

Las limitaciones de los modelos finitos de gramáticas para la definición de lenguajes

Esteban di Tada*

Resumen

El uso de reglas de reescritura para el estudio de lenguajes fue introducido por Axel Thue¹ y posteriormente profundizado por Avram Noam Chomsky² que creó el concepto de gramática generativa y estableció las cuatro jerarquías conocidas como *Jerarquías de Chomsky* en base a su capacidad de generación de lenguajes. El problema que planteamos aquí es el de demostrar que el modelo de gramática finita solo puede representar una cantidad muy reducida de lenguajes

Palabras clave: Lenguajes formales, gramáticas formales, Jerarquías de Chomsky

Abstract

The use of rewritten rules in the study of languages was used for the first time by Axel Thue and then formalized by Noam Chomsky who introduced the concept of generative grammars. He defined also the 4 hierarchies known as Chomsky Hierarchies based on the power to generate languages. The problem we will study here is to answer the question if with finite models defined over a given alphabet it is possible to represent all the possible languages that exist on this alphabet

Keywords: Formal Languages, formal grammars, Chomsky Hierarchies

Fecha de Recepción: octubre 2012 | Fecha de aceptación: noviembre 2012

• Universidad de Palermo, Decano de la Facultad de Ingeniería

1. Axel Thue, Matemático noruego nacido en 1863 y fallecido en 1922.

2. Noam Chomsky, lingüista norteamericano perfeccionador de las ideas originales de Thue, nacido en Filadelfia, EEUU en 1928

1. Introducción

1.1. Definiciones previas

- Se definirá como alfabeto a un conjunto finito de símbolos $\Sigma = \{\sigma_1, \sigma_2 \dots \sigma_n\}$
- Se definirá como cadena a una n-tupla ordenada de símbolos pertenecientes a un alfabeto y se lo designará por medio de letras minúsculas griegas. Si por ejemplo $\Sigma = \{a, b, 1, 4\}$, $\delta = ab4a$ será una cadena definida sobre el alfabeto Σ
- Se definirá la operación concatenación de dos cadena α, β a una nueva cadena γ que se forma agregando a continuación de los símbolos de α los símbolos de β . Si $\alpha = cd5r$ y $\beta = r5a$ entonces será $\gamma = cd5rr5a$. Se empleara la notación de la potenciación para representar la operación de multiplicar una cadena n veces por sí misma. Así $a^5 = aaaaa$
- Se definirá como longitud de una cadena a cantidad de símbolos por la que está compuesta. La longitud de la cadena δ del ejemplo del párrafo anterior será 4 y se la denotará $|\delta| = 4$
- Se definirá como Σ^n al conjunto de cadenas φ tal que $\varphi \in \Sigma^n \Leftrightarrow |\varphi| = n$. El conjunto Σ^0 contendrá un solo elemento que será la cadena λ que no contiene ningún símbolo y por definición será $|\lambda| = 0$.
- El conjunto de todas las cadenas se definirá como $\Sigma^* = \bigcup \Sigma^n$ para $n \geq 0$ y $\Sigma^+ = \bigcup \Sigma^n$ para $n \geq 1$
- Funciones sobre cadenas. Sean los conjuntos A y B y sea $f: A \rightarrow B$. Sea α una cadena de símbolos pertenecientes a A. La función $f: A^+ \rightarrow B^+$ se definirá recursivamente de la siguiente manera para $x \in A$ y $\alpha \in A^*$

$$f(x\alpha) = \begin{cases} f(x) & \text{si } |\alpha| = 0 \\ f(x)f(\alpha) & \text{si } |\alpha| > 0 \end{cases}$$

- Dado el conjunto arbitrario A se definirá el conjunto 2^A al conjunto de todos los conjuntos que se pueden formar con los elementos de A incluyendo al conjunto vacío \emptyset .
- Dado el alfabeto Σ se dirá que el subconjunto $\mathcal{L} \subseteq \Sigma^*$ es un lenguaje definido sobre Σ . La cardinalidad de Σ^* denotada $|\Sigma^*| = \aleph_0$ es decir que el conjunto es equipotente al de los naturales. Por el teorema de Cantor la cardinalidad de 2^{Σ^*} será \aleph_1 es decir la cardinalidad del continuo.

