

PROGRAMMING IN PYTHON

Gavrilut Dragos Course 8

Regular expression are implemented in Python in library "re".

Usual usage:

- Regular expression (string form) is compiled into a binary form (usually an automata)
- The binary form is used for the following:
 - Checks if a string matches a regular expression
 - Checks if a sub-string from a string can be identified using a regular expression
 - Replace substrings from a string based on a regular expression

Documentation:

• Python 3: https://docs.python.org/3/library/re.html

Regular expression special characters (here is a simple set of special characters)

Character	Match	Character	Match
•	All characters except new line	d/	Decimal characters 0,1,2,3,9
٨	Matches at the start of the string	$\setminus D$	All except decimal characters
\$	Matches at the end of the string	\s	Space, tab, new line (CR/LF) characters
*	>=0 repetition(s)	\ S	All except characters designated by $ackslash$ s
Ś	0 or 1 occurrence	\w	Word characters a-z, A-Z, 0-9 and _
+	>=1 repetition(s)	$\setminus W$	All except characters designated by $ackslash w$
{x}	Matches <x> times</x>	\setminus	Escape character
{x,y}	Matches between <x> and <y> times</y></x>	[^]	Not specified group of characters
[]	Group of characters	()	Grouping
1	Or condition	[]	'-' interval for a group of characters.

Usage:

- o use re.compile (regular_expression_string,flags) to compile a regular expression into its binary form
- Use the "match" method of the resulted object to check if a string matches the regular expression

Python 3.x	
import re	Output
r = re. <i>compile</i> ("07[0-9]{8}")	Match
if r. <i>match</i> ("0740123456"):	
<pre>print("Match")</pre>	

• The same result can be achieved by using the "match" function from the re module directly

```
Python 3.x
import re
if re.match("07[0-9]{8}", "0740123456"):
    print("Match")
```

Pattern	String that will be match
\w+\s+\w+	"Gavrilut Dragos", "Gavrilut Dragos Teodor"
^\w+\s+\w+\$	"Gavrilut Dragos"
[0-9]{1,3}\.[0-9]{1,3}\.[0-9]{1,3}	"192.168.0.1" , "999.999.999.999"
([0-9]{1,3}\.){3}[0-9]{1,3}	"192.168.0.1" , "999.999.999.999"
$ ^{((([0-9]) ([1-9][0-9]) (1[0-9]{2}) (2[0-4][0-9]) (25[0-5])) .){3}(([0-9]) ([1-9][0-9]) (1[0-9]{2}) (2[0-4][0-9]) (25[0-5]))$ } $	Will only match IP addresses
[12]\d{12}	CNP (but will not validate the correctness of the birth date
0x[0-9a-fA-F]+	A hex number
(if then else while continue break)	A special keyword

re.match starts the matching from the beginning of the string and stops once the matching ends and not when the string ends except for the case where regular expression pattern is using the "\$" character:

```
      Python 3.x

      import re
      Output

      if re.match("\d+", "123 USD"):
      Match

      print ("Match")
      NO Match

      if re.match("\d+", "Price is 123 USD"):
      print ("Match")

      if not re.match("\d+$", "123 USD"):
      print ("Match")
```

If you want to check if a regular expression pattern is matching a part of a string, the "search" method can be used:

Python 3.x	
import re	
if re.search("\d+"."Price is 123 USD").	Output
print ("Found")	Found

The same can be achieved using a compiled object:



search method stops after the first match is achieved.

The object returned by the **search** or **match** method is called a match object. A match object is always evaluated to **true**. If the search does not find any match, None is returned and will be evaluated to false. A match object has several members:

- group(index) → returns the substring that matches that specific group. If index is 0, the substring refers to the entire string that was matched by the regular expression
- lastindex → returns the index of the last object that was matched by the regular expression. To create a
 group within the regular expression, one must use (...).

Python 3.x	
import re	
result = re. <i>search(</i> "\d+","Price is 123 USD")	Output
<pre>if result:</pre>	123
<pre>print (result.group(0))</pre>	

In case of some operators (like * or +) they can be preceded by ?. This will specify a NON-greedy behavior.

Python 3.x	
import re	
<pre>result = re.search(".*(\d+)", "File size if 12345 bytes if result: print (result.group(1))</pre>	5")
<pre>result = re.search(".*?(\d+)", "File size if 12345 bytes if result: print (result.group(1))</pre>	5")
	5
	12345

In case of some operators (like * or +) they can be preceded by ?. This will specify a NON-greedy behavior.



