

Meta characters in regular expressions

Meta character	Description
\	Marks the next character as either a special character or a literal. For example, <code>n</code> matches the character <i>n</i> , whereas <code>\n</code> matches a newline character. The sequence <code>\\</code> matches <code>\</code> and <code>\(</code> matches <code>(</code> .
^	Matches the beginning of input.
\$	Matches the end of input.
*	Matches the preceding character zero or more times. For example, <code>zo*</code> matches either <i>z</i> or <i>zoo</i> .
+	Matches the preceding character one or more times. For example, <code>zo+</code> matches <i>zoo</i> but not <i>z</i> .
?	Matches the preceding character zero or one time. For example, <code>a?ve?</code> matches the <i>ve</i> in <i>never</i> .
.	Matches any single character except a newline character.
(pattern)	Matches a pattern and remembers the match. The matched substring can be retrieved from the resulting matches collection by using this code: <code>Item [0] . . . [n]</code> . To match parentheses characters <code>()</code> , use <code>\(</code> or <code>\)</code> .
<code>x y</code>	Matches either <i>x</i> or <i>y</i> . For example, <code>z wood</code> matches <i>z</i> or <i>wood</i> . <code>(z w)oo</code> matches <i>zoo</i> or <i>wood</i> .
<code>{n}</code>	<i>n</i> is a non-negative integer. Matches exactly <i>n</i> times. For example, <code>o{2}</code> does not match the <i>o</i> in <i>Bob</i> , but matches the first two <i>os</i> in <i>fooooood</i> .
<code>{n,}</code>	In this expression, <i>n</i> is a non-negative integer. Matches the preceding character at least <i>n</i> times. For example, <code>o{2,}</code> does not match the <i>o</i> in <i>Bob</i> and matches all the <i>os</i> in <i>fooooood</i> . The <code>o{1,}</code> expression is equivalent to <code>o+</code> and <code>o{0,}</code> is equivalent to <code>o*</code> .
<code>{n,m}</code>	The <i>m</i> and <i>n</i> variables are non-negative integers. Matches the preceding character at least <i>n</i> and at most <i>m</i> times. For example, <code>o{1,3}</code> matches the first three <i>os</i> in <i>fooooood</i> . The <code>o{0,1}</code> expression is equivalent to <code>o?</code> .
<code>[xyz]</code>	A character set. Matches any one of the enclosed characters. For example, <code>[abc]</code> matches the <i>a</i> in <i>plain</i> .
<code>[^xyz]</code>	A negative character set. Matches any character that is not enclosed. For example, <code>[^abc]</code> matches the <i>p</i> in <i>plain</i> .

Meta character	Description
<code>[a-z]</code>	A range of characters. Matches any character in the specified range. For example, <code>[a-z]</code> matches any lowercase alphabetic character in the English alphabet.
<code>[^m-z]</code>	A negative range of characters. Matches any character that is not in the specified range. For example, <code>[m-z]</code> matches any character that is not in the range <i>m</i> through <i>z</i> .
<code>\A</code>	Matches only at beginning of a string.
<code>\b</code>	Matches a word boundary, that is, the position between a word and a space. For example, <code>er\b</code> matches the <i>er</i> in <i>never</i> but not the <i>er</i> in <i>verb</i> .
<code>\B</code>	Matches a nonword boundary. The <code>ea*r\b</code> expression matches the <i>ear</i> in <i>never early</i> .
<code>\d</code>	Matches a digit character.
<code>\D</code>	Matches a non-digit character.
<code>\f</code>	Matches a form-feed character.
<code>\n</code>	Matches a newline character.
<code>\r</code>	Matches a carriage return character.
<code>\s</code>	Matches any white space including spaces, tabs, form-feed characters, and so on.
<code>\S</code>	Matches any non-white space character.
<code>\t</code>	Matches a tab character.
<code>\v</code>	Matches a vertical tab character.
<code>\w</code>	Matches any word character including underscore. This expression is equivalent to <code>[A-Za-z0-9_]</code> .
<code>\W</code>	Matches any non-word character. This expression is equivalent to <code>[^A-Za-z0-9_]</code> .
<code>\z</code>	Matches only the end of a string.
<code>\Z</code>	Matches only the end of a string, or before a newline character at the end.

Python RegEx

In this tutorial, you will learn about regular expressions (RegEx), and use Python's `re` module to work with RegEx (with the help of examples).

A **Regular Expression** (RegEx) is a sequence of characters that defines a search pattern. For example,

```
^a...s$
```

The above code defines a RegEx pattern. The pattern is: **any five letter string starting with `a` and ending with `s`.**

A pattern defined using RegEx can be used to match against a string.

