

Expresiones Regulares

Introducción a la Lógica
Fa.M.A.F., Universidad Nacional de Córdoba

28/10/15

Expresiones Regulares.

Autómatas Finitos y Expresiones Regulares

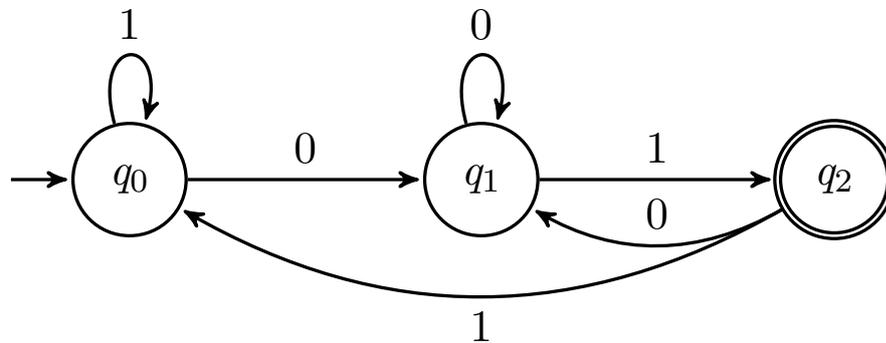
Lenguajes Regulares

Outline of the Talk

- 1 Expresiones Regulares.
 - Intuiciones
 - Operaciones sobre lenguajes
 - Algebra de Expresiones Regulares
- 2 Autómatas Finitos y Expresiones Regulares
 - De Expresiones regulares a Autómatas
 - De Autómatas a Expresiones Regulares
- 3 Lenguajes Regulares

Intuiciones I

- Sabemos que los autómatas finitos (DFA, NFA y NFA- ϵ) definen lenguajes.



$$L(M) = \{w \mid w \text{ termina en } 01\}$$

- Notación **machine-like** para especificar lenguajes.
- Convengamos que no es una notación muy cómoda!

Intuiciones II

Las expresiones regulares:

- Nos permiten dar una descripción algebraica de un lenguaje.

$$(0 + 1)^*01$$

- Introducidas en los 50s por Kleene.
- Se han vuelto una herramienta de uso común en muchos sistemas operativos y lenguajes de programación.
- Muchas aplicaciones de procesamiento de texto las utilizan.
 - Analizadores léxicos.
 - Búsqueda de textos.
- Estudiaremos una versión básica de las expresiones regulares.

Unión y Concatenación de Lenguajes

Unión

Unión de dos lenguajes L_1 y L_2 ($L_1 \cup L_2$): lenguaje formado por las cadenas que están en L_1 y las cadenas que están en L_2 .

Ejemplo: Sea $L_1 = \{00, 11\}$ y $L_2 = \{100, 111\}$ luego,

$$L_1 \cup L_2 = \{00, 11, 100, 111\}.$$

Concatenación

Concatenación de dos lenguajes L_1 y L_2 ($L_1 L_2$): lenguaje que contiene a todas las palabras que pueden ser formadas concatenando cualquier palabra en L_1 con cualquier palabra en L_2 , es decir, $L_1 L_2 = \{xy \mid x \in L_1, y \in L_2\}$.

Ejemplo: Sea $L_1 = \{00, 11\}$ y $L_2 = \{100, 111\}$, luego,

$$L_1 L_2 = \{00100, 00111, 11100, 11111\}.$$

Clausura de un lenguaje

- Definamos $L^0 = \{\epsilon\}$, $L^1 = L$ y $L^i = LL^{i-1}$, para $i \geq 1$.

Clausura de Kleene

La **clausura de un lenguaje** L (L^*) representa el conjunto de palabras que se pueden formar concatenando un nro. arbitrario de cadenas de L . Formalmente,

$$L^* = \bigcup_{i=0}^{\infty} L^i.$$

La **clausura positiva de** L (L^+) se define como

$$L^+ = \bigcup_{i=1}^{\infty} L^i.$$

Clausura de un Lenguaje

Ejemplo: Sea $L = \{0, 11\}$. Calculemos L^* .

