

Universidade da Beira Interior

Desenho de Linguagens de Programação e de Compiladores

Simão Melo de Sousa

Aula 5 - Análise Sintáctica Descendente

não é só o léxico que importa, a sintaxe e a morfologia da frase também



Man eating piranha mistakenly sold as pet fish

Babies are what the mother eats

o objectivo da análise sintáctica é reconhecer as frases pertencendo à sintaxe da linguagem

o que o fluxo de caracteres é para a análise léxica, o fluxo dos lexemas – resultado desta mesma análise – o são para a análise sintáctica

a saída da análise sintáctica é uma árvore de sintaxe abstracta

sequência de caracteres
 "fun x -> (x + 1)"

↓
análise léxica

↓
 sequência de lexemas

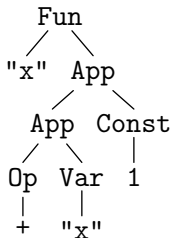
fun x -> (x + 1)

sequência de lexemas

fun x -> (x + 1)

↓
análise sintáctica

↓
 sintaxe abstracta



em particular, a análise sintáctica deve detectar os erros de sintaxe e

- localisá-los com precisão
- os identificar (na maioria dos casos assinalam-se por « Syntax Error » mas também « missing ending parenthesis », etc.)
- eventualmente, retomar a análise após detecção e identificação de erro para procurar outros potenciais erros

para a análise sintáctica vamos usar

- uma **gramática livre de contexto** (também conhecida como **gramática algébrica**, *context free grammar*, ou ainda CFG) para descrever a sintaxe
- um **autómato com pilha** para a reconhecer

é o análogo das expressões regulares e dos autómatos de estados finitos utilizados na análise léxica

Definição

Uma gramática algébrica é um 4-tuplo (N, T, S, R) onde

- N é um conjunto finito de **símbolos não terminais**
- T é um conjunto finito de **símbolos terminais**
- $S \in N$ é um símbolo de partida/inicial (designado de **axioma**)
- $R \subseteq N \times (N \cup T)^*$ é um conjunto finito de **regras de produção**

$N = \{E\}$, $T = \{+, *, (,), \text{int}\}$, $S = E$,
 e $R = \{(E, E+E), (E, E*E), (E, (E)), (E, \text{int})\}$

na prática, apresentamos as regras na forma seguinte:

$$\begin{array}{l}
 E \rightarrow E + E \\
 \quad | \quad E * E \\
 \quad | \quad (E) \\
 \quad | \quad \text{int}
 \end{array}$$

os terminais da gramática são os lexemas produzidos pela análise léxica

int designa aqui o lexema correspondendo a uma constante inteira
 (*i.e.* a sua natureza, e não o seu valor)

Definição

Uma palavra $u \in (N \cup T)^*$ **deriva-se** numa palavra $v \in (N \cup T)^*$, com a notação $u \rightarrow v$, se existe uma decomposição

$$u = u_1 X u_2$$

com $X \in N$, $X \rightarrow \beta \in R$ e

$$v = u_1 \beta u_2$$

exemplo :

$$\underbrace{E *}_{u_1} \left(\underbrace{E}_X \right) \underbrace{\quad}_{u_2} \rightarrow E * \left(\underbrace{E + E}_{\beta} \right)$$

uma sequência $w_1 \rightarrow w_2 \rightarrow \dots \rightarrow w_n$ é designada de **derivação**

falamos de **derivação esquerda** (resp. **direita**) se o não-terminal reduzido é sistematicamente o mais a esquerda *i.e.* $u_1 \in T^*$ (resp. o mais a direita *i.e.* $u_2 \in T^*$)

nota-se \rightarrow^* o fecho reflexivo transitivo de \rightarrow

$$\begin{aligned} E &\rightarrow E * E \\ &\rightarrow \text{int} * E \\ &\rightarrow \text{int} * (E) \\ &\rightarrow \text{int} * (E + E) \\ &\rightarrow \text{int} * (\text{int} + E) \\ &\rightarrow \text{int} * (\text{int} + \text{int}) \end{aligned}$$

em particular

$$E \rightarrow^* \text{int} * (\text{int} + \text{int})$$

Definição

A **linguagem** definida por uma gramática algébrica $G = (N, T, S, R)$ é o conjunto das palavras T^* derivadas do axioma, i.e.

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

no exemplo anterior : $\text{int} * (\text{int} + \text{int}) \in L(G)$

outro exemplo : para a gramática $S \rightarrow aSb \mid \epsilon$, temos

$$L = \{ a^n b^n \mid n \in \mathbb{N} \}$$

(exercício : demonstrar tal resultado)

Definição

A qualquer derivação $S \rightarrow^* w$, podemos associar uma **árvore de derivação**, cujos nodos possuem etiquetas definidas da seguinte forma