Expression	String	Matched?
	abs	No match
	alias	Match
<code>^a...s\$</code>	abyss	Match
	Alias	No match
	An abacus	No match

Python has a module named `re` to work with RegEx. Here's an example:

```
import re

pattern = '^a...s$'
test_string = 'abyss'
result = re.match(pattern, test_string)

if result:
    print("Search successful.")
else:
    print("Search unsuccessful.")
```

[Run Code](#)

Here, we used `re.match()` function to search `pattern` within the `test_string`. The method returns a match object if the search is successful. If not, it returns `None`.

There are other several functions defined in the `re` module to work with RegEx. Before we explore that, let's learn about regular expressions themselves.

If you already know the basics of RegEx, jump to [Python RegEx](#).

Specify Pattern Using RegEx

To specify regular expressions, metacharacters are used. In the above example, `^` and `$` are metacharacters.

MetaCharacters

Metacharacters are characters that are interpreted in a special way by a RegEx engine. Here's a list of metacharacters:

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

`[]` - Square brackets

Square brackets specifies a set of characters you wish to match.

Expression	String	Matched?
<code>[abc]</code>	<code>a</code>	1 match

Expression	String	Matched?
	<code>ac</code>	2 matches
	<code>Hey Jude</code>	No match
	<code>abc de ca</code>	5 matches

Here, `[abc]` will match if the string you are trying to match contains any of the `a`, `b` or `c`.

You can also specify a range of characters using `-` inside square brackets.

- `[a-e]` is the same as `[abcde]`.
- `[1-4]` is the same as `[1234]`.
- `[0-39]` is the same as `[01239]`.

You can complement (invert) the character set by using caret `^` symbol at the start of a square-bracket.

- `[^abc]` means any character except `a` or `b` or `c`.
- `[^0-9]` means any non-digit character.

`.` - Period

A period matches any single character (except newline `'\n'`).

Expression	String	Matched?
<code>..</code>	<code>a</code>	No match
	<code>ac</code>	1 match
	<code>acd</code>	1 match
	<code>acde</code>	2 matches (contains 4 characters)

`^` - Caret

The caret symbol `^` is used to check if a string **starts with** a certain character.

Expression	String	Matched?
<code>^a</code>	<code>a</code>	1 match
	<code>abc</code>	1 match
	<code>bac</code>	No match
<code>^ab</code>	<code>abc</code>	1 match
	<code>acb</code>	No match (starts with <code>a</code> but not followed by <code>b</code>)

`$` - Dollar

The dollar symbol `$` is used to check if a string **ends with** a certain character.

Expression	String	Matched?
<code>a\$</code>	<code>a</code>	1 match
	<code>formula</code>	1 match
	<code>cab</code>	No match

`*` - Star

The star symbol `*` matches **zero or more occurrences** of the pattern left to it.

Expression	String	Matched?
<code>ma*n</code>	<code>mn</code>	1 match
	<code>man</code>	1 match
	<code>maaan</code>	1 match
	<code>main</code>	No match (<code>a</code> is not followed by <code>n</code>)
	<code>woman</code>	1 match

`+` - Plus

The plus symbol `+` matches **one or more occurrences** of the pattern left to it.

Expression	String	Matched?
<code>ma+n</code>	<code>mn</code>	No match (no <code>a</code> character)
	<code>man</code>	1 match
	<code>maaan</code>	1 match
	<code>main</code>	No match (<code>a</code> is not followed by <code>n</code>)
	<code>woman</code>	1 match

`?` - Question Mark

The question mark symbol `?` matches **zero or one occurrence** of the pattern left to it.

Expression	String	Matched?
<code>ma?n</code>	<code>mn</code>	1 match
	<code>man</code>	1 match
	<code>maaan</code>	No match (more than one <code>a</code> character)
	<code>main</code>	No match (a is not followed by n)
	<code>woman</code>	1 match

`{}` - Braces

Consider this code: `{n,m}`. This means at least `n`, and at most `m` repetitions of the pattern left to it.

Expression	String	Matched?
<code>a{2,3}</code>	<code>abc dat</code>	No match
	<code>abc daat</code>	1 match (at <code>daat</code>)
	<code>aabc daaat</code>	2 matches (at <code>aabc</code> and <code>daaat</code>)
	<code>aabc daaaat</code>	2 matches (at <code>aabc</code> and <code>daaaat</code>)

Let's try one more example. This RegEx `[0-9]{2,4}` matches at least 2 digits but not more than 4 digits

Expression	String	Matched?
<code>[0-9]{2,4}</code>	<code>ab123csde</code>	1 match (match at <code>ab123csde</code>)
	<code>12 and 345673</code>	3 matches (<code>12</code> , <code>3456</code> , <code>73</code>)

Expression	String	Matched?
	1 and 2	No match

| - Alternation

Vertical bar `|` is used for alternation (or operator).