$$L^* = \bigcup_{i=0}^{\infty} L^i = L^0 \cup L^1 \cup L^2 \cup L^3 \cup \dots$$

$$L^0 = \{\epsilon\}$$

$$L^1 = L = \{0, 11\}$$

$$L^2 = LL^1 = \{00, 011, 110, 1111\}$$

$$L^3 = LL^2 = \{000, 0011, 0110, 01111, 1100, 11011, 11110, 111111\}$$

$$L^* = \{\epsilon, 0, 11, 00, 011, 110, 1111, \\ 000, 0011, 0110, 01111, 1100, \\ 11011, 11110, 111111, \dots\}$$

Sintaxis de Expresiones Regulares

Sintaxis

El conjunto de las expresiones regulares sobre el alfabeto Σ , ER_{Σ} es el menor conjunto tal que:

- 1 $\emptyset \in ER_{\Sigma}$.
- 2 $\epsilon \in ER_{\Sigma}$.
- 3 $a \in ER_{\Sigma}, \forall a \in \Sigma$.

Si r y s son expresiones regulares, entonces

- 4 $r + s \in ER_{\Sigma}$.
- 5 $rs \in ER_{\Sigma}$ (o $r \cdot s \in ER_{\Sigma}$).
- 6 $r^* \in ER_{\Sigma}$.
- 7 $(r) \in ER_{\Sigma}$.

Semántica de expresiones regulares

Lenguaje de una expresión regular

El lenguaje asociado a una expresión regular r , $L(r)$, se define recursivamente como:

- 1 $L(\emptyset) = \emptyset$.
- 2 $L(\epsilon) = \{\epsilon\}$.
- 3 $L(\mathbf{a}) = \{a\} \forall \mathbf{a} \in ER_{\Sigma}$.

Si r y s son expresiones regulares, entonces

- 4 $L(\mathbf{r} + \mathbf{s}) = L(\mathbf{r}) \cup L(\mathbf{s})$.
- 5 $L(\mathbf{rs}) = L(\mathbf{r})L(\mathbf{s})$.
- 6 $L(\mathbf{r}^*) = L(\mathbf{r})^*$.
- 7 $L((\mathbf{r})) = L(\mathbf{r})$.

Precedencia de Operadores

Precedencia

Sean $*$, $+$, \cdot los operadores utilizados en expresiones regulares. Definimos su precedencia de la siguiente forma:

$$* > \cdot > +$$

Ejemplo:

- $((\mathbf{0(1^*)}) + \mathbf{1}) = \mathbf{01^*} + \mathbf{1}$
- $((\mathbf{01})^*) + \mathbf{1} \neq \mathbf{01^*} + \mathbf{1}$

Ejemplos

① 00 : $L(00) = L(0)L(0) = \{00\}$.

② $r = (0 + 1)^*$:

$$L(r) = \{w \mid w \text{ es una cadena de 0s y 1s}\}.$$

③ $r = (0 + 1)^*00(0 + 1)^*$:

$$L(r) = \text{cadenas de 0s y 1s con al menos dos 0s consecutivos.}$$

④ $r = (0 + 1)^*011$:

$$L(r) = \{w \mid w \text{ es una cadena de 0s y 1s que termina en 011}\}.$$

⑤ $r = (1 + 10)^*$: $L(r) =$ cadenas de 0s y 1s que comienzan con 1 y no tienen dos 0s consecutivos.

⑥ $r = 0^*1^*2^*$: $L(r) =$ cadenas que tienen cierta cantidad de 0s seguido por cierta cantidad de 1s seguido por cierta cantidad de 2s.

Recordar!

Importante

No confundir una expresión regular con una palabra del lenguaje.

- $001 \neq 001$.
- Una expresión regular es un patrón que describe un conjunto de palabras.

Leyes Algebraicas

Leyes Algebraicas

Sean r, s y t expresiones regulares. Se cumple lo siguiente:

- $r + s = s + r.$
- $(r + s) + t = r + (s + t).$
- $(rs)t = r(st).$
- $\emptyset + r = r + \emptyset = r.$
- $\epsilon r = r\epsilon = r.$
- $\emptyset r = r\emptyset = \emptyset.$
- $r(s + t) = rs + rt.$
- $(s + t)r = sr + tr.$
- $r + r = r.$
- $(r^*)^* = r^*.$
- $\emptyset^* = \epsilon.$
- $\epsilon^* = \epsilon.$
- $r^+ = rr^* = r^*r.$
- $r^* = r^+ + \epsilon.$

Autómatas y Expresiones Regulares

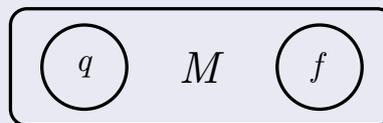
- Los autómatas finitos y las expresiones regulares definen lenguajes.
- Ahora vamos a ver que ambos definen los **mismos** lenguajes.
- Vamos a probar que partiendo de una expresión regular podemos construir un NFA- ϵ que define el mismo lenguaje.
 - Bastante sencillo!
- Vamos a probar que partiendo de un NFA- ϵ podemos construir una expresión regular que define el mismo lenguaje.
 - Mas complicado!.

