

# Table of contents

- [1.](#)
- [2.](#)
- [3.](#)
- [4.](#)
- [5.](#)
- [6.](#)
- [7.](#)
- [8.](#)
- [9.](#)
- [10.](#)
- [11.](#)
- [12.](#)
- [13.](#)
- [14.](#)
- [15.](#)
- [16.](#)
- [17.](#)
- [18.](#)
- [19.](#)
- [20.](#)
- [21.](#)
- [22.](#)
- [23.](#)

---

## 1. Import Module

([go to top](#))

```
In [3]: #import regex module
# see re documentation in last cell

import re

# x = re.compile('enter pattern here')
# y = x.search('enter text to be match')
# z, = y.group() # print out matched group
```

## 2. Overview

[\(go to top\)](#)

```
In [5]: # phone number pattern is all digits()
phoneNumRegex = re.compile(r'\d\d\d-\d\d\d-\d\d\d\d') #raw string
mo = phoneNumRegex.search('My number is 415-555-4242.')
print('Phone number found: ' + mo.group())
```

Phone number found: 415-555-4242

## 3. Grouping w Parentheses

[\(go to top\)](#)

```
In [6]: phoneNumRegex = re.compile(r'(\d\d\d)-(\d\d\d-\d\d\d\d)')
mo = phoneNumRegex.search('My number is 415-555-4242.')
```

```
In [7]: mo.group()
```

```
Out[7]: '415-555-4242'
```

```
In [10]: mo.group(0)
```

```
Out[10]: '415-555-4242'
```

```
In [17]: mo.groups()  
# returns a tuple of multiple values
```

```
Out[17]: ('415', '555-4242')
```

```
In [8]: mo.group(1)
```

```
Out[8]: '415'
```

```
In [9]: mo.group(2)
```

```
Out[9]: '555-4242'
```

```
In [19]: areacode, mainNum = mo.groups()  
print(areacode, 'and', mainNum)
```

```
415 and 555-4242
```

## 4. Match parenthesis in your text

- \ ( and \ )  
([go to top](#))

```
In [20]: phoneNumRegex = re.compile(r'(\d\d\d)-(\d\d\d-\d\d\d\d)')  
mo = phoneNumRegex.search('My number is (415)-555-4242.')
```

```
In [21]: mo.groups()
```

```
Out[21]: ('(415)', '555-4242')
```

```
In [22]: mo.group(1)
```

```
Out[22]: '(415)'
```

```
In [23]: mo.group(2)
```

```
Out[23]: '555-4242'
```

---

## 5. Special Regex Characters

([go to top](#))

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

to use these add a \ in front

---

## 6. Matching groups with the pipe |

- The regular expression `r'Batman|Tina Fey'` will match either 'Batman' or 'Tina Fey'.  
When both Batman and Tina Fey occur in the searched string, the first occurrence of matching text will be returned as the Match object

([go to top](#))

```
In [38]: heroRegex = re.compile(r'Batman|Tina Fey')
```

```
In [39]: mo1 = heroRegex.search('Batman and Tina Fey')
mo1.group()
```

```
Out[39]: 'Batman'
```

```
In [41]: mo2 = heroRegex.search('Tina Fey and Batman')
mo2.group()
```

```
Out[41]: 'Tina Fey'
```

---

---

- You can also use the pipe to match one of several patterns as part of your regex.  
For example, say you wanted to match any of the strings 'Batman', 'Batmobile', 'Batcopter', and 'Batbat'.

```
In [47]: batRegex = re.compile(r'Bat(man|mobile|copter|bat)')
```

```
In [50]: mo = batRegex.search('Batmobile lost a wheel')
```

```
In [51]: mo.group()
```

```
Out[51]: 'Batmobile'
```

```
In [52]: mo.group(1)
```

```
Out[52]: 'mobile'
```

-----  
-----

```
In [44]: mo3 = heroRegex.findall('Tina Fey and Batman')  
mo3
```

```
Out[44]: ['Tina Fey', 'Batman']
```

## 7. Optional Matching with Question Mark

- The ? character flags the group that precedes it as an optional part of the pattern  
That is, the regex should find a match regardless of whether that bit of text is there..