Python 3.x		
import re		
<pre>result = re.search("(\d+)[^\d]*(\d+)","Price is 123 USD aprox 110 EUR") if result: print (result.lastindex) for i in range(0, result.lastindex+1):</pre>		
<pre>print (i, "=>", result.group(i))</pre>	Output	
	2 0 => 123 USD, aprox. 110 1 => 123 2 => 110	



Python 3.x	
import re	
<pre>result = re.search("(\d+)[^\d]*(\d+)","Price is if result: print (result.lastindex) for i in range(0, result.lastindex+1):</pre>	s 123 USD aprox 110 EUR")
<pre>print (i, "=>", result.group(i))</pre>	Output
	2 0 => 123 USD, aprox. 110 1 => 123 2 => 110

Python 3.x			
import re			
result = re. <i>search</i> ("((\d+),(\d+))[^\d]*(\d+)", "Color from pixel 20,30 is 123")			
if result:	Output		
print (result. <i>lastindex</i>)			
<pre>for i in range(0, result.lastindex+1):</pre>	4		
<pre>print (i, "=>", result.group(i))</pre>	0 => 20,30 is 123		
	1 => 20,30		
	2 => 20		
	3 => 30		
	4 => 123		





Python 3.x	
import re	
<pre>result = re.search("((\d+),(\d+))[^\d]*(\d+)", "Color from pixel 20(30) is 3</pre>	123")
if result:	Output
print (result. <i>lastindex</i>)	
<pre>for i in range(0, result.lastindex+1):</pre>	4
<pre>print (i, "=>", result.group(i))</pre>	0 => 20,30 is 123
	1 => 20,30
	2 => 20
	3 => 30
	4 => 123

search stops after the first match. To find all substring that match a specific regular expression from a string, use **findall** method.

Python 3.x		
import re		
<pre>result = re.findall("\d+","Color from pixel 20,30 is 123") if result:</pre>		
print (result)	Output	
	['20', '30', '123']	

The result is a vector containing all substrings that matched the regular expression.

search stops after the first match. To find all substring that match a specific regular expression from a string, use **findall** method.

Using groups (...) is also allowed (in this case they will be converted to a tuple in each list element.

Python 3.x	
import re	
<pre>result = re.findall("(\d)(\d+)","Color from pixel 20,30 is 123") if recult.</pre>	
print (result)	

Output

[('2', '0'), ('3', '0'), ('1', '23')]

search stops after the first match. To find all substring that match a specific regular expression from a string, use **findall** method.

Using groups (...) is also allowed (in this case they will be converted to a tuple in each list element.



search stops after the first match. To find all substring that match a specific regular expression from a string, use **findall** method.

Using groups (...) is also allowed (in this case they will be converted to a tuple in each list element.



split method can be used to split a string using a regular expression.

The result is a vector with all elements that substrings that were obtained after the split occurred.

Python 3.x
import re
<pre>result = re.split("[aeiou]+","Color from pixel 20,30 is 123") print (result)</pre>

Outpu	ut													
['C',	'I',	'r fr',	'm p',	'x',	'I 20,30 ',	's 123']								
							- I - I				- T		1	
	1	r	flr	m	n i x	1 2	9 9	א	a	i	s	1	2	ר

Groups can also be used. In this case the split is done after each group that matches.

Python 3.x						
<pre>import re print (re.spl</pre>	<i>it</i> ("∖c	l\d","Colo	or from p	ixel 20,30) is 123"))	
Elements						
'Color from pixel '	· · · /	' is '	'3'			

Python 3.x											
<pre>import re print (re.sp</pre>	olit("	(\d)	(\d)'	',"Co	olor	from	pixel	20,	30 is	123 "))	
Flomonts											
Liemenis											
'Color from pixel	' '2'	'0'	'''	'3'	'0'	' is '	'1'	'2'	'3'		

Groups can also be used. In this case the split is done after each group that matches.

Py	thor	n 3. 2	K																					
in pr	por int	:t	re re	.sp	plit	t("	\d`	\d+	",'	Co	loı	f f	ron	n p	ixel	20	,30	is	5 1	23"	'))			
Οι	Output																							
['	Col	or	fr	ст	pi×	el	י נ	י ,	י נ	' j	ĹS	י ر	'']											

Groups can also be used. In this case the split is done after each group that matches.

Python 3.x	
<pre>import re print (re.split("\d\d+", "12345"))</pre>	<pre>import re print (re.split("\d", "12345"))</pre>
Output	
['', '']	['', '', '', '', '']

Python 3.x
<pre>import re print (re.split("(\d)", "12345"))</pre>
Output
['', '1', '', '2', '', '3', '', '4', '', '5', '']

Groups can also be used. In this case the split is done after each group that matches.

"If capturing parentheses are used in pattern, then the text of all groups in the pattern are also returned as part of the resulting list."

https://docs.python.org/3/library/re.html#re.split

Python 3.x
<pre>import re print (re.split("(\d\d+)", "Color from pixel 20,30 is 123"))</pre>
Output
['Color from pixel ', '20', ',', '30', ' is ', '123', '']

split method also allow flags and to specify how many times a split can be performed. The full format is *split* (*pattern*, *string*, *maxsplit=0*, *flags=0*)

Python 3.x

```
import re
s = "Today I'm having a python course"
print (re.split("[^a-z']+", s))
print (re.split("[^a-z']+", s, 2))
print (re.split("[^a-z']+", s, flags = re.IGNORECASE))
print (re.split("[^a-z']+", s, 2, flags = re.IGNORECASE))
print (re.split("[^a-z'A-Z]+", s))
```

