

07-semantic-analysis

Calcoliamo informazioni aggiuntive una volta che la struttura sintattica è conosciuta, queste informazioni sono al di là delle capacità delle grammatiche *context-free*, diciamo "vanilla". Tipicamente in questa fase:

- Popoliamo la tabella dei simboli (*symbol table*) dopo ogni dichiarazione.
- Facciamo dell'inferenza sui tipi e dei controlli su di essi nelle espressioni e dichiarazioni. Questo tipo di analisi si divide in due categorie:
- Analisi richiesta per stabilire la correttezza
- Analisi per aumentare l'efficienza del programma tradotto.

Il modo più semplice per implementare l'analisi semantica è identificare proprietà (attributi) di un simbolo della grammatica, e scrivere delle regole (regole semantiche) per descrivere come calcolare proprietà legate alle produzioni della grammatica. L'insieme dato da (attributi, regole, ecc....) viene detto grammatica attribuita o *syntax-directed definition*.

Come struttura la scelta ottimale ricade su un *abstract syntax tree* che risulta essere una rappresentazione compressa di un *derivation tree*.

Syntax-Directed Definitions

Introduciamo ora meglio il concetto di grammatica attribuita (Syntax-Directed Definitions) per gli amici *SDD*.

Sono delle grammatiche *context-free* arricchite con attributi e regole:

- **Attributi:** associati con dei simboli della grammatica, possono essere tipi, numeri, riferimenti alla tabella dei simboli ecc....
- **Regole semantiche:** associate con ogni produzione, tipicamente regolano il calcolo di attributi in funzione di attributi degli altri simboli della produzione.
Sia i simboli che le regole sono usati per dare senso (con l'analisi semantica) a quello che è espresso come un flusso di token.

Tipi di attributi

Gli attributi dei non-terminali sono suddivisi in due categorie:

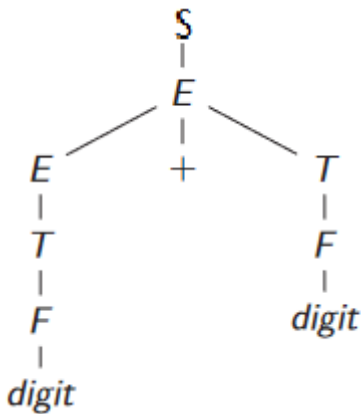
- **Sintetizzati:** gli attributi del driver sono definiti come una funzione degli attributi dei simboli della produzione.
- **Ereditati:** gli attributi dei non-terminali nel body sono definiti come funzione degli attributi dei simboli della produzione.
Gli attributi dei terminali possono essere solo sintetizzati poiché forniti dall'analizzatore lessicale e non esiste regola per calcolarli.

Esempio

Prendiamo la nostra grammatica per le espressioni aritmetiche e proviamo a scrivere delle regole e degli attributi per essa.

$$\mathcal{G}: \begin{cases} S \rightarrow E \\ E \rightarrow E + T | T \\ T \rightarrow T * F | F \\ F \rightarrow (E) | digit \end{cases}$$

Consideriamo ora il *parse tree* SLR(1), e diciamo che se un *digit* ha valore 3 e l'altro 4 il risultato deve essere 7.

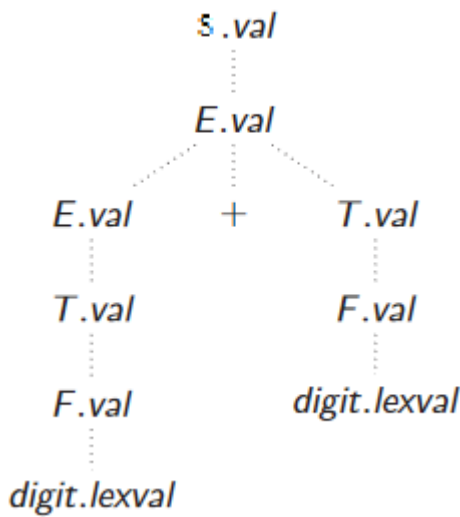


Vediamo quindi che forma deve avere il suo *SDD* per poter computare il valore finale in *S*.

$$\begin{aligned} S &\rightarrow E && \{S.val = E.val\} \\ E &\rightarrow E_1 + T && \{E.val = E_1.val + T.val\} \\ E &\rightarrow T && \{E.val = T.val\} \\ T &\rightarrow T_1 * F && \{T.val = T_1.val * F.val\} \\ T &\rightarrow F && \{T.val = F.val\} \\ F &\rightarrow (E) && \{F.val = E.val\} \\ F &\rightarrow digit && \{F.val = digit.lexval\} \end{aligned}$$