([go to top](#))

```
In [53]: batRegex = re.compile(r'Bat(wo)?man')
```

```
In [54]: mo1 = batRegex.search('The Adventures of Batman')
         mo1.group()
```

```
Out[54]: 'Batman'
```

```
In [58]: mo2 = batRegex.search('The Adventures of Batwoman')
         mo2.group()
```

```
Out[58]: 'Batwoman'
```

-----  
-----

```
In [59]: phoneRegex = re.compile(r'(\d\d\d-)?\d\d\d-\d\d\d\d')
```

```
In [60]: mo1 = phoneRegex.search('My number is 415-555-4242')
         mo1.group()
```

```
Out[60]: '415-555-4242'
```

```
In [62]: mo1.group(1)
```

```
Out[62]: '415-'
```

```
In [61]: mo2 = phoneRegex.search('My number is 555-4242')
         mo2.group()
```

```
Out[61]: '555-4242'
```

## 8. Zero or more matching with star \*

[\(go to top\)](#)

```
In [63]: batRegex = re.compile(r'Bat(wo)*man')
```

```
In [64]: mo1 = batRegex.search('The Adventures of Batman')
         mo1.group()
```

```
Out[64]: 'Batman'
```

```
In [65]: mo2 = batRegex.search('The Adventures of Batwoman')
mo2.group()
```

```
Out[65]: 'Batwoman'
```

```
In [70]: mo2.group(1)
```

```
Out[70]: 'wo'
```

```
In [68]: mo3 = batRegex.search('The Adventures of Batwowowowoman')
mo3.group()
```

```
Out[68]: 'Batwowowowoman'
```

```
In [69]: mo3.group(1)
```

```
Out[69]: 'wo'
```

## 9. One or more matching with plus +

- While \* means “match zero or more,” the + (or plus) means “match one or more.”

Unlike the star, which does not require its group to appear in the matched string, the group preceding a plus must appear at least once.

It is not optional.

[\(go to top\)](#)

```
In [72]: batRegex = re.compile(r'Bat(wo)+man')
```

```
In [77]: mo1 = batRegex.search('The Adventures of Batwoman')
mo1.group()
```

```
Out[77]: 'Batwoman'
```

```
In [76]: mo2 = batRegex.search('The Adventures of Batwowowoman')
mo2.group()
```

```
Out[76]: 'Batwowowoman'
```

- The regex `Bat(wo)+man` will not match the string `'The Adventures of Batman'`, because at least one `wo` is required by the plus sign.

```
In [78]: mo3 = batRegex.search('The Adventures of Batman')
mo3.group()
```

```
-----
AttributeError                                Traceback (most recent call
last)
<ipython-input-78-a63afd945b67> in <module>
      1 mo3 = batRegex.search('The Adventures of Batman')
----> 2 mo3.group()

AttributeError: 'NoneType' object has no attribute 'group'
```

## 10. Match Specific Repetitions with Braces

([go to top](#))

- For example, the regex `(Ha){3}` will match the string `'HaHaHa'`, but it will not match `'HaHa'`, since the latter has only two repeats of the `(Ha)` group. Instead of one number, you can specify a range by writing a minimum, a comma, and a maximum in between the braces. For example, the regex `(Ha){3,5}` will match `'HaHaHa'`, `'HaHaHaHa'`, and `'HaHaHaHaHa'`. You can also leave out the first or second number in the braces to leave the minimum or maximum unbounded. For example, `(Ha){3,}` will match three or more instances of the `(Ha)` group, while `(Ha){,5}` will match zero to five instances. Braces can help make your regular expressions shorter. These two regular expressions match identical patterns:



```
In [81]: haRegex = re.compile(r'(Ha){3}')
```

```
In [83]: mo1 = haRegex.search('HaHaHa')
mo1.group()
```

```
Out[83]: 'HaHaHa'
```

```
In [89]: mo2 = haRegex.search('Ha')
mo2 == None
```

```
Out[89]: True
```

```
In [91]: mo2.group()
```

```
-----
-----
AttributeError                                Traceback (most recent call
last)
<ipython-input-91-fe964b823f27> in <module>
----> 1 mo2.group()

AttributeError: 'NoneType' object has no attribute 'group'
```

