

Analisi sintattica efficiente

Giacomo Baldi
Elaborazione del Linguaggio
Naturale

Il problema

- Data una frase già segmentata in tokens si vuole costruire in modo efficiente la struttura della frase secondo una grammatica che descrive il linguaggio.
- Occorre quindi definire con accortezza la grammatica.

La grammatica

- La grammatica viene inizialmente data in BNF

S ::= NP VP | S Congiunzione S

NP ::= Pronome | Sostantivo | Articolo Sostantivo | NP PP | NP Relativa

VP ::= Verbo | VP NP | VP Aggettivo | VP PP | VP Avverbio

PP ::= Preposizione NP

Relativa ::= CongRel VP

Sostantivo ::= ...

Verbo ::= ...

→ Presenti nel dizionario

Aggettivo ::= ...

La grammatica aumentata

- Sovragerazione: "Me smells a flower"
- Una grammatica in BNF avrebbe un numero enorme di regole se si dovessero esplicitare tutte per garantire le concordanze.
- Si ricorre all'aumento delle categorie sintattiche

S ::= NP(**nominativo**) VP | ...

NP(**caso**) ::= Pronome(**caso**) | Nominativo | Articolo Nominativo ...

VP ::= VP NP(**accusativo**) | ...

Sottocategorizzazione dei verbi

- Ad ogni verbo viene associata la lista dei complementi che possono accompagnarlo.

	GIVE ME SOMETHING	[NP, NP]
	GIVE IT TO ME	[NP, PP]
	GIVE IT TO ME	
VP (cat) ::=		Verbo([NP,PP])
VP([NP, cat]) NP (accusativo)		VP([NP,PP])
VP([Aggettivo, cat]) Aggettivo		VP([NP,PP]) NP
VP([PP, cat]) PP		VP([PP])
Verbo (cat)		VP([PP]) PP
		VP([])

Sottocategorizzazione dei verbi (2)

Affinchè una frase sia riconosciuta almeno una lista di complementi del verbo deve essere soddisfatta

$S ::= NP(\text{Nominativo}) VP([])$

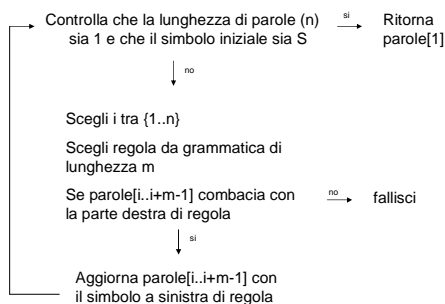
Per rendere espressiva la grammatica si aggiungono i **complementi esterni**

$VP(\text{cat}) ::= VP(\text{cat}) PP \mid VP(\text{cat}) \text{Avverbio}$



Parser nondeterministico bottom-up

Input: parole (vettore di simboli), grammatica (insieme di regole)



Parsing Top-Down

Un parser top-down (o predittivo) cerca di interpretare una frase con una derivazione leftmost della grammatica associata al linguaggio.

I DRINK TEA

$S \rightarrow NP VP \rightarrow \text{Pronome} VP \rightarrow \text{"I"} VP \rightarrow \text{"I"} \text{Verbo} NP \rightarrow \text{"I DRINK"} NP \rightarrow \text{"I DRINK TEA"}$

Parsing Bottom-Up

Un parser bottom-up cerca di interpretare una frase con una derivazione rightmost della grammatica associata al linguaggio.

I DRINK TEA

"I DRINK TEA" → Pronome "DRINK TEA" → NP "DRINK TEA"
→ NP Verbo "TEA" → NP Verbo NP → NP VP → S

Verso un'analisi efficiente

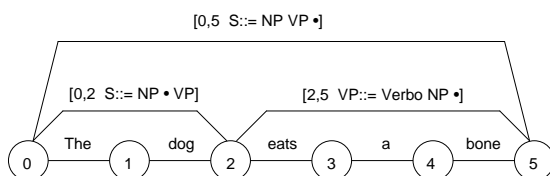
- Non fare due volte ciò che può essere fatto una sola volta
- Non fare ciò che si può evitare
- Non rappresentare ciò che non è necessario

Chart parsing

I risultati parziali vengono memorizzati in una struttura detta chart.

Un chart contiene nodi ed archi etichettati (completi e incompleti).

