

# Expresiones regulares y derivadas

Teoría de Lenguajes

1<sup>er</sup> cuatrimestre de 2020

# Plan de la clase

- 1 Expresiones regulares
  - Definiciones y propiedades
  - Ejemplos
  - Otras propiedades
  - Extensiones
- 2 Conversión de autómatas finitos a expresión regular
  - Definiciones y propiedades
  - Ejemplo
  - Conversión
- 3 Conversión de expresión regular a autómatas finitos
  - Derivada de un lenguaje
  - Derivada de una expresión regular
  - Ejemplos
  - Conversión

# Expresiones regulares

Las expresiones regulares son expresiones que se utilizan para *denotar* o *definir* lenguajes regulares. No sirven para denotar lenguajes menos restrictivos que los regulares (libres de contexto y demás).

Ejemplo: \*.txt

Las usan grep, awk, sed, perl, muchos editores de texto.

Se definen a partir de conjuntos simples y los que se pueden obtener de ellos a partir de las operaciones de unión, concatenación y clausura.

Cada expresión regular  $e$  va a denotar un lenguaje regular  $L(e)$ .

# Expresiones regulares

Se definen recursivamente de la siguiente manera:

- ① La expresión regular  $\emptyset$  define el lenguaje vacío.
- ② La expresión regular  $\lambda$  define el lenguaje  $\{\lambda\}$ .
- ③ La expresión regular  $a$  define el lenguaje  $\{a\}$ .
- ④ Sean:

$v$  la expresión regular que define el lenguaje  $L_v$

$w$  la expresión regular que define el lenguaje  $L_w$

entonces:

- ① la expresión regular  $(v.w)$  define el lenguaje  $L_v.L_w$
- ② la expresión regular  $(v|w)$  define el lenguaje  $L_v \cup L_w$
- ③ la expresión regular  $v^*$  define el lenguaje  $L_v^*$
- ④ la expresión regular  $v^+$  define el lenguaje  $L_v^+$

Precedencia: 1)  $*$  , 2)  $.$  , 3)  $|$

El operador  $.$  se puede omitir.

# Ejemplos

- La expresión regular  $(a|b)^+$  define el lenguaje de todas las cadenas de as y bs de longitud mayor o igual que 1.
- La expresión regular  $(1|\dots|9)(0|1|\dots|9)^*0$  define todos los números naturales sin ceros no significativos.
- Cadenas de 0s y 1s terminadas en 01:

$$(((0|1)^*0)1) = (0|1)^*01$$

- Cadenas que tienen como subcadena a 000:

$$((((0|1)^*0)0)0)(0|1)^* = (0|1)^*000(0|1)^*$$

# Semántica

## Definición

- $L(u.v) = L(u).L(v)$
- $L(u|v) = L(u) \cup L(v)$
- $L(u^*) = (L(u))^*$

## Definición

Dos expresiones regulares  $e$  y  $e'$  se dicen equivalentes, si  $L(e) = L(e')$ . Se suele escribir  $e = e'$ .

# Propiedades

La unión y la concatenación son asociativas:

## Asociatividad

$$(u|v)|w = u|(v|w) = u|v|w$$

$$(u.v).w = u.(v.w) = uvw$$

# Ejemplos

$$\Sigma = \{0, 1\}$$

Dar una expresión regular que defina:

- 1 Cadenas que terminan en 01:  $(0|1)^*01$
- 2 Cadenas con 0s y 1s alternados:  $(0|\lambda)(10)^*(1|\lambda)$
- 3 Complemento del primero (cadenas que NO terminan en 01):  
 $((0|1)^*(0|11))|1|\lambda$
- 4 Cadenas con cantidad par de ceros:  
 $1^*(01^*01^*)^*$  o equivalentemente:  $(1|01^*0)^*$



# Aclaraciones

- Por lo general no vamos a pedir que se demuestre que la ER denota el lenguaje pedido, pero sí tiene que haber una justificación de la misma.
- Como dijimos, el poder expresivo de las ER y los AF son el mismo. Pero hay veces que un formalismo es más intuitivo para definir uno que el otro. Por ejemplo, el ejemplo 4 es muy fácil de aceptar con AF.
- *Recordar que aceptar un lenguaje  $L$  (lo mismo, denotar un lenguaje) significa exactamente  $L$ , no alguno que lo incluya.*

# Ejercicios

- 1 Constantes reales con signo y notación exponencial.
- 2 Nombres de archivo con path (para DOS/Win, para Unix).
- 3 Direcciones de e-mail (sólo de Argentina o Uruguay).