1.2. Interpretación intuitiva

Dado que la cardinalidad del conjunto de lenguajes definidos sobre el alfabeto Σ es \aleph_1 , es decir la cardinalidad del continuo, ningún conjunto de modelos enumerable podrá contemplar la totalidad de los lenguajes posibles sobre Σ ya que si así lo fuera ambos conjuntos serian equipotentes por lo que existiría una función biyectiva h entre ambos. Esta función biyectiva mapearía el conjunto de los naturales en los reales lo que es una contradicción. Por lo tanto hay muchos más lenguajes que modelos enumerables. El objetivo del presente trabajo es demostrar formalmente esto para una clase general de gramáticas.

2. Definición de gramática

Se definirá una gramática como $G = \langle \Sigma, N, P, S \rangle$ en donde

- Σ es un conjunto finito de símbolos terminales $\sigma_k, k=1, \dots, n$
- N es un conjunto de símbolos no terminales $\eta_j, j=1, \dots, m$
- P es una relación llamada de producciones y definida de la siguiente manera: $P \subseteq (\Sigma \cup N)^* N (\Sigma \cup N)^* \times (\Sigma \cup N)^*$. Nótese que la parte de la izquierda de la relación debe contener obligatoriamente al menos un símbolo no terminal en tanto que la parte derecha puede ser la palabra vacía γ . Si $\alpha \in (\Sigma \cup N)^* N (\Sigma \cup N)^*$ y $\beta \in (\Sigma \cup N)^*$ y $\alpha, \beta \in P$ esto se denotará $\alpha \xrightarrow{G} \beta$ que se puede expresar enunciando que en la gramática G la subcadena α puede ser reemplazada por la cadena β . Cuando no pueda haber confusión se puede omitir la referencia a la gramática ($\alpha \rightarrow \beta$). En lenguaje corriente la formalización anterior podría expresarse diciendo que una producción es una regla de reescritura que permite reemplazar una cadena formada por símbolos terminales y no terminales pero que al menos contiene un no terminal por una cadena formada por símbolos terminales y no terminales (que puede ser la palabra vacía λ).
- $S \in N$ es el símbolo inicial.

3. Comportamiento de la Gramática

Se dirá que dadas dos cadenas $\alpha, \beta \in (\Sigma \cup N)^*$ β es derivable de α y que se denotará $\alpha \rightarrow \beta$ si y sólo si

- $\alpha = \gamma \zeta \omega$ donde $\gamma, \omega \in (\Sigma \cup N)^*$ y $\zeta \in (\Sigma \cup N)^* N (\Sigma \cup N)^*$ y
- $\beta = \gamma \eta \omega$ y
- $\zeta \rightarrow \eta \in N$

Esto se denota $\alpha \xRightarrow{G} \beta$. La referencia a la gramática puede omitirse cuando ello no pueda generar una confusión.

Se definirá la relación $\alpha \xrightarrow{n} \beta$ si y solo si existen $\gamma_0, \gamma_1, \dots, \gamma_n$ tal que $\gamma_{j-1} \rightarrow \gamma_j$ para $j=1, \dots, n$.

Se dirá que $\alpha \xrightarrow{*} \beta$ si existe un m finito tal que $\alpha \xrightarrow{m} \beta$

Ahora se está en condiciones de definir el lenguaje generado por una gramática. Dada la gramática

$$G = \langle \Sigma, N, P, S \rangle$$

el lenguaje L_G generado por G se definirá como el conjunto de cadenas α tal que

$$\alpha \in L_G \text{ si y solo si } S \xrightarrow{*} \alpha \text{ y } \alpha \in \Sigma^*$$

Esta definición puede expresarse diciendo que

El lenguaje generado por una gramática está constituido por todas las cadenas de símbolos terminales que se derivan del símbolo inicial.