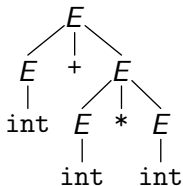
- a raíz é S
- as folhas formam a palavra w na ordem infixada
- todo o nodo interno X é um não-terminal cujos filhos tem por etiqueta $\beta \in (N \cup T)^*$ com $X \rightarrow \beta$ regra de derivação

de notar: não confundir com uma árvore de sintaxe abstracta, são coisas diferentes

a derivação esquerda

$$E \rightarrow E + E \rightarrow \text{int} + E \rightarrow \text{int} + E * E \rightarrow \text{int} + \text{int} * E \rightarrow \text{int} + \text{int} * \text{int}$$

devolve a árvore de derivação



mas a derivação direita

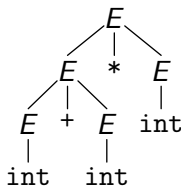
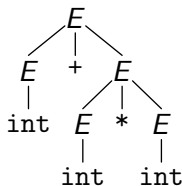
$$E \rightarrow E + E \rightarrow E + E * E \rightarrow E + E * \text{int} \rightarrow E + \text{int} * \text{int} \rightarrow \text{int} + \text{int} * \text{int}$$

também

Definição

Uma gramática é dita **ambígua** se uma palavra da sua linguagem admite mais do que uma árvore de derivação

exemplo : a palavra `int + int * int` admite as duas árvores de derivação

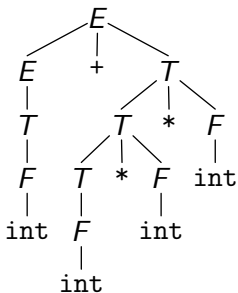


é no entanto e neste caso possível propor uma outra gramática, não ambígua, que define a mesma linguagem

$$\begin{array}{l} E \rightarrow E + T \\ \quad | \quad T \\ T \rightarrow T * F \\ \quad | \quad F \\ F \rightarrow (E) \\ \quad | \quad \text{int} \end{array}$$

esta nova gramática traduz a prioridade do produto sobre a soma e a escolha da associatividade esquerda para estas duas operações

assim, a palavra $\text{int} + \text{int} * \text{int} * \text{int}$ só tem doravante uma árvore de derivação, ou seja



que corresponde a derivação esquerda

$$\begin{aligned}
 E &\rightarrow E + T \rightarrow T + T \rightarrow F + T \rightarrow \text{int} + T \rightarrow \text{int} + T * F \\
 &\rightarrow \text{int} + T * F * F \rightarrow \text{int} + F * F * F \rightarrow \text{int} + \text{int} * F * F \\
 &\rightarrow \text{int} + \text{int} * \text{int} * F \rightarrow \text{int} + \text{int} * \text{int} * \text{int}
 \end{aligned}$$

determinar se uma gramática é ou não ambígua **não é decidível**

(lembrete :*decidível* significa que não podemos definir um algoritmo/programa que, para toda a entrada possível, termina e responde correctamente pela positiva ou pela negativa)

vamos utilizar **critérios de decisão suficientes** para garantir que uma gramática não é ambígua , e para as quais sabemos também decidir a pertença à linguagem de forma eficiente (com recurso a um autómato de pilha determinista)

as classes de gramáticas definidas por estes critérios chamam-se LL(1), LR(0), SLR(1), LALR(1), LR(1), etc.

execução por um autômato de pilha

reconhecer linguagens regulares = autómatos finitos

reconhecer linguagens livres de contexto = autómatos com pilha

num autómato finito, numa transição entre dois estados como por exemplo: $1 \xrightarrow{a} 2$, passar de 1 para 2 basta que a entrada forneça um a

num autómato com pilha, a diferença é que a transição é executada tendo em conta igualmente um **histórico** das sequências de caracteres já analisados: **pilha = histórico**

conjunto dos estados: Q

estado inicial: $q_0 \in Q$

estados finais: $F \subseteq Q$

histórico/pilha: Q^+

função/relação de transição: $\Delta \in Q^+ \times (A \cup \{\epsilon\}) \times Q^*$

princípio:

- A cada estado está associada uma palavra de Q^+ : a pilha.
- Uma transição = uma transformação no topo da pilha

consideremos a transição por x de $m_1 \in Q^+$ para $m_2 \in Q^*$ (notação $m_1 \xrightarrow{x} m_2$), então queremos que qualquer que seja a pilha com prefixo m_1 , digamos mm_1 , seja possível utilizar esta transição de cada vez que a entrada apresente x .

notação $mm_1 \xrightarrow{x} mm_2$

extendemos \rightarrow às palavras de $(A \cup \{\epsilon\})^*$

se $m_1 \xrightarrow{\omega_1} m_2$ e $m_2 \xrightarrow{\omega_2} m_3$ então $m_1 \xrightarrow{\omega_1\omega_2} m_3$

uma palavra ω é **reconhecida** por um autómato com pilha se $q_0 \xrightarrow{\omega} mf$ com $f \in F$ no topo da pilha

construção dum autómato com pilha a partir duma CFG

seja $G = (N, T, S, P)$ uma gramática livre de contexto

vamos considerar todos os pares possíveis (Y, α, β) , designados de **item**, tais que $Y \in N$, $\alpha, \beta \in (N \cup T)^*$ e $(Y ::= \alpha\beta) \in P$
notação $[Y \rightarrow \alpha.\beta]$ (poderemos omitir ocasionalmente os $[\]$)

alguns itens particulares:

se $\alpha = \epsilon$, $[Y \rightarrow .\beta]$

se $\beta = \epsilon$, $[Y \rightarrow \alpha.]$ este item é designado de **completo**

se $\alpha = \beta = \epsilon$, $[Y \rightarrow .]$

construção dum autómato com pilha a partir duma CFG

para construir o autómato com pilha correspondente introduzimos, se necessário, um novo símbolo inicial S' e a regra $S' \rightarrow S\#$. (Necessário se e a gramática não tiver esta configuração).

o símbolo $\#$ significa: marca de fim de buffer de entrada (ou ainda: fim de análise, fim de ficheiro etc...). Assim o buffer de entrada deverá considerar como último símbolo o símbolo $\#$

construção dum autómato com pilha a partir duma CFG

por exemplo:

$$E ::= E + E$$

dada a gramática $E ::= E * E$

$$E ::= x$$

$$\boxed{S' ::= E \#}$$

nova versão :

$$E ::= E + E$$

$$E ::= E * E$$

$$E ::= x$$

conjunto dos itens:

$$S' ::= . E \# \quad E ::= . E + E \quad E ::= . E * E \quad E ::= . x$$

$$S' ::= E . \# \quad E ::= E . + E \quad E ::= E . * E \quad E ::= x .$$

$$S' ::= E \# . \quad E ::= E + . E \quad E ::= E * . E$$

$$E ::= E + E . \quad E ::= E * E .$$

construção dum autómato com pilha a partir duma CFG

tendo uma gramática $G = (T, N, S, P)$, o autómato com pilha define-se então da seguinte forma:

Alfabeto: T .

