

Linguaggi di programmazione

Appunti

GioLaPalma

Contents

Chapter 1

Introduzione ai linguaggi normali: grammatiche

1.1 Linguaggi (naturali o artificiali)

La descrizione di un linguaggio avviene su 3 dimensioni:

- **Sintassi:** regole di formazione, ovvero la relazione tra segni
- **Semantica:** attribuzione di significato
- **Pregmatica:** in quale modo frasi corrette e sensate sono usate
- **Implementazione:** come eseguire una frase corretta rispettandone la semantica

Partiamo con il descrivere questi elementi

1.1.1 Sintassi

Definition 1.1.1: sintassi

La **sintassi** è la parte della grammatica che studia la struttura delle frasi e il modo in cui le parole si combinano per formare enunciati corretti e significativi

Si dirama in diversi aspetti:

- **Aspetto lessicale** che riguarda le parole che si possono usare, quindi:
 - descrizione del lessico, intuitivamente i dizionari per i linguaggi naturali assolvono a questo compito
 - errori dovuti a vocaboli inesistenti
- **Aspetto grammaticale** che si riferisce nel modo in cui è possibile costruire frasi corrette è possibile costruire con il lessico. Le frasi grammaticalmente corrette possono essere costruite grazie all'uso delle regole grammaticali, che sono in numero finito, mentre il numero di frasi generabili è infinito.
Un errore grammaticale è una frase scorretta, anche se il lessico utilizzato è corretto. Es. "La cane abbaiano"

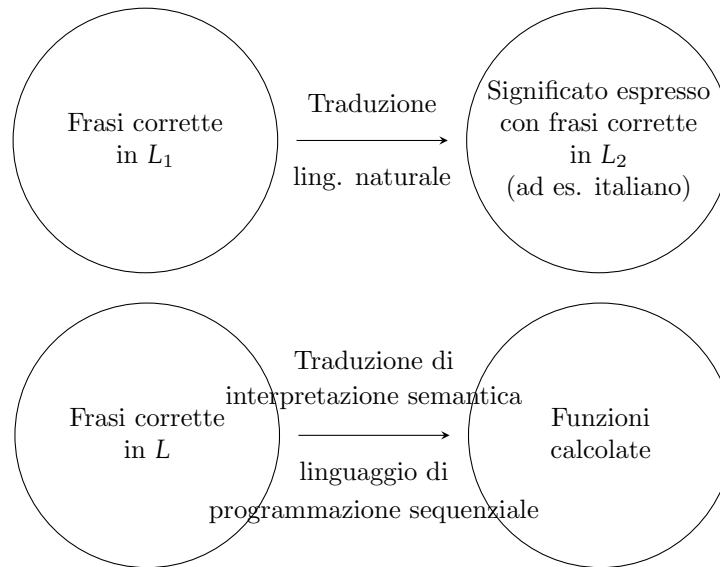
1.1.2 Semantica

Definition 1.1.2: Semantica

La **semantica** è la branca della linguistica che studia il significato delle parole, delle frasi e degli enunciati

- **Per il lessico** (quindi lo studio e il significato delle parole) bastano **i dizionari**
- **Per le frasi** è più complicato, **devo sapere** infatti:

1. a quale linguaggio appartiene la frase
2. su quale linguaggio basarmi per dare significato



1.1.3 Pragmatica

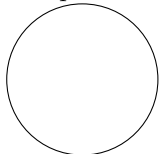
Definition 1.1.3: Pragmatica

La **pragmatica** è un insieme di regole che guidano l'uso e come i contesti influiscono sull'interpretazione di frasi sensate e corrette

Ad esempio quando e a chi dare del "tu" o del "lei" quando ci rivolgiamo a delle persone

1.1.4 Implementazione

L'implementazione è l'esecuzione di una frase sintatticamente corretta rispettandone la semantica



La semantica di P , è la funzione f che è pure la semantica del programma compilato Q , quindi l'implementazione Q di P preleva la semantica di P !

1.2 Lessico e frasi di un linguaggio

Innanzitutto diamo tre definizioni

Definition 1.2.1: alfabeto

Un **alfabeto** è un insieme (tipicamente) finito i cui elementi sono detti simboli

Definire l'alfabeto ci porta alla definizione di lessico:

Definition 1.2.2: Lessico

Il **lessico** è un insieme di sequenze finite costituite con caratteri o simboli dell'alfabeto

Il quale ci porta alla definizione di frase:

Definition 1.2.3: frase

Una **frase** è un insieme sequenze finite contruite con parole del lessico.

È, quindi, facile notare che **il lessico è un alfabeto per le frasi**
Si ci si può ora astrarre e definire un linguaggio formale:

Definition 1.2.4: linguaggio formale

Un **linguaggio formale** L su alfabeto A è un sottoinsieme di A^* ($L \subseteq A^*$), dove:

$$A^* = \bigcup_{n \geq 0} A^n \quad \text{dove } A^0 = \{\epsilon\}$$

e

$$A^{n+1} = A \cdot A^n \quad n \geq 0$$

con $A \cdot A^n = \{aw | a \in A \wedge w \in A^n\}$

Si osservi che A^* è un **insieme infinito contabile** dato un ordinamento \preceq sui simboli di A , possiamo elencare tutte le parole come segue:

1. elenco la parola vuota ϵ (A^0)
2. poi elenco le parole di lunghezza 1 (A^1) secondo l'ordinamento \preceq
Es. a, b, c, \dots
3. poi elenco le parole in A^2 secondo \preceq
Es. $aa, ab, ac, \dots, ba, bb, bc, \dots$
4. così via

Anche se l'alfabeto A fosse infinito (quindi $A = \{a_0, a_1, \dots\}$) A^* sarebbe ancora contabile, ovvero esisterebbe la possibilità di elencare tutte le possibili parole in A^* . Infatti, **esiste una biezione tra A e \mathbb{N}** e riguardo ai numeri naturali si sa che:

- $\mathbb{N} \times \mathbb{N}$ (prodotto cartesiano) è numerabile. La dimostrazione viene fatta attraverso il *dove-tailing*, una tecnica comune per dimostrare la numerabilità di coppie di numeri naturali. Questa dimostrazione introduce la cosiddetta **funzione di decodifica**

$$f^2(x_1, x_2) = \frac{(x_1 + x_2)(x_1 + x_2 + 1)}{2} + x_2$$

che presi due numeri in \mathbb{N} completa la seguente tabella ordinando i numeri naturali in "diagonale"

		x_1						
		0	1	2	3	4	5	...
x_2	0	0	1	3	6	10	15	
	1	2	4	7	11	16		
	2	5	8	12	17			
	3	9	13	18				
	4	14	19					
	5	20						

è pertanto vero, quindi, che

$$f^2 : \mathbb{N} \times \mathbb{N} \rightarrow \mathbb{N}$$

è biunivoca

- \mathbb{N}^k è numerabile. Infatti si può dimostrare attraverso questo algoritmo:

$$f^k(x_1, x_2, \dots, x_k) = \text{if } (k = 2) \text{ then } f^2(x_1, x_2) \quad \text{else } f^2(x_1, f^{k-1}(x_2, \dots, x_k))$$

Ovvero riduce una funzione con k variabili nel dominio ad una serie di funzioni matrisoska per ricondurla alla forma $f^2 : \mathbb{N} \times \mathbb{N} \rightarrow \mathbb{N}$

- $\mathbb{N}^* = \bigcup_{k \geq 0} \mathbb{N}^k$ è numerabile. Infatti:

$$f(x_1, \dots, x_k) = f^2(k, f^k(x_1, \dots, x_k))$$

1.3 Notazioni e definizioni ausiliarie

Vengono riportate qui alcune definizioni/notazioni

1.3.1 Lunghezza

Definition 1.3.1: lunghezza

La **lunghezza** di una parola o stringa è definita per induzione così:

- Caso $|\epsilon|$: 0
- Caso $|aw|$: $1 + |w|$

Es: $|abc| = 2$

1.3.2 Concatenazione

Definition 1.3.2: Concatenazione

La **concatenazione** xy tra una stringa x e y , è la parola ottenuta giustapponendo x e y . Formalmente:

$$w = xy \iff \begin{cases} |w| = |x| + |y| \\ w(j) = x(j) & \text{per } 1 \leq j \leq |x| \\ w(|x| + j) = y(j) & \text{per } 1 \leq j \leq |y| \end{cases}$$

Dove $w(j)$ indica il j -esimo simbolo di w

Questi sono le leggi della concatenazione

- Associatività: $x(yz) = (xy)z$
- Elemento neutro (ϵ): $x\epsilon = x = \epsilon x$

1.3.3 Sottostringa

Definition 1.3.3: Sottostringa

La stringa v si dice **sottostringa** di $w \iff \exists x, y \in A^*$ t.c. $w = xvy$ dove x e y possono essere ϵ

Si osservi, quindi, che:

- Ogni stringa è sottostringa di se stessa
- ϵ è sottostringa di ogni stringa

1.3.4 Suffisso

Definition 1.3.4: suffisso

v si dice **suffisso** di $w \iff \exists x \in A^*. w = xv$

1.3.5 Prefisso

Definition 1.3.5: prefisso

v si dice **prefisso** di $w \iff \exists x \in A^*. w = vx$

1.3.6 Potenza n-esima

Definition 1.3.6: potenza n-esima

Si dice **potenza n-esima** di una stringa w il valore $n \geq 0$ il cui significato è definito per induzione:

- Caso 0: $w^0 = \epsilon$
- Caso $n + 1$: $w^{n+1} = ww^n$

1.3.7 Linguaggio

Definition 1.3.7: linguaggio

Si dice **linguaggio** L su alfabeto A un sottoinsieme $L \subseteq A^*$

Vengono riportati qui alcuni esempi

Example 1.3.1

Se $A = \{a\}$, si possono avere:

- $\emptyset, \{\epsilon\}, \{a, aaa\}$ sono linguaggi finiti
- $L_1 = \{a^n \mid n \geq 1\} = \{a, aa, aaa, \dots\} = A^* \setminus \{\epsilon\}$
- $L_2 = \{a^{2n} \mid n \geq 0\} = \{\epsilon, aa, aaaa\}$

1.4 Operazione sui linguaggi

Qui sono elencati le varie operazioni

1.4.1 Complemento

Definition 1.4.1: complemento

È definito **complemento** il linguaggio completare ad un linguaggio L , ovvero:

$$\bar{L} = \{w \in A^* \mid w \notin L\} = A^* \setminus L$$

1.4.2 Unione e intersezione

Definition 1.4.2: unione e intersezione

Ovvi:

$$\begin{aligned} L_1 \cup L_2 &= \{w \mid w \in L_1 \vee w \in L_2\} \\ L_1 \cap L_2 &= \{w \mid w \in L_1 \wedge w \in L_2\} \end{aligned}$$

1.4.3 Concatenazione

Definition 1.4.3: concatenazione

È definita **concatenazione** tale operazione:

$$L_1 \cdot L_2 = \{w_1 w_2 \mid w_1 \in L_1 \wedge w_2 \in L_2\}$$

Ecco alcuni esempi:

Example 1.4.1

- $L_1 = \{a^n \mid n > 0\} \quad L_2 = \{b\}$
 $L_1 \cdot L_2 = \{a^n b \mid n > 0\}$
- $L_1 = \{a^{2n} \mid n > 0\} \quad L_2 = \{b^m \mid m > 0\}$
 $L_1 \cdot L_2 = \{a^{2n} b^m \mid n, m > 0\}$
- $L_1 = \{a^m b^n \mid n > 0\} \quad L_2 = \{b^m \mid n > 0\}$
 $L_1 \cdot L_2 = \{a^m b^{n+m} \mid n, m > 0\}$
 $= \{a^m b^n \mid m \geq n > 0\}$
- $L_1 = \{a^m \mid m > 1\} \quad L_2 = \{a^m b^m \mid n > 0\}$
 $L_1 \cdot L_2 = \{a^{m+n} b^m \mid m > 1, m > 0\}$
 $= \{a^m b^m \mid m > m > 0\}$

- $L_1 = \{ab^m \mid n \geq 1\}$ $L_2 = \{a, c\} \cup \{b^n \mid n \geq 1\}$
 $L_1 \cdot L_2 = \{ab^m a \mid m \geq 1\} \cup \{ab^m c \mid n \geq 1\}$
 $\cup \{ab^n \mid n \geq 2\}$
- $A = \{0, 1\}$
 $L_1 = \{w \in A^* \mid w \text{ contiene un numero pari di "0"}\}$
 $L_2 = \{w \in A^* \mid w = 0y \text{ e } y \in \{1^*\}\}$
 $L_1 \cdot L_2 = \{w \in A^* \mid w \text{ ha un numero dispari di "0"}\}$

1.4.4 Potenza di un linguaggio

Definition 1.4.4: Potenza di un linguaggio

La **potenza di un linguaggio** viene definita per induzione:

- Caso 0: $L^0 = \{\epsilon\}$
- Caso $n + 1$: $L \cdot L^n \quad \forall n \geq 0$

1.4.5 Stella di kleene

Definition 1.4.5: stella di kleene

Si dice **stella di kleene**:

$$L^* = \bigcup_{n \geq 0} L^n$$

Oppure

$$L^+ = \bigcup_{n \geq 1} L^n$$

Quest'ultima detta **chiusura positiva**

1.5 Definire finitamente un linguaggio

1.5.1 esempio 1: frasi palindrome

Una frase palindroma è una parola che letta da sx a dx è uguale a se stessa letta da dx a sx
 Es. "I topi non avevano nipoti"

- $A = \{a, b\}$ $L = \{\epsilon, a, b, aa, bb, aba, bab, dots\}$. Come si nota è piuttosto scomodo
- Una palindroma può essere:
 - o è la stringa ϵ
 - oppure a
 - oppure b
 - oppure a "palindroma" a
 - oppure b "palindroma" b
- rappresentazione tramite **Backus-aur form (BNF)**.

$$\langle P \rangle := \epsilon \mid a \mid b \mid a\langle P \rangle a \mid b\langle P \rangle b$$

- Come **grammatica**:

$$P \rightarrow \epsilon \mid a \mid b \mid aPa \mid bPb$$

- definizione ricorsiva in cui:
 - P è detto *simbolo non terminale*
 - a, b sono "simboli terminali"

1.5.2 Esempio 2

espressioni aritmetiche formate a partire dalle variabili a e b con gli operatori $\times, +$ e le parentesi $(,)$
 Una *expr* può essere:

- - la variabile a
 - la variabile b
 - $expr \times expr$
 - $expr + expr$
 - $(expr)$

- **bnf:**

$$\langle E \rangle ::= a \mid b \mid \langle E \rangle \times \langle E \rangle \mid \langle E \rangle + \langle E \rangle \mid (\langle E \rangle)$$

- Grammatica:

$$E \rightarrow a \mid b \mid E \times E \mid E + E \mid (E)$$

1.6 Grammatiche

Una grammatica è un insieme di regole che descrivono come le parole e le frasi possono essere combinate per formare espressioni valide. Queste regole determinano la struttura sintattica di un linguaggio, specificando come le unità di base (come le parole o i simboli) si connettono per formare frasi o espressioni più complesse. Ogni grammatica **segue lo stesso pattern definito** differenziandosi solo per come sono caratterizzate le produzioni. Quelle più utili sono le cosiddette grammatiche libere (in rapporto tra facilità di analisi ed espressività)

Definition 1.6.1: grammatiche libere

Una **grammatica libera** da contesto è una quadrupla (NT, T, R, S) dove:

- NT è un insieme finito di simboli non terminali
- T è un insieme finito di simboli terminali
- $S \in NT$ è detto simbolo iniziale
- R è un insieme finito di produzione (o regole) della forma:

$$V \rightarrow w \text{ dove } V \in NT \wedge w \in (T \cup NT)^*$$

Alcuni esempi:

Example 1.6.1

$$G = (\{S\}, \{a, b, +, \times\}, S, R)$$

Con

$$R = \{S \rightarrow a, S \rightarrow b, S \rightarrow S + S, S \rightarrow S \times S\}$$

1.7 Derivazioni

Definition 1.7.1: derivazione immediata

Data $G = (NT, T, R, S)$ libera dal contesto, diciamo che da v si **deriva immediatamente** w , e lo denotiamo con $v \Rightarrow w$, se:

$$\frac{v = xAy \quad (A \rightarrow z) \in R \quad w = xzy}{v \Rightarrow w}$$

Definition 1.7.2: derivazione

Diciamo che da v si **deriva** w (o anche " v si riscrive in w "), e lo denotiamo con $v \Rightarrow^* w$, se esiste una sequenza finita (eventualmente vuota) di derivazione immediate

$$v \Rightarrow w_0 \Rightarrow w_1 \Rightarrow \dots \Rightarrow w$$

Cioè:

$$\frac{}{v \Rightarrow^* v} \quad \frac{v \Rightarrow^* w \quad w \Rightarrow z}{v \Rightarrow^* z}$$

Dove \Rightarrow^* è la chiusa riflessiva e transitiva della relazione \Rightarrow

1.8 Linguaggio Generato

Definition 1.8.1: Linguaggio Generato

Il **linguaggio generato** da una grammatica $G = (NT, T, R, S)$ è l'insieme

$$L(G) = \{w \in T^* \mid S \Rightarrow^* w\}$$

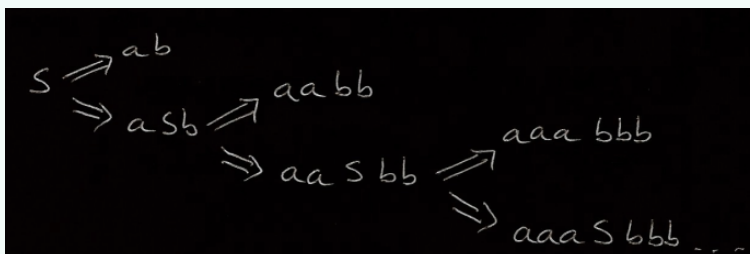
1.8.1 Algoritmo di Naif

Data una grammatica G è opportuno chiedersi come si fa a determinare un linguaggio $L(G)$ e a verificare se $w \in L(G)$. La domanda può essere complessa, tuttavia in casi semplici ci viene in aiuto l'**algoritmo di Naif** che consiste nel partire da S e provare ad applicare in tutti i modi possibili le produzioni (regole) per trovare una derivazione che generi w

