente par de producciones:

$$\begin{aligned} \langle A \rangle &\rightarrow \langle A \rangle \alpha \\ \langle A \rangle &\rightarrow \beta \end{aligned}$$

Las formas sentenciales que pueden generar son del tipo $\beta \alpha^n$ con $n \geq 0$. Podemos sustituirlas por las siguientes reglas, que generan las mismas formas sentenciales:

$$\begin{aligned} \langle A \rangle &\rightarrow \beta \langle A' \rangle \\ \langle A' \rangle &\rightarrow \alpha \langle A' \rangle \\ \langle A' \rangle &\rightarrow \lambda \end{aligned}$$

En algunos casos, la eliminación no es tan sencilla y se requiere expandir primero alguna de las reglas para después aplicar la transformación anterior. Existe un algoritmo para eliminar completamente la recursividad por la izquierda. Por desgracia, la gramática resultante no tiene ninguna relación con la original, lo que lo hace poco práctico. Afortunadamente, en la mayor parte de las ocasiones, la recursividad por la izquierda se puede eliminar utilizando gramáticas con partes derechas regulares.

EJERCICIO 25

Elimina la recursividad por la izquierda en la siguiente gramática:

$$\begin{aligned} \langle E \rangle &\rightarrow \langle E \rangle + \langle T \rangle \mid \langle E \rangle - \langle T \rangle \mid \langle T \rangle \\ \langle T \rangle &\rightarrow \langle T \rangle * \langle F \rangle \mid \langle T \rangle / \langle F \rangle \mid \langle F \rangle \\ \langle F \rangle &\rightarrow "(" \langle E \rangle ")" \mid \text{id} \end{aligned}$$

Calcula su tabla de análisis y analiza la cadena **id*(id+id)**.

EJERCICIO* 26

¿Cómo podemos eliminar la recursividad directa por la izquierda en las siguientes reglas?

$$\begin{aligned} \langle A \rangle &\rightarrow \langle A \rangle \alpha_1 \mid \langle A \rangle \alpha_2 \mid \dots \mid \langle A \rangle \alpha_n \\ \langle A \rangle &\rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_m \end{aligned}$$

Eliminación de prefijos comunes Otra de las causas de aparición de conflictos LL(1) es la existencia de prefijos comunes entre distintas partes derechas, lo que lleva a que los distintos conjuntos de primeros no sean disjuntos. Para eliminar estos conflictos, podemos aplicar el análogo de “sacar factor común”. Si tenemos las reglas:

$$\langle A \rangle \rightarrow \alpha \beta \mid \alpha \gamma$$

Podemos sustituirlas por:

$$\begin{aligned} \langle A \rangle &\rightarrow \alpha \langle A' \rangle \\ \langle A' \rangle &\rightarrow \beta \mid \gamma \end{aligned}$$

Puede suceder que los prefijos aparezcan “indirectamente”, entonces hay que expandir previamente algún no terminal. Por ejemplo:

$$\begin{aligned} \langle S \rangle &\rightarrow \text{id} := \langle E \rangle | \langle P \rangle \\ \langle P \rangle &\rightarrow \text{id} \text{ “(”} \langle E \rangle \text{ “)”} \end{aligned}$$

En este caso, deberemos expandir en primer lugar el no terminal $\langle P \rangle$:

$$\begin{aligned} \langle S \rangle &\rightarrow \text{id} := \langle E \rangle | \text{id “(”} \langle E \rangle \text{ “)”} \\ \langle P \rangle &\rightarrow \text{id “(”} \langle E \rangle \text{ “)”} \end{aligned}$$

Después aplicamos la factorización:

$$\begin{aligned} \langle S \rangle &\rightarrow \text{id } \langle S' \rangle \\ \langle S' \rangle &\rightarrow := \langle E \rangle | \text{ “(”} \langle E \rangle \text{ “)”} \\ \langle P \rangle &\rightarrow \text{id “(”} \langle E \rangle \text{ “)”} \end{aligned}$$

Finalmente, si $\langle P \rangle$ no aparece en ninguna parte derecha, podemos eliminarlo.