Ora le regole sembrano abbastanza semplici e non necessitano di ulteriori spiegazioni tranne per un paio di appunti:

- La prima produzione $S \rightarrow E$ è una produzione aggiuntiva non necessaria alla grammatica in se.
- Usiamo i numeri a pedice per differenziare lo stesso non-terminale nella stessa produzione.
- *digit.lexval* è il valore che viene trovato nella tabella dei simboli.
L'*abstract syntax tree* risulta quindi in:



Valutazione in ordine di un SDD

Non è sempre possibile che un SDD possa essere valutato, quindi definiamo un grafo (orientato) delle dipendenze degli attributi del nostro SDD e verifichiamo che non ci siano conflitti.

- Impostare un nodo del grafo delle dipendenze per ogni attributo associato con ogni nodo del *parse tree*.
- Per ogni attributo $X.x$ usato per definire l'attributo $Y.y$ creiamo un arco dal nodo di $X.x$ al nodo di $Y.y$.

Una volta creato il grafo delle dipendenze dobbiamo trovare un ordinamento topologico per il grafo, se l'ordinamento non esiste allora l'*SDD* non è valutabile.

Se invece un ordinamento esiste allora possiamo valutare l'*SDD* e abbiamo trovato anche un ordine in cui farlo.

Quando un *SDD* ha sia attributi ereditati che sintetizzati non ci sono garanzie che esista un ordinamento topologico, infatti potrebbe esserci un ciclo all'interno del grafo quindi nessun sorting è possibile.

$$A \rightarrow B \quad \{A.s = B.i; \quad B.i = A.s + 7\}$$

In questo caso l'attributo sintetizzato di A ha bisogno dell'attributo ereditato di B e viceversa.

Esistono due classi di *SDD* per le quali è garantita l'esistenza di un ordinamento topologico.

1. **S-attributed SDDs:** ci sono solo attributi sintetizzati quindi ci basta fare una visita in post-ordine.
2. **L-attributed SDDs:** attributi sia sintetizzati sia ereditati tali che:

Per ogni produzione $A \rightarrow X_1 \cdots X_n$ la definizione di ogni $X_j.i$ usa al più:

Attributi ereditati da A oppure

Attributi ereditati o sintetizzati dai fratelli a sinistra, ovvero X_1, \dots, X_{j-1}

Gli SDD* S-attribuiti sono ideali per il parsing bottom-up perchè l'albero può essere valutato mentre si fa il parsing.

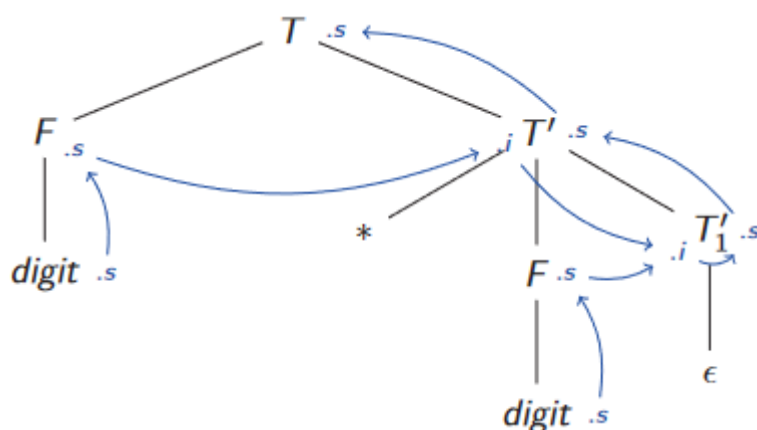
Gli L-attribuiti sono convenienti con il parsing top-down perchè ho attributi ereditato solo da sinistra e nel parsing top-down faccio derivazioni leftmost quindi ez.

