

Espressioni e linguaggi regolari

Prof. A. Morzenti

aa 2008-2009

ESPRESSIONI E LINGUAGGI REGOLARI

I LINGUAGGI REGOLARI sono la più semplice famiglia di linguaggi formali. Può essere definita in molti modi diversi:

- In modo algebrico (noi cominceremo da qui)
- Con grammatiche generative
- Con algoritmi di riconoscimento

$$\Sigma = \{a_1, a_2, \dots, a_i\}$$
$$\cdot \quad \cup \quad *$$
$$\{a_1\}, \{a_2\}, \dots, \{a_i\} \quad \emptyset$$

ESPRESSIONE REGOLARE (e.r.): stringa r costruita con i caratteri dell'alfabeto Σ e con i metasimboli $\emptyset \cup \cdot$ con le seguenti regole (dove s e t sono e.r.):

$$\begin{array}{lll} 1. r = \emptyset & 3. r = (s \cup t) & 5. r = (s)^* \\ 2. r = a, a \in \Sigma & 4. r = (s.t) \text{ o } r = (st) \end{array}$$

NB: al posto di \cup spesso si usa $|$

PRECEDENZA OPERATORI: stella '*', concatenamento '.', unione 'U'

Si possono usare anche ε definito come $\varepsilon = \emptyset^*$
 e^+ definito come $e.e^*$

IL SIGNIFICATO DI UNA e.r. r è un linguaggio L_r di alfabeto Σ secondo la seguente tabella

UN LINGUAGGIO E' DETTO
REGOLARE se è denotato da
una espressione regolare

espressione r	linguaggio L_r
1. ε	$\{\varepsilon\}$
2. $a \in \Sigma$	$\{a\}$
3. $s \cup t$ o $s t$	$L_s \cup L_t$
4. $s.t$ o st	$L_s.L_t$
5. s^*	L_s^*

ESEMPIO 1: L_e sia il linguaggio che contiene le sequenze di segnali multiple di tre

$$\begin{aligned}
 e &= (111)^* \\
 L_e &= \{\varepsilon, 111, 111111, \dots\} = \{1^n \mid n \bmod 3 = 0\} \\
 e_1 &= 11(1)^* \quad L_e \neq L_{e_1} \\
 L_{e_1} &= \{11, 111, 1111, 11111, \dots\} = \{111^n \mid n \geq 0\}
 \end{aligned}$$

ESEMPIO 2: Sia $\Sigma=\{+,-,d\}$ con d che denota una cifra decimale $0,1,\dots,9$.
Si definisce l'e.r. e che produce il linguaggio dei numeri interi con o senza segno.

$$e = (+ \cup - \cup \varepsilon) d d^*$$

$$L_e = \{+, -, \varepsilon\} \{d\} \{d\}^*$$

ESEMPIO 3: Il linguaggio di alfabeto $\{a,b\}$ è tale che il numero dei caratteri a è dispari e vi è almeno un b

$$A_p b A_d \mid A_d b A_p$$

$$A_p = b^* (a b^* a b^*)^* \quad A_d = b^* a b^* (a b^* a b^*)^*$$

LA FAMIGLIA DEI LINGUAGGI REGOLARI (REG)

È la collezione di tutti i linguaggi regolari

LA FAMIGLIA DEI LINGUAGGI FINITI (FIN)

E' la collezione di tutti i linguaggi aventi cardinalità finita

OGNI LINGUAGGI FINITO è REGOLARE perché è l'unione
di un numero finito di stringhe e ciascuna stringa è il concatenamento
di un numero finito di caratteri

$$(x_1 \cup x_2 \cup \dots \cup x_k) = (a_{1_1} a_{1_2} \dots a_{1_n} \cup \dots \cup a_{k_1} a_{k_2} \dots a_{k_m})$$

La famiglia dei linguaggi regolari contiene anche linguaggi di cardinalità non finita,
quindi l'inclusione è stretta: $\text{FIN} \subset \text{REG}$

SOTTOESPRESSIONE DI UNA E.R. (S.E.)

1. Consideriamo una e.r completamente parentesizzata
2. Produciamo una versione numerata della stessa
3. Produciamo le sottoespressioni

$$e = (a \cup (bb))^* (c^+ \cup (a \cup (bb)))$$

$$e_N = (a_1 \cup (b_2 b_3))^* (c_4^+ \cup (a_5 \cup (b_6 b_7)))$$

$(a_1 \cup (b_2 b_3))^*$	$c_4^+ \cup (a_5 \cup (b_6 b_7))$
$a_1 \cup (b_2 b_3)$	$c_4^+ \quad a_5 \cup (b_6 b_7)$
$a_1 \quad b_2 b_3$	$c_4 \quad a_5 \quad b_6 b_7$
$\quad b_2 \quad b_3$	$\quad \quad b_6 \quad b_7$

SCELTE: Gli operatori di unione e di ripetizione presenti in una e.r. corrispondono a possibili scelte. Fissando una scelta si ottiene una e.r. che Definisce un linguaggio più piccolo.

