

ALGORITMO DI EARLEY

COSTRUZIONE DI UN ANALIZZATORE SINTATTICO: L'ALGORITMO DI EARLEY

In questa lezione faremo vedere quali sono i passi che si devono compiere per costruire un analizzatore sintattico che verifica se una stringa appartiene a un dato linguaggio, espresso sotto forma di CFG.

In particolare **l'algoritmo di Earley** costruisce un analizzatore sintattico a partire da **una grammatica context free generica**.

Vedremo in seguito altri algoritmi che permettono di costruire parser a partire da alcune grammatiche context free particolari. La scelta di tali grammatiche permette una costruzione più efficiente e un riconoscimento più intuitivo

L'ALGORITMO DI EARLEY

L'idea è quella di costruire progressivamente tutte le possibili derivazioni leftmost (o rightmost) compatibili con la stringa in input.

Durante il procedimento si analizza la stringa in input da sinistra verso destra scartando via via le derivazioni in cui non vi sia corrispondenza tra i simboli derivati e quelli della stringa.

Se esistono due derivazioni leftmost (o rightmost) possibili, l'algoritmo restituisce i due alberi di derivazione.

La complessità di calcolo è proporzionale al cubo della lunghezza della stringa da analizzare e si riduce al quadrato se la grammatica non è ambigua e ancora di più se è deterministica.

IN DETTAGLIO

Data una stringa $x_1x_2\dots x_n$, l'algoritmo la scandisce da sinistra verso destra e per ogni x_i costruisce gli stati S_j . Uno stato è costituito da un insieme di coppie

(dotted_rule, puntatore).

Una **dotted rule** è una produzione di G avente sul lato destro un punto che ne marca una posizione.

Il **puntatore** è un intero che indica la posizione dell'input a partire dalla quale è iniziato l'esame della produzione contenuta nella dotted rule.

In altre parole un elemento dello stato S_j è del tipo:

$(A \rightarrow \alpha.\beta, i)$ con $0 \leq i \leq j$, $\alpha, \beta \in (N \cup T)^*$, $A \rightarrow \alpha\beta$ regola della grammatica

Ciò significa che:

- Si è iniziato l'esame della produzione $A \rightarrow \alpha\beta$ a partire dalla posizione $i+1$ dell'input;
- È già stata esaminata la parte α che precede il punto;
- È già stato verificato che α genera $x_{i+1} \dots x_j$

L'ALGORITMO

Per semplicità si aggiunge il simbolo \$ alla fine della stringa x e la produzione $S' \rightarrow S\$$.

input $x = x_1 \dots x_n$

$x_{n+1} = \$$

$S[0] = \{ (S' \rightarrow .S\$, 0) \}$ /*stato iniziale

for $j=0$ to n do

 elabora ogni coppia $(A \rightarrow \alpha.\beta, i)$ di $S[j]$ applicando una delle

 3 operazioni: scansione, predizione, completamento.

if $S[n+1] = \{ (S' \rightarrow S\$. , 0) \}$ then

 accetta

else rifiuta.

SCANSIONE (SCANNER)

Esaminiamo uno stato $(A \rightarrow \alpha . \beta, i) \in S_j$

SCANSIONE:

Si applica quando β inizia con un terminale a , cioè $\beta = a\beta'$ (cioè $A \rightarrow \alpha . a\beta'$), e $a = x_{j+1}$, allora si aggiunge la coppia $(A \rightarrow \alpha a . \beta', i)$ a S_{j+1} (cioè lo stato successivo)

La regola (che era stata applicata all' i -esimo passo) ha quindi riconosciuto il carattere $j+1$ -esimo, lo stato successivo si dovrà preoccupare di esaminare il prossimo carattere della produzione.

Se il carattere a è diverso dal carattere x_{j+1} letto nell'input, questa produzione fallisce e non va più avanti.

PREDIZIONE (PREDICTOR)

Esaminiamo $(A \rightarrow \alpha . \beta, i) \in S_j$

PREDIZIONE:

Si applica quando β inizia con un non terminale B , cioè $\beta = B\beta'$, allora, per ogni produzione $B \rightarrow \gamma$ della grammatica, aggiungi la coppia $(B \rightarrow . \gamma, j)$ a S_j (n.b. stesso stato in cui si trova $(A \rightarrow \alpha . \beta, i)$)

Significa che si predice l'espansione di B a partire dalla posizione j

COMPLETAMENTO (COMPLETED)

$$(A \rightarrow \alpha . \beta, i) \in S_j$$

COMPLETAMENTO:

Si applica quando $\beta = \epsilon$, $(A \rightarrow \alpha ., i)$ allora per ogni coppia $(C \rightarrow \eta . A \delta, h)$ di S_i (ossia ogni coppia in cui A appare a destra del punto nello stato i in cui è iniziata la predizione di A) si aggiunge $(C \rightarrow \eta A . \delta, h)$ a S_j

Ciò significa che la predizione $A \rightarrow . \alpha$ iniziata allo stato i ha avuto successo, quindi possiamo continuare ad analizzare le produzioni da cui questa è derivata (quelle del tipo $(C \rightarrow \eta . A \delta, h)$), a partire dal simbolo successivo. Si individuano allora in S_i tutti gli stati che avevano attivato la predizione $A \rightarrow . \alpha$

Significa che $\eta \Rightarrow^* x_{h+1} \dots x_i$, in più $A \rightarrow \alpha \Rightarrow^* x_{i+1} \dots x_j$, allora possiamo concludere che $\eta \alpha \Rightarrow^* x_{h+1} \dots x_j$

La derivazione (e quindi l'albero di derivazione) può essere ricostruita a ritroso.

Applicazione dell'algoritmo di Earley al riconoscimento di (a)\$

$S' \rightarrow E\$$

$E \rightarrow E + E$

$E \rightarrow E * E$

$E \rightarrow (E)$

$E \rightarrow a$

• **scanner** se β inizia con un terminale a , cioè $\beta = a\beta'$, allora se $a = x_{j+1}$, aggiungi la coppia $(A \rightarrow \alpha a \cdot \beta', i)$ a S_{j+1} .

• **predictor** : se β inizia con un non terminale B , cioè $\beta = B\beta'$, allora, per ogni produzione $B \rightarrow \gamma$, aggiungi la coppia $(B \rightarrow \cdot \gamma, j)$ a $S[j]$

• **completer** se $\beta = \epsilon$ allora data la coppia $(A \rightarrow \alpha \cdot, i)$, per ogni coppia $(C \rightarrow \eta \cdot A \delta, h)$ di S_i si aggiunge $(C \rightarrow \eta A \cdot \delta, h)$ a S_j

(a) \$

	S_1	S_2	S_3	S_4
$S' \rightarrow \cdot E \$, 0 \cdot$	$E \rightarrow (\cdot E), 0 \cdot$	$E \rightarrow a \cdot, 1 \cdot$	$E \rightarrow (E) \cdot, 0 \cdot$	$S' \rightarrow E \$ \cdot, 0$
$E \rightarrow \cdot E + E, 0 \cdot$	$E \rightarrow \cdot E + E, 1 \cdot$	$E \rightarrow (E \cdot), 0 \cdot$	$S' \rightarrow E \cdot \$, 0 \cdot$	
$E \rightarrow \cdot E * E, 0 \cdot$	$E \rightarrow \cdot E * E, 1 \cdot$	$E \rightarrow E \cdot + E, 1 \cdot$	$E \rightarrow E \cdot + E, 0 \cdot$	
$E \rightarrow \cdot (E), 0 \cdot$	$E \rightarrow \cdot (E), 1 \cdot$	$E \rightarrow E \cdot * E, 1 \cdot$	$E \rightarrow E \cdot * E, 0 \cdot$	
$E \rightarrow \cdot a, 0 \cdot$	$E \rightarrow \cdot a, 1 \cdot$			

Nell'operazione di scansione si inserisce un puntatore tra un insieme e un altro.
 Nell'operazione di completamento e di predizione all'interno dello stesso insieme.

Applicazione dell'algoritmo di Earley al riconoscimento di (a)\$

S_0	S_1	S_2	S_3	S_4
$S' \rightarrow .E\$, 0 \bullet$	$E \rightarrow (.E), 0 \bullet$	$E \rightarrow a., 1 \bullet$	$E \rightarrow (E)., 0 \bullet$	$S' \rightarrow E\$. , 0$
$E \rightarrow .E+E, 0 \bullet$	$E \rightarrow .E+E, 1 \bullet$	$E \rightarrow (E.), 0 \bullet$	$S' \rightarrow E.\$, 0 \bullet$	
$E \rightarrow .E * E, 0 \bullet$	$E \rightarrow .E * E, 1 \bullet$	$E \rightarrow E.+E, 1 \bullet$	$E \rightarrow E.+E, 0 \bullet$	
$E \rightarrow .(E), 0 \bullet$	$E \rightarrow .(E), 1 \bullet$	$E \rightarrow E.*E, 1 \bullet$	$E \rightarrow E.*E, 0 \bullet$	
$E \rightarrow .a, 0 \bullet$	$E \rightarrow .a, 1 \bullet$			