Esempio

Prendiamo la grammatica LL(1) per le operazioni aritmetiche:

$$\begin{cases} V \rightarrow E \\ E \rightarrow TE' \\ E' \rightarrow +TE' | \varepsilon \\ T \rightarrow FT' \\ T' \rightarrow *FT' | \varepsilon \\ F \rightarrow (E) | digit \end{cases}$$

Facciamo il parsing di $3 * 5$, come si vede partiamo da $V \Rightarrow E \Rightarrow TE'$, ora la Quaglia ha disegnato solo il sottoalbero che ha come radice T visto che dobbiamo $E' \Rightarrow \varepsilon$. Essendo un parsing top-down utilizziamo un albero L-attribuito:

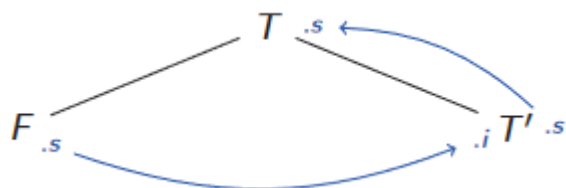


Ogni $.s$ sta per attributi sintetizzato mentre ogni $.i$ sta per ereditato. Iniziamo quindi con il dare le regole alle varie produzioni:

- Prima di tutto possiamo notare che quando T' viene copiato in T si ha già il valore della moltiplicazione.

$$T \rightarrow FT' \quad \{T'.i = F.s; \quad T.s = T'.s\}$$

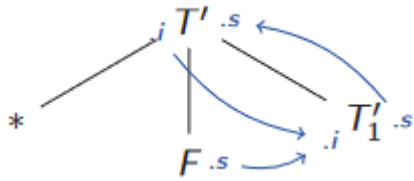
Quindi passiamo T' come valore ereditato $F.s$ e assegniamo a T sintetizzato il valore di $T'.s$.



- Ora concentriamoci nelle produzioni di T' , nel nostro caso abbiamo una moltiplicazione, quindi sappiamo che in $T'.i$ è memorizzato il valore del membro di sinistra per cui assegniamo a $T'_1.i$ il valore della moltiplicazione.

Dopodichè per poter risalire l'albero fino alla radice assegniamo a $T'.s$ il valore del risultato finale del ramo destro che sarà memorizzato in $T'_1.s$.

$$T' \rightarrow *FT'_1 \quad \{T'_1.i = F.s * T'.i; \quad T'.s = T'_1.s\}$$



Valutazione durante il parsing bottom-up

Il nostro obbiettivo è implementare la traduzione della parola durante il processo di parsing anzichè ottenere il parsing tree poi annotarlo ed infine valutarlo.

Il caso più semplice in cui questa elaborazione può essere fatta è durante l'algoritmo shift/reduce con un *SDD* S-attribuito.

L'idea è quella di tenere oltre ai due stack per gli stati *stSt* e per i simboli *symSt* un ulteriore stack per gli attributi *semSt*.

Ogni volta che faccio una uno shift (leggo una valore) faccio un push della prima istruzione e.g. `push digit.lexval` e ogni volta che faccio una riduzione faccio il `pop` dei valori nel body e poi il `push` del valore del driver.

N.B. devo mantenere lo stack *semSt* "allineato", questo vuol dire che devo inserire dei valori particolari quando leggo simboli che non sono direttamente dei valori che mi interessano, e.g. se ho letto $3 + 4$ nello stack semantico dovrò avere una sequenza tipo $3 - 0 - 4$ nella quale 0 simboleggia il mio valore speciale.

Caso studio

Pensiamo di dover tradurre un numero inserito come stringa in un intero base 10.

Una grammatica LALR(1) molto semplice potrebbe essere :

$$S \rightarrow Digits \quad \{print(D.val)\}$$

$$Digits \rightarrow Digits_1 d \quad \{Digits.val = Digits_1.val * 10 + d.lexval\}$$

$$Digits \rightarrow d \quad \{Digits.val = d.lexval\}$$

Notiamo che l'*SDD* è S-attribuito quindi l'albero è valutabile e siamo apposto.

Aggiungiamo ora un livello di difficoltà, diciamo che se la stringa è preceduta dal terminale *o* allora dobbiamo tradurre il numero in base 8. (La Quaglia ha un concetto strano di base 8).

Praticamente se abbiamo la *o* dobbiamo moltiplicare per 8.

Andiamo a tentativi per trovare una grammatica buona.

1.

$$\begin{cases} S \rightarrow Num \\ Num \rightarrow o Digits | Digits \\ Digits \rightarrow Digits d | d \end{cases}$$

Per essere sicuri di poter valutare l'albero gli attributi devono essere sintetizzati.

