

Регулярные выражения

1. Язык регулярных выражений в Python
2. Модуль для работы с регулярными выражениями

Литература

1. Марк Саммерфилд. Программирование на Python 3. (С. 525-552)
2. <https://docs.python.org/3/library/re.html>
3. <http://www.rexegg.com/>
4. <https://habr.com/ru/post/349860/>

- 1. Язык регулярных выражений в Python**
2. Модуль для работы с регулярными выражениями

Проблема!

- Тип Строки позволяет выполнять над строками различные операции, в том числе поиск, замену, вставку и удаление подстрок.
- Но есть классы задач по обработке символьной информации, где стандартных возможностей явно не хватает: проверить текст на соответствие определенному шаблону (например, адресу электронной почты), поиск по шаблону `pet*. ???` и др.

Регулярное выражение – что это?

- Регулярное выражение - это компактная форма записи набора строк, отвечающих требованиям регулярного выражения.
- Регулярное выражение - это текстовая строка, которая описывает **шаблон**, который используется механизмом регулярных выражений для поиска фрагмента текста (или позиций) в исходном тексте, как правило, для проверки, поиска, замены или разделения.
- В состав Python входит модуль **re**, с помощью которого можно создавать и использовать регулярные выражения.

Регулярное выражение	Значение
simple text	В точности текст «simple text»
<code>\d{5}</code>	Последовательности из 5 цифр, <code>\d</code> означает любую цифру, <code>{5}</code> — ровно 5 раз
<code>\d\d/\d\d/\d{4}</code>	Даты в формате ДД/ММ/ГГГГ (и прочие фрагменты, например, 98/76/5432)
<code>\b\w{3}\b</code>	Слова в точности из трёх букв <code>\b</code> означает границу слова (с одной стороны буква, а с другой — нет) <code>\w</code> — любая буква, <code>{3}</code> — ровно три раза
<code>[-+]? \d+</code>	Целое число, например, 7, +17, -42, 0013 (возможны ведущие нули), <code>[-+]?</code> — либо -, либо +, либо пусто

Регулярное выражение – для чего?

□ *Проверка:* проверка соответствия фрагментов текста некоторым критериям, например, наличие символа обозначения валюты и последующих за ним цифр.

□ *Поиск:* поиск подстрок, которые могут иметь несколько форм,

например, поиск подстрок «pet.png», «pet.jpg», «pet.jpeg» или «pet.svg», чтобы при этом не обнаруживались подстроки «carpet.png» и подобные ей.

Регулярное выражение – для чего?

- *Поиск и замена:* замена всего, что совпадает с регулярным выражением, и замена на указанную строку, например, поиск подстроки «устройство передвижения, движимое мускульной силой» и замена подстрокой «велосипед».
- *Разбиение строк:* разбиение строки по точкам совпадения с регулярным выражением, например, разбиение строки по подстроке «: » или «=».

Спецсимволы

✓ . ^ \$ * + ? { } [] \ | ()

✓ Для написания их просто как символов требуется их *экранировать*, для чего нужно поставить перед ними знак \.

Осторожно - регулярные выражения!

- ✓ Использовать только там, где они действительно приносят пользу, а не вред. Плохо написанные регулярные выражения работают медленно.
- ✓ Их сложно читать, особенно если *регулярка* написана не лично тобой пять минут назад.
- ✓ Даже небольшое изменение задачи (того, что требуется найти) приводит к значительному изменению выражения.
- ✓ Это *write only code* (код, который только пишут с нуля, но не читают и не правят).

Регулярные выражения: как?

- Регулярное выражение - это шаблон, по которому выполняется поиск соответствующего фрагмента текста.
- Язык описания регулярных выражений состоит из символов двух видов: **обычных символов** и **метасимволов**.
- Обычный символ представляет в выражении сам себя, а метасимвол - некоторый *класс символов*.

Наиболее употребительные МЕТАСИМВОЛЫ:

Use	To match any character
[set]	In that set
[^set]	Not in that set
[a-z]	In the a-z range
[^a-z]	Not in the a-z range
.	Any except \n (new line)
\char	Escaped special character

Примеры

Выражение	Фрагмент
<i>c.t</i>	<i>cat, cut, c#t, c{t</i>
<i>c[aui]t</i>	<i>cat, cut, cit</i>
<i>c[a-c]t</i>	<i>cat, cbt, cct</i>
<i>c[^aui]t</i>	<i>cbt, cct, c2t и т.д.</i>
<i>c[^a-c]t</i>	<i>cdt, cet, c%t</i>

УПРАВЛЯЮЩИЕ СИМВОЛЫ

Use	To match	Unicode
<code>\t</code>	Horizontal tab	<code>\u0009</code>
<code>\v</code>	Vertical tab	<code>\u000B</code>
<code>\b</code>	Backspace	<code>\u0008</code>
<code>\e</code>	Escape	<code>\u001B</code>
<code>\r</code>	Carriage return	<code>\u000D</code>
<code>\f</code>	Form feed	<code>\u000C</code>
<code>\n</code>	New line	<code>\u000A</code>
<code>\a</code>	Bell (alarm)	<code>\u0007</code>
<code>\c char</code>	ASCII control character	—

