



PYTHON — курс молодого бойца



RegEx — регулярные выражения в Python

👤 [Вадим В. Костерин](#)

🕒 [09/02/2020](#)

💬 [Оставьте комментарий](#)

🔔 *Сие есть перепечатка из Habra замечательной статьи в тему сайта [Регулярные выражения в Python от простого к сложному. Подробности, примеры, картинки, упражнения](#)*

Решил я давеча моим школьникам дать задачек на регулярные выражения для изучения. А к задачкам нужна какая-нибудь теория. И стал я искать хорошие тексты на русском. Пяток сносных нашёл, но всё не то. Что-то смято, что-то упущено. У этих текстов был не только *фатальный* недостаток. Мало картинок, мало примеров. И почти нет разумных задач. Ну неужели поиск IP-адреса — это самая частая задача для регулярных выражений? Вот и я думаю, что нет. Про разницу (?:...)/ (...) фиг найдёшь, а без этого знания в некоторых случаях можно только страдать.

Плюс в питоне есть немало регулярных плюшек. Например, `re.split` может добавлять тот кусок текста, по которому был разрез, в список частей. А в `re.sub` можно вместо шаблона для замены передать функцию. Это — реальные вещи, которые прямо очень нужны, но никто про это не пишет.

Так и родился этот достаточно многобуквенный материал с подробностями, тонкостями, картинками и задачами.



Надеюсь, вам удастся из него извлечь что-нибудь новое и полезное, даже если вы уже в ладах с регулярками.

PS. Решения задач школьники сдают в тестирующую систему, поэтому задачи оформлены в несколько формальном виде.

Содержание

[Регулярные выражения в Python от простого к сложному;](#)

[Содержание;](#)

[Примеры регулярных выражений;](#)

[Сила и ответственность;](#)

[Документация и ссылки;](#)

[Основы синтаксиса;](#)

[Шаблоны, соответствующие одному символу;](#)

[Квантификаторы \(указание количества повторений\);](#)

[Жадность в регулярках и границы найденного шаблона;](#)

[Пересечение подстрок;](#)

[Эксперименты в песочнице;](#)

[Регулярки в питоне;](#)

[Пример использования всех основных функций;](#)

[Тонкости экранирования в питоне \('\\\\\\\\\\\\\\\\foo'\);](#)

[Использование дополнительных флагов в питоне;](#)

[Написание и тестирование регулярных выражений;](#)

[Задачи — 1;](#)

[Скобочные группы \(?...\) и перечисления |;](#)

[Перечисления \(операция «ИЛИ»\);](#)

[Скобочные группы \(группировка плюс квантификаторы\);](#)

[Скобки плюс перечисления;](#)

[Ещё примеры;](#)

[Задачи — 2;](#)

[Группирующие скобки \(...\) и match-объекты в питоне;](#)

[Match-объекты;](#)

[Группирующие скобки \(...\);](#)

[Тонкости со скобками и нумерацией групп;](#)

[Группы и re.findall;](#)

[Группы и re.split;](#)

[Использование групп при заменах;](#)

[Замена с обработкой шаблона функцией в питоне;](#)

[Ссылки на группы при поиске;](#)

[Задачи — 3;](#)

[Шаблоны, соответствующие не конкретному тексту, а позиции;](#)

[Простые шаблоны, соответствующие позиции;](#)

[Сложные шаблоны, соответствующие позиции \(lookaround и Co\);](#)

[lookaround на примере королей и императоров Франции;](#)

[Задачи — 4;](#)

[Post scriptum;](#)

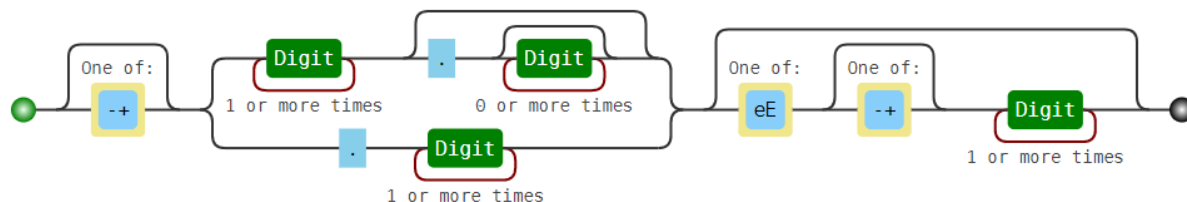
Регулярное выражение — это строка, задающая шаблон поиска подстрок в тексте. Одному шаблону может соответствовать много разных строчек. Термин «Регулярные выражения» является переводом английского словосочетания «Regular expressions». Перевод не очень точно отражает смысл, правильнее было бы «шаблонные выражения». Регулярное выражение, или коротко «регулярка», состоит из обычных символов и специальных командных последовательностей. Например, `\d` задаёт любую цифру, а `\d+` — задаёт любую последовательность из

одной или более цифр. Работа с регулярками реализована во всех современных языках программирования. Однако существует несколько «диалектов», поэтому функционал регулярных выражений может различаться от языка к языку. В некоторых языках программирования регулярками пользоваться очень удобно (например, в питоне), в некоторых — не слишком (например, в C++).

Примеры регулярных выражений

Регулярка	Её смысл
<code>simple text</code>	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 куски, на них похожие, например, 98/7 6/5432)
<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 <code>\d+</code> — 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.

