

Simple Parser

How to easily write parsers in Python

Author: Christophe Delord
Contact: cdelord@cdsoft.fr
Web site: <http://www.cdsoft.fr/sp.html>
Date: Sunday 07 March 2010
License: This software is released under the LGPL license.

Table of Contents

- 1 [Introduction and tutorial](#)
 - 1.1 [Introduction](#)
 - 1.2 [Installation](#)
 - 1.3 [Tutorial](#)
- 2 [SP reference](#)
 - 2.1 [Usage](#)
 - 2.2 [Grammar structure](#)
 - 2.3 [Lexer](#)
 - 2.4 [Parser](#)
 - 2.5 [Performances and memory consumption](#)
- 3 [Older Python versions](#)
 - 3.1 [Separators](#)
- 4 [SP mini language](#)
- 5 [Some examples to illustrate SP](#)
 - 5.1 [Newick format](#)
 - 5.2 [Infix/Prefix/Postfix notation converter](#)
 - 5.3 [Complete interactive calculator](#)

1 Introduction and tutorial

1.1 Introduction

SP (Simple Parser) is a Python¹ parser generator. It is aimed at easy usage rather than performance. SP produces [Top-Down Recursive descent](#) parsers. SP also uses [memoization](#) to optimize parsers' speed when dealing with ambiguous grammars.

License

SP is available under the GNU Lesser General Public:

Simple Parser: A Python parser generator

Copyright (C) 2009-2010 Christophe Delord

Simple Parser is free software: you can redistribute it and/or modify it under the terms of the GNU Lesser General Public License as published by the Free Software Foundation, either version 3 of the License, or (at your option) any later version.

Simple Parser is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the GNU Lesser General Public License for more details.

You should have received a copy of the GNU Lesser General Public License along with Simple Parser. If not, see <<http://www.gnu.org/licenses/>>.

Structure of the document

Introduction and tutorial starts smoothly with a gentle tutorial as an introduction. I think this tutorial may be sufficient to start with SP.

SP reference is a reference documentation. It will detail SP as much as possible.

Some examples to illustrate SP gives the reader some examples to illustrate SP.

1.2 Installation

Getting SP

SP is freely available on its web page (<http://www.cdsoft.fr/sp.html>).

Requirements

SP is a *pure Python* package. It may run on *any platform* supported by Python. The only requirement of SP is *Python 2.6* or newer². Python can be downloaded at <http://www.python.org>.

1.3 Tutorial

Introduction

This short tutorial presents how to make a simple calculator. The calculator will compute basic mathematical expressions (+, -, *, |) possibly nested in parenthesis. We assume the reader is familiar with regular expressions.

¹ Python is a wonderful object oriented programming language available at <http://www.python.org>

² Older *Python* versions may work (tested with Python 2.4 and 2.5). See the [Older Python versions](#) chapter.

Defining the grammar

Expressions are defined with a grammar. For example an expression is a sum of terms and a term is a product of factors. A factor is either a number or a complete expression in parenthesis.

We describe such grammars with rules. A rule describes the composition of an item of the language. In our grammar we have 3 items (expr, term, factor). We will call these items *symbols* or *non terminal symbols*. The decomposition of a symbol is symbolized with \rightarrow .

Grammar for expressions:

Grammar rule	Description
<code>expr -> term (('+' '-') term)*</code>	An expression is a term eventually followed with a plus (+) or a minus (-) sign and an other term any number of times (* is a repetition of an expression 0 or more times).
<code>term -> fact (('*' '/') fact)*</code>	A term is a factor eventually followed with a * or / sign and an other factor any number of times.
<code>fact -> ('+' '-') fact number '(' expr ')'</code>	A factor is either a factor preceded by a sign, a number or an expression in parenthesis.

We have defined here the grammar rules (i.e. the sentences of the language). We now need to describe the lexical items (i.e. the words of the language). These words - also called *terminal symbols* - are described using regular expressions. In the rules we have written some of these terminal symbols (+, -, *, /, (,)). We have to define **number**. For sake of simplicity numbers are integers composed of digits (the corresponding regular expression can be `[0-9]+`). To simplify the grammar and then the Python script we define two terminal symbols to group the operators (additive and multiplicative operators). We can also define a special symbol that is ignored by SP. This symbol is used as a separator. This is generally useful for white spaces and comments.

Terminal symbol definition for expressions:

Terminal symbol	Regular expression	Comment
number	<code>[0-9]+</code> or <code>\d+</code>	One or more digits
addop	<code>[+-]</code>	a + or a -
mulop	<code>[*/]</code>	a * or a /
spaces	<code>\s+</code>	One or more spaces

This is sufficient to define our parser with SP.

