```
Q
        3.12.0
                                                                                    Go
        self.maps[0][key] = value
         _delitem__(self, key):
        for mapping in self.maps:
            if key in mapping:
                del mapping[key]
                return
        raise KeyError(key)
>>> d = DeepChainMap({'zebra': 'black'}, {'elephant': 'blue'}, {'lion': 'yellow'})
>>> d['lion'] = 'orange'
                                # update an existing key two levels down
>>> d['snake'] = 'red'
                                 # new keys get added to the topmost dict
>>> del d['elephant']
                                 # remove an existing key one level down
                                 # display result
>>> d
DeepChainMap({'zebra': 'black', 'snake': 'red'}, {}, {'lion': 'orange'})
```

Counter objects

A counter tool is provided to support convenient and rapid tallies. For example:

```
class collections.Counter([iterable-or-mapping])
```

A Counter is a dict subclass for counting hashable objects. It is a collection where elements are stored as dictionary keys and their counts are stored as dictionary values. Counts are allowed to be any integer value including zero or negative counts. The Counter class is similar to bags or multisets in other languages.

Elements are counted from an iterable or initialized from another mapping (or counter):

Counter objects have a dictionary interface except that they return a zero count for missing items instead of raising a KeyError:



Setting a count to zero does not remove an element from a counter. Use del to remove it entirely:

```
>>> c['sausage'] = 0  # counter entry with a zero counter entry with a
```

New in version 3.1.

Changed in version 3.7: As a dict subclass, Counter inherited the capability to remember insertion order. Math operations on *Counter* objects also preserve order. Results are ordered according to when an element is first encountered in the left operand and then by the order encountered in the right operand.

Counter objects support additional methods beyond those available for all dictionaries:

elements()

Return an iterator over elements repeating each as many times as its count. Elements are returned in the order first encountered. If an element's count is less than one, elements() will ignore it.

```
>>> c = Counter(a=4, b=2, c=0, d=-2)
>>> sorted(c.elements())
['a', 'a', 'a', 'b', 'b']
```

most_common([n])

Return a list of the n most common elements and their counts from the most common to the least. If n is omitted or None, $most_common()$ returns all elements in the counter. Elements with equal counts are ordered in the order first encountered:

```
>>> Counter('abracadabra').most_common(3)
[('a', 5), ('b', 2), ('r', 2)]
```

subtract([iterable-or-mapping])

Elements are subtracted from an *iterable* or from another *mapping* (or counter). Like dict.update() but subtracts counts instead of replacing them. Both inputs and outputs may be zero or negative.

```
>>> c = Counter(a=4, b=2, c=0, d=-2)
>>> d = Counter(a=1, b=2, c=3, d=4)
>>> c.subtract(d)
>>> c
Counter({'a': 3, 'b': 0, 'c': -3, 'd': -6})
```

New in version 3.2.

total()

New in version 3.10.

The usual dictionary methods are available for Counter objects except for two which work differently for counters.

fromkeys(iterable)

This class method is not implemented for Counter objects.

```
update([iterable-or-mapping])
```

Elements are counted from an *iterable* or added-in from another *mapping* (or counter). Like dict.update() but adds counts instead of replacing them. Also, the *iterable* is expected to be a sequence of elements, not a sequence of (key, value) pairs.

Counters support rich comparison operators for equality, subset, and superset relationships: ==, !=, <, <=, >, >=. All of those tests treat missing elements as having zero counts so that Counter(a=1) == Counter(a=1, b=0) returns true.

New in version 3.10: Rich comparison operations were added.

Changed in version 3.10: In equality tests, missing elements are treated as having zero counts. Formerly, Counter(a=3) and Counter(a=3, b=0) were considered distinct.

Common patterns for working with Counter objects:

```
c.total()
                                 # total of all counts
c.clear()
                                 # reset all counts
list(c)
                                 # list unique elements
set(c)
                                 # convert to a set
dict(c)
                                 # convert to a regular dictionary
c.items()
                                 # convert to a list of (elem, cnt) pairs
Counter(dict(list of pairs))
                                # convert from a list of (elem, cnt) pairs
c.most_common()[:-n-1:-1]
                                 # n least common elements
                                 # remove zero and negative counts
+c
```

Several mathematical operations are provided for combining Counter objects to produce multisets (counters that have counts greater than zero). Addition and subtraction combine counters by adding or subtracting the counts of corresponding elements. Intersection and union return the minimum and maximum of corresponding counts. Equality and inclusion compare corresponding counts. Each operation can accept inputs with signed counts, but the output will exclude results with counts of zero or less.

Unary addition and subtraction are shortcuts for adding an empty counter or subtracting from an empty counter.

```
>>> c = Counter(a=2, b=-4)
>>> +c
Counter({'a': 2})
>>> -c
Counter({'b': 4})
```

New in version 3.3: Added support for unary plus, unary minus, and in-place multiset operations.

Note: Counters were primarily designed to work with positive integers to represent running counts; however, care was taken to not unnecessarily preclude use cases needing other types or negative values. To help with those use cases, this section documents the minimum range and type restrictions.

- The Counter class itself is a dictionary subclass with no restrictions on its keys and values. The values are intended to be numbers representing counts, but you *could* store anything in the value field.
- The most_common() method requires only that the values be orderable.
- For in-place operations such as c[key] += 1, the value type need only support addition and subtraction. So fractions, floats, and decimals would work and negative values are supported. The same is also true for update() and subtract() which allow negative and zero values for both inputs and outputs.
- The multiset methods are designed only for use cases with positive values. The inputs may be negative or zero, but only outputs with positive values are created. There are no type restrictions, but the value type needs to support addition, subtraction, and comparison.
- The elements() method requires integer counts. It ignores zero and negative counts.

See also:

- Bag class in Smalltalk.
- Wikipedia entry for Multisets.
- C++ multisets tutorial with examples.
- For mathematical operations on multisets and their use cases, see *Knuth, Donald. The Art of Computer Programming Volume II, Section 4.6.3, Exercise 19.*
- To enumerate all distinct multisets of a given size over a given set of elements, see

csv — CSV File Reading and Writing

Source code: Lib/csv.py

The so-called CSV (Comma Separated Values) format is the most common import and export format for spreadsheets and databases. CSV format was used for many years prior to attempts to describe the format in a standardized way in RFC 4180. The lack of a well-defined standard means that subtle differences often exist in the data produced and consumed by different applications. These differences can make it annoying to process CSV files from multiple sources. Still, while the delimiters and quoting characters vary, the overall format is similar enough that it is possible to write a single module which can efficiently manipulate such data, hiding the details of reading and writing the data from the programmer.

The csv module implements classes to read and write tabular data in CSV format. It allows programmers to say, "write this data in the format preferred by Excel," or "read data from this file which was generated by Excel," without knowing the precise details of the CSV format used by Excel. Programmers can also describe the CSV formats understood by other applications or define their own special-purpose CSV formats.

The csv module's reader and writer objects read and write sequences. Programmers can also read and write data in dictionary form using the DictReader and DictWriter classes.

See also:

PEP 305 - CSV File API

The Python Enhancement Proposal which proposed this addition to Python.

Module Contents

The csv module defines the following functions:

csv.reader(csvfile, dialect='excel', **fmtparams)

Return a reader object which will iterate over lines in the given *csvfile*. *csvfile* can be any object which supports the iterator protocol and returns a string each time its __next__() method is called — file objects and list objects are both suitable. If *csvfile* is a file object, it should be opened with newline=''. [1] An optional *dialect* parameter can be given which is used to define a set of parameters specific to a particular CSV dialect. It may be an instance of a subclass of the Dialect class or one of the strings returned by the list_dialects() function. The other optional *fmtparams* keyword arguments can be given to override individual formatting parameters in the current dialect. For full details about the dialect and formatting parameters, see section Dialects and Formatting Parameters.