$S' \rightarrow E\$$

$E \rightarrow E+E$

$E \rightarrow E * E$

$E \rightarrow (E)$

$E \rightarrow a$

RICOSTRUZIONE DELLA DERIVAZIONE

A partire dallo stato S_4 , vedo da dove ho ricavato $S' \rightarrow E\$. , 0$ ossia $S' \rightarrow E.\$, 0$ in S_3 . Questo a sua volta proviene da $E \rightarrow (E)., 0$. Questo proviene da $E \rightarrow (E.), 0$ in S_2 , che a sua volta proviene da $E \rightarrow a., 1$

$S' \rightarrow E\$ \rightarrow (E) \$ \rightarrow (a) \$$

Esercizio: provare il parser di a+a+a

	a	+	a	+	a	\$
S_0	S_1	S_2	S_3	S_4	S_5	S_6
$S' \rightarrow \cdot E \$, 0$	$E \rightarrow a \cdot, 0$	$E \rightarrow E + \cdot E, 0$	$E \rightarrow a \cdot, 2$	$E \rightarrow E + \cdot E, 2$	$E \rightarrow a \cdot, 4$	$S' \rightarrow E \$ \cdot, 0$
$E \rightarrow \cdot E + E, 0$	$S' \rightarrow E \cdot \$, 0$	$E \rightarrow \cdot E + E, 2$	$E \rightarrow E + E \cdot, 0$	$E \rightarrow E + \cdot E, 0$	$E \rightarrow E + E \cdot, 2$	
$E \rightarrow \cdot E * E, 0$	$E \rightarrow E \cdot + E, 0$	$E \rightarrow \cdot E * E, 2$	$E \rightarrow E \cdot + E, 2$	$E \rightarrow \cdot E + E, 4$	$E \rightarrow E + E \cdot, 0$	
$E \rightarrow \cdot (E), 0$	$E \rightarrow E \cdot * E, 0$	$E \rightarrow \cdot (E), 2$	$E \rightarrow E \cdot * E, 2$	$E \rightarrow \cdot E * E, 4$	$E \rightarrow E \cdot + E, 4$	
$E \rightarrow \cdot a, 0$		$E \rightarrow \cdot a, 2$	$S' \rightarrow E \cdot \$, 0$	$E \rightarrow \cdot (E), 4$	$E \rightarrow E \cdot * E, 4$	
			$E \rightarrow E \cdot + E, 0$	$E \rightarrow \cdot a, 4$	$E \rightarrow E \cdot + E, 2$	
			$E \rightarrow E \cdot * E, 0$		$E \rightarrow E \cdot * E, 2$	
					$S' \rightarrow E \cdot \$, 0$	
					$E \rightarrow E \cdot + E, 0$	
					$E \rightarrow E \cdot * E, 0$	

$S' \rightarrow E \$$

$E \rightarrow E * E$

$E \rightarrow E + E$

$E \rightarrow (E)$

$E \rightarrow a$

COMPLESSITÀ DELL'ALGORITMO

Ciascun insieme S_j può avere un numero di coppie che cresce linearmente con j , quindi $O(n)$;

Le operazioni di scansione e predizione su ogni coppia sono indipendenti da n ;

L'operazione di completamento richiede $O(j)$ per ogni coppia, quindi in totale $O(n^2)$

Sommando i passi per ogni i si ha $O(n^3)$.

In pratica l'algoritmo è più veloce: per molte Grammatiche è $O(n)$ e per ogni grammatica non ambigua è $O(n^2)$ Troppo per l'analisi sintattica di un Compilatore!

ESERCIZIO

Eseguire l'algoritmo di Earley per la grammatica

$S \rightarrow A \mid B$

$A \rightarrow aAb \mid ab$

$B \rightarrow aaBb \mid aab$

e la stringa $aabb$ e poi per la stringa $aaabb$

ESERCIZIO

Siano date la grammatica G e la stringa **aaaa**. Si esegua il riconoscimento della stringa utilizzando l'algoritmo di Earley.

$G: S \rightarrow \text{aaS} \mid \text{Saaa} \mid a \mid \varepsilon$

Nota: se ci sono produzioni $A \rightarrow \varepsilon$, si effettua subito il completamento

(La dot_rule è $A \rightarrow .$)

ESERCIZIO

Siano date la grammatica G e la stringa $aabb$. Si esegua il riconoscimento della stringa utilizzando l'algoritmo di Earley.

$S \rightarrow aAbB \mid C$

$A \rightarrow aA \mid a$

$B \rightarrow bB \mid b$

$C \rightarrow aCb \mid ab$