Parsificazione I

(analisi top-down)

Analisi sintattica I

Data una grammatica G l'analizzatore sintattico o parsificatore legge la stringa sorgente e se appartiene al linguaggio L(G) ne produce una derivazione o un albero sintattico, altrimenti si ferma segnalando l'errore.

Due classi importanti di analizzatori

Discendenti o top-down

Ascendenti o bottom-up

Come per l'analisi lessicale, anche per l'analisi sintattica sono stati sviluppati strumenti per la generazione automatica di parsificatori, sia per l'analisi bottom-up (Yacc) sia per quella top-down (Antlr).

Analisi sintattica deterministica

Un analizzatore sintattico deterministico leggendo 1 (o più) caratteri in input può eliminare le ambiguità e scegliere sempre la strada giusta che porta al riconoscimento della stringa. Ovviamente il modello su cui i rifanno tutti gli analizzatori sintattici (più o meno fedelmente) è l'automa a pila.

Le grammatiche che permettono parsing predittivo discendente sono chiamate LL(k), quelle che permettono parsing predittivo ascendente sono chiamate LR(k), dove k è il numero di simboli necessari per individuare la produzione senza ambiguità.

La famiglia **LL(k)** contiene tutti e soli i linguaggi che possono essere definiti da una grammatica LL(k) per un valore finito di k >= 1. Non tutti i linguaggi che hanno riconoscitori deterministici sono generabili da grammatiche LL(k), cioè la famiglia dei linguaggi LL(k) è strettamente contenuta nella famiglia dei linguaggi che hanno riconoscitori deterministici.

Parser LL(1) (iterative)

Un parser LL(1) iterativo per una grammatica LL(1) (definita esattamente in seguito) è costituito da:

- la stringa in input cui viene aggiunto un mark di fine stringa che denotiamo con \$. L'input viene letto sequenzialmente come negli automi.
- •una *pila* i cui elementi possono essere *terminali* o *non terminali* della grammatica

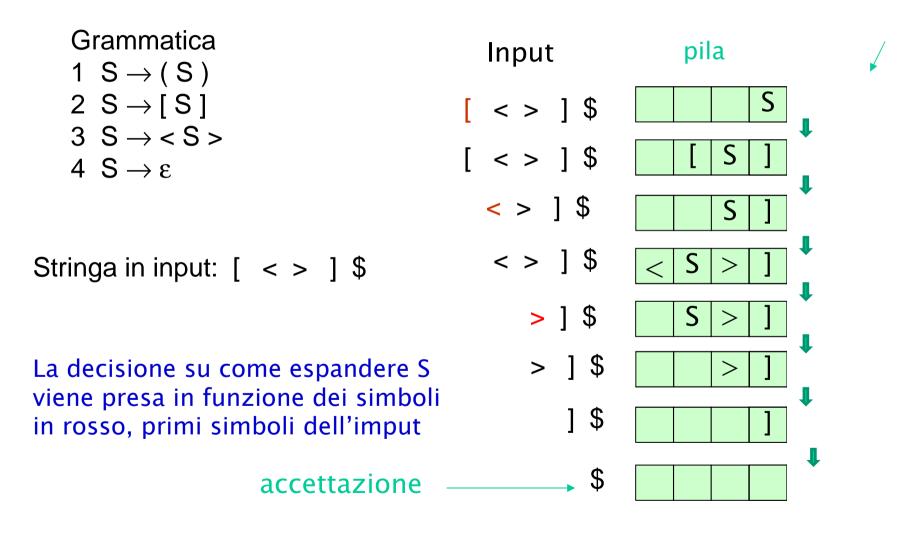
Durante l'analisi input e pila vengono modificati esattamente come l'input e la pila rappresentati nelle configurazioni istantanee di un automa a pila (non deterministico) che riconosce la grammatica G costruita nel modo standard.

Ma nell'algoritmo di parsificazione LL(1) la *decisione* su come espandere i nonterminali in cima alla pila viene presa in modo deterministico *guardando il primo simbolo dell'input (lookahead).*

Quando invece in cima alla pila si trova un terminale basterà verificare la corrispondenza di questo col simbolo in input e avanzare la testina di lettura (come nell'automa a pila).