Un poco de notación

- Queremos ver que a toda expresión regular r le corresponde un autómata NFA- ϵ M tal que $L(r) = L(M)$.

Representación de un autómata

Sea M un autómata con estado inicial q y estado final f . Lo representaremos de la siguiente forma:



Transformando expresiones regulares en autómatas

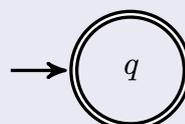
Algoritmo

Sean $r, s \in ER_{\Sigma}$. Definimos inductivamente el autómata M_r correspondiente a una expresión regular r de la siguiente forma:

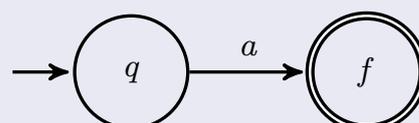
- 1 A \emptyset le corresponde el autómata:



- 2 A ϵ le corresponde el autómata:



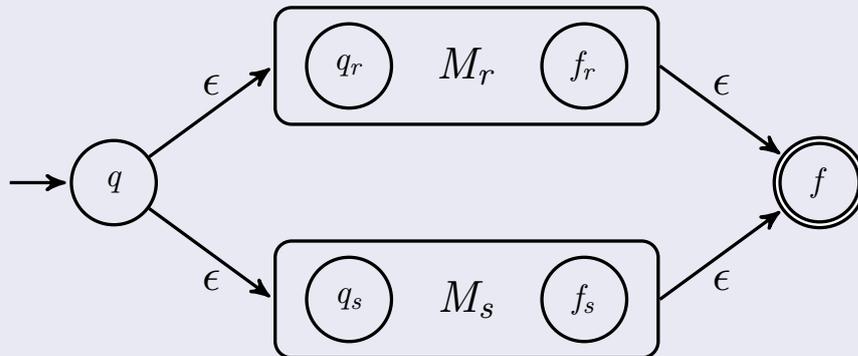
- 3 A toda $a \in ER_{\Sigma}$, le corresponde el autómata:



Transformando expresiones regulares en autómatas

Algoritmo (cont.)

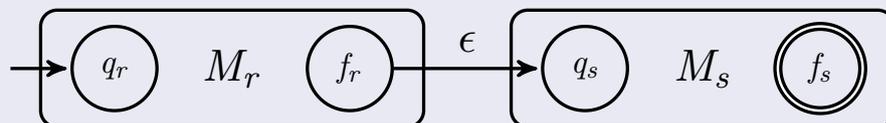
- 4 A $r + s$ le corresponde el autómata:



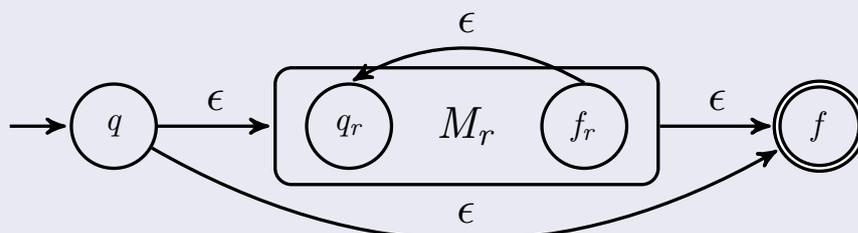
Transformando expresiones regulares en autómatas

Algoritmo (cont.)

- 5 A rs le corresponde el autómata:



- 6 A r^* le corresponde el autómata:

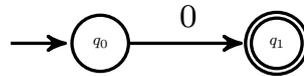


Ejemplo

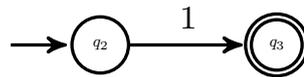
Construir el NFA- ϵ para $r = 01^* + 1$.

