

Регулярные выражения в Python: подробное руководство с примерами

Регулярные выражения (regex) — мощный инструмент для поиска и обработки текста в Python. Модуль `re` предоставляет все необходимые функции для работы с регулярными выражениями. В этой статье рассмотрим основные метасимволы и практические примеры их использования.

Основные метасимволы регулярных выражений

Метасимвол	Описание	Пример	Соответствие
<code>.</code>	Любой символ кроме новой строки	<code>a.c</code>	<code>abc</code> , <code>a1c</code> , <code>a@c</code>
<code>^</code>	Начало строки	<code>^hello</code>	<code>hello world</code> (в начале)
<code>\$</code>	Конец строки	<code>world\$</code>	<code>hello world</code> (в конце)
<code>*</code>	0 или более повторений	<code>ab*c</code>	<code>ac</code> , <code>abc</code> , <code>abbc</code>
<code>+</code>	1 или более повторений	<code>ab+c</code>	<code>abc</code> , <code>abbc</code> (не <code>ac</code>)
<code>?</code>	0 или 1 повторение	<code>ab?c</code>	<code>ac</code> , <code>abc</code>
<code>{n}</code>	Ровно <code>n</code> повторений	<code>a{3}</code>	<code>aaa</code>
<code>{n,}</code>	<code>n</code> или более повторений	<code>a{2,}</code>	<code>aa</code> , <code>aaa</code> , <code>aaaa</code>
<code>{n,m}</code>	От <code>n</code> до <code>m</code> повторений	<code>a{2,4}</code>	<code>aa</code> , <code>aaa</code> , <code>aaaa</code>
<code>[]</code>	Символьный класс	<code>[abc]</code>	<code>a</code> , <code>b</code> , <code>c</code>
<code>[^]</code>	Отрицание класса	<code>[^abc]</code>	Любой символ кроме <code>a</code> , <code>b</code> , <code>c</code>
<code> </code>	Альтернатива (ИЛИ)	<code>cat dog</code>	<code>cat</code> или <code>dog</code>

Метасимвол	Описание	Пример	Соответствие
()	Группировка	(ab)+	ab, abab, ababab
\	Экранирование	\\$	Буквальный символ \$

Символьные классы

Класс	Описание	Эквивалент
\d	Цифра	[0-9]
\D	Не цифра	[^0-9]
\w	Буква, цифра или подчёркивание	[a-zA-Z0-9_]
\W	Не буква, цифра или подчёркивание	[^a-zA-Z0-9_]
\s	Пробельный символ	[\t\n\r\f\v]
\S	Не пробельный символ	[^\t\n\r\f\v]

Квантификаторы (жадные и ленивые)

Жадный	Ленивый	Описание
*	*?	0 или более (минимум)
+	+?	1 или более (минимум)
?	??	0 или 1 (минимум)
{n,m}	{n,m}?	От n до m (минимум)

Границы слов

Метасимвол	Описание	Пример
\b	Граница слова	\bcat\b находит cat как отдельное слово
\B	Не граница слова	\Bcat\B находит cat внутри слова

Особые последовательности

Последовательность	Описание
\n	Символ новой строки
\t	Символ табуляции
\r	Символ возврата каретки
\f	Символ перевода страницы
\v	Символ вертикальной табуляции
\0	Нулевой символ
\xhh	Символ с шестнадцатеричным кодом hh

Примеры использования

```
import re

# Поиск email
email_pattern = r'\b[A-Za-z0-9._%+-]+@[A-Za-z0-9.-]+\.[A-Z|a-z]{2,}\b'

# Поиск телефона
phone_pattern = r'\+?[1-9]\d{1,14}'

# Поиск даты в формате DD.MM.YYYY
date_pattern = r'\d{2}\.\d{2}\.\d{4}'

# Поиск IP-адреса
ip_pattern = r'\b(?:[0-9]{1,3}\.){3}[0-9]{1,3}\b'
```

Точка (.) — универсальный символ

Метасимвол точка соответствует любому символу, кроме символа новой строки. Это один из наиболее часто используемых элементов регулярных выражений.

```
import re

text = "cat bat mat rat"
pattern = r".at" # Находим слова, оканчивающиеся на "at"
matches = re.findall(pattern, text)
print(matches) # Вывод: ['cat', 'bat', 'mat', 'rat']
```

Якорь начала строки (^)

Символ ^ указывает на начало строки. Он позволяет найти шаблон только в том случае, если он находится в самом начале текста.

```
import re

text = "hello world"
pattern = r"^hello" # Ищем слово, начинающееся с "hello"
match = re.search(pattern, text)
if match:
    print(match.group()) # Вывод: 'hello'
```

Якорь конца строки (\$)