RegExp: `/[-+]?(?:\d+(?:\.\d*)?|\.\d+)(?:[eE][-+]? \d+)?/`



Сила и ответственность

Регулярные выражения, или коротко, *регулярки* — это очень мощный инструмент. Но использовать их следует с умом и осторожностью, и только там, где они действительно приносят пользу, а не вред. Во-первых, плохо написанные

регулярные выражения работают медленно. Во-вторых, их зачастую очень сложно читать, особенно если регулярка написана не лично тобой пять минут назад. В-третьих, очень часто даже небольшое изменение задачи (того, что требуется найти) приводит к значительному изменению выражения. Поэтому про регулярки часто говорят, что это *write only code* (код, который только пишут с нуля, но не читают и не правят). А также шутят: *Некоторые люди, когда сталкиваются с проблемой, думают «Я знаю, я решу её с помощью регулярных выражений.» Теперь у них две проблемы.* Вот пример write-only регулярки (для проверки валидности e-mail адреса (не надо так делать!!!)):

```
(?:[a-z0-9!#$%&'*/=?^_`{|}~-]+(?:\.[a-z0-9!#$%&'*/=?^_`{|}~-]+)*|"(?:[\x01-\x08\x0b\x0c\x0e-\x1f\x21-\x5a\x5d-\x7f]|\\[\x01-\x08\x0b\x0c\x0e-\x1f\x21-\x7f])")@(?:(?:[a-z0-9](?:[a-z0-9-]{0,61}[a-z0-9])?\.)+(?:[a-z0-9-]{2,61}[a-z0-9])?|(?-?:25[0-5]|2[0-4][0-9]|0[0-9][0-9])?\.)+(?:[a-z0-9-]{2,61}[a-z0-9])?
```

А вот [здесь](#) более точная регулярка для проверки корректности email адреса стандарту RFC822. Если вдруг будете проверять email, то не делайте так! Если адрес вводит пользователь, то пусть вводит почти что угодно, лишь бы там была собака. Надёжнее всего отправить туда письмо и убедиться, что пользователь может его получить.

Документация и ссылки

- Оригинальная документация: docs.python.org/3/library/re.html;
- Очень подробный и обстоятельный материал: www.regular-expressions.info/;
- Разные сложные трюки и тонкости с примерами: <http://www.rexegg.com/>;
- Он-лайн отладка регулярок regex101.com (не забудьте поставить галочку Python в разделе FLAVOR слева);
- Он-лайн визуализация регулярок www.debuggex.com/ (не забудьте выбрать Python);
- Мощнейший текстовый редактор [Sublime text 3](https://www.sublimetext.com/3), в котором очень удобный поиск по регуляркам;

Основы синтаксиса

Любая строка (в которой нет символов `.^$*+?{}[]\|()`) сама по себе является регулярным выражением. Так, выражению **Хаха** будет соответствовать строка "Хаха" и только она. Регулярные выражения являются регистрозависимыми, поэтому строка "хаха" (с маленькой буквы) уже не будет соответствовать

выражению выше. Подобно строкам в языке Python, регулярные выражения имеют спецсимволы `.^$*+?{}[]\|()`, которые в регулярках являются управляющими конструкциями. Для написания их просто как символов требуется их *экранировать*, для чего нужно поставить перед ними знак `\`. Так же, как и в питоне, в регулярных выражениях выражение `\n` соответствует концу строки, а `\t` — табуляции.

Шаблоны, соответствующие одному символу

Во всех примерах ниже соответствия регулярному выражению выделяются бирюзовым цветом с подчёркиванием.

Шаблон	Описание	Пример	Применяем к тексту
.	Один любой символ, кроме новой строки <code>\n</code> .	<code>м.л.ко</code>	<u>молоко</u> , <u>малак</u> <u>о</u> , <u>Им0л0ко</u> Ихле б
<code>\d</code>	Любая цифра	<code>су\d\d</code>	<u>СУ35</u> , <u>СУ11</u> 1, А <u>ЛСУ14</u>
<code>\D</code>	Любой символ, кроме цифры	<code>926\D123</code>	<u>926</u>) <u>123</u> , <u>1926</u> - <u>1234</u>
<code>\s</code>	Любой пробельный символ (пробел, табуляция, конец строки и т.п.)	<code>бор\sод а</code>	<u>бор_ода</u> , <u>бор_ода</u> , борода
<code>\S</code>	Любой непробельный символ	<code>\S123</code>	<u>X123</u> , <u>я123</u> , <u>!12</u> <u>3456</u> , 1 + 1234 56
<code>\w</code>	Любая буква (то, что может быть частью слова), а также цифры и <code>_</code>	<code>\w\w\w</code>	<u>Год</u> , <u>f_3</u> , <u>qwert</u>
<code>\W</code>	Любая не-буква, не-цифра и не подчёркивание	<code>com\W</code>	<u>com!</u> , <u>com?</u>
<code>[...]</code>	Один из символов в скобках, а также любой символ из диапазона <code>a-b</code>	<code>[0-9][0-9A-Za-f]</code>	<u>12</u> , <u>1E</u> , <u>4B</u>
<code>[^...]</code>	Любой символ, кроме перечисленных	<code><[^>]></code>	<u><1></u> , <u><a></u> , <u><>></u>
<code>\d≈[0-9], \D≈[^0-9], \w≈[0-9a-zA-Za-яA-ЯёЁ], \s≈[\f\n\r\t\v]</code>	Буква "ё" не включается в общий диапазон букв! Вообще говоря, в <code>\d</code> включается всё, что в юникоде помечено как «цифра», а в <code>\w</code> — как буква. Ещё много всего!		
<code>[abc-], [-1]</code>	если нужен минус, его нужно указать последним или первым		
<code>[*[(+ \\)]\t]</code>	внутри скобок нужно экранировать только <code>]</code> и <code>\</code>		
<code>\b</code>	Начало или конец слова (слева пусто или не-буква, справа буква и наоборот). В отличие от предыдущих соответствует позиции, а не символу	<code>\bвал</code>	<u>вал</u> , перевал, Перевалка

Шаблон	Описание	Пример	Применяем к тексту
<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

Квантификаторы (указание количества повторений)

