

Teorema Chomsky-Schützenberger

Autómatas y lenguajes formales

17 de septiembre de 2021

Lenguaje de *Dyck*

Este tema en el plan del curso iba antes, pero lo salté, aquí lo recupero pues es necesario para lo que sigue.

Podemos definir un lenguaje para los paréntesis bien balanceados de distinto tipo¹, que llamaremos D_n de la siguiente forma:

$$S \rightarrow [S]_1|[S]_2|[S]_n|SS| \epsilon,$$

con esas producciones se pueden generar n pares de paréntesis bien anidados.

Un concepto importante en matemáticas y en las ciencias de la computación es la conversión de un problema a otro ya conocido², así si en el problema nuevo no encontramos solución, si lo podemos convertir a un viejo conocido con solución, pues ya solo convirtiendo algunas cosas tenemos la solución del nuevo. Esta es la base del teorema de Chomsky-Schützenberger.

Teorema de Chomsky-Schützenberger

Teorema 1 *Todo lenguaje independiente de contexto es una imagen homomorfa de la intersección de un lenguaje de paréntesis y un conjunto regular. En otras palabras, para todo lenguaje independiente de contexto A , hay una $n \geq 0$, un conjunto regular R y un homomorfismo h tales que:*

$$A = h(D_n \cap R)$$

Diciéndolo en otras palabras, todo lenguaje independiente de contexto es una ampliación a un lenguaje regular, ya que es equivalente a la unión de un lenguaje regular con un lenguaje de paréntesis de *Dyck* (que es un ejemplo muy sencillo de un lenguaje independiente de contexto).

El homomorfismo es un mapeo que va de un dominio a otro, en este caso $h : \Gamma^* \rightarrow \Lambda^*$, por ejemplo la función $h(x, y) = h(x)h(y)$, donde $x, y \in \Gamma^*$ y $h(x)h(y) \in \Lambda^*$. Como el nombre lo indica, un mapeo pasa punto por punto a

¹Como ya su experiencia en la matemáticas les habrá hecho ver, cuando usan más de un paréntesis, en ocasiones es útil usar paréntesis diferenciados, por ejemplo si en una misma expresión requieren tres tipos de paréntesis pueden usar $\{\{()\}$. Aunque en programación se suele usar sólo un tipo y la programadora/programador o el editor está pendiente de donde abre y cierra cada tipo.

²Incluso hay un chiste sobre cuántos matemáticos se necesitan para cambiar un foco, no me lo sé y para no arruinarlo no lo cuento.

algo muy parecido, como aquellos mapamundis que usaban en la escuela, que pasan el globo terráqueo a una superficie plana, sin dejar fuera ningún país. Una imagen para recordar un poco como es eso se muestra en la figura 1.

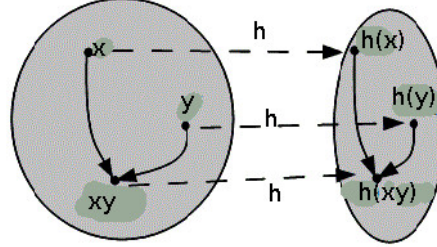


Figura 1: Ejemplo de un homomorfismo, piense que el conjunto de la izquierda es Γ^* y el de la derecha Σ^* . Imagen adaptada de una tomada de la página: <https://mathstrek.blog/2012/09/28/casual-introduction-to-group-theory-6/>

¿Cómo se hace este homomorfismo de la intersección de un lenguaje de *Dyck* y un lenguaje regular a un lenguaje independiente de contexto? Si recuerdan un lenguaje independiente de contexto se define por un cuarteto del tipo $G = (\Sigma, \Gamma, S, \rightarrow)$, con Σ el alfabeto de entrada, Γ el alfabeto de la memoria, S la variable inicial y \rightarrow el conjunto de producciones.