Parsing LL(1) (top-down o discendente

Consideriamo per esempio la seguente grammatica (LL(1)). Notare l'uso di "\$" per definire la fine stringa:



FIRST

Data una grammatica $G = \langle V, \Sigma, P, S \rangle$, l'insieme FIRST di una stringa α di variabili e terminali, è definito formalmente come:

FIRST
$$(\alpha) = \{a \mid \alpha \rightarrow^* a\beta\} \cup \{\epsilon \mid se \alpha \rightarrow^* \epsilon\}$$

E' l'insieme dei terminali con cui iniziano le stringhe derivabili da α nella grammatica G. **FIRST**(α) (abb. F(α)) soddisfa questa definizione (ricorsiva):

1.
$$F(\varepsilon) = \{\varepsilon\}$$

$$F(A) = \bigcup_{A \to \gamma_i \in P} F(\gamma_i)$$

N.B. A è annullabile se A \rightarrow * ϵ

N.B. $(A \rightarrow \gamma_1 | \gamma_2 | ... | \gamma_k \hat{e})$ l'insieme delle produzioni di A in G)

$$G = \langle \{X, Y, Z\}, \{a, c, d\}, P, Z \rangle$$

$$P: \quad Z \rightarrow d \mid XYZ$$

$$Y \rightarrow c \mid \varepsilon$$

$$X \rightarrow Y \mid a$$

$$FIRST(d) = \{d\}$$

$$FIRST(XYZ) = \{a, c, d\}$$

$$Infatti \quad Z \rightarrow XYZ \rightarrow aYZ$$

$$Z \rightarrow XYZ \rightarrow YZ \rightarrow cZ$$

$$Z \rightarrow XYZ \rightarrow YZ \rightarrow Z \rightarrow d$$

$$FIRST\{X\} = \{a, c, \varepsilon\}$$

$$Infatti \quad X \rightarrow a$$

$$X \rightarrow Y \rightarrow c$$

$$X \rightarrow Y \rightarrow \varepsilon$$

FOLLOW

Data una grammatica $G = \langle V, \Sigma, P, S \rangle$, l'insieme **FOLLOW** (insieme dei <u>seguiti</u>) di una variabile A è l'insieme dei terminali con cui iniziano le stringhe che seguono A nelle derivazioni della grammatica G (assumendo \$ in fine stringa). Formalmente:

FOLLOW(A) = {a |
$$S \rightarrow^* \alpha Aa\beta$$
} \cup {\$ |se $S \rightarrow^* \alpha A$ }

Notare che \$ appartiene <u>sempre</u> al Fw dell'assioma.

FOLLOW(A) (abb. Fw(A)) soddisfa la seguente equazione:

$$Fw(A) = \left[\bigcup_{B \to \alpha A\beta \in P} (F(\beta) - \{ E \}) \right] \cup \left[\bigcup_{B \to \alpha A\beta \in P \text{ tali che} \atop \beta \text{ annullabile e B} \neq A} Fw(B) \right] \cup \left[\bigcup_{B \to \alpha A\beta \in P \text{ tali che} \atop \beta \text{ annullabile e B} \neq A} Fw(B) \right] \cup \left[\bigcup_{B \to \alpha A\beta \in P \text{ tali che} \atop \beta \text{ annullabile e B} \neq A} Fw(B) \right] \cup \left[\bigcup_{B \to \alpha A\beta \in P \text{ tali che} \atop \beta \text{ annullabile e B} \neq A} Fw(B) \right] \cup \left[\bigcup_{B \to \alpha A\beta \in P \text{ tali che} \atop \beta \text{ annullabile e B} \neq A} Fw(B) \right] \cup \left[\bigcup_{B \to \alpha A\beta \in P \text{ tali che} \atop \beta \text{ annullabile e B} \neq A} Fw(B) \right] \cup \left[\bigcup_{B \to \alpha A\beta \in P \text{ tali che} \atop \beta \text{ annullabile e B} \neq A} Fw(B) \right] \cup \left[\bigcup_{B \to \alpha A\beta \in P \text{ tali che} \atop \beta \text{ annullabile e B} \neq A} Fw(B) \right] \cup \left[\bigcup_{B \to \alpha A\beta \in P \text{ tali che} \atop \beta \text{ annullabile e B} \neq A} Fw(B) \right] \cup \left[\bigcup_{B \to \alpha A\beta \in P \text{ tali che} \atop \beta \text{ annullabile e B} \neq A} Fw(B) \right] \cup \left[\bigcup_{B \to \alpha A\beta \in P \text{ tali che} \atop \beta \text{ annullabile e B} \neq A} Fw(B) \right] \cup \left[\bigcup_{B \to \alpha A\beta \in P \text{ tali che} \atop \beta \text{ annullabile e B} \neq A} Fw(B) \right]$$

