

Autómatas y Lenguajes Formales 2019-2

Facultad de Ciencias UNAM*

Nota de Clase 7: El lema del bombeo - Teorema de Myhill Nerode

Favio E. Miranda Perea

A. Liliana Reyes Cabello

Lourdes González Huesca

29 de marzo de 2019

1. Lenguajes regulares o ¿no?

Los lenguajes que hemos tratado han sido caracterizados por ser sencillos, es decir son conjuntos de cadenas cuyas propiedades son fáciles de abstraer mediante lenguajes básicos y operaciones como la concatenación, unión y cerradura de Kleene. A partir de ellos, las propiedades de cerradura nos permiten construir nuevos lenguajes regulares: Si L, M son lenguajes regulares entonces:

$$\begin{array}{ll} L \cup M \text{ es regular} & LM \text{ es regular} \\ L^* \text{ es regular} & L^+ \text{ es regular} \\ \overline{L} \text{ es regular} & L \cap M \text{ es regular} \\ L - M \text{ es regular} & \end{array}$$

Así mismo, hemos considerado la equivalencia entre lenguajes regulares y los autómatas finitos a través del teorema de Kleene. Pero ¿cómo es posible decidir si un lenguaje dado, sin una representación de expresión regular o autómatata que lo reconoce, es regular?

Un caso conocido de un lenguaje “sencillo” no regular es

$$L = \{a^i b^i \mid i \in \mathbb{N}\}$$

Veamos dos ejemplos de lenguajes no regulares al tratar de dar una expresión que los genere:

Ejemplo: Considere el lenguaje $L = \{a^i b^j \mid i \neq j, i, j \in \mathbb{N}\}$. Decida si es regular.

Para la demostración, supongamos primero que L es regular. Ahora procedemos a construir una expresión que lo represente: partimos del lenguaje $a^* b^*$ que claramente es regular.

Por propiedades de cerradura, el lenguaje $a^* b^* - L$ también debe ser regular.

¡Pero $a^* b^* - L = \{a^i b^i \mid i \in \mathbb{N}\}$ no es regular!

Por lo tanto L no puede ser regular.

*Material elaborado en el marco del proyecto PAPIIME PE102117

Ejemplo: Ahora considere el lenguaje $L = \{wb^n \mid |w| = n, n \geq 1\}$. Decida si es regular. Nuevamente, supongamos que L es regular. Y partimos otra vez de a^*b^* que claramente es regular. Usando las propiedades de cerradura $L \cap a^*b^*$ también debe ser regular. Pero $L \cap a^*b^* = \{a^ib^i \mid i \in \mathbb{N}\}$ no es regular. Por lo tanto L no puede ser regular.

1.1. ¿Cuántos lenguajes regulares hay?

Los lenguajes son variados y requerimos de una forma eficiente para decidir si un lenguaje es regular y entonces atrevemos a proporcionar una máquina que lo reconozca.

Consideremos el conjunto de todos los lenguajes regulares, dado un alfabeto:

$$REG = \{L \subseteq \Sigma^* \mid L \text{ es regular} \}$$

¿Cuál es la cardinalidad de REG ? Analicemos el conjunto:

- Dado que cualquier lenguaje L es un subconjunto de Σ^* , existen tantos lenguajes como elementos en $\mathcal{P}(\Sigma^*)$.
- Puesto que Σ^* es infinito numerable, es decir es del tamaño del conjunto \mathbb{N} de los números naturales, entonces $\mathcal{P}(\Sigma^*)$ es del tamaño del conjunto de los números reales \mathbb{R} .
- Existen sólo tantos lenguajes regulares como números naturales, $|REG| = |\mathbb{N}|$:
 - La idea de la prueba es enumerar lexicográficamente *todos* los AFD posibles con alfabeto de entrada Σ , es decir, primero los autómatas con un sólo estado, luego los de dos estados, etc.
 - Esto implica que el número de lenguajes regulares es a lo más numerable.
 - Además, claramente es numerable pues hay una infinidad numerable de lenguajes regulares, por ejemplo

$$\{a\}, \{aa\}, \{aaa\}, \dots$$

- De manera que el conjunto $\mathcal{P}(\Sigma^*) - REG$ no puede ser numerable, pues la unión de numerables sigue siendo numerable.
- Es decir, hay tantos lenguajes **no** regulares como números reales.