Символ \$ соответствует концу строки. Используется для поиска шаблонов в конце текста.

```
import re

text = "hello world"
pattern = r"world$" # Ищем слово, заканчивающееся на "world"
match = re.search(pattern, text)
if match:
    print(match.group()) # Вывод: 'world'
```

Квантификаторы в регулярных выражениях

Звездочка (*) — ноль или более вхождений

Квантификатор * соответствует нулю или более вхождениям предшествующего элемента.

```
import re

text = "ac abc abbc abbbc"
pattern = r"ab*c" # Ищем "a", за которым идет ноль или более "b", а затем "c"
matches = re.findall(pattern, text)
print(matches) # Вывод: ['ac', 'abc', 'abbc', 'abbbc']
```

Плюс (+) — одно или более вхождений

Квантификатор + требует наличия хотя бы одного вхождения предшествующего элемента.

```
import re

text = "ac abc abbc abbbc"
pattern = r"ab+c" # Ищем "a", за которым идет один или более "b", а затем "c"
matches = re.findall(pattern, text)
print(matches) # Вывод: ['abc', 'abbc', 'abbbc']
```

Знак вопроса (?) — ноль или одно вхождение

Квантификатор ? делает предшествующий элемент необязательным.

```
import re

text = "colour color"
pattern = r"colou?r" # Ищем "colou" с нулевым или одним "u", за которым следует "r"
matches = re.findall(pattern, text)
print(matches) # Вывод: ['colour', 'color']
```

Практические примеры использования regex

Валидация email-адресов

Поиск и валидация email-адресов — одна из наиболее распространенных задач при работе с регулярными выражениями.

```
import re

text = "Emails: test@example.com, another.test@mail.co.uk"
pattern = r"\b[A-Za-z0-9._%+-]+@[A-Za-z0-9.-]+\.[A-Z|a-z]{2,}\b"
matches = re.findall(pattern, text)
print(matches) # Вывод: ['test@example.com', 'another.test@mail.co.uk']
```

Извлечение URL из текста

Регулярные выражения помогают находить веб-адреса в тексте различных форматов.

```
import re

text = "Visit my website at https://www.example.com"
pattern = r"https?://(?:www\.)?([a-zA-Z0-9]+\.[a-zA-Z]{2,})"
matches = re.findall(pattern, text)
print(matches) # Вывод: ['example.com']
```

Поиск хэштегов в социальных сетях

Для анализа контента социальных сетей часто требуется извлечение хэштегов.

```
import re
```

```
text = "Check out #Python and #DataScience on Twitter!"
pattern = r"#\w+"
matches = re.findall(pattern, text)
print(matches) # Вывод: ['#Python', '#DataScience']
```

Поиск IP-адресов

Регулярные выражения эффективно находят IP-адреса в логах и конфигурационных файлах.

```
import re

text = "IP addresses: 192.168.1.1 and 10.0.0.1"
pattern = r"\b(?:\d{1,3}\.){3}\d{1,3}\b"
matches = re.findall(pattern, text)
print(matches) # Вывод: ['192.168.1.1', '10.0.0.1']
```

Парсинг дат в формате DD/MM/YYYY

Извлечение дат из текста с использованием групп захвата.

```
import re

text = "Dates: 20/01/2022 and 31/12/2023"
pattern = r"\b(0[1-9]|[12]\d|3[01])/(0[1-9]|1[0-2])/(\d{4})\b"
matches = re.findall(pattern, text)
print(matches) # Вывод: [('20', '01', '2022'), ('31', '12', '2023')]
```

Поиск HTML-тегов

Для обработки HTML-разметки используются ленивые квантификаторы.

```
import re

text = "<p>This is a paragraph</p> <a href='https://example.com'>Link</a>"
pattern = r"<.*?>" # Ленивый поиск всех символов между "<" и ">"
matches = re.findall(pattern, text)
print(matches) # Вывод: ['<p>', '</p>', "<a href='https://example.com'>", '</a>']
```

Продвинутые техники работы с regex

Поиск слов с заглавной буквы

Для выделения имен собственных и начала предложений.

```
import re

text = "Hello world, Python is Amazing"
pattern = r"\b[A-Z][a-z]*\b"
matches = re.findall(pattern, text)
print(matches) # Вывод: ['Hello', 'Python', 'Amazing']
```

Поиск палиндромов

Использование обратных ссылок для поиска слов, которые начинаются и заканчиваются одинаково.

```
import re
```

```
text = "level deed book car door"
pattern = r"\b(\w)\w*\1\b"
matches = re.findall(pattern, text, re.IGNORECASE)
print(matches) # Находит слова-палиндромы
```

Поиск повторяющихся символов

Выявление слов с двумя или более повторяющимися символами подряд.

```
import re

text = "ssss yyyy aa bb cccc"
pattern = r"\b\w*(\w)\1+\w*\b"
matches = re.findall(pattern, text)
print([match for match in re.finditer(r"\b\w*(\w)\1+\w*\b", text)])
```