In certi casi questa verifica è semplice

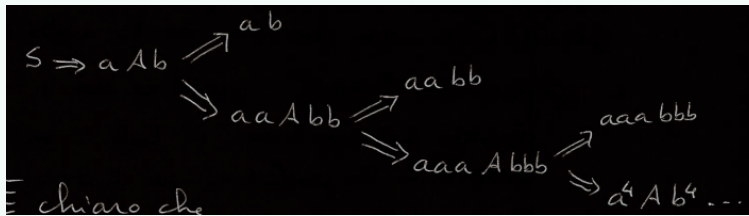
Example 1.8.1

- Sia G_3 con $S \rightarrow aSb|ab$



In questo esempio è facile determinare che $L(G_3) = \{a^n b^n \mid n \geq 1\}$

- Sia $G_1 \rightarrow aAb$ e $A \rightarrow aAb \mid \epsilon$



Quindi $L(G_1) = \{a^n b^n \mid n \geq 1\}$

Si può notare che G_1 e G_3 sono grammatiche equivalenti perché $L(G_1) = L(G_3)$

In generale **esistono grammatiche diverse che generano lo stesso linguaggio**

1.9 Alberi di derivazione

Per rappresentare graficamente e semplicemente una certa grammatica esiste uno strumento utilissimo, ovvero l'albero di derivazione

Definition 1.9.1: albero di derivazione

Data una grammatica libera $G = (NT, T, S, R)$, un albero di derivazione (o di parsing) è un albero ordinato in cui:

- Ogni nodo è etichettato con un simbolo in $NT \cup \{\epsilon\} \cup T$
- la radice è etichettata con S
- ogni nodo interno è etichettato con un simbolo in NT
- se il nodo n
 - ha etichetta $A \in NT$
 - i suoi figli sono nell'ordine m_1, \dots, m_k con etichetta x_1, \dots, x_k (in $NT \cup T$), allora

$A \rightarrow x_1, \dots, x_k$ è una produzione in R

- se il nodo n ha etichetta ϵ , allora n è una foglia, è figlio unico e, dato A suo padre, $A \rightarrow \epsilon$ è una produzione di R
- se inoltre ogni nodo foglia è etichettato su $T \cup \{\epsilon\}$ è una produzione di R
- se inoltre ogni nodo foglia è etichettato su $T \cup \{\epsilon\}$, allora l'alberello di derivazione corrisponde ad una derivazione completa

Un albero di derivazione, quindi, **riassume tante derivazioni diverse ma tutte equivalenti** (ovvero generano lo stesso albero)

Example 1.9.1

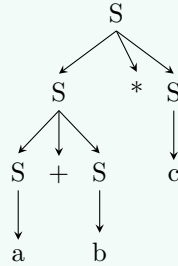
Consideriamo la grammatica

$$s \rightarrow a|b|c|S + S|S \times S$$

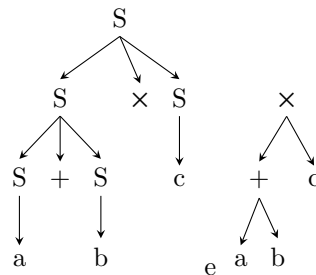
Inoltre si consideri la derivazione

$$S \Rightarrow \underline{S} + S \Rightarrow \underline{S} \times S + S \Rightarrow a \times \underline{S} + S \Rightarrow a \times b + \underline{S} \Rightarrow a \times b + c$$

Allora il suo albero di derivazione è:



Si osservi inoltre che l'albero di derivazione fornisce informazioni semantiche: "quali operandi per quali operatori" e possono essere riassunti nei cosiddetti **alberi sintattici**, tipo



Theorem 1.9.1

Una stringa $w \in T^*$ appartiene a $L(G)$ sse ammette un albero di derivazione completo (le cui foglie, lette da sx a dx danno la stringa w) cioè visita in ordine anticipato tralasciando i nonterminali

1.9.1 Ambiguità

Per spiegare cos'è l'ambiguità partiamo con un esempio

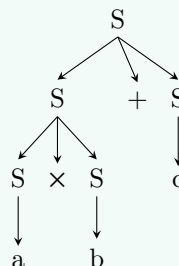
Example 1.9.2

Si consideri la seguente grammatica

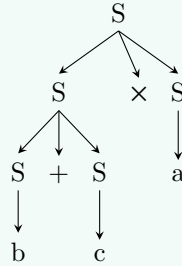
$$S \rightarrow a|b|c|S + S|S \times S$$

Si hanno due diverse derivazioni per $a \times b + c$:

1. $S \Rightarrow \underline{S} + S \Rightarrow \underline{S} \times S + S \Rightarrow a \times \underline{S} + S \Rightarrow a \times b + \underline{S} \Rightarrow a \times b + c$
con l'albero:



2. $S \Rightarrow \underline{S} \times S \Rightarrow a \times \underline{S} \Rightarrow a \times \underline{S} + s \Rightarrow a \times b + \underline{S} \Rightarrow a \times b + c$
Con l'albero:



Per questa grammatica, la stringa $a \times b + c$ ha più di un albero di derivazione in questi casi si dice che la grammatica è quindi **ambigua** e **inutilizzabile per dare semantica a $a \times b + c$** . Bisogna utilizzare grammatiche non ambigue, o manipolare grammatiche ambigue per disambiguarle

Grammatica ambigua

Definition 1.9.2: Grammatica ambigua

Una grammatica libera G è **ambigua** se $\exists w \in L(G)$ che ammette più alberi di derivazione

Definition 1.9.3: Linguaggio ambiguo

Un linguaggio L è **ambiguo** se tutte le grammatiche G , tali che $L(G) = L$, sono ambigue

Alcune grammatiche possono essere manipolate di modo da

- rimuovere l'ambiguità
- generare lo stesso linguaggio

In altre, invece, ti tieni l'ambiguità (non è possibile rimuoverla (cazzo))

1.9.2 Rimuovere l'ambiguità

Questa CAZZATA viene spiegata con solo degli esempi

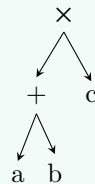
Example 1.9.3

Sia S la grammatica

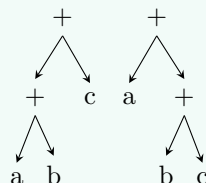
$$S \rightarrow a \mid b \mid c \mid S + S \mid S \times S \text{ è ambigua!}$$

Problemi:

- Precedenza del $*$ rispetto al $+$ in modo che $a \times b + c$ sia interpretato come:



- associatività del $+$ e del \times



$a + b + c$ ovvero bisogna scegliere l'associatività a dx o sx

Un modo per eliminare questa ambiguità è definire la grammatica così:

$$e \rightarrow E + T \mid T \quad T \rightarrow A \times T \mid A \quad A \rightarrow a \mid b \mid c \mid (E)$$

Chapter 2

Struttura di un compilatore, semantica statica, semantica dinamica

2.1 Vincoli contestuali

Definition 2.1.1: vincoli sintattici contestuali

I vincoli sintattici contestuali sono termini o parole riservate non esprimibili per mezzo di grammatiche libere (perché non possono descrivere vincoli che dipendono dal contesto) che bisogna evitare di considerare quando si esegue il codice

Tradizionalmente i vincoli sintattici contestuali appartengono alla sintassi, ma nel gergo dei *LP*, si intende:

- **Sintassi**: quello che si *scrive per mezzo di Grammatiche Libere*
- **semantica** tutto il resto ...

Pertanto *i vincoli contestuali sono dunque vincoli semantici, detti di semantica statica* cioè vincoli che possono essere verificati ispezionando il codice *senza mandare il programma in esecuzioni*

Il compilatore delega questi controlli di semantica statica alla cosiddetta **analisi semantica**

2.2 Semantica statica

Definition 2.2.1: Semantica statica

Per **semantica statica** si intende l'insieme di quei controlli che possono essere fatti sul testo del programma senza eseguirlo

Example 2.2.1

```
int A;  
bool B  
A := B (errore di tipo)
```

2.3 semantica dinamica

Per **semantica dinamica** si intende una rappresentazione formale dell'esecuzione del programma, la quale può mostrare errori durante l'esecuzione

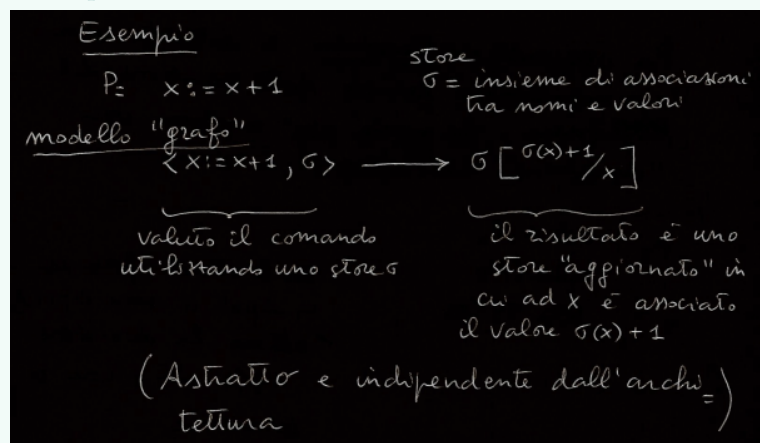
Example 2.3.1


```
read(A);  
B :=  $\frac{10}{A}$  (Se  $A = 0$  si dà un errore in esecuzione. F)
```

Staticamente non si può sapere l'errore perché la sua occorrenza **dipende dall'input dell'utente che fornirà durante l'esecuzione del programma**

Per implementare una semantica dinamica occorre fornire un modello matematico che descriva indipendentemente dall'architettura su cui il programma viene eseguito, il "comportamento del programma"

Example 2.3.2



2.3.1 Utilità della semantica dinamica

A chi serve la semantica dinamica?

- **Al programmatore:** *ANALISI DEL PROGRAMMA*
 - deve sapere esattamente cosa debba fare il suo programma
 - deve poter dimostrare proprietà del suo programma (ad es.: "termina sempre per ogni possibile input?")
- **Al progettista del linguaggio:**
 - strumento di specifica del linguaggio
 - deve poter dimostrare proprietà del linguaggio (ad es.: "è Turing-completo?")
- **All'implementatore del linguaggio:**
 - riferimento per dimostrare la correttezza dell'implementazione

Infatti **un compilatore è corretto quando preserva la semantica dinamica**, quindi per dimostrare che un compilatore è corretto serve avere una semantica per il linguaggio sorgente e per il linguaggio oggetto

2.3.2 definire la semantica

Per definire la semantica si utilizzano due tecniche principali:

- **operazionale:** (macchina astratta a stati e transizioni)
Ovvero si costruisce una specie di automa che, passo a passo, mostra l'effetto dell'esecuzione delle varie istruzioni. **vi è una maggiore enfasi su COME si calcola**
- **Denotazionale:** si associa ad ogni programma sequenziale una funzione da input ad output (incluse strutture ausiliarie e memoria). **vi è una maggiore enfasi su COSA si calcola**

2.4 Pragmatica nella descrizione di un linguaggio

Definition 2.4.1: Pragmatica nella descrizione di un linguaggio

si definisce **pragmatica nella descrizione di un linguaggio** insieme di regole sul modo in cui è meglio usare le istruzioni a disposizione

Esempietti:

Example 2.4.1

- evitare le istruzioni di salto quando possibile
- usare le variabili di controllo del `for` solo a quello scopo
- scelta della modalità più appropriata di passaggio di parametri ad una funzione
- scelta tra iterazione determinata (`for`) e indeterminata (`while`)

2.5 implementazione

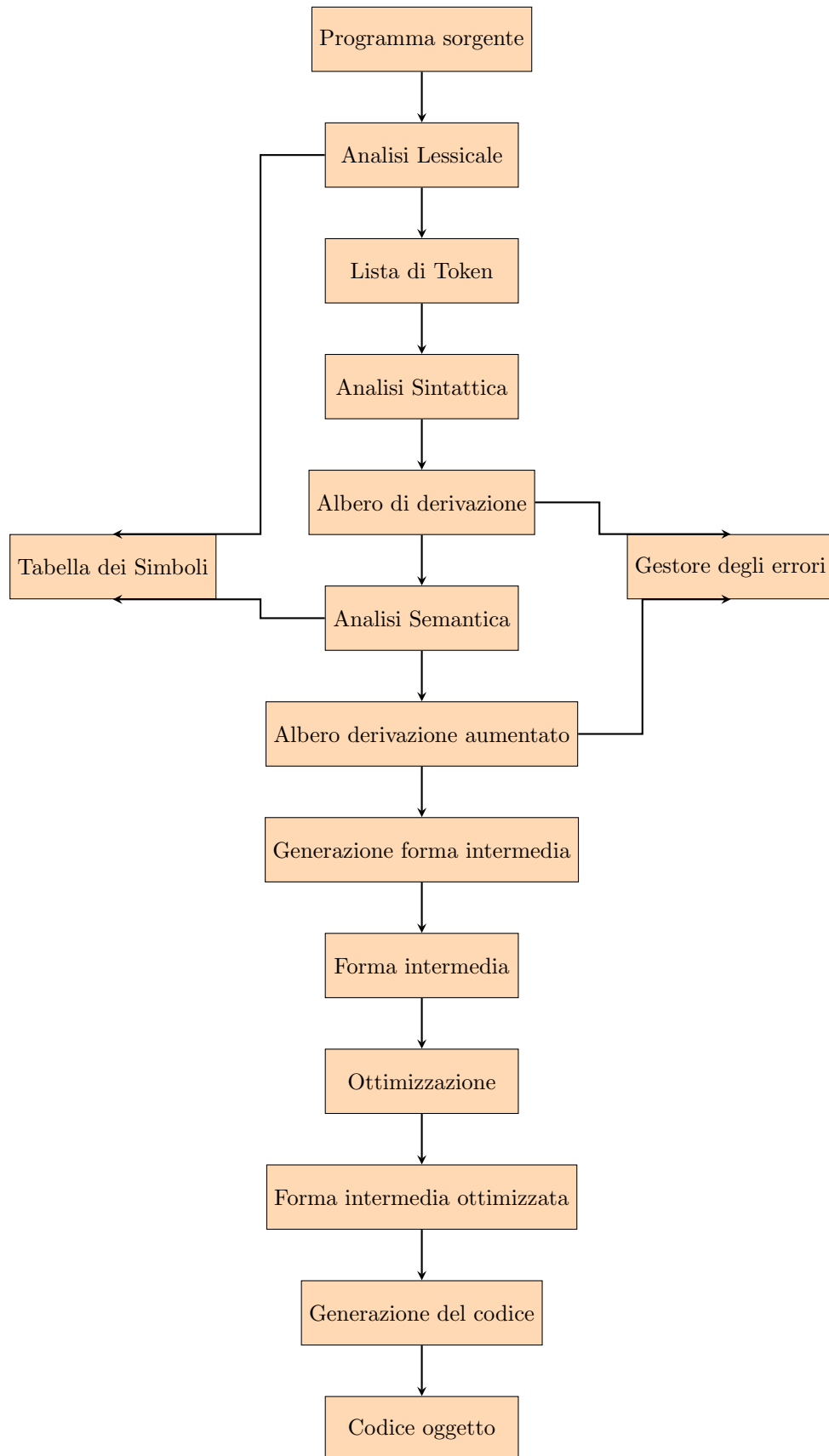
Definition 2.5.1: implementazione

Per **implementazione** si intende la scrittura di un compilatore per una macchina ospite già realizzata, costruendo così una macchina astratta per il linguaggio

2.5.1 Correttezza dell'implementazione

Per far sì che un compilatore sia corretto occorre **dimostrare che il programma preservi la semantica**, ovvero il programma sorgente e quello oggetto calcolino la stessa funzione

2.5.2 Struttura di un compilatore



2.6 fasi principali della compilazione

2.6.1 analisi lessicale (scanner)

L'analisi lessicale spezza il programma sorgente nei componenti sintattici primitivi chiamati "tokens" (identificatori, numeri, operatori, parametri, parole riservate)

- controlla solo che il lessico sia ammissibile
- riempie parzialmente la tabella dei simboli per gli identificatori di variabili, procedure funzioni ...

