

In questa lezione cominciamo ad analizzare il livello SINTATTICO nell'analisi del linguaggio.

GERARCHIA DI CHOMSKY ~ (10:00 min)

Data delle frasi P_1 e P_2

P_1 : Mario mangia una mela con il coltello.

P_2 : Mario mangia una mela con il verme.

Vogliamo trovare una rappresentazione delle STRUTTURE SINTATTICHE delle frasi P_1 e P_2 . Vogliamo quindi ESPlicitARE le relazioni che esistono tra le varie parole.

Una volta che abbiamo effettuato l'ANALISI GRAMMATICALE e assegnato ad ogni PAROLA la sua rispettiva CLASSE, non possiamo analizzare come le varie parole delle frasi sono collegate tra loro.

Per rappresentare queste relazioni possiamo utilizzare l'approccio di Chomsky, che consiste nel definire il concetto di GRAMMATICA FORMALE che ci permette di definire in LINGUAGGIO FORMALE.

Questi linguaggi nei sono stati classificati nelle
famora GERARCHIA DI CHOMSKY in base alle loro
COMPLESSITA', dove la completezza di un
linguaggio può essere vista nei seguenti due
modi :

- i) Come la completezza del macchinario (AUTOMA)
memoria e sufficiente per ACCETTARE il linguaggio.
- ii) Come la completezza delle REGOLE della
grammatica che GENERA il linguaggio.

La gerarchie di Chomsky può essere rappresentata
come segue :

LIVELLO	GRAMMATICA	AUTOMA
0	?	TURING MACHINE (INDECIDIBILE)
1	CONTEXT-SENSITIVE	TURING MACHINE (DECIDIBILE)
2	CONTEXT-FREE	PUSH DOWN AUTOMATON
3	REGOLARE	FINITE STATE AUTOMATON

Chomsky era interessato a capire in quale livello
potrebbe essere inserito il LINGUAGGIO NATURALE.

GRAMMATICHE CONTEXT-FREE

~ (28:00 min)

Lavorando con una grammatica CF è possibile rappresentare le informazioni sintattiche tramite delle STRUTTURE ALBOREE in quanto le regole non delle forme $\alpha \rightarrow \beta$, con $|\alpha| = 1$.

ESEMPIO:

$S \rightarrow AB \mid BC$

$A \rightarrow c \mid a$

$B \rightarrow BA \mid B$

$C \rightarrow c$

La frase cBa è ottenuta dalla seguente DERIVAZIONE

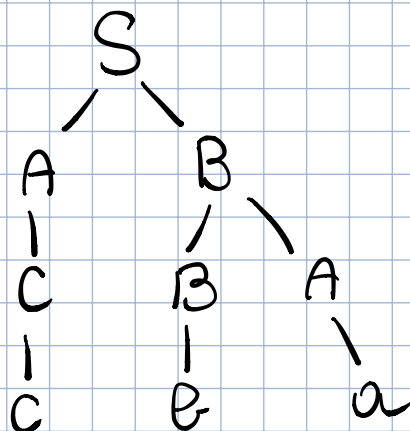
$S \rightarrow AB$

$\rightarrow CB$

$\rightarrow CBA$

$\rightarrow cBa$

Possiamo rappresentare la derivazione di cBa nel seguente modo



Andiamo adesso a rispondere alle seguenti domande:

- i) Perché ci servono le grammatiche CF e non ci bastano le grammatiche REGOLARI.
- ii) Perché ci bastano le CF e non ci servono le grammatiche CONTEXT-SENSITIVE.

Per rispondere a queste domande possiamo utilizzare il PUMPING LEMMA, che ci permette di capire a quale livello appartiene un linguaggio. Abbiamo due versioni del pumping lemma:

1) REGULAR \rightarrow CF

La prima versione è utilizzata per dimostrare che il linguaggio sotto analisi non è regolare.

Un particolare ci dice che se nel linguaggio non presenti stringhe della forma $a^n b^n$, allora il linguaggio non è regolare.

2) CF \longrightarrow CS

Utilizzato per dim. che un linguaggio non è CF.

Ci dice che stringhe delle forme $a^m b^n c^m$ non fanno parte del livello context-free.

Per capire dove si colloca il linguaggio naturale devo quindi analizzare le strutture prodotte e capire se ci sono strutture delle forme $a^m b^n$ o $a^m b^n c^m$.

Q: Le lingue naturali come l'italiano o l'inglese sono REGOLARI o no?

R: Non sono regolari in quanto l'italiano, ad esempio, quando trovo un "CHE" devo anche contare un VERBO. Tra il "CHE" e il VERBO non c'è una relazione fissa tra loro. Dovendo contare qualcosa, non concludere che l'italiano non è un linguaggio REGOLARE.

!

Potrei però finire un NUMERO MASSIMO DI INNESTI DA CONTARE e utilizzare in autonomia stati finiti per analizzare il CORPUS.