Поиск слов определенной длины

Для анализа текста по длине слов используются квантификаторы с фигурными скобками.

```
import re

text = "apple banana orange kiwi"
pattern = r"\b\w{6}\b" # Ищем слова длиной в 6 символов
matches = re.findall(pattern, text)
print(matches) # Вывод: ['banana', 'orange']
```

Дополнительные примеры для практики

Поиск номеров телефонов

```
import re

text = "Contact: +7-900-123-45-67 or 8(495)123-45-67"
pattern = r"(\+7|8)[\s\-\(\)]?(\d{3})[\s\-\(\)]?(\d{3})[\s\-\(\)]?(\d{2})[\s\-\(\)]?(\d{2})"
matches = re.findall(pattern, text)
print(matches)
```

Извлечение чисел из текста

```
import re

text = "Temperature: 25.5°C, Humidity: 60%, Pressure: 1013.25 hPa"
pattern = r"\d+\.\d*"
matches = re.findall(pattern, text)
print(matches) # Вывод: ['25.5', '60', '1013.25']
```

Поиск слов в кавычках

```
import re

text = 'He said "Hello world" and then "Goodbye"'
pattern = r'"([\^"]*)"'
matches = re.findall(pattern, text)
print(matches) # Вывод: ['Hello world', 'Goodbye']
```

Регулярные выражения в Python - Полная справочная таблица

Основные методы модуля re

Метод	Описание	Пример использования
<code>re.search(pattern, string, flags=0)</code>	Ищет первое совпадение в строке	<code>re.search(r'\d+', 'abc123def')</code>
<code>re.match(pattern, string, flags=0)</code>	Проверяет совпадение в начале строки	<code>re.match(r'\d+', '123abc')</code>
<code>re.fullmatch(pattern, string, flags=0)</code>	Проверяет полное совпадение всей строки	<code>re.fullmatch(r'\d+', '123')</code>
<code>re.findall(pattern, string, flags=0)</code>	Возвращает все совпадения как список	<code>re.findall(r'\d+', 'a1b2c3')</code>
<code>re.finditer(pattern, string, flags=0)</code>	Возвращает итератор объектов Match	<code>re.finditer(r'\d+', 'a1b2c3')</code>
<code>re.sub(pattern, repl, string, count=0, flags=0)</code>	Заменяет совпадения	<code>re.sub(r'\d+', 'X', 'a1b2c3')</code>
<code>re.subn(pattern, repl, string, count=0, flags=0)</code>	Заменяет и возвращает кортеж (результат, количество замен)	<code>re.subn(r'\d+', 'X', 'a1b2c3')</code>
<code>re.split(pattern, string, maxsplit=0, flags=0)</code>	Разбивает строку по паттерну	<code>re.split(r'\s+', 'a b c')</code>
<code>re.compile(pattern, flags=0)</code>	Компилирует регулярное выражение	<code>re.compile(r'\d+')</code>
<code>re.escape(string)</code>	Экранирует специальные символы	<code>re.escape('a.b*c')</code>

Метод	Описание	Пример использования
<code>re.purge()</code>	Очищает кеш регулярных выражений	<code>re.purge()</code>

Флаги (модификаторы)

Флаг	Константа	Описание	Пример
<code>re.IGNORECASE</code>	<code>re.I</code>	Игнорирует регистр	<code>re.search(r'hello', 'HELLO', re.I)</code>
<code>re.MULTILINE</code>	<code>re.M</code>	<code>^</code> и <code>\$</code> работают для каждой строки	<code>re.findall(r'^line', text, re.M)</code>
<code>re.DOTALL</code>	<code>re.S</code>	<code>.</code> включает символы новой строки	<code>re.search(r'a.b', 'a\nb', re.S)</code>
<code>re.VERBOSE</code>	<code>re.X</code>	Позволяет использовать комментарии и пробелы	<code>re.compile(r'(?x)\d+ # число')</code>
<code>re.ASCII</code>	<code>re.A</code>	<code>\w</code> , <code>\b</code> , <code>\s</code> и <code>\d</code> работают только с ASCII	<code>re.search(r'\w+', 'café', re.A)</code>
<code>re.LOCALE</code>	<code>re.L</code>	Использует текущую локаль (устарел)	<code>re.search(r'\w+', text, re.L)</code>
<code>re.UNICODE</code>	<code>re.U</code>	Включает Unicode (по умолчанию в Python 3)	<code>re.search(r'\w+', 'café', re.U)</code>

Флаг	Константа	Описание	Пример
re.DEBUG	-	Выводит отладочную информацию	re.compile(r'\d+', re.DEBUG)