$e_k, 1 \leq k \leq n$, è scelta dell'unione $e_1 \cup \dots \cup e_n$

$e^n, n \geq 1$, è scelta delle espressioni e^, e^+*

*ε è scelta dell'espressione e^**

Data una e.r. se ne può DERIVARE una seconda sostituendo al posto di una sottoespressione un'altra scelta da essa.

RELAZIONE DI DERIVAZIONE tra due e.r. e' ed e''

$e' \Rightarrow e''$ se le due e.r. si possono fattorizzare come

$$e' = \alpha\beta\gamma \quad e'' = \alpha\delta\gamma$$

dove: β è s.e. di e' , δ è s.e. di e'' , δ è una scelta di β

LA RELAZIONE DI DERIVAZIONE è applicabile più volte di seguito (cfr. definizione di **potenza** di una relazione, **chiusura transitiva e riflessiva**)

$$\overset{n}{e_0 \Rightarrow e_n} \text{ se } e_0 \Rightarrow e_1, e_1 \Rightarrow e_2, \dots, e_{n-1} \Rightarrow e_n$$

$$\overset{+}{e_0 \Rightarrow e_n} \quad e_0 \text{ deriva } e_n \text{ in } n \geq 1 \text{ passi}$$

$$\overset{*}{e_0 \Rightarrow e_n} \quad e_0 \text{ deriva } e_n \text{ in } n \geq 0 \text{ passi}$$

ESEMPI

Derivazioni immediate e derivazioni a più passi:

$$\begin{aligned} a^* \cup b^+ &\Rightarrow a^*, & a^* \cup b^+ &\Rightarrow b^+ \\ a^* \cup b^+ &\Rightarrow a^* \Rightarrow \varepsilon \text{ ossia } a^* \cup b^+ \xRightarrow{2} \varepsilon & \text{ o anche } a^* \cup b^+ \xRightarrow{+} \varepsilon \\ a^* \cup b^+ &\Rightarrow b^+ \Rightarrow bbb \text{ ossia } a^* \cup b^+ \xRightarrow{2} bbb & \text{ o anche } a^* \cup b^+ \xRightarrow{+} bbb \end{aligned}$$

Tra le e.r. ottenute per derivazione da una e.r. alcune contengono anche i metasimboli (operatori e parentesi), altre soltanto i simboli di Σ (detti anche **terminali**) e la ε .

Queste ultime costituiscono il linguaggio definito dalla e.r.

IL LINGUAGGIO DEFINITO DA UNA E.R. è

$$L(r) = \left\{ x \in \Sigma^* \mid r \xRightarrow{*} x \right\}$$

Due e.r sono dette EQUIVALENTI se definiscono lo stesso linguaggio

IL LINGUAGGIO DELLA E.R. DERIVATA E' INCLUSO IN QUELLO DELLA E.R. DERIVANTE

Ci possono essere molti ORDINI DISTINTI ma EQUIVALENTI per produrre una frase del linguaggio

ESEMPIO:

$$1.(ab)^* \Rightarrow abab$$

$$2.(ab \cup c) \Rightarrow ab$$

$$3.a(ba \cup c)^* d \Rightarrow ad$$

$$4.a(ba \cup c)^* d \Rightarrow a(ba \cup c)(ba \cup c)d$$

$$5.a^*(b \cup c \cup d)f^+ \Rightarrow aaa(b \cup c \cup d)f^+$$

$$6.a^*(b \cup c \cup d)f^+ \Rightarrow a^*cf^+$$

$$7.a^*(b \cup c \cup d)f^+ \Rightarrow aaacf^+ \text{ in 2 passi}$$

$$8.a^*(b \cup c \cup d)f^+ \Rightarrow aaacff^+ \text{ in 3 passi}$$

AMBIGUITA' DELLE ESPRESSIONI REGOLARI

Una frase può essere ottenuta con derivazioni che differiscono in modo più sostanziale del solo ordine di derivazione.

$$(a \cup b)^* a (a \cup b)^*$$

$$(a \cup b)^* a (a \cup b)^* \Rightarrow (a \cup b) a (a \cup b)^* \Rightarrow aa (a \cup b)^* \Rightarrow aa \varepsilon \Rightarrow aa$$

$$(a \cup b)^* a (a \cup b)^* \Rightarrow \varepsilon a (a \cup b)^* \Rightarrow \varepsilon a (a \cup b) \Rightarrow \varepsilon aa \Rightarrow aa$$

