

Teoria dei linguaggi e grammatiche

Un Alfabeto Σ è un insieme di elementi, detti simboli. Si denota con $|\Sigma|$ la cardinalità dell'alfabeto Σ .

Una stringa $\alpha = \alpha_1 \alpha_2 \dots \alpha_n$ è una sequenza di simboli. Si denota con $|\alpha|$ la lunghezza della stringa.

Un Linguaggio è un insieme di stringhe. Dato un linguaggio L si denota con $|L|$ la sua cardinalità (finita o infinita).

Dato un linguaggio L , una stringa $\alpha \in L$ è detta frase di L .

Operazioni sulle stringhe

Date due stringhe $\alpha = \alpha_1 \alpha_2 \dots \alpha_n$ e $\beta = \beta_1 \beta_2 \dots \beta_n$, $\alpha = \beta$ se e solo se $\alpha_i = \beta_i$ per $1 \leq i \leq n$.

Date due stringhe $\alpha = \alpha_1 \alpha_2 \dots \alpha_n$ e $\beta = \beta_1 \beta_2 \dots \beta_n$ si definisce concatenazione di α e β la stringa:

$$\gamma = \alpha\beta = \alpha_1 \alpha_2 \dots \alpha_n \beta_1 \beta_2 \dots \beta_n .$$

La concatenazione è un'operazione associativa ma non commutativa.

Si definisce stringa nulla la stringa ϵ tale che $|\epsilon| = 0$. Data una qualunque stringa α , $\alpha\epsilon = \epsilon\alpha$.

Date le stringhe α, β, γ e la stringa $\delta = \alpha\beta\gamma$, le stringhe α, β, γ sono dette sottostringhe di δ , e le stringhe α e γ sono dette rispettivamente prefisso e suffisso di δ .

Data la stringa $\alpha = \alpha_1 \alpha_2 \dots \alpha_n$ la stringa $\alpha^R = \alpha_n \alpha_{n-1} \dots \alpha_1$ è detta stringa riflessa di α .
Data la stringa α , $(\alpha^R)^R = \alpha$, $(\alpha\beta)^R = \beta^R \alpha^R$.

Data la stringa α , la potenza ennesima α^n di α è la stringa ottenuta concatenando $n-1$ volte α con se stessa; $\alpha^0 = \epsilon$.

Operazioni sui linguaggi

Dati due linguaggi L_1 e L_2 , si definisce concatenazione di L_1 e L_2 il linguaggio $L = L_1 L_2 = \{ \alpha\beta \mid \alpha \in L_1, \beta \in L_2 \}$

Si definisce il linguaggio vuoto Φ tale che $L\Phi = \Phi L; |\Phi| = 0$.

Data il linguaggio L , la potenza ennesima L^n di L è il linguaggio ottenuto concatenando $n-1$ volte L con se stesso; $L^0 = \{\epsilon\}$.

Dato il linguaggio L , si definisce chiusura riflessiva di L rispetto alla concatenazione il linguaggio $L^* = \bigcup_{i=0}^{\infty} L^i$.

Proprietà della chiusura riflessiva:

$$\begin{array}{ll} L^* \supseteq L & \text{monotonicità} \\ (x \in L^*) \wedge (y \in L^*) \Rightarrow xy \in L^* & \text{chiusura} \\ (L^*)^* = L^* & \text{idempotenza} \\ (L^*)^R = (L^R)^* & \text{commutatività} \end{array}$$

Si definisce linguaggio universale (o monoide libero) di alfabeto Σ il linguaggio $L = \Sigma^*$. Dato il linguaggio L , si definisce complemento di L (scritto $\neg L$) il suo complemento rispetto a Σ^* .

Dato il linguaggio L , si definisce chiusura non riflessiva di L rispetto alla concatenazione il linguaggio:

$$L^+ = \bigcup_{i=1}^{\infty} L^i.$$

Si noti che: $L^* \supseteq L^+ \cup \{\epsilon\}$, $L^+ = LL^+ = L^+L$.

Siano dati i linguaggi L ed L_b , con alfabeti rispettivamente Σ e Δ , e sia $b \in \Sigma$. Si definisce sostituzione di L_b a b in L il linguaggio L' tale che:

$$L' = \{ \alpha \in ((\Sigma - \{b\}) \cup \Delta)^* \mid \alpha = x_1 y_1 x_2 y_2 \dots y_{k-1} x_k, \text{ con } x_i, b x_i b \in L, x_i \in (S - \{b\})^*, y_i \in L_b, 1 \leq i \leq k \}$$

Dati i linguaggi L_1 e L_2 si definisce quoziente di L_1 rispetto a L_2 il linguaggio L tale che:

$$L = \{ \alpha \mid (\exists \beta \in L_2) \wedge (\alpha\beta \in L_1) \}$$

Espressioni e linguaggi regolari

Un linguaggio L_Σ di alfabeto $\Sigma = \{a_1, a_2, \dots, a_k\}$ è detto regolare se può essere espresso mediante operazioni di concatenazione, unione e chiusura riflessiva applicate un numero finito di volte ai linguaggi unitari (i.e. composti da una sola stringa) $L_1 = \{a_1\}$, $L_2 = \{a_2\}$, ... $L_k = \{a_k\}$, ed al linguaggio vuoto Φ .

Una espressione regolare (e.r.) x di alfabeto Σ è definita costruttivamente come segue:

1. $x = \Phi$
2. $x = \epsilon$
3. $x = a$, dove $a \in \Sigma$
4. $x = yz$, dove y e z sono e.r.
5. $x = y \cup z$ dove y e z sono e.r.
6. $x = y^*$ dove y è e.r.
7. nient'altro è un'espressione regolare

Ogni linguaggio finito è regolare, in quanto unione di un numero finito di linguaggi unitari.

NOTA: il linguaggio delle espressioni regolari può essere esteso con gli operatori potenza (x^n) e chiusura transitiva (x^+).

Una sottoespressione (s.e) di una e.r. è così definita:

x_k è una s.e. di $(x_1 \cup \dots \cup x_k \cup \dots \cup x_n)$

x_k è una s.e. di $x_1 \dots x_k \dots x_n$

x è una s.e. di x^* , x^+ , x^n

x è una s.e. di se stessa

ϵ è una s.e. di ogni espressione.

Date le e.r. x ed y , x implica y (scritto $x \Rightarrow y$) se il linguaggio definito da x include quello definito da y :

$(x_1 \cup \dots \cup x_k \cup \dots \cup x_n) \Rightarrow x_k \quad 1 \leq k \leq n$;

$(x)^* \Rightarrow xx \dots x$ (ripetuto k volte, $k \geq 0$)

$(x)^+ \Rightarrow xx \dots x$ (ripetuto k volte, $k > 0$)

$(x)^n \Rightarrow xx \dots x$ (ripetuto n volte)

Date due e.r. x_1 ed x_2 , x_1 diviene immediatamente x_2 (scritto $x_1 \rightarrow x_2$) se:

$$x_1 = \alpha\beta\gamma, \quad x_2 = \alpha\delta\gamma, \quad \beta \Rightarrow \delta$$

dove α, β, γ sono sottoespressioni di x_1

Nel modo solito si definiscono $x_1 \rightarrow^n x_2$, $x_1 \rightarrow^+ x_2$, $x_1 \rightarrow^* x_2$

Data una e.r. x , il linguaggio $L(x)$ prodotto da x è:

$$L(x) = \{ \alpha \in V_T^* \mid x \rightarrow^* \alpha \}$$

Grammatiche generative

Una grammatica è una quadrupla $G = \langle V_T, V_N, P, Z \rangle$ dove:

V_T = insieme di simboli terminali

V_N = insieme di simboli non terminali ($V_T \cap V_N = \emptyset$)

P = insieme di produzioni (o regole di riscrittura)

Z = assioma ($Z \in V_N$)

Le produzioni hanno la seguente struttura:

$$\alpha ::= \beta$$

dove α è detta parte sinistra e β è detta parte destra.