Como en los casos anteriores se puede obviar la mención de la gramática cuando esto no pueda traer indefiniciones.

Resulta conveniente hacer dos observaciones:

- Primero que la definición no garantiza que exista al menos una cadena perteneciente al lenguaje es decir que $L = \emptyset$ y segundo
- Que siempre de acuerdo con la definición de la palabra nula pertenecerá al lenguaje y que esta podría ser su único elemento, es decir $L = \{\lambda\}$. Para evitar esta posibilidad habría que modificar la definición anterior de lenguaje generado por G de la siguiente manera $\alpha \in L_G$ si y solo si $S \xrightarrow{*} \alpha$ y $\alpha \in \Sigma^+$.

Resulta conveniente ilustrar las definiciones con un ejemplo. Sea la gramática

$$G = \langle \Sigma, N, P, S \rangle \text{ donde}$$

$$\Sigma = \{a, b\}$$

$$N = \{Q\}$$

$$P = \{(Q, aQa), (Q, b)\}, \text{ o en el otro formato } P = \{Q \rightarrow aQa, Q \rightarrow b\}$$

$$S = Q$$

Esta gramática genera todas las cadenas del tipo $a^n b a^n$ para $n \geq 0$

4. Gramáticas isomorfas

Sean las gramática $G = \langle \Sigma, N, P, S \rangle$ y los conjuntos T y U equipotentes con Σ y P respectivamente y las funciones biyectivas $h: \Sigma \leftrightarrow T$ y $g: P \leftrightarrow U$. Si se define la función biyectiva $f: \Sigma \cup T \leftrightarrow N \cup U$

$$f(x) = \begin{cases} x \in \Sigma: h(x) \\ x \in N: g(x) \end{cases}$$

Cuya función inversa es (dado que h y g son biyectivas)

$$f^{-1}(x) = \begin{cases} x \in T: h^{-1}(x) \\ x \in U: g^{-1}(x) \end{cases}$$

Se puede crear una nueva gramática $H = \langle T, U, Q, R \rangle = \langle f(T), f(U), f(Q), f(R) \rangle$
Por razones de simplificación se denotará $H = f(G)$

Hay una función, que será definida como canónica, que jugará un especial rol en las demostraciones que siguen. Sea la gramática $G = \langle \Sigma, N, P, S \rangle$ en donde $\Sigma = \{\sigma_1, \sigma_2, \dots, \sigma_n\}$ y $N = \{\eta_1, \eta_2, \dots, \eta_m\}$ la función canónica será

$$\begin{aligned} c(\sigma_j) &= j \\ c(\eta_k) &= n+k \end{aligned}$$

Y la función inversa será par $x \in (\Sigma \cup N)$

$$c^{-1}(i) = \begin{cases} \sigma_i & \text{si } 0 < i \leq n \\ \eta_{i-n} & \text{si } n < i \leq n+m \end{cases}$$

De esta manera los símbolos terminales y no terminales son reemplazados por números naturales lo que facilitará la aplicación de métodos similares a los empleados en la demostración del teorema de Göedel.