Методы объекта Match

Метод	Описание	Пример
match.group(num=0)	Возвращает совпадение группы	match.group(1)
match.groups()	Возвращает все группы как кортеж	match.groups()
match.groupdict()	Возвращает именованные группы как словарь	match.groupdict()
match.start(group=0)	Возвращает начальную позицию	match.start()
match.end(group=0)	Возвращает конечную позицию	match.end()
match.span(group=0)	Возвращает кортеж (начало, конец)	match.span()
match.expand(template)	Расширяет шаблон с группами	match.expand(r'\1-\2')
match.string	Исходная строка	match.string
match.re	Объект регулярного выражения	match.re
match.pos	Позиция начала поиска	match.pos

Метод	Описание	Пример
<code>match.endpos</code>	Позиция конца поиска	<code>match.endpos</code>
<code>match.lastindex</code>	Индекс последней захваченной группы	<code>match.lastindex</code>
<code>match.lastgroup</code>	Имя последней захваченной группы	<code>match.lastgroup</code>

Методы скомпилированного объекта Pattern

Метод	Описание	Пример
<code>pattern.search(string, pos=0, endpos=len(string))</code>	Поиск в строке	<code>compiled.search('text')</code>
<code>pattern.match(string, pos=0, endpos=len(string))</code>	Совпадение с начала	<code>compiled.match('text')</code>
<code>pattern.fullmatch(string, pos=0, endpos=len(string))</code>	Полное совпадение	<code>compiled.fullmatch('text')</code>
<code>pattern.findall(string, pos=0, endpos=len(string))</code>	Все совпадения	<code>compiled.findall('text')</code>
<code>pattern.finditer(string, pos=0, endpos=len(string))</code>	Итератор совпадений	<code>compiled.finditer('text')</code>
<code>pattern.sub(repl, string, count=0)</code>	Замена	<code>compiled.sub('new', 'text')</code>
<code>pattern.subn(repl, string, count=0)</code>	Замена с подсчетом	<code>compiled.subn('new', 'text')</code>
<code>pattern.split(string, maxsplit=0)</code>	Разбиение строки	<code>compiled.split('text')</code>

Метод	Описание	Пример
<code>pattern.pattern</code>	Исходный паттерн	<code>compiled.pattern</code>
<code>pattern.flags</code>	Флаги компиляции	<code>compiled.flags</code>
<code>pattern.groups</code>	Количество групп	<code>compiled.groups</code>
<code>pattern.groupindex</code>	Словарь именованных групп	<code>compiled.groupindex</code>

Специальные последовательности для замены

Последовательность	Описание	Пример
<code>\g<number></code>	Группа по номеру	<code>re.sub(r'(\d+)', r'\g<1>', text)</code>
<code>\g<name></code>	Именованная группа	<code>re.sub(r'(?P<num>\d+)', r'\g<num>', text)</code>
<code>\1, \2, ...</code>	Обратная ссылка на группу	<code>re.sub(r'(\w+) (\w+)', r'\2 \1', text)</code>
<code>\g<0></code>	Вся совпадающая строка	<code>re.sub(r'\d+', r'[\g<0>]', text)</code>
<code>\\</code>	Литеральный обратный слеш	<code>re.sub(r'\d+', r'\\', text)</code>

Исключения

Исключение	Описание
re.error	Ошибка компиляции регулярного выражения

Примеры использования различных методов

```
import re

# Компиляция паттерна
pattern = re.compile(r'(\d{4})-(\d{2})-(\d{2})', re.IGNORECASE)

# Поиск
match = pattern.search('Date: 2023-12-25')
if match:
    print(match.group(0)) # 2023-12-25
    print(match.groups()) # ('2023', '12', '25')

# Замена с использованием групп
text = "Born: 1990-05-15, Died: 2020-10-30"
result = pattern.sub(r'\3/\2/\1', text)
print(result) # Born: 15/05/1990, Died: 30/10/2020

# Использование флагов
pattern_verbose = re.compile(r'''
    (\d{4})    # год
    -          # разделитель
    (\d{2})    # месяц
    -          # разделитель
    (\d{2})    # день
''', re.VERBOSE)

# Разбиение строки
text = "apple,banana;orange:grape"
fruits = re.split(r'[,:;]', text)
print(fruits) # ['apple', 'banana', 'orange', 'grape']

# Поиск всех совпадений
numbers = re.findall(r'\d+', 'a1b22c333d4444')
print(numbers) # ['1', '22', '333', '4444']

# Использование именованных групп
pattern_named = re.compile(r'(?P<year>\d{4})-(?P<month>\d{2})-(?P<day>\d{2})')
match = pattern_named.search('2023-12-25')
if match:
    print(match.groupdict()) # {'year': '2023', 'month': '12', 'day': '25'}
```