Più produzioni aventi la stessa parte sinistra possono essere raggruppate usando la notazione:

$$\alpha ::= \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$$

Date due stringhe $\alpha \in (V_T \cup V_N)^+$, $\beta \in (V_T \cup V_N)^*$ si dice che β deriva da α oppure che α si riduce a β

($\alpha \Rightarrow \beta$) se e solo se:

$$\alpha = \alpha_1 \gamma \alpha_2, \quad \beta = \beta_1 \delta \beta_2, \quad \gamma ::= \delta$$

Conseguentemente:

$\alpha \Rightarrow^n \beta$ se e solo se $\alpha \Rightarrow \gamma_1 \Rightarrow \gamma_2 \Rightarrow \dots \Rightarrow \gamma_{n-1} \Rightarrow \beta$

$\alpha \Rightarrow^+ \beta$ se e solo se $\alpha \Rightarrow^n \beta$ con $n > 0$

$\alpha \Rightarrow^* \beta$ se e solo se $\alpha \Rightarrow^n \beta$ con $n \geq 0$

Data una grammatica G , il linguaggio generato da G è: $L(G) = \{ \alpha \in V_T^* \mid Z \Rightarrow^* \alpha \}$.

Classi di linguaggi:

Tipo 0	$\alpha ::= \beta$, dove $\alpha \in V^*$ ($\alpha \in V_N^*$), $\beta \in V^*$ $V = V_T \cup V_N$	linguaggi ricorsivamente enumerabili	Macchina di Turing
Tipo 1	$\alpha ::= \beta$, dove $\alpha \in V^*$, $\beta \in V^*$ $ \alpha \leq \beta $ $V = V_T \cup V_N$	linguaggi contestuali (context sensitive, CS)	Macchina di Turing con nastro finito
Tipo 2	$A ::= \beta$, dove $A \in V_N^*$, $\beta \in V^*$ $V = V_T \cup V_N$	linguaggi non contestuali (context free, CF)	Automa a pila non deterministico
Tipo 3	$A ::= xB$, dove $A \in V_N^*$, $x \in V_T$ $B \in (V_N \cup \{\epsilon\})$	linguaggi regolari	Automa a stati finiti

Data una grammatica G ed un simbolo $A \in V_N$, il linguaggio generato da G partendo da A è:
 $L_A(G) = \{ \alpha \in V_T^* \mid A \Rightarrow^* \alpha \}$.

Data una grammatica G ed un simbolo $A \in V_N$, una forma di frase generata da G partendo da A è una stringa $\alpha \in (V_T \cup V_N)^*$ tale che $A \Rightarrow^* \alpha$.

Data una grammatica G , una frase generata da G è una stringa $\alpha \in V_T^*$ tale che $Z \Rightarrow^* \alpha$, (si noti che $\alpha \in L(G)$).

Linguaggi non contestuali

Una grammatica non contestuale (cf) G si dice ridotta se valgono le seguenti condizioni:

- 1) non consente derivazioni circolari del tipo $A \Rightarrow^+ A$ per qualunque $A \in V_N$
- 2) ogni non terminale A è raggiungibile dall'assioma:
 $Z \Rightarrow^* \alpha A \beta$, con $\alpha, \beta \in (V_T \cup V_N)^*$ per qualunque $A \in V_N$
- 3) ogni non terminale A genera un linguaggio non vuoto: $L_A(G) \neq \Phi$

Un simbolo $A \in V_N$ si dice ricorsivo se esiste una derivazione $A \Rightarrow^+ \alpha A \beta$, con $\alpha, \beta \in (V_T \cup V_N)^*$. La ricorsione è diretta se $A \Rightarrow \alpha A \beta$, altrimenti è indiretta.

La derivazione $A \Rightarrow^+ b_1 b_2 \dots b_n$, $b_i \in (V_T \cup V_N)$, $1 \leq i \leq n$, può essere descritta mediante un albero, detto albero sintattico, con radice A e foglie b_1, b_2, \dots, b_n (da sinistra a destra).

La derivazione $\alpha_0 \Rightarrow \alpha_1 \dots \Rightarrow \alpha_n$, dove $\alpha_i = \beta_i A \gamma_i$, $\alpha_{i+1} = \beta_i \delta_i \gamma_i$, con $\delta_i, \beta_i, \gamma_i \in (V_T \cup V_N)^*$, $A \in V_N$, è detta canonica sinistra (destra) se $\beta_i \in V_T$ ($\gamma_i \in V_T$) per ogni i : $0 \leq i \leq n-1$, ed è denotata con $\alpha_0 \xRightarrow{s}^+ \alpha_n$ ($\alpha_0 \xRightarrow{d}^+ \alpha_n$).

Una frase $\alpha \in L(G)$ è detta sintatticamente ambigua se è generabile da G con n ($n \geq 2$) alberi sintattici distinti (n è detto grado di ambiguità); in tal caso anche la grammatica è detta ambigua (proprietà indecidibile).

Un linguaggio è inherentemente ambiguo se tutte le grammatiche che lo generano sono ambigue.

Due grammatiche G_1 e G_2 sono debolmente equivalenti se $L(G_1) = L(G_2)$; la proprietà di equivalenza debole è indecidibile.

Trasformazioni sintattiche conservative (cf)

1) Sostituzione

Siano $B ::= \beta_1 | \beta_2 | \dots | \beta_n$ tutte le produzioni della grammatica G aventi il simbolo B come parte sinistra. Se la produzione $A ::= \alpha B \gamma$ di G è sostituita dalle produzioni $A ::= \alpha \beta_1 \gamma | \alpha \beta_2 \gamma | \dots | \alpha \beta_n \gamma$, la grammatica G' così ottenuta è debolmente equivalente a G .

2) Eliminazione delle ricorsioni sinistre

Un nonterminale A è ricorsivo a sinistra (destra) se esiste una derivazione $A \Rightarrow^+ A \beta$, ($A \Rightarrow^+ \beta A$) con $\beta \neq \epsilon$

Una produzione $A ::= Ab$ è detta immediatamente ricorsiva (i.r.) a sinistra.

Eliminazione delle produzioni i.r. a sinistra

Sia G una grammatica contenente le produzioni:

(1) $A ::= A\beta_1 | A\beta_2 | \dots | A\beta_n | \gamma_1 | \gamma_2 | \dots | \gamma_k$
con $\beta_1, \beta_2, \dots, \beta_n \neq \epsilon$

La grammatica G' ottenuta sostituendo le (1) con:

$A ::= \gamma_1 A' | \gamma_2 A' | \dots | \gamma_k A'$
 $A' ::= \beta_1 A' | \beta_2 A' | \dots | \beta_n A' | \epsilon$

è debolmente equivalente a G .

Eliminazione delle ricorsioni non immediate

Sia data la grammatica G ricorsiva a sinistra, con $V_N = \{A_1, A_2, \dots, A_n\}$ e priva di ϵ -regole (non esistono produzioni $A_i ::= \epsilon$).