Un arco consiste dell'indicazione dei nodi iniziali e terminali e della regola che deve venir completata.



Metodi per il chart parsing

Fare chart parsing significa aggiungere nodi ed archi al chart fino a terminare la frase e ottenere come soluzioni tutti gli archi completi che terminano sull'ultimo nodo del chart.

Inizializzatore: aggiunge al chart il primo nodo e il primo arco.

Predittore: scelto un arco incompleto che cerca una categoria C, aggiunge al chart un nuovo arco che se completato darebbe C.

Estensore: prende due archi. Uno incompleto che termina al nodo j e che cerca una C. L'altro completo che inizia al nodo j e aggiunge un arco in cui la C è stata trovata.

Analizzatore: Stesso comportamento dell'estensore tranne per il fatto che invece di prendere un arco completo prende una parola dallo stream di input.

Versione nondeterministica

Inizializzazione:

Chart $\leftarrow [0, 0, S' ::= \bullet S]$

While (altri archi possono essere aggiunti) do

arco \leftarrow scegli $[i, j, A ::= \alpha \bullet B \beta]$ dal chart

scegli uno fra i seguenti metodi che abbia successo:

Predittore:

scegli $(B ::= \gamma)$ dalla grammatica

aggiungi $[i, j, B ::= \bullet \gamma]$ al chart

Estensore:

scegli $[j, k, B ::= F \bullet]$ dal chart

aggiungi $[i, k, A ::= \alpha B \bullet \beta]$ al chart

Analizzatore:

se la parola corrente è di categoria B

allora aggiungi $[j, j+1, A ::= \alpha B \bullet \beta]$ al chart

Esempio

$_0 I_1 \text{SMELL}_2 \text{FOOD}_3$

- | | | |
|---|---------------|-----------------------------|
| 1) $[0, 0, S' ::= \bullet S]$ | \rightarrow | Inizializzatore |
| 2) $[0, 0, S ::= \bullet \text{NP VP}]$ | \rightarrow | predittore applicato a 1 |
| 3) $[0, 1, S ::= \text{NP} \bullet \text{VP}]$ | \rightarrow | analizzatore applicato a 2 |
| 4) $[1, 1, \text{VP} ::= \bullet \text{Verbo NP}]$ | \rightarrow | predittore applicato a 3 |
| 5) $[1, 2, \text{VP} ::= \text{Verbo} \bullet \text{NP}]$ | \rightarrow | analizzatore applicato a 4 |
| 6) $[1, 3, \text{VP} ::= \text{Verbo NP} \bullet]$ | \rightarrow | analizzatore applicato a 5 |
| 7) $[0, 3, S ::= \text{NP Verbo NP} \bullet]$ | \rightarrow | estensore applicato a 3 e 6 |

Versione deterministica

ChartParser(*stringa*, *grammatica*)

Aggiungi($[0, 0, S' ::= \bullet S]$)

Per ogni parola *x* in *stringa* chiama Analizzatore(*x*, *stringa*[*x*])

Aggiungi(arco):

se arco è già presente in chart[*fine*(arco)] non fa nulla altrimenti

inserisci arco in chart[*fine*(arco)]

se l'arco è completo allora Estensore(arco)

altrimenti Predittore(arco)

Analizzatore(*j*, *parola*):

per ogni arco $[i, \text{indice}, A ::= \alpha \bullet B \beta]$ in chart[*indice*]

se *parola* è di categoria B allora Aggiungi($[i, j+1, A ::= \alpha B \bullet \beta]$)

Predittore($[i, j, A ::= \alpha \bullet B \beta]$):

per ogni $(B ::= \gamma)$ in grammatica

Aggiungi($[i, j, B ::= \bullet \gamma]$)

Estensore ($[j, k, B ::= F \bullet]$):

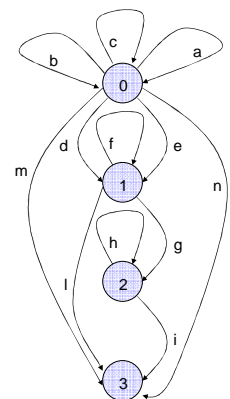
per ogni $[i, j, A ::= \alpha \bullet B \beta]$ in chart[*j*]

se $B=B'$ Aggiungi($[i, k, A ::= \alpha B' \bullet \beta]$)