Шаблон	Описание	Пример	Применяем к тексту
<code>{n}</code>	Ровно n повторений	<code>\d{4}</code>	1, 12, 123, <u>1234</u> , 12345
<code>{m,n}</code>	От m до n повторений включительно	<code>\d{2,4}</code>	1, <u>12</u> , <u>123</u> , <u>1234</u> , 12345
<code>{m,}</code>	Не менее m повторений	<code>\d{3,}</code>	1, 12, <u>123</u> , <u>1234</u> , <u>12345</u>
<code>{,n}</code>	Не более n повторений	<code>\d{,2}</code>	<u>1</u> , <u>12</u> , <u>123</u>
<code>?</code>	Ноль или одно вхождение, синоним <code>{0,1}</code>	<code>валы?</code>	<u>вал</u> , <u>валы</u> , <u>валов</u>
<code>*</code>	Ноль или более, синоним <code>{0,}</code>	<code>СУ\d*</code>	<u>СУ</u> , <u>СУ1</u> , <u>СУ12</u> , ...
<code>+</code>	Одно или более, синоним <code>{1,}</code>	<code>а\)+</code>	<u>а</u>), <u>а</u>)), <u>а</u>))), <u>ба</u>))
<code>*?</code> <code>+?</code> <code>??</code> <code>{m,n}?</code> <code>{,n}?</code> <code>{m,}?</code>	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>\(.*?)</code> <code>\)</code>	<u>(a+b)*(c+d)*(e+f)</u> <u>(a+b)*(c+d)*(e+f)</u>

Жадность в регулярках и границы найденного шаблона

Как указано выше, по умолчанию квантификаторы *жадные*. Этот подход решает очень важную проблему — проблему границы шаблона. Скажем, шаблон `\d+` захватывает максимально возможное количество цифр. Поэтому можно быть уверенным, что перед найденным шаблоном идёт не цифра, и после идёт не цифра. Однако если в шаблоне есть не жадные части (например, явный текст), то подстрока может быть найдена неудачно. Например, если мы хотим найти «слова», начинающиеся на `СУ`, после которой идут цифры, при помощи регулярки `СУ\d*`, то мы найдём и неправильные шаблоны:

`ПАСУ13 СУ12`, ЧТОБЫ СУ6ЕНИЕ УДАЛОСЬ.

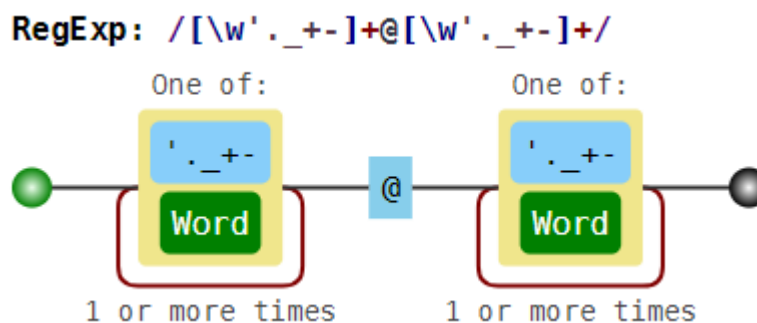
В тех случаях, когда это важно, условие на границу шаблона нужно обязательно добавлять в регулярку. О том, как это можно делать, будет дальше.

Пересечение подстрок

В обычной ситуации регулярки позволяют найти только непересекающиеся шаблоны. Вместе с проблемой границы слова это делает их использование в некоторых случаях более сложным. Например, если мы решим искать e-mail адреса при помощи неправильной регулярки `\w+@\w+` (или даже лучше, `[\w'._+-]+\@([\w'._+-]+)`), то в неудачном случае найдём вот что:

| foo@boo@goo@moo@roo@zoo

То есть это с одной стороны и не e-mail, а с другой стороны это не все подстроки вида **текст - собака - текст**, так как `boo@goo` и `moo@roo` пропущены.



Эксперименты в песочнице

Если вы впервые сталкиваетесь с регулярными выражениями, то лучше всего сначала попробовать [песочницу](#). Посмотрите, как работают простые шаблоны и квантификаторы. Решите следующие задачи для этого текста (возможно, к части придётся вернуться после следующей теории):

1. Найдите все натуральные числа (возможно, окружённые буквами);
2. Найдите все «слова», написанные капсом (то есть строго заглавными), возможно внутри настоящих слов (aaa**БББ**vvv);
3. Найдите слова, в которых есть русская буква, а когда-нибудь за ней цифра;
4. Найдите все слова, начинающиеся с русской или латинской большой буквы (`\b` — граница слова);
5. Найдите слова, которые начинаются на гласную (`\b` — граница слова);
6. Найдите все натуральные числа, не находящиеся внутри или на границе слова;

7. Найдите строчки, в которых есть символ `*` (`.` — это точно не конец строки!);
8. Найдите строчки, в которых есть открывающая и когда-нибудь потом закрывающая скобки;
9. Выделите одним махом весь кусок оглавления (в конце примера, вместе с тегами);
10. Выделите одним махом только текстовую часть оглавления, без тегов;
11. Найдите пустые строчки;

Регулярки в питоне

Функции для работы с регулярками живут в модуле `re`. Основные функции:

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> ;

Пример использования всех основных функций

```

1. import re
2.
3. match = re.search(r'\d\d\D\d\d', r'Телефон 123-12-12')
4. print(match[0] if match else 'Not found')
5. # -> 23-12
6. match = re.search(r'\d\d\D\d\d', r'Телефон 1231212')
7. print(match[0] if match else 'Not found')
8. # -> Not found
9.
10. match = re.fullmatch(r'\d\d\D\d\d', r'12-12')
11. print('YES' if match else 'NO')
```

```

12. # -> YES
13. match = re.fullmatch(r'\d\d\d\d\d', r'T. 12-12')
14. print('YES' if match else 'NO')
15. # -> NO
16.
17. print(re.split(r'\W+', 'Где, скажите мне, мои очки?!'))
18. # -> ['Где', 'скажите', 'мне', 'мои', 'очки', '']
19.
20. print(re.findall(r'\d\d\.\d\d\.\d{4}',
21.                 r'Эта строка написана 19.01.2018, а могла бы и
01.09.2017'))
22. # -> ['19.01.2018', '01.09.2017']
23.
24. 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'):
25.     print('Дата', m[0], 'начинается с позиции', m.start())
26. # -> Дата 19.01.2018 начинается с позиции 20
27. # -> Дата 01.09.2017 начинается с позиции 45
28.
29. print(re.sub(r'\d\d\.\d\d\.\d{4}',
30.              r'DD.MM.YYYY',
31.              r'Эта строка написана 19.01.2018, а могла бы и 01.09.2017'))
32. # -> Эта строка написана DD.MM.YYYY, а могла бы и DD.MM.YYYY

```

Тонкости экранирования в питоне ('\\\\\\\\foo')

Так как символ `\` в питоновских строках также необходимо экранировать, то в результате в шаблонах могут возникать конструкции вида `'\\\\par'`. Первый слеш означает, что следующий за ним символ нужно оставить «как есть». Третий также. В результате с точки зрения питона `'\\\\'` означает просто два слеша `\\`. Теперь с точки зрения движка регулярных выражений, первый слеш экранирует второй. Тем самым как шаблон для регулярки `'\\\\par'` означает просто текст `\\par`. Для того, чтобы не было таких нагромождений слешей, перед открывающей кавычкой нужно поставить символ `r`, что скажет питону «не рассматривай `\` как экранирующий символ (кроме случаев экранирования открывающей кавычки)». Соответственно можно будет писать `r'\\\\par'`.