Each row read from the csv file is returned as a list of strings. No automatic data type conversion is performed unless the QUOTE_NONNUMERIC format option is specified (in which case unquoted fields are transformed into floats).

csv.writer(csvfile, dialect='excel', **fmtparams)

Return a writer object responsible for converting the user's data into delimited strings on the given file-like object. <code>csvfile</code> can be any object with a write() method. If <code>csvfile</code> is a file object, it should be opened with newline='' [1]. An optional <code>dialect</code> parameter can be given which is used to define a set of parameters specific to a particular CSV dialect. It may be an instance of a subclass of the <code>Dialect</code> class or one of the strings returned by the <code>list_dialects()</code> function. The other optional <code>fmtparams</code> keyword arguments can be given to override individual formatting parameters in the current dialect. For full details about dialects and formatting parameters, see the <code>Dialects</code> and <code>Formatting</code> Parameters section. To make it as easy as possible to interface with modules which implement the DB API, the value <code>None</code> is written as the empty string. While this isn't a reversible transformation, it makes it easier to dump SQL NULL data values to CSV files without preprocessing the data returned from a <code>cursor.fetch*</code> call. All other non-string data are stringified with <code>str()</code> before being written.

A short usage example:

csv.register_dialect(name[, dialect[, **fmtparams]])

Associate *dialect* with *name*. *name* must be a string. The dialect can be specified either by passing a sub-class of Dialect, or by *fmtparams* keyword arguments, or both, with keyword arguments overriding parameters of the dialect. For full details about dialects and formatting parameters, see section Dialects and Formatting Parameters.

csv.unregister_dialect(name)

Delete the dialect associated with *name* from the dialect registry. An Error is raised if *name* is not a registered dialect name.

csv.get dialect(name)

Return the dialect associated with *name*. An Error is raised if *name* is not a registered dialect name. This function returns an immutable Dialect.

csv.list_dialects()

Return the names of all registered dialects.

```
csv.field_size_limit([new_limit])
```



Q

Go

The csv module defines the following classes:

```
class csv.DictReader(f, fieldnames=None, restkey=None, restval=None,
dialect='excel', *args, **kwds)
```

Create an object that operates like a regular reader but maps the information in each row to a dict whose keys are given by the optional *fieldnames* parameter.

The *fieldnames* parameter is a sequence. If *fieldnames* is omitted, the values in the first row of file f will be used as the fieldnames. Regardless of how the fieldnames are determined, the dictionary preserves their original ordering.

If a row has more fields than fieldnames, the remaining data is put in a list and stored with the fieldname specified by *restkey* (which defaults to None). If a non-blank row has fewer fields than fieldnames, the missing values are filled-in with the value of *restval* (which defaults to None).

All other optional or keyword arguments are passed to the underlying reader instance.

If the argument passed to *fieldnames* is an iterator, it will be coerced to a list.

Changed in version 3.6: Returned rows are now of type OrderedDict.

Changed in version 3.8: Returned rows are now of type dict.

A short usage example:

class csv.DictWriter(f, fieldnames, restval='', extrasaction='raise',
dialect='excel', *args, **kwds)

Create an object which operates like a regular writer but maps dictionaries onto output rows. The *fieldnames* parameter is a sequence of keys that identify the order in which values in the dictionary passed to the writerow() method are written to file *f*. The optional *restval* parameter specifies the value to be written if the dictionary is missing a key in *fieldnames*. If the dictionary passed to the writerow() method contains a key not found in *fieldnames*, the optional *extrasaction* parameter indicates what action to take. If it is set to 'raise', the default value, a ValueError is raised. If it is set to 'ignore', extra values in the dictionary are ignored. Any other optional or keyword arguments are passed to the underlying writer instance.

Note that unlike the DictReader class, the *fieldnames* parameter of the DictWriter class is not

Go



3.12.0

Q

If the argument passed to neldnames is an iterator, it will be coerced to a list.

A short usage example:

```
import csv

with open('names.csv', 'w', newline='') as csvfile:
    fieldnames = ['first_name', 'last_name']
    writer = csv.DictWriter(csvfile, fieldnames=fieldnames)

writer.writeheader()
    writer.writerow({'first_name': 'Baked', 'last_name': 'Beans'})
    writer.writerow({'first_name': 'Lovely', 'last_name': 'Spam'})
    writer.writerow({'first_name': 'Wonderful', 'last_name': 'Spam'})
```

class csv.Dialect

The Dialect class is a container class whose attributes contain information for how to handle doublequotes, whitespace, delimiters, etc. Due to the lack of a strict CSV specification, different applications produce subtly different CSV data. Dialect instances define how reader and writer instances behave.

All available Dialect names are returned by list_dialects(), and they can be registered with specific reader and writer classes through their initializer (__init__) functions like this:

class csv.excel

The excel class defines the usual properties of an Excel-generated CSV file. It is registered with the dialect name 'excel'.

class csv.excel tab

The excel_tab class defines the usual properties of an Excel-generated TAB-delimited file. It is registered with the dialect name 'excel-tab'.

class csv.unix dialect

The unix_dialect class defines the usual properties of a CSV file generated on UNIX systems, i.e. using '\n' as line terminator and quoting all fields. It is registered with the dialect name 'unix'.

New in version 3.2.

class csv.Sniffer

The Sniffer class is used to deduce the format of a CSV file.

The Sniffer class provides two methods:



3.12.0

Q

Go

the optional *delimiters* parameter is given, it is interpreted as a string containing possible valid delimiter characters.

has_header(sample)

Analyze the sample text (presumed to be in CSV format) and return True if the first row appears to be a series of column headers. Inspecting each column, one of two key criteria will be considered to estimate if the sample contains a header:

- the second through n-th rows contain numeric values
- the second through n-th rows contain strings where at least one value's length differs from that of the putative header of that column.

Twenty rows after the first row are sampled; if more than half of columns + rows meet the criteria, True is returned.

Note: This method is a rough heuristic and may produce both false positives and negatives.

An example for Sniffer use:

```
with open('example.csv', newline='') as csvfile:
    dialect = csv.Sniffer().sniff(csvfile.read(1024))
    csvfile.seek(0)
    reader = csv.reader(csvfile, dialect)
# ... process CSV file contents here ...
```

The csv module defines the following constants:

csv.QUOTE_ALL

Instructs writer objects to quote all fields.

csv.QUOTE MINIMAL

Instructs writer objects to only quote those fields which contain special characters such as *delimiter*, *quotechar* or any of the characters in *lineterminator*.

CSV.QUOTE NONNUMERIC

Instructs writer objects to quote all non-numeric fields.

Instructs reader objects to convert all non-quoted fields to type float.

csv.QUOTE NONE

Instructs writer objects to never quote fields. When the current *delimiter* occurs in output data it is preceded by the current *escapechar* character. If *escapechar* is not set, the writer will raise Error if any characters that require escaping are encountered.

Instructs reader objects to perform no special processing of quote characters.

csv.QUOTE NOTNULL

Instructs writer objects to quote all fields which are not None. This is similar to QUOTE_ALL,

Instructs reader objects to interpret an empty (unquoted) field as None and to otherwise behave as QUOTE_ALL.

csv.QUOTE_STRINGS

Instructs writer objects to always place quotes around fields which are strings. This is similar to QUOTE_NONNUMERIC, except that if a field value is None an empty (unquoted) string is written.

Instructs reader objects to interpret an empty (unquoted) string as None and to otherwise behave as QUOTE_NONNUMERIC.

The csv module defines the following exception:

exception csv.Error

Raised by any of the functions when an error is detected.

Dialects and Formatting Parameters

To make it easier to specify the format of input and output records, specific formatting parameters are grouped together into dialects. A dialect is a subclass of the Dialect class having a set of specific methods and a single validate() method. When creating reader or writer objects, the programmer can specify a string or a subclass of the Dialect class as the dialect parameter. In addition to, or instead of, the *dialect* parameter, the programmer can also specify individual formatting parameters, which have the same names as the attributes defined below for the Dialect class.

Dialects support the following attributes:

Dialect.delimiter

A one-character string used to separate fields. It defaults to ','.