Una e.r. f è AMBIGUA se e solo se, nel linguaggio definito dalla corrispondente e.r. marcata f' vi sono due diverse stringhe x e y , tali che, cancellando i numeri, esse vengono a coincidere.

$$f' = (a_1 \cup b_2)^* a_3 (a_4 \cup b_5)^*$$

definisce un linguaggio regolare di alfabeto

$$\{a_1, b_2, a_3, a_4, b_5\}$$

$a_1 a_3$ e $a_3 a_4$ del ling. di f' mostrano l'ambiguità

ESEMPIO (ambiguità)

$(aa \mid ba)^* a \mid b(aa \mid ba)^*$ è ambigua

infatti, marcando

$(a_1a_2 \mid b_3a_4)^* a_5 \mid b_6(a_7a_8 \mid b_9a_{10})^*$

si ottengono

$b_3a_4a_5 \quad b_6a_7a_8$

che si proiettano in modo ambiguo nella frase baa

APPLICAZIONE:

numeri reali in virgola (punto) mobile, con o senza segno ed esponente

$$\Sigma = \{+, -, \bullet, E, d\}$$

$$r = s.c.e$$

$s = (+ \cup - \cup \varepsilon)$ apporta il segno \pm opzionale

$c = (d^+ \bullet d^* \cup d^* \bullet d^+)$ genera costanti intere o frazionarie senza segno

$e = (\varepsilon \cup E(+ \cup - \cup \varepsilon)d^+)$ genera l'esponente opzionale preceduto da E

$$(+ \cup - \cup \varepsilon)(d^+ \bullet d^* \cup d^* \bullet d^+)(\varepsilon \cup E(+ \cup - \cup \varepsilon)d^+)$$

$$+ dd \bullet E - ddd \quad + 12 \bullet E - 341 \quad 12.10^{-341}$$

ALTRI OPERATORI

POTENZA: $a^h = aa...a$ (h volte)

RIPETIZIONE: da k a n > k: $[a]_k^n = a^k \cup a^{k+1} \cup ... \cup a^n$

OPZIONALITA': $[a] = (\varepsilon \cup a)$

INTERVALLO ORDINATO (0...9) (a...z) (A...Z)

Operazioni insiemistiche INTERSEZIONE, DIFFERENZA, COMPLEMENTO

Si parla di ESPRESSIONI REGOLARI ESTESE con operatori insiemistici

Si dimostra (studiando relazione con automi finiti cfr. cap. 3 testo) che op. insiemistiche non aumentano potenza espressiva delle e.r.

(sono solo una **comoda abbreviazione**)

INTERSEZIONE: comoda per definire linguaggi mediante congiunzione di condizioni

ESEMPIO: linguaggio $L \subseteq \{a,b\}^*$ delle stringhe di lunghezza pari e contenenti bb

Facile definire usando un'e.r. e con intersezione tale che $L = L_e$:

$$e = ((a \mid b)^* bb (a \mid b)^*) \cap ((a \mid b)^2)^*$$

frasi contenenti bb fr. Lungh. pari

Senza l'intersezione:

***bb* circondata da due stringhe pari o da due dispari**

$$((a \mid b)^2)^* bb ((a \mid b)^2)^* \mid (a \mid b) ((a \mid b)^2)^* bb (a \mid b) ((a \mid b)^2)^*$$

ESEMPIO (uso di e.r. estesa con operatore di complemento)

Definiamo linguaggio $L \subset \{a,b\}^*$ delle stringhe che **NON** contengono sottostringa aa

Facile definire il linguaggio complemento $\neg L = \{x \in (a|b)^* \mid x \text{ contiene la sottostringa aa}\}$

$$\neg L = ((a|b)^* aa (a|b)^*)$$

Quindi L definibile con e.r. estesa

$$L = \neg((a|b)^* aa (a|b)^*)$$

Definizione mediante e.r. **NON** estesa (forse meno leggibile)

$$L = (ab \mid b)^* (a \mid \varepsilon)$$

CHIUSURA DELLA FAMIGLIA REG RISPETTO ALLE OPERAZIONI / 1

Sia θ un operatore che applicato a un linguaggio o a una coppia di linguaggi
Ne produce un altro.

UNA FAMIGLIA DI LINGUAGGI SI DICE CHIUSA RISPETTO AD UN
OPERATORE θ SE IL LINGUAGGIO RISULTANTE DALL'APPLICAZIONE DI
 θ AI LINGUAGGI DELLA FAMIGLIA APPARTIENE ANCORA ALLA FAMIGLIA

PROPRIETA'. La famiglia REG dei linguaggi regolari è CHIUSA RISPETTO
AGLI OPERATORI di concatenamento, unione, stella (e quindi anche rispetto
agli operatori derivati croce ed elevamento a potenza)
E' un'ovvia conseguenza della definizione stessa di espressione regolare)

Quindi i linguaggi regolari possono essere combinati tra loro con detti operatori
senza pericolo di uscire dalla famiglia dei linguaggi definibili con e.r.