Chart

I EAT PIZZA

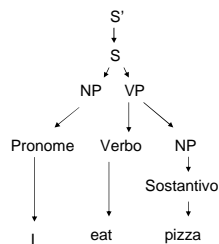
Arco	Metodo	Chart
a	Inizializzazione	$[0, 0, S' ::= \bullet S]$
b	Predittore(a)	$[0, 0, S ::= \bullet \text{NP VP}]$
c	Predittore(b)	$[0, 0, \text{NP} ::= \bullet \text{Pronome}]$
d	Analizzatore(1, I)	$[0, 1, \text{NP} ::= \text{Pronome} \bullet]$
e	Estensore(d) (b)	$[0, 1, S ::= \text{NP} \bullet \text{VP}]$
f	Predittore(e)	$[1, 1, \text{VP} ::= \bullet \text{Verbo NP}]$
g	Analizzatore(2, eat)	$[1, 2, \text{VP} ::= \text{Verbo} \bullet \text{NP}]$
h	Predittore(g)	$[2, 2, \text{NP} ::= \bullet \text{Sostantivo}]$
i	Analizzatore(3, pizza)	$[2, 3, \text{NP} ::= \text{Sostantivo} \bullet]$
l	Estensore(i) (g)	$[1, 3, \text{VP} ::= \text{Verbo NP} \bullet]$
m	Estensore(l) (e)	$[0, 3, S ::= \text{NP VP} \bullet]$
n	Estensore(m) (a)	$[0, 3, S' ::= S \bullet]$



Estrazione delle analisi

- Modificando opportunamente Estensore in modo che associ all'arco che produce, i due archi che fonde si possono estrarre gli alberi di analisi dal chart.

Arco	Metodo	Chart
a	Inizializzazione	[0,0, S' ::= • S]
b	Predittore(a)	[0,0, S ::= • NP VP]
c	Predittore(b)	[0,0, NP ::= • Pronome]
d	Analizzatore(1, I)	[0,1, NP ::= Pronome •]
e	Estensore(d) (b)	[0,1, S ::= NP • VP]
f	Predittore(e)	[1,1, VP ::= • Verbo NP]
g	Analizzatore(2, eat)	[1,2, VP ::= Verbo • NP]
h	Predittore(g)	[2,2, NP ::= • Sostantivo]
i	Analizzatore(3, pizza)	[2,3, NP ::= Sostantivo •]
l	Estensore(i) (g)	[1,3, VP ::= Verbo NP •]
m	Estensore(l) (e)	[0,3, S ::= NP VP •]
n	Estensore(m) (a)	[0,3, S' ::= S •]



Ambiguità

Fall leaves fall and spring leaves spring

- Le foglie autunnali cadono e le foglie primaverili nascono
- Le foglie autunnali cadono e la primavera lascia la primavera
- L'autunno lascia l'autunno e le foglie primaverili nascono
- L'autunno lascia l'autunno e la primavera lascia la primavera

Ambiguità derivante dalle interpretazioni molteplici di fall, leaves e spring.

Ogni sottofrase è ambigua in due modi per cui 4 interpretazioni.

N sottofrasi = 2^N interpretazioni

Impaccamento

- Si possono modificare Estensore e Aggiungi per ottenere una foresta impaccata di rappresentazioni.

$$[S [S \left\{ \left[\begin{array}{l} [NP \text{ fall leaves}] [VP \text{ fall}] \\ [NP \text{ fall}] [VP \text{ leaves fall}] \end{array} \right\} \text{ and } \left\{ \left[\begin{array}{l} [NP \text{ spring leaves}] [VP \text{ spring}] \\ [NP \text{ spring}] [VP \text{ leaves spring}] \end{array} \right\} \right]]$$

$2N$ invece di 2^N = ottimizzazione esponenziale dello spazio

Il parser più efficiente

- Il chart parser più efficiente è stato prodotto da Robert Moore per Microsoft

Preprocessing con fattorizzazione sinistra della grammatica:

$$A ::= BC \mid BD \rightarrow A ::= BA' \quad A' ::= C \mid D$$

Implementazione del chart come hash table

Rappresentazione solo della parte a destra del punto negli archi

$$A ::= \alpha X \cdot \beta \rightarrow A ::= \alpha \cdot \beta$$

Bottom-up filtering anticipato rispetto al top-down

	CT Grammar	ATIS Grammar	PT Grammar
Rules	24,456	4,592	15,039
Nonterminals	3,946	192	38
Terminals	1,032	357	47
# Test Sentences	162	98	30
Average Length	8.3	11.4	5.7
# Grammatical	150	70	30
Average # Parses	5.4	940	7.2×10^{27}