En esta nota nos dedicaremos a discutir dos métodos para probar que un lenguaje **no** es regular.

2. Lema de Bombeo para lenguajes regulares

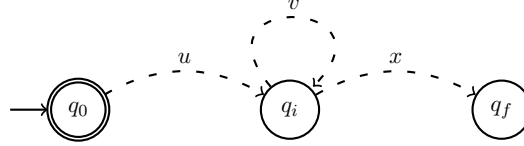
Lema 1 (Lema del Bombeo). *Si L es un lenguaje regular infinito entonces existe un número $n \in \mathbb{N}$, llamado constante de bombeo para L , tal que para cualquier cadena $w \in L$ con $|w| \geq n$ existen cadenas u, v, x tales que:*

1. $w = uvx$
2. $|uv| \leq n$

3. $v \neq \varepsilon$

4. $\forall m \in \mathbb{N} (uv^m x \in L)$.

Informalmente, la idea del lema se basa en que la característica de un lenguaje regular consiste en *repetir* o *ciclar* una subcadena y de ahí que exista un ciclo en un autómata que reconoce a L . Recordemos que un autómata es una máquina con una memoria corta, es decir sólo recuerda lo que está abstraído en un estado y no tiene memoria auxiliar. El siguiente diagrama resume el lema:



Veamos a continuación la prueba formal.

Demostración. Como L es regular, existe un AFD $M = \langle Q, \Sigma, q_0, \delta, F \rangle$ que reconoce a L . Sea $n = |Q|$, es decir n es el número de estados en Q . Sea $w \in L$ una cadena con $|w| \geq n$, digamos $w = a_1 \dots a_n z$ con $a_i \in \Sigma$ y $z \in \Sigma^*$. Ya que $w \in L$, entonces se tiene que $\delta^*(q_0, w) = q_f \in F$. Supongamos s.p.g. que para toda $0 \leq i < n$, $\delta(q_i, a_{i+1}) = q_{i+1}$. Entonces el proceso de cómputo para $\delta^*(q_0, w) = q_f$ es como sigue:

$$\begin{aligned}
 \delta^*(q_0, w) &= \delta^*(\mathbf{q_0}, a_1 \dots a_n z) \\
 &= \delta^*(\delta^*(q_0, a_1), a_2 \dots a_n z) \\
 &= \delta^*(\mathbf{q_1}, a_2 \dots a_n z) \\
 &= \delta^*(\delta^*(q_1, a_2), a_3 \dots a_n z) \\
 &= \delta^*(\mathbf{q_2}, a_3 \dots a_n z) \\
 &\vdots \\
 &= \delta^*(\delta^*(q_{n-2}, a_{n-1}), a_n z) \\
 &= \delta^*(\mathbf{q_{n-1}}, a_n z) \\
 &= \delta^*(\delta^*(q_{n-1}, a_n), z) \\
 &= \delta^*(\mathbf{q_n}, z) \\
 &= q_f \in F
 \end{aligned}$$

Obsérvese que en este proceso de ejecución de w figuran los estados q_0, \dots, q_n que son $n + 1$ estados y como Q cuenta sólo con n estados, por el principio del palomar, existen $j, k \in \{0, 1, \dots, n\}$

tales que $j < k$ y $q_j = q_k$. De aquí que el proceso anterior equivale a

$$\begin{aligned}
\delta^*(q_0, w) &= \delta^*(q_0, a_1 \dots a_j \dots a_k \dots a_n z) \\
&= \delta^*(\delta^*(q_0, a_1 a_2 \dots a_j), a_{j+1} \dots a_k \dots a_n z) \\
&= \delta^*(\mathbf{q}_j, a_{j+1} \dots a_k \dots a_n z) \\
&= \delta^*(\mathbf{q}_k, a_{k+1} \dots a_n z) \\
&= \delta^*(\mathbf{q}_j, a_{k+1} \dots a_n z) \\
&= \vdots \\
&= \delta^*(\delta^*(q_{n-1}, a_n), z) \\
&= \delta^*(q_n, z) \\
&= q_f \in F
\end{aligned}$$

Sean $u = a_1 \dots a_j$, $v = a_{j+1} \dots a_k$ y $x = a_{k+1} \dots a_n z$. Es claro que $w = uvx$ y que $|uv| = |a_1 \dots a_j a_{j+1} \dots a_k| \leq n$, de donde las condiciones 1 y 2 se cumplen. Por otra parte, se cumple que $v \neq \varepsilon$ (condición 3), pues $j < k$ y $|v| = |a_{j+1} \dots a_k| = k - j > 0$.