EJERCICIO 27

Elimina los prefijos comunes de estas producciones:

$$\begin{aligned} \langle S \rangle &\rightarrow \text{try} : \langle S \rangle \text{ except id} : \langle S \rangle \\ \langle S \rangle &\rightarrow \text{try} : \langle S \rangle \text{ except id} : \langle S \rangle \text{ finally} : \langle S \rangle \\ \langle S \rangle &\rightarrow \text{id} := \langle E \rangle ; | \langle E \rangle ; \\ \langle S \rangle &\rightarrow \{ \langle LS \rangle \} \\ \langle LS \rangle &\rightarrow \langle S \rangle | \langle S \rangle \langle LS \rangle | \lambda \\ \langle E \rangle &\rightarrow \text{id} | \text{cte} \end{aligned}$$

Manipulación directa de la tabla de análisis En algunos casos, la mejor manera de resolver un conflicto es cambiar directamente la tabla de análisis. Por ejemplo, sea la gramática:

$$\begin{aligned} \langle S \rangle &\rightarrow \text{if } \langle E \rangle \text{ then } \langle S \rangle \text{ else } \langle S \rangle \\ \langle S \rangle &\rightarrow \text{if } \langle E \rangle \text{ then } \langle S \rangle \\ \langle S \rangle &\rightarrow \text{sent} \\ \langle E \rangle &\rightarrow \text{id} \end{aligned}$$

En primer lugar, eliminamos prefijos comunes:

$$\begin{aligned} \langle S \rangle &\rightarrow \text{if } \langle E \rangle \text{ then } \langle S \rangle \langle E \rangle \\ \langle S \rangle &\rightarrow \text{sent} \\ \langle E \rangle &\rightarrow \text{id} \\ \langle E \rangle &\rightarrow \text{else } \langle S \rangle | \lambda \end{aligned}$$

Ahora escribimos la tabla de análisis:

	else	id	if	sent	then	\$
$\langle S \rangle$			if $\langle E \rangle$ then $\langle S \rangle$ $\langle E \rangle$	sent		
$\langle E \rangle$		id				
$\langle E \rangle$	else $\langle S \rangle$ λ					λ

Vemos que hay un conflicto en la celda correspondiente a $\langle E \rangle$ con el terminal **else**. Si analizamos en qué situaciones podemos llegar al conflicto, vemos que la única manera en la que podemos tener un $\langle E \rangle$ en el tope de la pila es por haber terminado de analizar la $\langle S \rangle$ de un **if-then**. Pero si seguimos la regla de asociar el **else** al **if** más cercano, lo que tenemos que hacer en esta situación es predecir la secuencia **else** $\langle S \rangle$. Así pues, modificamos la entrada de la tabla y obtenemos:

	else	id	if	sent	then	\$
$\langle S \rangle$			if $\langle E \rangle$ then $\langle S \rangle$ $\langle E \rangle$	sent		
$\langle E \rangle$		id				
$\langle E \rangle$	else $\langle S \rangle$					λ

5.6. Gramáticas RLL(1)

Si la gramática de la que partimos es una GPDR, tenemos dos opciones para construir un analizador. La más obvia es transformarla, utilizando los métodos que presentamos en el punto 4.2, en una gramática incontextual para la que después construiríamos el analizador. Sin embargo, vamos a ver cómo podemos construir los analizadores directamente a partir de las propias GPDR. Esto nos llevará a los analizadores RLL(1), en cierto sentido los equivalentes a los LL(1) para GPDR. Estudiaremos la construcción de manera intuitiva, sin entrar en una formalización rigurosa.

En primer lugar, debemos darnos cuenta de que ahora podrán aparecer en la pila de predicciones, y por lo tanto en la tabla, no sólo terminales y no terminales, sino también expresiones regulares en forma de clausuras o disyunciones.

Comenzaremos por un ejemplo para entenderlo mejor. Sea la GPDR:

$$\begin{aligned} \langle E \rangle &\rightarrow \langle T \rangle ((+|-)\langle T \rangle)^* \\ \langle T \rangle &\rightarrow \text{id} | \text{cte} | \text{"("} \langle E \rangle \text{"} \end{aligned}$$

Vamos a analizar la cadena **id+cte**. Utilizando nuestra representación habitual, empezamos por:

id	+cte\$	id	+cte\$	id	+cte\$
$\langle E \rangle$	\$	$\langle T \rangle$	$((+ -)\langle T \rangle)^* \mathbf{\$}$	id	$((+ -)\langle T \rangle)^* \mathbf{\$}$

Como coinciden los terminales, podemos avanzar en el análisis. Pero ahora nos encontramos con una importante diferencia, lo que avanza es la clausura, así que llegamos a:

$$\begin{array}{c|c} + & \text{cte\$} \\ \hline ((+|-)\langle T \rangle)^* & \mathbf{\$} \end{array}$$

Vemos que no hay problema en la predicción, ya que **+** es uno de los primeros de $((+|-)\langle T \rangle)^*$. Sin embargo, al encontrarnos frente a una clausura, debemos “expandirla”. Es más fácil verlo analizando cómo podemos “crear” el lenguaje de $(\rho)^*$. Podemos imaginar que primero se creará algo a partir de ρ y que después vendrá una nueva copia de la clausura, es decir, que nos encontraremos $\rho(\rho)^*$. En nuestro ejemplo, tendremos que predecir $(+|-)\langle T \rangle$ seguido de $((+|-)\langle T \rangle)^*$:

$$\begin{array}{c|c} + & \text{cte\$} \\ \hline (+|-) & \langle T \rangle ((+|-)\langle T \rangle)^* \mathbf{\$} \end{array}$$