Table 1: Grammars and test sets for parser evaluations

	CT Grammar	ATIS Grammar	PT Grammar
LC ₁	4.3	15.6	45.0
LC ₂	3.4	11.9	43.0
LC ₃	3.1	11.6	41.8
LC ₄	2.7	11.8	42.3

Table 2: LC parsing algorithm performance comparisons

	CT Grammar	ATIS Grammar	PT Grammar
LC ₁ UTF	4.3	15.6	45.0
LC ₁ FLF	7.4	63.5	timed out
LC ₁ PLF	6.2	66.2	timed out
LC ₁ BUPM	3.6	11.7	34.1
LC ₂ UTF	3.4	11.9	43.0
LC ₂ FLF	5.1	38.2	timed out
LC ₂ PLF	4.2	37.7	timed out
LC ₂ BUPM	3.1	7.0	27.0
LC ₃ UTF	3.1	11.6	41.8
LC ₃ FLF	4.2	12.3	45.4
LC ₃ PLF	3.8	12.1	43.6
LC ₃ BUPM	5.0	17.1	64.6
LC ₄ UTF	2.7	11.8	42.3
LC ₄ FLF	3.6	11.9	46.6
LC ₄ PLF	3.2	11.7	44.4
LC ₄ BUPM	3.2	14.7	63.6

Table 3: LC parsing grammar transformation performance comparisons

Parsing probabilistico

- Un parser probabilistico sceglie le regole da espandere in base a probabilità calcolate con numerosi meriti per arrivare il prima possibile ad un'analisi e restituirla come "più probabile"

Si cerca di massimizzare $p(N_{j,k}^i | t_{0,n})$ dove $t_{0,n}$ è la sequenza dei primi $n+1$ tokens e $N_{j,k}^i$ è la categoria i -esima che copre la frase da j a k .

Conoscendo queste probabilità si possono stabilire criteri di scelta delle categorie da espandere per raggiungere prima l'analisi voluta.

Parsing probabilistico

La probabilità che la categoria N^i copra da j a k data la stringa di caratteri da 0 a n è esprimibile in termini di α e β . α è dato dalla probabilità esterna mentre β è la probabilità interna.

Poiché α e il denominatore rappresentano l'influenza delle parole che circondano N^i , possiamo eliminarle e usare solo β . Questo metodo è detto β diretto. Può portare a parser inefficienti perché eliminando il denominatore non viene presa in considerazione l'influenza della lunghezza della sequenza, preferendo così sempre sequenze corte.

$$\begin{aligned}
 p(N_{j,k}^i | t_{0,n}) &= \frac{p(N_{j,k}^i, t_{0,n})}{p(t_{0,n})} \\
 &= \frac{p(N^i | t_{0,n}) p(t_{0,j}, t_{k,n})}{p(t_{0,n})} \\
 &= \frac{p(t_{0,j}, N_{j,k}^i, t_{k,n}) p(t_{j,k} | t_{0,j}, N_{j,k}^i, t_{k,n})}{p(t_{0,n})} \\
 p(N_{j,k}^i | t_{0,n}) &\approx \frac{p(t_{0,j}, N_{j,k}^i, t_{k,n}) p(t_{j,k} | N_{j,k}^i)}{p(t_{0,n})} \\
 &= \frac{p(t_{0,j}, N_{j,k}^i, t_{k,n}) \beta(N_{j,k}^i)}{p(t_{0,n})} \\
 p(N_{j,k}^i | t_{0,n}) &\approx \frac{\alpha(N_{j,k}^i) \beta(N_{j,k}^i)}{p(t_{0,n})}
 \end{aligned}$$

Altri metodi

β normalizzato: per tenere conto della lunghezza delle sequenze si utilizza come parametro per la scelta $\sqrt{\beta(N_{j,k}^i)}$.

α, β normalizzato: per tenere conto della lunghezza delle sequenze e delle parole già parse si usa $\sqrt{\alpha_L(N_{j,k}^i) \beta(N_{j,k}^i)}$.