Использование дополнительных флагов в питоне

Каждой из функций, перечисленных выше, можно дать дополнительный параметр `flags`, что несколько изменит режим работы регулярки. В качестве значения нужно передать сумму выбранных констант, вот они:

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: •١٢٣٤٥٦٧٨٩. <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

```

1. import re
2. print(re.findall(r'\d+', '12 + \v'))
3. # -> ['12', '\v']
4. print(re.findall(r'\w+', 'Hello, мир!'))
5. # -> ['Hello', 'мир']
6. print(re.findall(r'\d+', '12 + \v', flags=re.ASCII))
7. # -> ['12']
8. print(re.findall(r'\w+', 'Hello, мир!', flags=re.ASCII))
9. # -> ['Hello']
10. print(re.findall(r'[уеыаозяию]+', '0000 ааааа ррррр ыыыы яяяя'))
11. # -> ['ааааа', 'яяяя']
12. print(re.findall(r'[уеыаозяию]+', '0000 ааааа ррррр ыыыы яяяя',
13. flags=re.IGNORECASE))
14. # -> ['0000', 'ааааа', 'ыыыы', 'яяяя']
15.
16. text = r"""
17. Торт
18. с вишней1
19. вишней2
20. """
21. print(re.findall(r'Торт.с', text))
22. # -> []
23. print(re.findall(r'Торт.с', text, flags=re.DOTALL))
24. # -> ['Торт\nс']
25. print(re.findall(r'виш\w+', text, flags=re.MULTILINE))
26. # -> ['вишней1', 'вишней2']
27. print(re.findall(r'^виш\w+', text, flags=re.MULTILINE))
28. # -> ['вишней2']

```

Написание и тестирование регулярных выражений

Для написания и тестирования регулярных выражений удобно использовать сервис regex101.com (не забудьте поставить галочку Python в разделе FLAVOR слева) или текстовый редактор [Sublime text 3](#).

Задачи — 1

Задача 01. Регистрационные знаки транспортных средств

В России применяются регистрационные знаки нескольких видов.

Общего в них то, что они состоят из цифр и букв. Причём используются только 12 букв кириллицы, имеющие графические аналоги в латинском алфавите — А, В, Е, К, М, Н, О, Р, С, Т, У и Х.

У частных легковых автомобилях номера — это буква, три цифры, две буквы, затем две или три цифры с кодом региона. У такси — две буквы, три цифры, затем две или три цифры с кодом региона. Есть также и другие виды, но в этой задаче они не понадобятся.

Вам потребуется определить, является ли последовательность букв корректным номером указанных двух типов, и если является, то каким.

На вход даются строки, которые претендуют на то, чтобы быть номером. Определите тип номера. Буквы в номерах — заглавные русские. Маленькие и английские для простоты можно игнорировать.

Ввод		Вывод	
1.	C227HA777	1.	Private
2.	KY22777	2.	Taxi
3.	T22B7477	3.	Fail
4.	M227K19Y9	4.	Fail
5.	C227HA777	5.	Fail

Задача 02. Количество слов

Слово — это последовательность из букв (русских или английских), внутри которой могут быть дефисы.

На вход даётся текст, посчитайте, сколько в нём слов.

PS. Задача решается в одну строчку. Никакие хитрые техники, не упомянутые выше, не требуются.

Ввод		Вывод	
1.	Он --- серо-буро-малиновая редиска!!	1.	9
2.	>>>:->		
3.	А не кот.		
4.	www.kot.ru		

Задача 03. Поиск e-mailов

Допустимый формат e-mail адреса регулируется стандартом RFC 5322.

Если говорить вкратце, то e-mail состоит из одного символа @ (*at-символ* или *собака*), текста до собаки (*Local-part*) и текста после собаки (*Domain part*). Вообще в адресе может быть всякий беспредел (вкратце можно прочитать о нём в [википедии](#)). Довольно странные штуки могут быть валидным адресом, например:

```
"very.(),:;<>[]\".VERY.\"very@\\  
\"very\".unusual\"@[IPv6:2001:db8::1]  
\"()<>[:.,;@\\\"!#$%&'-/=?^_`{}| ~.a\"@(comment)exa-mple
```

Но большинство почтовых сервисов такой ад и вакханалию не допускают. И мы тоже не будем 😊

Будем рассматривать только адреса, имя которых состоит из не более, чем 64 латинских букв, цифр и символов `'._+,-`, а домен — из не более, чем 255 латинских букв, цифр и символов `.-`. Ни Local-part, ни Domain part не может начинаться или заканчиваться на `.-+,-`, а ещё в адресе не может быть более одной точки подряд.

Кстати, полезно знать, что часть имени после символа `+` игнорируется, поэтому можно использовать синонимы своего адреса (например, `shashkov+spam@179.ru` и `shashkov+vk@179.ru`), для того, чтобы упростить себе сортировку почты. (Правда не все сайты позволяют использовать «+», увы)

На вход даётся текст. Необходимо вывести все e-mail адреса, которые в нём встречаются. В общем виде задача достаточно сложная, поэтому у нас будет 3 ограничения:

две точки внутри адреса не встречаются;

две собаки внутри адреса не встречаются;

считаем, что e-mail может быть частью «слова», то есть в `boo@ya_ru` мы видим адрес `boo@ya`, а в `foo!boo@ya.ru` видим `boo@ya.ru`.