Per realizzare uno scanner avremo bisogno di studiare:

- **grammatiche regolari**
- **espressioni regolari**: un formalismo usato per descrivere i linguaggi generati da grammatiche regolari
- **automi a stati finiti**: uno strumento che permette di riconoscere i linguaggi regolari

2.6.2 analisi sintattica (parser)

A partire dalla lista di tokens, generata dallo scanner, il parser produce l'albero di derivazione del programma, riconoscendo se le frasi sono sintatticamente corrette

Ad esempio controlla che:

- le parentesi siano bilanciate: `((a)+b)))`
- che i comandi siano composti secondo le regole grammaticali `if(x=5) then then x:=3`

Per realizzare un Parser, avremo bisogno di:

- grammatiche libere dal contesto
- automi a pila

2.6.3 Analisi semantica

L'analisi semantica **esegue dei controlli di semantica statica** (ovvero sintattici contestuali) per rilevare eventuali errori semantici

Arricchisce l'albero di derivazione generato dal Parser con informazioni sui tipi, verifica i tipi negli assegnamenti, parametri attuali vs. formali, dichiarazione e uso di variabili e genera eventuali errori

2.6.4 Generazione della forma intermedia

Genera codice scritto in un **linguaggio intermedio** indipendente dall'architettura, facilmente traducibile nel linguaggio macchina di varie macchine diverse. Nel generare questo codice intermedio si esegue la struttura dell'albero sintattico, ricavato dall'albero di derivazione

2.6.5 Ottimizzazione

Si effettuano ottimizzazioni nel codice intermedio per renderlo più efficiente

- rimozione di codice inutile (dead code)
- espansione in linea di chiamate di funzioni
- fattorizzazione di sottoespressioni
- mettere fuori dai cicli sottoespressioni che non variano

Alla fine si ottiene un codice intermedio **ottimizzato**

2.6.6 Generazione del codice

Viene generato codice per una specifica architettura (include anche l'assegnazione dei registri e ottimizzazioni specifiche macchine)

2.6.7 Tabella dei simboli

Memorizza le informazioni sui nomi presenti nel programma (identificatori di variabili, funzioni, procedure)

Es: per le matrici mette, come attributo la dimensione e il tipo dei suoi elementi

2.7 semantica operativa strutturata

2.7.1 Definizione di un linguaggio a cui dare semantica

La **semantica operativa strutturata** È utilizzata per descrivere come ogni singola istruzione o espressione in un linguaggio modifica lo stato di un sistema in termini di transizioni di stato

Il suo **linguaggio** viene **definito tramite sintassi astratta semplice ed intuitiva, ma ambigua** ed una stringa viene sempre accoppiata ad un albero sintattico (non ambiguo)

Alcuni elementi fondamentali del linguaggio vengono definiti attraverso **insiemi di base**:

- **Booleani**: l'insieme dei valori booleani è composto da due valori: $\{\mathbf{tt}, \mathbf{ff}\}$. Le metavariable sono $t, t_1, t' \in \mathbf{T}$
- **numeri naturali**: $\{0, 1, 2, \dots\}$ $n, m, p \in \mathbf{N}$
- **variabili**: a, b, c, \dots, z $v \in \mathbf{Var}$

Per descrivere espressioni più complesse, vengono definiti alcuni **insiemi derivati** utilizzando la notazione BNF (Backus-Naur Form)

- **espressioni aritmetiche (exp)**:

$$e ::= m | v | e + e | e - e | e * e$$

- **espressioni booleane (Bexp)**:

$$b ::= t | e = e | b \text{ or } b | \neg b$$

- **Comandi Com**:

$$c ::= \text{skip} | v := e | c ; c | \text{while } b \text{ do } c | \text{if } b \text{ then } c \text{ else } c$$

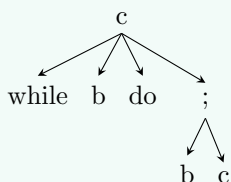
Questo tipo di sintassi è piuttosto semplice ma è ambigua, **per una sintassi non ambigua ne dovrei costruire una completa, ma molto più complicata** (dovrei gestire le precedenze, le parentesi ecc...) ma non serve nel dare una semantica in un linguaggio di programmazione perché un **parser (analizzatore sintattico) prende in input un programma scritto in sintassi concreta (non ambigua) e restituisce un albero sintattico di sintassi astratta** (quella che stiamo appena definendo), pertanto, nel dare semantica possiamo partire dagli alberi di sintassi astratta (ambigua) e ignorare la parte di analisi del parser

Example 2.7.1

Riportiamo qui un esempio di sintassi astratta.

Che tipo di albero sintattico vogliamo intendere con la seguente espressione?

while b do $c_1; c_2$



2.8 Dare semantica ad un linguaggio

Entriamo nel vivo del discorso, ma prima definiamo, per ogni categoria sintattica (cioè **Exp**, **Bexp**, **Com**) un modello detto **sistema di transizione** che è fondamentalmente un "grafo" di stati

Definition 2.8.1: sistema di transizione

Un **sistema di transizione** è una tripla $\langle \Gamma, T, \rightarrow \rangle$ dove

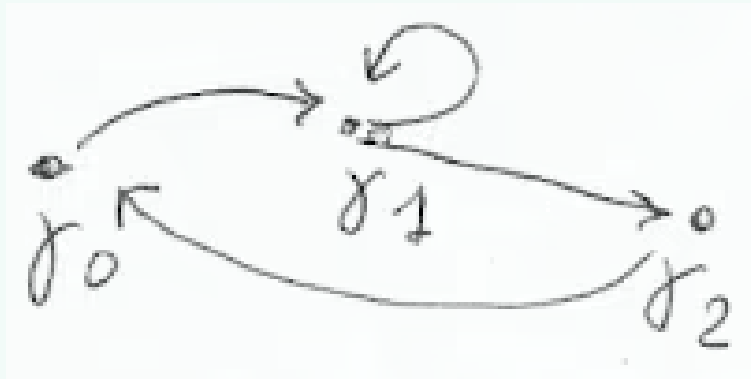
- Γ è l'insieme di stati (o configurazione)
- $T \subseteq \Gamma$ è l'insieme degli stati terminali (ovvero tutti quegli stati in cui il calcolo è stato terminato con successo)
- $\rightarrow \subseteq \Gamma \times \Gamma$ è la relazione di transazione che prende in input uno stato $\in \Gamma$ e restituisce un'altro stato $\in \Gamma$

Una computazione a partire dallo stato γ_0 è una sequenza $\gamma_0 \rightarrow \gamma_1 \rightarrow \gamma_2 \rightarrow \dots$ che può essere finita o infinita, invece con \rightarrow^* si indica la chiusura riflessiva e transitiva di \rightarrow , ovvero:

$$\frac{}{\gamma \rightarrow^* \gamma} \quad \frac{\gamma \rightarrow^* \gamma' \quad \gamma' \rightarrow \gamma''}{\gamma \rightarrow^* \gamma''}$$

ovvero si può raggiungere da uno stato γ uno stato γ'' in più passi

Example 2.8.1



Questa è una rappresentazione grafica di un grafo in cui i nodi sono gli stati e gli archi le transizioni

Se voglio definire la semantica (se voglio usare questo tipo di struttura) del linguaggio con la sintassi definita prima occorre definire uno stato di transazione specifico per **Exp**, per **Bexp** e per **Com**

Vi sono tuttavia **diversi problemmucci**, del tipo:

1. Γ è di solito un insieme infinito contabile, allora vi è la **necessità di trovare una rappresentazione finita ed implicita attraverso grammatiche**. Questo vuol dire che Γ coincide con uno dei linguaggi delle 3 categorie sintattiche (ovvero **Exp**, **Bexp**, **Com**)

Example 2.8.2

$$\Gamma_e = \{ \langle e, \sigma \rangle \mid e \in \text{Exp}, \sigma \in \text{Store} \}$$

Dove σ è una funzione che associa ad ogni variabile un numero naturale, perché lo stato del mio sistema è una coppia in cui la prima parte indica l'espressione che devo valutare, la seconda componente è lo store che indica il valore dell'espressione

2. $\rightarrow \subseteq \Gamma \times \Gamma$ è una relazione costituita da infinite coppie $\gamma \rightarrow \gamma'$, anche qui vi è la necessità di trovare una rappresentazione finita ed implicita come minima relazione che soddisfa **un certo insieme finito di assiomi e regole di inferenza**, quindi la semantica non è che un insieme di regole di inferenza che mi indicano in modo calcolare le transizioni che mi portano ad eseguire un certo comando
3. per dare significato alle variabili (che possono solo assumere valore su \mathbb{N}) è necessario introdurre uno **store** $\sigma : Var \rightarrow \mathbb{N}$, come funzione che associa ad ogni variabile un valore

$$\sigma = \{x_1/n_1, x_2/n_2, \dots, x_k/n_k\}$$

Se supponiamo che $var = \{x_1, x_2, \dots, x_k\}$

2.8.1 Semantica delle espressioni aritmetiche

Adesso introduciamo la **Semantica delle espressioni aritmetiche**, un tipo di semantica operativa. Deve ovviamente avere un sistema di transizione $\langle \Gamma_e, T_e, \rightarrow_e \rangle$ dove:

- $\Gamma_e = \{\langle e, \sigma \rangle \mid e \in Exp, \sigma \in Store\}$
- $T_e = \{\langle n, \sigma \rangle \mid n \in \mathbb{N}, \sigma \in Store\}$
- La relazione \rightarrow_e è definita come la minima relazione che soddisfa gli assiomi e le regole di inferenza qui sotto:

1. Variabile:

$$\overline{\langle v, \sigma \rangle \rightarrow_e \langle \sigma(v), \sigma \rangle}$$

Ovvero il tuo stato terminale sarà il numero della variabile v indicato dallo Store (inoltre σ rimane inalterato). Quindi valuto ciò che v vale in σ

2. Somma 1:

$$\frac{\langle e_0, \sigma \rangle \rightarrow_e \langle e'_0, \sigma' \rangle}{\langle e_0 + e_1, \sigma \rangle \rightarrow_e \langle e'_0 + e_1, \sigma' \rangle}$$

Nel momento in cui riesco a fare un passo di valutazione da e_0 a e'_0 alterando anche lo stato dello store da σ a σ' questa trasformazione si applica anche durante una somma, in altre parole l'espressione si semplifica o riduce (ad esempio, una variabile viene sostituita con il suo valore), e nel contempo lo stato della memoria potrebbe essere aggiornato se l'espressione stessa comporta una modifica ai valori delle variabili

3. Somma 2:

$$\frac{\langle e_1, \sigma \rangle \rightarrow_e \langle e'_1, \sigma' \rangle}{\langle m + e_1, \sigma \rangle \rightarrow_e \langle m + e'_1, \sigma' \rangle}$$

Stessa roba ma con un numero m

4. Somma 3:

$$\overline{\langle m + m', \sigma \rangle \rightarrow_e \langle P, \sigma \rangle} \quad \text{dove } P = m + m'$$

5. Sottrazione 1:

$$\frac{\langle e_0, \sigma \rangle \rightarrow_e \langle e'_0, \sigma' \rangle}{\langle e_0 - e_1, \sigma \rangle \rightarrow_e \langle e'_0 - e_1, \sigma' \rangle}$$

6. Sottrazione 2:

$$\frac{\langle e_1, \sigma \rangle \rightarrow_e \langle e'_1, \sigma' \rangle}{\langle m - e_1, \sigma \rangle \rightarrow_e \langle m - e'_1, \sigma' \rangle}$$

7. Sottrazione 3:

$$\overline{\langle m - m', \sigma \rangle \rightarrow_e \langle p, \sigma \rangle}$$

Si noti come la somma e la sottrazione prima valutano la sottoespressione di sinistra (e_0) con somma/sottrazione 1 poi, se questa s'è mutata in numero, valutano la sottoespressione di destra (e_1) con somma/sottrazione 2 ed infine, se questa s'è mutata in un numero, viene fatta la somma/sottrazione finale con somma/sottrazione 3

Example 2.8.3

Esempietto per valutare $\langle (x+2) - y, \{x/5, y/3\} \rangle$:

$$\begin{array}{l}
 \text{(Var)} \quad \frac{}{\langle x, \{x/5, y/3\} \rangle \rightarrow \langle 5, \{x/5, y/3\} \rangle} \\
 \text{(Sum}_1\text{)} \quad \frac{}{\langle x+2, \{x/5, y/3\} \rangle \rightarrow \langle 5+2, \{x/5, y/3\} \rangle} \\
 \text{(Sub}_1\text{)} \quad \frac{}{\underbrace{\langle (x+2) - y, \{x/5, y/3\} \rangle}_{\gamma_0} \rightarrow \underbrace{\langle (5+2) - y, \{x/5, y/3\} \rangle}_{\gamma_1}} \\
 \text{(Sum}_3\text{)} \quad \frac{}{\langle 5+2, \{x/5, y/3\} \rangle \rightarrow \langle 7, \{x/5, y/3\} \rangle} \\
 \text{(Sub}_1\text{)} \quad \frac{}{\underbrace{\langle (5+2) - y, \{x/5, y/3\} \rangle}_{\gamma_1} \rightarrow \underbrace{\langle 7 - y, \{x/5, y/3\} \rangle}_{\gamma_2}} \\
 \text{(Var)} \quad \frac{}{\langle y, \{x/5, y/3\} \rangle \rightarrow \langle 3, \{x/5, y/3\} \rangle} \\
 \text{(Sub}_2\text{)} \quad \frac{}{\underbrace{\langle 7 - y, \{x/5, y/3\} \rangle}_{\gamma_2} \rightarrow \underbrace{\langle 7 - 3, \{x/5, y/3\} \rangle}_{\gamma_3}} \\
 \text{(Sub}_3\text{)} \quad \frac{}{\underbrace{\langle 7 - 3, \{x/5, y/3\} \rangle}_{\gamma_3} \rightarrow \underbrace{\langle 4, \{x/5, y/3\} \rangle}_{\gamma_4}}
 \end{array}$$

Si ha, quindi, che $\gamma_0 \rightarrow \gamma_1 \rightarrow \gamma_2 \rightarrow \gamma_3 \rightarrow \gamma_4 \in T_e$ cioè $\langle (x+2) - y, \{x/5, y/3\} \rangle \rightarrow^* \langle 4, \{x/5, y/3\} \rangle$

Theorem 2.8.1

Vogliamo dimostrare che \rightarrow_e è deterministico, ovvero:

$$\gamma \rightarrow_e \gamma' \text{ e } \gamma \rightarrow_e \gamma'', \text{ allora } \gamma' = \gamma'' \quad \forall \gamma, \gamma', \gamma''$$

In altre parole significa che **da ogni transizione esce al più una transizione**, mai più di una

dimostrazione: mi riduco a dimostrare che $(\langle e, \sigma \rangle \rightarrow_e \gamma' \wedge \langle e, \sigma \rangle \rightarrow_e \gamma'') \Rightarrow \gamma' = \gamma''$

Procedo per induzione strutturale, con HP $(\langle e, \sigma \rangle \rightarrow_e \gamma' \wedge \langle e, \sigma \rangle \rightarrow_e \gamma'') \Rightarrow \gamma' = \gamma''$.

1. $e = m \in \mathbb{N}$: se $\langle e, \sigma \rangle \rightarrow_e e$ allora la conclusione è vera perché la premessa è falsa
2. $e = v \in Var$: Per la regola (Var), l'unica transizione derivabile per $\langle v, \sigma \rangle$ è $\langle v, \sigma \rangle \rightarrow_e \langle \sigma(v), \sigma \rangle$ poiché σ è una funzione (cioè $\sigma(v)$ è univoco) e la sola regola (Var) è applicabile allora per forza $\langle \sigma(v), \sigma \rangle = \sigma'$ e $\langle \sigma(v), \sigma \rangle = \sigma''$ quindi $\gamma' = \gamma''$
3. $e = e_0 + e_1$: Supponiamo che $\langle e_0 + e_1, \sigma \rangle \rightarrow \gamma'$ e $\langle e_0 + e_1, \sigma \rangle \rightarrow \gamma''$
Ci sono 3 sottocasi da esaminare in accordo nel modo in cui derivò $\langle e_0 + e_1, \sigma \rangle \rightarrow \gamma'$:
 - (a) $\langle e_0, \sigma \rangle \rightarrow \langle e'_0, \sigma' \rangle$ e $\gamma' = \langle e'_0 + e_1, \sigma' \rangle$ in questo caso ho che $e_0 \notin \mathbb{N}$ e la regola che ho applicato è **Somma** 1. Allora se $\langle e_0 + e_1, \sigma \rangle \rightarrow \gamma''$ è necessario che $\langle e_0, \sigma \rangle \rightarrow \langle e''_0, \sigma'' \rangle$ e che $\gamma'' = \langle e''_0 + e_1, \sigma'' \rangle$.
Tuttavia per (HP) si ha che $\langle e'_0, \sigma' \rangle = \langle e''_0, \sigma'' \rangle$ pertanto deve essere che $e'_0 = e''_0$ e $\sigma' = \sigma''$, da cui discende $\gamma' = \gamma''$

- (b) $e_0 = m \in \mathbb{N}$ ed $\langle e_1, \sigma \rangle \rightarrow \langle e'_1, \sigma' \rangle$ Caso analogo al precedente, dato che ho che $e_1 \notin \mathbb{N}$ e la regola che ho applicato è **Somma 2**. Allora se $\langle e_0 + e_1, \sigma \rangle \rightarrow \gamma''$, è necessario che $\langle e_1, \sigma \rangle \rightarrow \langle e''_1, \sigma'' \rangle$ e $\gamma'' = \langle e_0 + e''_1, \sigma'' \rangle$. Tuttavia per (HP) si ha che $\langle e'_1, \sigma' \rangle = \langle e''_1, \sigma'' \rangle$ pertanto deve essere che $e'_1 = e''_1$ e $\sigma' = \sigma''$, da cui discende $\gamma' = \gamma''$
- (c) $e_0 \in \mathbb{N}$ ed $e_1 \in \mathbb{N}$ In questo caso, solo **Somma 3** è applicabile, ottenendo una sola passibile transizione:

$$\langle e_0 + e_1, \sigma \rangle \rightarrow \langle P, \sigma \rangle \text{ dove } P = e_0 + e_1$$

Quindi la tesi segue:

$$\langle e_0 + e_1, \sigma \rangle \rightarrow \gamma' \wedge \langle e_0 + e_1, \sigma \rangle \rightarrow \gamma'' \implies \gamma' = \gamma''$$

4. $e = e_1 - e_2$: DEL TUTTO ANALOGO AL CASO PRECEDENTE

Q.e.d. ☺

Questo teorema ci porta ad un dio boia di corollario:

Corollary 2.8.1

poiché \rightarrow_e è deterministica, a partire da $\langle e, \sigma \rangle$ arriveremo su una sola configurazione terminale $\langle n, \sigma \rangle$: "n è il valore di e in σ "

È possibile perciò definire una funzione

$$eval : Expr \times Store \rightarrow \mathbb{N}$$

che da semantica alle espressioni

$$eval(e, \sigma) = \begin{cases} m & \text{se } \langle e, \sigma \rangle \rightarrow^* \langle m, \sigma \rangle \\ \text{indefinita} & \text{altrimenti} \end{cases}$$

Esempi:

Example 2.8.4

- $eval((x + 2) - y, \{x/5, y/3\}) = 4$ dato che

$$\langle (x + 2) - y, \{x/5, y/3\} \rangle \rightarrow^* \langle 4, \{x/5, y/3\} \rangle$$

- $eval((x + 2) - y, \{x/2, y/7\}) = \text{indefinito}$ dato che

$$\langle (x + 2) - y, \{x/2, y/7\} \rangle \rightarrow^* \langle 4 - 7, \{x/2, y/7\} \rangle \nrightarrow$$

Inoltre sia introdotta la definizione di equivalenza:

Definition 2.8.2: Equivalenza tra espressioni

Siano e ed e' due espressioni, allora si dicono **equivalenti** sse $\forall \sigma \in Store \quad eval(e, \sigma) = eval(e', \sigma)$
E si denota con $e \equiv e'$

Esempietto:

Example 2.8.5

$$v_1 + (v_2 + v_3) \equiv (v_1 + v_2) + v_3$$

Si osservi come Eval è definita rispetto alla disciplina di valutazione IS (interno destro), pertanto, rigorosamente, Eval è denotato come $Eval_{is}$. Si può, inoltre, dimostrare che anche per ID (interno destro), il risultato della valutazione è lo stesso:

$$Eval_{is} = Eval_{id}$$

Dove $Eval_{id}(e, \sigma) = \begin{cases} m & \text{se } \langle e, \sigma \rangle \rightarrow_{id}^* \langle m, \sigma \rangle \\ \text{indefinita} & \text{altrimenti} \end{cases}$

Come vedremo, è possibile definire anche altre discipline di valutazione come Esterna Sinistra, Esterne Destra, Esterna parallela.