Dialect.doublequote

Controls how instances of *quotechar* appearing inside a field should themselves be quoted. When True, the character is doubled. When False, the *escapechar* is used as a prefix to the *quotechar*. It defaults to True.

On output, if *doublequote* is False and no *escapechar* is set, Error is raised if a *quotechar* is found in a field.

Dialect.escapechar

A one-character string used by the writer to escape the *delimiter* if *quoting* is set to QUOTE_NONE and the *quotechar* if *doublequote* is False. On reading, the *escapechar* removes any special meaning from the following character. It defaults to None, which disables escaping.

Changed in version 3.11: An empty escapechar is not allowed.

Dialect.lineterminator

The string used to terminate lines produced by the writer. It defaults to '\r\n'.

Dialect.quotechar

A one-character string used to quote fields containing special characters, such as the *delimiter* or *quotechar*, or which contain new-line characters. It defaults to '"'.

Changed in version 3.11: An empty quotechar is not allowed.

Dialect.quoting

Controls when quotes should be generated by the writer and recognised by the reader. It can take on any of the QUOTE_* constants (see section Module Contents) and defaults to QUOTE_MINIMAL.

Dialect.skipinitialspace

When True, spaces immediately following the delimiter are ignored. The default is False.

Dialect.strict

When True, raise exception Error on bad CSV input. The default is False.

Reader Objects

Reader objects (DictReader instances and objects returned by the reader() function) have the following public methods:

```
csvreader.__next__()
```

Return the next row of the reader's iterable object as a list (if the object was returned from reader()) or a dict (if it is a DictReader instance), parsed according to the current Dialect. Usually you should call this as next(reader).

Reader objects have the following public attributes:

csvreader.dialect

A read-only description of the dialect in use by the parser.

csvreader.line_num

The number of lines read from the source iterator. This is not the same as the number of records returned, as records can span multiple lines.

DictReader objects have the following public attribute:

DictReader. fieldnames

If not passed as a parameter when creating the object, this attribute is initialized upon first access or when the first record is read from the file.

Writer Objects

Writer objects (DictWriter instances and objects returned by the writer() function) have the following public methods. A *row* must be an iterable of strings or numbers for Writer objects and a



3.12.0

Q

Go

cause some problems for other programs which read CSV files (assuming they support complex numbers at all).

```
csvwriter.writerow(row)
```

Write the *row* parameter to the writer's file object, formatted according to the current Dialect. Return the return value of the call to the *write* method of the underlying file object.

Changed in version 3.5: Added support of arbitrary iterables.

```
csvwriter.writerows(rows)
```

Write all elements in *rows* (an iterable of *row* objects as described above) to the writer's file object, formatted according to the current dialect.

Writer objects have the following public attribute:

```
csvwriter.dialect
```

A read-only description of the dialect in use by the writer.

DictWriter objects have the following public method:

```
DictWriter.writeheader()
```

Write a row with the field names (as specified in the constructor) to the writer's file object, formatted according to the current dialect. Return the return value of the csvwriter.writerow() call used internally.

New in version 3.2.

Changed in version 3.8: writeheader() now also returns the value returned by the csvwriter.writerow() method it uses internally.

Examples

The simplest example of reading a CSV file:

```
import csv
with open('some.csv', newline='') as f:
    reader = csv.reader(f)
    for row in reader:
        print(row)
```

Reading a file with an alternate format:

```
import csv
with open('passwd', newline='') as f:
    reader = csv.reader(f, delimiter=':', quoting=csv.QUOTE_NONE)
    for row in reader:
        print(row)
```

The corresponding simplest possible writing example is:

Go



Since open() is used to open a CSV file for reading, the file will by default be decoded into unicode using the system default encoding (see locale.getencoding()). To decode a file using a different encoding, use the encoding argument of open:

```
import csv
with open('some.csv', newline='', encoding='utf-8') as f:
    reader = csv.reader(f)
    for row in reader:
        print(row)
```

The same applies to writing in something other than the system default encoding: specify the encoding argument when opening the output file.

Registering a new dialect:

```
import csv
csv.register_dialect('unixpwd', delimiter=':', quoting=csv.QUOTE_NONE)
with open('passwd', newline='') as f:
    reader = csv.reader(f, 'unixpwd')
```

A slightly more advanced use of the reader — catching and reporting errors:

```
import csv, sys
filename = 'some.csv'
with open(filename, newline='') as f:
    reader = csv.reader(f)
    try:
        for row in reader:
            print(row)
    except csv.Error as e:
        sys.exit('file {}, line {}: {}'.format(filename, reader.line_num, e))
```

And while the module doesn't directly support parsing strings, it can easily be done:

```
import csv
for row in csv.reader(['one,two,three']):
    print(row)
```

Footnotes

1(1,2) If newline='' is not specified, newlines embedded inside quoted fields will not be interpreted correctly, and on platforms that use \r\n linendings on write an extra \r will be added. It should always be safe to specify newline='', since the csv module does its own (universal) newline handling.

re — Regular expression operations

Source code: Lib/re/

This module provides regular expression matching operations similar to those found in Perl.

Both patterns and strings to be searched can be Unicode strings (<u>str</u>) as well as 8-bit strings (<u>bytes</u>). However, Unicode strings and 8-bit strings cannot be mixed: that is, you cannot match a Unicode string with a bytes pattern or vice-versa; similarly, when asking for a substitution, the replacement string must be of the same type as both the pattern and the search string.

Regular expressions use the backslash character ('\') to indicate special forms or to allow special characters to be used without invoking their special meaning. This collides with Python's usage of the same character for the same purpose in string literals; for example, to match a literal backslash, one might have to write '\\\' as the pattern string, because the regular expression must be \\, and each backslash must be expressed as \\ inside a regular Python string literal. Also, please note that any invalid escape sequences in Python's usage of the backslash in string literals now generate a SyntaxWarning and in the future this will become a SyntaxError. This behaviour will happen even if it is a valid escape sequence for a regular expression.

The solution is to use Python's raw string notation for regular expression patterns; backslashes are not handled in any special way in a string literal prefixed with 'r'. So r"\n" is a two-character string containing '\' and 'n', while "\n" is a one-character string containing a newline. Usually patterns will be expressed in Python code using this raw string notation.

It is important to note that most regular expression operations are available as module-level functions and methods on <u>compiled regular expressions</u>. The functions are shortcuts that don't require you to compile a regex object first, but miss some fine-tuning parameters.

See also: The third-party <u>regex</u> module, which has an API compatible with the standard library <u>re</u> module, but offers additional functionality and a more thorough Unicode support.

Regular Expression Syntax

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 (or if a given regular expression matches a particular string, which comes down to the same thing).

Regular expressions can be concatenated to form new regular expressions; if A and B are both regular expressions, then AB is also a regular expression. In general, if a string p matches A and another string p matches B, the string pq will match AB. This holds unless A or B contain low precedence operations;

boundary conditions between *A* and *B*; or have numbered group references. Thus, complex expressions can easily be constructed from simpler primitive expressions like the ones described here. For details of the theory and implementation of regular expressions, consult the Friedl book [Frie09], or almost any text-book about compiler construction.

A brief explanation of the format of regular expressions follows. For further information and a gentler presentation, consult the <u>Regular Expression HOWTO</u>.

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'. (In the rest of this section, we'll write RE's in this special style, usually without quotes, and strings to be matched 'in single quotes'.)

Some characters, like '|' or '(', are special. Special characters either stand for classes of ordinary characters, or affect how the regular expressions around them are interpreted.

Repetition operators or quantifiers (*, +, ?, {m,n}, etc) cannot be directly nested. This avoids ambiguity with the non-greedy modifier suffix ?, and with other modifiers in other implementations. To apply a second repetition to an inner repetition, parentheses may be used. For example, the expression $(?:a\{6\})*$ matches any multiple of six 'a' characters.