Estados: Conjunto dos itens associado à gramática G

Transições: $T \cup \{\epsilon\}$.

Estado inicial: $S' \rightarrow .S\#$

significa: preparada para iniciar o reconhecimento duma frase gerada por S

Estado final (único): $S' \rightarrow S\#$.

significa: uma frase gerada por S foi reconhecido e a marca de fim de análise foi atingida

3 tipos de transições:

- **expansão:** para cada $Y ::= \alpha$ de P (lembrete: $\alpha \in (N \cup T)^*$).
 $[X \rightarrow \beta. Y\gamma] \xrightarrow{\epsilon} [X \rightarrow \beta. Y\gamma][Y \rightarrow . \alpha]$
- **leitura** de $a \in T$.
 $[X \rightarrow \beta. a\gamma] \xrightarrow{a} [X \rightarrow \beta a. \gamma]$
- **redução** da regra $Y ::= \alpha$.
 $[X \rightarrow \beta. Y\gamma][Y \rightarrow \alpha.] \xrightarrow{\epsilon} [X \rightarrow \beta Y. \gamma]$

reconhecimento de $aabb\#$ pela gramática

S	$::=$	aSb
S	$::=$	ϵ

	<i>start</i>	$[S' \rightarrow \cdot S\#]$
<i>expansão</i>	$S ::= aSb$	$[S' \rightarrow \cdot S\#][S \rightarrow \cdot aSb]$
	<i>leitura a</i>	$[S' \rightarrow \cdot S\#][S \rightarrow a \cdot Sb]$
<i>expansão</i>	$S ::= aSb$	$[S' \rightarrow \cdot S\#][S \rightarrow a \cdot Sb][S \rightarrow \cdot aSb]$
	<i>leitura a</i>	$[S' \rightarrow \cdot S\#][S \rightarrow a \cdot Sb][S \rightarrow a \cdot Sb]$
<i>expansão</i>	$S ::= \epsilon$	$[S' \rightarrow \cdot S\#][S \rightarrow a \cdot Sb][S \rightarrow a \cdot Sb][S \rightarrow \cdot]$
	<i>redução</i>	$[S' \rightarrow \cdot S\#][S \rightarrow a \cdot Sb][S \rightarrow aS \cdot b]$
	<i>leitura b</i>	$[S' \rightarrow \cdot S\#][S \rightarrow a \cdot Sb][S \rightarrow aSb \cdot]$
	<i>redução</i>	$[S' \rightarrow \cdot S\#][S \rightarrow aS \cdot b]$
	<i>leitura b</i>	$[S' \rightarrow \cdot S\#][S \rightarrow aSb \cdot]$
	<i>redução</i>	$[S' \rightarrow S \cdot \#]$
	<i>leitura #</i>	$[S' \rightarrow S\# \cdot]$

apresentar um reconhecimento de $id + id * id\#$ com base no autômato com pilha subjacente à gramática seguinte:

$$\begin{aligned} E &::= T \mid E + T \\ T &::= F \mid T * F \\ F &::= id \mid (E) \end{aligned}$$

este processo consegue lidar com todas as gramáticas livres de contexto, mas

os autómatos gerados por estes métodos são enormes (na prática, como vimos no exemplo anterior, nem os desenhamos)

e são igualmente não-deterministas

em termos conceptuais, um reconhecimento sintáctico baseado neste processo é completo mas, em prática, é particularmente ineficiente

que soluções? mais uma vez...visar um subconjunto das gramáticas livres de contexto representativo sobre o qual outros tipos de algoritmos de reconhecimento mais eficientes possam ser definidos: as gramáticas LL e LR

análise descendente

problema: tendo em conta o não determinismo do autómato com pilha, como escolher uma regra de expansão?

uma solução possível: olhar para o buffer de entrada **um ou mais** caracteres **a frente**.

contexto: tomemos por exemplo as produções $Y ::= \alpha_1 | \alpha_2 | \dots | \alpha_n$. Se cada α_i começar de forma distinta dos outros α_j ($i \neq j$) então ao ler um carácter a sabemos imediatamente que regras escolher: a regra $Y ::= \alpha_k$ com $\alpha_k \rightarrow^* a \dots$ (a na primeira posição).

$$\begin{array}{ll}
 1 & S ::= A d S \\
 2 & S ::= b \\
 3 & A ::= a A b \\
 4 & A ::= c
 \end{array}$$

vamos analisar a situação mais cuidadosamente

imaginemos que estamos numa situação em que pretendemos derivar algo de S

- se a entrada começar por a ou c então a única hipótese é escolher a regra 1
- se a entrada começar por b então a única possibilidade é escolher a regra 2
- qualquer outra entrada não pode reconhecida por S

relativamente a A

- se a entrada for a então podemos escolher a 3^{ra} regra
- se a entrada for c então podemos escolher a 4^{ta} regra
- qualquer outra entrada não pode reconhecida por A