2.8.2 Semantica delle espressioni booleane

Arriviamo alle espressioni booleane con la seguente grammatica:

$$b ::= t | e = e | b \text{ or } b | \neg b$$

(ricordo che t è una metavariable con un valore di verità true o false) E il seguente sistema di transazione:

$$\langle \Gamma_b, T_b, \rightarrow_b \rangle \text{ dove } \Gamma_b = \{ \langle b, \sigma \rangle | b \in Bexp, \sigma \in Store \} \text{ e } T_b = \{ \langle tt, \sigma \rangle, \langle ff, \sigma \rangle | \sigma \in Store \}$$

e \rightarrow_b è la minima relazione generata dai seguenti assiomi e regole di inferenza:

- **Eq1**

$$\frac{\langle e_0 = e_1, \sigma \rangle \rightarrow_b \langle e'_1, \sigma' \rangle}{\langle m = e_1, \sigma \rangle \rightarrow_b \langle m = e'_1, \sigma' \rangle}$$

- **Eq2**

$$\frac{\langle e_1, \sigma \rangle \rightarrow_e \langle e'_1, \sigma' \rangle}{\langle m = e_1, \sigma \rangle \rightarrow_b \langle m = e'_1, \sigma' \rangle}$$

- **Eq3**

$$\frac{}{\langle m = m, \sigma \rangle \rightarrow_b \langle t, \sigma \rangle} \text{ dove } t = \begin{cases} tt & \text{se } m = n \\ ff & \text{se } m \neq n \end{cases}$$

- **Or1**

$$\frac{\langle b_0, \sigma \rangle \rightarrow_b \langle b'_0, \sigma' \rangle}{\langle b_0 \text{ or } b_1, \sigma \rangle \rightarrow_b \langle b'_0 \text{ or } b_1, \sigma' \rangle}$$

- **Or2**

$$\frac{}{\langle tt \text{ or } b_1, \sigma \rangle \rightarrow_b \langle tt, \sigma \rangle}$$

- **Or3**

$$\frac{}{\langle ff \text{ or } b_1, \sigma \rangle \rightarrow_b \langle b_1, \sigma \rangle}$$

- **Neg1**

$$\frac{\langle b, \sigma \rangle \rightarrow_b \langle b', \sigma' \rangle}{\langle \neg b, \sigma \rangle \rightarrow_b \langle \neg b', \sigma' \rangle}$$

- **Neg2**

$$\frac{}{\langle \neg b, \sigma \rangle \rightarrow_b \langle t', \sigma \rangle} \text{ dove } t' = \begin{cases} tt & \text{se } t = ffff \\ ff & \text{se } t = tt \end{cases}$$

Si tenga presente che Eq1, Eq2 e Eq3 sono cosiddette **interne sinistre** perché inizio a valutare la sottoespressione di sinistra per poi restituire un valore di verità t sse ho ottenuto numeri in tutte e due le sottoespressioni mentre Or1, Or2 e Or3 sono **esterne sinistre** perché inizio a valutare la sottoespressione di sinistra per poi restituire un valore di verità t sse ho ottenuto numeri almeno in una sottoespressione. Quindi se nelle interne dovevo avere dei numeri in tutte le sottoespressioni per poi eseguire la valutazione finale nelle esterne per eseguire la valutazione finale mi basta avere una quantità sufficiente

Anche per i booleani si ha questo teorema:

Theorem 2.8.2

\rightarrow_b è deterministica, ovvero

$$(\gamma \rightarrow_b \gamma' \wedge \gamma \rightarrow_b \gamma'') \implies \gamma' = \gamma''$$

Che porta al seguente corollario:

Corollary 2.8.2

si può, quindi, definire:

$$eval_b(b, \sigma) = \begin{cases} t & \text{se } \langle b, \sigma \rangle \rightarrow^* \langle t, \sigma \rangle \\ \text{indefinita} & \text{altrimenti} \end{cases}$$

E si ha anche la seguente definizione:

Definition 2.8.3: Equivalenza booleani

Siano b ed b' due booleani, allora si dicono **equivalenti** sse $\forall \sigma \in Store \quad eval_b(b, \sigma) = eval_b(b', \sigma)$
E si denota con $b \equiv b'$

Example 2.8.6

$$\neg((3 = v) \vee (3 = 4)) = \neg(v = 3)$$

Si possono definire per b_0 or b_1 regole di valutazioni diverse da ES. Ad esempio ED o IS, ma non sono tutte equivalenti, si provi, ad esempio, con ED:

- Or1':

$$\frac{\langle b_1, \sigma \rangle \rightarrow_b \langle b'_1, \sigma' \rangle}{\langle b_0 \text{ or } b_1, \sigma \rangle \rightarrow_b \langle b_0 \text{ or } b'_1, \sigma' \rangle}$$

- Or2':

$$\overline{\langle b_0 \text{ or } tt, \sigma \rangle \rightarrow_b \langle tt, \sigma \rangle}$$

- Or3':

$$\overline{\langle b_0 \text{ or } ff, \sigma \rangle \rightarrow_b \langle b_0, \sigma \rangle}$$

Example 2.8.7

$$\gamma = \langle \rangle$$

Chapter 3

linguaggi liberi deterministici

3.1 PDA e linguaggi deterministici

3.1.1 PDA deterministici

Definition 3.1.1: PDA deterministico

Un PDA $N = (\Sigma, Q, \Gamma, \delta, q_0, \perp, F)$ si dice **deterministico** sse:

1. $\forall q \in Q, \forall z \in \Gamma, (\forall a \in \Sigma, (\delta(q, \epsilon, z) \neq \emptyset \implies \delta(q, a, z) = \emptyset))$
2. $\forall q \in Q, \forall z \in \Gamma, \forall a \in (\Sigma \cup \{\epsilon\}), (|\delta(q, a, z)| \leq 1)$

Ovvero un PDA è libero deterministico sse in ogni configurazione, il PDA ha al massimo una transizione possibile per un dato stato, simbolo di input, e simbolo in cima alla pila e se ha una transizione ϵ disponibile allora non ha altri tipi di transizioni.

Quindi un PDA:

- ha al massimo una transizione
- non ha conflitti tra transizioni ϵ e transizioni che leggono un simbolo

3.1.2 Definizione di linguaggi liberi deterministici

Definition 3.1.2: Linguaggio libero deterministico

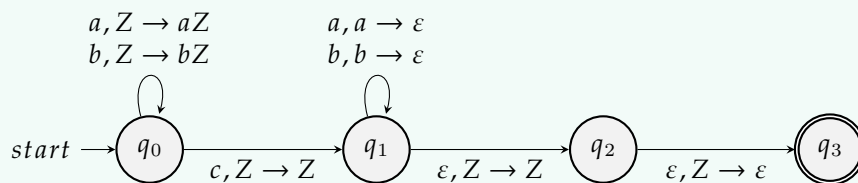
Un linguaggio è **libero deterministico** se è accettato per stato finale da un DPDA

Theorem 3.1.1

la classe dei linguaggi liberi deterministici è inclusa propriamente nella classe dei linguaggi liberi :)

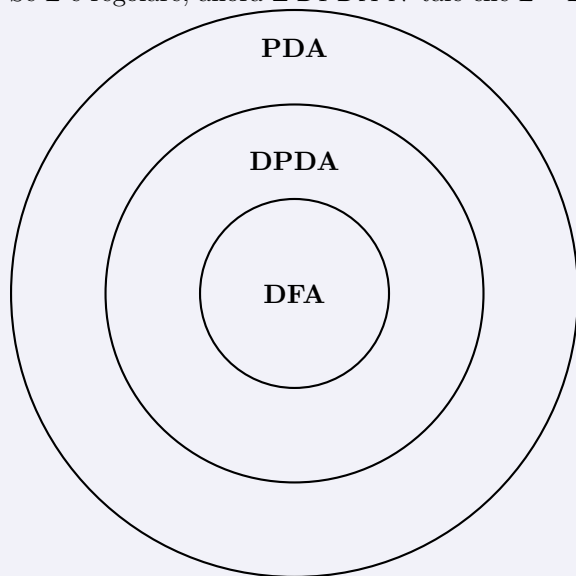
Example 3.1.1

- Sia $L_1 = \{ww^R \mid w \in \{a,b\}^*\}$ è libero, am si può dimostrare che non esiste un DPDA che lo riconosca, infatti con un DPDA non esiste un modo deterministico per riconoscere quando finisce w e inizia w^R
- $L_2 = \{wcw^R \mid w \in \{a,b\}^*\}$ è libero deterministico grazie al segnaposto c è possibile riconoscere quando inizia w^R , si ha infatti:



Theorem 3.1.2

Se L è regolare, allora $\exists DPDA N$ tale che $L = L[N]$ per stato finale



a^*b^*

wcw^R

ww^R

dimostrazione: Se L è regolare, allora $\exists DFA M$ tale che $L = L[M]$. A partire da M , posso costruire un DPDA N si compone come M senza mai manipolare lo stack, allora si che $L = L[N]$ per stato finale ☺

Prefix propriety

Si giunge così al seguente fatto:

Claim 3.1.1

Un linguaggio libero deterministico L è riconosciuto da un DPDA per pila vuota sse L gode della "prefix propriety", ovvero

$$\nexists x, y \in L : x \text{ è prefisso di } y$$

Pertanto si ha che:

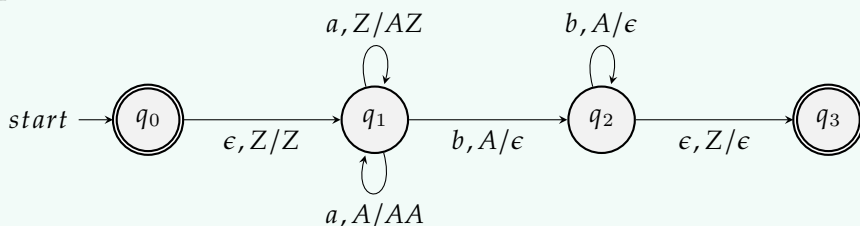
- Se L è non gode della prefix propriety non può essere riconosciuto da un DPDA per pila vuota
- Se L è libero deterministico gode della prefix propriety, allora può essere riconosciuto da un DPDA per pila vuota
- Se L è libero deterministico, allora $L\$ = \{w\$ \mid w \in L\}$ gode della prefix propriety, infatti $L\$$ può essere riconosciuto da un DPDA per pila vuota.

Dove $\$ \notin \Sigma$ ovvero $\$$ non è un simbolo dell'alfabeto di L , ma grazie ad esso alla fine di ogni linguaggio regolare vale la prefix propriety

Example 3.1.2

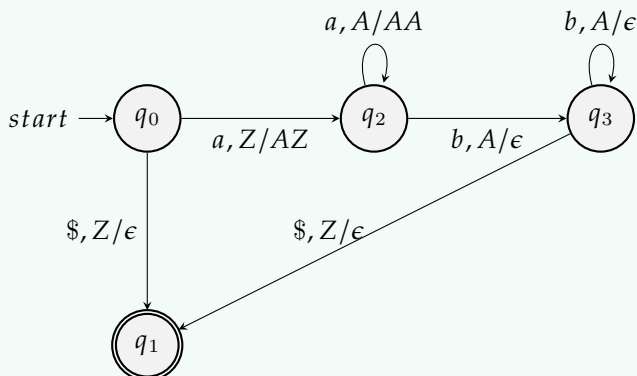
Sia $L_1 = \{a^n b^n \mid n \geq 0\}$ il seguente linguaggio che non gode della prefix propriety, in quanto $\epsilon \in L_1$ e ϵ è

prefisso di ab



Si ha quindi che il seguente linguaggio è riconosciuto da un DPDA *per stato finale* e non *per pila vuota*. Tuttavia il seguente linguaggio con il $\$ \notin \Sigma$ è riconosciuto da un DPDA *per pila vuota*

Sia $L_1\$$:

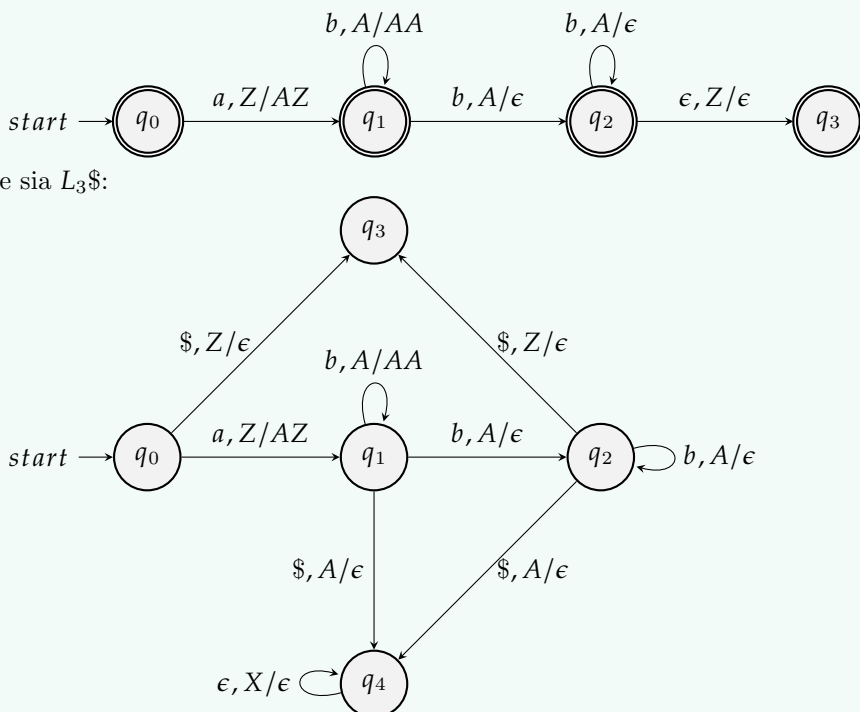


Si verifichi, infatti, come egli venga riconosciuto per pila vuota

Example 3.1.3

Sia $L_3 = \{a^n b^m \mid n \geq m \geq 0\}$ e dato che $\epsilon \in L_3$ tale linguaggio non gode della prefix propriety

e sia $L_3\$$:



non ambiguità dei linguaggi liberi deterministici

Proposition 3.1.1

Se L è libero deterministico, ovvero riconosciuto da un DPDA per *stato finale*, allora L è generabile da una grammatica libera non ambigua.

Si ha quindi che i linguaggi liberi deterministici non sono ambigui

Proprietà dei linguaggi liberi deterministici

I linguaggi liberi deterministici presentano le seguenti proprietà:

Proposition 3.1.2 chiusura solo per complementazione

Sia L un linguaggio libero deterministico, allora questo è chiuso per complementazione, ovvero:

$$(\exists \text{ DPDA } N : L = L(N)) \implies (\exists \text{ DPDA } N' : \bar{L} = L(N')) \quad \text{dove} \quad \bar{L} = \Sigma^* \setminus L$$

Dimostrazione: Bisogna rendere totale la δ di N , eventualmente aggiungendo stati non finali, e poi N' si ottiene da questo N "aumentato", semplicemente scambiando finali e non finali ☺

Proposition 3.1.3 Non chiusura per intersezione

Un linguaggio libero deterministico non è chiuso per intersezione

Example 3.1.4

$L_1 = \{a^n b^n c^m \mid n, m \geq 0\}$ è libero deterministico

$L_2 = \{a^m b^n c^n \mid n, m \geq 0\}$ è libero deterministico

ma

$L_1 \cap L_2 = \{a^n b^n c^n \mid n \geq 0\}$ non è libero!