Output

```
[", 'oday', ""m", 'having', 'a', 'python', 'course']
[", 'oday', "'m having a python course"]
['Today', "l'm", 'having', 'a', 'python', 'course']
```

['Today', "I'm", 'having a python course'] ['Today', "I'm", 'having', 'a', 'python', 'course']

Regexp can also be used to replace a matched string with another string using the method **sub**.

format is sub (pattern, replace, string, count=0, flags=0)

- o pattern is a regular expression to search for
- *replace* is either a string or a function
- o *string* is the string where you are going to search the pattern
- o count represents how many time the replacement can occur. If missing or 0 means for all matches.
- flags represents some flags (like re.IGNORECASE)

Python 3.x

import re

```
s = "Today I'm having a python course"
print (re.sub("having\s+a\s+\w+\s+course", "not doing anything", s))
```

Output

Today I'm not doing anything

Regexp can also be used to replace a match with another string using the method **sub**. format is *sub* (*pattern*, *replace*, *string*, *count=0*, *flags=0*)

If **replace** parameter is a string, there is a special operator (<**number**>) that if found within the replacement string will be replace with the group from the match search (for example \setminus 3 will be replaced with .group(3)).

Output

Today I'm not doing the python course

Regexp can also be used to replace a match with another string using the method **sub**. format is *sub* (pattern, replace, string, count=0, flags=0)

If **replace** parameter is a string, there is a special operator (<**number**>) that if found within the replacement string will be replace with the group from the match search (for example \setminus 3 will be replaced with .group(3)).



Regexp can also be used to replace a match with another string using the method **sub**. format is *sub* (*pattern*, *replace*, *string*, *count=0*, *flags=0*)

If **replace** parameter is a function it receives the match object. Usually that function will use .group(0) method to get the string that was matched and convert it to the replacement value.

Python 3.x						
import re						
	Output					
<pre>def ConvertToHex(s):</pre>	File size is 0x3039 bytes					
return hex(int(s.group(0)))						
<pre>s = "File size is 12345 bytes" print (re.sub("\d+",ConvertToHex, s))</pre>						

Python regular expressions supports extensions. The form of the extension is (?...)

• (**?P**<**name**>...) will set the name of a group to a given string. In case of a match, that group can be accessed based on its name.

Python 3.x import re s = "File size if 12345 bytes" result = re.search("(?P<file size>\d+)",s) if result: print ("Size is ",result.group("file_size"))

Output	
Size is 12345	

Python regular expressions supports extensions. The form of the extension is (?...)

 (?P<name>...) → The match object also has a groupdict method that returns a dictionary with all the keys and strings that match the specified regular expression

Python regular expressions supports extensions. The form of the extension is (?...)

- (?i)(...) ignore case will be applied for the current block match
- (?s)(...) "." (dot) will match everything

Python 3.x

Output		
abc		

Python regular expressions supports extensions. The form of the extension is (?...)

- (?=...) will match the previous expression only if next expression is ... (this is called look ahead assertion)
- (?!...) similarly, will match only if the next expression will **not** match ...

```
Python 3.x
import re
s = "Python Course"
result = re.search("(Python)\s+(?=Course)",s)
if result:
    print (result.group(1))
```

Output

Python

Python regular expressions supports extensions. The form of the extension is (?...)

• (?#...) represents a comment / information that can be added in the regular expression to reflect the purpose of a specific group

Python 3.x

```
import re
s = "Size is 1234 bytes"
result = re.search("(?# file size)(\d+)",s)
if result:
    print ("Size is ",result.group(1))
```

Output

Size is 1234

BUILDING A TOKENIZER

Python has a way to iterate through a string applying different regular expression. Because of this, a tokenizer can be built for different languages. Use method **finditer** for this.

Python 3.x	
import re	
number = "(?P <number>\d+)"</number>	Output
<pre>operation = "(?P<operation>[+\-*\/])"</operation></pre>	10 => number
braket = "(?P<braket></braket> [\(\)])"	=> space
<pre>space = "(?P<space>\s)"</space></pre>	* => operation
other = "(?P<other></other> .)"	=> space
<pre>r = re.compile(number+" "+operation+" "+braket+</pre>	(=> braket
" "+space+" "+other)	250 => number
expr = "10 * (250+3)"	+ => operation
<pre>for matchobj in r.finditer(expr):</pre>	3 => number
key = matchobj. <i>lastgroup</i>) => braket
print (matchobj. <i>group</i> (key)+" => "+key)	L

Recommendations:

- 1. If the same regular expression is used multiple times using it in the compile form will improve the performance of the script
- Even if Python recognizes some escape sequences (such as \d or \w) it is better to either use a raw string r"..." or to duplicate the escape character
 Instead of "\d" → use r"\d" or "\\d"
- 3. Regular expression need memory. If all you need is to search a substring within another substring or perform simple string operation, don't use regular expression for this.
- If you are trying to use the regular expression in a portable way, don't use some features like (?P=name) → other languages or regular expression engines might not support this.