

Шаблоны, соответствующие одному символу

Шаблон	Описание	Пример	Применяем к тексту
.	Один любой символ, кроме новой строки \n.	м.л.ко	молоко , малако , И мОлОко Ихлеб
\d	Любая цифра	су\d\d	СУ35 , СУ111 , АЛ СУ14
\D	Любой символ, кроме цифры	926\D123	926) 123 , 1 926-123 4
\s	Любой пробельный символ (пробел, табуляция, конец строки и т.п.)	бор\sода	бор ода , бор ода , борода
\S	Любой непробельный символ	\S123	X123 , я123 , !123 456, 1 + 123456
\w	Любая буква (то, что может быть частью слова), а также цифры и _	\w\w\w	Год , f_3 , qwert
\W	Любая не-буква, не-цифра и не подчёркивание	com\W	com! , com?
[...]	Один из символов в скобках, а также любой символ из диапазона a-b	[0-9][0-9A-Fa-f]	12 , 1F , 4B

Шаблон	Описание	Пример	Применяем к тексту
[^..]	Любой символ, кроме перечисленных	<[^>>	<1> , <a> , <>>
\ d=[0-9], \ D=[^0-9], \ w=[0-9a-zA-Z a-яA-яёЁ], \ s=[\f\n\r\t\v]	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

Квантификаторы (указание количества повторений)

Шаблон	Описание	Пример	Применяем к тексту
$\{n\}$	Ровно n повторений	<code>\d{4}</code>	1, 12, 123, <u>1234</u> , 12345
$\{m, n\}$	От m до n повторений включительно	<code>\d{2,4}</code>	1, <u>12</u> , <u>123</u> , <u>1234</u> , 12345
$\{m, \}$	Не менее m повторений	<code>\d{3,}</code>	1, 12, <u>123</u> , <u>1234</u> , <u>12345</u>
$\{, n\}$	Не более n повторений	<code>\d{,2}</code>	<u>1</u> , <u>12</u> , <u>123</u>
<code>?</code>	Ноль или одно вхождение, синоним $\{0, 1\}$	<code>валы?</code>	<u>вал</u> , <u>валы</u> , <u>вал</u> ов
<code>*</code>	Ноль или более, синоним $\{0, \}$	<code>су\d*</code>	<u>су</u> , <u>су1</u> , <u>су12</u> , ...
<code>+</code>	Одно или более, синоним $\{1, \}$	<code>а\)+</code>	<u>а</u>), <u>а</u>)), <u>а</u>))), б <u>а</u>))
<code>*?</code> <code>+?</code> <code>??</code> $\{m, n\}?$ $\{, n\}?$ $\{m, \}?$	По умолчанию квантификаторы <i>жадные</i> — захватывают максимально возможное число символов. Добавление <code>?</code> делает их <i>ленивыми</i> , они захватывают минимально возможное число символов	<code>\(.*\)</code> <code>\(.*?\)</code>	<u>(a + b) * (c + d) * (e + f)</u> <u>(a + b)</u> * (c + d) * (e + f)

Регулярки в питоне

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 по 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> ;

Пример использования всех основных функций

```
import re

match = re.search(r'\d\d\D\d\d', r'Телефон 123-12-12')
print(match[0] if match else 'Not found')
# -> 23-12

match = re.search(r'\d\d\D\d\d', r'Телефон 1231212')
print(match[0] if match else 'Not found')
# -> Not found

match = re.fullmatch(r'\d\d\D\d\d', r'12-12')
print('YES' if match else 'NO')
# -> YES

match = re.fullmatch(r'\d\d\D\d\d', r'T. 12-12')
print('YES' if match else 'NO')
# -> NO

print(re.split(r'\W+', 'Где, скажите мне, мои очки?!'))
# -> ['Где', 'скажите', 'мне', 'мои', 'очки', '']

print(re.findall(r'\d\d\.\d\d\.\d{4}',
                 r'Эта строка написана 19.01.2018, а могла бы и 01.09.2017'))
# -> ['19.01.2018', '01.09.2017']

for m in re.finditer(r'\d\d\.\d\d\.\d{4}', r'Эта строка написана 19.01.2018,
а могла бы и 01.09.2017'):
    print('Дата', m[0], 'начинается с позиции', m.start())
# -> Дата 19.01.2018 начинается с позиции 20
# -> Дата 01.09.2017 начинается с позиции 45

print(re.sub(r'\d\d\.\d\d\.\d{4}',
             r'DD.MM.YYYY',
             r'Эта строка написана 19.01.2018, а могла бы и 01.09.2017'))
# -> Эта строка написана DD.MM.YYYY, а могла бы и DD.MM.YYYY
```

Использование дополнительных флагов в питоне

Константа	Её смысл
<code>re.ASCII</code>	По умолчанию <code>\w</code> , <code>\W</code> , <code>\b</code> , <code>\B</code> , <code>\d</code> , <code>\D</code> , <code>\s</code> , <code>\S</code> соответствуют все юникодные символы с соответствующим качеством. Например, <code>\d</code> соответствуют не только арабские цифры, но и вот такие: <code>٠١٢٣٤٥٦٧٨٩</code> . <code>re.ASCII</code> ускоряет работу, если все соответствия лежат внутри ASCII.
<code>re.IGNORECASE</code>	Не различать заглавные и маленькие буквы. Работает медленнее, но иногда удобно
<code>re.MULTILINE</code>	Специальные символы <code>^</code> и <code>\$</code> соответствуют началу и концу каждой строки
<code>re.DOTALL</code>	По умолчанию символ <code>\n</code> конца строки не подходит под точку. С этим флагом точка — вообще любой символ

