

Semantic Web — Spring 2009
Homework #1

Simple XML Parser using Python Lex/Yacc

Behnam Esfahbod
behnam@sharif.edu

April 15, 2009

Listing 1: The Python program

```
1 #!/usr/bin/env python

import sys
from UserString import UserString

6 from ply import lex, yacc

#####
11 # LEXER

tokens = [

    # INITIAL
16     'CDATA',

    'OPENTAGOPEN',
    'CLOSETAGOPEN',

21     # tag
    'TAGATTRNAME',

    'TAGCLOSE',
    'LONETAGCLOSE',

26     'ATTRASSIGN',

    # attrvalue1
31     'ATTRVALUE1OPEN',
    'ATTRVALUE1STRING',
    'ATTRVALUE1CLOSE',

    # attrvalue2
36     'ATTRVALUE2OPEN',
    'ATTRVALUE2STRING',
    'ATTRVALUE2CLOSE',

]

41 # Regular expressions

re_digit      = r'([0-9])'
```

```

re_nondigit    = r'[_A-Za-z]',
re_identifier  = r'(' + re_nondigit + r'(' + re_digit + r'|' + re_nondigit + r')*)',
46
class SyntaxError(Exception):
    pass

51
class XmlLexer:
    # The XML Tokenizer

    # states:
56
    #
    # default:
    #     The default context, non-tag texts
    # tag:
    #     A document tag
61
    # string:
    #     Within quote-delimited strings inside tags

    states = (
        ('tag', 'exclusive'),
66
        ('attrvalue1', 'exclusive'),
        ('attrvalue2', 'exclusive'),
    )

    tokens = tokens

71

    # ANY

    def t_ANY_error(self, t):
76
        raise SyntaxError("Illegal character '%s'" % t.value[0])
        t.lexer.skip(1)
        pass

81

    # INITIAL

    t_ignore = ''

    def t_CLOSETAGOPEN(self, t):
86
        r'</'
        t.lexer.push_state('tag')
        return t

    def t_OPENTAGOPEN(self, t):
91
        r'<'
        t.lexer.push_state('tag')
        return t

    def t_CDATA(self, t):
96
        r'[<]+'
        return t

    # tag: name

101

    t_tag_ignore = ' \t'

    def t_tag_TAGATTRNAME(self, t):
        return t
106
    t_tag_TAGATTRNAME.__doc__ = re_identifier

    def t_tag_TAGCLOSE(self, t):
        r'>'
        t.lexer.pop_state()
111
        return t

```

```

116 def t_tag_LONETAGCLOSE(self, t):
    r'>'
    t.lexer.pop_state()
    return t

# tag: attr

121 t_tag_ATTRASSIGN = r'='

def t_tag_ATTRVALUE1OPEN(self, t):
    r'\''
    t.lexer.push_state('attrvalue1')
126 return t

def t_tag_ATTRVALUE2OPEN(self, t):
    r'"'
    t.lexer.push_state('attrvalue2')
131 return t

# attrvalue1

136 def t_attrvalue1_ATTRVALUE1STRING(self, t):
    r'[^\']*'
    t.value = unicode(t.value)
    return t

141 def t_attrvalue1_ATTRVALUE1CLOSE(self, t):
    r'\''
    t.lexer.pop_state()
    return t

146 t_attrvalue1_ignore = ''

# attrvalue2

151 def t_attrvalue2_ATTRVALUE2STRING(self, t):
    r'[^"]*'
    t.value = unicode(t.value)
    return t

156 def t_attrvalue2_ATTRVALUE2CLOSE(self, t):
    r'"'
    t.lexer.pop_state()
    return t

161 t_attrvalue2_ignore = ''

# misc

166 literals = '$%^'

def t_ANY_newline(self, t):
    r'\n'
    t.lexer.lineno += len(t.value)
171

# Build the lexer
def build(self, **kwargs):
    self.lexer = lex.lex(object=self, **kwargs)

176 # Test it output
def test(self, data):
    self.lexer.input(data)

```

```

181     _debug_header('LEXER')

    while 1:
        tok = self.lexer.token()
        if not tok: break
186     _debug_print_('LEXER', '%-12s' % (self.lexer.lexstate, tok))

    _debug_footer('LEXER')

    # XmlLexer ends
191

#####
# ESCAPE

196 _xml_escape_table = {
    "&": "&amp;",
    "'": "&quot;",
    '"': "&apos;",
    ">": "&gt;",
201    "<": "&lt;",
    }

def _xml_escape(text):
    L=[]
206    for c in text:
        L.append(_xml_escape_table.get(c,c))
    return "".join(L)

def _xml_unescape(s):
    rules = _xml_escape_table.items()
    rules.reverse()

    for x, y in rules:
        s = s.replace(y, x)
216
    return s

#####
221 # PARSER

tag_stack = []

# Customization
226
def parser_trace(x):
    _debug_print_('PARSER', '%-16s' % (sys._getframe(1).f_code.co_name, x))

def yacc_production_str(p):
231    #return "YaccProduction(%s, %s)" % (str(p.slice), str(p.stack))
    return "YaccP%s" % (str([i.value for i in p.slice]))

yacc.YaccProduction.__str__ = yacc_production_str

236 class ParserError(Exception):
    pass

# Grammar

241 def p_root_element(p):
    '''root : element
        | element CDATA
        , ,
    parser_trace(p)
246
    p[0] = p[1]