CHIUSURA DELLA FAMIGLIA REG RISPETTO ALLE OPERAZIONI / 2

PROPRIETA' PIU' FORTE. La famiglia REG dei linguaggi regolari è **la più piccola** famiglia di linguaggi che

(i) contiene tutti i linguaggi finiti ed

(ii) è chiusa rispetto a concatenamento, unione, stella.

Si puo` dare una facile dimostrazione per assurdo (cfr. testo §2.3.4)

REG è anche chiusa rispetto a INTERSEZIONE, COMPLEMENTO e RIFLESSIONE (useremo la teoria degli automi per dimostrarlo)

ASTRAZIONE LINGUISTICA

L'astrazione linguistica trasforma le frasi di un linguaggio reale, detto concreto, in una forma più semplice, detta rappresentazione astratta.

Essa trascura i simboli dell'alfabeto concreto e impiega al loro posto i caratteri di un altro alfabeto, quello astratto.

AL LIVELLO ASTRATTO, LE STRUTTURE DEI LINGUAGGI ARTIFICIALI SI POSSONO OTTENERE COME COMPOSIZIONE DI POCHI PARADIGMI ELEMENTARI, ATTRAVERSO LE OPERAZIONI DI CONCATENAMENTO, UNIONE E ITERAZIONE.

Dal linguaggio astratto a quello effettivo → scelta degli elementi lessicali (parole chiave, delimitatori, identificatori, ...)

LA COSTRUZIONE DEI COMPILATORI FA RIFERIMENTO ALLA STRUTTURA ASTRATTA

I linguaggi artificiali fanno uso di (poche) strutture astratte ricorrenti.
Tra queste, LE LISTE sono descrivibili con ESPRESSIONI REGOLARI.

LISTE ASTRATTE E CONCRETE

Una lista contiene un numero imprecisato di elementi e dello stesso tipo.
Essa è generata dalla e.r. e^+ o da e^* , se gli elementi possono mancare.

e può essere simbolo terminale o altro (stringa di un altro linguaggio formale, ...)

LISTE CON SEPARATORI E MARCHE DI APERTURA E CHIUSURA

ESEMPI

$$ie(se)^* f$$
$$i[e(se)^*]f$$
$$\overbrace{\text{begin}}^i \overbrace{\text{istr}_1}^e ; \overbrace{\text{istr}_2}^e ; \dots ; \overbrace{\text{istr}_n}^e \overbrace{\text{end}}^f$$
$$\overbrace{\text{procedure STAMPA}}^i (\overbrace{\text{par}_1}^e , \overbrace{\text{par}_2}^e , \dots , \overbrace{\text{par}_n}^e)$$
$$\overbrace{\text{array MATRICE}}^i [\overbrace{\text{int}_1}^e , \overbrace{\text{int}_2}^e , \dots , \overbrace{\text{int}_n}^e]$$

LISTE A PRECEDENZE O LIVELLI

Un elemento di una lista può essere

Una lista di livello inferiore

NB: il numero di livelli è limitato,
altrimenti servono notazioni più
potenti (grammatiche)

ESEMPI

livello 1: *begin istr₁;istr₂;...;istr_n end*

livello 2: *STAMPA(var₁, var₂, ..., var_n)*

$$lista_1 = i_1 lista_2 (s_1 lista_2)^* f_1$$

$$lista_2 = i_2 lista_3 (s_2 lista_3)^* f_2$$

...

$$lista_k = i_k e_k (s_k e_k)^* f_k$$

$$3 + \underbrace{5 \times 7 \times 4}_{\text{monomio1}} - \underbrace{8 \times 2 \div 5}_{\text{monomio2}} + 8 + 3$$

padre, madre, figlio e figlia

un padre forte, severo e giusto, una madre amorevole e fedele

un libro come lista di capitoli separati da pagine bianche,
chiusa tra due copertine

un capitolo come lista di sezioni

Si passa dalla forma astratta di un costrutto a quella concreta mediante SOSTITUZIONE, operazione semplice che rimpiazza un carattere terminale $b \in \Sigma$ di una stringa x di un primo linguaggio (detto sorgente) con le frasi di un secondo linguaggio $L_b \subseteq \Delta^*$ (detto pozzo)

La sostituzione del linguaggio L_b al posto di b nella stringa $x = a_1 a_2 \dots a_n$ produce un linguaggio di alfabeto $(\Sigma \setminus \{b\}) \cup \Delta$ così definito:

$$\{y \mid y = y_1 y_2 \dots y_n \wedge (\text{if } a_i \neq b \text{ then } y_i = a_i \text{ else } y_i \in L_b)\}$$