PS. Совсем не обязательно делать все проверки только регулярками. Регулярные выражения — это просто инструмент, который делает часть задач простыми. Не нужно делать их назад сложными 😊

Ввод		Вывод	
1.	Иван Иванович!	1.	ivanoff@ivan-chai.ru
2.	Нужен ответ на письмо от ivanoff@ivan-chai.ru.	2.	serge'o-lupin@mail.ru
3.	Не забудьте поставить в копию		
4.	serge'o-lupin@mail.ru- это важно.		
1.	N0: foo@ya.ru, foo@ya.ru	1.	boo@ya
2.	PARTLY: boo@ya_ru, -boo@ya.ru-, foo№boo@ya.ru	2.	boo@ya.ru
		3.	boo@ya.ru

Скобочные группы (?:...) и перечисления |

Перечисления (операция «ИЛИ»)

Чтобы проверить, удовлетворяет ли строка хотя бы одному из шаблонов, можно воспользоваться аналогом оператора **or**, который записывается с помощью символа **|**. Так, некоторая строка подходит к регулярному выражению **A|B** тогда и только тогда, когда она подходит хотя бы к одному из регулярных выражений **A** или **B**. Например, отдельные овощи в тексте можно искать при помощи шаблона **морковк|св[её]кл|картошк|редиск**.

Скобочные группы (группировка плюс квантификаторы)

Зачастую шаблон состоит из нескольких повторяющихся групп. Так, MAC-адрес сетевого устройства обычно записывается как шесть групп из двух шестнадцатиричных цифр, разделённых символами **-** или **:**. Например, **01:23:45:67:89:ab**. Каждый отдельный символ можно задать как **[0-9a-fA-F]**, и можно весь шаблон записать так:

```
[0-9a-fA-F]{2}[:-][0-9a-fA-F]{2}[:-][0-9a-fA-F]{2}[:-]
[0-9a-fA-F]{2}[:-][0-9a-fA-F]{2}[:-][0-9a-fA-F]{2}
```

Ситуация становится гораздо сложнее, когда количество групп заранее не зафиксировано.

Чтобы разрешить эту проблему в синтаксисе регулярных выражений есть группировка **(?:...)**. Можно писать круглые скобки и без значков **?:**, однако от этого у группировки значительно меняется смысл, регулярка начинает работать гораздо медленнее. Об этом будет написано ниже. Итак, если **REGEXP** — шаблон, то **(?:REGEXP)** — эквивалентный ему шаблон. Разница только в том, что теперь к **(?:REGEXP)** можно применять квантификаторы, указывая, сколько именно

строка вида **HH:MM:SS** или **HH:MM**, в которой **HH** — число от 00 до 23, а **MM** и **SS** — число от 00 до 59.

Ввод		Вывод	
1.	Уважаемые! Если вы к 09:00 не вернёт	1.	Уважаемые! Если вы к (TBD) не вернёт
2.	чемодан, то уже в 09:00:01 я за себя не отвечаю.	2.	чемодан, то уже в (TBD) я за себя не отвечаю.
3.	PS. С отношением 25:50 всё нормально!	3.	PS. С отношением 25:50 всё нормально!

Задача 05. Действительные числа в паскале

Паскаль требует, чтобы реальные константы имели либо десятичную точку, либо экспоненту (начиная с буквы e или E и официально называемую масштабным коэффициентом), либо обе, в дополнение к обычному набору десятичных цифр. Если десятичная точка включена, у нее должна быть хотя бы одна десятичная цифра с каждой стороны от нее. Как и ожидалось, знак (+ или -) может предшествовать целому числу или показателю степени, или обоим. Экспоненты могут не включать дробные цифры. Пробелы могут предшествовать или следовать за реальной константой, но они не могут быть встроены в нее. Обратите внимание, что синтаксические правила Паскаля для реальных констант не делают предположений о диапазоне действительных значений, как и эта проблема. Ваша задача в этой задаче состоит в том, чтобы определить действительные константы Паскаля.

Ввод		Вывод	
1.	1.2	1.	1.2 is legal.
2.	1.	2.	1. is illegal.
3.	1.0e-55	3.	1.0e-55 is legal.
4.	e-12	4.	e-12 is illegal.
5.	6.5E	5.	6.5E is illegal.
6.	1e-12	6.	1e-12 is legal.
7.	+4.1234567890E-99999	7.	+4.1234567890E-99999 is legal.
8.	7.6e+12.5	8.	7.6e+12.5 is illegal.
9.	99	9.	99 is illegal.

Задача 06. Аббревиатуры

Владимир устроился на работу в одно очень важное место. И в первом же документе он ничего не понял, там были сплошные *ФГУП НИЦ ГИДГЕО, ФГОУ ЧШУ АПК* и т.п. Тогда он решил

собрать все аббревиатуры, чтобы потом найти их расшифровки на <http://sokr.ru/>. Помогите ему.

Будем считать аббревиатурой слова только лишь из заглавных букв (как минимум из двух). Если несколько таких слов разделены пробелами, то они считаются одной аббревиатурой.

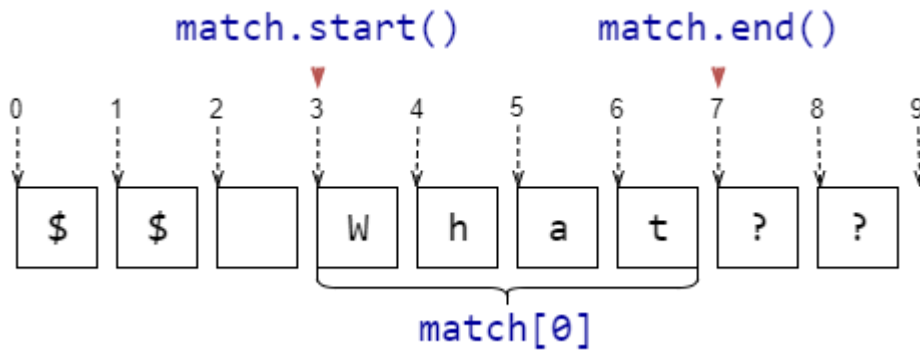
Ввод		Вывод	
1.	Это курс информатики соответствует ФГОС и ПОО	1.	ФГОС
	П,	2.	ПООП
2.	это подтверждено ФГУ ФНЦ НИИСИ РАН	3.	ФГУ ФНЦ НИИСИ Р АН

Группирующие скобки (. . .) и `match`-объекты в питоне

Match-объекты

Если функции `re.search`, `re.fullmatch` не находят соответствие шаблону в строке, то они возвращают `None`, функция `re.finditer` не выдаёт ничего. Однако если соответствие найдено, то возвращается `match`-объект. Эта штука содержит в себе кучу полезной информации о соответствии шаблону. Полный набор атрибутов можно посмотреть в [документации](#), а здесь приведём самое полезное.

