

# Tema 3

## Introducción al análisis sintáctico

### Bibliografía:

- Aho, A.V., Sethi, R., Ullman, J.D. (1990), *Compiladores: principios, técnicas y herramientas*, Tema 4, pág.: 163-186.
- Louden, K.C. (1997), *Compiler Construction: Principles and Practice*, Tema 3, páginas: 93-140.

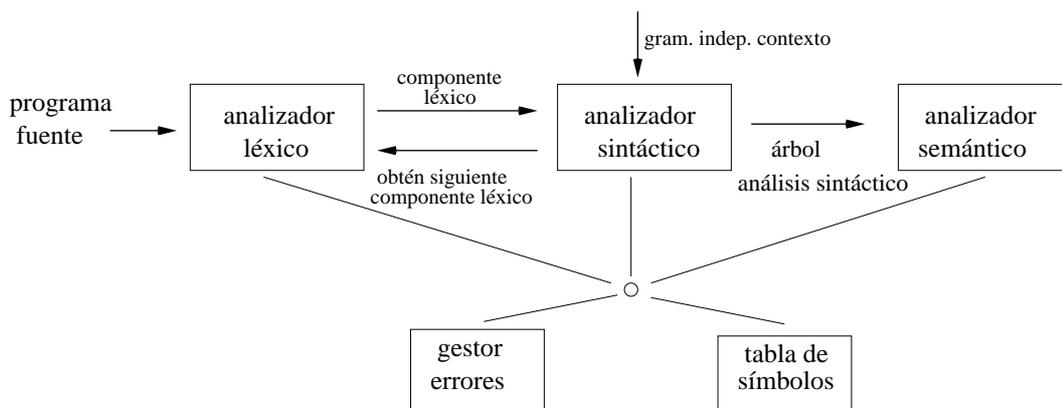
### Contenidos:

1. El proceso de análisis sintáctico.
2. Especificación sintáctica de los lenguajes de programación:
  - a) Gramáticas independientes del contexto versus expresiones regulares.
  - b) Gramáticas independientes del contexto versus gramáticas dependientes del contexto.
3. Derivaciones y árboles sintácticos.
4. Gramáticas limpias y bien formadas.
5. El problema de la ambigüedad en las gramáticas. Eliminación:
  - a) Mediante reglas de precedencia.
  - b) Mediante transformación.
6. Clasificación de los métodos de análisis sintáctico.

### 3.1. El proceso de análisis sintáctico

*Funciones del analizador sintáctico:*

- Comprobar si la cadena de componentes léxicos proporcionada por el analizador léxico puede ser generada por la gramática que define el lenguaje fuente (Gramática Independiente del Contexto, GIC).
- Construir el árbol de análisis sintáctico que define la estructura jerárquica de un programa y obtener la serie de derivaciones para generar la cadena de componentes léxicos. El árbol sintáctico se utilizará como representación intermedia en la generación de código.
- Informar de los errores sintácticos de forma precisa y significativa y deberá estar dotado de un mecanismo de recuperación de errores para continuar con el análisis.



El análisis sintáctico se puede considerar como una función que toma como entrada la secuencia de componentes léxicos producida por el análisis léxico y produce como salida el árbol sintáctico. En la realidad, el análisis sintáctico hace una petición al análisis léxico del componente léxico siguiente en la entrada (los símbolos terminales) conforme lo va necesitando en el proceso de análisis, conforme se mueve a lo largo de la gramática.

#### Reconocedor versus analizador sintáctico

Un reconocedor trata de determinar simplemente si la cadena puede o no ser generada por la gramática (salida booleana). Podríamos preguntar algo más. Si queremos reconocer las estructuras propias de los lenguajes de programación en el fichero de entrada para guiar la traducción, es necesario conocer su estructura jerárquica, el árbol de análisis sintáctico.

### 3.2. Especificación sintáctica de los lenguajes de programación

La mayoría de las construcciones de los lenguajes de programación se pueden representar con una gramática independiente del contexto (GIC). La mayoría de las construcciones de los lenguajes de programación implican recursividad y anidamientos.

$$G = \{S, V_T, V_{NT}, P\},$$

S: el axioma o símbolo de inicio

$V_T$ : conjunto de terminales, los componentes léxicos

$V_{NT}$ : conjunto de no-terminales

P: conjunto de reglas de producción de la forma

$V_{NT} \rightarrow X_1, \dots, X_n$ , con  $X_i \in (V_T \cup V_{NT})$

**Pregunta:** Si las expresiones regulares son un caso particular de GIC, **por qué no se incluye la especificación léxica como parte de la sintaxis?**

Los AFD son muy sencillos de implementar y son muy eficientes frente a los autómatas a pila necesarios para reconocer las GIC, la eficiencia del traductor se vería comprometida.

