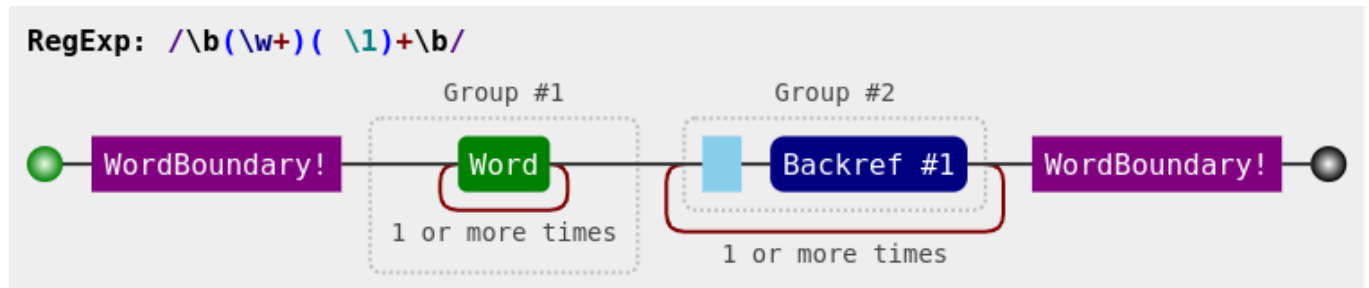


## Python regular expression cheatsheet and examples

2020-07-03



Above diagram created using *Regex*

From docs.python: re:

A regular expression (or RE) specifies a set of strings that matches it; the functions in this module let you check if a particular string matches a given regular expression

This blog post gives an overview and examples of regular expression syntax as implemented by the `re` built-in module (Python 3.8+). Assume ASCII character set unless otherwise specified. This post is an excerpt from my Python `re(gex)?` book.

### Elements that define a regular expression🔗

Anchors	Description
<code>\A</code>	restricts the match to the start of string
<code>\Z</code>	restricts the match to the end of string
<code>^</code>	restricts the match to the start of line
<code>\$</code>	restricts the match to the end of line
<code>\n</code>	newline character is used as line separator
<code>re.MULTILINE</code> or <code>re.M</code>	flag to treat input as multiline string
<code>\b</code>	restricts the match to the start/end of words
	word characters: alphabets, digits, underscore
<code>\B</code>	matches wherever <code>\b</code> doesn't match

`^`, `$` and `\` are metacharacters in the above table, as these characters have special meaning. Prefix a `\` character to remove the special meaning and match such characters literally. For example, `\^` will match a `^` character instead of acting as an anchor.

Feature	Description
	multiple RE combined as conditional OR
	each alternative can have independent anchors
(RE)	group pattern(s), also a capturing group
	<code>a(b c)d</code> is same as <code>abd acd</code>
(?:RE)	non-capturing group
(?P<name>pat)	named capture group
.	Match any character except the newline character <code>\n</code>
[]	Character class, matches one character among many

Greedy Quantifiers	Description
*	Match zero or more times
+	Match one or more times
?	Match zero or one times
{m,n}	Match <code>m</code> to <code>n</code> times (inclusive)
{m,}	Match at least <code>m</code> times
{,n}	Match up to <code>n</code> times (including <code>0</code> times)
{n}	Match exactly <code>n</code> times
pat1.*pat2	any number of characters between <code>pat1</code> and <code>pat2</code>
pat1.*pat2 pat2.*pat1	match both <code>pat1</code> and <code>pat2</code> in any order

Greedy here means that the above quantifiers will match as much as possible that'll also honor the overall RE.

Appending a `?` to greedy quantifiers makes them **non-greedy**, i.e. match as *minimally* as possible. Quantifiers can be applied to literal characters, groups, backreferences and character classes.