- 1 $S ::= A d S$
- 2 $S ::= b$
- 3 $A ::= a A b$
- 4 $A ::= c$

esta gramática tem então a particularidade de permitir escolher de forma certa a produção para a derivação, e isto com base na **antevisão** de um só carácter na entrada

ou seja, com tal gramática a análise sintáctica é facilitada: em cada situação sabemos exactamente o que fazer entre escolher leituras, expansões de determinadas regras ou reduções

- $$\begin{array}{ll}
 1 & S ::= A d S \\
 2 & S ::= b \\
 3 & A ::= a A b \\
 4 & A ::= c
 \end{array}$$

um exemplo: derivação da palavra $acbdb\#$

	Estado	antevisão da entrada	decisão e acção
0	$[\cdot S\#]$		Início
1	$[\cdot S\#]$	a	expansão regra 1
2	$[\cdot S\#][\cdot A d S]$	a	expansão regra 3
3	$[\cdot S\#][\cdot A d S][\cdot a A b]$	a	leitura a
4	$[\cdot S\#][\cdot A d S][a \cdot A b]$	c	expansão regra 4
5	$[\cdot S\#][\cdot A d S][a \cdot A b][\cdot c]$	c	leitura c
6	$[\cdot S\#][\cdot A d S][a \cdot A b][c \cdot]$	b	redução regra 4
7	$[\cdot S\#][\cdot A d S][aA \cdot b]$	b	leitura b
8	$[\cdot S\#][\cdot A d S][a A b \cdot]$	d	redução da regra 3
9	$[\cdot S\#][A \cdot d S]$	d	leitura d
10	$[\cdot S\#][A d \cdot S]$	b	expansão regra 2
11	$[\cdot S\#][A d \cdot S][\cdot b]$	b	leitura b
12	$[\cdot S\#][A d \cdot S][b \cdot]$	$\#$	redução da regra 2
13	$[\cdot S\#][A d S \cdot]$	$\#$	redução da regra 1
14	$[S \cdot \#]$	$\#$	leitura $\#$
15	$[S\# \cdot]$	ϵ	sucesso

esta análise corresponde a derivação esquerda seguinte:

$$S' \Rightarrow S\# \Rightarrow AdS\# \Rightarrow aAbdS\# \Rightarrow acbdS\# \Rightarrow acbdb\#$$

E se juntarmos mais uma regra?

1	$S ::= A d S$
2	$S ::= b$
3	$A ::= a A b$
4	$A ::= c$
5	$A ::= \epsilon$

o facto de A poder derivar o ϵ obriga aqui a um cuidado especial: temos de olhar para além de A para saber que regra escolher

relativamente a S

- Se a entrada começar por a , c ou d então a única hipótese é escolher a regra 1.
- Se a entrada começar por b então a única possibilidade é escolher a regra 2
- Qualquer outra entrada não pode reconhecida por S .

relativamente a A

- Se a entrada for a então podemos escolher a 3^{ra} regra.
- Se a entrada for c então podemos escolher a 4^{ta} regra.
- Se a entrada for b ou d então podemos escolher a 5^{ta} regra.

Para poder ter uma análise sintáctica nestes moldes é preciso conseguir determinar para toda a palavra m de $(N \cup T)^*$ (em particular não-terminais)

se esta pode **derivar a palavra vazia**,

que conjunto de caracteres (letras de T) **apareçam em primeira posição** das palavras geradas por m e finalmente

que conjunto de caracteres terminais **apareçam em primeira posição a seguir** a qualquer palavra gerada por m .

estes dados permitirão construir um analisador (descendente) eficiente e determinístico.

ideia : proceder por expansões sucessivas do não terminal o mais a esquerda (construímos assim uma derivação esquerda) partindo de S e servindo-se de uma **tabela** que indica, para um não terminal X por expandir e os k primeiros caracteres da entrada, qual a expansão $X \rightarrow \beta$ por realizar (em inglês falamos de *top-down parsing*)

supomos $k = 1$ no resto da apresentação e notemos $T(X, c)$ esta tabela

na prática, tendo em conta o símbolo terminal $\#$, a tabela indica também as expansões X quando o fim da entrada foi atingido

utilizamos uma pilha que é uma palavra de $(N \cup T)^*$; inicialmente a pilha está reduzida ao símbolo inicial

em cada instante, examinamos o topo da pilha e o primeiro caracter c da entrada

- se a pilha encontra-se vazia, paramos ; há sucesso se e só se $c = \#$
- se o topo da pilha é um terminal a , então a deve ser igual a c , tiramos a da pilha (*pop*) e consumimos c da entrada ; senão a análise falha
- se o topo da pilha é um não-terminal X , então substituímos X pela palavra $\beta = T(X, c)$ no topo da pilha: i.e. *pop* seguido de, quando apropriado, *push* para cada letra de β , começando pela última ; senão, a análise falha

consideremos mais uma gramática para as expressões aritméticas
 (assumimos implicitamente $S \rightarrow E\#$)
 e a tabela de expansão seguinte

$$\begin{array}{l}
 E \rightarrow T E' \\
 E' \rightarrow + T E' \\
 \quad | \epsilon \\
 T \rightarrow F T' \\
 T' \rightarrow * F T' \\
 \quad | \epsilon \\
 F \rightarrow (E) \\
 \quad | \text{int}
 \end{array}$$

	+	*	()	int	#
<i>E</i>			<i>TE'</i>		<i>TE'</i>	
<i>E'</i>	+ <i>TE'</i>			ϵ		ϵ
<i>T</i>			<i>FT'</i>		<i>FT'</i>	
<i>T'</i>	ϵ	* <i>FT'</i>		ϵ		ϵ
<i>F</i>			(<i>E</i>)		int	