Grammar of the expression recognizer:

```
def Calc():

    number = R(r'[0-9]+')
    addop = R(r'[+-]')
    mulop = R(r'[*/]')

    with Separator(r'\s+'):

        expr = Rule()
        fact = Rule()
        fact |= addop & fact
        fact |= '(' & expr & ')'
```

```

fact |= number
term = fact & ( mulop & fact )[:]
expr |= term & ( addop & term )[:]

return expr

```

`Calc` is the name of the Python function that returns a parser. This function returns `expr` which is the *axiom*³ of the grammar.

`expr` and `fact` are recursive rules. They are first declared as empty rules (`expr = Rule()`) and alternatives are later added (`expr |= ...`).

Slices are used to implement repetitions. `foo[:]` parses `foo` zero or more times, which is equivalent to `foo*` is a classical grammar notation.

The grammar can also be defined with the mini grammar language provided by SP:

```

def Calc():
    return compile("""
        number = r'[0-9]+' ;
        addop = r'[+-]' ;
        mulop = r'[*/]' ;

        separator: r'\s+' ;

        !expr = term (addop term)* ;
        term = fact (mulop fact)* ;
        fact = addop fact ;
        fact = '(' expr ')' ;
        fact = number ;
    """)

```

Here the *axiom*³ is identified by `!`.

With this small grammar we can only recognize a correct expression. We will see in the next sections how to read the actual expression and to compute its value.

Reading the input and returning values

The input of the grammar is a string. To do something useful we need to read this string in order to transform it into an expected result.

This string can be read by catching the return value of terminal symbols. By default any terminal symbol returns a string containing the current token. So the token `'('` always returns the string `'('`. For some tokens it may be useful to compute a Python object from the token. For example `number` should return an integer instead of a string, `addop` and `mulop`, followed by a number, should return a function corresponding to the operator. That's why we will add a function to the token and rule definitions. So we associate `int` to `number` and `op1` and `op2` to unary and binary operators.

`int` is a Python function converting objects to integers and `op1` and `op2` are user defined functions. `op1` and `op2` functions:

```

op1 = lambda f,x: {'+':pos, '-':neg}[f](x)
op2 = lambda f,y: lambda x: {'+': add, '-': sub, '*': mul, '/': div}[f](x,y)

# red applyies functions to a number

```

³ The axiom is the symbol from which the parsing starts

```
def red(x, fs):
    for f in fs: x = f(x)
    return x
```

To associate a function to a token or a rule it must be applied using / or * operators:

- / applies a function to an object returned by a (sub)parser.
- * applies a function to an tuple of objects returned by a sequence of (sub) parsers.

Token and rule definitions with functions:

```
number = R(r'[0-9]+') / int

fact |= (addop & fact) * op1
term = (fact & ( (mulop & fact) * op2 )[:]) * red

# R(r'[0-9]+') applied on "42" will return "42".
# R(r'[0-9]+') / int will return int("42")

# addop & fact applied on "+ 42" will return ('+', 42)
# (addop & fact) * op1 will return op1(*('+', 42)), i.e. op1('+', 42)
# so (addop & fact) * op1 returns +42

# (addop & fact) * op2 will return op2(*('+', 42)), i.e. op2('+', 42)
# so (addop & fact) * op2 returns lambda x: add(x, 42)

# fact & ( (mulop & fact) * op2 )[:] returns a number and a list of func-
tions
# for instance (42, [(lambda x:mul(x, 43)), (lambda x:mul(x, 44))])
# so (fact & ( (mulop & fact) * op2 )[:]) * red applied on "42*43*44"
# will return red(42, [(lambda x:mul(x, 43)), (lambda x:mul(x, 44))])
# i.e. 42*43*44
```

And with the SP language:

```
number = r'[0-9]+' : 'int' ;

addop = r'[+-]' ;
mulop = r'[*/]' ;

fact = addop fact :: 'op1' ;
term = fact (mulop fact :: 'op2')* :: 'red' ;

# r'[0-9]+' applied on "42" will return "42".
# r'[0-9]+' : 'int' will return int("42")

# "addop fact" applied on "+ 42" will return ('+', 42)
# "addop fact :: 'op1'" will return op1(*('+', 42)), i.e. op1('+', 42)
# so "addop fact :: 'op1'" returns +42

# "addop fact :: 'op2'" will return op2(*('+', 42)), i.e. op2('+', 42)
# so "addop fact :: 'op2'" returns lambda x: add(x, 42)

# "fact (mulop fact :: 'op2')*" returns a number and a list of functions
```

```
# for instance (42, [(lambda x:mul(x, 43)), (lambda x:mul(x, 44))])
# so "fact (mulop fact :: 'op2')* :: 'red'" applyied on "42*43*44"
# will return red(42, [(lambda x:mul(x, 43)), (lambda x:mul(x, 44))])
# i.e. 42*43*44
```

In the SP language, : (as /) applies a Python function (more generally a callable object) to a value returned by a sequence and :: (as *) applies a Python function to several values returned by a sequence. Here is finally the complete parser.