**¿Expresiones regulares o gramáticas independientes del contexto?**

Las expresiones regulares no permiten construcciones anidadas tan comunes en los lenguajes de programación: paréntesis equilibrados, concordancia de pares de palabras clave como `begin-end`, `do-while`, ...

Por ejemplo: consideremos el problema de los paréntesis equilibrados en una expresión aritmética. El hecho de que haya un

paréntesis abierto obliga a que haya un paréntesis cerrado. Este problema es similar a considerar el lenguaje

$$L = \{aba, aabaa, aaabaaa, \dots\} = \{a^n b a^n \mid n \geq 0\}$$

Si intentamos escribir una expresión regular, lo más próximo sería:  $\{a^* b a^*\}$  pero no se garantiza que el número de a's antes y después sea el mismo.

Las expresiones regulares NO SABEN contar. NO es posible especificar la estructura de un lenguaje de programación con sólo expresiones regulares.

**¿Son suficientes las gramáticas independientes del contexto?**

NO, existen ciertos aspectos de los lenguajes de programación que no se pueden representar con una GIC. A estos aspectos se les llama aspectos *semánticos*, son aspectos en general dependientes del contexto.

Por ejemplo, el problema de declaración de los identificadores antes de uso. Se podría representar mediante el lenguaje

$$L = \{w c w \mid w \in (a|b)^*\}$$

este lenguaje no puede ser representado por una GIC.

**Por qué no usar gramáticas contextuales?**

Su reconocimiento es demasiado complejo para ser eficiente en un problema práctico.

Será en el tema de análisis semántico donde se abordará estos problemas mediante el uso de acciones concretas (acciones semánticas) que no son propias de la gramática en sí. Será código escrito para hacer las comprobaciones semánticas necesarias (tipos, declaración antes de uso, ...)

Ejemplo de una gramática para un sencillo lenguaje:

---

Program →	<b>module id ; Declarations begin</b> StatementList <b>end id .</b>
Declarations →	<b>type id =</b> TypeDenoter ; TypeList Declarations   <b>var id IdList :</b> TypeDenoter ; VarList Declarations   <b>function id (</b> Parameter ParameterList <b>) :</b> SimpleType ; Declarations <b>begin</b> StatementList <b>end id ;</b> Declarations   $\epsilon$
TypeList →	<b>id =</b> TypeDenoter ; TypeList   $\epsilon$
VarList →	<b>id IdList :</b> TypeDenoter ; VarList   $\epsilon$
IdList →	, <b>id</b> IdList   $\epsilon$
Parameter →	TypeDenoter <b>id</b>   $\epsilon$
ParameterList →	, Parameter ParameterList   $\epsilon$
SimpleType →	<b>char   integer   real   boolean</b>
TypeDenoter →	SimpleType   <b>id</b>   CompositeType
CompositeType →	<b>array [ num_integer .. num_integer ] of</b> SimpleType   <b>record id</b> IdList : SimpleType ; FieldList <b>end</b>
FieldList →	<b>id</b> IdList : SimpleType ; FieldList   $\epsilon$
StatementList →	Statement ; StatementList   $\epsilon$
Statement →	AssignStm   IfStm   WhileStm   FunctionStm   ReturnStm   WriteStm   ReadStm
AssignStm →	LeftValue = Expression
LeftValue →	<b>id</b>   <b>id</b> [ Expression ]   <b>id . id</b>
IfStm →	<b>if</b> Expression <b>then</b> StatementList <b>else</b> ElseStm <b>end</b>
ElseStm →	<b>else</b> StatementList   $\epsilon$
WhileStm →	<b>while</b> Expression <b>do</b> StatementList <b>end</b>
FunctionStm →	<b>id</b> ( Arg ArgList )
ReturnStm →	<b>return</b> Expression
WriteStm →	<b>write</b> Expression
ReadStm →	<b>read</b> LeftValue
Arg →	Num   - Num   <b>id</b>   - <b>id</b>   $\epsilon$
Num →	<b>num_integer</b>   <b>num_real</b>
ArgList →	, Arg ArgList   $\epsilon$
Expression →	LeftValue   ( Expression )   <b>true</b>   <b>false</b>   Num   FunctionStm   Expression * Expression   Expression / Expression   Expression - Expression   Expression + Expression   - Expression   Expression <b>or</b> Expression   Expression <b>and</b> Expression   <b>not</b> Expression   Expression < Expression   Expression > Expression   Expression == Expression

