

Aula Teórica 6 (guião)

Semana de 20 a 24 de Outubro de 2025

José Carlos Ramalho

Sinopsis:

- Parsers Top-Down: o Recursivo Descendente;
 - A condição LL(1).
 - Imagens e algum conteúdo retirados da sebenta "Processamento de Linguagens: Reconhecedores Sintáticos" de José João Almeida e José Bernardo Barros, editada em Abril de 2022.
-

Gramáticas

- Gramáticas Regulares;
 - Gramáticas Independentes de Contexto;
 - Gramáticas Dependentes de Contexto;
 - Gramáticas Livres.
-

Gramáticas: especificação

Uma gramática G é um tuplo (N, T, S, P) , onde:

- N - conjunto de símbolos não terminais;
 - T - conjunto de símbolos terminais;
 - S - símbolo inicial;
 - P - as produções.
-

Exemplificar: uma turma de alunos

```
Turma PL2025 .  
Alunos{  
    ("A23876", "Ana Maria")  
    ("A78654", "Paulo Azevedo")  
    ...  
}
```

Qual seria a gramática?

```
T = {}  
S =
```

N =
P =

(Dar exemplos)

LISP + Markdown

```
(doc
  (tit "Primeiro exemplo")
  (subtit "Aula 6: 2025-03-14")
  "Este é um primeiro exemplo."
)
```

Expressões S

```
(+ 5 4)
(+ (- 7 2) (+ 12 19))
```

Revisitando os parentesis

```
()
((()))()
```

Expressões aritméticas básicas

```
4 + 15
7 - 12 * 3
7 + 8 * 2 / 3
```

Voltando às listas

Exercício: a linguagem das listas

Exemplos:

```
[]
[2]
```

```
[ 2, 4, 5]
[ 2, 4, [ 5, 7, 9], 6]
```

Símbolos terminais: $T = \{'[', ']', \text{num}\}$

Produções:

```
Lista --> '[' ']'
        | '[' Conteudo ']'

Conteudo --> num
          | num ',' Conteudo
```

Analizador Léxico

```
# listas_analex.py
# 2023-03-21 by jcr
# -----
import ply.lex as lex

tokens = ('NUM', 'PA', 'PF', 'VIRG')

t_NUM = r'[+\-]?\\d+'
t_PA = r'\\['
t_PF = r'\\]'
t_VIRG = r','

def t_newline(t):
    r'\\n+'
    t.lexer.lineno += len(t.value)

t_ignore = '\\t '

def t_error(t):
    print('Carácter desconhecido: ', t.value[0], 'Linha: ',
          t.lexer.lineno)
    t.lexer.skip(1)

lexer = lex.lex()
```

Programa exemplo

```
# listas_program.py
# 2023-03-21 by jcr
# -----
```

```
linha = input("Introduza uma lista: ")
rec_Parser(linha)
```

Reconhcedores (Parsers) Top-Down: o recursivo descendente

- Vamos escrever uma função para cada símbolo terminal (T) ou não terminal (N).
- Para os terminais basta verificar se o símbolo que é preciso reconhecer é o que está a seguir na string de entrada:

função $recT (t : T) : tipoValT$

$$\left\{ \begin{array}{l} \text{se } t = simbolo_{seg} \longrightarrow \left\{ \begin{array}{l} val \leftarrow simbolo_{seg}.valor \\ ABORT \end{array} \right. \\ \text{senão} \longrightarrow \\ \text{return } val \end{array} \right.$$

Para os não terminais a situação é um pouco mais complexa.

Para produções da forma: $A \rightarrow x_1 x_2 \dots x_n \mid \dots \mid y_1 y_2 \dots y_n$ é preciso decidir qual das regras se vai seguir.

No exemplo das listas, quais as produções com que devemos ter cuidado?

```
p1: Lista --> '[' ']'
p2:         | '[' Conteudo ']'

p3: Conteudo --> num
p4:           | num ',' Conteudo
```

No decurso da análise do parser, este terá de decidir entre **p1** e **p2** quando estiver a tentar derivar o símbolo **Lista** e entre **p3** e **p4** quando estiver a tentar derivar o símbolo **Conteúdo**.