OSS: In inglese WITH viene utilizzato per indicare che una frase è ACCESSORIA, mentre il THAT indica che la frase non è ACCESSORIA.

OSS: Gli LSTM non sono in grado di riconoscere i linguaggi naturali anche se non sono context-free.

Consideriamo ora la frase

Pino, Pino e Gino sono fratelli, figlio e madre di GAETANO

Notiamo che lavorando con una grammatica CF non siamo in grado di analizzare strutture del genere delle forme $(C_1 C_2)_1)_2$.

Dato però che questi casi sono pochi, l'idea è quella di rimanere nel livello CF e qualora dovessimo trovare frasi del genere, non insistere rimanendo la loro analisi a livelli successivi.

PARSING ~ (1:05:00 min)

Dato una grammatica $G = \langle R, NT, T, S \rangle$ CF e data una frase $w_1 \dots w_m$, $w_i \in T^*$, la procedura di costruire l'ALBERO SINTATTICO di $w_1 \dots w_m$ secondo G è detta PARSING.

Il PARSING ci permette di capire se la frase $w_1 \dots w_m$ è generata da G , e in caso positivo costruire l'albero sintattico.

Consideriamo le grammatiche

$S \rightarrow AB$

$S \rightarrow ASB$

$A \rightarrow a$

$B \rightarrow b$

e sia 'aab' la frase da analizzare.

Abbiamo varie tecniche per effettuare il parsing, tra cui:

- TOP-DOWN PARSING

L'idea è quella di partire tutte le regole della grammatica ogni singola volta.

$S \rightarrow AB \rightarrow aB$

$S \rightarrow ASB \rightarrow AAB \rightarrow aAB$

A questo punto lo superato la FRONTIERA e rigetto.

È possibile migliorare questa procedura utilizzando dei LOOK-AHEAD a partire dalla frontiera.

- BOTTOM-UP PARSING

Si parte dalle parole della frase per vedere se possiamo applicare delle regole della grammatica per ottenere le parole in analisi.

CYK - PARSING ~ (1:23:00 min)

SE CYK è un algoritmo di parsing di tipo BOTTOM-UP utilizzando la PROGRAMMAZIONE DINAMICA.

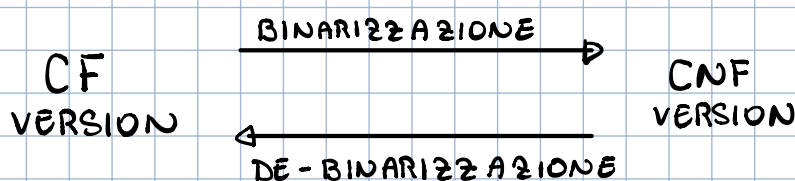
L'algoritmo lavora con una forma particolare della grammatica CF dette CNF (CHOMSKY NORMAL FORM), che è la seguente

$$\text{CNF} \begin{cases} NT \rightarrow NT_1 NT_2 \\ NT \rightarrow T \end{cases}$$

Mostriamo che gli alberi ottenuti da grammatiche di tipo CNF non sono ALBERI BINARI.

Mostriamo poi che le grammatiche CNF non (DEBOLMENTE) equivalenti a quelle CF, nel senso che generano le stesse frasi ma non gli stessi alberi.

Tramite dei processi di conversione possiamo passare dalle versioni CF generale a quelle in CNF date da grammatiche CF.



Consideriamo la seguente grammatica

G:	$S \rightarrow AB \mid BC$
	$A \rightarrow BC$
	$B \rightarrow BB \mid CC$
	$C \rightarrow a$
	$B \rightarrow b$
	$A \rightarrow a$

Andiamo a vedere come saranno le stringhe
 baa e ba

Per notare, il CYK costruisce una matrice che contiene le informazioni sul modo di PARSING. La matrice utilizzata è la seguente

ϵ	a	a	b	a
B	S, A	B	B	S, A
██████████	A, C	B	B	S, A
██████████	██████████	A, C	S	\emptyset
██████████	██████████	██████████	B	S, A
██████████	██████████	██████████	██████████	A, C

Ogni cella contiene i non terminali che permettono di ottenere una sottstringa della stringa in input.

Una possibile esecuzione riempie le diagonali della matrice a partire dalla diagonale principale. Ad ogni passo andiamo a confrontare delle coppie di celle riempite nel passo precedente.

Alla fine del processo se riusciamo ad inserire lo START SYMBOL S nelle celle in alto a destra, quella contiene l'intera stringa, allora la stringa è generata dalla grammatica.

Se la stringa è generata, per poterla costruire l'ALBERO SINTATTICO dobbiamo ogni volta mantenerci dei RIFERIMENTI alle celle utilizzate per produrre i non terminali.

In realtà molto spesso non abbiamo un solo albero, in quanto la grammatica utilizzata potrebbe essere ambigua AMBIGUA.

In caso di grammatiche esclusivamente ambigue non è possibile avere TANTI alberi sintattici. In questi casi capire quale albero scegliere può essere MOLTO DIFFICILE.