Expression recognizer and evaluator:

```
from sp import *

def Calc():

    from operator import pos, neg, add, sub, mul, truediv as div

    op1 = lambda f,x: {'+':pos, '-':neg}[f](x)
    op2 = lambda f,y: lambda x: {'+': add, '-':
': sub, '*': mul, '/': div}[f](x,y)

    def red(x, fs):
        for f in fs: x = f(x)
        return x

    number = R(r'[0-9]+') / int
    addop = R('[+-]')
    mulop = R('[*/]')

    with Separator(r'\s+'):

        expr = Rule()
        fact = Rule()
        fact |= (addop & fact) * op1
        fact |= '(' & expr & ')'
        fact |= number
        term = (fact & ( (mulop & fact) * op2 )[:]) * red
        expr |= (term & ( (addop & term) * op2 )[:]) * red

    return expr
```

Or with SP language:

```
from sp import *

def Calc():

    from operator import pos, neg, add, sub, mul, truediv as div

    op1 = lambda f,x: {'+':pos, '-':neg}[f](x)
    op2 = lambda f,y: lambda x: {'+': add, '-':
': sub, '*': mul, '/': div}[f](x,y)

    def red(x, fs):
        for f in fs: x = f(x)
```

```

    return x

return compile("""
    number = r'[0-9]+' : 'int' ;
    addop = r'[+-]' ;
    mulop = r'[*/]' ;

    separator: r'\s+' ;

    !expr = term (addop term :: 'op2')* :: 'red' ;
    term = fact (mulop fact :: 'op2')* :: 'red' ;
    fact = addop fact :: 'op1' ;
    fact = '(' expr ')' ;
    fact = number ;
""")

```

Embedding the parser in a script

A parser is a simple Python object. This example show how to write a function that returns a parser. The parser can be applied to strings by simply calling the parser.

Writting SP grammars in Python:

```

from sp import *

def MyParser():

    parser = ...

    return parser

# You can instanciate your parser here
my_parser = MyParser()

# and use it
parsed_object = my_parser(string_to_be_parsed)

```

To use this parser you now just need to instanciate an object.
Complete Python script with expression parser:

```

from sp import *

def Calc():

    ...

calc = Calc()
while True:
    expr = input('Enter an expression: ')
    try: print(expr, '=', calc(expr))
    except Exception as e: print("%s:%s" % e.__class__.__name__, e)

```

Conclusion

This tutorial shows some of the possibilities of SP. If you have read it carefully you may be able to start with SP. The next chapters present SP more precisely. They contain more examples to illustrate all the features of SP.

Happy SP'ing!

2 SP reference

2.1 Usage

SP is a package which main function is to provide basic objects to build a complete parser.

The grammar is a Python object.

Grammar embedding example:

```
def Foo():
    bar = R('bar')
    return bar
```

Then you can use the new generated parser. The parser is simply a Python object.

Parser usage example:

```
test = "bar"
my_parser = Foo()
x = my_parser(test)           # Parses "bar"
print x
```

2.2 Grammar structure

SP grammars are Python objects. SP grammars may contain two parts:

Tokens are built by the R or K keywords.

Rules are described after tokens in a **Separator** context.

Example of SP grammar structure:

```
def Foo():

    # Tokens
    number = R(r'\d+') / int

    # Rules
    with Separator(r'\s+'):
        S = number[:]

    return S

foo = Foo()
result = foo("42 43 44") # return [42, 43, 44]
```


2.3 Lexer

Regular expression syntax

The lexer is based on the *re*⁴ module. SP profits from the power of Python regular expressions. This document assumes the reader is familiar with regular expressions.

You can use the syntax of regular expressions as expected by the *re*⁵ module.

Predefined tokens

Tokens can be explicitly defined by the `R`, `K` and `Separator` keywords.