Para se poder escolher a produção a seguir é preciso calcular o conjunto de símbolos terminais que poderão dar início à derivação especificada nessa produção. A estes conjuntos chamamos **lookaheads**.

função $recA () : tipoValA$

$$\left\{ \begin{array}{l} \text{se } simbolo_{seg} \in lookahead(A \rightarrow x_1x_2...x_n) \longrightarrow \\ \quad \left\{ \begin{array}{l} val_1 \leftarrow recx_1; \\ val_2 \leftarrow recx_2; \\ \dots \\ val_n \leftarrow recx_n; \\ valA \leftarrow f(val_1, val_2, \dots val_n) \end{array} \right. \\ simbolo_{seg} \in lookahead(A \rightarrow \dots) \longrightarrow \\ \quad \{ \dots \\ simbolo_{seg} \in lookahead(A \rightarrow y_1y_2...y_m) \longrightarrow \\ \quad \left\{ \begin{array}{l} val_1 \leftarrow recy_1; \\ val_2 \leftarrow recy_2; \\ \dots \\ val_m \leftarrow recy_m; \\ valA \leftarrow g(val_1, val_2, \dots val_m) \end{array} \right. \\ \text{senão} \longrightarrow \\ \quad \{ [erro] \\ \text{return } valA \end{array} \right.$$

Cálculo de Lookaheads

O Cálculo do **lookahead** faz-se com a ajuda de duas funções auxiliares: **First** e **Follow**.

First

Esta função calcula o conjunto de símbolos terminais que são inícios válidos de uma determinada sequência de símbolos. Esta sequência pode iniciar-se por um símbolo terminal ou por um não terminal.

Para um sequência iniciada por um não terminal calcula-se o **FirstN** correspondente a esse símbolo:

função $first_nt (G : gramatica, A : NT) : T - set$

$$\left\{ \begin{array}{l} \text{para } (A \rightarrow rhs) \in Producoes(G) \text{ fazer} \\ \quad \left\{ \begin{array}{l} fir \leftarrow fir \cup first_str(rhs) \end{array} \right. \\ \text{return } fir \end{array} \right.$$

Para os restantes casos aplica-se a função seguinte:

função $first_str (G : gramatica : s : (T|NT) - seq) : T - set$

$$\left\{ \begin{array}{l} \text{se } s = \langle \rangle \longrightarrow \left\{ \begin{array}{l} fir \leftarrow \{ \} \end{array} \right. \\ s = \langle h.t \rangle \wedge h \in T \longrightarrow \left\{ \begin{array}{l} fir \leftarrow h \end{array} \right. \\ s = \langle h.t \rangle \wedge h \in NT_{anulaveis} \longrightarrow \left\{ \begin{array}{l} fir \leftarrow first_nt(h) \cup first_str(t) \end{array} \right. \\ s = \langle h.t \rangle \wedge h \notin NT_{anulaveis} \longrightarrow \left\{ \begin{array}{l} fir \leftarrow first_nt(h) \end{array} \right. \\ \text{return } fir \end{array} \right.$$

Follow

Follow ou terminal seguinte de um símbolo não terminal **A**, corresponde a determinar o conjunto de símbolos terminais que se podem seguir àquilo em que A derive. O cálculo do **Follow(A)** vai depender do contexto em que **A** possa estar inserido. Esses contextos vão ser lidos nas produções em que **A** aparece no lado direito.

Nomeadamente se existe uma produção: **B** -> ... **A** **cont**

então o **Follow(A)** vai incluir o início de **cont** ou seja **first_str(cont)**; Quando **cont** for anulável, então aquilo que se puder seguir a **B** também pertencerá a **Follow(A)**.

função $follow (G : gramatica, A : NT) : T - set$

$$\left\{ \begin{array}{l} \text{para } (B \rightarrow \dots A \text{ cont}) \in Producoes(G) \text{ fazer} \\ \quad \left\{ \begin{array}{l} \text{se } [\text{cont é anulável}] \longrightarrow \\ \quad \left\{ \begin{array}{l} fol \leftarrow fol \cup first_str(G, cont) \cup follow(G, B) \end{array} \right. \\ \text{senão} \longrightarrow \\ \quad \left\{ \begin{array}{l} fol \leftarrow fol \cup first_str(G, cont) \end{array} \right. \end{array} \right. \\ \text{return } fol \end{array} \right.$$

Condição LL(1)

Para quaisquer duas derivações com o mesmo lado esquerdo, a **interseção dos seus lookaheads tem que ser vazia**.