Algoritmo

for $i := 1$ to n do

begin

for $k := 1$ to $i-1$ do

begin

sostituire ogni produzione della forma

$$A_i ::= A_k \alpha$$

con le produzioni $A_i ::= \gamma_1 \alpha \mid \gamma_2 \alpha \mid \dots \mid \gamma_h \alpha$

dove $A_k ::= \gamma_1 \mid \gamma_2 \mid \dots \mid \gamma_h$ sono tutte le

produzioni presenti con parte sinistra A_k

end

eliminare le ricorsioni immediate a sinistra di A_i

end

La grammatica G' ottenuta applicando il precedente algoritmo alla grammatica G è debolmente equivalente a G.

Forme normali per linguaggi cf

Data una grammatica G possono essere costruite forme normali debolmente equivalenti.

1) Forma normale di Chomsky

Le produzioni sono di due tipi:

$$A ::= BC \quad \text{oppure} \quad A ::= a$$

con $A, B, C \in V_N$ e $a \in (V_T \cup \{\epsilon\})$

2) Forma normale di Greibach

Le produzioni sono del tipo:

$$A ::= a\alpha$$

con $a \in V_T$ e $\alpha \in (V_T \cup V_N)^*$

La forma di Greibach si applica a linguaggi privi di ϵ -regole e si ottiene applicando ripetutamente le trasformazioni di sostituzione ed eliminazione delle ricorsioni sinistre.

Grammatiche non contestuali estese (ecf)

Le produzioni di una grammatica ecf hanno la seguente struttura: $A ::= \alpha$, dove $A \in V_N$, α è una e.r. di alfabeto $(V_T \cup V_N)$.

Data la produzione $A ::= \alpha$, si dice che A produce β (scritto $A \Rightarrow \beta$) se $\alpha \Rightarrow \beta$, dove α e β sono e.r.

Grammatiche lineari

Una grammatica cf è detta lineare se ogni produzione ha la forma: $A ::= \alpha$, dove $\alpha \in V_T^* \times (V_N \cup \{\epsilon\}) \times V_T^*$.

Una grammatica cf è detta lineare destra (sinistra) se ogni produzione ha la forma: $A ::= \alpha B$ ($A ::= B\alpha$), dove $\alpha \in V_T^*$, $B \in (V_N \cup \{\epsilon\})$.

Una grammatica cf è detta strettamente lineare destra (sinistra) se ogni produzione ha la forma: $A ::= cB$ ($A ::= Bc$), dove $c \in V_T$, $B \in (V_N \cup \{\epsilon\})$.

Per ogni grammatica lineare destra (sinistra) G esiste una grammatica strettamente lineare G' debolmente equivalente.

La classe dei linguaggi generati dalle grammatiche lineari destre e lineari sinistre coincide con quella dei linguaggi regolari.

Linguaggi deterministici

Si definiscono linguaggi deterministici i linguaggi non contestuali riconoscibili da un automa push-down deterministico.

Siano:

L_{cf} l'insieme dei linguaggi non contestuali

L_d l'insieme dei linguaggi deterministici

L_r l'insieme dei linguaggi regolari

Si ha: $L_{cf} \supset L_d \supset L_r$

- (• il linguaggio $L_d = \{ a^n b^n \mid n > 0 \}$ è deterministico ma non regolare;
- il linguaggio $L_d = \{ a^n b^n \mid n > 0 \} \cup \{ a^n b^{2n} \mid n > 0 \}$ è non contestuale ma non deterministico)

Proprietà:

- la famiglia dei linguaggi deterministici non è chiusa rispetto all'unione, all'intersezione ed alla concatenazione;
- il complemento di un linguaggio deterministico è deterministico;
- se L_d è deterministico e L_r è regolare, L_d/L_r e $L_d L_r$ sono deterministici;

Analisi sintattica (parsing)

Si definiscono due metodologie principali:

- analisi ascendente (bottom-up): si costruisce la derivazione canonica destra della stringa procedendo dalle foglie verso la radice dell'albero sintattico (ordine riflesso);
- analisi discendente (top-down): si costruisce la derivazione canonica sinistra;

Analisi bottom-up LR(0)

Sia $Z \xRightarrow{\sigma} \beta A \gamma \xRightarrow{\delta} \beta \alpha \gamma$ una derivazione canonica destra, con $\langle A ::= \alpha \rangle \in P$, $\alpha = \alpha_1 \alpha_2$. La stringa $\beta \alpha_1$ è un prefisso ascendente e la stringa $\langle A ::= \alpha_1 \overset{\vee}{\alpha}_2 \rangle$ è detta candidata per $\beta \alpha_1$.

Dato il prefisso ascendente α , $C(\alpha)$ è l'insieme delle candidate per α .

Data una grammatica G , l'insieme dei prefissi ascendenti costituisce un linguaggio regolare $L_p(G)$ di alfabeto $\Sigma_p = V_T \cup V_N$ (riconoscibile mediante un automa a stati finiti).

Dati due insiemi di candidate C e C' , C' è detto estensione di C (denotato $C \mid\!-\! C'$) se e solo se: $C' = C \cup \{ \langle B ::= \overset{\vee}{\beta} \rangle \mid \langle B ::= \beta \rangle \in P \text{ e } \langle A ::= \alpha_1 \overset{\vee}{B} \alpha_2 \rangle \in C \}$

Dato l'insieme di candidate C si definisce chiusura di C l'insieme C^\oplus tale che:

$$C \mid\!-\!^+ C^\oplus \text{ e } C^\oplus = C', \forall C': C^\oplus \mid\!-\! C'$$

Data la grammatica $G = \langle V_T, V_N, P, Z \rangle$, l'automa $LR_0(G) = \langle Q_p, \Sigma_p, t_p, q_{p0} \rangle$ riconoscitore del linguaggio $L_p(G)$ con prospezione 0 è definito nel modo seguente:

- se $p \in Q_p$, $h \in \Sigma_p$ e $C = \{ \langle A ::= \alpha_1 \overset{\vee}{h} \alpha_2 \rangle \mid \langle A ::= \alpha_1 \overset{\vee}{h} \alpha_2 \rangle \in p \}$, allora $\exists q \in Q_p$: $q = K^\oplus$, dove $K = \{ \langle A ::= \alpha_1 h \overset{\vee}{\alpha}_2 \rangle \mid \langle A ::= \alpha_1 \overset{\vee}{h} \alpha_2 \rangle \in C \}$ e $t_p(p, h) = q$.
- se $C = \{ \langle Z ::= \overset{\vee}{\alpha} \rangle \mid \langle Z ::= \alpha \rangle \in P \}$, $q_{p0} = C^\oplus$.