El tope de la pila ahora es la disyunción, por lo que tendremos que predecir alguna de sus alternativas, lo que es fácil:

$$\begin{array}{c|c} + & \text{cte\$} \\ \hline + & \langle T \rangle ((+|-)\langle T \rangle)^* \mathbf{\$} \end{array}$$

Ahora avanzamos y expandimos $\langle T \rangle$ normalmente:

cte	\$	cte	\$	\$	
$\langle T \rangle$	$((+ -)\langle T \rangle)^* \mathbf{\$}$	cte	$((+ -)\langle T \rangle)^* \mathbf{\$}$	$((+ -)\langle T \rangle)^* \mathbf{\$}$	\$

Nuevamente encontramos la clausura. En este caso, no tenemos uno de sus primeros en la entrada, pero sí uno de sus siguientes, por lo que podemos sustituirla por la cadena vacía y llegamos a:

$$\frac{\$}{\$}$$

con lo que terminamos el análisis.

5.6.1. Cálculo de primeros

Ahora debemos adaptar el algoritmo presentado en el punto 5.3 para las GPDR. La adaptación es sencilla: basta con añadir los casos de la clausura y la disyunción.

¿Cuáles serán los primeros de una forma sentencial $(\rho)^*\beta$? Por un lado, serán los primeros de ρ ; además, la clausura puede producir la cadena vacía, así que también tendremos los primeros de β . En cuanto a la disyunción, tendremos que los primeros son la unión de los primeros de cada una de las opciones. Si alguna de las opciones es anulable, habrá que mirar el resto de la forma sentencial.

Teniendo esto en cuenta, el algoritmo queda:

Algoritmo primeros $(\alpha : (N \cup \Sigma)^*)$

$C = \emptyset$

si $\alpha = a\beta$ **entonces**

$C := \{a\}$;

si no si $\alpha = \langle A \rangle \beta$ **entonces**

para todo $\langle A \rangle \rightarrow \gamma \in P$ **hacer** $C := C \cup \text{primeros}(\gamma)$; **fin para**

si anulable $(\langle A \rangle)$ **entonces** $C := C \cup \text{primeros}(\beta)$; **fin si**

si no si $\alpha = (\rho)^*\beta$ **entonces**

$C := \text{primeros}(\rho) \cup \text{primeros}(\beta)$;

si no si $\alpha = (\rho_1|\rho_2|\dots|\rho_n)\beta$ **entonces**

$C := \bigcup_{i=1}^n \text{primeros}(\rho_i)$;

si anulable $(\rho_1|\rho_2|\dots|\rho_n)$ **entonces** $C := C \cup \text{primeros}(\beta)$; **fin si**

fin si

devolver C ;

fin primeros

5.6.2. Cálculo de anulables

Una modificación similar (figura 1) se puede hacer para el algoritmo que averigua si una forma sentencial es anulable. Basta con observar que la clausura siempre es anulable y que una disyunción lo es si lo es alguna de sus opciones.

5.6.3. Cálculo de siguientes

El cálculo de los siguientes se hace de manera similar al caso de las gramáticas incontextuales. Simplemente tendremos que aumentar las ecuaciones, teniendo en cuenta además que ahora tendremos que calcular los siguientes de clausuras y disyunciones.

En el resto de este punto, supondremos que ρ es un no terminal, una disyunción o una clausura. Las ecuaciones de la página 23 se pueden adaptar directamente:

- El fin de la entrada sigue al inicial:

$$\{\$ \} \subseteq \text{siguientes}(\langle S \rangle). \quad (4)$$

- Si $\langle A \rangle \rightarrow \alpha\rho\beta \in P$, entonces:

$$\text{primeros}(\beta) \subseteq \text{siguientes}(\rho); \quad (5)$$

```

Algoritmo anulable ( $\alpha : (N \cup \Sigma)^*$ )
  si  $\alpha = \lambda$  entonces
    devolver cierto;
  si no si  $\alpha = a\beta$  entonces
    devolver falso;
  si no si  $\alpha = \langle A \rangle \beta$  entonces
    si existe  $\gamma$  tal que  $\langle A \rangle \rightarrow \gamma \in P$  y anulable( $\gamma$ ) entonces
      devolver anulable( $\beta$ );
    si no
      devolver falso;
    fin si
  si no si  $\alpha = (\rho)^* \beta$  entonces
    devolver anulable( $\beta$ );
  si no si  $\alpha = (\rho_1 | \rho_2 | \dots | \rho_n) \beta$  entonces
    si existe  $i$  tal que anulable( $\rho_i$ ) entonces
      devolver anulable( $\beta$ );
    si no
      devolver falso;
    fin si
  fin si
fin anulable

```

Figura 1: Cálculo de anulables en GPDR.