Proposition 3.1.4 non chiusura per unione

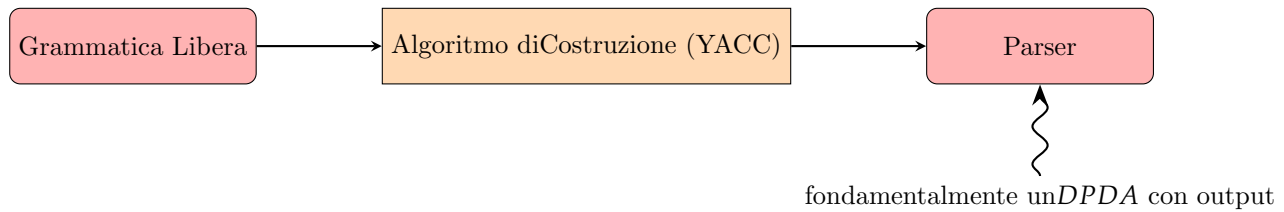
Un linguaggio libero deterministico non è chiuso per unione

Dimostrazione: Assumiamo per assurdo che un linguaggio libero deterministico sia chiuso per unione, allora:

$$L_1 \cap L_2 = \overline{\overline{L_1} \cup \overline{L_2}}$$

Per questo fatto e la proposizione precedente si è verificato un assurdo ☺

3.1.3 Analizzatori sintattici: parser



I parser possono essere:

- **nondeterministici**: se, durante la ricerca di una derivazione, si scopre che una scelta è improduttiva e non porta a riconoscere l'input, **il parser torna indietro** (*backtracking*), disfa parte della derivazione appena costruita e sceglie un'altra produzione, **tornando a leggere parte dell'input**
- **deterministici**: leggono l'input una sola volta ed **ogni loro decisione è definitiva**

entrambi cercano di sfruttare informazioni dall'input per guidare la ricerca della derivazione

introduzione al top-down parsing

Definition 3.1.3

Data $G = (NT, T, S, R)$ libera, costruiamo il $PDA M = (T, \{q\}, T \cup NT, \delta, q, S, \emptyset)$, che riconosce per pila vuota, dove $\delta : Q \times (\Sigma \cup \{\epsilon\}) \times \Gamma \times Q \times \Gamma^*$ è definita:

- **espandi**: $(q, \beta) \in \delta(q, \epsilon, A)$ se $A \rightarrow \beta \in R$
- **consuma**: $\forall a \in T ((q, \epsilon) \in \delta(q, a, a))$

Tale che $L(G) = P[M]$ (riconoscimento per pila vuota)

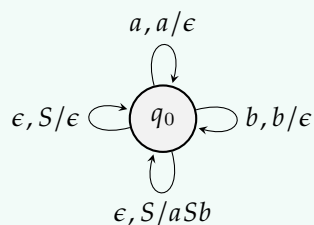
Quindi, grazie alla prima regola se il simbolo in cima alla pila è un non terminale A , l'automa può espanderlo sostituendolo con la produzione β , come descritto nelle regole R della grammatica, mentre secondo la regola di consumazione si ha che se il simbolo sulla cima della pila e il simbolo corrente della stringa in input sono entrambi a , allora l'automa può eliminare quel simbolo dalla pila e procedere nella lettura dell'input.

Example 3.1.5

Sia G la seguente grammatica:

$$S \rightarrow aSb \mid \epsilon$$

Si ha che $S \Rightarrow aSb \Rightarrow aaSbb \Rightarrow aabb$

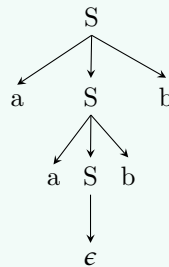


Si osservi che il seguente automa che riconosce il linguaggio della grammatica S non è un PDPA

Nella maggior parte delle configurazioni standard di un Pushdown Automaton (PDA) utilizzato per il parsing top-down, la pila viene inizializzata con il simbolo iniziale della grammatica, nel nostro caso, quindi, la serie di passaggi che porteranno a riconoscere il linguaggio sarà:

$$\begin{aligned}(q, aabb, S) &\vdash (q, aabb, aSb) \\ \vdash (q, abb, Sb) &\vdash (q, abb, aSbb) \\ \vdash (q, bb, Sbb) &\vdash (q, bb, bb) \\ \vdash (q, b, b) &\vdash (q, \epsilon, \epsilon)\end{aligned}$$

che costruisce il seguente albero di derivazione



Si ha, quindi:

- derivazione canonica a sx (leftmost)
- costruzione dell'albero dall'alto in basso

Tuttavia in questo esempio è facile verificare che il parser è nondeterministico, infatti l'automa è un *PDA*, tuttavia questo nondeterminismo può essere risolto scegliendo la produzione in base al simbolo di lettura nell'input (look-ahead), ad esempio:

- se leggo *a*, espando $S \rightarrow aSb$
- se leggo *b*, espando $S \rightarrow \epsilon$

input	stack	azione
<u>a</u> abb\$	S	leggo <i>a</i> e espando in <i>aSb</i>
	<u>a</u> Sb	consumo
<u>a</u> bb\$	Sb	leggo <i>a</i> e espando in <i>aSb</i>
	<u>a</u> Sbb	consumo
<u>b</u> b\$	Sbb	leggo <i>b</i> e espando in ϵ
	<u>b</u> b	consumo
<u>b</u> \$	<u>b</u>	leggo <i>b</i> e espando in ϵ
\$	ϵ	fin

Note:

Non tutte le grammatiche sono adatte per il top down-parser

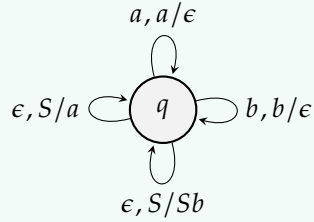
Example 3.1.6

Sia *G* la seguente grammatica:

$$S \rightarrow Sb \mid a$$

e $L(G) = ab^*$ (linguaggio regolare semplice)

L'automa che riconosce il linguaggio è:



Con i seguenti passaggi:

$$(q, ab, S) \vdash (q, ab, Sb) \vdash (q, ab, Sbb) \vdash \dots$$

Qui il determinismo non funziona, infatti se vogliamo espandere quando leggiamo a in input non possiamo anche consumare, pertanto l'automa espanderà all'infinito

Occorre manipolare la grammatica affinché non vi siano ricorsioni sinistre

Introduzione al botto-up parsing

Claim 3.1.2

Data una grammatica libera $G = (NT, T, R, S)$, costruiamo un $PDPA M = (T, \{q\}, T \cup NT \cup \{Z\}, \delta, q, Z, \emptyset)$ che riconosce $L(G).\$$ dove:

- **shift:** $\forall a \in T, \left(\forall X \in T \cup NT \cup Z, ((q, aX) \in \delta(q, a, Z)) \right)$

- **reduce:** $(A \rightarrow \alpha R) \implies (q, A) \in \delta(q, \epsilon, \alpha^R)$

Con α^R una generalizzazione dei PDA in cui si consuma una stringa sulla pila anziché solo il top

- **accept:** $(q, \epsilon) \in \delta(q, \$, SZ)$

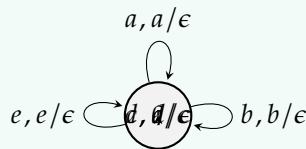
Con S che deve essere alla fine sulla pila e $\$$ simbolo di fine input

Cominciamo subito con un esempio

Example 3.1.7

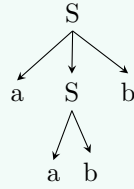
Sia G la seguente grammatica:

$$S \rightarrow aSb \mid ab$$



stack	input	azione
Z	aabb\$	shift
Za	abb\$	shift
Zaa	bb\$	shift
Zaab	b\$	reduce $S \rightarrow ab$
ZaS	b\$	shift
ZaSb	\$	reduce $S \rightarrow aSb$
ZS	\$	accept
ϵ	ϵ	

In passaggi, $S \implies aSb \implies aabb$. Viene prodotto il seguente albero di derivazione:



Si può verificare come la costruzione dell'input sia:

- **left-to-right** (come leggo l'input)
- **right-derivatione**
-

Tuttavia è facile verificare che c'è molto nondeterminismo:

- Conflitti: Shift-Reduce
 1. $zaab \quad b\$$ (può fare shift)
 2. $zaabb \quad \$$ (ma è un percorso infruttuoso)
- Conflitti: Reduce-Reduce
Più riduzioni possibili (non ci sono in quest'esempio, ma con grammatiche più complesse è possibile)

Per ottenere un DPDA, serve introdurre informazioni aggiuntive per risolvere i conflitti:

- **Più stati:** (o strutture particolari di supporto alle decisioni, come un DFA dei prefissi validi).
- **Look-ahead:** (guardare l'input in avanti).

Vediamo un esempio più corposo relativo alle grammatiche delle espressioni aritmetiche:

Example 3.1.8

Sia G tale grammatica:

$$\begin{aligned}
 E &\rightarrow T + E \mid T \\
 T &\rightarrow T \times A \mid A \\
 A &\rightarrow a \mid b \mid (E)
 \end{aligned}$$

Si ha:

G	M
(R1) $E \rightarrow T + E$	$t_0 : \delta(q, a, X) = (q, aX) \forall aX$ SHIFT
(R2) $E \rightarrow T$	$t_1 : \delta(q, \epsilon, E + T) = (q, E)$
(R3) $T \rightarrow T * A$	$t_2 : \delta(q, \epsilon, T) = (q, E)$
(R4) $T \rightarrow A$	$t_3 : \delta(q, \epsilon, A * T) = (q, T)$ REDUCE
(R5) $A \rightarrow (E)$	$t_4 : \delta(q, \epsilon, A) = (q, T)$
(R6) $A \rightarrow a$	$t_5 : \delta(q, \epsilon, (E)) = (q, A)$
(R7) $A \rightarrow b$	$t_6 : \delta(q, \epsilon, a) = (q, A)$
	$t_7 : \delta(q, \epsilon, b) = (q, A)$
	$t_8 : \delta(q, \$, EZ) = (q, \epsilon)$ ACCEPT

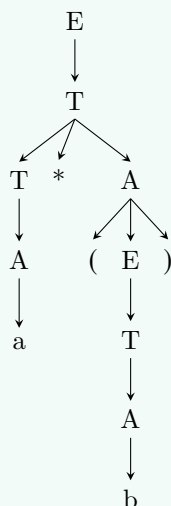
La cui sequenza è:

$$a * (b) \Leftarrow A * (b) \Leftarrow T * (b) \Leftarrow T * (A) \Leftarrow T * (T) \Leftarrow T * (E) \Leftarrow T * A \Leftarrow T \Leftarrow E$$

La cui pila-stack-input è:

Stack	Input	Action
Z	$a * (b) \$$	shift
Za	$*(b) \$$	reduce R_6
ZA	$*(b) \$$	reduce R_4
ZT	$*(b) \$$	shift
ZT*	$(b) \$$	shift
ZT * ($b) \$$	shift
ZT * (b)\$	shift
ZT * (b)	\$	reduce R_7
ZT * (A)\$	reduce R_4
ZT * (T)\$	reduce R_2
ZT * (E)\$	shift
ZT * (E)	\$	reduce R_5
ZT * A	\$	reduce R_3
ZT	\$	reduce R_2
ZE	\$	accept
ϵ	ϵ	

Il cui albero di derivazione:



L'automa è non deterministico perché:

1. Chi ha precedenza tra *shift* e *reduce*?

Ad esempio:

- 1) $Za \quad *(b) \$$ *shift*
- 2) $Za * \quad (b) \$$ non buono! Qui *reduce* è in conflitto con *shift*!

Altro esempio:

- (a) $ZT \quad *(b) \$$ *reduce R_2*
- (b) $ZE \quad *(b) \$$ non buono! Qui *shift* è in conflitto con *reduce*!

2. Chi ha precedenza tra due diverse *reduce*?

Ad esempio:

- (a) $ZT * A \quad \$$ *reduce R_3*
- (b) $ZT \quad \$$

- (c) $ZT * T$ \$ se faccio, *reduce* R_4
 Qui *reduce* R_3 è in conflitto con *reduce* R_4 !

È buona norma scegliere in modo tale che ciò che si trova nella pila sia un prefisso (a rovescio) di una parte destra di una produzione della grammatica.

Ad esempio:

- (a) Za
 (b) ZA con la “norma” produco chi coincide con ciò che è un prefisso di una parte destra che comincia con A .

problema con le pruzioni ϵ

Sia la produzione $S \rightarrow \epsilon$
 TODO! NON HO CAPITO

3.2 Semplificazione delle grammatiche

Per avere **PDA** efficienti e con minor non determinismo è necessario semplificare le grammatiche. Ad esempio:

- Eliminare le produzioni ϵ (del tipo $A \rightarrow \epsilon$) inadatte al bottom up parsing
- Eliminare le **produzione unitarie** (del tipo $A \rightarrow B$ che possono creare dei cicli $A \Rightarrow^+ A$)
- Eliminare **simboli inutili**, cioè quei terminali e non terminali che non sono raggiungibili/generabili a partire dal simbolo iniziale S
 es. gli stati d'errore
- Eliminare la **ricorsione sinistra** (del tipo $A \rightarrow A\alpha$), perché inadatte al top - down parsing
- **fattorizzare** le grammatiche, per ottenere grammatiche con meno non determinismo nel top-down parsing

3.2.1 Eliminare le produzioni ϵ

Per fare ciò si usi un algoritmo che ha:

- in **input**: una G libera con produzione ϵ
- in **output**: una G' libera senza produzione ϵ tale che $L(G') = L(G) \setminus \{\epsilon\}$

Note:

Se $\epsilon \in L(G)$ e si vuole ottenere una G'' t.c. $L(G) = L(G'')$, basta considerare $G' = (NT, T, S, R')$ e definire $G'' = G' \cup \{S' \rightarrow \epsilon | S\}$ t.c.

$$G'' = (NT \cup \{S'\}, T, S', R' \cup \{S' \rightarrow \epsilon | S\})$$

simboli annullabili

Per l'algoritmo occorre innanzi tutto definire i **simboli annullabili**

Definition 3.2.1: simboli annullabili

I **simboli annullabili** sono quei non terminali tale che possono risciversi in uno o più passi in ϵ , ovvero:

$$N(G) = \{A \in NT | A \Rightarrow^+ \epsilon\}$$

Dove $N(G)$ è l'insieme dei simboli annullabili e viene calcolato induttivamente come segue:

- $N_0(G) = \{A \in NT | A \rightarrow \epsilon\}$, questo è il caso in cui un non terminale A viene riscritto direttamente in ϵ tramite una produzione
- $N_{i+1}(G) = N_i(G) \cup \{B \in NT | B \rightarrow c_1, \dots, c_k \in \mathbb{R} \text{ e } c_1, \dots, c_k \in N_i(G)\}$ questo è il caso in cui un non terminale B possa essere ricondotto a ϵ in più passi

Ovviamente $\exists i_c$ tale che $N_{i_c}(G) = N_{i_c+1}(G)$, cioè che ad un certo punto non aggiungo nessun altro B all'insieme (NT è finito)

Theorem 3.2.1

L'insieme $N(G) = N_{i_c}(G)$ è esattamente l'insieme di tutti i simboli annullabili

algoritmo per il calcolo della grammatica

Una volta calcolato $N(G)$ per $G = (N, T, S, R)$, costruiamo la grammatica $G' = (N, T, S, R')$ dove per ogni produzione $A \rightarrow \alpha \in R$ con $\epsilon \notin \alpha$, in cui occorrono simboli annullabili Z_1, \dots, Z_k , mettiamo in R' tutte le produzioni del tipo $A \rightarrow \alpha'$ dove α' si ottiene da α cancellando tutti i possibili sottoinsiemi di Z_1, \dots, Z_k (incluso \emptyset), ad eccezione del caso in cui α' risulta ϵ , in altre parole si creano tutte le possibili combinazioni di α eliminando uno o più di questi simboli Z_1, \dots, Z_n :

- in G' non mettiamo produzioni $A \rightarrow \epsilon \in R$,
- in G' non introduciamo mai produzioni del tipo $A \rightarrow \epsilon$.

Theorem 3.2.2

Data una grammatica libera G , la grammatica G' determinata dall'algoritmo sopra non ha ϵ -produzioni, e $L(G') = L(G) \setminus \{\epsilon\}$

Example 3.2.1

Sia G una grammatica tale che:

$$G = \begin{cases} S \rightarrow AB \\ A \rightarrow aAA \mid \epsilon \\ B \rightarrow bBB \mid \epsilon \end{cases} \quad N_0(G) = \{A, B\} \text{ e } N_1(G) = \{A, B, S\} = N(G)$$

Quindi tutti sono simboli annullabili

Adesso procedo con l'algoritmo, procedo per ogni simbolo annullabili

- $S \rightarrow AB$: secondo l'algoritmo devo cancellare in tutti i modi possibili i simboli non terminali che compaiono nella parte destra della produzione. In questo caso dobbiamo considerare 4 casi:
 - $\emptyset \implies S \rightarrow AB$, rinuncio a cancellare
 - $\{B\} \implies S \rightarrow A$, se cancello B rimane A nella parte destra
 - $\{A\} \implies S \rightarrow B$
 - $\{A, B\} \implies S \rightarrow \epsilon$ dato che non devo mai introdurre produzioni del tipo $A \rightarrow \epsilon$ **non posso cancellare A e B**

Unisco i vari sottoinsiemi e si ha:

$$S \rightarrow AB|B|A$$

- $A \rightarrow aAA \in R$. dobbiamo considerare 4 casi:
 - $\emptyset \implies A \rightarrow aAA$
 - $\{A\} \implies A \rightarrow aA$
 - $\{A, A\} \implies A \rightarrow a$

Quindi si ha:

$$A \rightarrow aAA|aA|a$$

- si ha la stessa cosa con B , quindi cancellarehe verrà trasformato in

$$B \rightarrow bBB|bB|b$$

Così la nuova grammatica sarà:

$$G' = \begin{cases} S \rightarrow AB|A|B \\ A \rightarrow aAA|aA|a \\ B \rightarrow bBB|bB|b \end{cases}$$

3.2.2 Eliminazione delle produzioni unitarie

Definition 3.2.2: Produzione unitaria

Una produzione si dice **unitaria** quando $A \rightarrow B$ si ha che $A, B \in NT$

coppie unitarie

Per eliminare queste produzioni unitarie si deve però calcolare quelle che sono definite le "coppie unitarie"

Definition 3.2.3: Coppia unitaria

Una coppia (A, B) si dice **unitaria** quando $A \Rightarrow^* B$ (quindi quando A può risciversi in 0 o più passi nel non terminale B) usando solo produzioni unitarie.