$L(r) = \{w \mid w = 1 \text{ o } w \text{ empieza con } 0 \text{ seguida por cero o más } 1\text{'s}\}$

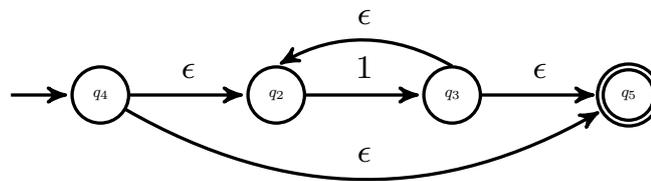
• **0:**



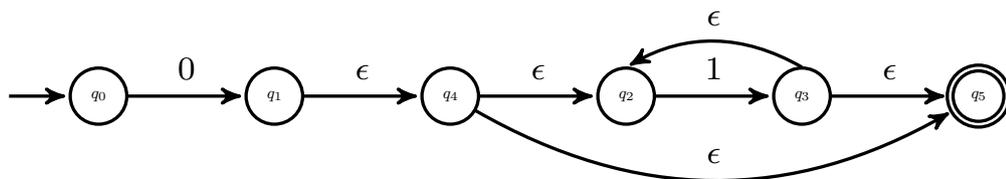
• **1:**



• **1*:**

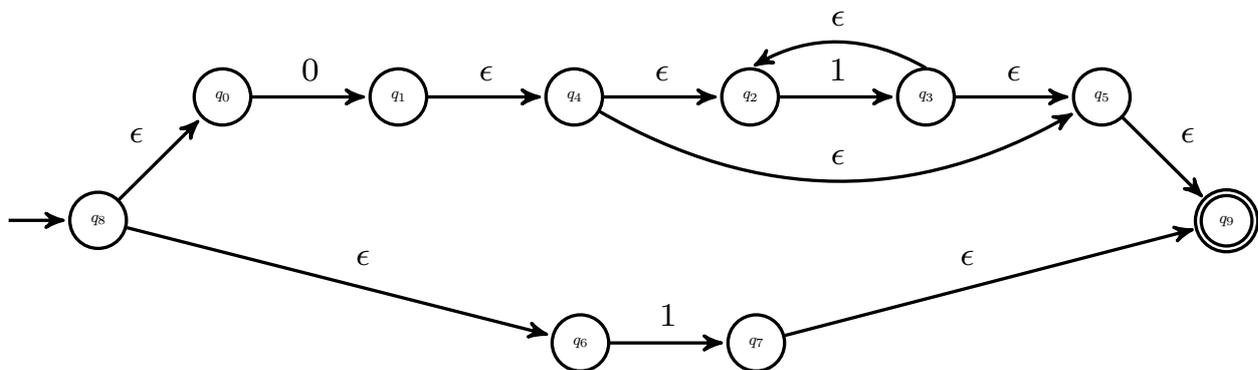


• **01*:**



Ejemplo (cont.)

• **01* + 1:**



Intuiciones

- Vamos a construir expresiones regulares que describan “caminos” dentro del autómata.
- Nos interesan los caminos que van desde el estado inicial a los estados finales.
- Camino = Palabra.
- En cada paso de la construcción vamos a considerar autómatas mas pequeños.

Un poco de notación

Dados $M = (Q, \Sigma, \delta, q_0, F)$, $R \subseteq Q$ y $q_n, q_m \in R$, definimos

$$M_{nm}(R) = (R, \Sigma, \delta|_R, q_n, \{q_m\}).$$

Algoritmo

Definición

Sea $M = (Q, \Sigma, \delta, q_0, F)$ un autómata. Sea $Q = \{q_0, \dots, q_r\}$ y $\Sigma = \{a_1, \dots, a_u\}$. Dados $0 \leq n, m \leq r$, $R \subseteq Q$, definimos recursivamente las siguientes expresiones regulares.

$$L_{nm}(R) = \begin{cases} \emptyset & \text{si } q_n \circ q_m \notin R. \\ I_n(R)^* & \text{si } n = m. \\ I_n(R)^* F_{nm}(R) & \text{si } n \neq m. \end{cases}$$

$$I_n(R) = \sum_{\substack{q_n \xrightarrow{a} q_t \\ q_s \xrightarrow{b} q_n}} \mathbf{a} L_{ts}(R \setminus \{q_n\}) \mathbf{b} + \sum_{q_n \xrightarrow{c} q_n} \mathbf{c}.$$

$$F_{nm}(R) = \sum_{q_n \xrightarrow{a} q_t} \mathbf{a} L_{tm}(R \setminus \{q_n\}).$$