Expression	Usage
<code>R</code>	defines a regular token. The token is defined with a regular expression and returns a string (or a tuple of strings if the regular expression defines groups).
<code>K</code>	defines a token that returns nothing (useful for keywords for instance). The keyword is defined by an identifier (in this case word boundaries are expected around the keyword) or another string (in this case the pattern is not considered as a regular expression). The token just recognizes a keyword and returns nothing.
<code>Separator</code>	is a context manager used to define separators for the rules defined in the context. The token is defined with a regular expression and returns nothing.

A token can be defined by:

a name which identifies the token. This name is used by the parser.

a regular expression which describes what to match to recognize the token.

an action which can translate the matched text into a Python object. It can be a function of one argument or a non callable object. If it is not callable, it will be returned for each token otherwise it will be applied to the text of the token and the result will be returned. This action is optional. By default the token text is returned.

Token definition examples:

```
integer = R(r'\d+') / int
identifier = R(r'[a-zA-Z]\w*\b')
boolean = R(r'(True|False)\b') / (lambda b: b=='True')

spaces = K(r'\s+')
comments = K(r'#. *')

with Separator(spaces|comments):
    # rules defined here will use spaces and comments as separators
    atom = '(' & expr & ')'
```

There are two kinds of tokens. Tokens defined by the `R` or `K` keywords are parsed by the parser and tokens defined by the `Separator` keyword are considered as separators (white spaces or comments for example) and are wiped out by the lexer.

The word boundary `\b` can be used to avoid recognizing “True” at the beginning of “Truexyz”.

If the regular expression defines groups, the parser returns a tuple containing these groups:

```
couple = R('<(\d+)-(\d+)>')

couple("<42-43>") == ('42', '43')
```

⁴ *re* is a standard Python module. It handles regular expressions. For further information about *re* you can read <http://docs.python.org/library/re.html> 9

⁵ Read the Python documentation for further information: <http://docs.python.org/library/re.html#re-syntax>

If the regular expression defines only one group, the parser returns the value of this group:

```
first = R('<(\d+)-\d+>')
```

```
first("<42-43>") == '42'
```

Unwanted groups can be avoided using `(?:...)`.

A name can be given to a token to make error messages easier to read:

```
couple = R('<(\d+)-(\d+)>', name="couple")
```

Regular expressions can be compiled using specific compilation options. Options are defined in the `re` module:

```
token = R('...', flags=re.IGNORECASE|re.DOTALL)
```

`re` defines the following flags:

I (IGNORECASE) Perform case-insensitive matching.

L (LOCALE) Make `\w`, `\W`, `\b`, `\B`, dependent on the current locale.

M (MULTILINE) `"^"` matches the beginning of lines (after a newline) as well as the string. `"$"` matches the end of lines (before a newline) as well as the end of the string.

S (DOTALL) `"."` matches any character at all, including the newline.

X (VERBOSE) Ignore whitespace and comments for nicer looking RE's.

U (UNICODE) Make `\w`, `\W`, `\b`, `\B`, dependent on the Unicode locale

Inline tokens

Tokens can also be defined on the fly. Their definition are then inlined in the grammar rules. This feature may be useful for keywords or punctuation signs.

In this case tokens can be written without the `R` or `K` keywords. They are considered as keywords (as defined by `K`).

Inline token definition examples:

```
IfThenElse = 'if' & Cond &
              'then' & Statement &
              'else' & Statement
```

2.4 Parser

Declaration

A parser is declared as a Python object.

Grammar rules

Rule declarations have two parts. The left side declares the symbol associated to the rule. The right side describes the decomposition of the rule. Both parts of the declaration are separated with an equal sign (`=`).

Rule declaration example:

```
SYMBOL = (A & B) * (lambda a, b: f(a, b))
```

Sequences

Sequences in grammar rules describe in which order symbols should appear in the input string. For example the sequence `A & B` recognizes an `A` followed by a `B`.

For example to say that a `sum` is a `term` plus another `term` you can write:

```
Sum = Term & '+' & Term
```

Alternatives

Alternatives in grammar rules describe several possible decompositions of a symbol. The infix pipe operator (`|`) is used to separate alternatives. `A | B` recognizes either an `A` or a `B`. If both `A` and `B` can be matched only the first longest match is considered. So the order of alternatives may be very important when two alternatives can match texts of the same size.

For example to say that an `atom` is an *integer* or an *expression in parenthesis* you can write:

```
Atom = integer | '(' & Expr & ')'
```

Repetitions

Repetitions in grammar rules describe how many times an expression should be matched.