Vi è qui ripostata la definizione induttiva:

- $U_0(G) = \{(A, A) | A \in NT\}$, quindi ogni non terminale fa coppia con se stesso
- $U_{i+1}(G) = U_i(G) \cup \{(A, C) | (A, B) \in U_i(G) \text{ e } B \rightarrow C \in R, \text{ quindi è l'insieme delle coppie al passo } i \text{ unito alle coppie alle coppie } (A, C) \text{ tali che } (A, B) \text{ sono coppie presenti nell'insieme dell'iterazione precedente e } B \rightarrow C \in R\}$

Anche in questo caso $\exists i_c$ t.c. $U_{i_c}(G) = U_{i_c+1}(G)$ dato che NT è finito. Pertanto per definizione si ha che $U(G) = U_{i_c}(G)$, detto insieme di tutte le coppie unitarie

algoritmo per il calcolo dell'eliminazione delle produzioni unitarie

Data $G = (N, T, R, S)$ libera, si definisce una $G' = (N, T, R', S)$ dove, per ogni $(A, B) \in U(G)$, R' contiene tutte le produzioni $A \rightarrow \alpha$, dove $B \rightarrow \alpha \in R$ e non è unitaria

Note:

Poiché, per ogni $A \in N$, la coppia $(A, A) \in U(G)$, R' contiene tutte le produzioni non unitarie di R e in aggiunta un po' di altre.

Theorem 3.2.3

Sia $G = (NT, T, R, S)$ libera e sia $U(G)$ l'insieme delle sue coppie unitarie. Sia $G' = (NT, T, R', S)$ la grammatica ottenuta dall'algoritmo G' non ha produzioni unitarie e $L(G) = L(G')$

Esempietto:

Example 3.2.2

Prediamo con esempio la grammatica non ambigua E delle espressioni aritmetiche:

$$E = \begin{cases} E \rightarrow E + T | T \\ T \rightarrow T * A | A \\ A \rightarrow a | b | (E) \end{cases}$$

Si noti subito che ha 2 produzioni unitarie: $E \rightarrow T$ e $T \rightarrow A$. Iniziamo a calcolare l'insieme delle coppie unitarie:

- $U_0(G) = \{(E, E), (T, T), (A, A)\}$ e grazie al cuzzo
- $U_1(G) = U_0(G) \cup \{(E, T), (T, A)\}$
- $U_2(G) = U_1(G) \cup \{(E, A)\} = U_3(G) = U(G)$ dato che da E si arriva ad T e si arriva A

Per calcolare la grammatica G' devo prendere tutte le produzioni non unitarie della grammatica originale, ovvero:

$$G' = \begin{cases} E \rightarrow E + T \\ T \rightarrow T \times A \\ A \rightarrow a | b | (E) \end{cases}$$

In aggiunta:

$$\begin{cases} E \rightarrow E + T & \text{perché } (E, T) \in U_1(G) \\ T \rightarrow a | b | (E) & \text{perché } (T, A) \in U_1(G) \\ E \rightarrow a | b | (E) & \text{perché } (E, A) \in U_2(G) \end{cases}$$

Pertanto G' sarà:

$$G' = \begin{cases} E \rightarrow E + T | T \times A | a | b | (E) \\ T \rightarrow T \times A | a | b | (E) \\ A \rightarrow a | b | (E) \end{cases}$$

Sia ha che non contiene produzioni unitarie ed è equivalente a G

3.2.3 Rimuovere i simboli inutili

Definition 3.2.4: Simboli generatori, raggiungibili e utili

Un simbolo $X \in T \cup NT$ è

- Un **generatore** $\iff \exists w \in T^*$ con $x \implies *w$

Quindi un generatore è o un terminale (un simbolo può risciversi in se stesso) oppure un non terminale che in uno o più passi. è definito induttivamente come segue:

- $G_0(G) = T$ se $a \in T, a \implies *a$ (quindi tutti i terminali sono generatori)
- $G_{i+1}(G) = G_i(G) \cup \{B \in NT | B \rightarrow C_1, \dots, C_k \in R \wedge C_1, \dots, C_k \in G_i(G)\}$

- Un **raggiungibile** $\iff (\exists \alpha, \beta \in (T \cup NT)^*. (S \implies *\alpha X \beta))$. Sono definiti induttivamente:

- $R_0(G) = \{S\}$
- $R_{i+1}(G) = R_i(G) \cup \{x_1, \dots, x_k\} \forall B \in R_i(G), B \rightarrow x_1, \dots, x_k \in R$

- **utile** sse è sia un generatore e sia raggiungibile, ovvero se $S \implies *\alpha X \beta \implies *x \in L(G)$ cioè X compare in almeno una derivazione di una stringa $z \in L(G)$

algoritmo per l'eliminazione dei simboli inutili

1. Prima di tutto elimino tutti i non-generatori (e tutte le produzioni che usano almeno uno di questi)

2. Poi dalla nuova grammatica, elimino tutti i non - raggiungibili (E tutte le produzioni che li usano)

Theorem 3.2.4

Sia $G = (NT, T, R, S)$ una grammatica libera t.c. $L(G) \neq \emptyset$

- Sia G_1 la grammatica che si ottiene da G eliminando tutti i simboli che non appartengono a $G(G)$ (insieme dei generatori), e tutte le produzioni che fanno uso di almeno uno di tali simboli
- Sia G_2 la grammatica che si ottiene da G_1 eliminando tutti i simboli che non appartengono a $R(G)$, e tutte le produzioni che fanno uso di almeno uno di tali simboli

Allora G_2 non ha simboli inutili e $L(G_2) = L(G)$

Dimostrazione: La dimostrazione si divide nelle due parti dell'enunciato:

- $L(G_2) \subseteq L(G)$ è ovvio, dato che G_2 contiene meno produzioni di G
- $L(G) \subseteq L(G_2)$: dobbiamo dimostrare che $S \Rightarrow_G^* w$ (ovvero se S deriva w usando le produzioni di G) allora $S \Rightarrow_{G_2}^* w$

Si ha che ogni simbolo usato in $S \Rightarrow_G^* w$ è, ovviamente, sia raggiungibile sia generatore

Quindi quelle derivazioni sono anche una derivazione per G_2

Q.e.d.



Note:

L'ordine dei due generatori è importante!

- prima elimino i non-generatori
- poi i non - raggiungibili

ma se inverte l'ordine, allora può capitare che non elimino tutti i simboli inutili

Esempietto di eliminazione di tutti quei simboli non utili (inutili)

Example 3.2.3

Si parta da questa grammatica:

$$G = \begin{cases} S \rightarrow AB|a \\ B \rightarrow b \end{cases}$$

Poiché $a \Rightarrow^* a, b \Rightarrow^* b, S \Rightarrow^* a, B \Rightarrow^* b$ si ha che i generatori saranno $\{S, B, a, b\}$ (dove manca A). Possiamo così eliminare tutte le produzioni che includono A :

$$G' = \begin{cases} S \rightarrow a \\ B \rightarrow b \end{cases}$$

Adesso posso eliminare tutti i non raggiungibili da S , che in questo caso l'unico è solo B . Si ha che:

$$G'' = S \rightarrow a$$

Si ha che G'' è equivalente a G , ma non contiene simboli inutili

Esempio secondo:

Example 3.2.4

$$G = \begin{cases} S \rightarrow aC \\ A \rightarrow a \\ B \rightarrow bB \\ C \rightarrow b \mid AC \\ D \rightarrow a \mid aS \end{cases}$$

Poiché $G(G) = \{S, a, C, b, A, D\}$ e solo B non è generatore, possiamo eliminare tutte le produzioni che includono B . Si ha quindi:

$$G_1 = \begin{cases} S \rightarrow aC \\ A \rightarrow a \\ C \rightarrow b \mid AC \\ D \rightarrow a \mid aS \end{cases}$$

A questo punto, notiamo che solo D non è raggiungibile, quindi possiamo eliminarlo. Si ottiene:

$$G_2 = \begin{cases} S \rightarrow aC \\ A \rightarrow a \\ C \rightarrow b \mid AC \end{cases}$$

G_2 è la grammatica semplificata, equivalente a G , senza simboli inutili.

In questo esempio si ha che $L(G_2) = \{ab, aab, aaab, \dots\} = a^+b$

Si osservi però che è possibile trovare una grammatica più semplice per il linguaggio a^+b :

$$S \rightarrow aS \mid ab$$

3.2.4 mettere insieme le cose

Se, nel semplificare la grammatica G , **seguiamo questo ordine**:

- Eliminare le ϵ -produzioni
- Eliminare le produzioni unitarie (ovvero i cicli)
- eliminare i simboli inutili

allora la grammatica risultante **è garantita non avere nè ϵ -produzioni, ne produzioni unitarie, ne simboli inutili ed è equivalente a quella di partenza.**

Note:

Si presti attenzione all'ordine poiché alcune delle costruzioni possono interagire tra di loro durante la fase di eliminazione delle ϵ -produzioni, potremmo introdurre produzioni unitarie, pertanto le ϵ -produzioni vanno eliminate prima della fase di eliminazione delle produzioni unitarie

Esempietto:

Example 3.2.5

$$G = \begin{cases} S \rightarrow aAa \mid aa \\ A \rightarrow C \\ C \rightarrow S \mid \epsilon \end{cases}$$

1. Togliere le ϵ -produzioni

$$N(G) = \{S, C, A\} \Rightarrow G' = \begin{cases} S \rightarrow aAa & | \quad aa \\ A \rightarrow C \\ C \rightarrow S \end{cases}$$

2. Togliere le produzioni unitarie

$$U(G') = \{(A, A), (C, C), (S, S), (A, C), (C, S), (A, S)\}$$

$$G'' = \begin{cases} S \rightarrow aAa & | \quad aa & \text{perché } (S, S) \in U(G') \\ C \rightarrow aAa & | \quad aa & \text{perché } (C, S) \in U(G') \\ A \rightarrow aAa & | \quad aa & \text{perché } (A, S) \in U(G') \end{cases}$$

3. Rimuovere i simboli inutili

$$G(G'') = \{S, a, A, C\} \quad \text{tutti i generatori}$$

$$R(G'') = \{S, a, A\} \quad \text{ma non } C$$

$$G''' = \begin{cases} S \rightarrow aAa & | \quad aa \\ A \rightarrow aAa & | \quad aa \end{cases}$$

In questo esempio si ha che $L(G''') = \{aa, aaaa, \dots\} = (aa)^+$

Si osservi però che è possibile trovare una grammatica più semplice per il linguaggio $(aa)^+$:

$$S \rightarrow aSa|aa$$

o anche

$$S \rightarrow aaS|aa$$

3.2.5 forme normali

le **forme normali** sono particolari configurazioni di rappresentazione di un linguaggio formale o di un'espressione logica che rispettano determinate regole e strutture. Ne studieremo di due tipi:

- **Chomsky:** Una grammatica è in forma normale di Chomsky se ogni produzione ha la forma $A \rightarrow BC$ o $A \rightarrow a$, dove $A, B, C \in NT$ e $a \in T$. Ogni produzione deriva quindi o una coppia di variabili o un singolo terminale. Questa forma è utile, per esempio, negli algoritmi di parsing
- **Greibach:** Una grammatica è in forma normale di Greibach se ogni produzione ha la forma $A \rightarrow a\alpha$, dove $A \in NT$, $a \in T$ e α è (eventualmente) una stringa di variabili. La GNF è usata in particolare per costruire parser discendenti

Forma normale di Chomsky

Definition 3.2.5: Forma normale di Chomsky

Una grammatica si dice in **forma normale di Chomsky** se sono nella forma:

$$\begin{aligned} A &\rightarrow BC \\ A &\rightarrow a \end{aligned}$$

Dove ϵ è trattato a parte $S \rightarrow \epsilon|BC$ e S non compare mai a destra in una produzione

Note:

se G è libera in forma normale di Chomsky, allora:

- non ha ϵ -produzioni
- non ha produzioni unitarie

Note:

ogni grammatica libera G può essere trasformata in una equivalente G' in forma normale di Chomsky

Forma normale di Greibach**Definition 3.2.6: forma normale di Greibach**

Una grammatica si dice in **forma normale di Greibach** se sono nella forma:

$$\begin{aligned} A &\rightarrow aBC \\ A &\rightarrow aB \\ A &\rightarrow a \end{aligned}$$

Dove ϵ è trattato a parte $S \rightarrow \epsilon|BC$ e S non compare mai a destra in una produzione

Note:

se G è libera in forma normale di Greibach, allora:

- non ha ϵ -produzioni
- non ha produzioni unitarie
- non è ricorsiva a sinistra
- ogni produzione applicata in una derivazione allunga il prefisso di terminali \implies il parser costruito a partire dalla forma normale di Greibach sono meno non deterministici

Note:

ogni grammatica libera G può essere trasformata in una equivalente G' in forma normale di Greibach

3.2.6 Eliminare la ricorsione a sinistra

L'eliminazione della ricorsione a sinistra è un problema tipico dei parser top-down

Definition 3.2.7: produzione ricorsiva a sinistra

Si definisce **una produzione ricorsiva a sinistra** una produzione del tipo

$$A \rightarrow A\alpha \in R$$

Definition 3.2.8: grammatica ricorsiva a sinistra

Si definisce una **grammatica ricorsiva a sinistra** una grammatica G del tipo:

$$A \implies {}^+A\alpha \text{ per qualche } A \in NT, \alpha \in (T \cup NT)^*$$

Una tipica ricorsione a sinistra è:

$$A \rightarrow A_{\alpha_1} | \dots | A_{\alpha_n} | \beta_1 | \dots | \beta_n$$

Dove le stringhe β_i non cominciano per A . Queste produzioni possono essere rimpiazzate da

$$\begin{aligned} A &\rightarrow \beta_1 A' | \dots | \beta_m A' \\ A' &\rightarrow \alpha_1 A' | \dots | \alpha_n A' | \epsilon \end{aligned}$$

Se nella grammatica originale avviamo la derivazione

$$A \Rightarrow A\alpha_{i_1} \Rightarrow A\alpha_{i_2}\alpha_{i_1} \Rightarrow \dots \Rightarrow A\alpha_{i_k}\dots\alpha_{i_2}\alpha_{i_1} \Rightarrow \beta_i\alpha_{i_k}\dots\alpha_{i_2}\alpha_{i_1}$$

Con la nuova grammatica si ha:

$$A \Rightarrow \beta_i A' \Rightarrow \beta_i \alpha_{i_k} A' \Rightarrow \dots \Rightarrow \beta_i \alpha_{i_k} \dots \alpha_{i_2} A' \Rightarrow \beta_i \alpha_{i_k} \dots \alpha_{i_1} A' \Rightarrow \beta_i \alpha_{i_k} \dots \alpha_{i_1}$$

Esempi concreti:

Example 3.2.6

	$A \rightarrow Aa b$
	\Rightarrow
	$A \rightarrow bA'$
	$A' \rightarrow aA' \epsilon$
Poi	
	$A \rightarrow Ab Ac d$
	\Rightarrow
	$A \rightarrow dA'$
	$A' \rightarrow bA' cA' \epsilon$

Note:

Se $G = [A \rightarrow Aa]$, non si può applicare l'algoritmo perché mancano le produzioni di base da cui partire ($A \rightarrow \beta_1 | \dots | \beta_m$). Infatti, $L(G) = \emptyset$ e la grammatica corrispondente non ha produzioni

3.2.7 Ricorsione sx non-immediata

Consideriamo

$$G = \begin{cases} S \rightarrow Ba|b \\ B \rightarrow Bc|Sc|d \end{cases}$$

In G c'è ricorsione sx immediata ($B \rightarrow Bc$) ma anche non immediata ($S \Rightarrow Ba \Rightarrow Sca$)

algoritmo per il calcolo della ricorsione non immediata

Algorithm 1: ricorsione non immediata

Input: una G libera senza ϵ -prod, senza produzioni unitarie, ma con ricorsione sc non immediata
Output: una G libera senza ϵ -prod, senza produzioni unitarie e senza alcuna ricorsione a sx

- 1 Let $NT = \{A_1, A_2, \dots, A_n\}$ in un ordine fissato;
- 2 **for** $i = 1$ **to** n **do**
- 3 **for** $j = 1$ **to** $i - 1$ **do**
- 4 Sostituisci ogni produzione della forma $A_i \rightarrow A_j \alpha$ con le produzioni $A_i \rightarrow \beta_1 \alpha | \dots | \beta_k \alpha$, dove
 $A_j \rightarrow \beta_1 | \dots | \beta_k$ sono produzioni correnti per A_j ;
- 5 Elimina la ricorsione immediata su A_i ;

- L'obiettivo dell'algoritmo è che, alla fine, ogni produzione del tipo $A_i \rightarrow A_k \alpha$ sia tale che $i < k$, in modo che sia impossibile avere ricorsione sx non immediata.
- Quando $i = 1$, l'unica cosa che viene fatta è l'istruzione 2), che rimuove l'eventuale ricorsione sx immediata. Al termine, $A_1 \rightarrow A_k \alpha$ avremo $i < k$.
- Alla i -esima iterazione del **for** esterno, tutti i non-terminali A_m con $m < i$ hanno produzioni con la proprietà desiderata.

Ora il ciclo **for** interno (istruzione 1) aumenta progressivamente l'indice del non-terminali in prima posizione; finché, al termine del ciclo ($j = i - 1$), avremo che ogni produzione $A_i \rightarrow A_k \alpha$ è tale che $i < k$.

Ora l'istruzione 2) rimuove l'eventuale ricorsione sx immediata da A_i , sicché ogni produzione $A_i \rightarrow A_k \alpha$ è tale che $i < k$.

- Quindi al termine dell'algoritmo, avremo che ogni produzione $A_i \rightarrow A_k \alpha$ è tale che $i < k$, garantendo l'impossibilità di creare ricorsione sx non immediata.