Простые шаблоны, соответствующие позиции

Шаблон	Описание	Пример	Применяем к тексту
<code>^</code>	Начало всего текста или начало строки текста, если <code>flag=re.MULTILINE</code>	<code>^Привет</code>	
<code>\$</code>	Конец всего текста или конец строки текста, если <code>flag=re.MULTILINE</code>	Будь здоров! <code>\$</code>	
<code>\A</code>	Строго начало всего текста		
<code>\Z</code>	Строго конец всего текста		
<code>\b</code>	Начало или конец слова (слева пусто или не-буква, справа буква и наоборот)	<code>\bвал</code>	<u>вал</u> , перевал, Перевалка
<code>\B</code>	Не граница слова: либо и слева, и справа буквы, либо и слева, и справа НЕ буквы	<code>\Bвал</code>	пере <u>вал</u> , вал, Пере <u>вал</u> ка
		<code>\Bвал\B</code>	перевал, вал, Пере <u>вал</u> ка

Шаблон

Применяем к тексту

(?:\w\w\d\d)+

Есть миг29а, ту15б. Некоторые делают даже миг29ту154ил86.

(?:\w+\d+)+

Есть миг29а, ту154б. Некоторые делают даже миг29ту154ил86.

(?:\+7|8)(?:-\d{2,3}){4}

+7-926-123-12-12, 8-926-123-12-12

(?:[Xx][аоеи]+)+

Муха — хахахехо, ну хааахоооохе, да хахахехохиии! Хам трамвайный.

\b(?:[Xx][аоеи]+)+\b

Муха — хахахехо, ну хааахоооохе, да хахахехохиии! Хам трамвайный.

Сложные шаблоны, соответствующие позиции (*lookaround* и Co)

Шаблон	Описание	Пример	Применяем к тексту
(?=...)	<i>lookahead assertion</i> , соответствует каждой позиции, сразу после которой начинается соответствие шаблону ...	Isaac (?=Asimov)	Isaac Asimov, Isaac other
(?!...)	<i>negative lookahead assertion</i> , соответствует каждой позиции, сразу после которой НЕ может начинаться шаблон ...	Isaac (?!Asimov)	Isaac Asimov, Isaac other
(?<=...)	<i>positive lookbehind assertion</i> , соответствует каждой позиции, которой может заканчиваться шаблон ... Длина шаблона должна быть фиксированной, то есть <code>abc</code> и <code>a b</code> — это ОК, а <code>a*</code> и <code>a{2,3}</code> — нет.	(?<=abc)def	abc def , bcdef
(?<!...)	<i>negative lookbehind assertion</i> , соответствует каждой позиции, которой НЕ может заканчиваться шаблон ...	(?<!abc)def	abcdef, bc def

lookaround на примере королей и императоров Франции

Людовик(?=VI) — Людовик, за которым идёт VI
КарлIV, КарлIX, КарлV, КарлVI, КарлVII, КарлVIII,
ЛюдовикIX, [ЛюдовикVI](#), [ЛюдовикVII](#), [ЛюдовикVIII](#), ЛюдовикX, ..., ЛюдовикXVIII,
ФилиппI, ФилиппII, ФилиппIII, ФилиппIV, ФилиппV, ФилиппVI

Людовик(?!VI) — Людовик, за которым идёт не VI

КарлIV, КарлIX, КарлV, КарлVI, КарлVII, КарлVIII,
[ЛюдовикIX](#), ЛюдовикVI, ЛюдовикVII, ЛюдовикVIII, [ЛюдовикX](#), ..., [ЛюдовикXVIII](#),
ФилиппI, ФилиппII, ФилиппIII, ФилиппIV, ФилиппV, ФилиппVI

(?<=Людовик)VI — «шестой», но только если Людовик

КарлIV, КарлIX, КарлV, КарлVI, КарлVII, КарлVIII,
ЛюдовикIX, Людовик[VI](#), Людовик[VII](#), Людовик[VIII](#), ЛюдовикX, ..., ЛюдовикXVIII,
ФилиппI, ФилиппII, ФилиппIII, ФилиппIV, ФилиппV, ФилиппVI

(?!Людовик)VI — «шестой», но только если не Людовик

КарлIV, КарлIX, КарлV, Карл[VI](#), Карл[VII](#), Карл[VIII](#),
ЛюдовикIX, ЛюдовикVI, ЛюдовикVII, ЛюдовикVIII, ЛюдовикX, ..., ЛюдовикX[VIII](#),
ФилиппI, ФилиппII, ФилиппIII, ФилиппIV, ФилиппV, Филипп[VI](#)

Шаблон	Комментарий	Применяем к тексту
(?<!\d)\d(?!\d)	Цифра, окружённая не-цифрами	Text ABC 123 A1B2C3!
(?<=#START#).*?(?=#END#)	Текст от #START# до #END#	text from #START# till #END#
\d+(?=_(?!_))	Цифра, после которой идёт ровно одно подчёркивание	12 _34__56
^(?:(?!boo).)*?\$	Строка, в которой нет boo (то есть нет такого символа, перед которым есть boo)	a foo and boo and zoo and others

Шаблон	Комментарий	Применяем к тексту
<code>^(?:(!boo)(?!foo).)*?\$</code>	Строка, в которой нет ни boo, ни foo	a foo and boo and zoo <u>and others</u>