- si, además, anulable(β) entonces:

$$\text{siguientes}(\langle A \rangle) \subseteq \text{siguientes}(\rho). \quad (6)$$

Además debemos considerar los casos en que ρ esté en una clausura o una disyunción. Podemos entender mejor las ecuaciones que siguen si observamos que $(\alpha\rho\beta)^*$ tiene dos posibles reescrituras: la cadena vacía y $(\alpha\rho\beta)(\alpha\rho\beta)^*$. En este segundo caso, lo que aparece tras ρ es β , por lo que sus primeros son siguientes de ρ . Si β es anulable, detrás de ρ podrá aparecer cualquier cosa que pueda seguir a la clausura o, si consideramos la reescritura $(\alpha\rho\beta)(\alpha\rho\beta)(\alpha\rho\beta)^*$, los primeros de $(\alpha\rho\beta)$. En resumen:

- Si $(\alpha\rho\beta)^*$ aparece en la parte derecha de alguna regla, entonces:

$$\text{primeros}(\beta) \subseteq \text{siguientes}(\rho) \quad (7)$$

- y si anulable(β) entonces

$$\text{primeros}(\alpha\rho\beta) \cup \text{siguientes}((\alpha\rho\beta)^*) \subseteq \text{siguientes}(\rho). \quad (8)$$

Cuando ρ aparece en una disyunción, el cálculo es más sencillo:

- Si $(\alpha_1 | \dots | \beta\rho\gamma | \dots | \alpha_n)$ aparece en la parte derecha de una regla, entonces:

$$\text{primeros}(\gamma) \subseteq \text{siguientes}(\rho) \quad (9)$$

- y si anulable(γ) entonces

$$\text{siguientes}((\alpha_1 | \dots | \beta\rho\gamma | \dots | \alpha_n)) \subseteq \text{siguientes}(\rho). \quad (10)$$

Si estas reglas te resultan engorrosas, siempre puedes convertir la GPDR en una gramática incontextual y calcular sobre ella los siguientes.

5.6.4. Tabla de análisis

La construcción de la tabla de análisis es análoga a la que hacíamos en el caso de las gramáticas LL(1). La única diferencia reseñable es que tendremos filas para las clausuras y disyunciones presentes en la gramática.

Los pasos que seguimos son:

- Averiguar qué no terminales y disyunciones son anulables.
- Calcular primeros de las partes derechas de las reglas y de los cuerpos de las clausuras y las disyunciones.
- Calcular siguientes de los no terminales anulables, de las clausuras y de las disyunciones anulables.
- Por cada regla $\langle A \rangle \rightarrow \alpha$, poner α en las celdas $M[\langle A \rangle, \mathbf{a}]$ tales que:
 - $\mathbf{a} \in \text{primeros}(\alpha)$
 - o $\text{anulable}(\alpha)$ y $\mathbf{a} \in \text{siguientes}(\langle A \rangle)$.
- Por cada clausura $(\rho)^*$:
 - Introducir $\rho(\rho)^*$ en las celdas $M[(\rho)^*, \mathbf{a}]$ tales que $\mathbf{a} \in \text{primeros}(\rho)$.
 - Si ρ es anulable, introducir $\rho(\rho)^*$ en las celdas $M[(\rho)^*, \mathbf{a}]$ tales que $\mathbf{a} \in \text{siguientes}((\rho)^*)$.
 - Introducir λ en las celdas $M[(\rho)^*, \mathbf{a}]$ tales que $\mathbf{a} \in \text{siguientes}((\rho)^*)$.
- Por cada disyunción $(\rho_1 | \dots | \rho_n)$:
 - Introducir ρ_i en las celdas $M[(\rho_1 | \dots | \rho_n), \mathbf{a}]$ tales que $\mathbf{a} \in \text{primeros}(\rho_i)$.
 - Introducir ρ_i en las celdas $M[(\rho_1 | \dots | \rho_n), \mathbf{a}]$ tales que $\mathbf{a} \in \text{siguientes}(\rho_1 | \dots | \rho_n)$ y $\text{anulable}(\rho_i)$.

Si no hay conflictos en la tabla, decimos que la GPDR es RLL(1).