The special characters are:

- (Dot.) In the default mode, this matches any character except a newline. If the <u>DOTALL</u> flag has been specified, this matches any character including a newline. (?s:.) matches any character regardless of flags.
- (Caret.) Matches the start of the string, and in <u>MULTILINE</u> mode also matches immediately after each newline.
- Matches the end of the string or just before the newline at the end of the string, and in <u>MULTILINE</u> mode also matches before a newline. foo matches both 'foo' and 'foobar', while the regular expression foo\$ matches only 'foo'. More interestingly, searching for foo.\$ in 'foo1\nfoo2\n' matches 'foo2' normally, but 'foo1' in <u>MULTILINE</u> mode; searching for a single \$ in 'foo\n' will find two (empty) matches: one just before the newline, and one at the end of the string.
- Causes the resulting RE to match 0 or more repetitions of the preceding RE, as many repetitions as are possible. ab* will match 'a', 'ab', or 'a' followed by any number of 'b's.
- Causes the resulting RE to match 1 or more repetitions of the preceding RE. ab+ will match 'a' followed by any non-zero number of 'b's; it will not match just 'a'.

2 of 29

\$

*

+

?

Causes the resulting RE to match 0 or 1 repetitions of the preceding RE. ab? will match either 'a' or 'ab'.

*?, +?, ??

The '*', '+', and '?' quantifiers are all greedy; they match as much text as possible. Sometimes this behaviour isn't desired; if the RE <.*> is matched against '<a> b <c>', it will match the entire string, and not just '<a>'. Adding ? after the quantifier makes it perform the match in non-greedy or minimal fashion; as few characters as possible will be matched. Using the RE <.*?> will match only '<a>'.

*+, ++, ?+

Like the '*', '+', and '?' quantifiers, those where '+' is appended also match as many times as possible. However, unlike the true greedy quantifiers, these do not allow back-tracking when the expression following it fails to match. These are known as *possessive* quantifiers. For example, a*a will match 'aaaa' because the a* will match all 4 'a's, but, when the final 'a' is encountered, the expression is backtracked so that in the end the a* ends up matching 3 'a's total, and the fourth 'a' is matched by the final 'a'. However, when a*+a is used to match 'aaaa', the a*+ will match all 4 'a', but when the final 'a' fails to find any more characters to match, the expression cannot be backtracked and will thus fail to match. x*+, x++ and x?+ are equivalent to (?>x*), (?>x+) and (?>x?) correspondingly.

Added in version 3.11.

 $\{m\}$

Specifies that exactly m copies of the previous RE should be matched; fewer matches cause the entire RE not to match. For example, a{6} will match exactly six 'a' characters, but not five.

 $\{m,n\}$

Causes the resulting RE to match from m to n repetitions of the preceding RE, attempting to match as many repetitions as possible. For example, a{3,5} will match from 3 to 5 'a' characters. Omitting m specifies a lower bound of zero, and omitting n specifies an infinite upper bound. As an example, a{4,}b will match 'aaaab' or a thousand 'a' characters followed by a 'b', but not 'aaab'. The comma may not be omitted or the modifier would be confused with the previously described form.

 $\{m,n\}$?

Causes the resulting RE to match from m to n repetitions of the preceding RE, attempting to match as few repetitions as possible. This is the non-greedy version of the previous quantifier. For example, on the 6-character string 'aaaaaa', a{3,5} will match 5 'a' characters, while a{3,5}? will only match 3 characters.

 ${m,n}+$

Causes the resulting RE to match from m to n repetitions of the preceding RE, attempting to match as many repetitions as possible without establishing any backtracking points. This is the possessive version of the quantifier above. For example, on the 6-character string 'aaaaaa', a{3,5}+aa attempt

to match 5 'a' characters, then, requiring 2 more 'a's, will need more characters than available and thus fail, while a $\{3,5\}$ aa will match with a $\{3,5\}$ capturing 5, then 4 'a's by backtracking and then the final 2 'a's are matched by the final aa in the pattern. $x\{m,n\}$ + is equivalent to $(?>x\{m,n\})$.

Added in version 3.11.

Either escapes special characters (permitting you to match characters like '*', '?', and so forth), or signals a special sequence; special sequences are discussed below.

If you're not using a raw string to express the pattern, remember that Python also uses the backslash as an escape sequence in string literals; if the escape sequence isn't recognized by Python's parser, the backslash and subsequent character are included in the resulting string. However, if Python would recognize the resulting sequence, the backslash should be repeated twice. This is complicated and hard to understand, so it's highly recommended that you use raw strings for all but the simplest expressions.

[]

\

Used to indicate a set of characters. In a set:

- Characters can be listed individually, e.g. [amk] will match 'a', 'm', or 'k'.
- Ranges of characters can be indicated by giving two characters and separating them by a '-', for example [a-z] will match any lowercase ASCII letter, [0-5] [0-9] will match all the two-digits numbers from 00 to 59, and [0-9A-Fa-f] will match any hexadecimal digit. If is escaped (e.g. [a\-z]) or if it's placed as the first or last character (e.g. [-a] or [a-]), it will match a literal '-'.
- Special characters lose their special meaning inside sets. For example, [(+*)] will match any of the literal characters '(', '+', '*', or ')'.
- Character classes such as \w or \S (defined below) are also accepted inside a set, although the characters they match depend on the <u>flags</u> used.
- Characters that are not within a range can be matched by *complementing* the set. If the first character of the set is '^', all the characters that are *not* in the set will be matched. For example, [^5] will match any character except '5', and [^^] will match any character except '^'. ^ has no special meaning if it's not the first character in the set.
- To match a literal ']' inside a set, precede it with a backslash, or place it at the beginning of the set. For example, both [()[\]{}] and []()[{}] will match a right bracket, as well as left bracket, braces, and parentheses.
- Support of nested sets and set operations as in <u>Unicode Technical Standard #18</u> might be added in the future. This would change the syntax, so to facilitate this change a <u>FutureWarning</u> will be raised in ambiguous cases for the time being. That includes sets starting with a literal '[' or containing literal character sequences '--', '&&', '~~', and '||'. To avoid a warning escape them with a backslash.

Changed in version 3.7: FutureWarning is raised if a character set contains constructs that will

4 of 29

change semantically in the future.

A|B, where A and B can be arbitrary REs, creates a regular expression that will match either A or B. An arbitrary number of REs can be separated by the '|' in this way. This can be used inside groups (see below) as well. As the target string is scanned, REs separated by '|' are tried from left to right. When one pattern completely matches, that branch is accepted. This means that once A matches, B will not be tested further, even if it would produce a longer overall match. In other words, the '|' operator is never greedy. To match a literal '|', use \|, or enclose it inside a character class, as in [|].

 (\ldots)

Matches whatever regular expression is inside the parentheses, and indicates the start and end of a group; the contents of a group can be retrieved after a match has been performed, and can be matched later in the string with the \number special sequence, described below. To match the literals '(' or ')', use \((or \)), or enclose them inside a character class: [(], [)].

(?...)

This is an extension notation (a '?' following a '(' is not meaningful otherwise). The first character after the '?' determines what the meaning and further syntax of the construct is. Extensions usually do not create a new group; (?P<name>...) is the only exception to this rule. Following are the currently supported extensions.

(?aiLmsux)

(One or more letters from the set 'a', 'i', 'L', 'm', 's', 'u', 'x'.) The group matches the empty string; the letters set the corresponding flags for the entire regular expression:

- re.A (ASCII-only matching)
- re.I (ignore case)
- re.L (locale dependent)
- re.M (multi-line)
- re.S (dot matches all)
- re.U (Unicode matching)
- re.X (verbose)

(The flags are described in <u>Module Contents</u>.) This is useful if you wish to include the flags as part of the regular expression, instead of passing a *flag* argument to the <u>re.compile()</u> function. Flags should be used first in the expression string.

Changed in version 3.11: This construction can only be used at the start of the expression.

(?:...)