---

Cuadro 3.1: Ejemplo de especificación sintáctica

### 3.3. Derivaciones. Árboles de anal. sintáctico

Dada  $\alpha, \beta \in (V_T \cup V_{NT})^*$  cadenas arbitrarias de símbolos gramaticales, se dice que

$\alpha A \beta \xRightarrow{*} \alpha \gamma \beta$  es una *derivación*, si existe una producción  $A \rightarrow \gamma$ .

Sea  $S \xRightarrow{*} \alpha$ ,

si  $\alpha \in (V_T \cup V_{NT})^*$  se dice que  $\alpha$  es una *forma de frase o forma sentencial*

si  $\alpha \in (V_T)^*$  se dice que  $\alpha$  es una *frase o sentencia*.

Existen dos tipos de derivaciones:

- *Derivación más a la izquierda*: se sustituye en cada paso el no-terminal más a la izquierda de la forma de frase. La denotaremos por:

$$\alpha \xRightarrow{*mi} \beta$$

- *Derivación más a la derecha*: se sustituye en cada paso el no-terminal más a la derecha de la forma de frase. La denotaremos por:

$$\alpha \xRightarrow{*md} \beta$$

Por ejemplo para la gramática:

$$E \rightarrow \text{id} \mid \text{num} \mid E + E \mid ( E ) \mid - E$$

Si queremos derivar la frase  $-(id + id)$

$$E \Rightarrow -E \Rightarrow -(E + E) \Rightarrow -(id + E) \Rightarrow -(id + id)$$

derivación más a la izquierda.

$$E \Rightarrow -E \Rightarrow -(E + E) \Rightarrow -(E + id) \Rightarrow -(id + id)$$

derivación más a la derecha.

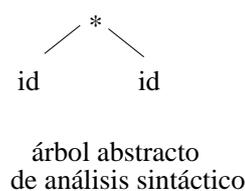
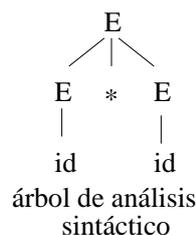
### 3.3.1. Árboles de análisis sintáctico

Un árbol de análisis sintáctico indica cómo a partir del axioma de la gramática se deriva una frase (cadena) del lenguaje. Dada una gramática independiente del contexto, un *árbol de análisis sintáctico* es un árbol tal que:

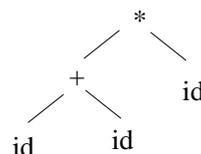
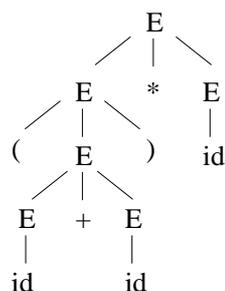
1. La raíz está etiquetada con el símbolo inicial.
2. Cada hoja está etiquetada con un componente léxico. Las hojas de izquierda a derecha forman la frase (el programa fuente).
3. Cada nodo interior está etiquetado con un no-terminal.
4. Si  $A$  es un no-terminal y  $X_1, X_2, \dots, X_n$  son sus hijos de izquierda a derecha, entonces existe la producción  $A \rightarrow X_1, X_2, \dots, X_n$ , con  $X_i \in (V_T \cup V_{NT})$ .

El árbol de análisis sintáctico contiene en general mucha más información que la estrictamente necesaria para generar el código. Se puede construir una estructura más sencilla, los *árboles abstractos de análisis sintáctico*. Ejemplo: expresiones aritméticas (igual semántica, menor complejidad)

Entrada: id \* id



Entrada: (id +id)\* id



### 3.4. Gramáticas limpias y bien formadas

Llamamos *símbolo vivo* al símbolo a partir del cual se puede derivar una cadena de terminales. Llamamos *símbolo muerto* a los símbolos no-vivos, no generan una cadena del lenguaje. Llamamos *símbolo inaccesible* si nunca aparece en la parte derecha de una producción. A las gramáticas que contienen estos tipos de símbolos se les llama *gramáticas sucias*.