Finalmente, observemos que

$$\begin{aligned}
\delta^*(q_j, v) &= \delta^*(q_j, a_{j+1} \dots a_k) \\
&= \delta^*(q_{j+1}, a_{j+2} \dots a_k) \\
&= \vdots \\
&= \delta^*(q_{k-1}, a_k) \\
&= \delta^*(\mathbf{q}_k, \varepsilon) \\
&= \delta^*(\mathbf{q}_j, \varepsilon) = q_j
\end{aligned}$$

de donde, por inducción sobre m , se sigue que $\forall m \in \mathbb{N} (\delta^*(q_j, v^m) = q_j)$. La propiedad 4 resulta ahora inmediata puesto que:

$$\delta^*(q_0, uv^m x) = \delta^*(\delta^*(q_0, u), v^m x) = \delta^*(q_j, v^m x) = \delta^*(\delta^*(q_j, v^m), x) = \delta^*(q_j, x) = \delta^*(q_n, z) = q_f \in F$$

□

Obsérvese que de la prueba del lema del bombeo se sigue que la constante de bombeo no es única, pues hay una por cada AFD que reconozca al lenguaje en cuestión. Más aún, la constante del bombeo es meramente existencial pues no se conoce cómo es el autómata particular que la define, de modo que no es posible saber cuál es el número de estados de dicho autómata.

Otra observación importante es que el lema del bombeo sólo proporciona una condición necesaria para que un lenguaje sea regular, pero no una condición suficiente, es decir, existen lenguajes que satisfacen el lema del bombeo pero no son regulares. El siguiente es un ejemplo cuya verificación se deja al lector.

Ejemplo 1. El lenguaje $L \subseteq \{a, b, c\}^*$ definido como:

$$L = \{a^i b^j c^j \mid i \geq 1, j \geq 0\} \cup \{b^j c^k \mid j, k \geq 0\}$$

tiene una constante de bombeo, a saber $n = 1$. Sin embargo L no es un lenguaje regular.

De manera que el lema del bombeo no sirve para mostrar que un lenguaje es regular. Su principal aplicación es Para probar que un lenguaje L **no** es regular, para lo cual se procede por *contradicción* como sigue:

1. Si L fuera regular entonces existiría una constante de bombeo n .
2. Y cualquier palabra $w \in L$ con longitud mayor o igual a n se descompone como $w = uvx$ donde $v \neq \varepsilon$, $|uv| \leq n$.
3. Se llega a una contradicción como sigue:

por el lema del bombeo la cadena $uv^i x$ debe pertenecer a L , para toda $i \geq 0$. Pero por la definición particular de L , se puede mostrar alguna i tal que $uv^i x \notin L$.

Debemos observar que encontrar la i adecuada depende del problema particular y no hay un método general, pero usualmente basta con valores pequeños de i , incluyendo $i = 0$.

Ejemplo 2. Veamos que $L = \{a^i b^i \mid i \in \mathbb{N}\}$ **no** es regular usando el Lema del bombeo.

Supóngase que L es regular y sea n una constante de bombeo. Ahora considere una palabra cualquiera del lenguaje, es decir tiene la forma $w = a^n b^n$. Una descomposición de w en uvx es la siguiente: u, v contienen sólo a 's, digamos

$$u = a^k, v = a^\ell, k \geq 0, \ell \geq 1$$

Con esto aseguramos que se cumple que $v \neq \varepsilon$, $|uv| \leq n$ y además $x = a^{n-k-\ell} b^n$.

Si tomamos $i = 2$, por el lema del bombeo la cadena $uv^2 x$ debe pertenecer a L . Pero, realmente sucede que

$$uv^2 x = a^k a^\ell a^\ell a^{n-k-\ell} b^n = a^{n+\ell} b^n \notin L$$

Por lo tanto, L no es regular.

Ejemplo 3. Demuestre que $L_1 = \{w \in \{a, b\}^* \mid w = w^R\}$ **no** es regular.