CHARACTER CLASSES

Use	To match character
<code>\p{ctgry}</code>	In that Unicode category or block
<code>\P{ctgry}</code>	Not in that Unicode category or block
<code>\w</code>	Word character. Любая буква (то, что может быть частью слова), а также цифры и <code>_</code>
<code>\W</code>	Non-word character. Любая не-буква, не-цифра и не подчёркивание
<code>\d</code>	Decimal digit
<code>\D</code>	Not a decimal digit. Любой символ, кроме цифры
<code>\s</code>	White-space character. Любой пробельный символ (пробел, табуляция, конец строки и т.п.)
<code>\S</code>	Non white space char. Любой непобельный символ

Примеры

<i>Выражение</i>	<i>Фрагмент</i>
<i>c wt</i>	<i>cbt, cct, c2t</i> и т.д., но не <i>c%t, c{t</i> и т.д.
<i>c Wt</i>	<i>c%t, c{t, c. t</i> , но не <i>cbt, cct, c2t</i> и т.д.
<i> s w w w s</i>	Любое слово из трех букв, окруженное пробельными символами.
<i> s S S S s</i>	Любые три непробельных символа, окруженные пробельными.
<i>c dt</i>	<i>c1t, c2t, c3t</i>
<i>c Dt</i>	Не <i>c1t, c2t, c3t</i> и т.д.

УТОЧНЯЮЩИЕ СИМВОЛЫ

Use	To specify position
<code>^</code>	At start of string or line
<code>\A</code>	At start of string (многострочная)
<code>\z</code>	At end of string
<code>\Z</code>	At end (or before <code>\n</code> at end) of string
<code>\$</code>	At end (or before <code>\n</code> at end) of string or line
<code>\G</code>	Where previous match ended
<code>\b</code>	On word boundary
<code>\B</code>	Not on word boundary

Повторители или квантификаторы

Greedy	Lazy	Matches
*	*?	0 or more times
+	+?	1 or more times
?	??	0 or 1 time
{n}	{n}?	Exactly n times
{n,}	{n,}?	At least n times
{n,m}	{n,m}?	From n to m times

Про квантификаторы

По умолчанию квантификаторы *жадные* — захватывают максимально возможное число символов.

Добавление *?* делает их *ленивыми*, они захватывают минимально возможное число символов.

Выражение		Фрагмент	
ca*t		ct, cat, caat, caaat	
ca+t		cat, caat, caaat	
ca?t		ct, cat	
ca{3}t		caaat	
(cat){2}		catcat	
ca{3, }t		caaat, caaaat, caaaaaaat	
(cat){2, }		catcat, catcatcat	
ca{2, 4}t		caat, caaat, caaaat	

1. Язык регулярных выражений в Python

2. Модуль для работы с регулярными выражениями

pattern, string

- Два параметра в функциях модуля регулярных выражений *re* :
 - Шаблон регулярного выражения для определения текста.
 - Текст, который будет проанализирован на соответствие шаблону регулярного выражения.

Функция	Её смысл
<code>re.search(pattern, string)</code>	Найти в строке <code>string</code> первую строчку, подходящую под шаблон <code>pattern</code> ;
<code>re.fullmatch(pattern, string)</code>	Проверить, подходит ли строка <code>string</code> под шаблон <code>pattern</code> ;
<code>re.split(pattern, string, maxsplit=0)</code>	Аналог <code>str.split()</code> , только разделение происходит по подстрокам, подходящим под шаблон <code>pattern</code> ;

`re.findall(pattern, string)`

Найти в строке `string` все непересекающиеся шаблоны `pattern`;

`re.finditer(pattern, string)`

Итератор по всем непересекающимся шаблонам `pattern` в строке `string` (выдаются `match-объекты`);

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

Заменить в строке `string` все непересекающиеся шаблоны `pattern` на `repl`;

Примеры регулярных выражений:

- 1) Слово math **math, «math»**
- 2) Номер телефона в формате xxx-xx-xx
@"\d\d\d-\d\d-\d\d" или r"\d{3}(-\d\d){2}
- 3) Номер автомобиля - **r"[A-Z]\d{3}[A-Z]{2}\d{2,3}RUS"**

Задание. Запишите регулярное выражение, соответствующее:

1. дате в формате дд.мм.гг или дд.мм.гггг
2. времени в формате чч.мм или чч:мм
3. целому числу (со знаком и без)
4. вещественному числу (со знаком и без, с дробной частью и без, с целой частью и без)
5. Найдите все натуральные числа (возможно, окружённые буквами);
6. Найдите все «слова», написанные капсом (то есть строго заглавными), возможно внутри настоящих слов (aaa**БББ**vvv);
7. Найдите слова, в которых есть русская буква, а когда-нибудь за ней цифра;
8. Найдите все слова, начинающиеся с русской или латинской большой буквы (`\b` — граница слова);
9. Найдите слова, которые начинаются на гласную (`\b` — граница слова).

```
match = re.findall(r'[0-9]+\w', r'Телефон 123-12-12')
```