Sean la gramática $G = \langle \Sigma, N, P, S \rangle$ dos cadenas $\alpha, \beta \in (\Sigma \cup N)^*$ tal que $\alpha \rightarrow \beta$ luego se verificará por definición lo siguiente:

- $\alpha = \gamma \zeta \omega$ donde $\gamma, \omega \in (\Sigma \cup N)^*$ y $\zeta \in (\Sigma \cup N)^* N (\Sigma \cup N)^*$ y
- $\beta = \gamma \eta \omega$ y
- $\zeta \rightarrow \eta \in N$

Por lo tanto se verificará que $c(\alpha) \rightarrow c(\beta)$ dado que:

- $c(\alpha) = c(\gamma)c(\zeta)c(\omega)$ donde $c(\gamma), (\omega) \in c((\Sigma \cup N)^*)$ y $c(\zeta) \in c((\Sigma \cup N)^* N (\Sigma \cup N)^*)$
y
- $c(\beta) = c(\gamma)c(\eta)c(\omega)$ y
- $c(\zeta) \rightarrow c(\eta) \in c(N)$

Dado que c es una función biyectiva se puede demostrar entonces que $x \in L_G$ si y solo si $c(x) \in c(L_G)$.

Se definirá la gramática $c(L_G)$ como la forma canónica de la gramática G

En otras palabras los símbolos terminales y no terminales fueron reemplazados por números enteros mediante la función c y las cadenas generadas por esta nueva gramática (que serán cadenas de naturales) restituyen la cadena de G mediante la aplicación de c^{-1} . La cardinalidad de L_G y $L_{c(G)}$ es la misma.

5. Teorema de Schröder-Bernstein-Cantor

Un teorema fundamental para la teoría de los números transfinitos de Cantor es el demostrado por Schröder, Bernstein y Cantor que establece que dados dos conjuntos A y B y dos funciones sobreyectivas $f: A \rightarrow B$ y $g: B \rightarrow A$ entonces existe una función biyectiva $h: A \leftrightarrow B$. El teorema resulta bastante evidente para conjuntos finitos pero no lo es así para conjuntos no finitos.

No se incluirá aquí la demostración de este teorema pero para una mejor comprensión se presentará un ejemplo. Se desea demostrar que existe una función biyectiva $h: [0, 1] \leftrightarrow [0, 1)$ – entre el intervalo cerrado $0 \leq x \leq 1$ y el intervalo semiabierto $0 \leq x < 1$. Es fácil construir las funciones sobreyectivas f y g de la siguiente manera:

$$f(x) = x/2: [0, 1] \rightarrow [0, 1) \text{ y } g(x) = x: [0, 1) \rightarrow [0, 1]$$

La construcción de la función biyectiva se hará en forma iterativa. Si se elige la función $f_0(x) = x$ se ve que no satisface las condiciones requeridas para ser una función ya que $f_0(1) = 1$ no pertenece a $[0, 1)$ que, por definición, no contiene el punto 1. Para resolver este problema se modifica la función de la siguiente manera:

$$f_1(x) = \left\{ \begin{array}{l} x \iff x \neq 1 \\ 1/2 \iff x = 1 \end{array} \right\}$$

Con esta solución queda resuelto el problema anterior pero esta función presenta ahora el problema que tanto para $x=1$ como para $x=1/2$ la función vale $1/2$. Para solucionar este problema se modifica la definición de la siguiente manera

$$f_2(x) = \left\{ \begin{array}{l} x \Leftrightarrow x \notin \{1, 1/2\} \\ 1/2 \Leftrightarrow x \in \{1\} \\ 1/4 \Leftrightarrow x \in \{1/2\} \end{array} \right\} = \left\{ \begin{array}{l} x \Leftrightarrow x \notin \{1, 1/2\} \\ x/2 \Leftrightarrow x \in \{1, 1/2\} \end{array} \right\}$$

El problema ahora se trasladó al punto $x = 1/4$. Para resolverlo defínase la nueva función

$$f_3(x) = \left\{ \begin{array}{l} x \Leftrightarrow x \notin \{1, 1/2, 1/4\} \\ x/2 \Leftrightarrow x \in \{1, 1/2, 1/4\} \end{array} \right\}$$