Suponer que L es regular y sea n una constante de bombeo. Ahora considere que la palabra $w = a^n b^n a^n$ en L_1 tiene una descomposición en $w = uvx$ con $v \neq \varepsilon$, $|uv| \leq n$ y en donde u, v tienen sólo a 's digamos $u = a^k$, $v = a^\ell$, $\ell \geq 1$.

Por tanto $x = a^{n-k-\ell} b^n a^n$. Al hacer $i = 2$ debemos tener a $uv^2 x \in L$, por el lema del bombeo. Pero por otra parte se tiene que:

$$uv^2 x = a^k a^\ell a^\ell a^{n-k-\ell} b^n a^n = a^{n+\ell} b^n a^n \notin L$$

Y por lo tanto L no es regular.

Así el lema anterior *no* permite demostrar que un lenguaje es regular, generalmente se usa como método de demostración para probar que un lenguaje **no** es regular. Veamos a continuación otro método, el cual permite decidir si un lenguaje es regular.

3. El Teorema de Myhill-Nerode

Considérense las siguientes relaciones de equivalencia sobre Σ^* relacionadas a un lenguaje dado L y a un autómata finito determinista dado M , sean u, v dos cadenas:

- $u \equiv_L v$ si y sólo si

$$\forall w \in \Sigma^* (uw \in L \Leftrightarrow vw \in L)$$

Se dice que u, v son cadenas *indistinguibles* para L .

- $u \equiv_M v$ si y sólo si

$$\delta^*(q_0, u) = \delta^*(q_0, v)$$

Es decir, u, v son cadenas *indistinguibles* según M .

Ejemplo 4. Dado el lenguaje $L = \{a^i b^i \mid i \in \mathbb{N}\}$ sobre $\Sigma = \{a, b\}$ se tiene que

- $a^4 b^3 \equiv_L a^3 b^2$ pues

$$\forall z \in \Sigma^* (a^4 b^3 z \in L \Leftrightarrow a^3 b^2 z \in L)$$

En particular con $z = b$, ambas concatenaciones están en L . En otro caso no lo están pero esto también hace que se cumpla la equivalencia.

- $a^2 b^2 \not\equiv_L a^3 b^2$ pues para $z = \varepsilon$ se tiene

$$a^2 b^2 z \in L \text{ y } a^3 b^2 z \notin L$$

- $a^4 b \not\equiv_L a^3 b^2$ pues para $z = b$, se tiene

$$a^4 b^2 z \notin L \text{ y } a^3 b^2 z \in L$$

- $abb \equiv_L baba$ pues ambas cadenas no pertenecen a L y, debido a la definición de L , $abbb$ y $babav$ tampoco están en L para cualquier z .

- La relación \equiv_L tiene una infinidad de clases de equivalencia, por ejemplo:

$$[\varepsilon], [a], [a^2], \dots, [a^n], \dots$$

Todas estas clases son diferentes pues si $i \neq j$ entonces $a^i \not\equiv_L a^j$ ya que si consideramos $z = b^i$ entonces $a^i z \in L$ pero $a^j z \notin L$.

Por lo general no hay relación alguna entre un lenguaje L y un autómata M . Más aún, la relación \equiv_L puede definirse para cualquier lenguaje L aún cuando este no sea regular. Sin embargo, en el caso particular en que $L = L(M)$ se cumple que \equiv_M es un refinamiento de \equiv_L .

Proposición 1. Para cualesquiera $x, y \in \Sigma^*$, Si $x \equiv_M y$ entonces $x \equiv_{L(M)} y$.

Demostración. Ejercicio

□

□

Esta proposición nos deja ver la más importante limitación de los autómatas finitos, el hecho de que carecen de memoria más allá de lo que recuerde el estado actual: si $x \equiv_M y$ entonces $x \equiv_{L(M)} y$, por lo que ninguna cadena w procesada después de x o y permitirá que M determine cuál de x o y se procesó anteriormente.

Definición 1. Una relación de equivalencia \equiv sobre Σ^* es invariante por la derecha si y sólo si

$$\forall x, y, w \in \Sigma^* (x \equiv y \rightarrow xw \equiv yw).$$

Así la relación \equiv_L es invariante por la derecha.

Lema 2 (Lema de continuación). Sean $x, y \in \Sigma^*$. Si $\delta^*(q_0, x) = \delta^*(q_0, y)$ entonces para cualquier $z \in \Sigma^*$, se cumple que

$$\delta^*(q_0, xz) = \delta^*(q_0, yz)$$

Demostración. Ejercicio. □

Del lema anterior se sigue que la relación \equiv_M es invariante por la derecha.