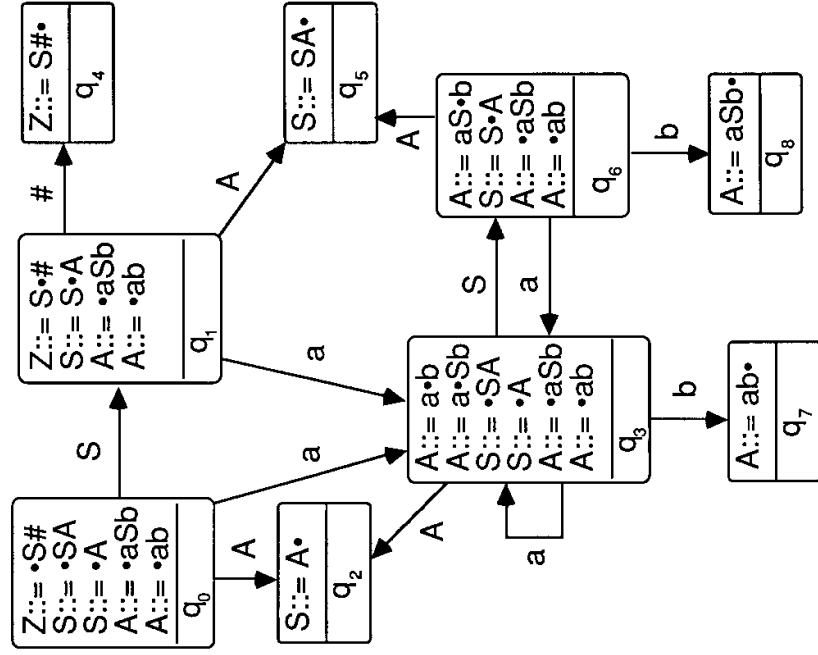
Dato l'automa $LR_0(G)$ uno stato q si definisce:

- stato di riduzione se contiene la candidata di riduzione $\langle A ::= \alpha \rangle$
- stato di spostamento se contiene la candidata di spostamento $\langle A ::= \alpha_1 \overset{\vee}{\alpha}_2 \rangle$

Uno stato di riduzione che contenga più di una candidata, oppure uno stato che sia contemporaneamente di spostamento e di riduzione è detto stato inadeguato per l'analisi $LR(0)$.

Esempio

Data la grammatica G_1 con $P = \{ Z ::= S\#, S ::= SA \mid A, A ::= aSb \mid ab \}$ e assioma Z , l'automa $LR_0(G_1)$ è:



Nessuno stato è inadeguato.

Una grammatica G è detta LR(0) (left-to-right con proiezione 0) se l'automa $LR_0(G)$ non contiene stati inadeguati.

Data una grammatica G di tipo LR(0), l'automa $A_{LR(0)}(G) = \langle Q, I, \Gamma, t, Z_0, q_0 \rangle$ riconoscitore del linguaggio $L(G)$ è definito nel modo seguente:

- $Q = \{q_0, q_F\}$
 - $I = V_T$
 - $\Gamma = Q_p \cup \Sigma_p$
 - $Z_0 = q_{p0}$
 - $\langle q_0, \alpha\alpha, h_1 \rangle \vdash \langle q_0, \alpha, h_k a h_1 \rangle$ se $h_i \in Q_p$ è uno stato di spostamento e $t_p(h_i, a) = h_k$
 - $\langle q_0, \alpha, h_i b_n \dots h_1 b_1 h_{i0} \rangle \vdash \langle q_0, \alpha, h_k A h_{i0} \rangle$ se $h_i \in Q_p$ è uno stato di riduzione con candidata $\langle A ::= b_1 b_2 \dots b_n \rangle$ ($A \neq Z$) e $t_p(h_{i0}, A) = h_k$
 - $\langle q_0, \varepsilon, h_i b_n \dots h_1 b_1 Z_0 \rangle \vdash \langle q_F, \varepsilon, Z_0 \rangle$ se $h_i \in Q_p$ è uno stato di riduzione con candidata $\langle Z ::= b_1 b_2 \dots b_n \rangle$
- Z è l'assioma di G .

L'insieme $L_{LR(0)}$ dei linguaggi generati da grammatiche LR(0) e' costituito da linguaggi deterministici privi di prefissi (se $L \in L_{LR(0)}$ e $\alpha \in L$, $\alpha\beta \notin L$ per $\forall \beta \neq \varepsilon$); quindi $L_G \supset L_{LR(0)}$

Ogni linguaggio deterministico può essere trasformato in un equivalente linguaggio privo di prefissi terminando ogni frase con la marca di fine stringa "#".
Tuttavia una grammatica LR(0) può risultare complessa ed inefficiente dal punto di vista del trattamento degli errori .

Analisi bottom-up LR(1)

Ad ogni candidata LR(0) viene associato un insieme di simboli terminali detto insieme di prospezione .

Una candidata LR(1) ha quindi la forma:

$$\langle A ::= \alpha_1 \overset{\vee}{\alpha_2} \{a_1, a_2, \dots, a_n\} \rangle \quad a_i \in V_T \cup \{\#\}$$

dove $\{a_1, a_2, \dots, a_n\}$ è l'insieme di prospezione.

$\langle A ::= \alpha_1 \overset{\vee}{\alpha_2} \Pi \rangle$ è candidata di un prefisso ascendente $\beta\alpha_1$ se per ogni $a_i \in \Pi$ esistono le derivazioni destre:

$$Z \overset{d}{\Rightarrow}^* \beta A \gamma \overset{d}{\Rightarrow} \beta \alpha_1 \alpha_2 \gamma \quad \text{e} \quad \gamma = a_i \delta \quad (\gamma = \varepsilon \text{ se } a_i = \#)$$

L'automa $LR_1(G)$ riconoscitore dei prefissi ascendenti può essere calcolato in modo analogo al caso LR(0) :

- se $p \in Q_p$, $h \in \Sigma_p$ e $C = \{ \langle A ::= \alpha_1 \overset{\vee}{h\alpha_2}, \Pi \rangle \mid \langle A ::= \alpha_1 \overset{\vee}{h\alpha_2}, \Pi \rangle \in p \}$, allora $\exists q \in Q_p$: $q = K^\oplus$, dove $K = \{ \langle A ::= \alpha_1 h \overset{\vee}{\alpha_2}, \Pi \rangle \mid \langle A ::= \alpha_1 \overset{\vee}{h\alpha_2}, \Pi \rangle \in C \}$ e $t_p(p, h) = q$.
- se $C = \{ \langle Z ::= \overset{\vee}{\alpha_1} \{\#\} \rangle \mid \langle Z ::= \overset{\vee}{\alpha} \rangle \in P \}$, $q_{p0} = C^\oplus$.

ridefinendo il concetto di estensione nel modo seguente:

- dati due insiemi di candidate C e C' , C' è detto estensione di C (denotato $C \mid - C'$) se e solo se:

$$C' = C \cup \{ \langle B ::= \overset{\vee}{\beta}, \Pi' \rangle \mid \langle B ::= \beta \rangle \in C \}$$

$$\langle A ::= \alpha_1 \overset{\vee}{B} \alpha_2, \{a_1, a_2, \dots, a_n\} \rangle \in C'$$

dove: $b \in \Pi'$ se esiste la derivazione $\alpha_2 \overset{d}{\Rightarrow}^+ b \gamma$

$a_i \in \Pi'$ se esiste la derivazione $\alpha_2 \overset{d}{\Rightarrow}^+ \varepsilon$

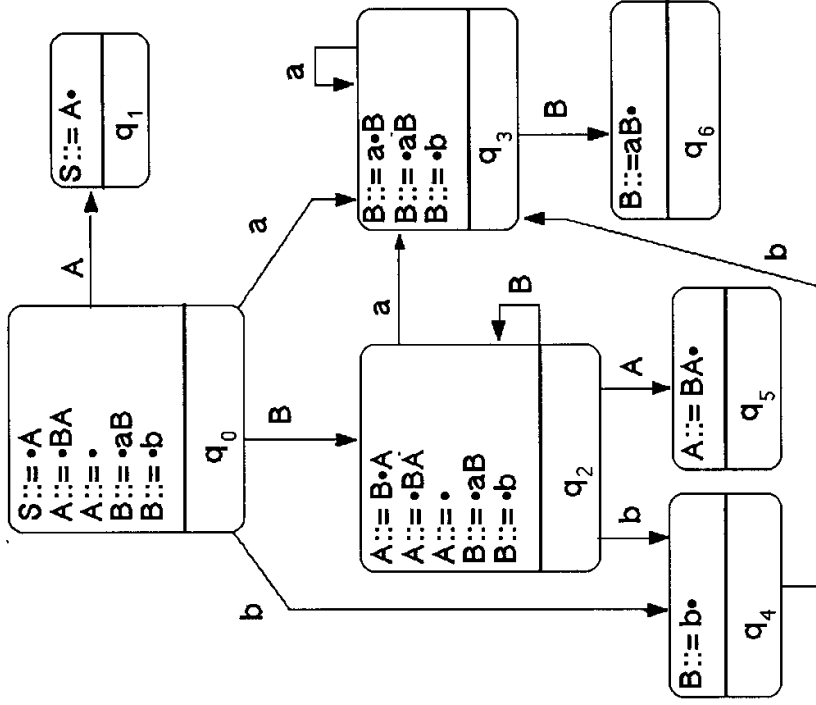
Una grammatica G è detta LR(1) se per ogni stato q dell'automa $LR_1(G)$ contenente una candidata di riduzione $\langle A ::= \alpha \overset{\vee}, \Pi_1 \rangle$ valgono entrambe le seguenti condizioni:

- $t(q, a)$ non è definito per alcun $a \in \Pi_1$
- se in q vi è un'altra candidata di riduzione $\langle B ::= \beta \overset{\vee}, \Pi_2 \rangle$, $\Pi_1 \cap \Pi_2 = \Phi$

Uno stato che non soddisfa questa condizione è detto inadeguato per l'analisi LR(1).

Esempio

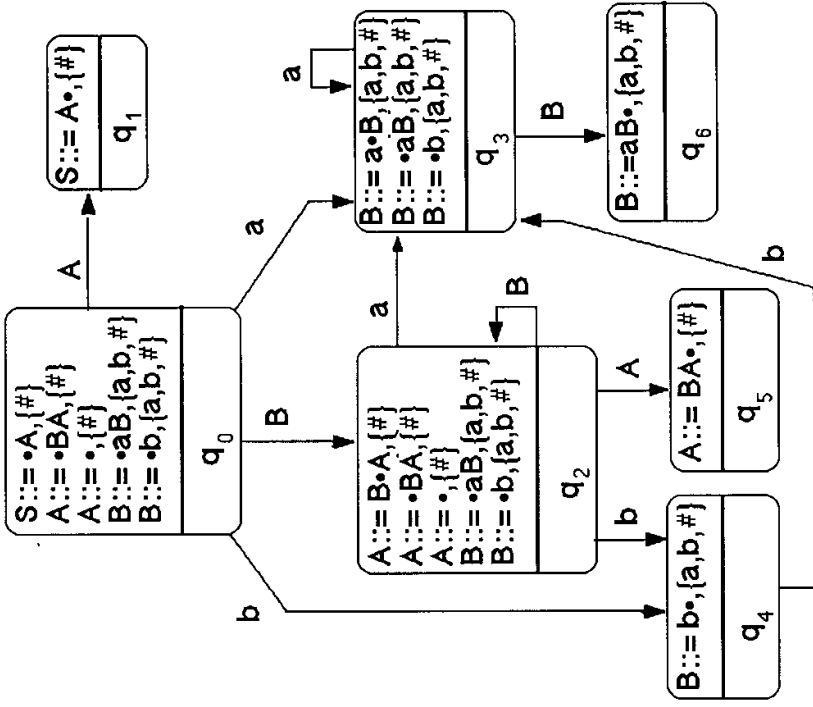
Data la grammatica G_2 con $P = \{ S ::= A, A ::= BA \mid \epsilon, B ::= aB \mid b \}$ e assioma S , l'automa $LR_0(G_2)$ è:



Gli stati q_0 e q_2 sono inadeguati nell'analisi LR(0)

Esempio

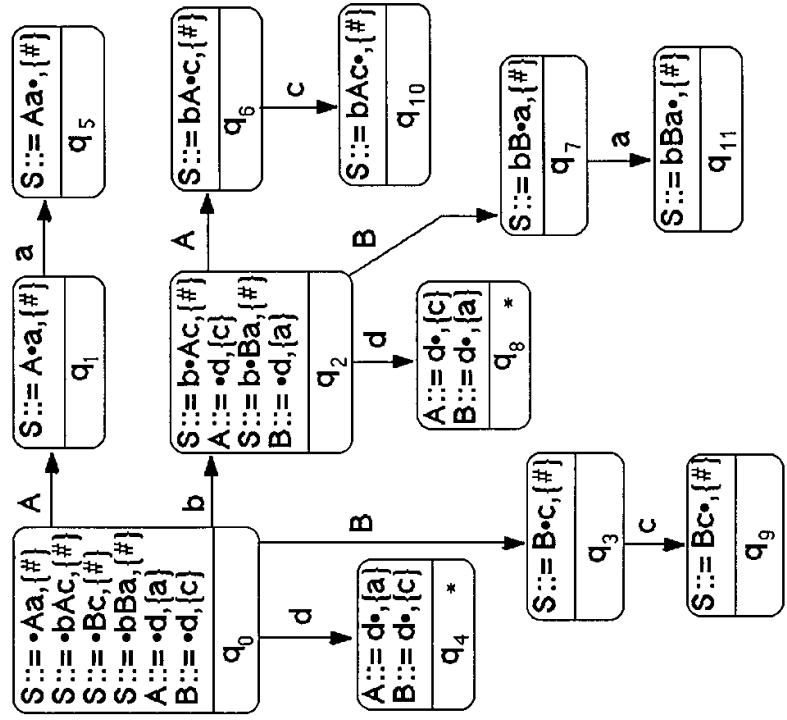
automa $LR_1(G_2)$ con insiemi di prospezione per analisi LR(1)



Gli stati q_0 e q_2 sono adeguati per l'analisi LR(1)

Esempio

Data la grammatica G_3 con $P = \{ S ::= Aa \mid bAc \mid Bc \mid bBa, A ::= d, B ::= d \}$ e assioma S , l'automa $LR_1(G_3)$ è:



NOTA: l'automa per il riconoscimento $LR(0)$ non distinguerebbe gli stati q_4 e q_8

Data una grammatica G di tipo $LR(1)$, l'automa $A_{LR(1)}(G) = \langle Q, I, \Gamma, t, Z_0, q_0 \rangle$ riconoscitore del linguaggio $L(G)$ è definito nel modo seguente:

- $Q = \{q_0, q_F\}$
- $I = V_T$
- $\Gamma = Q_p \cup \Sigma_p$
- $Z_0 = q_{p0}$
- $\langle q_0, \alpha\alpha, h_i \gamma \rangle \vdash \langle q_0, \alpha, h_k a h_i \gamma \rangle$ se $h_i \in Q_p$ contiene la candidata $\langle A ::= \beta_1 \overset{\vee}{a} \beta_2, \Pi \rangle$ e $t_p(h_i, a) = h_k$
- $\langle q_0, \alpha\alpha, h_i b_n \dots h_{i-1} b_1 h_{i_0} \rangle \vdash \langle q_0, \alpha\alpha, h_k A h_{i_0} \rangle$ se $h_i \in Q_p$ contiene la candidata $\langle A ::= b_1 b_2 \dots b_n, \Pi \rangle$
- $\langle A \neq Z \rangle, t_p(h_{i_0}, A) = h_k$ e $a \in \Pi$
- $\langle q_0, \#, h_{i_n} b_n \dots h_{i_1} b_1 Z_0 \rangle \vdash \langle q_F, \#, Z_0 \rangle$ se $h_{i_n} \in Q_p$ contiene la candidata $\langle Z ::= b_1 b_2 \dots b_n, \{ \# \} \rangle$ Z è l'assioma di G .