## 11. Greedy or nah?

- Since `(Ha){3,5}` can match three, four, or five instances of `Ha` in the string `'HaHaHaHaHa'`, you may wonder why the `Match` object's call to `group()` in the previous brace example returns `'HaHaHaHaHa'` instead of the shorter possibilities. After all, `'HaHaHa'` and `'HaHaHaHa'` are also valid matches of the regular expression `(Ha){3,5}`.

Python's regular expressions are greedy by default, which means that in ambiguous situations they will match the longest string possible. The non-greedy (also called lazy) version of the braces, which matches the shortest string possible, has the closing brace followed by a question mark.

Enter the following into the interactive shell, and notice the difference between the greedy and non-greedy forms of the braces searching the same string:

([go to top](#))

```
In [92]: greedyRegex = re.compile(r'(Ha){3,5}')
```

```
In [93]: mo1 = greedyRegex.search('HaHaHaHaHa')
mo1.group()
```

```
Out[93]: 'HaHaHaHaHa'
```

```
In [94]: mo1.group(1)
```

```
Out[94]: 'Ha'
```

-----

```
In [96]: nonGreedyRegex = re.compile(r'(Ha){3,5}?')
```

```
In [98]: mo2 = nonGreedyRegex.search('HaHaHaHaHa')
mo2.group()
```

```
Out[98]: 'HaHaHa'
```

## 12. findall() method

[\(go to top\)](#)

```
In [115]: phoneNumRegex = re.compile(r'\d\d\d-\d\d\d-\d\d\d\d')
```

```
In [116]: mo = phoneNumRegex.search('Cell: 415-555-999 Work: 212-555-0000')
mo.group()
```

```
Out[116]: '212-555-0000'
```

-----

```
In [117]: phoneNumRegex = re.compile(r'\d\d\d-\d\d\d-\d\d\d\d')
```

```
In [118]: mo1 = phoneNumRegex.findall('Cell: 415-555-9999 Work: 212-555-0000')
mo1
```

```
Out[118]: ['415-555-9999', '212-555-0000']
```

If there are groups in the regular expression, then findall() will return a list of tuples.

Each tuple represents a found match, and its items are the matched strings for each group in the regex.

To see findall() in action, enter the following into the interactive shell (notice that the regular expression being compiled now has groups in parentheses):

```
In [119]: phoneNumRegex = re.compile(r'(\d\d\d)-(\d\d\d)-(\d\d\d\d)')
```

```
In [120]: mo2 = phoneNumRegex.findall('Cell: 415-555-9999 Work: 212-555-0000')
mo2
```

```
Out[120]: [('415', '555', '9999'), ('212', '555', '0000')]
```

```
In [ ]:
```

---

## 13. Character Classes

[\(go to top\)](#)

`\d` Any numeric digit from 0 to 9.

`\D` Any character that is not a numeric digit from 0 to 9.

`\w` Any letter, numeric digit, or the underscore character. (Think of this as matching “word” characters.)

`\W` Any character that is not a letter, numeric digit, or the underscore character.

`\s` Any space, tab, or newline character. (Think of this as matching “space” characters.)