El problema surge ahora en el punto $1/8$ que corresponde a dos valores de x que son $x = 1/4$ y $x = 1/8$. Continuando con este proceso, puede comprobarse que la función biyectiva será

$$f(x) = \lim_{n \rightarrow \infty} f_n(x) = \left\{ \begin{array}{l} x \Leftrightarrow x \notin \{1, 1/2, 1/4, 1/8, \dots\} \\ x/2 \Leftrightarrow x \in \{1, 1/2, 1/4, 1/8, \dots\} \end{array} \right\}$$

Si se define el conjunto $X_n = \{1, 1/2, 1/4, 1/8, \dots, 1/n\}$ y

$$X = \lim_{n \rightarrow \infty} X_n = \lim_{n \rightarrow \infty} \{1, 1/2, 1/4, 1/8, \dots, 1/n\}$$

La función puede redefinirse de la siguiente manera

$$f(x) = \left\{ \begin{array}{l} x \Leftrightarrow x \notin X \\ x/2 \Leftrightarrow x \in X \end{array} \right\}$$

Este teorema permite demostrar fácilmente que existe una función biyectiva entre los pares ordenados $(x,y) \in \mathbb{N}^2$ y $x \in \mathbb{N}$. Demostraremos que existe una función

sobreyectiva $f: \mathbb{N}^2 \rightarrow \mathbb{N}$. La función es sencilla $f(x,y)=2^x 3^y$. Esta función está bien definida porque dado que el teorema fundamental de la aritmética garantiza que la descomposición en los factores primos de un natural es única salvo el orden de los factores. La función g se puede definir $g(x)=(x,1): \mathbb{N} \rightarrow \mathbb{N}^2$. Por lo tanto existe una función biyectiva entre $h: \mathbb{N} \leftrightarrow \mathbb{N}^2$.

El conjunto de gramáticas $G=\langle \Sigma, N, P, S \rangle$ tiene la cardinalidad \aleph_0

Se demostrara que existe una función biyectiva entre el conjunto Ψ de todas las gramáticas de la forma $G=\langle \Sigma, N, P, S \rangle$ y el de los números naturales \mathbb{N} donde G está en la forma canónica.

Para ello se empleará el teorema de Schröder-Bernstein-Cantor. El primer paso será demostrar que existe una función sobreyectiva $f(G): \Psi \rightarrow \mathbb{N}$ que mapea cada gramática G en un numero natural

Sea la función f definida de la siguiente manera:

$$f(G)=2^{|\Sigma|} 3^{|N|} 5^q$$

- Donde $|\Sigma|$ es la cantidad de símbolos del alfabeto o cantidad de terminales y
- $|N|$ es la cantidad de símbolos no terminales y
- q es un producto de potencias de números primos $q=p_1 r^1 p_2 r^2 \dots p_{|N|} r^{|N|}$ donde
 - p_i es el i -ésimo número primo y
 - la producción $(\gamma_j, \varphi_j) \in P$ es la j -ésima producción de la gramática y
 - $r_j = 2^{a_j} 3^{b_j}$ donde $a_j = \prod_{i=1}^{|\gamma_j|} p_i^{r_i^j}$ y $b_j = \prod_{i=1}^{|\sigma_j|} p_i^{\sigma_i^j}$ siendo $|\gamma_j|$ y $|\sigma_j|$ la longitud de la parte de la izquierda y de la parte de la derecha de la j -ésima producción respectivamente.

Se demostrará que la función $f(G)$ está bien definida. Dado $m=f(G)$ es posible reconstruir la gramática canónica. La demostración de basa en el teorema fundamental de la aritmética por el cual la descomposición de un número natural en sus factores primos es única. En efecto dado m de su descomposición en sus factores primos se pueden determinar $|\Sigma|$, $|N|$ y q . Continuando con un procedimiento análogo dado q se pueden determinar los $r_j, j=1, \dots, |N|$. Dados los valores de r_j y dado que $r_j = 2^{a_j} 3^{b_j}$ se pueden determinar los valores de a_j y b_j .

