

Regular Expressions 2







Sandeep Sadanandan (TU, Munich)

э June 5, 2009 1/14

< ∃⇒

Opening URLs

- The module used for opening URLs is urllib2
- The method used is similar to the file open in syntax
- Returns a handler to the URL, which could be used as a handle to a file (readlines, read etc.)

```
1>>> import urllib2
2 >>> r = urllib2.urlopen('http://python.org/')
_{3} >>> html = r.read(300)
4 >>> print html
5 <!DOCTYPE html PUBLIC "-//W3C//DTD XHTML 1.0
6 Transitional//EN" "http://www.w3.org/TR/xhtml1/
7 DTD/xhtml1-transitional.dtd">
8
9
10 <html xmlns="http://www.w3.org/1999/xhtml" xml:</pre>
11 lang="en" lang="en">
12
13 < head>
   <meta http-equiv="content-type"
14
15 COntent="text/html; charset=utf-8" />
16 <title >Python Programming Language ---
17 Official Website</title>
```

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ のQで

General Way

- Not all urls can be opened this way.
- There could be complicated operations such as communicating with he cgi-bin of the server; or some ftp server; etc.
- For that purpose, there are Requests and Opener objects
 - Requests can send along extra data to the server
 - Opener can be used for complicated operations.

1>>> from urllib2 import Request 2 >>> reg = Request('http://www.google.com/') 3 >>> brwser = 'Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.0)' 4 5 >>> reg.add_header('User-Agent', brwser) 6>>> opener = urllib2.build_opener() $\gamma >>>$ opened = opener.open(reg) » >>> print opened.read(150) o <!doctype html><head><meta http-equiv=content-typ</p> 10 Content="text/html; charset=UTF-8"><title>hallo -II Google Search</title><script>window.google={kEl: 12 >>>

On Error?

- In case of errors, one can use the exception to show the error messages
- Two Exceptions which come handy are HTTPError and URLError
- They have to be used in the same order when you write the code. Because HTTPError is a subclass of URLError
- See the example below.

1

- 2 from urllib2 import Request, urlopen, URLError, H
 3 req = Request(someurl)
- 4 <u>try</u>:
- 5 response = urlopen(req)
- ₀ except HTTPError, e:
- print 'The server didn't fulfill the req.'
- 8 print 'Error code: ', e.code

```
    except URLError, e:
```

- ¹⁰ print 'We failed to reach a server.'
- print 'Reason: ', e.reason

12 else:

print 'everything is fine'

E nar

Regular Expressions - A recap

- What are they? A means to find out string patters, To match strings, To find substrings and so forth
- When not to use them? When they are unavoidable. In normal cases where one needs to check whether a string is a substring of another, then is could be easier and more understandable and perhaps more efficient to use the normal string methods.
- When to use them?
 When you know they must be.

A D b 4 A b

Regular Expressions in Theory

- Finite Automata NFA and DFA, Alphabets
- Books on Compilers give a good account of these
- Limitations : (*aⁿbⁿ*), palindromes

Image: Image:

Meta Characters

- If you want to search for ' 'test'', then easy.
- What if you don't know what you want to search for. For example a telephone number? (Which you don't know)
- There are some classes of characters which are dedicated to make the using of regular expressions possible.
- Normal characters match for themselves.
 E.g. t matches t.
- Some special characters don't match themselves.

. ^ \$ * + ? { [] \ | ()

[and] : These can be used to specify a class of characters.

[a-z] : stands for all the lowercase characters. The literal '-' has special meaning inside the square brackets.

[abc\$] stands for the characters 'a', 'b', 'c' and the dollar sign.

Even though \$ has special meaning in RE context, but inside [and]

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ のQで

^ : For negation/complementing a set

[^a-z] means everything which is not lowercase.

\ is perhaps the most important metacharacter.

It is used when a meta-character is to be matched.

= 990

- \d : Every decimal digit. [0-9]
- \D : Everything non-digit; [^0-9]
- \s : Any whitespace; [\t\n\r\f\b]
- \S : Any nonwhitespace character
- \w : Any alpha-numeric; [a-zA-Z0-9_]
- \W : Any non-alpha-numeric-character

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQ@

Importance of DOT

The character "." matches everything but a newline.

Even that can be done using a different mode of the RE module, using re.DOTALL

Repeating Things

- *: ca*t would match ct, cat, caat, caaaat, ...
- +: ca+t would match all of them except for ct
- ? : ca?t would match only ct or cat
- {m,n}: Minimum m times, maximum n times.

 $ca{2,4}t$ would match caat, caaat and caaaat. But not anything else.

