

Regex + scrapping

1.10.2020

regex

Регулярные выражения - язык описания шаблонов (patterns) для извлечения информации из текста.

Регулярка	Её смысл
<code>simple text</code>	В точности текст «simple text»
<code>\d{5}</code>	Последовательности из 5 цифр \d означает любую цифру {5} — ровно 5 раз
<code>\d\d/\d\d/\d{4}</code>	Даты в формате ДД/ММ/ГГГГ (и прочие куски, на них похожие, например, 98/76/5432)
<code>\b\w{3}\b</code>	Слова в точности из трёх букв \b означает границу слова (с одной стороны буква, а с другой — нет) \w — любая буква, {3} — ровно три раза
<code>[-+]? \d+</code>	Целое число, например, 7, +17, -42, 0013 (возможны ведущие нули) [-+]? — либо -, либо +, либо пусто \d+ — последовательность из 1 или более цифр
<code>[-+]?(?:\d+(?:\.\d*)? \.\d+)(?:[eE][-+]? \d+)?</code>	Действительное число, возможно в экспоненциальной записи Например, 0.2, +5.45, -.4, 6e23, -3.17E-14. См. ниже картинку.

Write only code

```
(?:[a-z0-9!#$%&'*/=?^_`{|}~-]+(?:\. [a-z0-9!#$%&'*/=?^_`{|}~-]+)*|"(?:[\x01-\x08\x0b\x0c\x0e-\x1f\x21\x23-\x5b\x5d-\x7f]|\\[\x01-\x09\x0b\x0c\x0e-\x7f])*")@(?:(?:[a-z0-9](?:[a-z0-9-]*[a-z0-9])?\.\.)+[a-z0-9](?:[a-z0-9-]*[a-z0-9])?|\\[(?:[0-5]|2[0-4])[0-9]|[01]?[0-9][0-9]?)\.\.){3}(?:25[0-5]|2[0-4][0-9]|[01]?[0-9][0-9]?|[a-z0-9-]*[a-z0-9]:(?:[\x01-\x08\x0b\x0c\x0e-\x1f\x21-\x5a\x53-\x7f]|\\[\x01-\x09\x0b\x0c\x0e-\x7f])+\))\)
```

regex

- Поиск точного совпадения.
- Поиск шаблонного совпадения.
- Возможность введения переменных.
- Жадность.
- Метасимволы.

regex

Точное совпадение:

Строка на вход «aaa bbb ccc»

шаблон: r"a"

на выходе: [a, a, a]

regex

[] - СИМВОЛ «ИЛИ»

Строка на вход «aaa bbb ccc»

шаблон: r"[ab]"

на выходе: [a, a, a, b, b, b]

regex

() - группировка

Строка на вход «aaa bbb ccc»

шаблон: r"(aaa) bbb"

на выходе: [aaa]

regex

Если группа не нужна, то в группе ставим :?

Строка на вход «aaa bbb ccc»

шаблон: `r"(aaa) (?:bbb)"`

на выходе: `[aaa]`

regex

Метасимволы

- `\d` - все цифры [0-9]
- `\w` - все буквы [a-Za-Я]
- `\s` - пробельные символы
- `\b` - слова

regex

Метасимволы

- `\D` - Не цифры `[0-9]`
- `\W` - Не буквы `[a-zA-Z]`
- `\S` - Не пробельные символы
- `\b` - Не слова

regex

Шаблон	Описание	Пример	Применяем к тексту
.	Один любой символ, кроме новой строки \n.	м.л.ко	<u>молоко</u> , <u>малако</u> , И <u>м0л0ко</u> Ихлеб
\d	Любая цифра	су\d\d	<u>СУ35</u> , <u>СУ111</u> , АЛ <u>СУ14</u>
\D	Любой символ, кроме цифры	926\D123	<u>926</u>) <u>123</u> , <u>1926-</u> <u>1234</u>
\s	Любой пробельный символ (пробел, табуляция, конец строки и т.п.)	бор\sода	<u>бор_ода</u> , <u>бор</u> <u>ода</u> , борода
\S	Любой непробельный символ	\S123	<u>X123</u> , <u>я123</u> , <u>!123</u> 456, 1 + 123456
\w	Любая буква (то, что может быть частью слова), а также цифры и _	\w\w\w	<u>Год</u> , <u>f_3</u> , <u>qwert</u>
\W	Любая не-буква, не-цифра и не подчёркивание	com\W	<u>com!</u> , <u>com?</u>