# Conmutatividad y distributividad

Conmutatividad de  $|$

$$u|v = v|u$$

(notar que  $.$  no es conmutativo)

Distributividad de  $.$  respecto de  $|$

$$u(v|w) = uv|uw$$

$$(u|v)w = uw|vw$$

# Neutro y absorbente

Elemento neutro de  $\cdot$

$$u \cdot \lambda = \lambda \cdot u = u$$

Elemento absorbente de  $\cdot$

$$u \cdot \emptyset = \emptyset \cdot u = \emptyset$$

Elemento neutro de  $|$

$$u | \emptyset = \emptyset | u = u$$

## Otros operadores

Los operadores  $.$ ,  $|$ ,  $*$  y  $+$  se comportan de la forma que vimos.

Si  $u$  es una expresión regular,  $n$  y  $m$  naturales:

- $u? = (u|\lambda)$
- $u\{n\} = u^n$
- $u\{n, \} = u^n.u^*$
- $u\{n, m\} = u^n|u^{n+1}|\dots|u^m$
- $() = \lambda$
- $.$  = cualquier carácter
- $^$  = comienzo de línea
- $\$$  = fin de línea
- $\backslash$ seguida de alguno de los caracteres  $^.[\$()—*+?\{\\$  representa ese carácter literal

# Conjuntos de caracteres

- $[abc] = a|b|c$
- $[a - z] = a|...|z$
- $^$  al comienzo de la lista indica complemento  $[^abc] =$  cualquier caracter menos a, b o c.

Hay clases de caracteres predefinidas: (en posix) *alnum*, *alpha*, *digit*, *space*...

# Conversión de autómata finito a expresión regular

Sea

$A = (Q = \{q_0, q_1, \dots, q_n\}, \Sigma = \{a_1, a_2, \dots, a_m\}, F \subseteq Q, q_0, \delta_A)$   
un autómata finito determinístico.

$\delta_A$  es la función de transición y está definida en:

$$\delta_A : Q \times \Sigma \rightarrow Q$$

Esta función se puede extender para cadenas de cualquier longitud sobre  $\Sigma$ , definiendo la función  $\delta$ :

$$\delta : Q \times \Sigma^* \rightarrow Q \text{ tal que } \begin{cases} \delta(q, ax) &= \delta(\delta_A(q, a), x), \text{ para } x \in \Sigma^* \\ \delta(q, \lambda) &= q \end{cases}$$

## Generalización

El lenguaje aceptado por el autómata se define como:

$$L(A) = \{x \in \Sigma^* \mid \delta(q_0, x) \in F\}$$

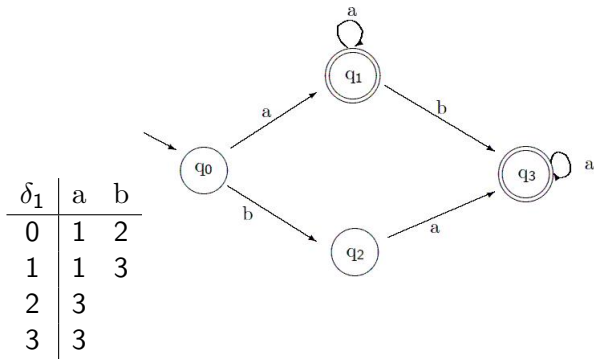
De la misma forma que definimos el lenguaje  $L$  como el conjunto de cadenas aceptadas por el autómata partiendo del estado  $q_0$ , podemos definir un lenguaje  $L_i$  para un estado  $q_i$  dado del autómata, que corresponde al lenguaje aceptado por el autómata si partiéramos, no de  $q_0$  sino de ese estado  $q_i$ :

$$L_i = \{x \in \Sigma^* \mid \delta(q_i, x) \in F\}$$



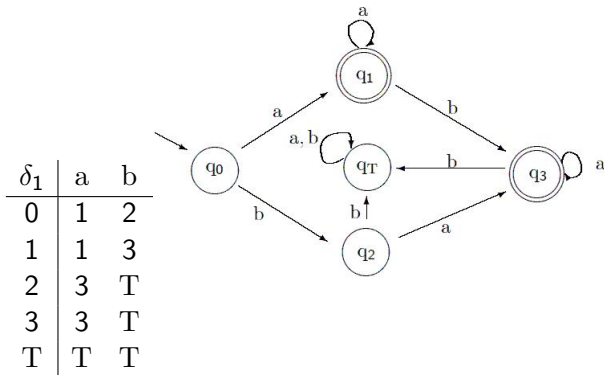
## Ejemplo

Dado el autómata  $A = (\{q_0, q_1, q_2, q_3\}, \{a, b\}, \{q_1, q_3\}, q_0, \delta_1)$