A non-capturing version of regular parentheses. Matches whatever regular expression is inside the parentheses, but the substring matched by the group *cannot* be retrieved after performing a match or referenced later in the pattern.

```
(?aiLmsux-imsx:...)
```

(Zero or more letters from the set 'a', 'i', 'L', 'm', 's', 'u', 'x', optionally followed by '-' followed by one or more letters from the 'i', 'm', 's', 'x'.) The letters set or remove the corresponding flags for the part of the expression:

- <u>re.A</u> (ASCII-only matching)
- re.I (ignore case)
- re.L (locale dependent)
- re.M (multi-line)
- re.S (dot matches all)
- re.U (Unicode matching)
- re.X (verbose)

(The flags are described in Module Contents.)

The letters 'a', 'L' and 'u' are mutually exclusive when used as inline flags, so they can't be combined or follow '-'. Instead, when one of them appears in an inline group, it overrides the matching mode in the enclosing group. In Unicode patterns (?a:...) switches to ASCII-only matching, and (?u:...) switches to Unicode matching (default). In bytes patterns (?L:...) switches to locale dependent matching, and (?a:...) switches to ASCII-only matching (default). This override is only in effect for the narrow inline group, and the original matching mode is restored outside of the group.

Added in version 3.6.

Changed in version 3.7: The letters 'a', 'L' and 'u' also can be used in a group.

```
(?>...)
```

Attempts to match ... as if it was a separate regular expression, and if successful, continues to match the rest of the pattern following it. If the subsequent pattern fails to match, the stack can only be unwound to a point *before* the (?>...) because once exited, the expression, known as an *atomic group*, has thrown away all stack points within itself. Thus, (?>.*). would never match anything because first the .* would match all characters possible, then, having nothing left to match, the final . would fail to match. Since there are no stack points saved in the Atomic Group, and there is no stack point before it, the entire expression would thus fail to match.

Added in version 3.11.

```
(?P<name>...)
```

Similar to regular parentheses, but the substring matched by the group is accessible via the symbolic group name *name*. Group names must be valid Python identifiers, and in bytes patterns they can only contain bytes in the ASCII range. Each group name must be defined only once within a regular expression. A symbolic group is also a numbered group, just as if the group were not named.

Named groups can be referenced in three contexts. If the pattern is (?P<quote>['"]).*?(? P=quote) (i.e. matching a string quoted with either single or double quotes):

Context of reference to group "quote"	Ways to reference it		
in the same pattern itself	(?P=quote) (as shown)\1		
when processing match object m	m.group('quote')m.end('quote') (etc.)		
in a string passed to the <i>repl</i> argument of re.sub()	\g<quote></quote>\g<1>\1		

Changed in version 3.12: In <u>bytes</u> patterns, group name can only contain bytes in the ASCII range ($b' \times 00' - b' \times 7f'$).

(?P=name)

A backreference to a named group; it matches whatever text was matched by the earlier group named *name*.

(?#...)

A comment; the contents of the parentheses are simply ignored.

(?=...)

Matches if ... matches next, but doesn't consume any of the string. This is called a *lookahead assertion*. For example, Isaac (?=Asimov) will match 'Isaac ' only if it's followed by 'Asimov'.

(?!...)

Matches if ... doesn't match next. This is a negative lookahead assertion. For example, Isaac (?! Asimov) will match 'Isaac ' only if it's not followed by 'Asimov'.

(?<=...)

Matches if the current position in the string is preceded by a match for ... that ends at the current position. This is called a *positive lookbehind assertion*. (?<=abc)def will find a match in 'abcdef', since the lookbehind will back up 3 characters and check if the contained pattern matches. The contained pattern must only match strings of some fixed length, meaning that abc or a|b are allowed, but a* and a{3,4} are not. Note that patterns which start with positive lookbehind assertions will not match at the beginning of the string being searched; you will most likely want to use the search() function rather than the match() function:

```
>>> import re
>>> m = re.search('(?<=abc)def', 'abcdef')
>>> m.group(0)
'def'
```

This example looks for a word following a hyphen:

```
>>> m = re.search(r'(?<=-)\w+', 'spam-egg')
```

```
>>> m.group(0)
'egg'
```

Changed in version 3.5: Added support for group references of fixed length.

```
(?<!...)
```

Matches if the current position in the string is not preceded by a match for This is called a *negative lookbehind assertion*. Similar to positive lookbehind assertions, the contained pattern must only match strings of some fixed length. Patterns which start with negative lookbehind assertions may match at the beginning of the string being searched.

```
(?(id/name)yes-pattern|no-pattern)
```

Will try to match with yes-pattern if the group with given id or name exists, and with no-pattern if it doesn't. no-pattern is optional and can be omitted. For example, $(<)?(\w+@\w+(?:\\w+)+)(?(1)>|\$)$ is a poor email matching pattern, which will match with '<user@host.com>' as well as 'user@host.com', but not with '<user@host.com' nor 'user@host.com>'.

Changed in version 3.12: Group id can only contain ASCII digits. In <u>bytes</u> patterns, group name can only contain bytes in the ASCII range ($b' \times 00' - b' \times 7f'$).

The special sequences consist of '\' and a character from the list below. If the ordinary character is not an ASCII digit or an ASCII letter, then the resulting RE will match the second character. For example, \\$ matches the character '\$'.

\number

Matches the contents of the group of the same number. Groups are numbered starting from 1. For example, (.+) \1 matches 'the the' or '55 55', but not 'thethe' (note the space after the group). This special sequence can only be used to match one of the first 99 groups. If the first digit of *number* is 0, or *number* is 3 octal digits long, it will not be interpreted as a group match, but as the character with octal value *number*. Inside the '[' and ']' of a character class, all numeric escapes are treated as characters.

****A

Matches only at the start of the string.

****b

Matches the empty string, but only at the beginning or end of a word. A word is defined as a sequence of word characters. Note that formally, $\$ is defined as the boundary between a $\$ w and a $\$ character (or vice versa), or between $\$ and the beginning or end of the string. This means that $\$ r'\bat\b' matches 'at', 'at.', '(at)', and 'as at ay' but not 'attempt' or 'atlas'.

The default word characters in Unicode (str) patterns are Unicode alphanumerics and the underscore, but this can be changed by using the <u>ASCII</u> flag. Word boundaries are determined by the current locale if the <u>LOCALE</u> flag is used.

Note: Inside a character range, \b represents the backspace character, for compatibility with

8 of 29

Python's string literals.

\B

Matches the empty string, but only when it is *not* at the beginning or end of a word. This means that r'at\B' matches 'athens', 'atom', 'attorney', but not 'at', 'at.', or 'at!'. \B is the opposite of \b, so word characters in Unicode (str) patterns are Unicode alphanumerics or the underscore, although this can be changed by using the <u>ASCII</u> flag. Word boundaries are determined by the current locale if the <u>LOCALE</u> flag is used.

Note: Note that \B does not match an empty string, which differs from RE implementations in other programming languages such as Perl. This behavior is kept for compatibility reasons.

\d

For Unicode (str) patterns:

Matches any Unicode decimal digit (that is, any character in Unicode character category [Nd]). This includes [0–9], and also many other digit characters.

Matches [0–9] if the ASCII flag is used.

For 8-bit (bytes) patterns:

Matches any decimal digit in the ASCII character set; this is equivalent to [0-9].

\D

Matches any character which is not a decimal digit. This is the opposite of \d.

Matches [^0-9] if the ASCII flag is used.

\s

For Unicode (str) patterns:

Matches Unicode whitespace characters (as defined by $\underline{str.isspace()}$). This includes [$t\n\r\f\v]$, and also many other characters, for example the non-breaking spaces mandated by typography rules in many languages.

Matches [$\t \n\r\f\v$] if the ASCII flag is used.

For 8-bit (bytes) patterns:

Matches characters considered whitespace in the ASCII character set; this is equivalent to $[\t \r \]$.

\S

Matches any character which is not a whitespace character. This is the opposite of \s.

Matches [$^ \t n\r\f\v$] if the ASCII flag is used.