Proposición 2. Dado un AFD $M = \langle Q, \Sigma, q_0, \delta, F \rangle$ se cumple lo siguiente:

- La relación \equiv_M es invariante por la derecha.
- La relación \equiv_M es de índice¹ finito.
- $L(M)$ es la unión de algunas de las clases de equivalencia de la relación \equiv_M .

Demostración.

Que \equiv_M es invariante por la derecha es consecuencia del lema de continuación 2.

- Que \equiv_M es de índice finito es claro pues cada clase de equivalencia se define como:

$$[x]_{\equiv_M} = \{y \in \Sigma^* \mid \delta^*(q_0, x) = \delta^*(q_0, y)\}$$

de modo que la definición de cada clase de equivalencia depende de una igualdad entre dos estados. Así, como Q es finito entonces sólo hay un número finito de igualdades entre estados, de donde sólo puede haber un número finito de clases de equivalencia.

- Veamos que $L(M) = \bigcup \{[x]_{\equiv_M} \mid \delta^*(q_0, x) \in F\}$. La contención (\subseteq) es clara. Para (\supseteq), tomemos $y \in [x]_{\equiv_M}$ con $\delta^*(q_0, x) \in F$. Entonces como $y \equiv_M x$, tenemos que $\delta^*(q_0, y) = \delta^*(q_0, x) \in F$, es decir $y \in L(M)$. □

¹Recordemos que el índice de una relación de equivalencia \equiv es el número de clases de equivalencia generadas por \equiv .

Las propiedades de la proposición 2 resultan características de un lenguaje regular y esto es precisamente lo que nos dice el siguiente

Teorema 1 (Myhill-Nerode). *Sea $L \subseteq \Sigma^*$. Las siguientes condiciones son equivalentes:*

1. L es regular.
2. Existe una relación de equivalencia \equiv sobre Σ^* , invariante por la derecha y de índice finito, tal que L es la unión de algunas de las clases de equivalencia de \equiv .
3. La relación de equivalencia \equiv_L tiene índice finito.

Demostración.

(1) \Rightarrow (2). Esto es consecuencia directa de la proposición 2.

- (2) \Rightarrow (3). Sean \equiv una relación de equivalencia de índice finito e invariante por la derecha y $w_0, \dots, w_n \in \Sigma^*$ tales que

$$L = \bigcup_{i=0}^n [w_i]_{\equiv}.$$

Para mostrar que \equiv_L es de índice finito, basta ver que \equiv es un refinamiento de \equiv_L , puesto que así el índice de \equiv_L es menor² que el índice de \equiv , que sabemos finito. Supongamos pues que $x \equiv y$ y mostremos que $x \equiv_L y$. Para esto sea $z \in \Sigma^*$

$$\begin{aligned} xz \in L \text{ syss } xz &\in [w_i]_{\equiv} \quad (\text{para alguna } i) \\ \text{syss } xz &\equiv w_i \\ \text{syss } yz &\equiv w_i \quad (\equiv \text{ es invariante por la derecha y } x \equiv y) \\ \text{syss } yz &\in [w_i]_{\equiv} \quad (\text{para alguna } i) \\ \text{syss } yz &\in L \end{aligned}$$

así hemos probado que $\forall x \in \Sigma^* (xz \in L \Leftrightarrow yz \in L)$, es decir $x \equiv_L y$.

- (3) \Rightarrow (1). Supongamos que \equiv_L es de índice finito y digamos que $[w_0], \dots, [w_n]$ son todas las clases de equivalencia inducidas por \equiv_L , donde s.p.g. $w_0 \equiv_L \varepsilon$. Para mostrar que L es regular definimos el siguiente AFD M :

- $Q = \Sigma^* / \equiv_L = \{[w_0], \dots, [w_n]\}$
- $q_0 = [\varepsilon]$
- $F = \{[w_i] \mid w_i \in L\}$
- Para cualesquiera $x \in \Sigma^*$ y $a \in \Sigma$, la función de transición δ se define como:

$$\delta([w_i], a) = [w_i a]$$

Es fácil ver que δ está bien definida.