## Ejemplo completado con trampa

Dado el autómata  $A = (\{q_0, q_1, q_2, q_3\}, \{a, b\}, \{q_1, q_3\}, q_0, \delta_1)$



Autómata completado con el estado trampa. Trampa es un estado no final, adonde van las transiciones no definidas, y que cicla sobre el mismo con cada símbolo de  $\Sigma$ .

## Observación

Partiendo del estado  $q_0$ , tendremos el lenguaje  $L = L_0$  aceptado por el autómata. Ahora, si en vez de partir del estado  $q_0$  partiéramos del estado  $q_1$ , tendríamos el lenguaje  $L_1$ , y así con los demás. En particular, el estado trampa define el lenguaje vacío, ya que partiendo de él no se puede reconocer ninguna cadena, i.e.  $L_T = \emptyset$ .

## Ecuaciones

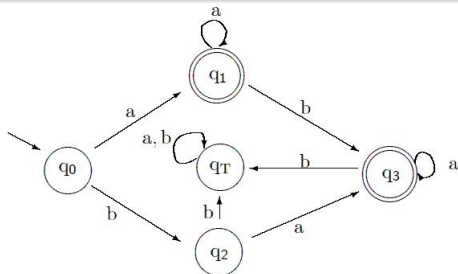
En el autómata definido antes, veamos que para cada estado  $q_i$  dado ( $1 \leq i \leq n$ ), podemos considerar las siguientes igualdades, donde  $\Sigma = \{a_1, \dots, a_m\}$ :

$$\begin{aligned}\delta(q_i, a_1) &= p_{k_1} \\ \delta(q_i, a_2) &= p_{k_2} \\ &\dots \\ \delta(q_i, a_m) &= p_{k_m}\end{aligned}\tag{1}$$

y a partir de ellas expresar el lenguaje  $L_i$  mediante la siguiente ecuación:

$$L_i = a_1.L_{k_1} | a_2.L_{k_2} | \dots | a_m.L_{k_m} | \epsilon(L_i)\tag{2}$$

## Ecuaciones



En el autómata del ejemplo, estas ecuaciones serán:

$$L_0 = a.L_1 | b.L_2$$

$$L_1 = a.L_1 | b.L_3 | \lambda$$

$$L_2 = a.L_3 | b.L_T$$

$$L_3 = a.L_3 | \lambda$$

$$L_T = a.L_T | b.L_T | \emptyset$$

Ahora, hay que resolverlas.

## Ecuaciones entre lenguajes

Imaginen ejemplos como

$$\begin{array}{ll} L = aL & L = La \\ L = La|aL & L = a|bL \\ L = LL & L = LaL \\ L = L^* & L = L \\ \dots & \end{array}$$

¿Cómo se resuelven? Es decir, dada una ecuación, ¿cómo encontramos un lenguaje  $L$  tal que al sustituirlo a ambos lados de la igualdad esta se cumpla? (Más allá de la solución trivial  $\emptyset$  que sirve en muchos casos...) Es complicado dar un método general. Pero nos interesan particularmente las ecuaciones de la forma  $L = \alpha L | \beta$  donde  $\alpha, \beta$  son expresiones regulares sobre  $\Sigma$ .

# Propiedad

Podemos resolver el sistema de ecuaciones anterior (asociado al autómata) usando la siguiente propiedad de los lenguajes definidos sobre  $\Sigma$  (junto a otras propiedades de e.r.):

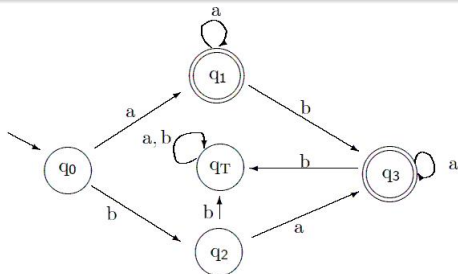
## Proposición

Dados  $R, \alpha, \beta$  lenguajes sobre  $\Sigma$ ,

$$\text{si } R = \alpha.R|\beta \wedge \lambda \notin \alpha, \text{ entonces } R = \alpha^*.\beta$$

Si no se cumpliera la hipótesis  $\lambda \notin \alpha$ , podría haber más soluciones... incluso no regulares.