(veremos mais adiante como construir tal tabela)

ilustremos a análise descendente da palavra

int + int * int

	+	*	()	int	#
<i>E</i>			<i>TE'</i>		<i>TE'</i>	
<i>E'</i>	<i>+TE'</i>			ϵ		ϵ
<i>T</i>			<i>FT'</i>		<i>FT'</i>	
<i>T'</i>	ϵ	<i>*FT'</i>		ϵ		ϵ
<i>F</i>			<i>(E)</i>		int	

pilha	entrada
<i>E</i>	int + int * int #
<i>E'T</i>	int + int * int #
<i>E'T'F</i>	int + int * int #
<i>E'T'int</i>	int + int * int #
<i>E'T'</i>	+int * int #
<i>E'</i>	+int * int #
<i>E'T+</i>	+int * int #
<i>E'T</i>	int * int #
<i>E'T'F</i>	int * int #
<i>E'T'int</i>	int * int #
<i>E'T'</i>	*int #
<i>E'T'F*</i>	*int #
<i>E'T'F</i>	int #
<i>E'T'int</i>	int #
<i>E'T'</i>	#
<i>E'</i>	#
ϵ	#

programa-se um analisador descendente com alguma facilidade, introduzindo uma função para cada símbolo não-terminal da gramática

cada função examina a entrada e, consoante o caso, consome esta ou invoca recursivamente as funções correspondentes a outros não-terminais, conforme a tabela de expansão

programação de um analisador descendente

façamos a escolha de uma programação puramente aplicativa, onde a entrada é uma lista de lexemas do tipo

```
type token = Tplus | Tmult | Tleft | Tright | Tint | Teof
```

vamos assim escrever cinco funções que « consumam » a lista de entradas

```
val e : token list -> token list  
val e': token list -> token list  
val t : token list -> token list  
val t': token list -> token list  
val f : token list -> token list
```

e o reconhecimento de uma entrada pode assim fazer-se por

```
let recognize l =  
  e l = [Teof]
```

programação de um analisador descendente

as funções operam por *pattern matching* sobre a entrada em conformidade com a tabela

	+	*	()	int	#
<i>E</i>			<i>TE'</i>		<i>TE'</i>	

```
let rec e = function
  | (Tleft | Tint) :: _ as m -> e' (t m)
  | _ -> error ()
```

	+	*	()	int	#
<i>E'</i>	+ <i>TE'</i>			ε		ε

```
and e' = function
  | Tplus :: m -> e' (t m)
  | (Tright | Teof) :: _ as m -> m
  | _ -> error ()
```

programação de um analisador descendente

	+	*	()	int	#
<i>T</i>			<i>FT'</i>		<i>FT'</i>	

```
and t = function
```

```
  | (Tleft | Tint) :: _ as m -> t' (f m)
  | _ -> error ()
```

	+	*	()	int	#
<i>T'</i>	ϵ	<i>*FT'</i>		ϵ		ϵ

```
and t' = function
```

```
  | (Tplus | Tright | Teof) :: _ as m -> m
  | Tmult :: m -> t' (f m)
  | _ -> error ()
```

programação de um analisador descendente

	+	*	()	int	#
<i>F</i>			(<i>E</i>)		int	

```
and f = function
| Tint :: m -> m
| Tleft :: m -> begin match e m with
  | Tright :: m -> m
  | _ -> error ()
end
| _ -> error ()
```

algumas notas

- a tabela de expansão não está explícita : está embutida no código de cada funções
- a pilha também não : ela é materializada pela pilha de chamadas
- em alternativa à implementação cá exposta, poderíamos tê-las tornadas explícitas no programa
- poderíamos ter preferido uma programação mais imperativa

```
val next_token : unit -> token
```


falta-nos abordar uma questão importante : como construir a tabela de expansão ?

a ideia é simples : para decidir se realizamos uma expansão $X \rightarrow \beta$ quando o primeiro caracter da entrada é c , vamos procurar determinar se c faz parte dos **primeiros** caracteres das palavras reconhecidos por β

como já vimos, surge uma dificuldade neste processo quando existam produções de tipo $X \rightarrow \epsilon$,
o primeiro caracter que pode ser gerado neste ponto pertence ao conjunto dos caracteres que podem **seguir** X

para determinar os primeiros caracteres e os caracteres seguintes é preciso também determinar se uma palavra se pode derivar em ϵ

Definição (null)

Seja $\alpha \in (T \cup N)^*$. $\text{null}(\alpha)$ é verdade se e só se podemos derivar ϵ partindo de α i.e. $\alpha \rightarrow^* \epsilon$.

Definição (first)

Seja $\alpha \in (T \cup N)^*$. $\text{first}(\alpha)$ é o conjunto de todos os primeiros terminais das palavras derivadas de α , i.e. $\{a \in T \mid \exists w. \alpha \rightarrow^* aw\}$.