`\S` Any character that is not a space, tab, or newline.`

Character classes are nice for shortening regular expressions. The character class `[0-5]` will match only the numbers 0 to 5; this is much shorter than typing `(0|1|2|3|4|5)`. Note that while `\d` matches digits and `\w` matches digits, letters, and the underscore, there is no shorthand character class that matches only letters. (Though you can use the `[a-zA-Z]` character class, as explained next.)

```
In [124]: xmasRegex = re.compile(r'\d+\s\w+')
xmasRegex.findall('12 drummers, 11 pipers, 10 lords, 9 ladies, 8 maids
['12 drummers', '11 pipers', '10 lords', '9 ladies', '8 maids', '7 swa
```

```
Out[124]: ['12 drummers',
           '11 pipers',
           '10 lords',
           '9 ladies',
           '8 maids',
           '7 swans',
           '6 geese',
           '5 rings',
           '4 birds',
           '3 hens',
           '2 doves',
           '1 partridge']
```

## 14. Make Your own Character Classes

- There are times when you want to match a set of characters but the short-hand character classes (\d, \w, \s, and so on) are too broad.

You can define your own character class using square brackets. For example, the character class [aeiouAEIOU] will match any vowel, both lowercase and uppercase. Enter the following into the interactive shell:

[\(go to top\)](#)

```
In [130]: vowelRegex = re.compile(r'[aeiouAEIOU]')
```

```
In [131]: vowelRegex.findall('Robocop eats baby food')
```

```
Out[131]: ['o', 'o', 'o', 'e', 'a', 'a', 'o', 'o']
```

```
In [135]: vowelRegex2 = re.compile(r'[a-zA-Z0-9]')
```

```
In [137]: vowelRegex2.findall('Robocop eats baby food 345')
```

```
Out[137]: ['R',  
           'o',  
           'b',  
           'o',  
           'c',  
           'o',  
           'p',  
           'e',  
           'a',  
           't',  
           's',  
           'b',  
           'a',  
           'b',  
           'y',  
           'f',  
           'o',  
           'o',  
           'd',  
           '3',  
           '4',  
           '5']
```

- 
- 
- Note that inside the square brackets, the normal regular expression symbols are not interpreted as such. This means you do not need to escape the ., \*, ?, or () characters with a preceding backslash. For example, the character class [0-5.] will match digits 0 to 5 and a period. You do not need to write it as [0-5\\.]. By placing a caret character (^) just after the character class's opening bracket, you can make a negative character class. A negative character class will match all the characters that are not in the character class. For example:

```
In [138]: consonantRegex = re.compile(r'^aeiouAEIOU')  
consonantRegex.findall('RoboCop eats baby food. BABY FOOD.')
```

```
Out[138]: ['R',  
           'b',  
           'C',  
           'p',  
           't',  
           's',  
           'b',  
           'b',  
           'y',  
           'f',  
           'd',  
           'B',  
           'B',  
           'Y',  
           'F',  
           'D',  
           '.']
```

---

## 15. Title

[\(go to top\)](#)

---

## 16. Title

[\(go to top\)](#)

---

## 17. Title

[\(go to top\)](#)

---

## 18. Title

[\(go to top\)](#)

---

## 19. Title

[\(go to top\)](#)

---

## 20. Title

[\(go to top\)](#)

---

---



In [ ]:

In [ ]:

Support for regular expressions (RE).

This module provides regular expression matching operations similar to those found in Perl. It supports both 8-bit and Unicode strings; both the pattern and the strings being processed can contain null bytes and characters outside the US ASCII range.

Regular expressions can contain both special and ordinary characters.

Most ordinary characters, like "A", "a", or "0", are the simplest regular expressions; they simply match themselves. You can concatenate ordinary characters, so `last` matches the string `'last'`.