## Conversión



$L_T = (a|b).L_T|\emptyset$ , y como  $\lambda \notin \{a, b\}$ , entonces  $L_T = (a|b)^*.\emptyset = \emptyset$

$L_3 = a.L_3|\lambda$ , y como  $\lambda \notin \{a\}$ , entonces  $L_3 = a^*.\lambda = a^*$

$L_2 = a.L_3|b.L_T = a.a^*|\emptyset = a^+$

$L_1 = a.L_1|b.L_3|\lambda = a.L_1|(b.a^*|\lambda)$ , y como  $\lambda \notin \{a\}$ ,  
entonces  $L_1 = a^*.(b.a^*|\lambda)$

$L_0 = a.L_1|b.L_2 = a^+(b.a^*|\lambda)|ba^+$



# Derivada de un lenguaje

Dado un lenguaje  $L$  sobre un alfabeto  $\Sigma$  definimos

## Derivada

- 1  $\partial\lambda(L) = L$
- 2  $\partial a(L) = \{\alpha : a\alpha \in L\}$  donde  $a \in \Sigma, \alpha \in \Sigma^*$
- 3  $\partial w(L) = \partial u(\partial a(L))$  si  $w = a.u, a \in \Sigma, u \in \Sigma^*$

Notaremos  $\partial_{ab} L = \partial b(\partial a(L))$

# Ejemplos

Por ejemplo, si

$L = \{a^{3n} : n \in \mathbf{N}_{\geq 1}\} = \{aaa, aaaaaa, aaaaaaaaa, \dots\}$ , entonces:

$$\partial_a L = \{a^{3n-1} : n \in \mathbf{N}_{\geq 1}\} = \{aa, aaaaa, aaaaaaaaa, \dots\}$$

$$\partial_b L = \emptyset$$

$$\partial_{aaa} L = L \cup \{\lambda\}$$

$$\partial_{aaaa} L = \partial_a L$$

## Derivada de una expresión regular

Podemos definir la derivada  $\partial_w(u)$  de una expresión regular  $u$  con respecto a una cadena  $w \in \Sigma^*$ , mediante las siguientes reglas:

- 1  $\partial_\lambda(u) = u$
- 2 Para  $a \in \Sigma$ :

### Derivadas

$$\begin{aligned}\partial_a(\emptyset) &= \emptyset \\ \partial_a(\lambda) &= \emptyset \\ \partial_a(b) &= \begin{cases} \lambda & \text{si } a = b \\ \emptyset & \text{si } a \neq b \end{cases}\end{aligned}$$

## Derivada de una expresión regular

Si  $u$  y  $v$  son expresiones regulares, entonces:

### Derivada

$$\begin{aligned}\partial a(u|v) &= \partial a(u)|\partial a(v) \\ \partial a(u.v) &= \partial a(u).v|\epsilon(u).\partial a(v) \\ \partial a(u^*) &= \partial a(u).u^*\end{aligned}$$

donde:

### $\epsilon(\bullet)$

$$\epsilon(u) = \begin{cases} \emptyset & \text{si } \lambda \notin L_u \\ \lambda & \text{si } \lambda \in L_u \end{cases}$$

Para  $a \in \Sigma \wedge w \in \Sigma^*$ :

$$\partial aw(u) = \partial w(\partial a(u))$$

# Ejemplos

$$\partial a(a) = \lambda$$

$$\partial a(a|b) = \partial a(a)|\partial a(b) = \lambda|\emptyset = \lambda$$

$$\partial a((a|b).c) = \partial a(a|b).c|\epsilon(a|b).\partial a(c) = \lambda.c|\emptyset.\emptyset = c$$

$$\partial a(a^*) = a^*$$

$$\partial a(a^+) = \partial a(a.a^*) = \partial a(a).a^*|\epsilon(a).\partial a(a^*) = \lambda.a^*|\emptyset.a^* = a^*$$

$$\partial a(ba^{20}) = \dots = \emptyset$$

## Conversión

Podemos formular un método para, dada una expresión regular que define un lenguaje, encontrar un autómata finito que acepta ese lenguaje.