Trigram: Si fanno stime sulla frequenza di gruppi di tre parole e non di tutto il contesto precedente $\frac{p(N^i) \beta(N_{j,k}^i)}{p(t_{j,k} | t_{j-2,j})}$.

Prefix: Si fanno stime sulla frequenza del gruppo di parole precedente normalizzato $\frac{\alpha_L(N_{j,k}^i) \beta(N_{j,k}^i)}{p(t_{0,k})}$.

Risultati

Figure of Merit	%E	%non-0 E	%popped
straight β	97.6	97.5	93.8
normalized β	34.7	31.6	61.5
normalized $\alpha_L\beta$	39.7	36.4	57.3
trigram estimate	25.2	21.7	44.3
prefix estimate	21.8	17.4	38.3

Figure of Merit	CPU time
straight β	3966
normalized β	1631
normalized $\alpha_L\beta$	68660
trigram estimate	1547
prefix estimate	26520

I modelli migliori (prefix e $\alpha_L\beta$) risultano costosi poiché il calcolo dei parametri è impegnativo

Le tabelle riassumono i risultati dei vari metodi.

%E indica la percentuale di archi inseriti rispetto ad un parser esaustivo.

%non-0 E sono le produzioni non vuote

%popped è la percentuale degli archi utilizzati per arrivare ad una soluzione.

Parsing incrementale

- 1) Possibilità di eseguire operazioni su parti di stringa già parseate in modo da garantire una riorganizzazione efficiente degli alberi di analisi.
- 2) L'input viene fornito al parser man mano che l'utente lo inserisce
- 3) Il parser, monitorando il processo di input, guida l'utente nella composizione del testo dando suggerimenti o chiedendo chiarificazioni.

Operazioni disponibili

Si definisce una relazione di dipendenza D fra archi del chart tale che un arco a è in relazione D con un arco b (a D b) se a è stato inserito nel chart o tramite predizione da b o tramite estensione con b.

Inserzione:

se avviene all'estrema destra dello stream di input viene trattata come una semplice chiamata ad Analizzatore .

se avviene in posizione x precedente allora vengono

- 1) rimossi gli archi dipendenti da quello in posizione x
- 2) aggiornati tutti gli archi con posizioni maggiori di x ai nuovi valori shiftati
- 3) rieseguito Analizzatore e aggiornato il chart

Operazioni disponibili (2)

Delezione:

se avviene all'estrema destra dello stream di input vengono rimossi tutti gli archi in relazione D* con l'ultimo aggiunto.

se avviene all'interno della stringa già parseata in posizione x:

- 1) Si rimuovono tutti gli archi in relazione D* con l'arco preterminale in posizione x
- 2) si shiftano a sinistra tutti gli archi con posizioni maggiori di x

Sostituzione:

si può realizzare attraverso combinazione di delezione e inserzione ma viene spesso implementata separatamente per motivi di efficienza.

Limiti e Benefici

- 1) Il parser incrementale così costruito non tiene conto delle relazioni fra periodi diversi, trattando ogni frase come indipendente. Questa limitazione può essere superata perdendo però notevolmente in efficienza.
- 2) Mantenere le dipendenze causa un aumento considerevole della memoria occupata, per questo si è pensato a varianti che utilizzino meno informazione peggiorando lievemente il tempo d'esecuzione
- 3) Un parser incrementale è un tool utilissimo in contesti come linguaggi di programmazione (individuazione di bug prima della compilazione), e linguaggio naturale (composizione testi assistita).

Conclusioni

In definitiva il problema dell'analisi sintattica non ha una soluzione unica ma viene adattato di volta in volta al problema che si vuole affrontare.

In ambito statistico sono infatti nati i parser probabilistici, in uno più applicativo i parser incrementali.

I passi da fare sono quelli di unificare gli sforzi verso un unico modello di parser che sia efficiente, probabilistico e incrementale.

In più si dovrebbe riuscire a compiere durante l'analisi sintattica una parte importante di analisi semantica per guidare soprattutto la disambiguazione e la scelta dell'analisi nel chart

Bibliografia

- Russel-Norving: **Intelligenza Artificiale**
- Robert Moore: **Improved Left-Corner Chart Parsing for Large Context-Free Grammars**
- A. Caraballo, E. Charniak: **Figures of Merit for Best-First Probabilistic Chart Parsing**
- Mats Wirèn: **Interactive Incremental Chart Parsing**