

Autómatas y Lenguajes Formales 2019-II
Facultad de Ciencias UNAM
Nota de Clase 11, Ambigüedad, gramáticas en lenguajes de
programación

Favio E. Miranda Perea A. Liliana Reyes Cabello Lourdes González Huesca

17 de mayo de 2019

1. Ambigüedad

Puede ser que dos derivaciones distintas tengan un mismo árbol de derivación y también puede suceder que haya más de un árbol de derivación para una cadena. Lo ideal es que cada cadena tenga sólo un árbol asociado. Si sucede lo anterior, implica que el lenguaje **no es ambiguo**. Desafortunadamente existen lenguajes ambiguos. Veamos una definición formal de este fenómeno:

Definición 1 Una gramática se dice **ambigua** si existe una palabra w con dos o más árboles de derivación distintos. En general una palabra puede tener más de una derivación, pero un sólo árbol y en tal caso no hay ambigüedad.

Algunas veces se puede suprimir la ambigüedad directamente. Sin embargo **no** hay un algoritmo para remover ambigüedad. Es decir el problema de eliminar la ambigüedad es indecidible.

Ejemplo: Considere la siguiente gramática:

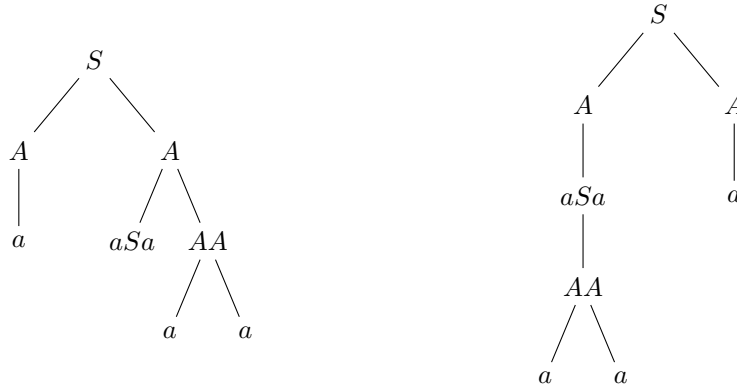
$$S \rightarrow AA \quad A \rightarrow aSa \mid a$$

La palabra a^5 tiene las siguientes derivaciones diferentes y por la izquierda:

$$S \rightarrow AA \rightarrow aA \rightarrow aaSa \rightarrow aaAAa \rightarrow aaaAa \rightarrow aaaaa$$

$$S \rightarrow AA \rightarrow aSaA \rightarrow aAAaA \rightarrow aaAaA \rightarrow aaaaA \rightarrow aaaaa$$

Las dos derivaciones generan árboles distintos:



Definición 2 Un lenguaje L es ambiguo si existe una gramática ambigua G que genera a L .
Un lenguaje es **inherentemente** ambiguo si todas las gramáticas que lo generan son ambiguas.

Ejemplo: El lenguaje $L = \{a^{2+3i} \mid i \geq 0\}$ es ambiguo por ser generado por la siguiente gramática ambigua:

$$S \rightarrow AA \quad A \rightarrow aSa \mid a$$

Sin embargo este lenguaje también es generado por una gramática no ambigua:

$$S \rightarrow aa \mid aaU \quad U \rightarrow aaaU \mid aaa$$

En este caso la derivación de a^5 es:

$$S \rightarrow aaU \rightarrow aaaaa$$

y por lo tanto L no es un lenguaje inherentemente ambiguo.

Ejemplo: El lenguaje

$$L = \{a^n b^n c^m d^m \mid n, m \geq 1\} \cup \{a^n b^m c^m d^n \mid n, m \geq 1\}$$

es generado por la gramática:

$$\begin{aligned} S &\rightarrow AB \mid C & A &\rightarrow aAb \mid ab & B &\rightarrow cBd \mid cd \\ C &\rightarrow aCd \mid aDd & D &\rightarrow bDc \mid bc \end{aligned}$$

La cadena $aabbccdd$ tiene dos derivaciones por la izquierda:

$$S \rightarrow AB \rightarrow aAbB \rightarrow aabbB \rightarrow aabbcBd \rightarrow aabbccdd$$

$$S \rightarrow C \rightarrow aCd \rightarrow aaDdd \rightarrow aabDcdd \rightarrow aabbccdd$$

Como se comentó antes, la ambigüedad no puede ser eliminada con un método y en general, probar la ambigüedad inherente es complicado.

Veamos ahora algunos ejemplos de gramáticas ambiguas en relación a lenguajes de programación

Ejemplo 1 Lenguaje de paréntesis balanceados:

$$S \rightarrow \varepsilon \mid (S) \mid SS$$

- *Ambigüedad:* $()()()$ tiene dos árboles de derivación
- *Gramática no ambigua equivalente:*

$$S \rightarrow \varepsilon \mid (S)S$$

Ejemplo 2 *Expresiones aritméticas:*

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

donde a es un terminal que representa a los identificadores y constantes.

- *Ambigüedad:* $a + a * a$
- *Gramática no ambigua equivalente:* se obtiene modelando la precedencia de operadores como sigue:

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

2. Gramáticas libres de contexto en lenguajes de programación

- El estudio formal de los lenguajes de programación se divide en sintaxis, pragmática y semántica.
- La semántica se encarga de definir el significado de las expresiones, enunciados y unidades de programa.
- La pragmática define la implementación del lenguaje basada en la metodología y estrategias de programación deseadas
- La sintaxis se encarga de definir la forma de las expresiones y enunciados de un lenguaje y se sirve fundamentalmente de los conceptos y herramientas de nuestro curso.
- Antes del proceso de evaluación, un compilador e intérprete necesita realizar los procesos de análisis léxico y sintáctico, descritos a continuación a grandes rasgos.