Vamos a calcular la tabla de análisis de la GPDR del ejemplo:

$$\begin{aligned} \langle E \rangle &\rightarrow \langle T \rangle ((+|-)\langle T \rangle)^* \\ \langle T \rangle &\rightarrow \text{id|cte|“(”}\langle E \rangle\text{”} \end{aligned}$$

El único elemento anulable es la clausura $((+|-)\langle T \rangle)^*$. Necesitaremos saber los primeros de cada no terminal y de las clausuras y disyunciones:

Elemento	Primeros
$\langle E \rangle$	{id, cte, “(”}
$\langle T \rangle$	{id, cte, “(”}
$(+ -)$	{+, -}
$((+ -)\langle T \rangle)^*$	{+, -}

Para calcular los siguientes, comenzamos con la tabla:

Elemento	Siguientes
$\langle E \rangle$	{“)”, \$}
$\langle T \rangle$	{+, -}
$(+ -)$	{id, cte, “(”}
$((+ -)\langle T \rangle)^*$	\emptyset

Extendemos los siguientes de $\langle E \rangle$:

Elemento	Siguientes
$\langle E \rangle$	{“”, \$}
$\langle T \rangle$	{+, -, “”, \$}
(+ -)	{id, cte, “(”}
$((+ -)\langle T \rangle)^*$	{“”, \$}

Y ahora ya podemos calcular la tabla (por problemas de espacio, la dibujamos en dos partes):

	id	cte	+
$\langle E \rangle$	$\langle T \rangle ((+ -)\langle T \rangle)^*$	$\langle T \rangle ((+ -)\langle T \rangle)^*$	
$\langle T \rangle$	id	cte	
(+ -)			+
$((+ -)\langle T \rangle)^*$			$(+ -)\langle T \rangle ((+ -)\langle T \rangle)^*$

	-	“(”	“)”	\$
$\langle E \rangle$		$\langle T \rangle ((+ -)\langle T \rangle)^*$		
$\langle T \rangle$		“(” $\langle E \rangle$ “)”		
(+ -)				
$((+ -)\langle T \rangle)^*$	$(+ -)\langle T \rangle ((+ -)\langle T \rangle)^*$		λ	λ

EJERCICIO 28 — Utiliza la tabla anterior para analizar la sentencia $(id + cte) - cte$.

6. Implementación del analizador sintáctico

6.1. Implementación mediante tablas

En principio, podemos implementar el analizador sintáctico utilizando el algoritmo de análisis que presentamos en 5.2. Bastaría con buscar una implementación adecuada de la tabla. Sin embargo, este tipo de implementación dificulta la integración de acciones semánticas junto con el análisis, por lo que estudiaremos una implementación alternativa en la que la gramática se codifica directamente en forma de programa.

6.2. Implementación de analizadores descendentes recursivos

Podemos reflejar la tabla de análisis (o la gramática) directamente en un programa. La idea básica es tener una función (o un método si decidimos hacer que el analizador sea un objeto) por cada uno de los no terminales. Así al no terminal $\langle A \rangle$ le asociaríamos la función (método) `Analiza_A`. De momento, no utilizaremos parámetros. El resultado de la llamada será el valor cierto si hemos encontrado en la entrada una secuencia que se pueda derivar de $\langle A \rangle$ y falso en caso contrario. El análisis se hará mediante una llamada a `Analiza_S`: si devuelve cierto, se comprobará que hemos llegado al final de la entrada; si no, ha habido un error.

Como en el caso del algoritmo de análisis por tabla, utilizaremos la variable a para guardar la componente que estamos analizando en este momento. Si nuestro analizador es un objeto, podemos utilizar un atributo para a . Si el analizador es un conjunto de funciones, tendremos que hacer que a sea una variable global. En cualquier caso, antes de comenzar (antes de la llamada a `Analiza_S`) habrá que realizar la asignación $a := siguiente()$ (lógicamente, si el analizador léxico es un objeto, la llamada sería algo parecido a $a := alex.siguiente()$).

Ahora veremos una serie de esquemas que nos permitirán construir el analizador. Hemos utilizado la abreviatura $acceptables_\beta$ para referirnos al conjunto:

$$acceptables_\beta(\alpha) = \begin{cases} primeros(\alpha), & \text{si no anulable}(\alpha), \\ primeros(\alpha) \cup siguientes(\beta), & \text{si anulable}(\alpha), \end{cases}$$

donde se supone que estamos analizando β , que será un no terminal, una disyunción o una clausura. Si te das cuenta, $\text{acceptables}_{\beta}(\alpha)$ es el conjunto de aquellos terminales t para los cuales la celda $M[\beta, t] = \alpha$.