< A ▶

Repeating Things

- It is easy to see that * is nothing but {0, }
- Similarly, + is nothing but {1,} and
- ? is {0,1}

A 3 3 4

(日)

- a b matches a or b.
- ^, \$ match the beginning and ending of a line.
- \A, \Z match the beginning and end of a string
- '\A[abc]*\Z' matches all strings
 which are combinations of a, b and c
- \b matches word boundaries: 'class\b' match 'class next Thursday' 'class\b' doesn't match 'classified'

a[bcd]*b against 'abcbd'

a The a in the RE matches.

abcbd The engine matches [bcd]*, going as far as it can, which is to the end of the string.

Failure The engine tries to match b, but the current position is at the end of the string, so it fails.

abcb Back up, so that [bcd]* matches one less character.

Failure Try b again, but the current position is at the last character, which is a "d".

abc Back up again, so that [bcd]* is only matching "bc".

abcb Try b again. This time but the character at the current position is "b", so it succeeds.

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ のQで

Using Them

- Compile them
- Match them
 - match(): Determine if the rematches the string
 - search(): Scan and find the matches
 - findall(): Find all the matches
 - finditer(): Return and iterator
- Use them

1 >>> <u>import</u> re

2>>> p = re.compile('[a-z]+')

з >>> р

- 4 <_sre.SRE_Pattern object at 80c3c28>
- 5 >>> p.match("")
- 6 >>> print p.match("")

7 None

- »>>> m = p.match('tempo')
- ∘>>> <u>print</u> m
- 10 <_sre.SRE_Match object at 80c4f68>

Using Them

- group(): The string matched
- start(): Start of the string
- end(): The End of the string
- span(): A tuple with (start, end)

```
\rightarrow >> m.group()
2 'tempo'
_{3} >>> m. start(), m. end()
4 (0, 5)
_5 >>> m.span()
₀ (0, 5)
>>> print p.match('::: message')
» None
>>> m = p.search('::: message') ; print m
10 < re. MatchObject instance at 80c9650>
11 >>> m.group()
12 'message'
13 >>> m. span()
_{14}(4, 11)
```

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三 のQ@

```
p = re.compile(...)
2 m = p.match( 'string goes here' )
3 if m:
      print 'Match found: ', M.group()
5 <u>else</u>:
6 print 'No match'
s >>> p = re.compile(' \d+')

9 >>> p.findall('12 drummers drumming,

                  11 pipers piping,
10
                  10 lords a-leaping')
11
12 ('12', '11', '10')
13
14 >>> iterator = p.finditer('12 drummers drumming,
                            11 \dots 10 \dots ')
15
_{16} >>>  iterator
17 < callable-iterator object at 0x401833ac>
18 >>> for match in iterator:
                                    ▲□▶▲□▶▲□▶▲□▶ □ のQで
```

<u>print</u> match.span()

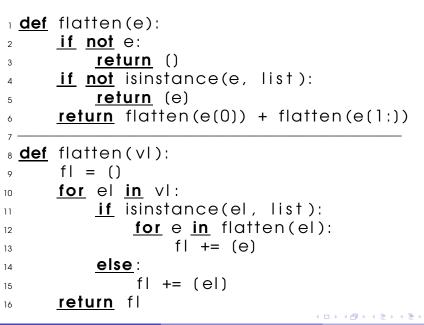
- 20 . . .
- 21 (0, 2)
- 22 (22, 24)
- 23 (29, 31)

▲□▶ ▲□▶ ▲□▶ ▲□▶ = 三 のへで

Tree Deletion

1	def delete(self , item):
2	if item < self.data:
2	if self.lchild:
4	self.lchild = self.lchild.delete(item)
5	<u>elif</u> item > self.data:
6	<u>if</u> self.rchild:
7	<pre>self.rchild = self.rchild.delete(item)</pre>
8	<u>else</u> :
9	<u>if</u> self.isLeaf() :
10	<u>return</u> None
11	<u>if</u> self.singleDad():
12	<u>return</u> self.singleDad()
13	Large = self.leastLarger()
14	self.data, lLarge.data = lLarge.data, self.
15	<pre>self.rchild = self.rchild.delete(item)</pre>
16	<u>return</u> self
	◆□> ◆聞> ◆問> ◆問> 「問」 ひんの

Flatten A List



Flatten A List

```
def reverserec(wh):
           if not wh: return wh
2
           return reverserec(wh(1:)) + wh(0)
3
4
5 def reverserecM(wh):
           if len(wh) == 1: return wh
6
           m = len(wh) / 2
7
           return reverserecM(wh(m;))
8
                   + reverserecM(wh(:m))
9
10
n def reverseStrNotWord(str):
           wh = str.split()
12
           whr = reverserecM(wh)
13
           return ' '.join(whr)
14
Sandeep Sadanandan (TU, Munich)
                                              June 5, 2009
                                                       13/14
```



- Nodes
- Edges
- In Python

・ロト ・ 聞 ト ・ ヨ ト ・ ヨ ト

Crawler/Spider

- Open an URL
- Write a re.query
- Make a graph
- Nodes are pages