Ma in questo caso non è possibile perchè noi prima facciamo prima una riduzione

$Digits \rightarrow d$ e poi tutte le riduzioni $Digits \rightarrow Digit\ d$, quindi non sappiamo ancora fino all'ultima riduzione $Num \rightarrow o\ Digits|Digits$ la base in cui convertire, unlucky.

2. Proviamo con la grammtica:

$$\left\{ \begin{array}{ll} S \rightarrow Num & \{print(Num.v)\} \\ Num \rightarrow oO & \{Num.v = O.v\} \\ Num \rightarrow D & \{Num.v = D.v\} \\ O \rightarrow O_1d & \{O.v = O_1.v * 8 + d.lexval\} \\ O \rightarrow d & \{O.v = d.lexval\} \\ D \rightarrow D_1d & \{D.v = D_1.v * 10 + d.lexval\} \\ D \rightarrow d & \{D.v = d.lexval\} \end{array} \right.$$

Buona per attributi sintetizzati e LALR(1) ma c'è troppa ridondanza.

3. Proviamo quindi a sfoltire un po' la grammatica:

$$\left\{ \begin{array}{ll} S \rightarrow Num & \{print(Num.v)\} \\ Num \rightarrow Octal\ Digits & \{Num.v = Digits.v\} \\ Num \rightarrow Decimal\ Digits & \{Num.v = Digits.v\} \\ Octal \rightarrow o & \{base = 8\} \\ Decimal \rightarrow \varepsilon & \{base = 10\} \\ Digits \rightarrow Digits_1\ d & \{D.v = D_1.v * base + d.lexval\} \\ Digits \rightarrow d & \{Digits.v = d.lexval\} \end{array} \right.$$

Abbiamo ancora margine di miglioramento perchè stiamo usando una variabile globale `base`, come in programmazione usiamo il meno possibile delle variabili globali.

4. Proviamo quindi a unire il terminale che gestisce le basi, così che ne venga subito fatta una riduzione.

$$\left\{ \begin{array}{ll} S \rightarrow D & \{print(D.v)\} \\ D \rightarrow D_1\ d & \{D.v = D_1.v * D_1.base + d.lexval; \ D.base = D_1.base\} \\ D \rightarrow B\ d & \{D.v = d.lexval; \ D.base = B.val\} \\ B \rightarrow o & \{B.val = 8\} \\ B \rightarrow \varepsilon & \{B.val = 10\} \end{array} \right.$$

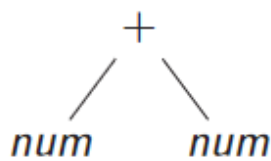
questo sembra un buon compromesso perchè se parsiamo la stringa *odd* la prima cosa che facciamo è leggere *o* e ridurlo in *B* quindi sappiamo la base, se invece abbiamo una parola tipo *dd* vuol dire che faremo la riduzione $B \rightarrow \varepsilon$ e quindi sappiamo che la base è 10.

Abstract syntax trees

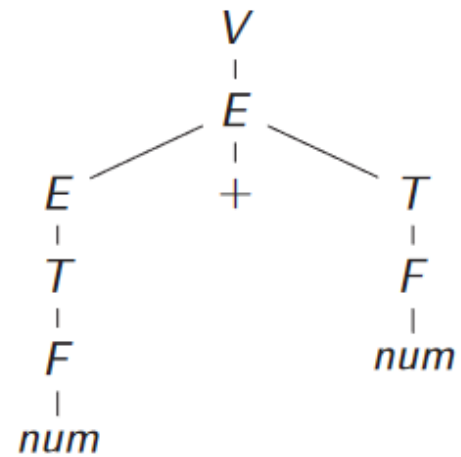
Una rappresentazione più compatta del *parse tree*, spesso usato come rappresentazione intermedia.

Non esistono regole generali per costruirli ma dipende tutto dalla grammatica che andiamo ad analizzare e dalle scelte implementative.

Devono contenere tutte le informazioni per portare avanti l'analisi.



instead of



D'ora in poi per semplicità li chiameremo *AST*.

Memorizzazione di un AST

Di fatto stiamo andando a creare in memoria un albero nel quale ogni nodo è un passaggio intermedio del processo di analisi di una parola.

I nodi dell'*AST* sono così strutturati:

- **Foglie:** contengono dei link alla tabella dei simboli, viene creata una foglia per ogni valore diretto (e.g. valore numero) trovato nella parola.
- **Nodi:** sono la struttura dell'albero e rappresentano i passaggi effettuati.