Si el no terminal $\langle A \rangle$ tiene las reglas $\langle A \rangle \rightarrow \alpha_1, \langle A \rangle \rightarrow \alpha_2, \dots, \langle A \rangle \rightarrow \alpha_n$, la estructura de `Analiza_A` es:

```

Función Analiza_A
  opción a
     $\text{acceptables}_{\langle A \rangle}(\alpha_1)$ : // Código de  $\alpha_1$ 
     $\text{acceptables}_{\langle A \rangle}(\alpha_2)$ : // Código de  $\alpha_2$ 
    ...
     $\text{acceptables}_{\langle A \rangle}(\alpha_n)$ : // Código de  $\alpha_n$ 
  si no // hemos encontrado un error
    devolver falso;
  fin opción
  devolver cierto;
fin Analiza_A

```

Las acciones de las partes derechas dependerán de la forma que tengan estas partes derechas. Si $\alpha = \beta_1\beta_2 \dots \beta_n$, tendremos que realizar el reconocimiento secuencialmente:

```

// Código de  $\alpha$ :
  // Código de  $\beta_1$ 
  // Código de  $\beta_2$ 
  ...
  // Código de  $\beta_n$ 

```

Si $\alpha = \langle A \rangle$, tenemos que llamar a la función de análisis correspondiente y comprobar que no ha habido error. La acción queda entonces:

```

// Código de  $\langle A \rangle$ :
si Analiza_A()= falso entonces
  devolver falso // transmitimos el error
fin si

```

Si $\alpha = t$ es un terminal, tendremos que comprobar que es el que hemos leído en ese momento. En caso de que no lo sea, hemos encontrado un error:

```

// Código de  $t$ 
si  $a \neq t$  entonces
  devolver falso // hemos encontrado un error
fin si
 $a := \text{siguiente}()$ ;

```

La acción correspondiente a λ es la acción nula (`pass` en Python).

Si $\alpha = (\rho)^*$, tenemos que entrar en un bucle y mantenernos en él mientras estemos en uno de los primeros de ρ . Al salir debemos comprobar que hemos encontrado uno de los siguientes de ρ^2 :

```

// Código de  $(\rho)^*$ 
mientras  $a \in \text{primeros}((\rho)^*)$  hacer
  // Código de  $\rho$ 
fin mientras
si  $a \notin \text{siguientes}((\rho)^*)$  entonces
  devolver falso // hemos encontrado un error
fin si

```

²En realidad, no es estrictamente necesario hacer esa comprobación. Si no se hace, el error también será detectado, pero más tarde, lo que puede hacer más difícil dar un mensaje adecuado.

Finalmente, si $\alpha = (\rho_1|\rho_2|\dots|\rho_n)$, tenemos que elegir la opción adecuada entre los ρ_i :

```
// Código de  $(\rho_1|\rho_2|\dots|\rho_n)$ 
opción a
  aceptables $_{(\rho_1|\dots|\rho_n)}$ ( $\rho_1$ ): // Código de  $\rho_1$ 
  aceptables $_{(\rho_1|\dots|\rho_n)}$ ( $\rho_2$ ): // Código de  $\rho_2$ 
  ...
  aceptables $_{(\rho_1|\dots|\rho_n)}$ ( $\rho_n$ ): // Código de  $\rho_n$ 
si no // hemos encontrado un error
  devolver falso
fin opción
```

Para la gramática de nuestro ejemplo, el analizador sería el que se ve en la figura 2. Este código es el que obtenemos directamente de los esquemas. Lógicamente, se puede mejorar con facilidad para eliminar muchas de las comprobaciones redundantes. El resultado lo puedes ver en la figura 3.

6.3. Tratamiento de errores

El tratamiento de errores es una de las partes de la compilación que más depende del ingenio del programador. Siempre es posible encontrar mejoras a la estrategia de recuperación que se sigue y obtener mensajes de error más ajustados al verdadero problema. Lo que veremos son algunas ideas generales que después tendrás que ajustar a tus propias necesidades.

6.3.1. Conceptos generales

Esperamos del analizador que haga algo más que aceptar/rechazar un programa:

- debe informar de la presencia de errores con tanta claridad y precisión como sea posible,
- y debe tratar de recuperarse del error y seguir con el análisis para detectar otros (eventuales) errores, aprovechando así una única ejecución del compilador para detectar el mayor número posible de errores.

Pero, en cualquier caso, la detección de errores y su recuperación no deben afectar significativamente a la complejidad del analizador sintáctico.

6.3.2. ¿Dónde se detecta un error?

Los analizadores LL(1) tienen la propiedad del **prefijo viable**: los errores se detectan en el momento en que se ha leído el prefijo más corto de la entrada que no es válido para ninguna cadena del lenguaje. Por lo tanto, este es el primer punto donde se puede anunciar un error analizando la entrada símbolo a símbolo de izquierda a derecha.