```

```

def p_root_cdata_element(p):
    '''root : CDATA element
    251         | CDATA element CDATA
    '''
    parser_trace(p)

    p[0] = p[2]
256
def p_element(p):
    '''element : opentag children closetag
    261         | lonetag
    '''
    parser_trace(p)

    if len(p) == 4:
        p[1].children = p[2]

266    p[0] = p[1]

# tag
def p_opentag(p):
    '''opentag : OPENTAGOPEN TAGATTRNAME attributes TAGCLOSE
    271         '''
    parser_trace(p)

    tag_stack.append(p[2])
    p[0] = DOM.Element(p[2], p[3])
276
def p_closetag(p):
    '''closetag : CLOSETAGOPEN TAGATTRNAME TAGCLOSE
    281         '''
    parser_trace(p)

    n = tag_stack.pop()
    if p[2] != n:
        raise ParserError('Close tag name ("%s") does not match the corresponding
        open tag ("%s").' % (p[2], n))

286
def p_lonetag(p):
    '''lonetag : OPENTAGOPEN TAGATTRNAME attributes LONETAGCLOSE
    291         '''
    parser_trace(p)

    p[0] = DOM.Element(p[2], p[3])

# attr
def p_attributes(p):
    296     '''attributes : attribute attributes
    296         | nothing
    '''
    parser_trace(p)

    if len(p) == 3:
    301         if p[2]:
            p[1].update(p[2])
            p[0] = p[1]
        else:
            p[0] = p[1]
    306     else:
        p[0] = {}

def p_attribute(p):
    311     '''attribute : TAGATTRNAME ATTRASSIGN attrvalue
    311         '''
    parser_trace(p)

    p[0] = {p[1]: p[3]}

```

```

316 def p_attrvalue(p):
    '''attrvalue : ATTRVALUE1OPEN ATTRVALUE1STRING ATTRVALUE1CLOSE
                  | ATTRVALUE2OPEN ATTRVALUE2STRING ATTRVALUE2CLOSE
    '''
    parser_trace(p)
321     p[0] = _xml_unescape(p[2])

# child
def p_children(p):
326     '''children : child children
                  | nothing
    '''
    parser_trace(p)

331     if len(p) > 2:
        if p[2]:
            p[0] = [p[1]] + p[2]
        else:
            p[0] = [p[1]]
336     else:
        p[0] = []

def p_child_element(p):
    '''child : element'''
341     parser_trace(p)

    p[0] = p[1]

def p_child_cdata(p):
346     '''child : CDATA'''
    parser_trace(p)

    p[0] = DOM.Text(p[1])

351 # nothing
def p_nothing(p):
    '''nothing :'''
    pass

356 # Error rule for syntax errors
def p_error(p):
    raise ParserError("Parse error: %s" % (p,))
    pass

361 #####
# DOM

class DOM:
366     class Element:
        # Document object model
        #
        # Parser returns the root element of the XML document

371         def __init__(self, name, attributes={}, children=[]):
            self.name = name
            self.attributes = attributes
            self.children = children

376         def __str__(self):
            attributes_str = ''
            for attr in self.attributes:
                attributes_str += ' %s="%s"' % (attr, _xml_escape(self.attributes[
                    attr]))

381         children_str = ''

```

```

        for child in self.children:
            if isinstance(child, self.__class__):
                children_str += str(child)
            else:
386                 children_str += child

        return '<%s%s>%s</%s>'%(self.name, attributes_str, children_str, self.
            name)

    def __repr__(self):
391         return str(self)

class Text(UserString):
    pass

396 #####
# Methods

def parse(data):
401
    _debug_header('INPUT')
    _debug_print_('INPUT', data)
    _debug_footer('INPUT')

406 # Tokenizer
    xml_lexer = XmlLexer()
    xml_lexer.build()

    xml_lexer.test(data)

411 # Parser
    yacc.yacc(method="SLR")

    _debug_header('PARSER')
416 root = yacc.parse(data)
    _debug_footer('PARSER')

    _debug_header('OUTPUT')
    _debug_print_('OUTPUT', root)
421 _debug_footer('OUTPUT')

    return root

426 def tree(node, level=0, prefix=''):
    'Returns a tree view of the XML data'

    s_node = prefix + node.name + ':'

431 s_children = ''

    children = node.children
    children.reverse()

436 if len(children) == 1 and not 'name' in children[0].__dict__:
    s_node += ' %s' % node.children[0] + '\n'

    else:
        first = True
        for i in xrange(len(children)):
441             if 'name' in node.children[i].__dict__:
                p = ' '
                s_children = tree(node.children[i], level+1, prefix+p) + s_children

446         s_node += '\n'

    return s_node + s_children

```

```

451 #####
# Debug
#
DEBUG = {
456     'INPUT': False,
        'LEXER': False,
        'PARSER': False,
        'OUTPUT': False,
    }
461
def _debug_header(part):
    if DEBUG[part]:
        print '-----'
        print '%s:' % part
466
def _debug_footer(part):
    if DEBUG[part]:
        pass
471
def _debug_print_(part, s):
    if DEBUG[part]:
        print s
476 #####
# MAIN
def main():
    data = open(sys.argv[1]).read()
481    root = parse(data)
    print tree(root)
if __name__ == '__main__':
    main()

```