Expression	String	Matched?
	cde	No match
a b	ade	1 match (match at <u>a</u> de)
	acdbea	3 matches (at <u>a</u> <u>c</u> <u>d</u> <u>b</u> <u>e</u> <u>a</u>)

Here, `a|b` match any string that contains either `a` or `b`

() - Group

Parentheses `()` is used to group sub-patterns. For

example, `(a|b|c)xz` match any string that matches either `a` or `b` or `c` followed by `xz`

Expression	String	Matched?
	ab xz	No match
(a b c)xz	abxz	1 match (match at <u>a</u> <u>b</u> <u>x</u> <u>z</u>)
	axz cabxz	2 matches (at <u>a</u> <u>x</u> <u>z</u> <u>b</u> <u>c</u> <u>a</u> <u>b</u> <u>x</u> <u>z</u>)

`\` - Backslash

Backslash `\` is used to escape various characters including all metacharacters. For example,

`\$a` match if a string contains `$` followed by `a`. Here, `$` is not interpreted by a RegEx engine in a special way.

If you are unsure if a character has special meaning or not, you can put `\` in front of it. This makes sure the character is not treated in a special way.

Special Sequences

Special sequences make commonly used patterns easier to write. Here's a list of special sequences:

`\A` - Matches if the specified characters are at the start of a string.

Expression	String	Matched?
<code>\Athe</code>	the sun	Match
	In the sun	No match

`\b` - Matches if the specified characters are at the beginning or end of a word.

Expression	String	Matched?
<code>\bfoo</code>	football	Match
	a football	Match
	afootball	No match

Expression	String	Matched?
	the foo	Match
foo\b	the afoo test	Match
	the afootest	No match

`\B` - Opposite of `\b`. Matches if the specified characters are **not** at the beginning or end of a word.

Expression	String	Matched?
	football	No match
\Bfoo	a football	No match
	afootball	Match
	the foo	No match
foo\B	the afoo test	No match
	the afootest	Match

`\d` - Matches any decimal digit. Equivalent to `[0-9]`

Expression	String	Matched?
	12abc3	3 matches (at <u>1</u> <u>2</u> abc <u>3</u>)
\d	Python	No match

`\D` - Matches any non-decimal digit. Equivalent to `[^0-9]`

Expression	String	Matched?
<code>\D</code>	1ab34"50	3 matches (at 1ab34"50)
	1345	No match

`\s` - Matches where a string contains any whitespace character. Equivalent to `[\t\n\r\f\v]`.

Expression	String	Matched?
<code>\s</code>	Python RegEx	1 match
	PythonRegEx	No match

`\S` - Matches where a string contains any non-whitespace character. Equivalent to `[^\t\n\r\f\v]`.

Expression	String	Matched?
<code>\S</code>	a b	2 matches (at a b)
		No match

`\w` - Matches any alphanumeric character (digits and alphabets). Equivalent to `[a-zA-Z0-9_]`. By the way, underscore `_` is also considered an alphanumeric character.

Expression	String	Matched?
<code>\w</code>	<code>12&" : ;c</code>	3 matches (at <code>12&" : ;c</code>)
	<code>%"> !</code>	No match

`\W` - Matches any non-alphanumeric character. Equivalent to `[^a-zA-Z0-9_]`

Expression	String	Matched?
<code>\W</code>	<code>1a2%c</code>	1 match (at <code>1a2%c</code>)
	<code>Python</code>	No match

`\Z` - Matches if the specified characters are at the end of a string.

Expression	String	Matched?
<code>Python\Z</code>	<code>I like Python</code>	1 match
	<code>I like Python Programming</code>	No match
	<code>Python is fun.</code>	No match

Tip: To build and test regular expressions, you can use RegEx tester tools such as [regex101](https://regex101.com/). This tool not only helps you in creating regular expressions, but it also helps you learn it.

Now you understand the basics of RegEx, let's discuss how to use RegEx in your Python code.

Python RegEx

Python has a module named `re` to work with regular expressions. To use it, we need to import the module.

```
import re
```

The module defines several functions and constants to work with RegEx.

re.findall()

The `re.findall()` method returns a list of strings containing all matches.

Example 1: re.findall()

```
# Program to extract numbers from a string

import re

string = 'hello 12 hi 89. Howdy 34'
pattern = '\d+'

result = re.findall(pattern, string)
```

```
print(result)
```

```
# Output: ['12', '89', '34']
```

If the pattern is not found, `re.findall()` returns an empty list.

re.split()

The `re.split()` method splits the string where there is a match and returns a list of strings where the splits have occurred.

Example 2: re.split()

```
import re

string = 'Twelve:12 Eighty nine:89.'
pattern = '\d+'

result = re.split(pattern, string)
print(result)

# Output: ['Twelve:', ' Eighty nine:', '.']
Run Code
```

If the pattern is not found, `re.split()` returns a list containing the original string.