El compilador debe señalar el lugar donde se detecta el error. Típicamente se indica la *línea* donde se ha detectado el error y un código o *mensaje informativo* del tipo de error.

6.3.3. ¿Qué acciones debemos emprender al detectar un error?

En principio, podemos pensar en abandonar la compilación (implícitamente, eso hemos hecho hasta ahora). Sin embargo esto no es adecuado más que en caso de entornos muy interactivos. En general, compilar es un proceso costoso y debemos aprovechar la compilación para detectar tantos errores como sea posible.

Otra solución es intentar modificar el programa fuente con la esperanza de adivinar lo que quiso escribir el programador. Generaríamos entonces código de acuerdo con nuestra hipótesis. Esto es muy peligroso:

- En general, es muy difícil saber qué quería exactamente el programador.

```

Objeto Analizador
Atributos a, alex
Método inicialización(al:AnalizadorLéxico)
  alex := al;
  a := alex.siguiente();
  si analiza_E() y a = $ entonces
    devolver cierto
  si no
    devolver falso
  fin si
fin inicialización
Método analiza_E()
  opción a
  id, cte, “(“: //⟨E⟩ → ⟨T⟩((+|-)⟨T⟩)*
    si analiza_T()=falso entonces
      devolver falso
    fin si
    mientras a ∈ {+, -} hacer
      opción a:
        +:
          si a ≠ + entonces
            devolver falso
          fin si
          a := alex.siguiente()
        -:
          si a ≠ - entonces
            devolver falso
          fin si
          a := alex.siguiente()
      otro: devolver falso
    fin opción
    si analiza_T()= falso entonces
      devolver falso
    fin si
  fin mientras
  si a ∉ {“(”, “)”}:
    devolver falso
  fin si
si no
  devolver falso
fin opción
devolver cierto
fin analiza_E()

Método analiza_T()
  opción a:
  id: //⟨T⟩ → id
    si a ≠ id entonces
      devolver falso
    fin si
    a := alex.siguiente()
  cte: //⟨T⟩ → cte
    si a ≠ cte entonces
      devolver falso
    fin si
    a := alex.siguiente()
  “(“: //⟨T⟩ → “(”⟨E⟩“”
    si a ≠ “(” entonces
      devolver falso
    fin si
    a := alex.siguiente()
    si analiza_E()=falso entonces
      devolver falso
    fin si
    si a ≠ “)” entonces
      devolver falso
    fin si
    a := alex.siguiente()
  si no
    devolver falso
  fin opción
  devolver cierto
fin analiza_T()
fin Analizador

```

Figura 2: Analizador para la gramática de ejemplo.

<pre> Objeto Analizador Atributos a, F Método inicialización(al:AnalizadorLéxico) alex := al; a := alex.siguiete(); devolver analiza_E() y a = \$; fin inicialización Método analiza_E() si a ∈ {id, cte, "("} entonces si analiza_T()=falso entonces devolver falso fin si mientras a ∈ {+, -} hacer a := alex.siguiete() si analiza_T()=falso entonces devolver falso fin si fin mientras si a ∉ {"", "\$}: devolver falso fin si si no devolver falso fin si devolver cierto fin analiza_E() </pre>	<pre> Método analiza_T() opción a: id,cte: //⟨T⟩ → id cte a := alex.siguiete() "(": //⟨T⟩ → "("⟨E⟩" a := alex.siguiete() si analiza_E()=falso entonces devolver falso fin si si a ≠ ")" entonces devolver falso fin si a := alex.siguiete() si no devolver falso fin opción devolver cierto fin analiza_T() fin Analizador </pre>
--	---

Figura 3: El analizador de la figura 2 simplificado.

- Puede pasar inadvertido que el código generado ha sido creado a partir de un fuente con errores, lo que puede traer malas consecuencias.
- En el mejor de los casos, será una pérdida de tiempo porque tras corregirse el error se volverá a compilar.

Así pues, la generación de código debe detenerse.

Podemos resumir lo anterior diciendo que debemos recuperarnos el error y seguir analizando, pero sin generar código.

Hay diversas estrategias de recuperación de errores:

- Entrar en **modo pánico**: se descartan todos los componentes léxicos de entrada desde la detección del error hasta encontrar un *componente léxico de sincronización*.
- Modelar los errores mediante **producciones de error**: se enriquece la gramática con producciones que modelan los errores más frecuentes y se asocia a éstas una acción de tratamiento de error.
- Tratar el error a **nivel de frase**: se efectúa una corrección local, sustituyendo un prefijo del resto del programa por una secuencia que permita seguir analizando.
- Efectuar una **corrección global**: se busca el programa correcto más próximo al erróneo.