Vamos a desglosar un poco este último conjunto, las producciones son varias, en este caso las llamaremos a partir de letras griegas minúsculas. Sólo para identificarlas con las funciones en el nuevo dominio.

$$\{\pi, \rho, \sigma, \dots\} \in \rightarrow,$$

esta elección de llamar al conjunto de producciones \rightarrow al momento puede ser medio confuso, pero espero conforme avancemos no lo sea tanto. Otra manera de verlo es como si cada producción tuviera un nombre:

$$\{\rightarrow_{\pi}, \rightarrow_{\rho}, \rightarrow_{\sigma}, \dots\} \in \rightarrow,$$

pero sólo las identificaremos por las letras que las designan.

Lo que sigue es transformar las producciones al lenguaje de paréntesis, por ejemplo si se tiene el lenguaje independiente de contexto:

$$\begin{aligned} A &\rightarrow BC \\ A &\rightarrow a, \end{aligned}$$

Las podemos juntar en una sola producción:

$$A \rightarrow BC|a$$

como pueden ver ya en forma normal de Chomsky, a esta le llamamos nuestra producción π . En el lenguaje de *Dyck* se transforma a:

$$A \rightarrow \begin{matrix} 1 & 12 & 2 \\ [B] & [C] \\ \pi & \pi\pi & \pi \end{matrix}$$

$$A \rightarrow \begin{matrix} 1122 \\ [][] \\ \pi\pi\pi\pi \end{matrix},$$

respectivamente. Los números de arriba indican el tipo de paréntesis (un paréntesis del tipo 1 no cierra con uno del tipo 2) y la letra griega minúscula inferior hace referencia a la producción, en este caso π .

Para hacerlo más claro completemos el lenguaje:

$$\begin{aligned} A &\rightarrow BC \\ A &\rightarrow a \\ B &\rightarrow b \\ C &\rightarrow c, \end{aligned}$$

aprovechando que ya está en forma normal de Chomsky. Cada regla de producción que involucre a la misma variable del lado izquierdo se marca con el mismo símbolo, así la producción π designa a las primeras dos líneas (con la variable A), la producción de B la llamaremos ρ , y la de C será σ . Por suerte estas últimas dos solo tienen una regla.

Ponerlo en forma normal de Chomsky lo reduce a que tengamos sólo dos tipos distintos de paréntesis por cada producción, sólo habrá dos tipos para π , dos tipos para ρ y dos para σ , pero un paréntesis de π no puede ser cerrado con uno de σ o ρ , son reglas de producción distintas. El lenguaje se transforma en:

$$A \rightarrow \begin{matrix} 1 & 12 & 2 \\ [B] & [C] \\ \pi & \pi\pi & \pi \end{matrix}$$

$$A \rightarrow \begin{matrix} 1122 \\ [][] \\ \pi\pi\pi\pi \end{matrix}$$

$$B \rightarrow \begin{matrix} 1122 \\ [][] \\ \rho\rho\rho\rho \end{matrix}$$

$$C \rightarrow \begin{matrix} 1122 \\ [][] \\ \sigma\sigma\sigma\sigma \end{matrix}.$$

De aquí se pueden ver unas reglas para el acomodo de los paréntesis y de que forma pueden ser identificados unos con variables y otros con terminales, eso ya se lo dejo un poco a su análisis, pero éstas son las reglas. Por ejemplo las cadenas en ese lenguaje independiente de contexto son a y bc , en el nuevo lenguaje se transforman en $\begin{matrix} 1122 \\ [][] \\ \pi\pi\pi\pi \end{matrix}$ y $\begin{matrix} 111221211222 \\ [][] [] [] [] \\ \pi\rho\rho\rho\rho\pi\sigma\sigma\sigma\sigma\pi \end{matrix}$ respectivamente.

Este lenguaje de puros paréntesis es equivalente al anterior, aunque no lo parezca, esos monstruos de paréntesis ya son un lenguaje independiente de contexto, todos los símbolos son el alfabeto del lenguaje de *Dyck*, pero las reglas para que hagan sentido es la intersección con un lenguaje regular (esta construcción podría representarse por un autómata finito).

Referencias

- [1] Kozen, Dexter C. “Automata and Computability” Springer (1997)