Character class	Description
[aeiou]	Match any vowel
[^aeiou]	<code>^</code> inverts selection, so this matches any consonant
[a-f]	<code>-</code> defines a range, so this matches any of abcdef characters
\d	Match a digit, same as <code>[0-9]</code>
\D	Match non-digit, same as <code>[^0-9]</code> or <code>[^\d]</code>
\w	Match word character, same as <code>[a-zA-Z0-9_]</code>

Character class	Description
<code>\W</code>	Match non-word character, same as <code>[^a-zA-Z0-9_]</code> or <code>[^\w]</code>
<code>\s</code>	Match whitespace character, same as <code>[\ \t\n\r\f\v]</code>
<code>\S</code>	Match non-whitespace character, same as <code>[^\ \t\n\r\f\v]</code> or <code>[^\s]</code>

Lookarounds	Description
lookarounds	custom assertions, zero-width like anchors
<code>(?!pat)</code>	negative lookahead assertion
<code>(?&lt;!pat)</code>	negative lookbehind assertion
<code>(?=pat)</code>	positive lookahead assertion
<code>(?&lt;=pat)</code>	positive lookbehind assertion
<code>(?!pat1)(?=pat2)</code>	multiple assertions can be specified in any order
	as they mark a matching location without consuming characters
<code>((?!pat).)*</code>	Negate a grouping, similar to negated character class

Flags	Description
<code>re.IGNORECASE</code> or <code>re.I</code>	flag to ignore case
<code>re.DOTALL</code> or <code>re.S</code>	allow <code>.</code> metacharacter to match newline character
<code>flags=re.S re.I</code>	multiple flags can be combined using <code> </code> operator
<code>re.MULTILINE</code> or <code>re.M</code>	allow <code>^</code> and <code>\$</code> anchors to match line wise
<code>re.VERBOSE</code> or <code>re.X</code>	allows to use literal whitespaces for aligning purposes
	and to add comments after the <code>#</code> character
	escape spaces and <code>#</code> if needed as part of actual RE
<code>re.ASCII</code> or <code>re.A</code>	match only ASCII characters for <code>\b</code> , <code>\w</code> , <code>\d</code> , <code>\s</code>
	and their opposites, applicable only for Unicode patterns
<code>re.LOCALE</code> or <code>re.L</code>	use locale settings for byte patterns and 8-bit locales
<code>(?#comment)</code>	another way to add comments, not a flag
<code>(?flags:pat)</code>	inline flags only for this <code>pat</code> , overrides <code>flags</code> argument
	flags is <code>i</code> for <code>re.I</code> , <code>s</code> for <code>re.S</code> , etc, except <code>L</code> for <code>re.L</code>
<code>(?-flags:pat)</code>	negate flags only for this <code>pat</code>
<code>(?flags-flags:pat)</code>	apply and negate particular flags only for this <code>pat</code>
<code>(?flags)</code>	apply flags for whole RE, can be used only at start of RE
	anchors if any, should be specified after <code>(?flags)</code>

Matched portion	Description
<code>re.Match</code> object	details like matched portions, location, etc
<code>m[0]</code> or <code>m.group(0)</code>	entire matched portion of <code>re.Match</code> object <code>m</code>
<code>m[n]</code> or <code>m.group(n)</code>	matched portion of <i>n</i> th capture group
<code>m.groups()</code>	tuple of all the capture groups' matched portions
<code>m.span()</code>	start and end+1 index of entire matched portion
	pass a number to get span of that particular capture group
	can also use <code>m.start()</code> and <code>m.end()</code>
<code>\N</code>	backreference, gives matched portion of <i>N</i> th capture group
	applies to both search and replacement sections
	possible values: <code>\1</code> , <code>\2</code> up to <code>\99</code> provided no more digits
<code>\g&lt;N&gt;</code>	backreference, gives matched portion of <i>N</i> th capture group
	possible values: <code>\g&lt;0&gt;</code> , <code>\g&lt;1&gt;</code> , etc (not limited to 99)
	<code>\g&lt;0&gt;</code> refers to entire matched portion
<code>(?P&lt;name&gt;pat)</code>	named capture group
	refer as <code>'name'</code> in <code>re.Match</code> object
	refer as <code>(?P=name)</code> in search section
	refer as <code>\g&lt;name&gt;</code> in replacement section
<code>groupdict</code>	method applied on a <code>re.Match</code> object
	gives named capture group portions as a <code>dict</code>

`\0` and `\100` onwards are considered as octal values, hence cannot be used as backreferences.