L'insieme $L_{LR(1)}$ dei linguaggi generati da grammatiche $LR(1)$ coincide con l'insieme dei linguaggi deterministici.

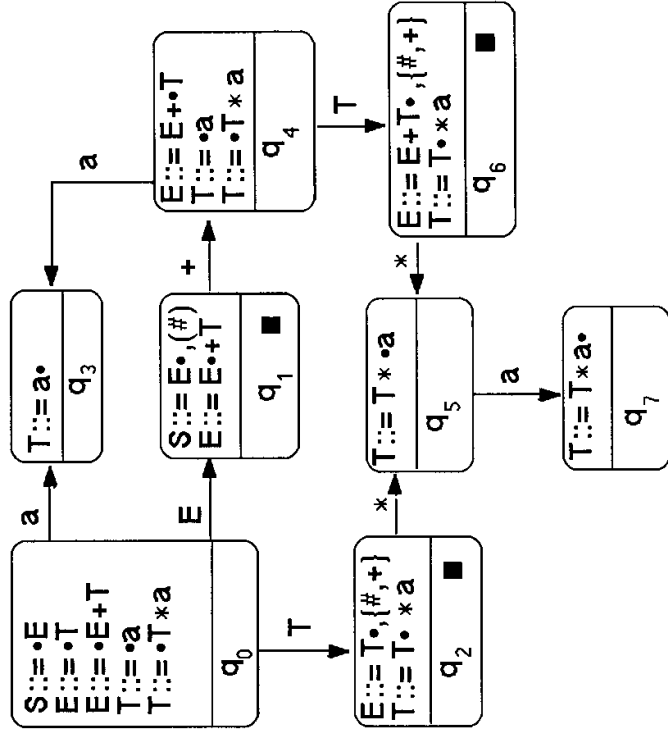
Analisi bottom-up SLR(1)

Sia A un simbolo nonterminale della grammatica G ; si definisce l'insieme $S(A)$ dei seguiti di A nel modo seguente:
 $S(A) = \{ a \in V_T \mid Z \Rightarrow^* \alpha A a \beta \}, \# \in S(A)$ se $Z \Rightarrow^* \alpha A$

Per ogni candidata $\langle A ::= \alpha^v, \Pi \rangle$ si avrà $\Pi = S(A)$.
 L'insieme di prospezione è calcolato solo per le candidate di riduzione appartenenti a stati inadeguati per l'analisi LR(0).

Esempio

Data la grammatica G_4 con $P = \{ S ::= E, E ::= T \mid E+T, T ::= a \mid T * a \}$ e assioma S , l'automa $SLR_1(G_4)$ è:



Gli stati q_1, q_2, q_6 sono inadeguati per l'analisi LR(0) e adeguati per l'analisi SLR(1).

Una grammatica è detta SLR(1) se per ogni candidata di riduzione $\langle A ::= \alpha^v, S(A) \rangle$ in uno stato q inadeguato per l'analisi LR(0) valgono le seguenti condizioni:

- in q non vi è un'altra candidata $\langle B ::= \beta^v a \delta \rangle$ con $a \in S(A)$
- se in q vi è un'altra candidata $\langle B ::= \alpha^v, S(B) \rangle$, $S(A) \cap S(B) = \Phi$

Per il calcolo di $S(A)$ si possono usare le seguenti regole:

- $\# \in S(Z)$, dove Z è l'assioma
- $a \in S(A)$ se $\langle B ::= \alpha A \beta \rangle \in P$ e $a \in I(\beta)$
- $a \in S(A)$ se $\langle B ::= \alpha A \rangle \in P$ ($A \neq B$) e $a \in S(B)$
- $a \in S(A)$ se $\langle B ::= \alpha A \beta \rangle \in P$ ($A \neq B$), $\beta \Rightarrow^+ \epsilon$ e $a \in S(B)$

dove $I(\beta) = \{ b \in V_T \mid \beta \Rightarrow^+ b \gamma \}$ (insieme degli inizi di β) è calcolabile come segue:

- $a \in I(\beta)$ se $\beta = a \gamma$
- $a \in I(\beta)$ se $\beta = A \gamma$, $\langle A ::= \delta \rangle \in P$ e $a \in I(\delta)$
- $a \in I(\beta)$ se $\beta = A \gamma$, $A \Rightarrow^+ \epsilon$ e $a \in I(\gamma)$

Analisi bottom-up LALR(1)

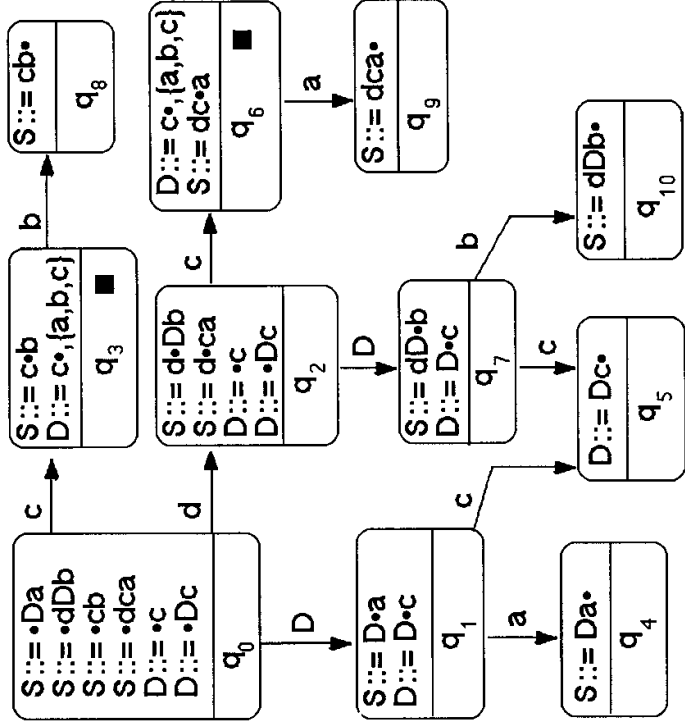
Per ogni candidata $\langle A ::= \alpha^v, \Pi \rangle$ appartenente allo stato q si avrà:

$\Pi = \{ a \in V_T \mid Z \Rightarrow^* \gamma A a \beta \Rightarrow \gamma \alpha a \beta \text{ e } t_p(q_{po}, \gamma \alpha) = q \}$
 dove t_p è la funzione di transizione di $LR_0(G)$.

L'insieme di prospezione è calcolato solo per le candidate di riduzione appartenenti a stati inadeguati per l'analisi LR(0).

Esempio

Data la grammatica G_5 con $P = \{ S ::= Da \mid dDb \mid cb \mid dca, D ::= c \mid Dc \}$ e assioma S , l'automa $SLR_1(G_5)$ è:



Gli stati q_3 e q_6 sono inadeguati per l'analisi LR(0) e SLR(1).