 \cup {\$} se A è lo start symbol di G

$$Z \rightarrow d \mid XYZ$$

 $Y \rightarrow c \mid \epsilon$
 $X \rightarrow Y \mid a$

FOLLOW(Y) = {a, d, c}
Infatti:
$$Z \rightarrow XYZ \rightarrow XYXYZ \rightarrow XYaYZ$$

 $Z \rightarrow XYZ \rightarrow XYd$
 $Z \rightarrow XYZ \rightarrow YYZ \rightarrow YcZ$
FOLLOW(X) = {c, d, a}
Infatti: $Z \rightarrow XYZ \rightarrow XcZ$
 $Z \rightarrow XYZ \rightarrow XYd \rightarrow Xd$
 $Z \rightarrow XYZ \rightarrow XYZ \rightarrow XXYZ \rightarrow XaYZ$

 $FOLLOW(Z) = \{\$\}$ perchè Z è l'assioma.

Insiemi guida

Data una grammatica G, **l'insieme guida** di una produzione della grammatica $A \to \alpha$ - **Gui** $(A \to \alpha)$ - è l'insieme dei terminali (o ϵ se ci si trova a fine parola) con cui *iniziano* le stringhe generabili a partire dalla produzione stessa :

Gui (A
$$\rightarrow \alpha$$
) = {a | S $\rightarrow_{\underline{c}}^*$ wA $\beta \rightarrow$ w $\alpha \beta \rightarrow_{\underline{lm}}$ wa γ }

L'insieme Gui (A $\rightarrow \alpha$) si può esprimere usando F() e FW():

$$Gui(A \to \alpha) = \begin{cases} F(\alpha) & \text{se } \alpha \text{ non è annullabile} \\ \\ (F(\alpha) - \{\epsilon\}) \cup FW(A) & \text{se } \alpha \text{ è annullabile} \end{cases}$$

Esempio:

$$\begin{array}{ll} \text{Gui } (Z \rightarrow d) = \{d\} & \text{Gui } (Y \rightarrow \epsilon) = \{a,\,c,\,d\} \\ \text{Gui } (Z \rightarrow XYZ) = \{a,\,c,\,d\} & \text{Gui } (X \rightarrow Y) = \{a,\,c,\,d\} \\ \text{Gui } (Y \rightarrow c) = \{c\} & \text{Gui } (X \rightarrow a) = \{a\} \end{array}$$

Grammatiche LL(1)

Una **grammatica** è **LL(1)** se *per ogni non terminale* A e *per ogni coppia* di produzioni $A \rightarrow \alpha$ e $A \rightarrow \beta$, gli insiemi guida sono disgiunti:

$$Gui(A \rightarrow \alpha) \cap Gui(A \rightarrow \beta) = \Phi$$

Esempio:

La grammatica $\{S\}$, $\{(, [, <,),], >\}$, $P = \{S \rightarrow (S) \mid [S] \mid <S> \mid \epsilon\}$, $S> \in LL(1)$.

 $Gui(S \rightarrow \varepsilon) = \{\$, \}, \}$

$$F(S) = \{(, [, <\} \\ Gui(S \rightarrow (S)) = \{(\} \\ Gui(S \rightarrow [S]) = \{[\} \\ FW(S) = \{\$, \}, \}, \}\}$$

$$Gui(S \rightarrow (S)) = \{(\} \\ Gui(S \rightarrow$$

Insiemi GUIDA: esempio di calcolo

$$\begin{split} Z \to d \mid XYZ \\ Y \to c \mid & \epsilon \\ X \to Y \mid a \\ \\ F(Z) &= \{a, c, d\} \quad F(X) = \{a, c, \epsilon\} \quad F(Y) = \{c, \epsilon\} \\ Fw(Z) &= \{\$\}, \quad Fw(X) = \{a, c, d\}, \quad Fw(Y) = \{a, c, d\} \\ Gui (Z \to d) &= F(d) = \{d\} \\ Gui (Z \to XYZ) &= F(XYZ) = (F(X) - \{\epsilon\}) \cup F(YZ) = \\ &= (F(X) - \{\epsilon\}) \cup (F(Y) - \{\epsilon\}) \cup F(Z) = \{a, c, d\} \\ Gui (Y \to c) &= \{c\} \\ Gui (Y \to \epsilon) &= (F(\epsilon) - \{\epsilon\}) \cup FW(Y) = \{a, c, d\} \\ Gui (X \to Y) &= \{a, c, d\} \\ Gui (X \to a) &= \{a\} \end{split}$$

Parsificazione top-down: esercizio

Data la seguente Grammatica:

Produzione

- 1. $S \rightarrow PQ$
- 2. $Q \rightarrow \&PQ$
- 3. $Q \rightarrow \epsilon$
- 4. $P \rightarrow aPb$
- 5. $P \rightarrow bPa$
- 6. $P \rightarrow c$

Insieme guida

{a, b, c}

{&}

{\$}

{a}

{b}

{C}

Verificare gli insiemi guida calcolando I FIRST e i FOLLOW usando direttamente le definizioni.