Exemplo: cálculo dos lookaheads para o caso das listas

Gramática:

```
p1: Lista --> '[' ']'
p2:      | '[' Conteudo ']'

p3: Conteudo --> num
p4:      | num ',' Conteudo
```

Lookaheads:

```
la(p1) = First('[' ']) = {'['}
la(p2) = First('[' Conteudo ']') = {'['}

la(p1) ∩ la(p2) = {'['} ==> Conflito LL(1)

la(p3) = First(num) = {num}
la(p4) = First(num ',' Conteudo) = {num}

la(p3) ∩ la(p4) = {num} ==> Conflito LL(1)
```

Quando existem conflitos não é possível reconhecer a linguagem com um parser TopDown que esteja a ver apenas um símbolo à frente. A solução passa por transformar a gramática.

Exercício (da sebenta):

Considere a seguinte gramática para um subconjunto do LISP:

```
p1: S    -> Exp '.'
p2: Exp -> int
p3:      | '(' Funcao ')'
p4: Funcao -> '+' Lista
p5:      | '*' Lista
p6: Lista  -> Exp Lista
p7:      | ε
```

Identifique os conjuntos de produções potencialmente problemáticos e verifique a condição LL(1) para esses casos.

Analizador Sintático: recursivo descendente

```

# listas_anasin.py
# 2023-03-21 by jcr
# -----
import listas_analex

prox_simb = ('Erro', '', 0, 0)

# Lista --> '[' ']'
#         | '[' Conteudo ']'
def rec_Lista():
    global prox_simb
    ... Temos um problema!!!

def rec_Parser(data):
    global prox_simb
    lexer.input(data)
    prox_simb = lexer.token()
    rec_Lista()
    print("That's all folks!")

```

Com a alteração das produções de Lista

```

def parserError(simb):
    print("Erro sintático, token inesperado: ", simb)

def rec_term(simb):
    global prox_simb
    if prox_simb.type == simb:
        prox_simb = lexer.token()
    else:
        parserError(prox_simb)

# P4: Conteudo --> num
# P5:          | num ',' Conteudo
# É preciso alterar para:
# P4: Conteudo --> num Cont2
# P5: Cont2     -->
# P6: Cont2     --> ',' Conteudo

def rec_Cont2():
    global prox_simb
    if prox_simb.type == 'VIRG':
        rec_term('VIRG')
        rec_Conteudo()
        print("Reconheci P6: Cont2     --> ',' Conteudo")
    elif prox_simb.type == '???':
        print("Reconheci P5: Cont2 -->")
    else:
        parserError(prox_simb)

```



```

def rec_Conteudo():
    rec_term('NUM')
    rec_Cont2()
    print("Reconheci P4: Conteudo --> num Cont2")

def rec_LCont():
    global prox_simb
    if prox_simb.type == 'PF':
        rec_term('PF')
        print("Reconheci P2: LCont --> ']'")
    elif prox_simb.type == 'NUM':
        rec_Conteudo()
        rec_term('PF')
        print("Reconheci P3: LCont --> Conteudo ']'")
    else:
        parserError(prox_simb)

# P1: Lista --> '[' LCont
# P2: LCont --> ']'
# P3:      | Conteudo ']'
def rec_Lista():
    global prox_simb
    rec_term('PA')
    rec_LCont()
    print("Reconheci P1: Lista --> '[' LCont")

```

Gramática Concreta Final

```

P1: Lista --> '[' LCont
P2: LCont --> ']'
P3:      | Conteudo ']'
P4: Conteudo --> num Cont2
P5: Cont2    -->
P6: Cont2    --> ',' Conteudo

```

Gramática Abstrata

```

P1: Lista -->
P2:      | Conteudo
P3: Conteudo --> num Conteudo
P4:      |

```

Modelo em Python

```

class Lista

```