Finalmente dado que $a_j = \prod_{i=1}^{|Y_j|} p_i^{\gamma_j^i}$ y $b_j = \prod_{i=1}^{|\sigma_j|} p_i^{\sigma_j^i}$ se pueden determinar los valores de γ_j^i y σ_j^i con lo que se reconstruye el conjunto de producciones.

La demostración que existe una función sobreyectiva g que va de los naturales al conjunto de gramáticas es sencillo ya que dado un numero natural n cualquiera se le puede asociar la gramática $G = \langle \Sigma, N, P, S \rangle$ donde

$$\Sigma = \{a\}, N = \{T\}, P = \{(R, a^n)\} \quad S = T.$$

El lenguaje generado por esta gramática contiene una sola cadena que consiste en n símbolos a seguidos.

Por lo tanto existen las funciones sobreyectivas $f: \Psi \rightarrow \mathbb{N}$ y $g: \mathbb{N} \rightarrow \Psi$ y por lo tanto debe existir de acuerdo con el teorema de Schröder-Bernstein-Cantor una función biyectiva $h: \Psi \leftrightarrow \mathbb{N}$.

Con esto queda demostrado que la cardinalidad de Ψ es \aleph_0 .

Sea nuevamente la gramática $G = \langle \Sigma, N, P, S \rangle$ donde

$$\Sigma = \{a, b\}$$

$$N = \{Q\}$$

$$P = \{(Q, aQa), (Q, b)\}, \text{ o en el otro formato } Q \rightarrow aAa, Q \rightarrow b$$

$$S = Q$$

La gramática canónica correspondiente sería $G' = \langle \Sigma', N', P', S' \rangle$ donde

$$\Sigma' = \{1, 2\}$$

$$N' = \{3\}$$

$$P' = \{(3, 131), (3, 2)\}$$

$$S' = 3$$

$$a_1 = 2^3 = 8, \quad b_1 = 2^1 3^3 5^1 = 270$$

$$a_2 = 2^3 = 8, \quad b_2 = 2^2 = 4$$

Por lo que será

$$r1 = 2^8 3^{270}$$

$$r2 = 2^8 3^4$$

El valor de q resulta

$$q = 2^{2^{83}270} 3^{2^{83}4}$$

Y finalmente

$$fG = 2^2 3^{15} 5^{2^{2^{83}270} 3^{2^{83}4}}$$

Este número, para esta gramática tan sencilla, es enormemente grande y casi imposible de poder calcular. Pero lo importante es que existe.

7. Conclusiones

De lo demostrado previamente se concluye que la cardinalidad del conjunto de gramáticas $G = \langle \Sigma, N, P, S \rangle$ es \aleph_0 en tanto que la cardinalidad del conjunto de gramáticas que existen sobre Σ^* es la de 2^{Σ^*} que por el teorema de Cantor tiene la cardinalidad \aleph_1 del continuo. Por lo tanto existen infinitas gramáticas que no pueden ser representadas con dicho modelo. Ello explica la gran dificultad que existe en la representación de lenguajes naturales por medio de modelos finitos.

Bibliografía

- [1] K. R. Chowdhari: “Fundamental of Discrete Mathematical”, Second Edition, Structures PHI Learning Private Limited, January 2012
- [2] Philip Edward Bertrand Jourdain, “Selected essays on the History of Set Theory and Logics”, CLUEB 1991
- [3] Joseph Warren Dauben, “Georg Cantor His Mathematics and Philosophy of the infinity, Princeton University Press 1979
- [4] John E. Hopcroft, “Introduction To Automata Theory, Languages, And Computation”, Pearson Education, 2008

[5] Rebecca Goldstein, “Gödel Paradoja y vida” , Antoni Bosh Editor España, 2005

[6] Ernest Nagel, James R. Newman, “Gödel’s Proof”, New York University Press, 2001

[7] Georg Cantor, “Contribution to the Founding of the theory of transfinite Numbers, Cosimo Inc. 2007