Calcolo di FIRST(X) (abbreviato F(X))

Un metodo per calcolare gli F(X) per una grammatica $G = \langle V, \Sigma, P, S \rangle$ dove $X \in V \cup \Sigma$:

- 1. Si pone $F(a) = \{a\}$ per ogni $a \in \Sigma$.
- 2. Si pone inizialmente

$$F(A) = \{a \mid A \rightarrow a \ \alpha \in P\} \cup \{\epsilon \mid se \ A \rightarrow \epsilon \in P\}.$$

Altrimenti si inizializza $F(A) = \{\}$ (insieme vuoto)

- 3. *per ogni* produzione $A \rightarrow Y_1...Y_k$:
 - 1. si aggiunge $F(Y_1)$ -{ ε } a F(A).
 - 2. se $\varepsilon \in F(Y_1)$ (tipicamente se $Y_1 \rightarrow \varepsilon \in P$) si aggiunge $F(Y_2) \{\varepsilon\}$ a F(A).
 - 3. se $\varepsilon \in F(Y_1)$ e $\varepsilon \in F(Y_2)$ si aggiunge $F(Y_3) \{\varepsilon\}$ a F(A).
 - 4.
 - 5. se $\varepsilon \in F(Y_1),...,F(Y_k)$ si aggiunge ε a F(A).
- 4. si ripete il passo 3, fino a che gli insiemi F(A) non cambiano più

$$G = \langle \{X, Y, Z\}, \{a, c, d\}, P, Z \rangle$$

$$P: \quad Z \rightarrow d \mid XYZ$$

$$Y \rightarrow c \mid \varepsilon$$

$$X \rightarrow Y \mid a$$

Valori iniziali: F(a) = {a}, F(c)= {c}, F{d} = {d} (non cambiano più)
$$F(Z) = {d}, F(Y) = {\epsilon, c}, F(X) = {a}$$

Esaminiamo le produzioni nell'ordine in cui sono scritte:

Dopo una *prima* passata: $F(Z) = \{d, c,a\}, F(Y) = \{\epsilon,c\}, F(X\} = \{\epsilon,c,a\}$ Questi sono i valori *definitivi*.

Per esempio $c \in F(Z)$. Infatti $Z \to XYZ \to YYZ \to YZ \to cZ$ Esempio 2 : $F(XYZ) = \{c, a\} \cup \{c\} \cup \{d, c, a\} = \{d, c, a\}$

Calcolo dei Follow(A) (abbreviato Fw(A))

Sia data una grammatica $G = \langle V, \Sigma, P, S \rangle$. Si eseguono i seguenti passi

- 1. Per ogni $A \in V$ Si calcolano F(A). Si suppone inizialmente $Fw(S) = \{\$\}$ e Fw(A) vuoto per ogni altra $A \in V$.
- 2. per ogni produzione $A \rightarrow \alpha B\beta \in P$, e per ogni β si *aggiunge* $F(\beta) \{\epsilon\}$ a Fw(B)
- 3. per ogni produzione $A \rightarrow aB\beta \in P$ tale che $\varepsilon \in F(\beta)$ (ovvero $\beta \in V^*$ è annullabile) si aggiunge Fw(A) a Fw(B);
- 4. si ripete il passo 3. fino a che gli insiemi non si cambiano più.

Note: conviene esaminare le produzioni in un ordine fisso.

$$Z \rightarrow d \mid XYZ$$

 $Y \rightarrow c \mid \varepsilon$
 $X \rightarrow Y \mid a$

Ricordiamo che $F(Z) = \{d,c,a\}, F(Y) = \{\epsilon,c\}, F(X) = \{\epsilon,c,a\}$

Valori *iniziali*: $Fw(Z)=\{\$\}$, $Fw(Y)=\{\}$, $Fw(X)=\{\}$

Dopo il passo 2.: $Fw(Z)=\{\$\}$, $Fw(Y)=\{d,c,a\}$, $Fw(X)=\{d,a,c\}$.

Non ci sono le condizioni per eseguire il passo 3. quindi si salta.

Questi sono i valori definitivi.

Per esempio $a \in Fw(Y)$. Infatti $Z \rightarrow XYZ \rightarrow XYXYZ \rightarrow XYaYZ$