Метод	Описание	Пример
<code>match[0]</code> , <code>match.group()</code>	Подстрока, соответствующая шаблону	<pre>match = re.search(r'\w+', r'\$\$ What??') match[0] # -> 'What'</pre>
<code>match.start()</code>	Индекс в исходной строке, начиная с которого идёт найденная подстрока	<pre>match = re.search(r'\w+', r'\$\$ What??') match.start() # -> 3</pre>
<code>match.end()</code>	Индекс в исходной строке, который следует сразу за найденной подстрокой	<pre>match = re.search(r'\w+', r'\$\$ What??') match.end() # -> 7</pre>



Группирующие скобки (...)

Если в шаблоне регулярного выражения встречаются скобки (...) без `?:`, то они становятся *группирующими*. В `match`-объекте, который возвращают `re.search`, `re.fullmatch` и `re.finditer`, по каждой такой группе можно получить ту же информацию, что и по всему шаблону. А именно часть подстроки, которая соответствует (...), а также индексы начала и окончания в исходной строке. Достаточно часто это бывает полезно.

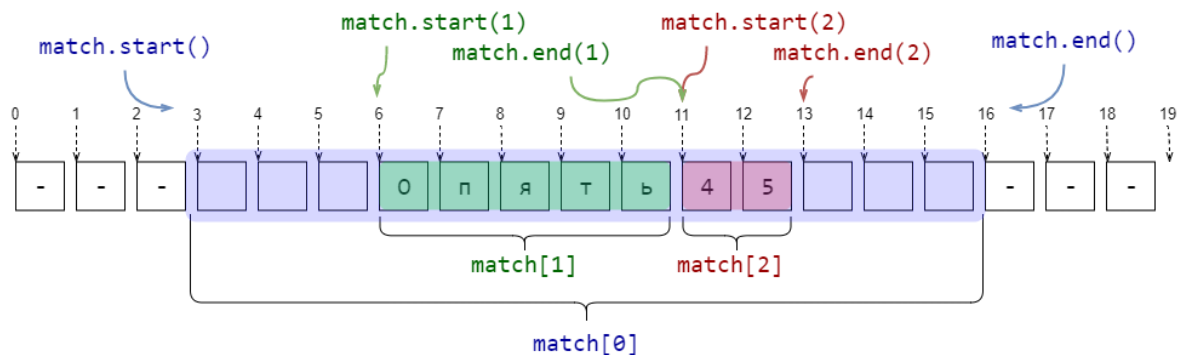
```

1. import re
2. pattern = r'\s*([А-Яа-яЁё]+)(\d+)\s*'
3. string = r'--- Опять45 ---'
4. match = re.search(pattern, string)
5. print(f'Найдена подстрока >{match[0]}< с позиции {match.start(0)} до {match.end(0)}')
6. print(f'Группа букв >{match[1]}< с позиции {match.start(1)} до {match.end(1)}')
7. print(f'Группа цифр >{match[2]}< с позиции {match.start(2)} до {match.end(2)}')
8. ###
9. # -> Найдена подстрока >   Опять45   < с позиции 3 до 16
10. # -> Группа букв >Опять< с позиции 6 до 11
11. # -> Группа цифр >45< с позиции 11 до 13

```



«>



Тонкости со скобками и нумерацией групп.

Если к группирующим скобкам применён квантификатор (то есть указано число повторений), то подгруппа в match-объекте будет создана только для последнего соответствия. Например, если бы в примере выше квантификаторы были снаружи от скобок `'\s*([А-Яа-яЁё])+(\d)+\s*'`, то вывод был бы таким:

1. # -> Найдена подстрока > Оять45 < с позиции 3 до 16
2. # -> Группа букв >ь< с позиции 10 до 11
3. # -> Группа цифр >5< с позиции 12 до 13

Внутри группирующих скобок могут быть и другие группирующие скобки. В этом случае их нумерация производится в соответствии с номером появления открывающей скобки с шаблоне.

```

1. import re
2. pattern = r'((\d)(\d))((\d)(\d))'
3. string = r'123456789'
4. match = re.search(pattern, string)
5. print(f'Найдена подстрока >{match[0]}< с позиции {match.start(0)} до {match.end(0)}')
6. for i in range(1, 7):
7.     print(f'Группа №{i} >{match[i]}< с позиции {match.start(i)} до {match.end(i)}')
8. ###
9. # -> Найдена подстрока >1234< с позиции 0 до 4
10. # -> Группа №1 >12< с позиции 0 до 2
11. # -> Группа №2 >1< с позиции 0 до 1
12. # -> Группа №3 >2< с позиции 1 до 2
13. # -> Группа №4 >34< с позиции 2 до 4
14. # -> Группа №5 >3< с позиции 2 до 3
15. # -> Группа №6 >4< с позиции 3 до 4

```

Группы и `re.findall`

Если в шаблоне есть группирующие скобки, то вместо списка найденных подстрок будет возвращён список кортежей, в каждом из которых только соответствие каждой группе. Это не всегда происходит по плану, поэтому обычно нужно использовать негруппирующие скобки `(?:...)`.

```

1. import re
2. print(re.findall(r'([a-z]+)(\d*)', r'foo3, im12, go, 24buz42'))
3. # -> [('foo', '3'), ('im', '12'), ('go', ''), ('buz', '42')]