The special characters are:

<code>"."</code>	Matches any character except a newline.
<code>"^"</code>	Matches the start of the string.
<code>"\$"</code>	Matches the end of the string or just before the newline at the end of the string.
<code>"*"</code>	Matches 0 or more (greedy) repetitions of the preceding RE. Greedy means that it will match as many repetitions as possible.
<code>"+"</code>	Matches 1 or more (greedy) repetitions of the preceding RE.
<code>"?"</code>	Matches 0 or 1 (greedy) of the preceding RE.
<code>*? , +? , ??</code>	Non-greedy versions of the previous three special characters.
<code>{m,n}</code>	Matches from m to n repetitions of the preceding

RE.

`{m,n}?` Non-greedy version of the above.

`"\"` Either escapes special characters or signals a special sequence.

`[]` Indicates a set of characters.

A `"^"` as the first character indicates a complementing set.

`"|"` `A|B`, creates an RE that will match either A or B.

`(...)` Matches the RE inside the parentheses.

The contents can be retrieved or matched later in the string.

`(?aiLmsux)` The letters set the corresponding flags defined below.

`(?:...)` Non-grouping version of regular parentheses.

`(?P<name>...)` The substring matched by the group is accessible by name.

`(?P=name)` Matches the text matched earlier by the group named name.

`(?#...)` A comment; ignored.

`(?=...)` Matches if ... matches next, but doesn't consume the string.

`(?!...)` Matches if ... doesn't match next.

`(?<=...)` Matches if preceded by ... (must be fixed length).

`(?<!...)` Matches if not preceded by ... (must be fixed length).

`(?(id/name)yes|no)` Matches yes pattern if the group with id/name matched,

the (optional) no pattern otherwise.

The special sequences consist of `"\"` and a character from the list

below. If the ordinary character is not on the list, then the resulting RE will match the second character.

`\number` Matches the contents of the group of the same number.

`\A` Matches only at the start of the string.

`\Z` Matches only at the end of the string.

`\b` Matches the empty string, but only at the start or end of a word.

`\B` Matches the empty string, but not at the start or end of a word.

`\d` Matches any decimal digit; equivalent to the set `[0-9]` in bytes patterns or string patterns with the ASCII flag. In string patterns without the ASCII flag, it will match the whole range of Unicode digits.

`\D` Matches any non-digit character; equivalent to `[\^ \d]`.

`\s` Matches any whitespace character; equivalent to `[\t\n\r\f\v]` in bytes patterns or string patterns with the ASCII flag. In string patterns without the ASCII flag, it will match the whole range of Unicode whitespace characters.

`\S` Matches any non-whitespace character; equivalent to `[\^ \s]`.

`\w` Matches any alphanumeric character; equivalent to `[a-zA-Z0-9_]` in bytes patterns or string patterns with the ASCII flag. In string patterns without the ASCII flag, it will match the range of Unicode alphanumeric characters (letters plus digits plus underscore). With LOCALE, it will match the set `[0-9_]` plus characters defined as letters for the current locale.

`\W` Matches the complement of `\w`.

`\\` Matches a literal backslash.

This module exports the following functions:

`match` Match a regular expression pattern to the beginning of a string.

`fullmatch` Match a regular expression pattern to all of a string.

`search` Search a string for the presence of a pattern.

`sub` Substitute occurrences of a pattern found in a string.

`subn` Same as `sub`, but also return the number of substitutions made.

stitutions made.

split	Split a string by the occurrences of a pattern.
findall	Find all occurrences of a pattern in a string.
finditer	Return an iterator yielding a Match object for each match.
compile	Compile a pattern into a Pattern object.
purge	Clear the regular expression cache.
escape	Backslash all non-alphanumerics in a string.

Each function other than purge and escape can take an optional 'flags' argument

consisting of one or more of the following module constants, joined by "|".

A, L, and U are mutually exclusive.

A	ASCII	For string patterns, make \w, \W, \b, \B, \d, \D match the corresponding ASCII character categories (rather than the whole Unicode categories, which is the default).
---	-------	---

For bytes patterns, this flag is the only available behaviour and needn't be specified.

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	For compatibility only. Ignored for string patterns (it is the default), and forbidden for bytes patterns.
---	---------	--