*Teorema 1:* si todos los símbolos de la parte derecha de una producción son símbolos vivos, entonces el símbolo de la parte izquierda también lo es.

*Algoritmo para detectar símbolos muertos*

1. Hacer una lista de no-terminales que tengan al menos una producción con sólo símbolos terminales en la parte derecha.
2. Dada una producción, si todos los no-terminales de la parte derecha pertenecen a la lista, entonces podemos incluir en la lista al no-terminal de la parte izquierda de la producción.
3. Cuando ya no se puedan incluir más símbolos en la lista mediante la aplicación del paso 2, la lista contendrá los símbolos no-terminales vivos y el resto serán símbolos muertos.

*Teorema 2:* si el símbolo no-terminal de la parte izquierda de una producción es accesible, entonces todos los símbolos de la parte derecha también lo son.

*Algoritmo para detectar símbolos inaccesibles*

1. Se inicializa una lista de no-terminales que sabemos que son accesibles con el axioma.
2. Si la parte izquierda de una producción está en la lista, entonces se incluye en la misma al no-terminal que aparece en la parte derecha de la producción.
3. Cuando ya no se pueden incluir más símbolos a la lista mediante la aplicación del paso 2, entonces la lista contendrá todos los símbolos accesibles y el resto de los no-terminales serán innaccesibles.

Para limpiar una gramática primero se eliminan los símbolos muertos y después los símbolos inaccesibles.

Una gramática *está bien formada* si es limpia y además no contiene producciones- $\epsilon$ . Para eliminar producciones- $\epsilon$ , remitirse a la signatura de TALF.

Ejemplo: limpiar la gramática

$$\begin{aligned}
 S &\rightarrow a A B \mid A \\
 A &\rightarrow c B d \\
 B &\rightarrow e \mid f S \\
 C &\rightarrow g D \mid h D t \\
 D &\rightarrow x \mid y \mid z
 \end{aligned}$$

### 3.5. Gramáticas ambiguas

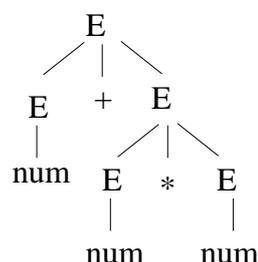
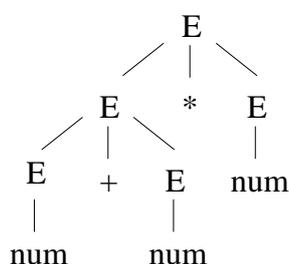
Una gramática es *ambigua* cuando para una determinada sentencia produce más de un árbol de derivación.

La gramática siguiente es ambigua:

$$E \rightarrow \text{id} \mid \text{num} \mid E + E \mid E * E \mid ( E ) \mid - E$$

Supongamos la sentencia  $\text{id} + \text{id} * \text{id}$

Entrada:  $\text{num} + \text{num} * \text{num}$



El significado semántico es DIFERENTE. No existe una única traducción posible. Se genera código diferente.

No existe un algoritmo que determine con certeza en un plazo finito de tiempo si una gramática es ambigua o no. A lo sumo que se ha llegado en el estudio de la ambigüedad es que hay algunas condiciones que de cumplirse determinan la ambigüedad, pero en el caso de no cumplirse no se puede decir que la gramática es no ambigua.

Necesidad de evitar las gramáticas ambiguas. **Cómo?**

- transformando la gramática o
- estableciendo precedencias entre operadores y de asociatividad.

Se puede eliminar la ambigüedad transformando la gramática agrupando todos los operadores de igual precedencia en grupos y asociando a cada uno una regla, de forma que los que tengan menor precedencia aparezcan más cercanos al símbolo de inicio, precedencia en cascada. Esto conlleva el aumento de la complejidad de la gramática y con ello en la del árbol sintáctico. La gramática deja de ser intuitiva.