## re module functions

Function	Description
<code>re.search</code>	Check if given pattern is present anywhere in input string
	Output is a <code>re.Match</code> object, usable in conditional expressions
	r-strings preferred to define RE
	Use byte pattern for byte input
	Python also maintains a small cache of recent RE
<code>re.fullmatch</code>	ensures pattern matches the entire input string
<code>re.compile</code>	Compile a pattern for reuse, outputs <code>re.Pattern</code> object
<code>re.sub</code>	search and replace

Function	Description
<code>re.sub(r'pat', f, s)</code>	function <code>f</code> with <code>re.Match</code> object as argument
<code>re.escape</code>	automatically escape all metacharacters
<code>re.split</code>	split a string based on RE
	text matched by the groups will be part of the output
	portion matched by pattern outside group won't be in output
<code>re.findall</code>	returns all the matches as a list
	if 1 capture group is used, only its matches are returned
	1+, each element will be tuple of capture groups
	portion matched by pattern outside group won't be in output
<code>re.finditer</code>	iterator with <code>re.Match</code> object for each match
<code>re.subn</code>	gives tuple of modified string and number of substitutions

The function definitions are given below:

```

re.search(pattern, string, flags=0)
re.fullmatch(pattern, string, flags=0)
re.compile(pattern, flags=0)
re.sub(pattern, repl, string, count=0, flags=0)
re.escape(pattern)
re.split(pattern, string, maxsplit=0, flags=0)
re.findall(pattern, string, flags=0)
re.finditer(pattern, string, flags=0)
re.subn(pattern, repl, string, count=0, flags=0)

```

## Regular expression examples

As a good practice, always use **raw strings** to construct RE, unless other formats are required. This will avoid clash of special meaning of backslash character between RE and normal quoted strings.

- examples for `re.search`

```

>>> sentence = 'This is a sample string'

# need to load the re module before use
>>> import re
# check if 'sentence' contains the pattern described by RE argument
>>> bool(re.search(r'is', sentence))
True
# ignore case while searching for a match
>>> bool(re.search(r'this', sentence, flags=re.I))
True

```

```
>>> bool(re.search(r'xyz', sentence))
```

```
False
```

```
# re.search output can be directly used in conditional expressions
```

```
>>> if re.search(r'ring', sentence):
```

```
...     print('mission success')
```

```
...
```

```
mission success
```

```
# use raw byte strings if input is of byte data type
```

```
>>> bool(re.search(rb'is', b'This is a sample string'))
```

```
True
```

- difference between string and line anchors

```
# string anchors
```

```
>>> bool(re.search(r'\Ahi', 'hi hello\ntop spot'))
```

```
True
```

```
>>> words = ['surrender', 'up', 'newer', 'do', 'ear', 'eel', 'pest']
```

```
>>> [w for w in words if re.search(r'er\Z', w)]
```

```
['surrender', 'newer']
```

```
# line anchors
```

```
>>> bool(re.search(r'^par$', 'spare\npar\ndare', flags=re.M))
```

```
True
```

- examples for `re.findall`

```
# whole word par with optional s at start and optional e at end
```

```
>>> re.findall(r'\bs?pare?\b', 'par spar apparent spare part pare')
```

```
['par', 'spar', 'spare', 'pare']
```

```
# numbers >= 100 with optional leading zeros
```

```
>>> re.findall(r'\b0*[1-9]\d{2,}\b', '0501 035 154 12 26 98234')
```

```
['0501', '154', '98234']
```

```
# if multiple capturing groups are used, each element of output
```

```
# will be a tuple of strings of all the capture groups
```

```
>>> re.findall(r'([^\s]+)/([^\s/]+),?', '2020/04,1986/Mar')
```

```
[('2020', '04'), ('1986', 'Mar')]
```

```
# normal capture group will hinder ability to get whole match
```

```
# non-capturing group to the rescue
```

```
>>> re.findall(r'\bw*(?:st|in)\b', 'cost akin more east run')
```

```
['cost', 'akin', 'east']
```

```
# useful for debugging purposes as well
```

```
>>> re.findall(r't.*?a', 'that is quite a fabricated tale')
```

```
['tha', 't is quite a', 'ted ta']
```

- examples for `re.split`

# split based on one or more digit characters

```
>>> re.split(r'\d+', 'Sample123string42with777numbers')
['Sample', 'string', 'with', 'numbers']
```