regex

[..]	Один из символов в скобках, а также любой символ из диапазона a-b	[0-9][0-9A-Fa-f]	12 , 1F , 4B
[^..]	Любой символ, кроме перечисленных	<[^>]>	<1> , <a> , <>>
<code>\d≈[0-9],</code> <code>\D≈[^0-9],</code> <code>\w≈[0-9a-zA-Z</code> <code>a-яA-ЯёЁ],</code> <code>\s≈[</code> <code>\f\n\r\t\v]</code>	Буква "ё" не включается в общий диапазон букв! Вообще говоря, в \d включается всё, что в юникоде помечено как «цифра», а в \w — как буква. Ещё много всего!		
[abc-], [-1]	если нужен минус, его нужно указать последним или первым		
[* [(+\\)]\t]	внутри скобок нужно экранировать только] и \		
\b	Начало или конец слова (слева пусто или не-буква, справа буква и наоборот). В отличие от предыдущих соответствует позиции, а не символу	\bвал	вал , перевал, Перевалка
\B	Не граница слова: либо и слева, и справа буквы, либо и слева, и справа НЕ буквы	\Bвал	пере вал , вал, Пере вал ка
		\Bвал\B	перевал, вал, Пере вал ка

[Check here](#)

regex

Квантификаторы

Шаблон	Описание	Пример
<code>{n}</code>	Ровно n повторений	<code>\d{4}</code>
<code>{m,n}</code>	От m до n повторений включительно	<code>\d{2,4}</code>
<code>{m,}</code>	Не менее m повторений	<code>\d{3,}</code>
<code>{,n}</code>	Не более n повторений	<code>\d{,2}</code>

regex

Пример кватификатора

Строка на вход «На дворе - трава, на траве - дрова.»

шаблон: `r"[а-я]{3,}"`

на выходе: ['дворе', 'трава', 'траве', 'дрова']

regex

Синонимы кватификаторов

Шаблон	Описание	Пример
<code>?</code>	Ноль или одно вхождение, синоним $\{0,1\}$	<code>валы?</code>
<code>*</code>	Ноль или более, синоним $\{0,\}$	<code>су\d*</code>
<code>+</code>	Одно или более, синоним $\{1,\}$	<code>а\)+</code>

regex

Жадность - по умолчанию регулярные выражения захватывают максимум символов, которые помещаются под шаблон.

Шаблон	Описание	Пример
?	По умолчанию квантификаторы <i>жадные</i> — захватывают максимально возможное число символов.	\(.\)
+?		\(.*?\)
??		
{m,n}?	Добавление ? делает их <i>ленивыми</i> , они захватывают минимально возможное число символов	
{,n}?		
{m,}?		

Python

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> ;
<code>re.findall(pattern, string)</code>	Найти в строке <code>string</code> все непересекающиеся шаблоны <code>pattern</code> ;
<code>re.finditer(pattern, string)</code>	Итератор всем непересекающимся шаблонам <code>pattern</code> в строке <code>string</code> (выдаются <code>match</code> -объекты);
<code>re.sub(pattern, repl, string, count=0)</code>	Заменить в строке <code>string</code> все непересекающиеся шаблоны <code>pattern</code> на <code>repl</code> ;

scraping

Скрапинг - процесс сбора информации с веб страниц.

Парсинг - процесс обработки текста, часто подразумевается разбор текста на составные части.

scraping

- requests + bs4 - низкоуровневый поиск.
- selenium - имитация клиента.
- scrapy - готовое решение для скрапинга, основано на асинхронных запросах.

scraping

Из чего состоит веб страница

- html
- CSS
- js

scraping

Возможные проблемы:

- Скорость...
- Тайминги, фриззы.
- AJAX страницы.