**Ejemplo 1:** gramática ambigua
$$E \rightarrow \mathbf{num} \mid E + E \mid E - E \mid E * E \mid E / E$$

Si la transformamos, esta gramática ya no es ambigua.

$$\begin{aligned} \text{exp} &\rightarrow \text{exp} + \text{term} \mid \text{exp} - \text{term} \mid \text{term} \\ \text{term} &\rightarrow \text{term} * \text{factor} \mid \text{term} / \text{factor} \mid \text{factor} \\ \text{factor} &\rightarrow ( \text{exp} ) \mid \mathbf{num} \end{aligned}$$

**IMPORTANTE:** Una regla recursiva a izquierdas, hace que el operador sea asociativo por la izquierda. Una regla recursiva a derechas hace que el operador sea asociativo por la derecha.

Ejemplo 2: el problema del *dangling else* y el convenio de la asociación al `if` más cercano.

$$\begin{aligned} \text{statement} &\rightarrow \text{if\_stmt} \mid \mathbf{other} \\ \text{if\_stmt} &\rightarrow \mathbf{if} ( \text{exp} ) \text{statement} \\ &\quad \mid \mathbf{if} ( \text{exp} ) \text{statement} \mathbf{else} \text{statement} \\ \text{exp} &\rightarrow \mathbf{0} \mid \mathbf{1} \end{aligned}$$

Por ejemplo: para la sentencia `if (0) if (1) other else other` tenemos dos árboles de derivación.

Crear los dos árboles ...

### 3.6. Clasificación métodos de análisis sintáctico

Atendiendo a la forma en que se construye el árbol de análisis sintáctico los métodos de análisis sintáctico se clasifican:

- *Métodos descendentes*: se parte del símbolo de inicio de la gramática que se coloca en la raíz y se va construyendo el árbol desde arriba hacia abajo hasta las hojas, eligiendo la derivación que da lugar a una concordancia con la cadena de entrada. Se basa en la idea de *predice* una derivación y establece una *concordancia* con el símbolo de la entrada (*predict/match*). El análisis sintáctico descendente corresponde con un recorrido prefijo del árbol de análisis sintáctico (primero expandimos el nodo que visitamos y luego procesamos los hijos).

S → Sujeto Verbo Objeto

Sujeto → el Nombre | un Nombre

Objeto → el Nombre | un Nombre

Verbo → el Nombre | un Nombre

Nombre → perro | gato | león

Ejemplo: El león cazó un gato

Crear el árbol...

Problemas: recursión a izquierdas, diferentes alternativas en una producción, ¿cómo conseguir un coste lineal?

- *Métodos ascendentes*: se construye el árbol de análisis sintáctico desde las hojas hasta la raíz. En las hojas está la cadena a analizar y se intenta reducirla al símbolo de inicio de la gramática que está en la raíz. Se trata de *desplazar-se* en la cadena de entrada y encontrar una subcadena para aplicar una *reducción* (una regla de producción a la inversa), (*shift-reduce*). El análisis sintáctico ascendente corresponde con un recorrido postorden del árbol (primero reconocemos los hijos y luego mediante una reducción reconocemos el padre).

Insertar un ejemplo de como efectivamente los dos tipos de análisis corresponden con los tipos de recorrido prefijo y postfijo.

Atendiendo a la forma en que procesan la cadena de entrada se clasifican en:

- *Métodos direccionales*: procesan la cadena de entrada símbolo a símbolo de izquierda a derecha.
- *Métodos no-direccionales*: acceden a cualquier lugar de la cadena de entrada para construir el árbol. Necesitan tener toda la cadena de componentes léxicos en memoria. Más costosos, no se suelen implementar.

Atendiendo al número de alternativas posibles en una derivación se clasifican en:

- *Métodos deterministas*: dado un símbolo de la cadena de entrada se puede decidir en cada paso cuál es la alternativa/derivación adecuada a aplicar, sólo hay una posible. No se produce retroceso y el coste es lineal.
- *Métodos no-deterministas*: en cada paso de la construcción del árbol se deben probar diferentes alternativas/derivaciones para ver cual es la adecuada, con el correspondiente aumento del coste.

Importancia de los métodos direccionales y deterministas en el diseño de traductores. El determinismo por la necesidad de eficiencia (coste lineal frente a exponencial que le haría prohibitivo en el diseño práctico de traductores) y porque las herramientas para la generación de traductores asumen esta condición. Los

	descendentes		ascendentes
no direccionales	Unger		CYK
direccionales	no deterministas	Predice/Concuerda 1° en anchura 1° en profundida	Desplaza/Reduce 1° en anchura 1° en profundida
	deterministas	Predice/Concuerda Gram. LL(1)	Desplaza/Reduce Gram. LR(k) LR(0), SLR(1), LALR(1)

métodos direccionales por la propia naturaleza secuencial en que se va procesando los símbolos del fichero fuente de izquierda a derecha.