```

Группы и `re.split`

Если в шаблоне нет группирующих скобок, то `re.split` работает очень похожим образом на `str.split`. А вот если группирующие скобки в шаблоне есть, то между каждыми разрезанными строками будут все соответствия каждой из подгрупп.

```

1. import re
2. print(re.split(r'(\s*)([+*/-])(\s*)', r'12 + 13*15 - 6'))
3. # -> ['12', ' ', '+', ' ', '13', '*', ' ', '15', ' ', '-', ' ', '6']

```

В некоторых ситуациях эта возможность бывает чрезвычайно удобна! Например, достаточно из предыдущего примера убрать лишние группы, и польза сразу станет очевидна!

```

1. import re
2. print(re.split(r'\s*([+*/-])\s*', r'12 + 13*15 - 6'))
3. # -> ['12', '+', '13', '*', '15', '-', '6']

```

Использование групп при заменах

Использование групп добавляет замене (`re.sub`, работает не только в питоне, а почти везде) очень удобную возможность: в шаблоне для замены можно сослаться на соответствующую группу при помощи `\1`, `\2`, `\3`, Например, если нужно даты из неудобного формата ММ/ДД/ГГГГ перевести в удобный ДД.ММ.ГГГГ, то можно использовать такую регулярку:

```

1. import re
2. text = "We arrive on 03/25/2018. So you are welcome after 04/01/2018."
3. print(re.sub(r'(\d\d)/(\d\d)/(\d{4})', r'\2.\1.\3', text))
4. # -> We arrive on 25.03.2018. So you are welcome after 01.04.2018.

```

Если групп больше 9, то можно сослаться на них при помощи конструкции вида `\g<12>`.

Замена с обработкой шаблона функцией в питоне

Ещё одна питоновская фишка для регулярных выражений: в функции `re.sub` вместо текста для замены можно передать функцию, которая будет получать на вход `match`-объект и должна возвращать строку, на которую и будет произведена

замена. Это позволяет не писать ад в шаблоне для замены, а использовать удобную функцию. Например, «зацензурируем» все слова, начинающиеся на букву «Х»:

```
1. import re
2. def repl(m):
3.     return '>censored(' + str(len(m[0])) + '>'
4.     text = "Некоторые хорошие слова подозрительны: хор, хоровод,
хороводоводовед."
5.     print(re.sub(r'\b[хХxX]\w*', repl, text))
6.     # -> Некоторые >censored(7)< слова подозрительны: >censored(3)<,
>censored(7)<, >censored(15)<.
```

Ссылки на группы при поиске

При помощи `\1`, `\2`, `\3`, ... и `\g<12>` можно сослаться на найденную группу и при поиске. Необходимость в этом встречается довольно редко, но это бывает полезно при обработке простых xml и html.

Только пообещайте, что не будете парсить сложный xml и тем более html при помощи регулярок! Регулярные выражения для этого не подходят. Используйте другие инструменты. Каждый раз, когда неопытный программист парсит html регулярками, в мире умирает котёнок. Если кажется «Да здесь очень простой html, напишу регулярку», то сразу вспоминайте шутку про две проблемы. Не нужно пытаться парсить html регулярками, даже Пётр Митричев не сможет это сделать в общем случае 😊 Использование регулярных выражений при парсинге html подобно залатыванию резиновой лодки шилом. Закон Мёрфи для парсинга html и xml при помощи регулярок гласит: парсинг html и xml регулярками иногда работает, но в точности до того момента, когда правильность результата будет *очень* важна.

Используйте [lxml](#) и [beautiful soup](#).

```
1. import re
2. text = "SPAM <foo>Here we can <boo>find</boo> something interesting</foo>
SPAM"
3. print(re.search(r'<(\w+?)>.*?</\1>', text)[0])
4. # -> <foo>Here we can <boo>find</boo> something interesting</foo>
5. text = "SPAM <foo>Here we can <foo>find</foo> OH, NO MATCH HERE!</foo>
SPAM"
6. print(re.search(r'<(\w+?)>.*?</\1>', text)[0])
7. # -> <foo>Here we can <foo>find</foo>
```

Задачи — 3

Задача 07. Шифровка

Владимиру потребовалось срочно запутать финансовую документацию. Но так, чтобы это было обратимо.

Он не придумал ничего лучше, чем заменить каждое целое число (последовательность цифр) на его куб. Помогите ему.

Ввод		Вывод	
1.	Было закуплено 12 единиц техники	1.	Было закуплено 1728 единиц техники
2.	по 410.37 рублей.	2.	по 68921000.50653 рублей.

Задача 08. То ли акrostих, то ли акроним, то ли апроним

Акrostих — осмысленный текст, сложенный из начальных букв каждой строки стихотворения.

Акроним — вид аббревиатуры, образованной начальными звуками (напр. НАТО, вуз, НАСА, ТАСС), которое можно произнести слитно (в отличие от аббревиатуры, которую произносят «по буквам», например: КГБ — «ка-гэ-бэ»).

На вход даётся текст. Выведите слитно первые буквы каждого слова. Буквы необходимо выводить заглавными.

Эту задачу можно решить в одну строчку.

Ввод		Вывод	
1.	Московский государственный институт международных отношений	1.	МГИМО
1.	микоян авиацию снабдил алкоголем,	1.	МАСАНД
2.	народ доволен работой авиаконструктора	РА	

Задача 09. Хайку

Хайку — жанр традиционной японской лирической поэзии века, известный с XIV века.

Оригинальное японское хайку состоит из 17 слогов, составляющих один столбец иероглифов. Особыми разделительными словами — кирэдзи — текст хайку делится на части из 5, 7 и снова 5 слогов. При переводе хайку на западные языки традиционно вместо разделительного слова использую разрыв строки и, таким образом, хайку записываются как трёхстишия.