# split based on digit or whitespace characters

```
>>> re.split(r'[\d\s]+', '**1\f2\n3star\t7 77\r**')
['**', 'star', '**']
```

# to include the matching delimiter strings as well in the output

```
>>> re.split(r'(\d+)', 'Sample123string42with777numbers')
['Sample', '123', 'string', '42', 'with', '777', 'numbers']
```

# use non-capturing group if capturing is not needed

```
>>> re.split(r'hand(?:y|ful)', '123handed42handy777handful500')
['123handed42', '777', '500']
```

- backreferencing within search pattern

# whole words that have at least one consecutive repeated character

```
>>> words = ['effort', 'flee', 'facade', 'oddball', 'rat', 'tool']
```

```
>>> [w for w in words if re.search(r'\b\w*(\w)\1\w*\b', w)]
['effort', 'flee', 'oddball', 'tool']
```

- working with matched portions

```
>>> re.search(r'b.*d', 'abc ac adc abbbc')
<re.Match object; span=(1, 9), match='bc ac ad'>
# retrieving entire matched portion, note the use of [0]
>>> re.search(r'b.*d', 'abc ac adc abbbc')[0]
'bc ac ad'
```

# capture group example

```
>>> m = re.search(r'a(.*?)d(.*a)', 'abc ac adc abbbc')
# to get matched portion of second capture group
>>> m[2]
'c a'
# to get a tuple of all the capture groups
>>> m.groups()
('bc ac a', 'c a')
```

- examples for `re.finditer`

# numbers < 350

```
>>> m_iter = re.finditer(r'[0-9]+', '45 349 651 593 4 204')
>>> [m[0] for m in m_iter if int(m[0]) < 350]
['45', '349', '4', '204']
```

```
# start and end+1 index of each matching portion
>>> m_iter = re.finditer(r'ab+c', 'abc ac adc abbbc')
>>> for m in m_iter:
...     print(m.span())
...
(0, 3)
(11, 16)
```

- examples for `re.sub`

```
>>> ip_lines = "catapults\nconcatenate\ncat"
>>> print(re.sub(r'^', r'* ', ip_lines, flags=re.M))
* catapults
* concatenate
* cat
```

```
# replace 'par' only at start of word
>>> re.sub(r'\bpar', r'X', 'par spar apparent spare part')
'X spar apparent spare Xt'
```

```
# same as: r'part|parrot|parent'
>>> re.sub(r'par(en|ro)?t', r'X', 'par part parrot parent')
'par X X X'
```

```
# remove first two columns where : is delimiter
>>> re.sub(r'\A(?:[^\:]+:){2}', r'', 'foo:123:bar:baz', count=1)
'bar:baz'
```

- backreferencing in replacement section

```
# remove consecutive duplicate words separated by space
>>> re.sub(r'\b(\w+)(\1)+\b', r'\1', 'aa a a a 42 f_1 f_1 f_13.14')
'aa a 42 f_1 f_13.14'
```

```
# add something around the matched strings
>>> re.sub(r'\d+', r'(\g<0>)', '52 apples and 31 mangoes')
'(520) apples and (310) mangoes'
```

```
# swap words that are separated by a comma
>>> re.sub(r'(\w+),(\w+)', r'\2,\1', 'good,bad 42,24')
'bad,good 24,42'
```

- using functions in replacement section of `re.sub`

```
>>> from math import factorial
>>> numbers = '1 2 3 4 5'
>>> def fact_num(n):
...     return str(factorial(int(n[0])))
```



```
...
>>> re.sub(r'\d+', fact_num, numbers)
'1 2 6 24 120'

# using lambda
>>> re.sub(r'\d+', lambda m: str(factorial(int(m[0]))), numbers)
'1 2 6 24 120'
```

- examples for lookarounds

```
# change 'foo' only if it is not followed by a digit character
# note that end of string satisfies the given assertion
# foofoo has 2 matches as the assertion doesn't consume characters
>>> re.sub(r'foo(?!\d)', r'baz', 'hey food! foo42 foot5 foofoo')
'hey bazd! foo42 bazt5 bazbaz'
```

```
# change whole word only if it is not preceded by : or -
>>> re.sub(r'(?