### 3.7. Ejercicios

1. (0.2 ptos) Dadas la siguiente gramática, elimina los símbolos muertos e inaccesibles.

$$\begin{aligned} S &\rightarrow a A B C \mid D d \\ A &\rightarrow a B C \\ B &\rightarrow b \\ C &\rightarrow c \\ D &\rightarrow d \\ E &\rightarrow g h \end{aligned}$$

2. (0.2 ptos) Dada la siguiente gramática, escribe la derivación más a la izquierda, el árbol de análisis sintáctico y el árbol sintáctico abstracto para las siguientes expresiones: (a)  $3+4*5-6$ ; (b)  $3*(4-5+6)$

$$\begin{aligned} \text{exp} &\rightarrow \text{exp addop term} \mid \text{term} \\ \text{addop} &\rightarrow + \mid - \\ \text{term} &\rightarrow \text{term mulop factor} \mid \text{factor} \\ \text{mulop} &\rightarrow * \mid / \\ \text{factor} &\rightarrow ( \text{exp} ) \mid \mathbf{num} \end{aligned}$$

3. (0.2 ptos) Dada la siguiente gramática simplificada de una expresión en LISP: (a) escribe las derivaciones más a la derecha y más a la izquierda de la entrada  $(a\ 23\ (m\ x\ y))$ . Dibuja el árbol de análisis sintáctico para esa sentencia.

$$\begin{aligned} \text{lexp} &\rightarrow \text{atom} \mid \text{list} \\ \text{atom} &\rightarrow \mathbf{number} \mid \mathbf{identifier} \\ \text{list} &\rightarrow ( \text{lexp-seq} ) \\ \text{lexp-seq} &\rightarrow \text{lexp-seq lexp} \mid \text{lexp} \end{aligned}$$

4. (0.2 ptos) Dada la siguiente gramática, construya el árbol de análisis sintáctico para la frase `not (true or false)`. ¿Es una gramática ambigua?

$$\begin{aligned} \text{bexp} &\rightarrow \text{bexp } \mathbf{or} \text{ bterm } | \text{ bterm} \\ \text{bterm} &\rightarrow \text{bterm } \mathbf{and} \text{ bfactor } | \text{ bfactor} \\ \text{bfactor} &\rightarrow \mathbf{not} \text{ bfactor } | ( \text{ bexp } ) | \mathbf{true} | \mathbf{false} \end{aligned}$$

5. (0.2 ptos) Escribe una gramática para expresiones booleanas que incluya las constantes `true` y `false`, los operadores `and`, `or` y `not` y los paréntesis. Se ha de tener en cuenta que el operador `or` tiene menor precedencia que el operador `and` y éste menor que el operador `not`. Construye la gramática de forma que los operadores `and` y `or` sean asociativos por la izquierda y el operador `not` por la derecha.
6. (0.2 ptos) Comprueba que se puede resolver el problema de la ambigüedad del *dangling else* transformando la gramática a la forma:

$$\begin{aligned} \text{statement} &\rightarrow \text{matched-stmt} | \text{unmatched-stmt} \\ \text{matched-stmt} &\rightarrow \mathbf{if} ( \text{exp} ) \text{matched-stmt } \mathbf{else} \text{matched-stmt} | \mathbf{other} \\ \text{unmatched-stmt} &\rightarrow \mathbf{if} ( \text{exp} ) \text{statement} \\ &\quad | \mathbf{if} ( \text{exp} ) \text{matched-stmt } \mathbf{else} \text{unmatched-stmt} \\ \text{exp} &\rightarrow 0 | 1 \end{aligned}$$

Justificar usando un ejemplo.

7. (0.2 ptos) Comprueba que el intento de resolver el problema de la ambigüedad del *dangling else* no es adecuado:

$$\begin{aligned} \text{statement} &\rightarrow \text{if} ( \text{exp} ) \text{statement} | \text{matched-stmt} \\ \text{matched-stmt} &\rightarrow \text{if} ( \text{exp} ) \text{matched-stmt } \mathbf{else} \text{statement} | \mathbf{other} \\ \text{exp} &\rightarrow 0 | 1 \end{aligned}$$

Justificar usando un ejemplo.