Una grammatica è detta LALR(1) se per ogni candidata di riduzione $\langle A ::= \alpha^v, \Pi_A \rangle$ in uno stato q inadeguato per l'analisi LR(0) valgono le seguenti condizioni:

- in q non vi è un'altra candidata $\langle B ::= \beta^v a \delta \rangle$ con $a \in \Pi_A$
- se in q vi è un'altra candidata $\langle B ::= \alpha^v, \Pi_B \rangle$, $\Pi_A \cap \Pi_B = \Phi$

Sia G_x l'insieme delle grammatiche di tipo x ; si ha:

- $G_{LR(1)} \supset G_{LALR(1)} \supset G_{SLR(1)} \supset G_{LR(0)}$

Analisi top-down LL(1)

Data una grammatica G e una produzione $\langle A ::= \alpha \rangle$ si definisce insieme guida di $\langle A ::= \alpha \rangle$ l'insieme:

$$G(\langle A ::= \alpha \rangle) = I(\alpha), \quad \text{se non esiste } \alpha \Rightarrow^* \epsilon$$

$$G(\langle A ::= \alpha \rangle) = I(\alpha) \cup S(A), \quad \text{se esiste } \alpha \Rightarrow^* \epsilon$$

Una grammatica G è LL(1) se, per ogni coppia di produzioni $\langle A ::= \alpha \rangle$ e $\langle A ::= \beta \rangle$:

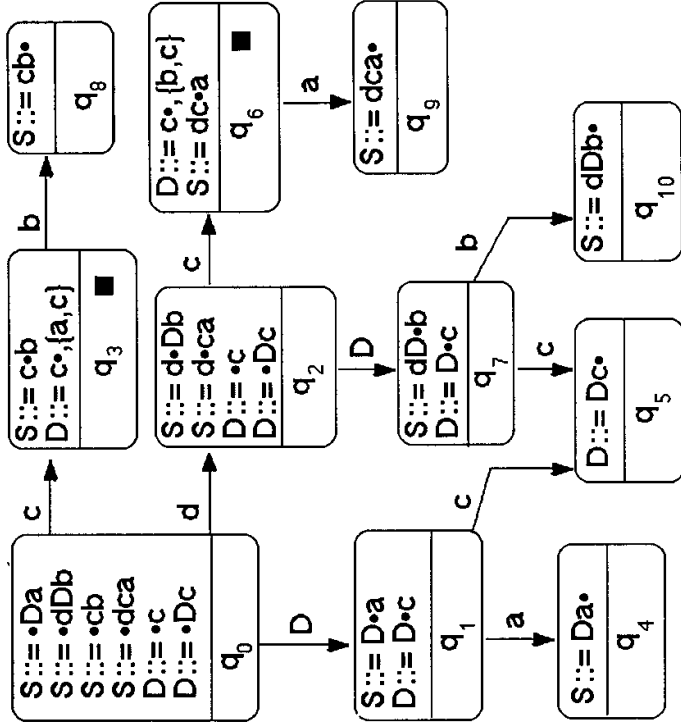
$$G(\langle A ::= \alpha \rangle) \cap G(\langle A ::= \beta \rangle) = \Phi$$

NOTA: una grammatica LL(1) non può avere produzioni ricorsive a sinistra.

Le grammatiche LL(1) sono un sottoinsieme delle grammatiche LR(1): $G_{LR(1)} \supset G_{LL(1)}$

Esempio

Data la grammatica G_5 con $P = \{ S ::= Da \mid dDb \mid cb \mid dca, D ::= c \mid Dc \}$ e assioma S , l'automa $LALR_1(G_5)$ è:



Gli stati q_3 e q_6 sono adeguati per l'analisi LALR(1).

Data una grammatica G di tipo $LL(1)$, l'automa $A_{LL(1)}(G) = \langle Q, I, \Gamma, t, Z_0, q_0 \rangle$ riconoscitore del linguaggio $L(G)$ è definito nel modo seguente:

- $Q = \{q_0, q_1, q_F\}$
- $I = V_T$
- $\Gamma = V_T \cup V_N$
- $\langle q_0, \alpha, Z_0 \rangle \vdash \langle q_1, \alpha, ZZ_0 \rangle$, dove Z è l'assioma di G
- $\langle q_1, a\alpha, A\gamma \rangle \vdash \langle q_1, a\alpha, \beta\gamma \rangle$ se $\langle A ::= \beta \rangle$ è una produzione, $a \in G(\langle A ::= \beta \rangle)$
- $\langle q_1, a\alpha, a\gamma \rangle \vdash \langle q_1, \alpha, \gamma \rangle$
- $\langle q_1, \#, Z_0 \rangle \vdash \langle q_F, \#, Z_0 \rangle$

Esempio

Le seguenti grammatiche sono $LL(1)$:

- $G_6 = \{ S ::= Ab \mid a,b, S ::= cS \mid c, A ::= aA \mid a, A ::= \epsilon \mid b \}$
- $G_7 = \{ S ::= aBb \mid a, S ::= BS \mid bc, B ::= b \mid b, S ::= c \mid c \}$

Alcune grammatiche non $LL(1)$ possono essere trasformate in grammatiche $LL(1)$ debolmente equivalenti con operazioni di:

- fattorizzazione sinistra
- sostituzione sinistra
- eliminazione di ricorsione a sinistra

Trattamento degli errori sintattici

Obiettivi:

- consentire la prosecuzione dell'analisi del testo
- suggerire correzioni

Qualunque trasformazione di una stringa in un'altra può essere ottenuta mediante una o più applicazioni delle seguenti operazioni di edizione :

- cancellazione del simbolo x nella posizione i : $[x]_i$.
- inserimento del simbolo x nella posizione i : $[x]_{i+}$.
- sostituzione di x con y nella posizione i : $[x/y]_i$.
- scambio di due simboli adiacenti nelle posizioni i e $i+1$: $[x\leftrightarrow y]_i$.

Date due stringhe α e β , si definisce distanza di edizione $d(\alpha,\beta)$ il numero di operazioni di edizione necessario per trasformare α in β (eventualmente si attribuisce un peso ad ogni operazione e si considera la somma pesata).

- Esempio* : $\alpha = \text{"loop"}, \beta = \text{"loot"}, \gamma = \text{"cloo"}, \delta = \text{"loto"}$
- $d(\alpha,\beta) = 1$: $[p/t]_4$
 - $d(\alpha,\gamma) = 2$: $[c]_{1+}, [p]_5$.
 - $d(\alpha,\delta) = 2$: $[p/t]_4, [\alpha\leftrightarrow t]_3$ oppure $[o/t]_3, [p/o]_4$

Si definisce correzione a distanza minima (c.d.m) la sostituzione di una stringa errata (i.e. non appartenente al linguaggio) con una scelta fra quelle corrette a distanza minima.

La c.d.m. è adeguata per la correzione di errori ortografici (individuazione della parola corretta dato il dizionario) e può essere utilizzata in fase di analisi lessicale.

Un analizzatore sintattico in grado di effettuare la c.d.m. risulta in genere non deterministico, con complessità $O(n^3)$, quindi non accettabile.

Correzione con analizzatori deterministici

Si definisce errore di prefisso , nell'analisi da sinistra a destra, il primo simbolo che fa sì che la parte di stringa già analizzata non sia prefisso di alcuna frase del linguaggio.

Esempio

- frase corretta: IF A = B THEN
- frase errata: IFA = B THEN
prefisso valido: IFA = B
errore di prefisso: THEN

Gli analizzatori deterministici segnalano l'errore di prefisso e tentano una correzione locale.

Analisi deterministica bottom-up

Situazione dell'analizzatore sintattico al momento della segnalazione dell'errore:

- contenuto dello stack: $Z_0 a_1 h_1 a_2 h_2 \dots a_n h_n$
- stringa residua in ingresso: $b_0 b_1 \dots \#$
- non esiste transizione $t_p(h_n, b_0)$
- se l'analizzatore ha prospezione > 0 , b_0 non appartiene al contesto di alcuna riduzione

Strategia 1:

si rimuovono (se necessario) gli elementi dallo stack fino ad incontrare uno stato h_i tale che esista la transizione $t_p(h_i, A) = h_k$, con A non terminale, e da h_k possa proseguire l'analisi con b_0 ; si simula il riconoscimento di A e si prosegue. Se non esiste alcun stato h_i , si ripristina lo stack e si rimuove b_0 , ripetendo l'operazione precedente con il simbolo b_1 , e così via fino alla fine della stringa.