\w

For Unicode (str) patterns:

Matches Unicode word characters; this includes all Unicode alphanumeric characters (as defined by str.isalnum()), as well as the underscore (_).

Matches [a-zA-Z0-9_] if the ASCII flag is used.

For 8-bit (bytes) patterns:

Matches characters considered alphanumeric in the ASCII character set; this is equivalent to [a–zA–Z0–9_]. If the <u>LOCALE</u> flag is used, matches characters considered alphanumeric in the current locale and the underscore.

\W

Matches any character which is not a word character. This is the opposite of \w. By default, matches non-underscore (_) characters for which str.isalnum() returns False.

Matches [^a-zA-Z0-9_] if the ASCII flag is used.

If the <u>LOCALE</u> flag is used, matches characters which are neither alphanumeric in the current locale nor the underscore.

\Z

Matches only at the end of the string.

Most of the <u>escape sequences</u> supported by Python string literals are also accepted by the regular expression parser:

∖a	\b	\f	\n
\N	\r	\t	\u
\U	\v	\f \t \x	11

(Note that \b is used to represent word boundaries, and means "backspace" only inside character classes.)

'\u', '\U', and '\N' escape sequences are only recognized in Unicode (str) patterns. In bytes patterns they are errors. Unknown escapes of ASCII letters are reserved for future use and treated as errors.

Octal escapes are included in a limited form. If the first digit is a 0, or if there are three octal digits, it is considered an octal escape. Otherwise, it is a group reference. As for string literals, octal escapes are always at most three digits in length.

Changed in version 3.3: The '\u' and '\U' escape sequences have been added.

Changed in version 3.6: Unknown escapes consisting of '\' and an ASCII letter now are errors.

Changed in version 3.8: The '\N{name}' escape sequence has been added. As in string literals, it expands to the named Unicode character (e.g. '\N{EM DASH}').

Module Contents

10/6/25, 11:20 AM

at a time needn't worry about compiling regular expressions.

re.search(pattern, string, flags=0)

Scan through *string* looking for the first location where the regular expression *pattern* produces a match, and return a corresponding <u>Match</u>. Return None if no position in the string matches the pattern; note that this is different from finding a zero-length match at some point in the string.

The expression's behaviour can be modified by specifying a *flags* value. Values can be any of the <u>flags</u> variables, combined using bitwise OR (the | operator).

re.match(pattern, string, flags=0)

If zero or more characters at the beginning of *string* match the regular expression *pattern*, return a corresponding <u>Match</u>. Return None if the string does not match the pattern; note that this is different from a zero-length match.

Note that even in <u>MULTILINE</u> mode, <u>re.match()</u> will only match at the beginning of the string and not at the beginning of each line.

If you want to locate a match anywhere in *string*, use <u>search()</u> instead (see also <u>search()</u> vs. match()).

The expression's behaviour can be modified by specifying a *flags* value. Values can be any of the <u>flags</u> variables, combined using bitwise OR (the | operator).

re.fullmatch(pattern, string, flags=0)

If the whole *string* matches the regular expression *pattern*, return a corresponding <u>Match</u>. Return None if the string does not match the pattern; note that this is different from a zero-length match.

The expression's behaviour can be modified by specifying a *flags* value. Values can be any of the flags variables, combined using bitwise OR (the | operator).

Added in version 3.4.

re.split(pattern, string, maxsplit=0, flags=0)

Split *string* by the occurrences of *pattern*. 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. If *maxsplit* is nonzero, at most *maxsplit* splits occur, and the remainder of the string is returned as the final element of the list.

```
>>> re.split(r'\W+', 'Words, words, words.')
['Words', 'words', '']
>>> re.split(r'(\W+)', 'Words, words, words.')
['Words', ', ', 'words', ', ', 'words', '.', '']
>>> re.split(r'\W+', 'Words, words, words.', 1)
['Words', 'words, words.']
>>> re.split('[a-f]+', '0a3B9', flags=re.IGNORECASE)
['0', '3', '9']
```

If there are capturing groups in the separator and it matches at the start of the string, the result will

start with an empty string. The same holds for the end of the string:

```
>>> re.split(r'(\W+)', '...words, words...')
['', '...', 'words', ', ', 'words', '...', '']
```

That way, separator components are always found at the same relative indices within the result list.

Empty matches for the pattern split the string only when not adjacent to a previous empty match.

```
>>> re.split(r'\b', 'Words, words, words.')
['', 'Words', ', ', 'words', '.']
>>> re.split(r'\W*', '...words...')
['', '', 'w', 'o', 'r', 'd', 's', '', '']
>>> re.split(r'(\W*)', '...words...')
['', '...', '', 'w', '', 'o', '', 'r', '', 'd', '', 's', '...', '', '']
```

The expression's behaviour can be modified by specifying a *flags* value. Values can be any of the <u>flags</u> variables, combined using bitwise OR (the | operator).

Changed in version 3.1: Added the optional flags argument.

Changed in version 3.7: Added support of splitting on a pattern that could match an empty string.

```
re.findall(pattern, string, flags=0)
```

Return all non-overlapping matches of *pattern* in *string*, as a list of strings or tuples. The *string* is scanned left-to-right, and matches are returned in the order found. Empty matches are included in the result.

The result depends on the number of capturing groups in the pattern. If there are no groups, return a list of strings matching the whole pattern. If there is exactly one group, return a list of strings matching that group. If multiple groups are present, return a list of tuples of strings matching the groups. Non-capturing groups do not affect the form of the result.

```
>>> re.findall(r'\bf[a-z]*', 'which foot or hand fell fastest') ['foot', 'fell', 'fastest'] >>> re.findall(r'(\w+)=(\d+)', 'set width=20 and height=10') [('width', '20'), ('height', '10')]
```

The expression's behaviour can be modified by specifying a *flags* value. Values can be any of the <u>flags</u> variables, combined using bitwise OR (the | operator).

Changed in version 3.7: Non-empty matches can now start just after a previous empty match.

re.finditer(pattern, string, flags=0)

Return an <u>iterator</u> yielding <u>Match</u> objects over all non-overlapping matches for the RE *pattern* in *string*. The *string* is scanned left-to-right, and matches are returned in the order found. Empty matches are included in the result.

The expression's behaviour can be modified by specifying a *flags* value. Values can be any of the <u>flags</u> variables, combined using bitwise OR (the | operator).

Changed in version 3.7: Non-empty matches can now start just after a previous empty match.

```
re.sub(pattern, repl, string, count=0, flags=0)
```

Return the string obtained by replacing the leftmost non-overlapping occurrences of *pattern* in *string* by the replacement *repl*. If the pattern isn't found, *string* is returned unchanged. *repl* can be a string or a function; if it is a string, any backslash escapes in it are processed. That is, \n is converted to a single newline character, \r is converted to a carriage return, and so forth. Unknown escapes of ASCII letters are reserved for future use and treated as errors. Other unknown escapes such as \& are left alone. Backreferences, such as \6, are replaced with the substring matched by group 6 in the pattern. For example:

If *repl* is a function, it is called for every non-overlapping occurrence of *pattern*. The function takes a single Match argument, and returns the replacement string. For example:

```
>>> def dashrepl(matchobj):
...    if matchobj.group(0) == '-': return ' '
...    else: return '-'
...
>>> re.sub('-{1,2}', dashrepl, 'pro----gram-files')
'pro--gram files'
>>> re.sub(r'\sAND\s', ' & ', 'Baked Beans And Spam', flags=re.IGNORECASE)
'Baked Beans & Spam'
```

The pattern may be a string or a Pattern.

The optional argument *count* is the maximum number of pattern occurrences to be replaced; *count* must be a non-negative integer. If omitted or zero, all occurrences will be replaced. Empty matches for the pattern are replaced only when not adjacent to a previous empty match, so sub('x*', '-', 'abxd') returns '-a-b--d-'.

