

Formal languages and compilers

Projects 2018

Henri LEFEBVRE

February 18, 2019

1 Top-down parser : begin...end

2 Bottom-up parser : NLP with multiple clauses

Introduction

Problem statement

Parse input strings s respecting the following rules :

- s must start with "begin" and end with "end"
- both "begin...end" and "{...}" define block structures which must not overlap in s

Valid example

```
begin
  <text>
  {
    <text>
  }
  <text>
end
```

Invalid example

```
begin
  <text>
  {
    begin <text>
  }
  end <text>
end
```

Grammar definition :

$$\langle \text{pgm} \rangle ::= \langle \text{stmts} \rangle \$$$

$$\langle \text{stmts} \rangle ::= \text{begin } \langle \text{txt} \rangle \text{ end}$$

$$\langle \text{txt} \rangle ::= \langle \text{letter} \rangle \langle \text{txt} \rangle \mid \langle \text{stmts} \rangle \langle \text{txt} \rangle \mid \{ \langle \text{txt} \rangle \} \langle \text{txt} \rangle \mid \epsilon$$

$$\langle \text{letter} \rangle ::= a \mid \dots \mid z$$

Predictive parsing :

NT symbol	[a-z]	{	}	begin	end	\$	€
$\langle \text{pgm} \rangle$				$\langle \text{stmts} \rangle$			
$\langle \text{stmts} \rangle$				$\text{begin } \langle \text{txt} \rangle \text{ end}$			
$\langle \text{txt} \rangle$	$\langle \text{letter} \rangle \langle \text{txt} \rangle$	$\{ \langle \text{txt} \rangle \}$	€	$\langle \text{stmts} \rangle \langle \text{txt} \rangle$	€		
$\langle \text{letter} \rangle$	€						

- Implemented in C++
 - *Token* class to represent tokens
 - *ParsingTable* class to link token and their predicted successors
 - *Buffer* class to manage efficient reading of input string
 - *Interpreter* class to take actions when rules are recognized
 - *SyntacticAnalyser* class to generate token from the *Buffer* class
- see <https://github.com/hlefebvr/unige-top-down-parser>

1 Top-down parser : begin...end

2 Bottom-up parser : NLP with multiple clauses

Problem statement

- Parse in-english input strings
- Decompose multiple-clause sentences into several single-clause sentences

Worked example

The boy [who is tall, skinny and strange] smiles and laughs.

- The boy is tall
- The boy is skinny
- The boy is strange
- The boy smiles
- The boy laughs

```
< whole_sentence > ::= < sentence > .  
  < sentence > ::= < np? > < vp >  
    < np > ::= < ng > < rel_clause? >  
    < ng > ::= DET? CARD? ORD? QUANT? < ap >? NOUN  
      | < np > COORD < np >  
  < rel_clause > ::= REL_WORD < sentence >  
    < vp > ::= < single_vp > | < add_vp >  
  < single_vp > ::= < conjugated_verb > < np? >  
    | < conjugated_verb > < ap? >  
  < add_vp > ::= < vp > COORD < vp >  
  < conjugated_verb > ::= NEGATION? CONJUGATED_BE  
    | AUXILIARY NEGATION? VERB?  
    | VERB  
  < ap > ::= ADJ | < ap > < ap > | < ap > COORD < ap >
```

An external dictionary defines COORD, NOUN, ADJ, NEGATION, CONJUGATED_BE, VERB, ... which are token returned by the lexer

- Python 3.0
- Lark parsing library (<https://github.com/lark-parser/lark>)
 - with a custom lexer making use of custom dictionary
- Idea : transform the parsing tree into a simpler tree where, given one level of exploration (from bottom to up), we end up treating structures like $\langle np \rangle \langle rel_clause? \rangle \langle vp \rangle$
- see <https://github.com/hlefebvr/nlp-multiple-clauses-to-single>

he can speak italian and english.

- he can speak italian.
- he can speak english.

the boy and the cat who is tall, skinny and strange smile and laugh.

- the cat is tall.
- the cat is skinny.
- the cat is strange.
- the boy smile.
- the boy laugh.
- the cat smile.
- the cat laugh.

the dog eat a man who is tall.

- a man is tall.
- the dog eat a man.

the boy whom I met yesterday cried and laughed.

- I met yesterday the boy
- the boy cried.
- the boy laughed.

- Using a context free grammar enforces us to make choices.
Example : "this" recognized as noun or determinant ?
 - this dog is eating. ("this" as determinant)
 - this is strange. ("this" as noun)