Expression	Usage
<code>A[:1]</code>	recognizes zero or one <code>A</code> .
<code>A[:]</code>	recognizes zero or more <code>A</code> .
<code>A[1:]</code>	recognizes one or more <code>A</code> .
<code>A[m:n]</code>	recognizes at least <code>m</code> and at most <code>n</code> <code>A</code> .
<code>A[m:n:s]</code>	recognizes at least <code>m</code> and at most <code>n</code> <code>A</code> using <code>s</code> as a separator.

Repetitions are greedy. Repetitions are implemented as Python loops. Thus whatever the length of the repetitions, the Python stack will not overflow.

The separator is useful to parse lists. For instance a comma separated parameter list is `parameter[:,',']`.

Precedence and grouping

The following table lists the different structures in increasing precedence order. To override the default precedence you can group expressions with parenthesis.

Precedence in SP expressions:

Structure	Example
Alternative	<code>A B</code>
Sequence	<code>A & B</code>
Repetitions	<code>A[x:y]</code>
Symbol and grouping	<code>A and (...)</code>

Actions

Grammar rules can contain actions as Python functions.
Functions are applied to parsed objects using / or *.

Expression	Value
<code>parser / function</code>	returns <i>function(result of parser)</i> .
<code>parser * function</code>	returns <i>function(*result of parser)</i> .

* can be used to analyse the result of a sequence.

Abstract syntax trees

An abstract syntax tree (AST) is an abstract representation of the structure of the input. A node of an AST is a Python object (there is no constraint about its class). AST nodes are completely defined by the user.

AST example (parsing a couple):

```
class Couple:
    def __init__(self, a, b):
        self.a = a
        self.b = b

def Foo():
    couple = ((' & item & ',' & item & ')) * Couple
    return couple
```

Constants

It is sometimes useful to return a constant. C defines a parser that matches an empty input and returns a constant.

Constant example:

```
number = ( '1' & C("one")
          | '2' & C("two")
          | '3' & C("three")
          )
```

Position in the input string

To know the current position in the input string, the At() parser returns an object containing the current index (attribute `index`) and the corresponding line and column numbers (attributes `line` and `column`):

```
position = At() / 'lambda p: (p.line, p.column)'  
rule = ... & pos & ...
```

2.5 Performances and memory consumption

Backtracking has a cost. The parser may often try to parse again the same string at the same position. To improve the speed of the parser, some time consuming functions are *memoized*. This drastically

fasten the parser but requires more memory. If a lot of string are parsed in a single script this mechanism can slow down the computer because of heavy swap disk usage or even lead to a memory error.

To avoid such problems it is recommended to clean the memoization cache by calling the `sp.clean` function:

```
import sp

...

for s in a_lot_of_strings:
    parse(s)
    sp.clean()
```

3 Older Python versions

This document describes the usage of SP with Python 2.6. Grammars need some adaptations to work with Python 2.5. or older.

3.1 Separators

Separators use context managers which don't exist in Python 2.4. Context managers have been introduced in Python 2.5 (from `__future__` import `with_statement`) and in Python 2.6 (as a standard feature). When the context managers are not available, it may be possible to call the `__enter__` and `__exit__` method explicitly (tested for Python 2.4).

Python 2.6 and later:

```
number = R(r'\d+') / int
with Separator('\s+'):
    coord = number & ', ' & number
```

Python 2.5 with `with_statement`:

```
from __future__ import with_statement

number = R(r'\d+') / int
with Separator('\s+'):
    coord = number & ', ' & number
```

Python 2.5 or 2.4 (or older but not tested) without `with_statement`:

```
sep = Separator('\s+')

number = R(r'\d+') / int
sep.__enter__()
coord = number & ', ' & number
sep.__exit__()
```

4 SP mini language

Instead of using Python expressions that can sometimes be difficult to read, it's possible to write grammars in a cleaner syntax and compile these grammar with the `sp.compile` function. This function takes the grammar as a string parameter. The `sp.compile_file` function reads the grammar in a separate file.

Here the equivalence between Python expressions and the SP mini language:

SP Python expressions	SP mini language	Description
<code>R("regular expression")</code> <code>R("regexpr", name="name")</code>	<code>r"regular expression"</code> <code>name.r"regexpr"</code>	Token defined by a regular expression
<code>K("plain text")</code> <code>K("plain text", name="name")</code>	<code>"plain text"</code> <code>name."plain text"</code>	Keyword defined by a non interpreted string
<code>t = R('...', flags=re.I re.S)</code>	<code>lexer: I S; t = r'...'</code>	Regular expression options
<code>with Separator(...):</code>	<code>separator: ... ;</code>	Separator definition
<code>C(object)</code>	<code>'object'</code>	Parses nothing and returns object
<code>... / function</code>	<code>... : 'function'</code>	Parses ... and apply the result to <code>function(function(...))</code>
<code>... * function</code>	<code>... :: 'function'</code>	Parses ... and apply the result (multiple values) to <code>function(function(*...))</code>
<code>... & At() & ...</code>	<code>... @ ...</code>	Position in the input string
<code>(...)[:]</code>	<code>(...)*</code>	Zero or more matches
<code>(...)[1:]</code>	<code>(...)+</code>	One or more matches
<code>(...)[:1]</code>	<code>(...)?</code>	Zero or one matche
<code>(...)[::S]</code>	<code>[.../S]*</code>	Zero or more matches separated by S
<code>(...)[1::S]</code>	<code>[.../S]+</code>	One or more matches separated by S
<code>A & B & C</code>	<code>A B C</code>	Sequence
<code>A B C</code>	<code>A B C</code>	Alternative
<code>(...)</code>	<code>(...)</code>	Grouping
<code>rule_name = ...</code>	<code>rule_name = ... ;</code>	Rule definition
<code>axiom_name = ...</code>	<code>!axiom_name = ... ;</code>	Axiom definition

5 Some examples to illustrate SP

5.1 Newick format

In mathematics, Newick tree format (or Newick notation or New Hampshire tree format) is a way to represent graph-theoretical trees with edge lengths using parentheses and commas. It was created by James Archie, William H. E. Day, Joseph Felsenstein, Wayne Maddison, Christopher Meacham, F. James Rohlf, and David Swofford, at two meetings in 1986, the second of which was at Newick's restaurant in Dover, New Hampshire, USA.

—Wikipedia, the free encyclopedia

The grammar given by Wikipedia is:

```

Tree --> Subtree ";" | Branch ";"
Subtree --> Leaf | Internal
Leaf --> Name
Internal --> "(" BranchSet ")" Name
BranchSet --> Branch | Branch "," BranchSet
Branch --> Subtree Length
Name --> empty | string
Length --> empty | ":" number

```

With very few transformation, this grammar can be converted to a Simple Parser grammar. Only BranchSet is rewritten to use a comma separated list parser:

```

Tree = Subtree ';' | Branch ';' ;
Subtree = Leaf | Internal ;
Leaf = Name ;
Internal = '(' [Branch/',' ]+ ')' Name ;
Branch = Subtree Length ;
Name = r'[^;:,()]*';
Length = '' | ':' r'[0-9.]+' ;

```

Here is the complete parser (newick.py):

```

#!/usr/bin/env python

# Simple Parser
# Copyright (C) 2009-2010 Christophe Delord
# http://www.cdsoft.fr/sp.html

# This file is part of Simple Parser.
#
# Simple Parser is free software: you can redistribute it and/or modify
# it under the terms of the GNU Lesser General Public License as published
# by the Free Software Foundation, either version 3 of the License, or
# (at your option) any later version.
#
# Simple Parser is distributed in the hope that it will be useful,
# but WITHOUT ANY WARRANTY; without even the implied warranty of
# MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the
# GNU Lesser General Public License for more details.
#
# You should have received a copy of the GNU Lesser General Public License
# along with Simple Parser. If not, see <http://www.gnu.org/licenses/>.

# from http://en.wikipedia.org/wiki/Newick_format

import sp

EXAMPLES = """\
(,,(,));                no nodes are named
(A,B,(C,D));            leaf nodes are named
(A,B,(C,D)E)F;          all nodes are named
(:0.1,:0.2,(0.3,:0.4):0.5); all but root node have a dis-
tance to parent
(:0.1,:0.2,(0.3,:0.4):0.5):0.0; all have a distance to parent

```



```

(A:0.1,B:0.2,(C:0.3,D:0.4):0.5);      distances and leaf names (popular)
(A:0.1,B:0.2,(C:0.3,D:0.4)E:0.5)F;    distances and all names
((B:0.2,(C:0.3,D:0.4)E:0.5)F:0.1)A;   a tree rooted on a leaf node (rare)
"""

class Leaf:
    def __init__(self, name): self.name = name
    def __str__(self): return self.name
    def nb_leaves(self): return 1

class Internal:
    def __init__(self, subtrees, name): self.subtrees, self.name = sub-
trees, name
    def __str__(self): re-
turn "(%s)%s" % ('','.join(str(st) for st in self.subtrees), self.name)
    def nb_leaves(self): return sum(st.nb_leaves() for st in self.subtrees)

class Branch:
    def __init__(self, subtree, length): self.subtree, self.length = sub-
tree, length
    def __str__(self): return "%s:%s"%(self.subtree, self.length)
    def nb_leaves(self): return self.subtree.nb_leaves()

parser = sp.compile(r"""
!Tree = Subtree ';' | Branch ';' ;
Subtree = Leaf | Internal ;
Leaf = Name : 'Leaf' ;
Internal = '(' [Branch/',' ]+ ')' Name :: 'Internal' ;
Branch = Subtree Length :: 'Branch' ;
Name = r'[^;:,()]*';
Length = ':' r'[0-9.]+': 'float' | '0.0' ;
""")

for example in EXAMPLES.splitlines():
    example, description = example.split(' ', 1)
    description = description.strip()
    tree = parser(example)
    print "%s:"%description
    print "-"*len(description)
    print "    Input :", example
    print "    Parsed:", tree
    print "    Leaves:", tree.nb_leaves()
    print