Definição (follow)

Seja $X \in N$. $\text{follow}(X)$ é o conjunto de todos os terminais que podem aparecer a seguir a X numa derivação, i.e. $\{a \in T \mid \exists u, w. S \rightarrow^* uXaw\}$.

para calcular $\text{null}(\alpha)$, basta determinar $\text{null}(X)$ para todo o $X \in N$

$\text{null}(X)$ é verdade se e só se

- existe uma produção $X \rightarrow \epsilon$,
- ou existe uma produção $X \rightarrow Y_1 \dots Y_m$ onde $\text{null}(Y_i)$ para todo i

(note que é falso em qualquer outra situação, em particular quando há sempre terminais na parte direita das produções)

problema : trata-se de um conjunto de equação mutuamente recursivas

dito de outra forma,

se $N = \{X_1, \dots, X_n\}$ e se $\vec{V} = (\text{null}(X_1), \dots, \text{null}(X_n))$,

vamos procurar a **menor solução** de uma equação da forma

$$\vec{V} = F(\vec{V})$$

Theorema (existência de um menor ponto fixo (Tarski))

Seja A um conjunto finito equipado de uma relação de ordem \leq e de um menor elemento ε . Toda a função $f : A \rightarrow A$ crescente, i.e. tal que $\forall x, y. x \leq y \Rightarrow f(x) \leq f(y)$, admite um menor ponto fixo.

demonstração :

como ε é o menor elemento, temos $\varepsilon \leq f(\varepsilon)$

sendo f crescente, temos $f^k(\varepsilon) \leq f^{k+1}(\varepsilon)$ para todo k

sendo A finito, existe assim um menor k_0 tal que $f^{k_0}(\varepsilon) = f^{k_0+1}(\varepsilon)$

$a_0 = f^{k_0}(\varepsilon)$ é assim um menor ponto fixo de f

seja b um outro ponto fixo de f

temos $\varepsilon \leq b$ e então $f^k(\varepsilon) \leq f^k(b)$ para todo k

em particular $a_0 = f^{k_0}(\varepsilon) \leq f^{k_0}(b) = b$

a_0 é assim o menor ponto fixo de f



o teorema de Tarski apresenta condições **suficientes** para a existência do menor ponto fixo, mas não apresenta condições **necessárias**

no caso do cálculo de null, temos $A = \text{Bool} \times \dots \times \text{Bool}$ com $\text{Bool} = \{\text{false}, \text{true}\}$

podemos equipar Bool da ordem $\text{false} \leq \text{true}$ e A da ordem *produto*

$$(x_1, \dots, x_n) \leq (y_1, \dots, y_n) \quad \text{se e só se} \quad \forall i. x_i \leq y_i$$

o teorema aplica-se então tomando

$$\varepsilon = (\text{false}, \dots, \text{false})$$

porque a função calculando $\text{null}(X)$ a partir dos $\text{null}(X_i)$ é crescente

para calcular os $\text{null}(X_i)$, podemos começar por

$$\text{null}(X_1) = \text{false}, \dots, \text{null}(X_n) = \text{false}$$

e aplicar as equações até obter o ponto fixo *i.e.* até os valores de $\text{null}(X_i)$ não se alterarem mais

$$\begin{array}{l}
 E \rightarrow T E' \\
 E' \rightarrow + T E' \\
 \quad | \quad \epsilon \\
 T \rightarrow F T' \\
 T' \rightarrow * F T' \\
 \quad | \quad \epsilon \\
 F \rightarrow (E) \\
 \quad | \quad \text{int}
 \end{array}$$

$$\begin{aligned}
 \text{null}(E) &= \text{null}(T) \wedge \text{null}(E') \\
 \text{null}(E') &= \text{null}(+ T E') \vee \text{null}(\epsilon) \\
 \text{null}(T) &= \text{null}(F) \wedge \text{null}(T') \\
 \text{null}(T') &= \text{null}(* F T') \vee \text{null}(\epsilon) \\
 \text{null}(F) &= \text{null}((E)) \vee \text{null}(\text{int})
 \end{aligned}$$

E	E'	T	T'	F
false	false	false	false	false
false	true	false	true	false
false	true	false	true	false

porque procuramos o **menor** ponto fixo ?

- ⇒ por recorrência sobre o número de passos do cálculo anterior, mostra-se que se $\text{null}(X) = \text{true}$ então $X \rightarrow^* \epsilon$
- ⇐ por recorrência sobre o número de passos da derivação $X \rightarrow^* \epsilon$, mostra-se que $\text{null}(X) = \text{true}$ pelo cálculo anterior

de forma semelhante, as equações que definam first são mutuamente recursivas

$$\text{first}(X) = \bigcup_{X \rightarrow \beta} \text{first}(\beta)$$

e

$$\text{first}(\epsilon) = \emptyset$$

$$\text{first}(a\beta) = \{a\}$$

$$\text{first}(X\beta) = \text{first}(X), \quad \text{se } \neg \text{null}(X)$$

$$\text{first}(X\beta) = \text{first}(X) \cup \text{first}(\beta), \quad \text{se } \text{null}(X)$$

mais uma vez, procedemos por cálculo de ponto fixo sobre o produto cartesiano $A = \mathcal{P}(T) \times \cdots \times \mathcal{P}(T)$ equipado da ordem produto sobre \subseteq e com $\epsilon = (\emptyset, \dots, \emptyset)$

null

E	E'	T	T'	F
false	true	false	true	false

$E \rightarrow T E'$
 $E' \rightarrow + T E'$
 $\quad \mid \epsilon$
 $T \rightarrow F T'$
 $T' \rightarrow * F T'$
 $\quad \mid \epsilon$
 $F \rightarrow (E)$
 $\quad \mid \text{int}$