Las dos últimas opciones se pueden considerar únicamente de interés teórico por su elevado coste computacional y el poco beneficio que dan en la práctica.

La segunda opción tiene más interés, aunque obtener unas buenas reglas es bastante costoso. Un ejemplo sería añadir a la gramática para las expresiones aritméticas una regla como $\langle F \rangle \rightarrow "(" \langle E \rangle$ y después emitir el mensaje de error correspondiente a la falta de paréntesis cuando se detectara este error. Como puedes imaginar, este tipo de reglas dan lugar a todo tipo de conflictos por lo que la escritura del compilador se complica enormemente. A cambio, los mensajes de error que se pueden emitir son muy precisos y la recuperación muy buena.

El método que utilizaremos será el tratamiento en modo pánico.

6.3.4. Tratamiento de errores en modo pánico

Como se ha comentado arriba, el tratamiento de errores en modo pánico consiste en detener el análisis sintáctico e ir leyendo componentes léxicos hasta encontrar algún componente que permita la “sincronización”. Si lo consideramos desde la perspectiva del analizador, la sincronización puede conseguirse de tres maneras:

- Podemos proseguir el análisis como si hubiéramos encontrado la categoría léxica que esperábamos. Este es el tipo de tratamiento que haríamos por ejemplo al no encontrar un punto y coma donde debería estar. En este caso podemos sincronizarnos con los siguientes del elemento que esperábamos encontrar y con el propio elemento. Si encontramos el propio elemento al sincronizarnos, lo eliminaremos antes de continuar el análisis.
- Podemos proseguir el análisis como si hubiésemos encontrado la estructura (no terminal, clausura o disyunción) en la que hemos encontrado el error. Esta estrategia supone sincronizarse con alguno de los siguientes de la estructura y devolver (en caso del no terminal) el control como si todo hubiera ido bien. Podríamos aplicarla, por ejemplo, si estamos intentando encontrar una expresión en una llamada a función.
- Finalmente, podemos reintentar analizar la estructura en caso de que encontremos alguno de los primeros de la estructura que estamos analizando.

En cualquier caso, el fin de la entrada debe de ser siempre una de las categorías con las que nos podamos sincronizar.

Para todo esto, es útil disponer de una función de sincronización a la que se le pasa un conjunto S de categorías con las que podemos sincronizarnos. La función sería:

```

Función sincroniza( $S \subseteq \Sigma \cup \{\$\}$ )
   $S := S \cup \{\$\}$  // Añadimos el fin de la entrada
  mientras  $a \notin S$  hacer
     $a := \text{alex.siguiente}()$ 
  fin mientras
fin sincroniza

```

En nuestro ejemplo, una posible estrategia de tratamiento de errores sería:

- Si no encontramos el primero de una $\langle E \rangle$, hay tres posibilidades:
 - Hemos encontrado el fin de fichero. Emitimos error “Fin de fichero inesperado” y salimos.
 - Hemos encontrado un paréntesis cerrado. Emitimos error “Paréntesis cerrado no esperado”, sincronizamos con **id**, **cte** o paréntesis abierto y continuamos el análisis.
 - Hemos encontrado un $+$ o un $-$. Emitimos error “Falta operando”, sincronizamos con **id**, **cte** o paréntesis abierto y continuamos el análisis.
- Si después de analizar $((+|-)\langle T \rangle)^*$, no hay un paréntesis ni fin de fichero, emitimos error “Falta operador”, y nos sincronizamos con **id**, **cte** o paréntesis abierto (que, además, es lo único que puede haber en la entrada).
- Los casos para $\langle T \rangle$ son análogos a los de $\langle E \rangle$, salvo que estemos esperando el paréntesis cerrado. Entonces, emitimos error “Falta paréntesis cerrado” y nos sincronizamos con él o con cualquiera de los siguientes de $\langle T \rangle$.

7. Resumen del tema

- El analizador sintáctico encuentra las estructuras de la entrada.
- Especificamos el analizador mediante gramáticas incontextuales.
- El árbol de análisis muestra cómo derivar la entrada.
- Debemos trabajar con gramáticas no ambiguas o dar reglas para resolver las ambigüedades.
- Extendemos las gramáticas para utilizar partes derechas regulares.
- Dos tipos de análisis: ascendente y descendente.
- Empleamos análisis LL(1) (o RLL(1) para GPDR):
 - En cada momento, tomamos la decisión en base a un solo símbolo de la entrada.
 - No puede haber recursividad a izquierdas.
 - Las distintas alternativas tienen que diferenciarse por sus primeros (o siguientes si son anulables).
- Implementación: mediante analizadores descendentes recursivos.
- Tratamiento de errores:
 - Detiene la generación de código.
 - Recuperación en modo pánico.