In string-type *repl* arguments, in addition to the character escapes and backreferences described above, \g<name> will use the substring matched by the group named name, as defined by the (? P<name>...) syntax. \g<number> uses the corresponding group number; \g<2> is therefore equivalent to \2, but isn't ambiguous in a replacement such as \g<2>0. \20 would be interpreted as a reference to group 20, not a reference to group 2 followed by the literal character '0'. The backreference \g<0> substitutes in the entire substring matched by the RE.

The expression's behaviour can be modified by specifying a *flags* value. Values can be any of the flags variables, combined using bitwise OR (the | operator).

Changed in version 3.1: Added the optional flags argument.

Changed in version 3.5: Unmatched groups are replaced with an empty string.

Changed in version 3.6: Unknown escapes in pattern consisting of '\' and an ASCII letter now are errors.

Changed in version 3.7: Unknown escapes in repl consisting of '\' and an ASCII letter now are errors.

Changed in version 3.7: Empty matches for the pattern are replaced when adjacent to a previous non-empty match.

Changed in version 3.12: Group id can only contain ASCII digits. In <u>bytes</u> replacement strings, group *name* can only contain bytes in the ASCII range (b'\x00'-b'\x7f').

```
re.subn(pattern, repl, string, count=0, flags=0)
```

Perform the same operation as sub(), but return a tuple (new_string, number_of_subs_made).

Changed in version 3.1: Added the optional flags argument.

Changed in version 3.5: Unmatched groups are replaced with an empty string.

The expression's behaviour can be modified by specifying a *flags* value. Values can be any of the <u>flags</u> variables, combined using bitwise OR (the | operator).

re.escape(pattern)

Escape special characters in *pattern*. This is useful if you want to match an arbitrary literal string that may have regular expression metacharacters in it. For example:

```
>>> print(re.escape('https://www.python.org'))
https://www\.python\.org
>>> legal_chars = string.ascii_lowercase + string.digits + "!#$%&'*+-.^_`|~:"
>>> print('[%s]+' % re.escape(legal_chars))
[abcdefghijklmnopqrstuvwxyz0123456789!\#\$%\&'\*\+\-\.\^_`\|\~:]+
>>> operators = ['+', '-', '*', '/', '**']
>>> print('|'.join(map(re.escape, sorted(operators, reverse=True))))
/|\-|\+|\*\*|\*
```

This function must not be used for the replacement string in <u>sub()</u> and <u>subn()</u>, only backslashes should be escaped. For example:

```
>>> digits_re = r'\d+'
>>> sample = '/usr/sbin/sendmail - 0 errors, 12 warnings'
>>> print(re.sub(digits_re, digits_re.replace('\\', r'\\'), sample))
/usr/sbin/sendmail - \d+ errors, \d+ warnings
```

Changed in version 3.3: The '_' character is no longer escaped.

re.purge()

Clear the regular expression cache.

Exceptions

```
exception re.error(msg, pattern=None, pos=None)
```

Exception raised when a string passed to one of the functions here is not a valid regular expression (for example, it might contain unmatched parentheses) or when some other error occurs during compilation or matching. It is never an error if a string contains no match for a pattern. The error instance has the following additional attributes:

msg

The unformatted error message.

pattern

The regular expression pattern.

pos

The index in pattern where compilation failed (may be None).

lineno

The line corresponding to pos (may be None).

colno

The column corresponding to pos (may be None).

Changed in version 3.5: Added additional attributes.

Regular Expression Objects

class re.Pattern

Compiled regular expression object returned by re.compile().

Changed in version 3.9: re.Pattern supports [] to indicate a Unicode (str) or bytes pattern. See Generic Alias Type.

```
Pattern.search(string[, pos[, endpos]])
```

Scan through *string* looking for the first location where this regular expression produces a match, and return a corresponding <u>Match</u>. Return None if no position in the string matches the pattern; note that this is different from finding a zero-length match at some point in the string.

The optional second parameter *pos* gives an index in the string where the search is to start; it defaults to 0. This is not completely equivalent to slicing the string; the '^' pattern character matches at the real beginning of the string and at positions just after a newline, but not necessarily at the index where the search is to start.

The optional parameter endpos limits how far the string will be searched; it will be as if the string is endpos characters long, so only the characters from pos to endpos -1 will be searched for a match. If endpos is less than pos, no match will be found; otherwise, if rx is a compiled regular expression object, rx.search(string, 0, 50) is equivalent to rx.search(string[:50], 0).

```
>>> pattern = re.compile("d")
>>> pattern.search("dog")  # Match at index 0
<re.Match object; span=(0, 1), match='d'>
>>> pattern.search("dog", 1)  # No match; search doesn't include the "d"
```

Pattern.match(string[, pos[, endpos]])

If zero or more characters at the *beginning* of *string* match this regular expression, return a corresponding <u>Match</u>. Return None if the string does not match the pattern; note that this is different from a zero-length match.

The optional pos and endpos parameters have the same meaning as for the search() method.

```
>>> pattern = re.compile("o")
>>> pattern.match("dog")  # No match as "o" is not at the start of "dog".
>>> pattern.match("dog", 1)  # Match as "o" is the 2nd character of "dog".
<re.Match object; span=(1, 2), match='o'>
```

If you want to locate a match anywhere in *string*, use <u>search()</u> instead (see also <u>search()</u> vs. match()).

```
Pattern.fullmatch(string[, pos[, endpos]])
```

If the whole *string* matches this regular expression, return a corresponding <u>Match</u>. Return None if the string does not match the pattern; note that this is different from a zero-length match.

The optional pos and endpos parameters have the same meaning as for the search() method.

```
>>> pattern = re.compile("o[gh]")
>>> pattern.fullmatch("dog")  # No match as "o" is not at the start of "dog"
>>> pattern.fullmatch("ogre")  # No match as not the full string matches.
>>> pattern.fullmatch("doggie", 1, 3)  # Matches within given limits.
<re.Match object; span=(1, 3), match='og'>
```

Added in version 3.4.

Pattern.split(string, maxsplit=0)

Identical to the split() function, using the compiled pattern.

```
Pattern.findall(string[, pos[, endpos]])
```

Similar to the <u>findall()</u> function, using the compiled pattern, but also accepts optional *pos* and *endpos* parameters that limit the search region like for search().

Pattern.finditer(string[, pos[, endpos]])

Similar to the <u>finditer()</u> function, using the compiled pattern, but also accepts optional *pos* and *endpos* parameters that limit the search region like for search().

```
Pattern.sub(repl, string, count=0)
```

Identical to the sub() function, using the compiled pattern.

```
Pattern.subn(repl, string, count=0)
```

Identical to the subn() function, using the compiled pattern.

Pattern.flags

The regex matching flags. This is a combination of the flags given to <u>compile()</u>, any (?...) inline flags in the pattern, and implicit flags such as <u>UNICODE</u> if the pattern is a Unicode string.

Pattern.groups

The number of capturing groups in the pattern.

Pattern.groupindex

A dictionary mapping any symbolic group names defined by (?P<id>) to group numbers. The dictionary is empty if no symbolic groups were used in the pattern.

Pattern.pattern

The pattern string from which the pattern object was compiled.

Changed in version 3.7: Added support of copy.copy() and copy.deepcopy(). Compiled regular expression objects are considered atomic.

Match Objects

Match objects always have a boolean value of True. Since $\underline{match()}$ and $\underline{search()}$ return None when there is no match, you can test whether there was a match with a simple if statement:

```
match = re.search(pattern, string)
if match:
    process(match)
```

class re.Match

Match object returned by successful matches and searches.

Changed in version 3.9: re.Match supports [] to indicate a Unicode (str) or bytes match. See Generic Alias Type.

Match.expand(template)

Return the string obtained by doing backslash substitution on the template string template, as done

20 of 29

math — Mathematical functions

This module provides access to the mathematical functions defined by the C standard.