Primero observemos que las ecuaciones (1) se pueden formular en términos de las derivadas del lenguaje  $L_i$ :

### Ecuaciones

$$\begin{aligned}\partial a_1(L_i) &= L_{k_1} \\ \partial a_2(L_i) &= L_{k_2} \\ &\dots \\ \partial a_m(L_i) &= L_{k_m}\end{aligned}$$

Reemplazando estos valores en (2) nos queda:

$$L_i = a_1.\partial a_1(L_i)|a_2.\partial a_2(L_i)|\dots|a_m.\partial a_m(L_i)|\epsilon(L_i)$$

## Observaciones

Si tenemos la expresión regular que define  $L$ , tenemos la expresión regular que define  $L_0$ . Si derivamos esta expresión regular con respecto a un símbolo  $a_j \in \Sigma$ , obtendremos la expresión regular que define el lenguaje aceptado por el autómata comenzando por el estado  $\delta(q_0, a_j)$  (puede ser el mismo  $L_0$  u otro lenguaje  $L_k$ ). Así iremos obteniendo una cantidad finita de expresiones regulares, que iremos adoptando como estados del autómata.

## Conversión

Veamos cómo obtenemos el autómata que corresponde a la expresión regular  $a^+(b.a^*|\lambda)|ba^+ = L_0$ :

$$\begin{aligned}
 \partial a L_0 &= \partial a(a^+(b.a^*|\lambda))|\partial a(ba^+) = \\
 &\quad a^*(b.a^*|\lambda)|\emptyset.\partial a(ba^*|\lambda)|\emptyset = a^*(b.a^*|\lambda) = L_1 \\
 \partial b L_0 &= \partial b(a^+(b.a^*|\lambda)|ba^+) = \\
 &\quad \partial b(a^+).(b.a^*|\lambda)|\emptyset.\partial b(b.a^*|\lambda)|\partial b(ba^+) = a^+ = L_2 \\
 \partial a L_1 &= \partial a(a^*(b.a^*|\lambda)) = \partial a(a^*).(b.a^*|\lambda)|\lambda.\partial a(b.a^*|\lambda) = \\
 &\quad a^*(b.a^*|\lambda)|\lambda.\emptyset = a^*(b.a^*|\lambda) = L_1 \\
 \partial b L_1 &= \partial b(a^*(b.a^*|\lambda)) = \emptyset.(b.a^*|\lambda)|\lambda.\partial b(b.a^*|\lambda) = a^* = L_3 \\
 \partial a L_2 &= \partial a(a^+) = a^* = L_3 \\
 \partial b L_2 &= \partial b(a^+) = \emptyset = L_T \\
 \partial a L_3 &= \partial a(a^*) = a^* = L_3 \\
 \partial b L_3 &= \partial b(a^*) = \emptyset = L_T \\
 \partial a L_T &= \partial a \emptyset = \emptyset \\
 \partial b L_T &= \partial b \emptyset = \emptyset
 \end{aligned}$$



# Conversión

$$\begin{aligned}\partial_a L_0 &= a^*(b.a^*|\lambda) = L_1 \\ \partial_b L_0 &= a^+ = L_2 \\ \partial_a L_1 &= a^*(b.a^*|\lambda) = L_1 \\ \partial_b L_1 &= a^* = L_3 \\ \partial_a L_2 &= a^* = L_3 \\ \partial_b L_2 &= \emptyset = L_T \\ \partial_a L_3 &= a^* = L_3 \\ \partial_b L_3 &= \emptyset = L_T \\ \partial_a L_T &= \emptyset = L_T \\ \partial_b L_T &= \emptyset = L_T\end{aligned}$$

## Conclusión

Con las derivadas calculadas podemos construir el autómata correspondiente. Los estados serán las expresiones regulares  $L_i$  obtenidas, y  $\delta(L_i, a) = L_j$  cada vez que  $\partial_a L_i = L_j$ . Serán estados finales aquellos  $L_i$  tales que  $\lambda \in L_i$ .

El estado trampa  $L_T$  siempre corresponde a la expresión regular  $\emptyset$ . En el caso del ejemplo, son finales  $L_1$  y  $L_3$ .