```

5.2 Infix/Prefix/Postfix notation converter

Introduction

In the previous example, the parser computes the value of the expression on the fly, while parsing. It is also possible to build an abstract syntax tree to store an abstract representation of the input. This may be useful when several passes are necessary.

This example shows how to parse an expression (infix, prefix or postfix) and convert it in infix, prefix and postfix notation. The expression is saved in a tree. Each node of the tree correspond to

an operator in the expression. Each leaf is a number. Then to write the expression in infix, prefix or postfix notation, we just need to walk through the tree in a particular order.

Abstract syntax trees

The AST of this converter has three types of node:

class Op is used to store operators (+, -, *, /, ^). It has two sons associated to the sub expressions.

class Atom is an atomic expression (a number or a symbolic name).

class Func is used to store functions.

These classes are instantiated by the init method. The infix, prefix and postfix methods return strings containing the representation of the node in infix, prefix and postfix notation.

Grammar

Lexical definitions

```
ident = r'\w+' : 'Atom' ;

func1 = r'sin' | r'cos' | r'tan' ;
func2 = r'min' | r'max' ;

op = op_add | op_mul | op_pow ;
op_add = r'[+-]' ;
op_mul = r'[*/]' ;
op_pow = r'\^' ;
```

Infix expressions The grammar for infix expressions is similar to the grammar used in the previous example:

```
expr = term (op_add term :: 'lambda op, y: lambda x: Op(op, x, y)')* :: 'red' ;
term = fact (op_mul fact :: 'lambda op, y: lambda x: Op(op, x, y)')* :: 'red' ;
fact = atom (op_pow fact :: 'lambda op, y: lambda x: Op(op, x, y)')? :: 'red' ;
atom = ident ;
atom = '(' expr ')' ;
atom = func1 '(' expr ')' :: 'Func' ;
atom = func2 '(' expr ',' expr ')' :: 'Func' ;
```

Prefix expressions The grammar for prefix expressions is very simple. A compound prefix expression is an operator followed by two subexpressions, or a binary function followed by two subexpressions, or a unary function followed by one subexpression:

```
expr_pre = ident ;
expr_pre = op expr_pre expr_pre :: 'lambda op, x, y: Op(op, x, y)' ;
expr_pre = func1 expr_pre :: 'Func' ;
expr_pre = func2 expr_pre expr_pre :: 'Func' ;
```

Postfix expressions At first sight postfix and infix grammars may be very similar. Only the position of the operators changes. So a compound postfix expression is a first expression followed by a second one and an operator. This rule is left recursive. As SP is a descendant recursive parser, such rules are forbidden to avoid infinite recursion. To remove the left recursion a classical solution is to rewrite the grammar like this:

```
expr_post = ident sexpr_post :: 'lambda x, f: f(x)' ;

sexpr_post = expr_post op :: 'lambda y, op: lambda x: Op(op, x, y)' ;
sexpr_post = expr_post func2 :: 'lambda y, f: lambda x: Func(f, x, y)' ;
sexpr_post = func1 :: 'lambda f: lambda x: Func(f, y)' ;
```

The parser searches for an atomic expression and builds the AST corresponding to the remaining subexpression. `sexpr_post` returns a function that builds the complete AST when applied to the first atomic expression. This is a way to simulate inherited attributes.