Example 3.2.7

Come esempio si consideri la grammatica di prima, ovvero

$$G = \begin{cases} S \rightarrow Ba|b \\ B \rightarrow Bc|Scd \end{cases}$$

Si segua passo-passo l'algoritmo j3:

- $i=1$ (ovvero S): il ciclo interno non viene eseguito e, siccome non c'è ricorsione immediata per S , non viene fatto nulla
- $i=2$ (cioè $A_i = B$): il ciclo interno (j da 1 a 1) si esegue solo per $A_j = A_1 = S$.

Allora la produzione $B \rightarrow Sc$ viene rimpiazzata con:

$$B \rightarrow Bac|bc$$

Ora le produzioni complessive per B sono :

$$B \rightarrow Bc|Bac|bc|d$$

Dalla quale dobbiamo eliminare la ricorsione immediata, il risultato è:

$$\begin{aligned} B &\rightarrow bcB'|sB' \\ B' &\rightarrow cB'|acB'|\epsilon \end{aligned}$$

Pertanto la gigagrammatica risultante è:

$$\begin{aligned} S &\rightarrow Ba|B \\ B &\rightarrow bcB'|sB' \\ B' &\rightarrow cB'|acB'|\epsilon \end{aligned}$$

3.2.8 Fattorizzazione a sinistra

Si prendi in esempio la seguente grammatica:

$$A \rightarrow aBbC|aBd$$

Se, in un top-down parsing, sulla pila ha A e leggo in input a , non sono in grado di determinare quale produzione scegliere tipico del nondeterminismo, pertanto occorre **raccogliere la parte comune (aB) alle 2 produzioni e introduco un nuovo nonterminale per rappresentare il resto delle produzioni**, quindi:

$$\begin{aligned} A &\rightarrow aBA' \\ A' &\rightarrow bC|d \end{aligned}$$

algoritmo per il calcolo della fattorizzazione

Algorithm 2: Fattorizzazione LU

Input: Grammatica G non fattorizzata
Output: Grammatica G' fattorizzata

- 1 Let N be a new variable;
- 2 $N \leftarrow NT$;
- 3 **while** *è possibile modificare a N o all'insieme delle produzioni* **do**
- 4 **foreach** $A \in N$ **do**
- 5 Sia α il prefisso più lungo comune alle parti destre di alcune produzioni di A ;
- 6 **if** $\alpha \neq \epsilon$ **then**
- 7 Sia A un nuovo non terminale ;
- 8 $N \leftarrow N \cup \{A\}$;
- 9 rimpiazza tutte le produzioni per A del tipo
$$A \rightarrow \alpha\beta_1 | \dots | \alpha\beta_k | \gamma_1 | \dots | \gamma_h$$

con le produzioni:
$$A \rightarrow \alpha A' | \gamma_1 | \dots | \gamma_h$$
$$A' \rightarrow \beta_1 | \dots | \beta_k$$

Example 3.2.8

Riparto una grammatica da fattorizzare:

$$\begin{aligned} E &\rightarrow T | T + E | T - E \\ T &\rightarrow A | A * T \\ A &\rightarrow a | b | (E) \end{aligned}$$

Dove sia E che T si possono fattorizzare, perciò diventa:

$$\begin{aligned} E &\rightarrow TE' \\ E' &\rightarrow \epsilon | + E | - E \\ T &\rightarrow AT' | \\ T' &\rightarrow \epsilon | * T \\ A &\rightarrow a | b | (E) \end{aligned}$$

Chapter 4

Parser Top-Down

Un **parser Top-Down** è un tipo di analizzatore sintattico per analizzare strutture gerarchiche, come le frasi di una lingua o la struttura di un codice. **Funziona esplorando e costruendo l'albero sintattico partendo dalla radice e procedendo verso le foglie**, quindi "dall'alto verso il basso"

Adesso presentiamo un primo esempio di parser Top-Down **nondeterministico** che usa implicitamente una pila per gestire le chiamate ricorsive

4.1 Parser a discesa ricorsiva

Data una grammatica libera $G = (NT, T, S, R), \forall A \in NT$ con produzioni:

$$A \rightarrow X_1^1 \dots X_{n_1}^1 \mid \dots \mid X_k^1 \dots X_{n_k}^k$$

Definisce la funzione

Algorithm 3: $A()$

```
1 scegli non deterministicamente  $h$  tra 1 e  $k$ , ovvero una produzione  $X_k^h \dots X_{n_h}^h$ ;  
2 for  $i = 1$  to  $n_h$  do  
3   if  $X_i^h \in NT$  then  
4      $\lfloor$  Nothing();  
5   else if  $X_i^h = \text{simbolo corrente dell'input}$  then  
6      $\lfloor$  avanza di un simbolo nell'input;  
7   else  
8      $\lfloor$  Fail() ; // backtracking! si torna alla riga 2 e si sceglie un'altra produzione  
9      $\lfloor$  return;
```

Si comincia invocando la funzione per il simbolo iniziale S

Example 4.1.1

Sia G la grammatica:

$$S \rightarrow ac|aSb$$

Col linguaggio:

$$L = \{a^{n+1}cb^n | n \geq 0\}$$

e sia $aacb$ un input. Si ha

4.2.1 First

Definition 4.2.1: First

Data una grammatica libera G e $\alpha \in (T \cup NT)^*$, si definisce **First**(α) come l'insieme dei terminali che possono stare in prima posizione in una stringa che si deriva da α

- per $a \in T, a \in \text{First}(\alpha) \iff \alpha \implies *a\beta$ per $\beta \in (T \cup NT)^*$
- inoltre $(\alpha \implies *\epsilon) \implies \epsilon \in \text{First}(\alpha)$

Note:

Sia la grammatica

$$A \rightarrow \alpha_1 | \alpha_2$$

Se $\text{First}(\alpha_1) \cap \text{First}(\alpha_2) = \emptyset$ la scelta della produzione è deterministica

Qui vi è riportato un esempietto:

Example 4.2.1

$$A \rightarrow aB | bC$$

Si ha che:

$$\begin{aligned}\text{First}(aB) &= \{a\} \\ \text{First}(bC) &= \{b\}\end{aligned}$$

Pertanto abbiamo del determinismo con un solo carattere in lettura

algoritmo per calcolare il first

Algorithm 4: First()

Input: Una grammatica G

Output: $\text{First}(G)$

```
1 for  $x \in T$  do
2    $\text{First}(x) \leftarrow \{x\};$                                 // un terminale è il primo elemento di se stesso
3 for  $X \in NT$  do
4    $\text{First}(X) \leftarrow \emptyset;$                             // per ogni  $x$  non terminale si inizializza il suo first a ""
5 while almeno un  $\text{First}(X)$  può essere modificato in una iterazione do
6   foreach  $x \rightarrow Y_1, \dots, Y_k$  do
7     foreach  $i = 1$  to  $k$  do
8       // se ciascuno di questi simboli  $y_1, \dots, y_{i-1}$  può derivare la stringa vuota  $\epsilon$ 
9       if  $Y_1, \dots, Y_{i-1} \in N(G)$  then
10         $\text{First}(X) \leftarrow \text{First}(X) \cup (\text{First}(Y_i) \setminus \{\epsilon\});$  // allora è possibile aggiungere gli elementi di
11         $\text{FIRST}(Y_i)$  a  $\text{FIRST}(X)$  per la produzione  $y_1, \dots, y_k$ 
12        // Se invece uno dei simboli da  $Y_1$  a  $Y_{i-1}$  non è annullabile, si interrompe la ricerca per quella
13        produzione, perché non possiamo "saltare" i simboli non annullabili per arrivare a  $Y_i$ 
14
15 foreach  $X \in N(G)$  do
16    $\text{First}(X) = \text{First}(X) \cup \{\epsilon\};$ 
```

In generale per una stringa α si ha che:

- Se $\alpha = \epsilon$, allora $\text{FIRST}(\alpha) = \{\epsilon\}$.
- Se $\alpha = X\beta$ e $X \notin N(G)$, allora $\text{FIRST}(X\beta) = \text{FIRST}(X)$.
- Se $\alpha = X\beta$ e $X \in N(G)$, allora $\text{FIRST}(X\beta) = (\text{FIRST}(X) \setminus \{\epsilon\}) \cup \text{FIRST}(\beta)$

In pratica se

$$A \rightarrow \alpha_1 | \dots | \alpha_k$$

si ha che

$$First(A) = First(\alpha_1) \cup \dots \cup First(\alpha_k)$$

Example 4.2.2

Si ossrvi la seguente grammatica:

$$\begin{aligned} S &\rightarrow Ab|c \\ A &\rightarrow aA|\epsilon \end{aligned}$$

$$\begin{aligned} FIRST(S) &= FIRST(Ab) \cup FIRST(c) \\ &= (FIRST(A) \setminus \{\epsilon\}) \cup FIRST(b) \cup \{\epsilon\} \\ &= \{a\} \cup \{b\} \cup \{c\} = \{a, b, c\} \end{aligned}$$

$$\begin{aligned} FIRST(A) &= FIRST(aA) \cup FIRST(\epsilon) \\ &= \{a\} \cup \{\epsilon\} = \{a, \epsilon\} \end{aligned}$$

4.2.2 Follow

Definition 4.2.2: Follow

Data una grammatica libera G e $A \in NT$, definiamo che $Follow(A)$ è l'insieme dei terminali che possono comparire immediatamente a destra di A in una forma sentenziale.

- Per ogni $a \in T$, $a \in Follow(A)$ se $S \Rightarrow^* \alpha A a \beta$ per qualche α e $\beta \in (T \cup NT)^*$.
- $\$ \in Follow(A)$ se $S \Rightarrow^* \alpha A$ (Poiché $S \Rightarrow^* S$, allora $\$ \in Follow(S)$!)

Riporto qui un esempio

Example 4.2.3

$$\begin{aligned} S &\rightarrow Ab \mid c \\ A &\rightarrow aA \mid \epsilon \end{aligned}$$

$$Follow(S) = \{\$ \} \quad Follow(A) = \{b\}$$

- $S \Rightarrow^* S$
- $S \Rightarrow^* Ab$

algoritmo per calcolare il Follow

Algorithm 5: Follow()

Input: Una grammatica credo
Output: bho

```

1 foreach  $X \in NT$  do
2    $First(X) \leftarrow \emptyset;$  // per ogni X non terminale si inizializza il suo first a "0"
3  $Follow(S) \leftarrow \{\$ \};$ 
4 while almeno un  $Follow(X)$  può essere modificato in una iterazione do
5   foreach  $X \rightarrow \alpha Y \beta$  do
6      $Follow(Y) \leftarrow Follow(Y) \cup (First(\beta) \setminus \{\epsilon\});$ 
7   foreach  $X \rightarrow \alpha Y$  do
8     foreach  $X \rightarrow \alpha \beta, \epsilon \in First(\beta)$  do
9        $Follow(Y) \leftarrow Follow(Y) \cup Follow(X);$ 

```

in pratica, occorre cercare tutte le produzioni in cui $Y \in NT$ appare e, per ognuna di esse, applicare la 1 o la 2 sopra

Example 4.2.4

(10)

$$\begin{aligned}
 E &\rightarrow TE' \\
 E' &\rightarrow \epsilon \mid +E \mid -E \\
 T &\rightarrow AT' \\
 T' &\rightarrow \epsilon \mid *T \\
 A &\rightarrow a \mid b \mid (E)
 \end{aligned}$$

	First	Follow
E	a, b, (\$,)
E'	$\epsilon, +, -$	\$,)
T	a, b, (\$,), +, -
T'	$\epsilon, *$	\$,), +, -
A	a, b, (\$,), +, -, *

$Follow(E) = \{ \$,) \}$ perché E è il simbolo iniziale
 perché $A \rightarrow (E)$
 ? $E' \rightarrow +E \mid -E$ richiediamo che $Follow(E') \subseteq Follow(E)$

$Follow(E') = Follow(E)$ poiché $E \rightarrow TE'$, deve essere $Follow(E) \subseteq Follow(E')$
 $\Rightarrow Follow(E) = Follow(E')$

$Follow(T) = \{ \$,), +, - \}$ perché $E \rightarrow TE'$ \Rightarrow include $First(E') \setminus \{\epsilon\}$
 + poiché $\epsilon \in First(E')$ include anche $Follow(E)$!
 $T' \rightarrow *T \Rightarrow Follow(T') \subseteq Follow(T)$

$Follow(T') = Follow(T)$ perché $T \rightarrow AT'$ $\Rightarrow Follow(T) \subseteq Follow(T') \Rightarrow Follow(T) = Follow(T')$

$Follow(A) = \{ \$,), +, -, * \}$ perché $T \rightarrow AT'$ \Rightarrow include $First(T') \setminus \{\epsilon\}$
 + poiché $\epsilon \in First(T')$ \Rightarrow include $Follow(T)$

Adesso che abbiamo introdotto i le procedure First e Follow occorre fare un passo in più per definire i parser

4.2.3 Parser per linguaggi $LL(1)$

tabella di parsing $LL(1)$

La tabella di parsing $LL(1)$ è una struttura di dati usata nei parser sintattici molto utili per risolvere il non determinismo. Questi parser leggono l'input da sinistra a destra (da qui il primo "L" di "LL"), costruendo una derivazione sinistra, o leftmost (da qui il secondo "L") e usano un solo simbolo di lookahead (da cui il "(1)"). Questa tabella è formata da una **matrice bidimensionale** M che è formata da:

- **righe**: non-terminali
- **colonne**: terminali (incluso \$)
- **casella** (A, a) : $M[A, a]$ contiene le produzioni che possono essere scelte dal parser mentre tenta di espandere A e l'input corrente è a .

Se ogni casella contiene **al più una produzione**, allora il parser è **deterministico**!

Per riempire la tabella occorre procedere in questo modo:

Per ogni produzione $A \rightarrow \alpha$:

1. per ogni $a \in T$ e $a \in \text{First}(\alpha)$, inserisci $A \rightarrow \alpha$ nella casella $M[A, a]$
2. se $\epsilon \in \text{First}(\alpha)$, inserisci $A \rightarrow \alpha$ in tutte le caselle $M[A, x]$ per $x \in \text{Follow}(A)$ (x può essere \$)

Ogni casella vuota, dopo aver elaborato tutte le produzioni, è un errore (cioè la funzione ricorsiva chiama 'fail')

grammatica $LL(1)$

Definition 4.2.3: grammatica $LL(1)$

Una grammatica si definisce $LL(1)$ sse ogni casella della tabella di parsing $LL(1)$ contiene al più una produzione, ovvero non presenta conflitti

Si ha che se $G = LL(1)$ allora il parser è **predittivo e deterministico**, questo perché il parser ricostruisce l'albero di derivazione per l'input w , in modo top-down, predicendo quale produzione usare (tra le molte possibili) guardando il prossimo carattere dell'input

Theorem 4.2.1

G è $LL(1)$ sse per ogni coppia di produzioni distinte con la stessa testa

$$A \rightarrow \alpha | \beta$$

si ha che

1. $\text{First}(\alpha) \cap \text{First}(\beta) = \emptyset$
2. (a) $(\epsilon \in \text{First}(\alpha)) \implies (\text{First}(\beta) \cap \text{Follow}(A) = \emptyset)$
(b) $(\epsilon \in \text{First}(\beta)) \implies (\text{First}(\alpha) \cap \text{Follow}(A) = \emptyset)$

Dimostrazione: Se sono soddisfatte le condizione 1 e 2 per ogni coppia di produzioni distinte con medesima testa allora la tabella di parsing $LL(1)$ contiene al più una prodizione in ogni cassella. Ma vale anche viceversa! ☺

Linguaggio $LL(1)$

Definition 4.2.4: Linguaggio $LL(1)$

Un linguaggio si definisce $LL(1)$ $\iff \exists G'$ grammatica $= LL(1)$ che lo genera

Example 4.2.5

Sia G la seguente grammatica:

$$\begin{aligned} S &\rightarrow A|B \\ A &\rightarrow ab|cd \\ B &\rightarrow ad|cb \end{aligned}$$

Si può notare che G non è $LL(1)$ dato che $S \rightarrow A|B$ e

$$\begin{aligned} First(A) &= \{a, c\} \\ First(B) &= \{a, c\} \\ First(A) \cap First(B) &= \{a, c\} \end{aligned}$$

Dal teorema sopra fornito si può dimostrare che non è $LL(1)$

Tuttavia si può manipolarla per farla diventare $LL(1)$, quindi espando S :

$$S \rightarrow ab|cd|ad|cb$$

$$\begin{aligned} S &\rightarrow aT|cT' \\ T &\rightarrow b|d \\ T' &\rightarrow b|d \end{aligned}$$

Poi osservo che T e T' sono identici, sia quindi G' la nuova grammatica:

$$\begin{aligned} S &\rightarrow aT|cT \\ T &\rightarrow b|d \end{aligned}$$

Si può dimostrare che è $LL(1)$, pertanto, per la definizione di linguaggio $LL(1)$ e nonostante G non sia $LL(1)$, si ha che $L(G) = \{ab, cd, ad, cb\}$ è un linguaggio $LL(1)$ perché G' che lo genera è una grammatica $LL(1)$

Theorem 4.2.2

Ogni linguaggio regolare è generabile da una grammatica G di classe $LL(1)$

Dimostrazione: Sia L un linguaggio regolare, allora \exists DFA $M = (Q, \Sigma, \delta, q_0, F) : L = [M]$.