Creazione di un AST

Potremmo crearlo dopo aver trovato il *parse tree*, ma il nostro obiettivo è farlo mentre leggiamo la parola quindi durante il parsing.

Dobbiamo avere una grammatica S-attribuita per poterlo fare durante l'analisi sintattica (parsing).

Facciamo un paio di assunzioni su 2 funzioni a nostra disposizione:

- `newLeaf(label, val)` crea una nuova foglia con due campi, la label del nodo e il valore.
- `newNode(label, c1, ..., ck)` crea un nodo interno che può avere fino a k figli, label è l'identificatore del nodo mentre i successivi argomenti sono riferimenti ai figli.

Esempio

Prendiamo la grammatica LALR(1) delle espressioni aritmetiche:

$$\left\{ \begin{array}{ll} E \rightarrow E_1 + T & \{E.\text{node} = \text{newNode}('+', E_1.\text{node}, T.\text{node})\} \\ E \rightarrow E_1 - T & \{E.\text{node} = \text{newNode}('-', E_1.\text{node}, T.\text{node})\} \\ E \rightarrow T & \{E.\text{node} = T.\text{node}\} \\ T \rightarrow (E) & \{T.\text{node} = E.\text{node}\} \\ T \rightarrow id & \{T.\text{node} = \text{newLeaf}(id, id.\text{entry})\} \\ T \rightarrow num & \{T.\text{node} = \text{newLeaf}(num, num.\text{lexval})\} \end{array} \right.$$

Ora proviamo a fare il parsing del lessema $a - 4 + c$ che risulta essere $id - num + id$.
Essendo noi dei parser umani possiamo fare le riduzioni a occhio senza la tabella di parsing:

1. $T \rightarrow id$

2. $E \rightarrow T$

3. $T \rightarrow num$

4. $E \rightarrow E - T$

5. $T \rightarrow id$

6. $E \rightarrow E + T$

Ora vediamo graficamente le regole eseguite ogni riduzione:

7. $T \rightarrow id \quad \{T.node = newLeaf(id, id.entry)\}$



8. $E \rightarrow T \quad \{E.node = T.node\}$



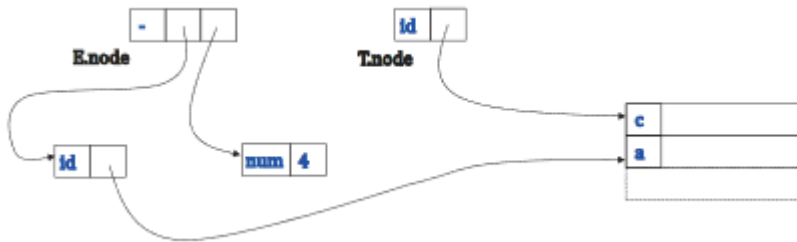
9. $T \rightarrow num \quad \{T.node = newLeaf(num, num.lexval)\}$



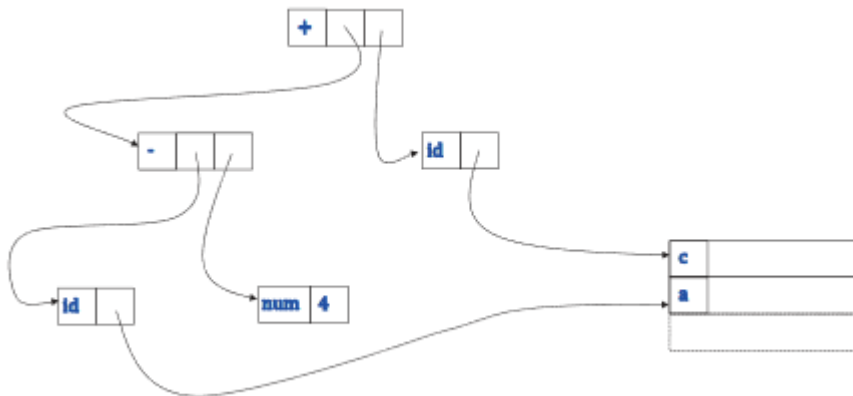
10. $E \rightarrow E_1 - T \quad \{E.node = newNode('-', E_1.node, T.node)\}$



11. $T \rightarrow id \quad \{T.node = newLeaf(id, id.entry)\}$



12. $E \rightarrow E_1 + T \quad \{E.node = newNode('+', E_1.node, T.node)\}$



Abstract syntax trees per LL(1)

Dovendo fare un parsing top-down ci conviene usare una grammatica L-attribuita, quindi possiamo avere attributi ereditati, come facciamo a costruire un AST per loro?