$first(E) = first(T E') = first(T)$
 $first(E') = first(+ T E') \cup first(\epsilon) = first(+) \cup \emptyset = \{+\}$
 $first(T) = first(F T') = first(F)$
 $first(T') = first(* F T') \cup first(\epsilon) = first(*) \cup \emptyset = \{*\}$
 $first(F) = first((E)) \cup first(\text{int}) = \{(, \text{int})\}$

first

E	E'	T	T'	F
\emptyset	\emptyset	\emptyset	\emptyset	\emptyset
\emptyset	$\{+\}$	\emptyset	$\{*\}$	$\{(, \text{int})\}$
\emptyset	$\{+\}$	$\{(, \text{int})\}$	$\{*\}$	$\{(, \text{int})\}$
$\{(, \text{int})\}$	$\{+\}$	$\{(, \text{int})\}$	$\{*\}$	$\{(, \text{int})\}$
$\{(, \text{int})\}$	$\{+\}$	$\{(, \text{int})\}$	$\{*\}$	$\{(, \text{int})\}$

mais uma vez, as equações definindo follow são mutuamente recursivas

$$\text{follow}(X) = \bigcup_{Y \rightarrow \alpha X \beta} \text{first}(\beta) \cup \bigcup_{Y \rightarrow \alpha X \beta, \text{null}(\beta)} \text{follow}(Y)$$

procedemos por cálculo de ponto fixo, sobre o mesmo domínio usado para o cálculo de first

nota : é necessário introduzir # no calculo dos caracteres seguintes do símbolo inicial

(podemos fazê-lo directamente ou então juntando a regra $S' \rightarrow S\#$)

null

E	E'	T	T'	F
false	true	false	true	false

first

E	E'	T	T'	F
{(, int}	{+}	{(, int}	{*}	{(, int}

$E \rightarrow T E'$
 $E' \rightarrow + T E'$
 $\quad \mid \epsilon$
 $T \rightarrow F T'$
 $T' \rightarrow * F T'$
 $\quad \mid \epsilon$
 $F \rightarrow (E)$
 $\quad \mid \text{int}$

$\text{follow}(E) = \{\#\} \cup \text{first}() = \{\#,)\}$
 $\text{follow}(E') = \text{first}(\epsilon) \cup \text{first}(\epsilon) \cup \text{follow}(E) \cup \text{follow}(E')$
 $\quad = \text{follow}(E) \cup \text{follow}(E')$
 $\text{follow}(T) = \text{first}(E') \cup \text{follow}(E) \cup \text{follow}(E')$
 $\quad = \{+\} \cup \text{follow}(E) \cup \text{follow}(E')$
 $\text{follow}(T') = \text{first}(\epsilon) \cup \text{follow}(T) \cup \text{follow}(T')$
 $\quad = \text{follow}(T) \cup \text{follow}(T')$
 $\text{follow}(F) = \text{first}(T') \cup \text{follow}(T) \cup \text{follow}(T')$
 $\quad = \{*\} \cup \text{follow}(T) \cup \text{follow}(T')$

follow

E	E'	T	T'	F
{#}	\emptyset	\emptyset	\emptyset	\emptyset
{#,)}	{#}	{+, #}	\emptyset	{*}
{#,)}	{#,)}	{+, #,)}	{+, #}	{*, +, #}
{#,)}	{#,)}	{+, #,)}	{+, #,)}	{*, +, #,)}
{#,)}	{#,)}	{+, #,)}	{+, #,)}	{*, +, #,)}

com base nestes três conceitos, construímos a tabela de expansão $T(X, a)$ da forma seguinte

para cada produção $X \rightarrow \beta$,

- consideramos $T(X, a) = \beta$ para todo o $a \in \text{first}(\beta)$
- se $\text{null}(\beta)$, consideramos também $T(X, a) = \beta$ para todo o $a \in \text{follow}(X)$

$E \rightarrow TE'$	first	E	E'	T	T'	F
$E' \rightarrow +TE'$		{(, int)}	{+}	{(, int)}	{*}	{(, int)}
ϵ						
$T \rightarrow FT'$						
$T' \rightarrow *FT'$	follow	E	E'	T	T'	F
ϵ		{#,)}	{#,)}	{+, #,)}	{+, #,)}	{*, +, #,)}
$F \rightarrow (E)$						
int						

	+	*	()	int	#
E			TE'	TE'	
E'	$+TE'$			ϵ	ϵ
T			FT'	FT'	
T'	ϵ	$*FT'$		ϵ	ϵ
F			(E)	int	

Definição (gramática LL(1))

Uma gramática é dita LL(1) se, na sua tabela de expansão, há no máximo uma produção em cada célula.

LL significa « **L**eft to right scanning, **L**eftmost derivation »

nem todas as gramáticas livres de contexto são LL(1), como tal é frequente ser necessário transformar uma gramática para obter uma gramática equivalente que seja LL(1)

vejamos alguns **critérios suficientes** para saltar fora do âmbito das gramáticas LL(1)

como caracterizar gramáticas $LL(1)$?

já sabemos que são gramáticas que permitam uma análise descendente eficiente e simples

sabemos igualmente $LL(1) \iff$ uma entrada no máximo em cada célula da tabela de expansão

mas não haverá uma forma informal ou intuitiva de caracterizar tais gramáticas? Sim:

- é possível prever, sempre, que regra escolher só com base nos k primeiros caracteres.
- ou seja, para $k = 1$ e um não-terminal A tal que $A \rightarrow \alpha_1|\alpha_2|\dots|\alpha_n$ e para i e j distintos :
 1. no máximo um dos α_i é tal que $\text{null}(\alpha_i) = \text{true}$
 2. $\text{first}(\alpha_i) \cap \text{first}(\alpha_j) = \emptyset$
 3. $\exists i. \text{null}(\alpha_i) = \text{true} \implies \forall j. \text{first}(\alpha_j) \cap \text{follow}(A) = \emptyset$

construir uma gramática LL(1)

uma gramática **recursiva esquerda (direta)**, i.e. contendo uma produção da forma

$$X \rightarrow X\alpha$$

nunca será LL(1)

torná-la LL(1) passará sempre por remover a recursividade esquerda (direta ou indireta)

de igual forma é necessário factorizar as produções que começam pelo mesmo prefixo (**fatorização esquerda**)

suficiente?