You can pass `maxsplit` argument to the `re.split()` method. It's the maximum number of splits that will occur.

```
import re

string = 'Twelve:12 Eighty nine:89 Nine:9.'
pattern = '\d+'

# maxsplit = 1
# split only at the first occurrence
result = re.split(pattern, string, 1)
print(result)

# Output: ['Twelve:', ' Eighty nine:89 Nine:9.']
Run Code
```

By the way, the default value of `maxsplit` is 0; meaning all possible splits.

re.sub()

The syntax of `re.sub()` is:

```
re.sub(pattern, replace, string)
```

The method returns a string where matched occurrences are replaced with the content of `replace` variable.

Example 3: re.sub()

```
# Program to remove all whitespaces
import re

# multiline string
string = 'abc 12\
de 23 \n f45 6'

# matches all whitespace characters
pattern = '\s+'
```



```
# empty string
replace = ''

new_string = re.sub(pattern, replace, string)
print(new_string)
```

```
# Output: abc12de23f456
Run Code
```

If the pattern is not found, `re.sub()` returns the original string.

You can pass `count` as a fourth parameter to the `re.sub()` method. If omitted, it results to 0. This will replace all occurrences.

```
import re

# multiline string
string = 'abc 12\
de 23 \n f45 6'

# matches all whitespace characters
pattern = '\s+'
replace = ''

new_string = re.sub(r'\s+', replace, string, 1)
print(new_string)

# Output:
# abc12de 23
# f45 6
```

re.subn()

The `re.subn()` is similar to `re.sub()` except it returns a tuple of 2 items containing the new string and the number of substitutions made.

Example 4: re.subn()

```
# Program to remove all whitespaces
import re

# multiline string
string = 'abc 12\
de 23 \n f45 6'

# matches all whitespace characters
pattern = '\s+'

# empty string
replace = ''

new_string = re.subn(pattern, replace, string)
print(new_string)

# Output: ('abc12de23f456', 4)
Run Code
```

re.search()

The `re.search()` method takes two arguments: a pattern and a string. The method looks for the first location where the RegEx pattern produces a match with the string.

If the search is successful, `re.search()` returns a match object; if not, it returns `None`.

```
match = re.search(pattern, str)
```

Example 5: re.search()

```
import re

string = "Python is fun"

# check if 'Python' is at the beginning
match = re.search('\APython', string)

if match:
    print("pattern found inside the string")
else:
    print("pattern not found")

# Output: pattern found inside the string
Run Code
```

Here, `match` contains a match object.

Match object

You can get methods and attributes of a match object using [dir\(\)](#) function. Some of the commonly used methods and attributes of match objects are:

match.group()

The `group()` method returns the part of the string where there is a match.

Example 6: Match object

```
import re

string = '39801 356, 2102 1111'
```

```
# Three digit number followed by space followed by two digit number
pattern = '(\d{3}) (\d{2})'

# match variable contains a Match object.
match = re.search(pattern, string)

if match:
    print(match.group())
else:
    print("pattern not found")

# Output: 801 35
Run Code
```

Here, `match` variable contains a match object.

Our pattern `(\d{3}) (\d{2})` has two subgroups `(\d{3})` and `(\d{2})`. You can get the part of the string of these parenthesized subgroups. Here's how:

```
>>> match.group(1)
'801'

>>> match.group(2)
'35'

>>> match.group(1, 2)
('801', '35')

>>> match.groups()
('801', '35')
```

`match.start()`, `match.end()` and `match.span()`

The `start()` function returns the index of the start of the matched substring. Similarly, `end()` returns the end index of the matched substring.

```
>>> match.start()
2

>>> match.end()
```

The `span()` function returns a tuple containing start and end index of the matched part.

```
>>> match.span()
(2, 8)
```

match.re and match.string

The `re` attribute of a matched object returns a regular expression object. Similarly, `string` attribute returns the passed string.

```
>>> match.re
re.compile('(\d{3}) (\d{2})')

>>> match.string
'39801 356, 2102 1111'
```

We have covered all commonly used methods defined in the `re` module. If you want to learn more, visit [Python 3 re module](#).

Using r prefix before RegEx

When `r` or `R` prefix is used before a regular expression, it means raw string. For example, `'\n'` is a new line whereas `r'\n'` means two characters: a backslash `\` followed by `n`.

Backslash `\` is used to escape various characters including all metacharacters. However, using `r` prefix makes `\` treat as a normal character.

Example 7: Raw string using `r` prefix

```
import re

string = '\n and \r are escape sequences.'

result = re.findall(r'[\n\r]', string)
print(result)

# Output: ['\n', '\r']
```