

Expresiones regulares y derivadas

Teoría de Lenguajes

1^{er} cuatrimestre de 2002

1 Expresiones regulares

Las expresiones regulares son expresiones que se utilizan para denotar lenguajes regulares. No sirven para denotar lenguajes menos restrictivos que los regulares (libres de contexto y demás). Se definen recursivamente de la siguiente manera:

1. La expresión regular \emptyset define el lenguaje vacío.

La expresión regular λ define el lenguaje $\{\lambda\}$.

La expresión regular a define el lenguaje $\{a\}$.

2. Sean:

v la expresión regular que define el lenguaje L_v

w la expresión regular que define el lenguaje L_w

entonces:

- (a) la expresión regular $(v.w)$ define el lenguaje $L_v.L_w$
- (b) la expresión regular $(v|w)$ define el lenguaje $L_v \cup L_w$
- (c) la expresión regular v^* define el lenguaje L_v^*
- (d) la expresión regular v^+ define el lenguaje L_v^+

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.

2 Derivada de un lenguaje

Dado un lenguaje L sobre un alfabeto Σ definimos las siguientes derivadas:

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 $\partial ab L = \partial b(\partial a(L))$

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

$$\begin{aligned}\partial a L &= \{a^{3n-1} : n \in \mathbf{N}_{\geq 1}\} = \{aa, aaaaa, aaaaaaaaa, \dots\} \\ \partial b L &= \emptyset \\ \partial aaa L &= L \cup \{\lambda\} \\ \partial aaaa L &= \partial a L\end{aligned}$$

Para un lenguaje L definimos también $\epsilon(L)$, que indica si la cadena nula pertenece o no al lenguaje.

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

3 Derivada de una expresion regular

De la misma manera que definimos la derivada de cadenas de un lenguaje, podemos definir la derivada $\partial x(u)$ de una expresión regular u con respecto a una cadena $x \in \Sigma^*$, mediante las siguientes reglas:

1. $\partial\lambda(u) = u$
2. Para $a \in \Sigma$:

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

3. Sean u y v expresiones regulares que denotan los lenguajes L_u y L_v respectivamente, entonces:

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

donde:

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

4. Para $a \in \Sigma \wedge x \in \Sigma^*$:

$$\partial a x u = \partial x (\partial a u)$$

Ejemplos:

$$\begin{aligned}\partial a(a) &= \lambda \\ \partial a(a|b) &= \partial a(a)|\partial a(b) = \lambda|\emptyset = \lambda \\ \partial a((a|b).c) &= \partial a(a|b).c|\epsilon(a|b).\partial a(c) = \lambda.c|\emptyset.\emptyset = c \\ \partial a(a^*) &= a^* \\ \partial a(a^+) &= \partial a(a.a^*) = \partial a(a).a^*|\epsilon(a).\partial a(a^*) = \lambda.a^*|\emptyset.a^* = a^*\end{aligned}$$

Dada una expresión regular, ésta denota un lenguaje regular, para el cual existe un autómata finito que lo reconoce.

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.

δ_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 Σ , definiendo la función δ :

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

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 :

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

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

δ_1	a	b		a	b
0	1	2	0	1	2
1	1	3	1	1	3
2	3		2	3	T
3	3		3	3	T
			T	T	T

Del lado derecho figura el autómata completado con el estado trampa, que es un estado que se puede agregar al autómata para completar las transiciones no definidas. Trampa es un estado no final, adonde van a parar las transiciones no definidas y que cicla sobre el mismo con cada símbolo de Σ^* .

Partiendo del estado q_0 , tenemos el lenguaje $L = L_0$ aceptado por el autómata, que está definido por la expresión regular $(a^+(ba^*|\lambda)|ba^+)$. Ahora, si en vez de partir del estado q_0 partiéramos del estado q_1 , tendríamos el lenguaje L_1 , que corresponde a la expresión regular a^*ba^* , y si partimos del estado q_2 tenemos el lenguaje L_2 con expresión regular a^+ . Análogamente, $L_3 = a^*$ y el estado trampa define el lenguaje vacío, ya que partiendo de él no se puede reconocer ninguna cadena, $L_T = \emptyset$.

Volviendo al autómata general definido más arriba, veamos que para un estado q_i dado ($1 \leq i \leq n$), y para cada símbolo $a_j \in \Sigma$ ($1 \leq j \leq m$), podemos armar el siguiente sistema de ecuaciones:

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

Basándonos en estas ecuaciones, podemos definir el lenguaje L_i de la siguiente forma:

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

En el autómata del ejemplo, estas ecuaciones corresponderían a:

$$\begin{aligned} L_0 &= a.L_1|b.L_2 \\ L_1 &= a.L_1|b.L_3|\lambda \\ L_2 &= a.L_3|b.L_T \\ L_3 &= a.L_3|\lambda \\ L_T &= a.L_T|b.L_T|\emptyset \end{aligned}$$

Podemos resolver este sistema de ecuaciones utilizando la siguiente propiedad de los lenguajes definidos sobre Σ : (dados R, S, T lenguajes sobre Σ)

$$\text{Si } R = S.R|T \wedge \lambda \notin S, \text{ entonces } R = S^*.T$$

Utilizando esta propiedad y otras de las expresiones regulares, despejamos las ecuaciones planteadas arriba:

$$\begin{aligned} L_T &= (a|b).L_T|\emptyset, \text{ y como } \lambda \notin \{a, b\}, \text{ entonces } L_T = (a|b)^*.\emptyset = \emptyset \\ L_3 &= a.L_3|\lambda, \text{ entonces } L_3 = a^*.\lambda = a^* \\ L_2 &= a.L_3|b.L_T = a.a^*|\emptyset = a^+ \\ L_1 &= a.L_1|b.L_3|\lambda = a.L_1|(b.a^*|\lambda), \text{ y como } \lambda \notin \{a\}, \text{ entonces } L_1 = a^*.(b.a^*|\lambda) \\ L_0 &= a.L_1|b.L_2 = a^+(b.a^*|\lambda)|ba^+ \end{aligned}$$

De esta manera, tenemos un método para, dado un autómata finito, hallar la expresión regular que define el lenguaje aceptado por el autómata. Veamos que utilizando estos mismos conceptos, podemos formular un método para, dada la expresión regular que define un lenguaje, encontrar el autómata finito que acepta ese lenguaje.

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

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

Reemplazando estos valores en (2) nos queda:

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

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 L_δ)

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

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

Con estas derivadas podemos construir el autómata correspondiente. Serán estados finales aquellos que estén representados por expresiones regulares que contengan la cadena nula. En este caso son finales L_1 y L_3 . El estado trampa corresponde a la expresión regular \emptyset .