Entendiendo el algoritmo - $L_{nm}(R)$ (1)

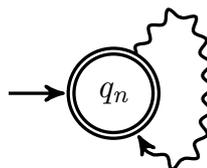
$$L_{nm}(R) = \begin{cases} \emptyset & \text{si } q_n \circ q_m \notin R. \\ I_n(R)^* & \text{si } n = m. \\ I_n(R)^* F_{nm}(R) & \text{si } n \neq m. \end{cases}$$

- $L_{nm}(R)$ es el lenguaje formado por las cadenas que llevan del estado q_n al estado q_m en el autómata $M_{nm}(R)$.
- $L_{nm}(R)$ depende de n , m y R :
- Si $q_n, q_m \notin R$: $L_{nm}(R)$ es el lenguaje vacío.

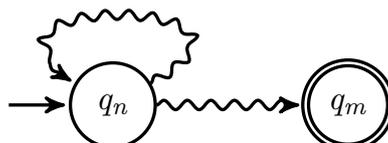
Entendiendo el algoritmo - $L_{nm}(R)$ (2)

$$L_{nm}(R) = \begin{cases} \emptyset & \text{si } q_n \circ q_m \notin R. \\ I_n(R)^* & \text{si } n = m. \\ I_n(R)^* F_{nm}(R) & \text{si } n \neq m. \end{cases}$$

- Si $q_n, q_m \in R$ y $q_n = q_m$: $L_{nm}(R)$ es el lenguaje formado por las cadenas que salen de q_n y vuelven a q_n , sin pasar por q_n .



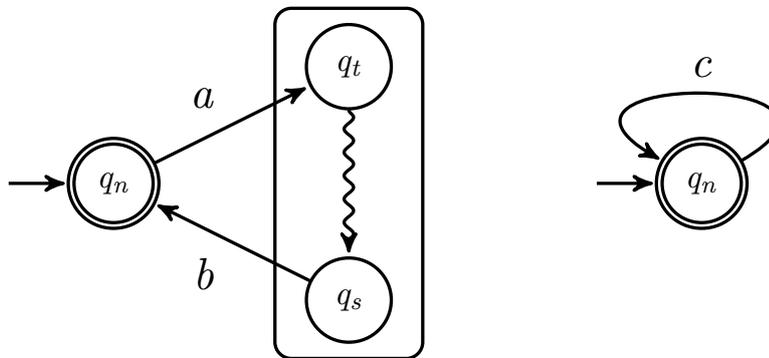
- Si $q_n, q_m \in R$ y $q_n \neq q_m$: $L_{nm}(R)$ es el lenguaje que resulta de concatenar el lenguaje de las cadenas que salen de q_n y vuelven a q_n , sin pasar por q_n , con el lenguaje de las cadenas que van de q_n a q_m .



Entendiendo el algoritmo - $I_n(R)$

$$I_n(R) = \sum_{q_n \xrightarrow{a} q_t, q_s \xrightarrow{b} q_n} \mathbf{a}L_{ts}(R \setminus \{q_n\})\mathbf{b} + \sum_{q_n \xrightarrow{c} q_n} \mathbf{c}.$$

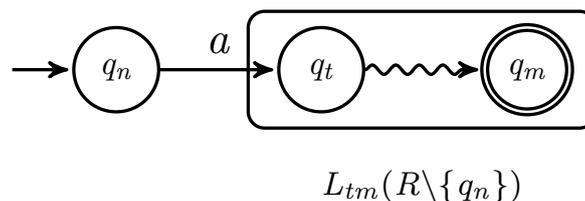
- $I_n(R)$ es el lenguaje formado por las cadenas que salen del estado q_n y vuelven a q_n , sin pasar por q_n , en el autómata $M_{nn}(R)$.
- $I_n(R)$ es la unión de 2 lenguajes.
- Las cadenas que **salen** de q_n y **regresan** a q_n .
- Las cadenas de longitud 1 (loops) que salen y entran a q_n .