Strategia 2:

si inserisce prima di b_0 un simbolo terminale a tale che esista la transizione $t_p(h_i, a) = h_k$ e da h_k possa proseguire l'analisi con b_0

Strategia 3:

si aggiungono alla grammatica una o più produzioni di errore, della forma:

$A ::= \langle \text{error} \rangle a$

- dove $\langle \text{error} \rangle$ è considerato terminale e $a \in V_T \cup \{\epsilon\}$. Nel corso dell'analisi, in presenza di un errore si procede come segue:
- viene inserito $\langle \text{error} \rangle$ prima di b_0 ;
 - si rimuovono (se necessario) gli elementi dallo stack fino ad incontrare uno stato h_i tale che esista la transizione $t_p(h_i, \langle \text{error} \rangle) = h_k$
 - si rimuovono (se necessario) i simboli della stringa in ingresso fino ad incontrare il simbolo $b_1 = a$
 - si effettua la riduzione secondo la produzione di errore e si rimuovono (se necessario) i simboli della stringa in ingresso fino ad incontrare un simbolo che consenta di riprendere l'analisi.

Analisi deterministica LL(1)

Dato l'automa $A_{LL(1)}(G)$ si possono avere le seguenti situazioni di errore:

- la configurazione corrente di $A_{LL(1)}(G)$ è $\langle q_1, b_0 b_1 \dots \#, A \gamma \rangle$, $b_0 \notin G(\langle A ::= \beta \rangle)$ per qualunque produzione avente A come parte sinistra.
- la configurazione corrente di $A_{LL(1)}(G)$ è $\langle q_1, b_0 b_1 \dots \#, a \gamma \rangle$, $b_0 \notin a$.

Si definisce sincrotripila la terna $\langle b, c, A \rangle$, con $A \in V_N$, $b \in MA(A)$, $c \in MC(A)$.

Metodo delle sincrotripile: ad ogni sottoalbero dell'albero sintattico può essere associata una lista di sincrotripile attivabili, detta sincrotabella. L'analizzatore sintattico può trovarsi in uno *stato di analisi* o in uno stato di errore: entrambi gli stati possono essere attivati ricorsivamente.

Stato di analisi attivato dalla marca a_a della sincrotripila $\langle a_a, a_c, A \rangle$:

- se incontra la marca a_c esce dallo stato di analisi e ritorna nello stato di errore precedente (termina se non esistono stati di errore precedenti)
- se incontra un errore e b_a è l'ultima marca di apertura incontrata, attiva un stato di errore con sincrotripila corrente $\langle b_a, b_c, B \rangle$.

Stato di errore; ignora tutti i simboli in ingresso finché:

- incontra la marca di chiusura della sincrotripila corrente e ritorna nello stato di analisi precedente
- incontra la marca di apertura della sincrotripila $\langle c_a, c_c, C \rangle$ ed attiva un nuovo stato di analisi
- incontra il simbolo '#' e termina

Vincoli sulla definizione delle sincrotripile della stessa sincrotabella

- 1) due sincrotripile non possono avere la stessa marca di apertura
- 2) la marca di apertura di una sincrotripila non può essere marca di chiusura della stessa o di un'altra sincrotripila
- 3) due sincrotripile non possono essere parzialmente sovrapposte: data le sincrotripile $\langle a_a, a_c, A \rangle$ e $\langle b_a, b_c, B \rangle$, per ogni stringa α tale che $A \Rightarrow^* a_a \alpha a_c$ $b_a \in \alpha$ se e solo se $b_c \in \alpha$.

Traduzioni sintattiche

Dati due linguaggi L_1 e L_2 , la traduzione da L_1 (linguaggio sorgente) a L_2 (linguaggio oggetto) può essere definita come una relazione binaria $L_1 \tau L_2$.

Date due grammatiche G_1 e G_2 , tali che:

- $V_{N1} = V_{N2} = V_N$, $Z_1 = Z_2 = Z$
- per ogni $\langle A::= \alpha_1 \rangle \in P_1$ esiste una $\langle A::= \alpha_2 \rangle \in P_2$ tale che $\pi_{V_N}(\alpha_1) = \pi_{V_N}(\alpha_2)$ [dove $\pi_{V_N}(\alpha)$ è la proiezione di α su V_N] e viceversa. Le coppie di produzioni $\langle A::= \alpha_1, A::= \alpha_2 \rangle$ si dicono corrispondenti.

L'insieme delle coppie di produzioni corrispondenti costituisce uno schema di traduzione.

Le due grammatiche possono essere fuse in un'unica grammatica di traduzione $G_1 = \langle V_{T1} \cup V_{T2}, V_N, P, Z \rangle$ tale che:

- $P_1 = \pi_{V_T}(P_1)$, $P_2 = \pi_{V_T}(P_2)$
- $L(G_1) = \{ \alpha \mid \gamma \in L(G_1) \text{ e } \alpha = \pi_{V_T}(\gamma) \}$
- $L(G_2) = \{ \alpha \mid \gamma \in L(G_2) \text{ e } \alpha = \pi_{V_T}(\gamma) \}$
- $\beta = \tau(\alpha)$ se e solo se esiste $\gamma \in L(G_1)$ tale che: $\gamma \in L(G_1)$, $\alpha = \pi_{V_T}(\gamma)$, $\beta = \pi_{V_T}(\tau(\gamma))$ dove $V_1 = V_{T1} \cup V_N$ e $V_2 = V_{T2} \cup V_N$

Esempio :

Schema di traduzione	Grammatica di trad.
G_1 $E ::= E + T$ $E ::= T$ $T ::= T * F$ $T ::= F$ $F ::= (E)$ $F ::= i$	G_2 $E ::= \oplus E T$ $E ::= T$ $T ::= \otimes T * F$ $T ::= F$ $F ::= (E)$ $F ::= id$
	G_1 $E ::= \oplus E + T$ $E ::= T$ $T ::= \otimes T * F$ $T ::= F$ $F ::= (E)$ $F ::= id$

Traduzioni postfixe

Uno schema di traduzione è detto postfixo se le produzioni della grammatica oggetto sono del tipo:

$$A ::= \gamma \omega \text{ dove } \gamma \in (V_{T1} \cup V_N)^*, \omega \in V_{T2}$$

Dato uno schema di traduzione non postfixo, è sempre possibile definirne uno equivalente postfixo; la produzione $A ::= \alpha \beta$, con $\alpha, \beta \in (V_{T1} \cup V_{T2} \cup V_N)^*$, $b \in V_{T2}$, è equivalente alla coppia: $A ::= \alpha B \beta$, $B ::= \epsilon b$

Traduttori deterministici bottom-up

Richiedono grammatiche di traduzione basate su schemi postfissi. Un traduttore si ottiene modificando un analizzatore come segue:

- in corrispondenza dell'operazione di riduzione secondo la candidata $A ::= \gamma^* \omega$, dove $\gamma \in (V_{T1} \cup V_N)^*$, $\omega \in V_{T2}^*$ si produce la stringa ω .

Traduttori deterministici top-down

Possono operare anche con schemi di traduzione non postfissi. Un traduttore si ottiene modificando un analizzatore come segue:

- se in cima allo stack vi è un simbolo $b \in V_{T2}$, il traduttore lo estrae dallo stack e lo aggiunge alla stringa di uscita.