2.1. Análisis léxico y sintáctico

- **Análisis léxico:** se encarga de transformar el programa fuente en una lista de unidades sintácticas de bajo nivel llamadas lexemas, los cuales se clasifican en distintas categorías llamadas *tokens*, como pueden ser identificadores, constantes, separadores, etc.
- El análisis léxico se sirve fundamentalmente de expresiones regulares para su definición y reconocimiento.
- **Análisis sintáctico:** se encarga de transformar la lista de lexemas en un programa objeto, el cual es una expresión válida de la llamada sintaxis abstracta del lenguaje. Este programa es esencialmente un árbol de derivación dictado por una gramática libre de contexto que define al lenguaje de programación.
- Por lo tanto el análisis sintáctico es esencialmente una forma del problema de la pertenencia en gramáticas libres de contexto.

2.2. Forma de Backus-Naur

- Las gramáticas libres de contexto para lenguajes de programación suelen escribirse en la forma de Backus-Naur o BNF.
- Este método de definición de gramáticas fue introducido por John Backus para el lenguaje ALGOL 58 en 1959 y fue mejorado por Peter Naur para la definición de ALGOL 60.
- Este sistema notacional para definir lenguajes libres de contexto sigue las siguientes convenciones:
 - El símbolo de reescritura \rightarrow se reemplaza con $::=$.
 - El símbolo $|$ significa *o* y abreviar la definición de producciones de una misma variable.
 - Las variables se escriben entre paréntesis triangulares y por lo general utilizan nombres largos que ayuden a la descripción de las categorías del lenguaje.

Ejemplo 3 Algunos ejemplos de gramáticas en FBN son:

- *Lenguaje de paréntesis balanceados*

$$\begin{aligned} \langle \text{parent_balanc} \rangle & ::= \varepsilon \mid \\ & \quad (\langle \text{parent_balanc} \rangle) \mid \\ & \quad \langle \text{parent_balanc} \rangle \langle \text{parent_balanc} \rangle \end{aligned}$$

- *Expresiones aritméticas:*

$$\begin{aligned} \langle \text{expr} \rangle & ::= \langle \text{expr} \rangle \langle \text{op} \rangle \langle \text{expr} \rangle \mid \\ & \quad (\langle \text{expr} \rangle) \mid \langle \text{id} \rangle \\ \langle \text{op} \rangle & ::= + \mid - \mid * \mid / \\ \langle \text{id} \rangle & ::= a \mid b \mid c \end{aligned}$$

- *Bloques de asignación de expresiones aritméticas:*

$$\begin{aligned} \langle \text{programa} \rangle & ::= \text{begin} \langle \text{sec_enunc} \rangle \text{end} \\ \langle \text{sec_enunc} \rangle & ::= \langle \text{enunc} \rangle \mid \langle \text{enunc} \rangle ; \langle \text{sec_enunc} \rangle \\ \langle \text{enunc} \rangle & ::= \langle \text{id} \rangle := \langle \text{expr} \rangle \\ \langle \text{expr} \rangle & ::= \langle \text{expr} \rangle \langle \text{op} \rangle \langle \text{expr} \rangle \mid \\ & \quad (\langle \text{expr} \rangle) \mid \langle \text{id} \rangle \\ \langle \text{op} \rangle & ::= + \mid - \mid * \mid / \\ \langle \text{id} \rangle & ::= a \mid b \mid c \end{aligned}$$

Terminamos con un ejemplo importante en la historia de los lenguajes de programación:

- Expresiones condicionales:

$$\begin{aligned} \langle \text{enunc} \rangle &::= \langle \text{condicional} \rangle \mid \langle \text{otro} \rangle \\ \langle \text{condicional} \rangle &:= \text{if } \langle \text{expr} \rangle \text{ then } \langle \text{enunc} \rangle \mid \\ &\quad \text{if } \langle \text{expr} \rangle \text{ then } \langle \text{enunc} \rangle \\ &\quad \text{else } \langle \text{enunc} \rangle \end{aligned}$$

- Problema del **else** colgante (*dangling else*):

```
if false then if false then 0 else 1
```

- Dos significados:

```
if false then (if false then 0) else 1
if false then (if false then 0 else 1)
```

- Una gramática equivalente no ambigua es:

$$\begin{aligned} \langle \text{enunc} \rangle &::= \langle \text{condicional} \rangle \mid \langle \text{otro} \rangle \\ \langle \text{condicional} \rangle &:= \langle \text{cond_doble} \rangle \mid \langle \text{cond_final} \rangle \\ \langle \text{cond_doble} \rangle &:= \text{if } \langle \text{expr} \rangle \text{ then } \langle \text{cond_doble} \rangle \text{ else } \langle \text{cond_doble} \rangle \\ \langle \text{cond_final} \rangle &:= \text{if } \langle \text{expr} \rangle \text{ then } \langle \text{enunc} \rangle \mid \\ &\quad \text{if } \langle \text{expr} \rangle \text{ then } \langle \text{cond_simple} \rangle \text{ else } \langle \text{cond_final} \rangle \end{aligned}$$

- Idea: *C1* genera condicionales dobles (if-then-else) balanceados; *C2* representa condicionales simples (if-then) y condicionales dobles pero de forma que un if-then sólo figura colgando al final (en el else).