A partire da M si può costruire una grammatica regolare $G = (NT, T, S, R)$ basata sul seguente automa M :

- $NT = \{[q] | q \in Q\}$, cioè un non terminale per ogni stato q
- $T = \Sigma$ cioè un terminale per ogni simbolo dell'alfabeto
- $S = [q_0]$ simbolo iniziale lo stato iniziale
- R (insieme delle produzioni) è definito come:

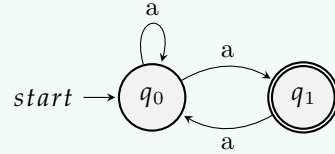
– se $\delta(q, a) = q'$, allora $[q] \rightarrow a[q'] \in R$

Infatti $\delta(q, a) = q'$ vuol dire che l'automata si trova allo stato q e legge in input a allora arriverà allo stato q' , che viene "tradotto" nella grammatica $[q] \rightarrow a[q']$ che corrisponde ad una produzione in cui il non terminale $[q]$ produce il terminale a e il non terminale $[q']$, per passare al non terminale $[q']$ occorre, infatti, fare match con a

– se $q \in F$, allora $[q] \rightarrow \epsilon \in R$

Poi che M è deterministico, $\forall q \in Q \forall a \in \Sigma \quad \exists! q'. q \xrightarrow{a} q'$, cioè $[q]$ avrà una sola produzione $[q] \rightarrow a[q']$ che "inizia" per a e dato che se q è finale, allora $[q] \rightarrow \epsilon$ è applicabile solo per i $\text{Follow}([q]) = \{\$ \} \implies$ nessun conflitto, dato che nessuna produzione genera $\$$, si ha che G è $LL(1)$ ☺

Example 4.2.6



Le produzioni corrispondenti sono:

$$\begin{aligned} [q_0] &\rightarrow a[q_1] \mid a[q_0] \\ [q_1] &\rightarrow a[q_0] \mid \epsilon \end{aligned}$$

La grammatica G è $LL(1)$, perché:

$$\text{First}(a[q_0]) \cap \text{First}(\epsilon) = \emptyset, \quad \text{First}(a[q_1]) \cap \text{First}(\epsilon) = \emptyset$$

Riporto qui un esempio di tabella di parsing $LL(1)$:

Example 4.2.7

Sia G la seguente grammatica:

$$\begin{aligned} E &\rightarrow TE' \\ E' &\rightarrow \epsilon \mid +TE' \mid -TE' \\ T &\rightarrow AT' \\ T' &\rightarrow \epsilon \mid *T \\ A &\rightarrow a \mid b \mid (E) \end{aligned}$$

Si può costruire la seguente tabella di First e Follow:

Produzione	First	Follow
E	$a, b, ($	$\$,)$
E'	$+, -, \epsilon$	$\$,)$
T	$a, b, ($	$+, -, \$,)$
T'	$*, \epsilon$	$+, -, \$,)$
A	$a, b, ($	$*, +, -, \$,)$

Ed ecco a voi la tabella di parsing:

	a	b	$($	$+$	$*$	$\$$
E	$E \rightarrow TE'$	$E \rightarrow TE'$	$E \rightarrow TE'$			
E'				$E' \rightarrow +TE'$	$E' \rightarrow -TE'$	$E' \rightarrow \epsilon$
T	$T \rightarrow AT'$	$T \rightarrow AT'$	$T \rightarrow AT'$			
T'				$T' \rightarrow \epsilon$	$T' \rightarrow *T$	$T' \rightarrow \epsilon$
A	$A \rightarrow a$	$A \rightarrow b$	$A \rightarrow (E)$			

<p>input: $a * (b + a) \\$</p> <p>Stack: E</p> <p>TE'</p> <p>$AT'E'$</p> <p>$aT'E'$</p> <p>$T'E'$</p> <p>$*T'E'$</p> <p>$(b+a) \\$</p> <p>TE'</p> <p>$AT'E'$</p> <p>$(E)T'E'$</p> <p>$b+a) \\$</p> <p>$E)T'E'$</p> <p>$TE')T'E'$</p> <p>$AT'E')T'E'$</p> <p>$bT'E')T'E'$</p>	<p>input: $+ a) \\$</p> <p>Stack: $T'E'$</p> <p>$T'E')T'E'$</p> <p>$E')T'E'$</p> <p>$+E)T'E'$</p> <p>$E)T'E'$</p> <p>$TE')T'E'$</p> <p>$AT'E')T'E'$</p> <p>$aT'E')T'E'$</p> <p>$T'E')T'E'$</p> <p>$E')T'E'$</p> <p>$)T'E'$</p> <p>$T'E'$</p> <p>E'</p> <p>ϵ</p> <p>Esistono: costruire l'albero di derivazione!</p>	<p>input: $a) \\$</p> <p>Stack: $T'E'$</p> <p>$T'E')T'E'$</p> <p>$E')T'E'$</p> <p>$+E)T'E'$</p> <p>$E)T'E'$</p> <p>$TE')T'E'$</p> <p>$AT'E')T'E'$</p> <p>$aT'E')T'E'$</p> <p>$T'E')T'E'$</p> <p>$E')T'E'$</p> <p>$)T'E'$</p> <p>$T'E'$</p> <p>E'</p> <p>ϵ</p> <p>Avuto la pile è a posto!</p>
---	--	--

Algoritmo Per il calcolo di un parser $LL(1)$

Algorithm 6: Parser $LL(1)$

```

Input: Stringa  $w$ 
Output: Niente
1  $Pila \leftarrow S\$;$                                      // cima della pila a sinistra
2  $X \leftarrow S;$                                        // top della pila
   // lettura in input
3  $input \leftarrow w\$;$ 
4  $i_c =$  primo carattere dell'input;
   // viene eseguito il ciclo While finché la pila non è vuota o l'input non è stato consumato
5 while  $X \neq \$$  do
6   if  $X$  è un terminale then
7     // si controlla se  $X$  fa "match" con  $i_c$ 
8     if  $X = i_c$  then
9       Pop  $X$  dalla pila;                                // rimuovo  $X$  dalla pila
10      avanza  $i_c$  sull'input;
11    else
12      Errore();                                         // no match
13  else
14    // se nella tabella di Parsing  $M$  esiste una regola  $X \rightarrow Y_1, \dots, Y_n$ 
15    if  $M[X, i_c] = X \rightarrow Y_1, \dots, Y_n$  then
16      Pop  $X$  dalla pila;
17      Push  $Y_1, \dots, Y_n$  sulla pila;                // mette sulla pila i simboli  $Y_1, \dots, Y_n$  con  $Y_1$  in cima
18      In output la produzione  $X \rightarrow Y_1, \dots, Y_n$ ;
19    else
20      Errore ();                                       // se non esiste una regola in  $M[X, i_c]$  si genera un errore, viene definito caso
21      "bianco"
22     $X \leftarrow$  top della pila;                        // Aggiorna  $X$  con il nuovo simbolo in cima alla pila
23 if  $i_c \neq \$$  then
24   Errore();                                           // pila svuotata ma vi è ancora dell'input da leggere

```

Un altro esempio:

Example 4.2.8

Sia G la seguente grammatica:

$$S \rightarrow aAB \mid bS$$

$$A \rightarrow a$$

$$B \rightarrow b$$

Che genera il seguente linguaggio

$$L(G) = L(b^*aab)$$

Si ha che questa grammatica G è $LL(1)$, perché:

$$First(aAB) \cap First(bS) = \emptyset, \quad \{a\} \cap \{b\} = \emptyset$$

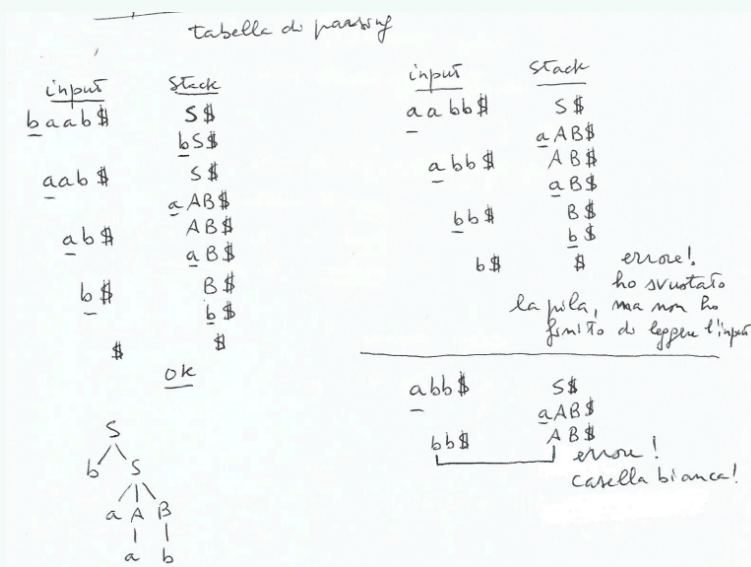
Si può proseguire con la tabella First e follow:

Simbolo	First	Follow
S	a, b	$\$$
A	a	b
B	b	$\$$

Da cui si può costruire la seguente tabella di parsing:

	a	b	$\$$
S	$S \rightarrow aAB$	$S \rightarrow bS$	
A	$A \rightarrow a$		
B		$B \rightarrow b$	

Viene qui descritto il funzionamento del parser con pila con diversi input:



4.2.4 Parser per linguaggi $LL(K)$

grammatiche $LL(K)$

Le grammatiche $LL(k)$ sono "un'estensione" del concetto di grammatiche $LL(1)$, dove il parser ha la capacità di guardare in avanti fino a k simboli per determinare le scelte di parsing

Per questi tipi di grammatica gli insiemi *First* e *Follow* assumo significati diversi rispetto a alle grammatiche $LL(K)$

Definition 4.2.5: First $LL(K)$

L'insieme $First_k(\alpha)$ contiene tutte le stringhe di lunghezza k o minore derivabili dall'inizio di una produzione con α , in particolare $w \in First_k(\alpha) \iff \alpha \Rightarrow^* w\beta$ con $|w| = k, w \in T^*, \beta \in (T \cup NT)^*$ oppure $\alpha \Rightarrow^* w$ con $|w| \leq k$ e $w \in T^*$

Definition 4.2.6: Follow $LL(K)$

L'insieme $Follow_k(A)$ definisce quali stringhe possono apparire immediatamente dopo un simbolo non terminale A in una derivazione a partire dal simbolo iniziale S . In particolare: $w \in Follow_k(A)$ se $S \Rightarrow^* \alpha Aw\beta$ con $|w| = k, w \in T^*$, e $\alpha, \beta \in (T \cup NT)^*$, oppure, $S \Rightarrow^* \alpha Aw$ con $|w| \leq k$ e $w \in T^*$.

Tabella di Parsing $LL(K)$

- **Righe:** non terminali.
- **Colonne:** $\{w \in T^* \mid |w| \leq k\}$ (solo quelle necessarie).

Per ogni produzione $A \rightarrow \alpha$, la tabella $M[A, w]$ contiene:

- $A \rightarrow \alpha$, per ogni $w \in \text{First}_k(\alpha)$ ($w \neq \varepsilon$);
- $w \in \text{Follow}_k(A)$ se $\varepsilon \in \text{First}_k(\alpha)$.

Ogni entrata/casella contiene al più una produzione. Se non esistono w_1 e w_2 tali che w_1 prefisso di w_2 con le due entrate corrispondenti su una riga entrambe riempite, allora G è una grammatica $LL(k)$.

Nota Bene: Le colonne sono tante quante sono le stringhe w che appartengono a $\text{First}_k(\alpha)$ per $A \rightarrow \alpha$ o a $\text{Follow}_k(A)$ per $A \rightarrow \alpha$. Questo va verificato per tutte le produzioni:

$$w \in \text{First}_k(\alpha).$$

Example 4.2.9

Sia G la seguente grammatica

$$S \rightarrow aSb \mid ab \mid c$$

Si ha che $L(G) = \{a^n b^n \mid n \geq 1\} \cup \{a^n c b^n \mid n \geq 0\}$, inoltre:

- $\text{First}_2(aSb) = \{aa, ac\}$
- $\text{First}_2(ab) = \{ab\}$
- $\text{First}_2(c) = \{c\}$

Si può dimostrare che G è $LL(2)$ dato che:

- $\text{First}_2(aSb) \cap \text{First}_2(ab) = \emptyset$
- $\text{First}_2(aSb) \cap \text{First}_2(c) = \emptyset$
- $\text{First}_2(ab) \cap \text{First}_2(c) = \emptyset$

Da cui si può ricavare la tabella di parsing:

Simbolo	aa	ab	ac	c
S	$S \rightarrow aSb$	$S \rightarrow ab$	$S \rightarrow aSb$	$S \rightarrow c$

Theorem 4.2.3

- Una grammatica ricorsiva sinistra non è $LL(K)$ per nessun K
- Una grammatica ambigua non è $LL(K)$
- Se G è $LL(K)$ per qualche k , allora G non è ambigua
- Se G è $LL(K)$, allora $L(G)$ è libero deterministico
- esiste L libero deterministico tale che non esiste G di classe $LL(K)$ per nessun K , tale che $L = L(G)$

Viene qui riportato il funzionamento del parser con pila:

Linguaggio $LL(K)$

Definition 4.2.7: Linguaggio $LL(K)$

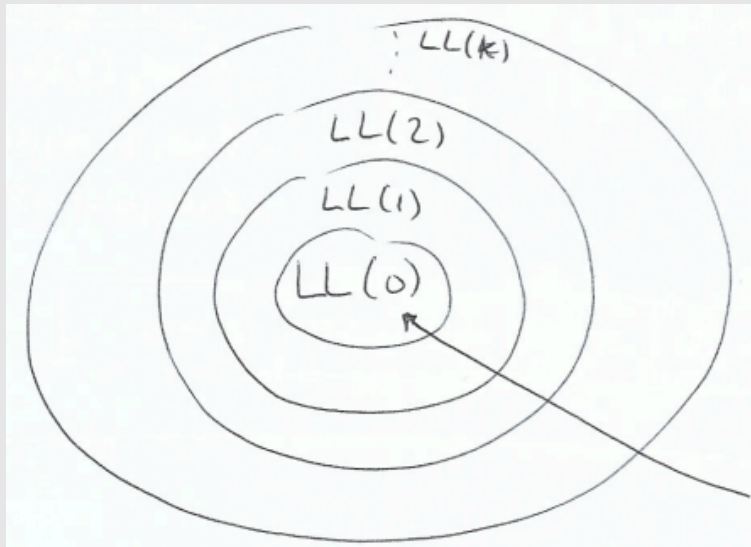
un linguaggio L è di classe $LL(K)$ se G di classe $LL(K)$ tale che $L = L(G)$

Theorem 4.2.4

$\forall k \geq 0$, la classe dei linguaggi $LL(k)$ contiene strettamente la classe dei linguaggi $LL(k)$

Note:

$\forall k \geq 0$, la classe dei linguaggi $LL(k+1)$ contiene strettamente la classe dei linguaggi $LL(k)$



Si ha che: $(\forall A \in NT, \exists! \alpha \in (T \cup NT)^*, A \rightarrow \alpha \implies G \text{ è } LL(0)) \implies L(G) = \{w\}$ ovvero una sola parola al massimo

Nella pratica tuttavia si usano solo $LL(1)$, spesso la si può manipolare trasformandola in $LL(1)$

Example 4.2.10

$$S \rightarrow Asb \mid ab \mid c$$

G è un $LL(2)$ ma non $LL(1)$ Si fattorizza

$$\begin{aligned} S &\rightarrow aT \mid c \\ T &\rightarrow Sb \mid b \end{aligned}$$

Ottenuta fattorizzando G è $LL(1)$ infatti:

- $First(aT) \cap First(c) = \emptyset$
- $First(Sb) \cap First(b) = \emptyset$

Si ha che $L = \{a^n b^n \mid n \geq 1\} \cup \{a^n c b^n \mid n \geq 0\}$ è un linguaggio di classe $LL(1)$ perché G' è $LL(1)$

Casi speciali

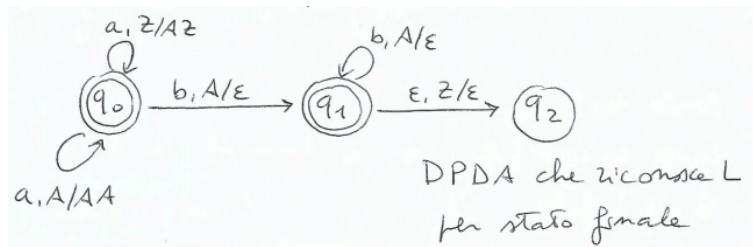
Sia

$$L = \{a^i b^j \mid i \geq j\} \text{ è libero deterministico}$$

Ma non è $LL(K)$ per nessun K !

Adesso, mostriamo una grammatica per L e dimostriamo che non è $LL(K)$ per nessun K . Sia G la seguente grammatica:

$$\begin{aligned} S &\rightarrow aS \mid B \\ B &\rightarrow aBb \mid \epsilon \end{aligned}$$



Sia G una grammatica libero per L

$$S \rightarrow aS \mid BB \rightarrow aBb \mid \epsilon$$

e poniamo $L = L(G)$

Per scegliere tra $S \rightarrow aS$ e $S \rightarrow B$ dovrei leggere fino in fondo l'input per sapere quante b in meno di a ci sono nella stringa! Allora G non può essere $LL(k)$ per nessun k . Infatti quanto possa essere grande il K posso trovare una stringa più lunga che richiede di leggere più di k simboli di lookahead

Non è tuttavia possibile alcuna G' e k tali che

$$L(G') = L \text{ e } G' \in L(K)$$

La dimostrazione non verrà illustrata

Chapter 5

Bottom up parser

Il **parser bottom up** è un tipo di analizzatore sintattico che costruisce l'albero di derivazione partendo dalle foglie, viene anche detto **Shift-Reduce** in quanto:

- un simbolo terminale viene spostato dall'input alla pila
- una serie di simboli (terminali e non terminali) sulla cima della pila corrisponde al "reverse" di una parte destra di una produzione, ovvero:

$$A \rightarrow \alpha \in R - \alpha^R \text{ sulla pila}$$

La stringa α^T viene rimossa dalla pila e sostituita con A , quindi α viene ridotta ad A

Possono essere di diversi tipi

5.1 Parser $L(R)$

Il parser LR è un parser bottom up in cui

- L sta per "leggo da sx a dx"
- R sta per "derivazione rightmost"

Theorem 5.1.1 GioLaPalma-Bonzo

Ciao darione nazionale

Dimostrazione: Per LEMMA dimostrato



zio pera