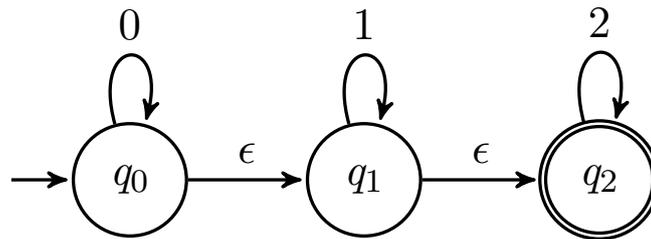
Entendiendo el algoritmo - $F_{nm}(R)$

$$F_{nm}(R) = \sum_{q_n \xrightarrow{a} q_t} \mathbf{a}L_{tm}(R \setminus \{q_n\}).$$

- $F_{nm}(R)$ es lenguaje formado por las cadenas que **salen** del estado q_n y llevan al estado q_m .



Ejemplo



$$L_{02}(Q) = I_0(Q)^* F_{02}(Q) = \mathbf{0}^*(\mathbf{0} + \epsilon)\mathbf{1}^*(\mathbf{1} + \epsilon)\mathbf{2}^* = \mathbf{0}^*\mathbf{1}^*\mathbf{2}^*$$

$$I_0(Q) = \mathbf{0}$$

$$F_{02}(Q) = (\mathbf{0} + \epsilon)L_{12}(Q \setminus \{q_0\}) = (\mathbf{0} + \epsilon)\mathbf{1}^*(\mathbf{1} + \epsilon)\mathbf{2}^*$$

$$L_{12}(Q \setminus \{q_0\}) = I_1(Q \setminus \{q_0\})^* F_{12}(Q \setminus \{q_0\}) = \mathbf{1}^*(\mathbf{1} + \epsilon)\mathbf{2}^*$$

$$I_1(Q \setminus \{q_0\}) = \mathbf{1}$$

$$F_{12}(Q \setminus \{q_0\}) = (\mathbf{1} + \epsilon)L_{22}(Q \setminus \{q_0, q_1\}) = (\mathbf{1} + \epsilon)\mathbf{2}^*$$

$$L_{22}(Q \setminus \{q_0, q_1\}) = I_2(Q \setminus \{q_0, q_1\})^* = \mathbf{2}^*$$

$$I_2(Q \setminus \{q_0, q_1\}) = \mathbf{2}$$

Equivalencia entre expresiones regulares y autómatas

Teorema

Sea r una expresión regular y $L(r)$ el lenguaje asociado. Entonces existe un NFA- ϵ M_r tal que $L(M_r) = L(r)$.

Prueba: Por inducción en la estructura de r .

Teorema de Kleene

Sea M un NFA con ϵ -transiciones y $L(M)$ el lenguaje asociado a M . Entonces existe una expresión regular r tal que $L(M_r) = L(r)$.

Prueba: Hay que probar 2 lemas intermedios.

Lenguajes Reconocidos

- Los DFA, NFA, NFA- ϵ y las expresiones regulares son equivalentes.
 - Reconocen los mismos lenguajes: **lenguajes regulares**.
- ¿Todos los lenguajes pueden ser descritos de esta forma?
- NO. Existen lenguajes que no pueden ser descritos por autómatas finitos y/o expresiones regulares.

Ejemplo

Considere el lenguaje $L_{01} = \{0^n 1^n \mid n \geq 1\}$. ¿Es regular?

Prueba: Si L_{01} es regular existe un DFA M con k estados tq $L(M) = L_{01}$. Sabemos que si $|\alpha| = k$ entonces hay $k + 1$ prefijos. Supongamos que M recibe k 0's. Por cada prefijo, M esta en un estado, como hay $k + 1$ prefijos y k estados, por el principio del pidgeonhole, existen $i, j \leq k + 1$ y $i \neq j$ tq $q_0 \xrightarrow{0^i} p$ y $q_0 \xrightarrow{0^j} p$. Si luego de recibir el prefijo 0^i (ó 0^j) M empieza a recibir 1's, luego de i 1's M debería aceptar la cadena $(p \xrightarrow{1^i} q_f, q_f \in F)$ solo si antes se recibieron i 0's (prefijo 0^i) y no aceptar la cadena si se recibieron j 0's. Sin embargo como estamos en el estado p , no podemos distinguir si hemos recibido el prefijo 0^i ó 0^j . Por lo tanto, el autómatata podría aceptar la cadena $0^j 1^i$ con $i \neq j$. Luego L_{01} no es un lenguaje regular.