Dobbiamo passare riferimenti ai nodi anzichè valori e fare il giro lungo, ovvero produrre il *parse tree* annotato poi definire un grafo delle dipendenze ed infine costruire l'AST.

Esempio

$$\begin{cases}
 E \rightarrow TE' & \{E.node = E'.node; E'.i = T.node\} \\
 E' \rightarrow +TE'_1 & \{E'.node = E'_1.node; E'_1.i = newNode('+', E'.i, T.node)\} \\
 E' \rightarrow -TE'_1 & \{E'.node = E'_1.node; E'_1.i = newNode('-', E'.i, T.node)\} \\
 E' \rightarrow \varepsilon & \{E'.node = E'.i\} \\
 T \rightarrow (E) & \{T.node = E.node\} \\
 T \rightarrow id & \{T.node = newLeaf(id, id.entry)\} \\
 T \rightarrow num & \{T.node = newLeaf(num, num.lexval)\}
 \end{cases}$$

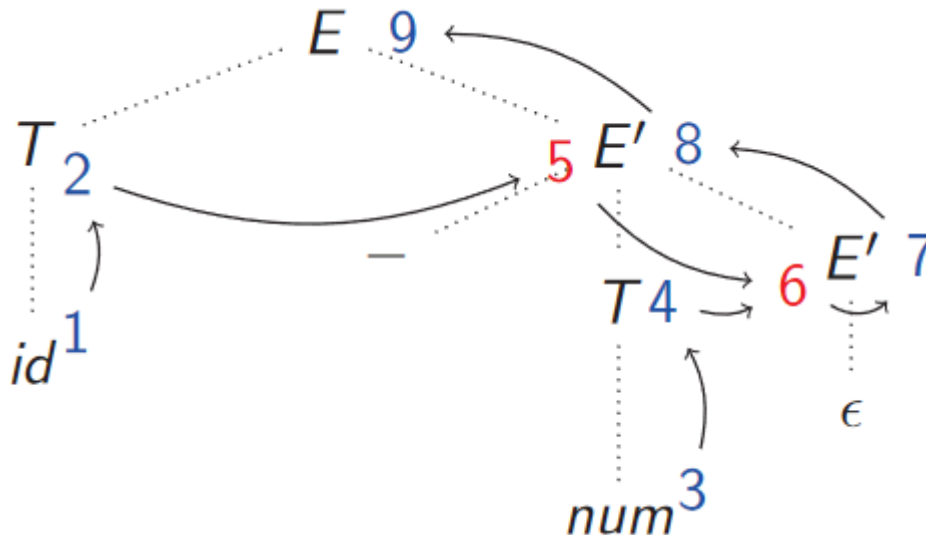
Proviamo quindi a fare il parsing del lessema $a - 4$ ovvero $id - num$.

Essendo un parsing top-down con derivazioni leftmost le derivazioni saranno:

$$E \Rightarrow TE' \Rightarrow idE' \Rightarrow id - TE' \Rightarrow id - numE' \Rightarrow id - num$$

L'ordinamento topologico ci restituirà il seguente ordine di valutazione:

//Nei passaggi successivi creeremo dei nodi in seguito ci riferiremo agli stessi nodi con le loro *label*



1. $id.entry$ (riferimento alla entry di a nella *symbol table*)
2. $T.node = newLeaf(2, id.entry)$
vado in E' e valuto quello che sta sotto
3. $num.lexval = 4$
4. $T.node = newLeaf(4, num.lexval)$
5. Sfruttiamo $E'.i = T.node = 2$ della produzione $E \rightarrow TE'$ per portarci il numero contenuti in id .
6. Ora possiamo scendere in con $E'_1.i = newNode('-', 2, 4)$ nel punto 6.
7. Fortunatamente abbiamo una produzione che va in ϵ quindi creiamo un riferimento al nodo della sottrazione con $E'.node = E'.i$
8. Creiamo un altro riferimento al nodo della sottrazione con $E'.node = E'_1.node$ della produzione $E' \rightarrow -TE'_1$
9. Finisco di risalire l'albero creando un riferimento allo stesso nodo con $E.node = E'.node$ della prima produzione.