

プログラミング言語処理

第2回 (平成15年度9月9日)

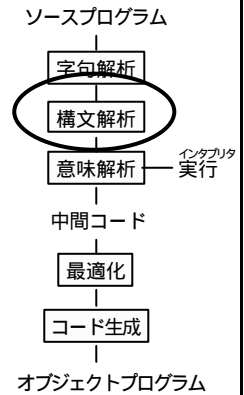
top-down parserの作り方 tiny Cの概要とデータ構造

筑波大学 佐藤

プログラミング言語処理

言語処理系の基本構成

- ◆ 字句解析 (lexical analysis): 文字列を言語の要素 (トークン、token) の列に分解する。
- ◆ 構文解析 (syntax analysis): token列を意味を反映した構造に変換。この構造は、しばしば、木構造で表現されるので、抽象構文木 (abstract syntax tree) と呼ばれる。ここまでの言語を認識する部分を言語の parser と呼ぶ。
- ◆ 意味解析 (semantics analysis): 構文木の意味を解析する。インタプリタでは、ここで意味を解析し、それに対応した動作を行う。コンパイラでは、この段階で内部的なコード、中間コードに変換する。



プログラミング言語処理

課題 1

- ◆ 掛け算、割り算の優先度を入れたインタプリタを作りなさい。tokenの種類に*や/に対応した演算子が増えることになる。入力として、
 $12*3 + 3*4 - 10$
をいれて、正しく実行できることを確認しなさい。
- ◆ さらに、括弧をいれた式が正しく処理できるよう拡張せよ。tokenの種類に括弧に対応するものが増えることになる。入力として、
 $12*(3+13) - 10$
をいれて、正しく実行できることを確認しなさい。できたプログラムを提出すること。

プログラミング言語処理

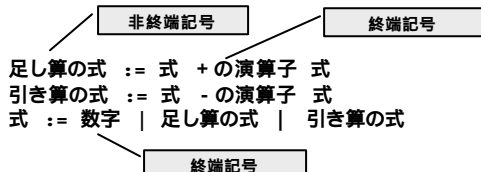
top-down parser

- ◆ これから、読み込むべきものを仮定して、それに合致するかを調べていく構文解析法
- 構文規則の上位から下位に向かって、parserを構成する。
- ◆ 構文規則から、直感的なプログラムを構成する。
- ◆ 人間にわかりやすい。人手でparserを書く場合に有効な方法。

プログラミング言語処理

構文規則のおさらい

- ◆ BNFによる数式の構文規則
- ◆ 構文規則の最終的な要素 = 終端記号 (terminal)
- ◆ ほかの構文規則によって定義される記号 = 非終端記号 (non-terminal)



プログラミング言語処理

構文規則の定義

- ◆ 定義
 1. 構文規則は、非終端記号 $\langle T \rangle$ に対して、 $\langle T \rangle ::= e$ (e は構文規則)で、表現される。これは、非終端記号 $\langle T \rangle$ は、構文規則 e によって置き換えられることを意味する。
 2. e は空でもよい。
 3. e は、非終端記号、終端記号、もしくは $e_1 \mid e_2$ 、 $e_1 e_2$ 、 e_1^* のいずれか。 $e_1 \mid e_2$ は、 e_1 もしくは e_2 であることを意味し、 $e_1 e_2$ は e_1 の次に e_2 が現れることを意味し、 e_1^* は、 e_1 の0個以上の繰り返しを意味する。
- (...)は、構文規則のまとまりを示す
- $e_1 \mid e_2 \mid e_3$ は、 $((e_1 \mid e_2) \mid e_3)$ と同じ
- $e_1 e_2 e_3$ は、 $((e_1 e_2) e_3)$ と同じである。
- 正規表現と似ている

構文規則と文脈

- ◆ 構文規則の左辺が、一つの終端記号 $\langle T \rangle$ だけという文法を考える。これはすなわち、どのような場合でも

$\langle T \rangle ::= e$

の規則を使って、右辺に置き換えることができることを意味し、このような文法を文脈自由文法とよぶ

- ◆ e_1 e_2 の間に構文要素 $\langle T \rangle$ が現れたときだけ、置き換えることができるというような文法が考えられるが、このような文法を文脈依存文法と呼ぶ

$e_1 \langle T \rangle e_2 ::= \dots$

top-down parser

- ◆ これから、読み込むべきものを仮定して、それに合致するかを調べていく構文解析法
 - 構文規則の上位から下位に向かって、parserを構成する。
 - 再帰的下向き構文解析 (recursive decent parsing)

- ◆ $\langle A \rangle ::= a \langle B \rangle c$ という構文規則に対し $\langle A \rangle$ のための構文解析関数 readA は以下のように作ることができる

$\langle A \rangle ::= a \langle B \rangle c$

```

readA()
{
    aを読み込む;
    readB(); /* Bを読み込むための関数を呼ぶ*/
    bを読み込む;
}
    
```

top down parserの構成法

- ◆ 非終端記号 $\langle T \rangle$ に対して、 $\langle T \rangle$ を読む関数 readT という関数を作る。

- ◆ token を先読みしておく。

- ◆ $\langle T \rangle ::= a b$ ならば、a に対する処理、b に対する処理を書く。

readT() { aに対する処理; bに対する処理; }

- ◆ $\langle T \rangle ::= a | b$ ならば、先読みしている token を使って、a か b かを選択

readT() { if(a or b) aに対する処理 else bに対する処理; }

- ◆ a が、終端記号ならば、先読みしている token が a であるかを判定する。

数式の例

```

<expression> ::= <expression> <expr_op> <term>
<term> ::= <term> <term_op> <factor>
<factor> ::= number | '(' <expression> ')'
<expr_op> ::= '+' | '-'
<term_op> ::= '*'
    
```

$\langle \text{expression} \rangle ::= \langle \text{expression} \rangle \langle \text{expr_op} \rangle \langle \text{term} \rangle$

```

readExpr()
{
    readExpr();
    readExprOp();
    readTerm();
}
    
```

readExpr() が永遠にとまらない!?
左再帰性の問題

左再帰性の問題

- ◆ 次の規則を持つ場合にうまくいかない (は空)
 - 最終的に $\langle T \rangle ::= \langle T \rangle$ になるもの

```

readExpr() {
    readTerm();
    while(readExprOp() is OK)
        readTerm();
}
    
```

- 以下のように書き換えたのと同様

```

<expression> ::= <term> <expression1>
<expression1> ::= <expr_op> <term> <expression1> |
    
```

左再帰性の問題

- ◆ $\langle T \rangle ::= b (c |)$
 - $(c |)$ は、c を読むかなにもしないか

- ◆ $\langle T \rangle$ の次に何がくるかによって、c を読むかどうかが決まるので、top-down parser では、処理ができないことになってしまう