Source code

Here is the complete source code (notation.py):

```
#!/usr/bin/env python3

# Simple Parser
# Copyright (C) 2009-2010 Christophe Delord
# http://www.cdsoft.fr/sp.html

# This file is part of Simple Parser.
#
# Simple Parser is free software: you can redistribute it and/or modify
# it under the terms of the GNU Lesser General Public License as published
# by the Free Software Foundation, either version 3 of the License, or
# (at your option) any later version.
#
# Simple Parser is distributed in the hope that it will be useful,
# but WITHOUT ANY WARRANTY; without even the implied warranty of
# MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the
# GNU Lesser General Public License for more details.
#
# You should have received a copy of the GNU Lesser General Public License
# along with Simple Parser. If not, see <http://www.gnu.org/licenses/>.

# Infix/prefix/postfix expression conversion

import sp

try:
    import readline
except ImportError:
    pass

class Op:
    """ Binary operator """
    precedence = {'+':1, '-':1, '*':2, '/':2, '^':3}
    def __init__(self, op, a, b):
```

```

        self.op = op                    # operator ("+", "-",
", "*", "/", "^")
        self.prec = Op.precedence[op]  # precedence of the operator
        self.a, self.b = a, b          # operands
    def infix(self):
        a = self.a.infix()
        if self.a.prec < self.prec: a = "(%s)"%a
        b = self.b.infix()
        if self.b.prec <= self.prec: b = "(%s)"%b
        return "%s %s %s"%(a, self.op, b)
    def prefix(self):
        a = self.a.prefix()
        b = self.b.prefix()
        return "%s %s %s"%(self.op, a, b)
    def postfix(self):
        a = self.a.postfix()
        b = self.b.postfix()
        return "%s %s %s"%(a, b, self.op)

class Atom:
    """ Atomic expression """
    def __init__(self, s):
        self.a = s
        self.prec = 99
    def infix(self): return self.a
    def prefix(self): return self.a
    def postfix(self): return self.a

class Func:
    """ Function expression """
    def __init__(self, name, *args):
        self.name = name
        self.args = args
        self.prec = 99
    def infix(self):
        args = [a.infix() for a in self.args]
        return "%s(%s)"%(self.name, ",".join(args))
    def prefix(self):
        args = [a.prefix() for a in self.args]
        return "%s %s"%(self.name, " ".join(args))
    def postfix(self):
        args = [a.postfix() for a in self.args]
        return "%s %s"(" ".join(args), self.name)

# Grammar for arithmetic expressions

def red(x, fs):
    for f in fs: x = f(x)
    return x

parser = sp.compile(r"""

    ident = r'\w+' : 'Atom' ;

```

```

func1 = r'sin' | r'cos' | r'tan' ;
func2 = r'min' | r'max' ;

op = op_add | op_mul | op_pow ;
op_add = r'[+-]' ;
op_mul = r'[*/]' ;
op_pow = r'\^' ;

separator: r'\s+' ;

!axiom = expr      "infix"
        | expr_pre "prefix"
        | expr_post "postfix"
        ;

# Infix expressions

expr = term (op_add term :: 'lambda op, y: lambda x: Op(op, x, y)')* :: 'red' ;
term = fact (op_mul fact :: 'lambda op, y: lambda x: Op(op, x, y)')* :: 'red' ;
fact = atom (op_pow fact :: 'lambda op, y: lambda x: Op(op, x, y)')? :: 'red' ;
atom = ident ;
atom = '(' expr ')' ;
atom = func1 '(' expr ')' :: 'Func' ;
atom = func2 '(' expr ',' expr ')' :: 'Func' ;

# Prefix expressions

expr_pre = ident ;
expr_pre = op expr_pre expr_pre :: 'Op' ;
expr_pre = func1 expr_pre :: 'Func' ;
expr_pre = func2 expr_pre expr_pre :: 'Func' ;

# Postfix expressions

expr_post = ident sexpr_post :: 'lambda x, f: f(x)' ;

sexpr_post = expr_post op :: 'lambda y, op: lambda x: Op(op, x, y)' ;
sexpr_post = expr_post func2 :: 'lambda y, f: lambda x: Func(f, x, y)' ;
sexpr_post = func1 :: 'lambda f: lambda x: Func(f, y)' ;

"""

while 1:
    e = input(":")
    if e == "": break
    try:
        expr, t = parser(e)
    except Exception as e:
        print(e)
    else:
        print(e, "is a", t, "expression")
        print("\tinfix    :", expr.infix())

```

```
print("\tprefix  :", expr.prefix())
print("\tpostfix :", expr.postfix())
```

5.3 Complete interactive calculator

This chapter presents an extension of the calculator described in the [tutorial](#). This calculator has more functions and a memory.

The grammar has been rewritten using the SP language.

New functions

The calculator has memories. A memory cell is identified by a name. For example, if the user types `pi = 3.14`, the memory cell named `pi` will contain the value of `pi` and `2*pi` will return `6.28`.

The variables are saved in a dictionary.

Source code

The complete source code is available in the example directory of the archive.