

# Язык программирования



python

Лекция № 9. Регулярные выражения в Python

Евгений Сергеевич Чухланцев



Регулярное выражение `[A-Za-z]\w+` означает, что первый символ должен быть алфавитным, т.е. должен относиться к диапазону A-Z или a-z, а за ним следует по крайней мере один (+) алфавитно-цифровой символ (`\w`).

Обозначение	Описание	Пример регулярно-го выражения
<b>Знаки</b>		
<code>literal</code>	Строка содержит символьный литерал <code>literal</code>	<code>foo</code>
<code>re1 re2</code>	Строка содержит регулярные выражения <code>re1</code> или <code>re2</code>	<code>foo bar</code>
<code>.</code>	Соответствует <i>любому символу</i> (кроме <code>\n</code> )	<code>b.b</code>
<code>^</code>	Соответствует <i>началу строки</i>	<code>^Dear</code>
<code>\$</code>	Соответствует <i>концу строки</i>	<code>/bin/*sh\$</code>
<code>*</code>	Предшествующее регулярное выражение встречается в строке <i>любое количество</i> (или не встречается вообще)	<code>[A-Za-z0-9]*</code>
<code>+</code>	Предшествующее регулярное выражение встречается в строке <i>не менее одного раза</i>	<code>[a-z]+\ .com</code>
<code>?</code>	Предшествующее регулярное выражение встречается в строке <i>только один раз или вообще в ней не встречается</i>	<code>goo?</code>

{N}	Предшествующее регулярное выражение встречается в строке N раз	[0-9]{3}
{M,N}	Предшествующее регулярное выражение встречается в строке от M до N раз	[0-9]{5,9}
[...]	Соответствует любому отдельному символу из класса символов	[aeiou]
[..x-y..]	Соответствует любому отдельному символу в диапазоне от x до y	[0-9], [A-Za-z]
[^...]	Не соответствует ни одному символу из класса символов, включая любые диапазоны, если они заданы	[^aeiou], [^A-Za-z0-9_]
(* + ? { })?	Предусматривает применение "нежадных" версий приведенных выше знаков вхождения/повторения (*, +?, {})	. *?[a-z]
(...)	Соответствует регулярному выражению, заключенному в круглые скобки, и сохраняет его в памяти как подгруппу	([0-9]{3})?, f(oo u)bar

## Специальные символы

<code>\d</code>	Соответствует любой <i>десятичной цифре</i> , так же, как [0-9] ( <code>\D</code> является обратным по отношению к <code>\d</code> : не соответствует ни одной цифре)	<code>data\d+.txt</code>
<code>\w</code>	Соответствует любому <i>алфавитно-цифровому символу</i> , эквивалентно [A-Za-z0-9_] ( <code>\W</code> является обратным по отношению к <code>\w</code> )	<code>[A-Za-z_] \w+</code>
<code>\s</code>	Соответствует любому <i>пробельному символу</i> , эквивалентно [ <code>\n\t\r\v\f</code> ] ( <code>\S</code> является обратным по отношению к <code>\s</code> )	<code>of\sthe</code>
<code>\b</code>	Позиция, соответствующая <i>границе слова</i> ( <code>\B</code> является обратным по отношению к <code>\b</code> ),	<code>\bThe\b</code>
<code>\N</code>	Соответствует <i>сохраненной подгруппе N</i> (см. (...) выше)	<code>price: \16</code>
<code>\c</code>	Буквально соответствует <i>любому специальному символу c</i> (т.е. игнорирует особое значение этого литерала)	<code>\., \\, \*</code>
<code>\A (\Z)</code>	Позиция, соответствующая <i>началу (концу) строки</i> (см. также <code>^</code> и <code>\$</code> выше)	<code>\ADear</code>

## Расширенная система обозначений

(?iLmsux)	Внедряет один или несколько специальных "флагов" непосредственно в регулярное выражение (применяется вместо способа, основанного на использовании функции/метода)	(?x), (?im)
(?:...)	Обозначает группу, содержимое которой <i>не сохраняется в памяти</i>	(?:\w+\.)*
(?P<name>...)	Обозначает группу, заданную именем, а не числовым идентификатором	(?P<data>)
(?P=name)	Сопоставляется с текстом, который был перед этим сгруппирован оператором (?P<name>) в той же строке	(?P=data)
(?#...)	Задаёт примечание; все содержимое примечания игнорируется	(?#comment)
(?=...)	Обеспечивает сопоставление, если далее следует шаблон ..., без сохранения в памяти сопоставленной части входной строки; эта операция именуется <i>положительной опережающей проверкой</i> (positive lookahead assertion)	(?= .com)