```

```
# match digits only if it is preceded by - and followed by ; or :
>>> re.findall(r'(?<=-)\d+(?=[:;])', 'fo-5, ba3; x-83, y-20: f12')
['20']
```

```
# words containing 'b' and 'e' and 't' in any order
>>> words = ['sequoia', 'questionable', 'exhibit', 'equation']
>>> [w for w in words if re.search(r'(?=.*b)(?=.*e).*t', w)]
['questionable', 'exhibit']
```

```
# match if 'do' is not there between 'at' and 'par'
>>> bool(re.search(r'at((?!do).)*par', 'fox,cat,dog,parrot'))
False
# match if 'go' is not there between 'at' and 'par'
>>> bool(re.search(r'at((?!go).)*par', 'fox,cat,dog,parrot'))
True
```

- examples for `re.compile`

Regular expressions can be compiled using `re.compile` function, which gives back a `re.Pattern` object. The top level `re` module functions are all available as methods for this object. Compiling a regular expression helps if the RE has to be used in multiple places or called upon multiple times inside a loop (speed benefit). By default, Python maintains a small list of recently used RE, so the speed benefit doesn't apply for trivial use cases.

```
>>> pet = re.compile(r'dog')
>>> type(pet)
<class 're.Pattern'>
>>> bool(pet.search('They bought a dog'))
True
>>> bool(pet.search('A cat crossed their path'))
```

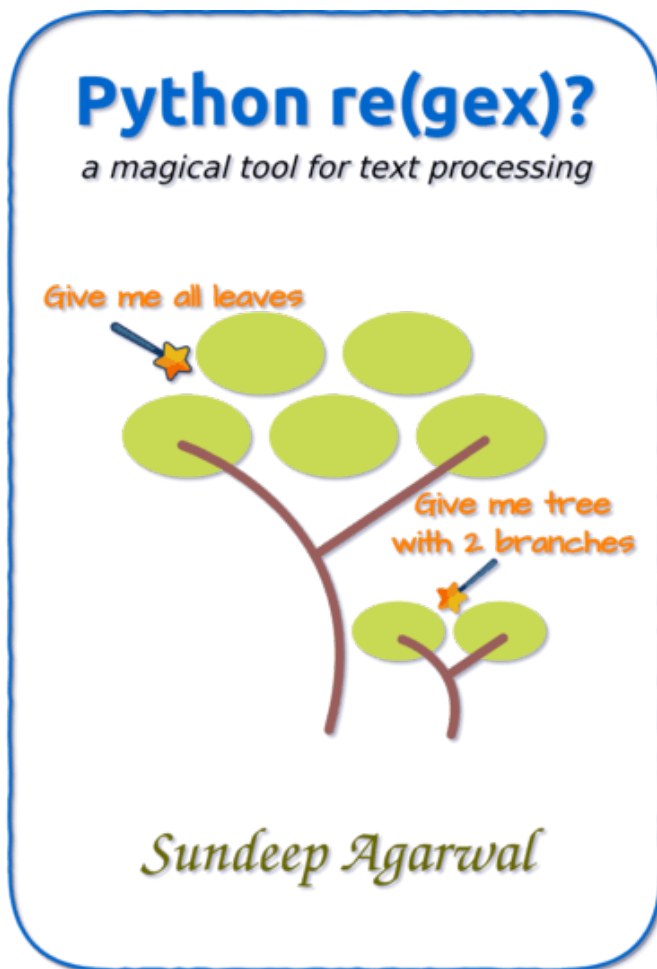
False

```
>>> pat = re.compile(r'\([^)]*\)')
>>> pat.sub('', 'a+b(addition) - foo() + c%d(#modulo)')
'a+b - foo + c%d'
>>> pat.sub('', 'Hi there(greeting). Nice day(a(b))')
'Hi there. Nice day'
```

## Python re(gex)? book

Visit my repo [Python re\(gex\)?](#) for details about the book I wrote on Python regular expressions. The ebook uses plenty of examples to explain the concepts from the very beginning and step by step introduces more advanced concepts. The book also covers the third party module `regex`. The cheatsheet and examples presented in this post are based on contents of this book.

Use this [leanpub](#) link for a discounted price.



#python #regular-expressions #cheatsheet #re-module #examples

◀ **Example driven book on Python regular expressions**

**JavaScript regular expressions cheatsheet and examples** ▶