em geral, não

como proceder então?

completar a transformação da gramática consoante diagnósticos feitos às anomalias presentes na tabela de transição (soluções *Ad-Hoc*)

por exemplo a redução da gramática por remoção de não-terminais.

o que é a recursividade esquerda?

uma gramática é recursiva esquerda se existir uma derivação tal $X \rightarrow^+ X\alpha$

uma gramática recursiva esquerda não pode ser $LL(1)$

se pretendemos obter uma gramática $LL(1)$ então é preciso eliminar a recursividade esquerda.

situação:

$$\begin{array}{l}
 A \rightarrow A \alpha_1 \mid A \alpha_2 \mid \cdots \mid A \alpha_n \\
 A \rightarrow \beta_1 \mid \beta_2 \mid \cdots \mid \beta_m \quad \text{os } \beta_i \text{ não começam por } A
 \end{array}$$

transformar em

$$\begin{array}{l}
 A \rightarrow \beta_1 A' \mid \beta_2 A' \mid \cdots \mid \beta_m A' \\
 A' \rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \cdots \mid \alpha_n A' \\
 A' \rightarrow \epsilon
 \end{array}$$

caso particular: caso surgem regras do tipo $A \rightarrow A$, estas regras podem ser removidas

consideremos esta gramática:

$$S \rightarrow A a$$

$$S \rightarrow b$$

$$A \rightarrow A c$$

$$A \rightarrow S d$$

$$S \rightarrow^+ S d a$$

como proceder para remover esta recursividade esquerda?

remoção da recursividade esquerda indireta

organizar os não-terminais de 1 a n (no caso de $|N| = n$): A_1, A_2, \dots, A_n

mas: **a ordem escolhida tem a sua importância!** Preferir atribuir os índices menores aos não-terminais do fundo da gramática e o índice final ao símbolo inicial

processo em n etapas:

etapa 1: Eliminar a recursividade directa de A_1

etapa i :

- enquanto existir uma produção $A_i \rightarrow A_j\alpha$ (com $i > j$) olhar para todas as produções tais que $A_j \rightarrow \beta$ e substituir $A_i \rightarrow A_j\alpha$ por $A_i \rightarrow \beta\alpha$
- retirar então a recursividade esquerda directa (quando $i = j$)
- no final, se β começa por um A_k então necessariamente $k > i$.

o processo iterativo anterior só acaba quando temos produções $A_i \rightarrow A_j\alpha$ tal que $i < j$

ordem: S, A (contrária à sugestão feita...)

começemos por S : $S \rightarrow A a \quad S \rightarrow b$

$S \rightarrow A a$
 $S \rightarrow b$
 $A \rightarrow A c$
 $A \rightarrow S d$

para A , em $A \rightarrow S d$ substituímos S por $A a \quad \epsilon \quad b$ logo:
 $A \rightarrow A c \quad A \rightarrow A a d \quad A \rightarrow b d$

podemos então eliminar a recursividade esquerda directa:

	$A' \rightarrow c A'$
$A \rightarrow b d A'$	$A' \rightarrow a d A'$
	$A' \rightarrow \epsilon$

conflito por prefixo comum

$$\begin{aligned}
 A &\rightarrow \underline{\alpha}\beta_1 \\
 A &\rightarrow \underline{\alpha}\beta_2 \\
 &\vdots \\
 A &\rightarrow \underline{\alpha}\beta_n
 \end{aligned}$$

neste caso transformar tais regras em:

$$\begin{aligned}
 A &\rightarrow \underline{\alpha}A' \\
 A' &\rightarrow \beta_1 \\
 A' &\rightarrow \beta_2 \\
 &\vdots \\
 A' &\rightarrow \beta_n
 \end{aligned}$$

remoção prévia de Não-Terminais

pode acontecer que duas produções de um não terminal Y tenham conjuntos first não disjuntos sem no entanto apresentarem explicitamente prefixo comum

pode assim ser interessante tentar explicitar estes prefixos e proceder à transformação anterior

$$Y \rightarrow \beta \underline{X} \gamma \quad (X \neq Y)$$

neste caso substituir tal regra pelas regras:

$$Y \rightarrow \beta \underline{\alpha} \gamma$$

para todo α tal que $(X \rightarrow \alpha) \in P$.

proceder em seguida à remoção de factorização esquerda caso necessário

conclusão

os analisadores LL(1) são relativamente simples de implementar
(ver aulas práticas)

mas obrigam à escrita de gramáticas que podem não serem *naturais*

iremos ver na próxima aula um outro tipo de solução

- aulas práticas
 - analisadores LL(1)

- aulas teóricas
 - análise sintáctica (segunda parte)

estes acetatos resultam essencialmente de uma adaptação do material pedagógico gentilmente cedido pelo Jean-Christophe Filliâtre ([link1](#), [link2](#))

adicionalmente poderá consultar as obras seguintes

- **Modern Compilers: Principles, Techniques, and Tools**, Alfred V. Aho, Monica S. Lam, Ravi Sethi et Jeffrey D. Ullman
- **Types and Programming Languages**, Benjamin C. Pierce
- **Modern Compiler Implementation**, Andrew W. Appel (3 versões: ML, C, Java)