(?!...)	Обеспечивает сопоставление, если далее не следует шаблон . . . , без сохранения в памяти сопоставленной части входной строки; эта операция именуется <i>отрицательной опережающей проверкой</i> (negative lookahead assertion)	(?! .net)
(?<=...)	Обеспечивает сопоставление, если далее находится шаблон . . . , без сохранения в памяти сопоставленной части входной строки; эта операция именуется <i>положительной ретроспективной проверкой</i> (positive lookbehind assertion)	(?<=800-)
(?<!...)	Обеспечивает сопоставление, если далее не находится шаблон . . . , без сохранения в памяти сопоставленной части входной строки; эта операция именуется <i>отрицательной ретроспективной проверкой</i> (positive lookbehind assertion)	(?<!192\.168\.)
(? (id/name) Y N)	Обеспечивает сопоставление регулярного выражения Y по условию, если группа с указанным идентификатором или именем существует; в противном случае возвращает N; часть  N является необязательной	(? (1) y x)

В регулярных выражениях операция чередования обозначается с помощью символа канала ( | ), который представлен на клавиатуре вертикальной чертой (pipeline symbol). Символ канала используется для отделения друг от друга разных регулярных выражений.

---

**Шаблон регулярного выражения****Сопоставленные строки**

---

at|home

at, home

r2d2|c3po

r2d2, c3po

bat|bet|bit

bat, bet, bit

---

Знак точки ( . ) обеспечивает сопоставление с любым отдельным символом, кроме \n. Знак точки сопоставляется с любой буквой, цифрой, пробельным символом (не включая " \n"), печатаемым или непечатаемым знаком.

---

**Шаблон регулярного выражения****Сопоставленные строки**

---

`f.o`

В этом примере точка между "f" и "o" представляет любой символ, например, как в строках fao, f9o, f#o и т.д.

`..`

Любая пара символов

`.end`

Любой символ перед строкой end

---

Для сопоставления с шаблоном, начиная с начала строки, необходимо использовать знак вставки (^) или специальный символ \A (прописная буква "A", которая следует за обратной косой чертой). Последний вариант применяется в основном на компьютерах с клавиатурой, на которой отсутствует знак вставки (такой как международная клавиатура). Аналогичным образом знак доллара (\$) или специальный символ \Z применяется для сопоставления с шаблоном, начиная с конца строки.

<b>Шаблон регулярного выражения</b>	<b>Сопоставленные строки</b>
<code>^From</code>	Любая строка, которая начинается с подстроки <code>From</code>
<code>/bin/tcsh\$</code>	Любая строка, которая заканчивается подстрокой <code>/bin/tcsh</code>
<code>^Subject: hi\$</code>	Любая строка, полностью совпадающая со строкой <code>Subject: hi</code>

<b>Шаблон регулярного выражения</b>	<b>Сопоставленные строки</b>
<code>the</code>	Любая строка, содержащая подстроку <code>the</code>
<code>\bthe</code>	Любое слово, которое начинается с подстроки <code>the</code>
<code>\bthe\b</code>	Сопоставляется только со словом <code>the</code>
<code>\Bthe</code>	Любая строка, которая содержит подстроку <code>the</code> , но с этой подстроки не начинается слово

Безусловно, знак точки хорошо подходит для тех случаев, когда необходимо обеспечить сопоставление с любым знаком, но иногда требуется провести сопоставление лишь с конкретным набором символов. По этой причине была предусмотрена возможность применения в шаблонах знаков квадратных скобок ( [ ] ).

<b>Шаблон регулярного выражения</b>	<b>Сопоставленные строки</b>
<code>b[aeiu]t</code>	<code>bat, bet, bit, but</code>
<code>[cr][23][dp][o2]</code>	Строка из четырех символов: в начале следует "с" или "r", затем — "2" или "3", после этого — "d" или "p" и, наконец, "o" или "2". Примерами могут служить подстроки <code>c2do, r3p2, r2d2, c3po</code> и т.д.

Квадратные скобки позволяют задавать не только наборы из отдельных символов, но и диапазоны символов. Для обозначения диапазона символов применяется пара символов, заключенных в квадратные скобки, между которыми проставлен знак дефиса. В качестве примера можно указать диапазоны A-Z, a-z и 0-9, применяемые для обозначения прописных букв, строчных букв и цифровых знаков соответственно.