Más aún, de la siguiente propiedad, cuya demostración dejamos como ejercicio,

$$\text{Para cualesquiera } w_i, x \in \Sigma^*, \quad \delta^*([w_i], x) = [w_i x]$$

²Cada clase de \equiv_L es unión de clases de \equiv , a saber $[x]_{\equiv_L} = \bigcup \{[y]_{\equiv} \mid x \equiv_L y\}$ y como por hipótesis \equiv es de índice finito entonces \equiv_L también.

se sigue que M acepta exactamente a L :

$$\delta^*([\varepsilon], x) \in F \text{ syss } [x] \in F \text{ syss } x \in L$$

Por lo tanto L es un lenguaje regular. □

El teorema anterior permite demostrar que un lenguaje L **no** es regular al mostrar que L no es de índice finito. Es decir, basta ver que \equiv_L tiene una infinidad de clases de equivalencia. Esto se hace explícito mediante el siguiente lema que es una consecuencia directa del teorema de Myhill-Nerode.

Lema 3 (Lema del índice finito). *Sea $L \subseteq \Sigma^*$ un lenguaje regular infinito. Cualquier conjunto $S \subseteq \Sigma^*$ suficientemente grande contiene al menos dos cadenas distintas, $x, y \in S$ tales que $x \equiv_L y$.*

Demostración. Como L es regular entonces, por el teorema de Myhill-Nerode, la relación \equiv_L es de índice finito, digamos n . Si $S \subseteq L$ y $|S| > n$ entonces, el principio del palomar implica que existen dos cadenas distintas $x, y \in S$ tales que tanto x como y pertenecen a la misma clase de equivalencia. Es decir, $x \equiv_L y$. □

Las pruebas de no regularidad se sirven de la contrapositiva del lema del índice finito. Es decir, hallando un conjunto $S \subseteq \Sigma^*$ suficientemente grande³ que no cumpla la propiedad del lema, habremos probado que L no es regular. La siguiente definición de conjuntos estafadores⁴ servirá para encontrar los conjuntos que no cumplen la propiedad, así para mostrar que un lenguaje L **no** es regular basta construir un conjunto estafador para L .

Definición 2. *Un conjunto infinito $S \subseteq \Sigma^*$ es un conjunto estafador para L si y sólo si para cualesquiera $x, y \in S$ existe una cadena $z \in \Sigma^*$ tal que una y sólo una de xz y de yz pertenece a L . Es decir, S es un conjunto estafador para L si y sólo si*

$$\forall x, y \in S (x \not\equiv_L y).$$

Veamos algunos ejemplos de pruebas de no regularidad usando este concepto.

Ejemplo 5. *El lenguaje $L = \{a^i b^i \mid i \in \mathbb{N}\}$ **no** es regular.*

Basta hallar un conjunto estafador para L . Sea $S = \{a^k \mid k \in \mathbb{N}\}$, veamos que S es un conjunto estafador: Sean $a^i, a^j \in S$ con $i \neq j$, entonces claramente $a^i b^i \in L$ y $a^i b^j \notin L$. Por lo tanto $a^i \not\equiv_L a^j$ y así S es un conjunto estafador para L .

Ejemplo 6. *Decidir si el lenguaje $L = \{a^{i^2} \mid i \in \mathbb{N}\}$ es regular⁵. Tratemos de encontrar un conjunto estafador para el lenguaje. Sea S exactamente L y consideremos $a^{i^2}, a^{j^2} \in S$ con $i < j$.*

- Por un lado tenemos que $a^{i^2} a^{2i+1} = a^{i^2+2i+1} = a^{(i+1)^2} \in L$
- Por otra parte, $a^{j^2} a^{2i+1} = a^{j^2+2i+1} \notin L$ puesto que $j^2 < j^2 + 2i + 1 < j^2 + 2j + 1 = (j+1)^2$.

³Por lo general S es infinito

⁴En inglés *fooling set*.

⁵Analizando informalmente, se puede ver que este lenguaje requiere memoria más allá de la que ofrece un autómata finito pues se necesita contabilizar i^2

Por lo tanto $a^{i^2} \not\equiv_L a^{j^2}$ y S es un conjunto estafador para L .

El siguiente ejemplo se deja al lector y constata que el lenguaje del ejemplo 1 no es regular.

Ejemplo 7. El lenguaje L del ejemplo 1 no es regular, puesto que el conjunto

$$S = \{ab^n \mid n \in \mathbb{N}\}$$

es un estafador para L .