These functions cannot be used with complex numbers; use the functions of the same name from the cmath module if you require support for complex numbers. The distinction between functions which support complex numbers and those which don't is made since most users do not want to learn quite as much mathematics as required to understand complex numbers. Receiving an exception instead of a complex result allows earlier detection of the unexpected complex number used as a parameter, so that the programmer can determine how and why it was generated in the first place.

The following functions are provided by this module. Except when explicitly noted otherwise, all return values are floats.

Number-theoretic and representation functions

math.ceil(x)

Return the ceiling of x, the smallest integer greater than or equal to x. If x is not a float, delegates to x.__ceil__, which should return an Integral value.

math.comb(n, k)

Return the number of ways to choose *k* items from *n* items without repetition and without order.

```
Evaluates to n! / (k! * (n - k)!) when k \le n and evaluates to zero when k > n.
```

Also called the binomial coefficient because it is equivalent to the coefficient of k-th term in polynomial expansion of $(1 + x)^n$.

Raises <u>TypeError</u> if either of the arguments are not integers. Raises <u>ValueError</u> if either of the arguments are negative.

```
Added in version 3.8.
```

math.copysign(x, y)

Return a float with the magnitude (absolute value) of x but the sign of y. On platforms that support signed zeros, copysign(1.0, -0.0) returns -1.0.

math.fabs(x)

Return the absolute value of x.

math.factorial(n)

Return *n* factorial as an integer. Raises ValueError if *n* is not integral or is negative.

Return True if the values a and b are close to each other and False otherwise.

Whether or not two values are considered close is determined according to given absolute and relative tolerances. If no errors occur, the result will be: $abs(a-b) \le max(rel_tol * max(abs(a), abs(b)), abs_tol)$.

rel_tol is the relative tolerance – it is the maximum allowed difference between a and b, relative to the larger absolute value of a or b. For example, to set a tolerance of 5%, pass rel_tol=0.05. The default tolerance is 1e-09, which assures that the two values are the same within about 9 decimal digits. rel_tol must be nonnegative and less than 1.0.

abs_tol is the absolute tolerance; it defaults to 0.0 and it must be nonnegative. When comparing x to 0.0, isclose(x, 0) is computed as abs(x) <= rel_tol * abs(x), which is False for any nonzero x and rel_tol less than 1.0. So add an appropriate positive abs_tol argument to the call.

The IEEE 754 special values of NaN, inf, and -inf will be handled according to IEEE rules. Specifically, NaN is not considered close to any other value, including NaN. inf and -inf are only considered close to themselves.

Added in version 3.5.

See also: PEP 485 – A function for testing approximate equality

math.isfinite(x)

Return True if x is neither an infinity nor a NaN, and False otherwise. (Note that 0.0 is considered finite.)

Added in version 3.2.

math.isinf(x)

Return True if x is a positive or negative infinity, and False otherwise.

math.**isnan**(x)

Return True if x is a NaN (not a number), and False otherwise.

math.isqrt(n)

Return the integer square root of the nonnegative integer n. This is the floor of the exact square root of n, or equivalently the greatest integer a such that $a^2 \le n$.

For some applications, it may be more convenient to have the least integer a such that $n \le a^2$, or in other words the ceiling of the exact square root of n. For positive n, this can be computed using $a = 1 + i \operatorname{sqrt}(n - 1)$.

Added in version 3.8.

math.lcm(*integers)

```
the first float bigger than x is x + ulp(x).
```

ULP stands for "Unit in the Last Place".

See also math.nextafter() and sys.float_info.epsilon.

```
Added in version 3.9.
```

Note that <u>frexp()</u> and <u>modf()</u> have a different call/return pattern than their C equivalents: they take a single argument and return a pair of values, rather than returning their second return value through an 'output parameter' (there is no such thing in Python).

For the <u>ceil()</u>, <u>floor()</u>, and <u>modf()</u> functions, note that *all* floating-point numbers of sufficiently large magnitude are exact integers. Python floats typically carry no more than 53 bits of precision (the same as the platform C double type), in which case any float x with abs(x) >= 2**52 necessarily has no fractional bits.

Power and logarithmic functions

```
math.cbrt(x)
```

Return the cube root of x.

Added in version 3.11.

math.exp(x)

Return e raised to the power x, where e = 2.718281... is the base of natural logarithms. This is usually more accurate than math.e ** x or pow(math.e, x).

math.exp2(x)

Return 2 raised to the power x.

Added in version 3.11.

math.expm1(x)

Return e raised to the power x, minus 1. Here e is the base of natural logarithms. For small floats x, the subtraction in $\exp(x) - 1$ can result in a <u>significant loss of precision</u>; the $\exp(x)$ function provides a way to compute this quantity to full precision:

```
>>> from math import exp, expm1
>>> exp(1e-5) - 1 # gives result accurate to 11 places
1.0000050000069649e-05
>>> expm1(1e-5) # result accurate to full precision
1.0000050000166668e-05
```

Added in version 3.2.

```
math.log(x[, base])
```

With one argument, return the natural logarithm of x (to base e).

With two arguments, return the logarithm of x to the given base, calculated as log(x)/log(base).

math.log1p(x)

Return the natural logarithm of 1+x (base e). The result is calculated in a way which is accurate for x near zero.

math.log2(x)

Return the base-2 logarithm of x. This is usually more accurate than log(x, 2).

Added in version 3.3.

See also: <u>int.bit_length()</u> returns the number of bits necessary to represent an integer in binary, excluding the sign and leading zeros.

math.log10(x)

Return the base-10 logarithm of x. This is usually more accurate than log(x, 10).

math.pow(x, y)

Return x raised to the power y. Exceptional cases follow the IEEE 754 standard as far as possible. In particular, pow(1.0, x) and pow(x, 0.0) always return 1.0, even when x is a zero or a NaN. If both x and y are finite, x is negative, and y is not an integer then pow(x, y) is undefined, and raises ValueError.

Unlike the built-in ** operator, math.pow() converts both its arguments to type float. Use ** or the built-in pow() function for computing exact integer powers.

Changed in version 3.11: The special cases pow(0.0, -inf) and pow(-0.0, -inf) were changed to return inf instead of raising ValueError, for consistency with IEEE 754.

math.sqrt(x)

Return the square root of x.

Trigonometric functions

math.acos(x)

Return the arc cosine of x, in radians. The result is between 0 and pi.

math.asin(x)

Return the arc sine of x, in radians. The result is between -pi/2 and pi/2.

math.atan(x)

Return the arc tangent of x, in radians. The result is between -pi/2 and pi/2.

math.atan2(y, x)

Return atan(y / x), in radians. The result is between -pi and pi. The vector in the plane from the

7 of 11

math.lgamma(x)

Return the natural logarithm of the absolute value of the Gamma function at x.

Added in version 3.2.

Constants

math.pi

The mathematical constant π = 3.141592..., to available precision.

math.e

The mathematical constant e = 2.718281..., to available precision.

math.tau

The mathematical constant τ = 6.283185..., to available precision. Tau is a circle constant equal to 2π , the ratio of a circle's circumference to its radius. To learn more about Tau, check out Vi Hart's video \underline{Pi} is (still) Wrong, and start celebrating \underline{Tau} day by eating twice as much pie!

Added in version 3.6.

math.inf

A floating-point positive infinity. (For negative infinity, use -math.inf.) Equivalent to the output of float('inf').

Added in version 3.5.

math.nan

A floating-point "not a number" (NaN) value. Equivalent to the output of float('nan'). Due to the requirements of the IEEE-754 standard, math.nan and float('nan') are not considered to equal to any other numeric value, including themselves. To check whether a number is a NaN, use the Isnan()) function to test for NaNs instead of is or ==. Example:

```
>>> import math
>>> math.nan == math.nan
False
>>> float('nan') == float('nan')
False
>>> math.isnan(math.nan)
True
>>> math.isnan(float('nan'))
True
```

Added in version 3.5.

Changed in version 3.11: It is now always available.

CPython implementation detail: The <u>math</u> module consists mostly of thin wrappers around the platform C math library functions. Behavior in exceptional cases follows Annex F of the C99 standard where appro-