Шаблон регулярного выражения	Сопоставленные строки
<code>z.[0-9]</code>	Буква "z", за которой следует любой символ, за которым следует одна цифра
<code>[r-u][env-y][us]</code>	Буква "r", "s", "t" или "u", за которой следует буква "e", "n", "v", "w", "x" или "y", за которой следует буква "u" или "s"
<code>[^aeiou]</code>	Знак, отличный от гласной буквы ( <b>Упражнение.</b> Можно ли в этой формулировке вместо "знак, отличный от гласной буквы" применить "знак согласной буквы"?)
<code>[^\t\n]</code>	Знак, отличный от знака табуляции или \n
<code>["-a]</code>	В системе ASCII — все символы, которые расположены между " " и "a", т.е. между знаками с порядковыми номерами 34 и 97

# Использование операторов

См. также

Шаблон регулярного выражения	Сопоставленные строки
<code>[dn]ot?</code>	Буквы "d" или "n", за которыми следует буква "o", а затем, самое большее, — одна буква "t". Таким образом, это регулярное выражение сопоставляется со словами do, no, dot, not
<code>0?[1-9]</code>	Любой цифровой символ, которому может предшествовать цифра "0". В качестве примера можно указать множество числовых обозначений месяцев от января до сентября, представленных с помощью одной или двух цифр
<code>[0-9]{15,16}</code>	Пятнадцать или шестнадцать цифр (например, номера кредитных карточек)
<code>&lt;/?[&gt;]+&gt;</code>	Строки, которые сопоставляются со всеми допустимыми (и недопустимыми) тегами HTML
<code>[KQRBNP][a-h][1-8]-[a-h][1-8]</code>	Допустимый шахматный ход в алгебраической нотации (только ходы; взятия фигур, шахи и др. не представлены); т.е. строки, которые начинаются с любой из букв "K", "Q", "R", "B", "N" или "P", обозначающих шахматные фигуры, за которой следует разделенная дефисами пара обозначений клеток шахматной доски от "a1" до "h8" (и всего, что находится между ними), в которой первое обозначение клетки указывает прежнее местонахождение фигуры, а второе — новое местонахождение

# Специальные символы, обозначающие наборы символов

---

## Шаблон регулярного выражения

## Сопоставленные строки

`\w+-\d+`

Алфавитно-цифровая строка и число, разделенные дефисом

`[A-Za-z]\w*`

Первый символ — алфавитный; следующие символы (если они присутствуют) могут быть алфавитно-цифровыми. Это регулярное выражение почти эквивалентно выражению, которое описывает множество допустимых идентификаторов Python (см. примеры)

`\d{3}-\d{3}-\d{4}`

Номера телефонов в формате, принятом в США, с префиксом кода города, как в примере 800-555-1212

`\w+@\w+\.com`

Простые адреса электронной почты в форме XXX@YYY.com

---

Пара круглых скобок ( ( ) ) в регулярном выражении позволяет решить любую из следующих задач (или обе эти задачи).

- Выполнить группирование регулярных выражений.
- Провести сопоставление с подгруппами.

---

**Шаблон регулярного выражения****Сопоставленные строки**

---

`\d+(\.\d*)?`

Строки, представляющие числа с плавающей точкой в простом формате, т.е. любое количество цифр, за которым следует необязательная отдельная десятичная точка, затем нуль или большее количество цифровых символов, например, "0.004", "2", "75." и т.д.

`(Mr?s?\. )?[A-Z]  
[a-z]* [A-Za-z-]+`

Имя и фамилия, в которых имя может быть представлено сокращенно (должно начинаться с прописной буквы, за которой следуют строчные буквы остальной части имени, если она задана), фамилия представлена полностью, а впереди находится необязательное обращение, "Mr.", "Mrs.", "Ms." или "M.". При указании фамилии предусмотрена возможность задавать несколько ее компонентов с использованием дефисов, прописных и строчных букв

---

Шаблон регулярного выражения	Определение расширенных обозначений
<code>(?:\w+\.)*</code>	Строки, которые оканчиваются точкой, такие как "google.", "twitter.", "facebook.". Обнаруженные при этом сопоставления не сохраняются для дальнейшего использования и не могут быть извлечены
<code>(?#comment)</code>	Эта конструкция не определяет сопоставление и применяется исключительно для ввода комментария
<code>(?=\.com)</code>	Сопоставление происходит, только если далее следует подстрока ".com"; обработчик не изымает ни одной части целевой строки.
<code>(?!\.net)</code>	Сопоставление происходит, только если далее не следует подстрока ".net"
<code>(?&lt;=800-)</code>	Сопоставление происходит, только если данной строке предшествует подстрока "800-"; эта конструкция может применяться, допустим, для поиска номеров телефонов. В этом случае обработчик не изымает ни одной части входной строки
<code>(?&lt;!192\.168\.)</code>	Сопоставление происходит, только если данной строке не предшествует подстрока "192.168.". Эта конструкция может применяться, например, для фильтрации IP-адресов класса C
<code>(?(1)y x)</code>	Если сопоставленная группа 1 ( $\Lambda 1$ ) существует, провести сопоставление с y; в противном случае — с x

# Спасибо за внимание !

## Домашнее задание

Продолжаем читать книгу: Лутц М. "Изучаем Python" (4-е издание, в 2-х томах) (2011, PDF) !