Перед вами трёхстишия, которые претендуют на то, чтобы быть хайку. В качестве разделителя строк используются символы `/`. Если разделители делят текст на строки, в которых 5/7/5 слогов, то выведите «Хайку!». Если число строк не равно 3, то выведите строку «Не хайку. Должно быть 3 строки.» Иначе выведите строку вида «Не хайку. В i строке слогов не s, а j.», где строка `i` — самая ранняя, в которой количество слогов неправильное.

Для простоты будем считать, что слогов ровно столько же, сколько гласных, не задумываясь о тонкостях.

Ввод	Вывод
Вечер за окном. / Еще один день прожит. / Жизнь скоротечна...	Хайку!
Просто текст	Не хайку. Должно быть 3 строки.
Как вишня расцвела! / Она с коня согнала / И князя-гордеца.	Не хайку. В 1 строке слогов не 5, а 6.
На голой ветке / Ворон сидит одиноко... / Осенний вечер!	Не хайку. В 2 строке слогов не 7, а 8.
Тихо, тихо ползи, / Улитка, по склону Фудзи, / Вверх, до самых высот!	Не хайку. В 1 строке слогов не 5, а 6.
Жизнь скоротечна... / Думает ли об этом / Маленький мальчик.	Хайку!

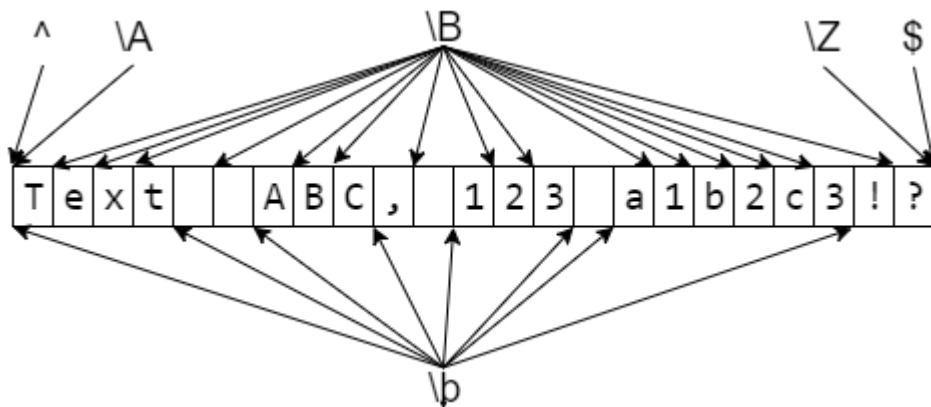
Шаблоны, соответствующие не конкретному тексту, а позиции

Отдельные части регулярного выражения могут соответствовать не части текста, а позиции в этом тексте. То есть такому шаблону соответствует не подстрока, а некоторая позиция в тексте, как бы «между» буквами.

Простые шаблоны, соответствующие позиции

Для определённости строку, в которой мы ищем шаблон будем называть *всем текстом*. Каждую строчку *всего текста* (то есть каждый максимальный кусок без символов конца строки) будем называть *строчкой текста*.

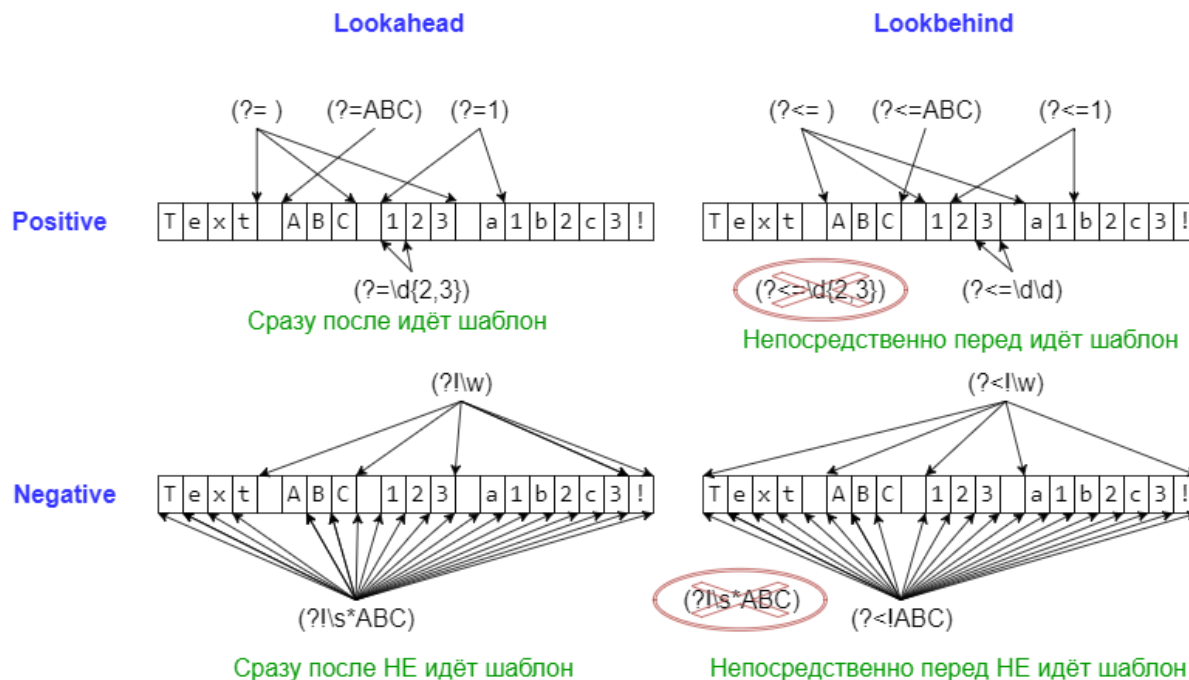
Шаблон	Описание	Пример	Применяем к тексту
<code>^</code>	Начало всего текста или начало строки текста, если <code>flag=re.MULTILINE</code>	<code>^Привет</code>	
<code>\$</code>	Конец всего текста или конец строки текста, если <code>flag=re.MULTILINE</code>	<code>Будь здоров!\$</code>	
<code>\A</code>	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, Перевалка
		<code>\Bвал\b</code>	перевал, вал, Перевалка



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

Следующие шаблоны применяются в основном в тех случаях, когда нужно уточнить, что должно идти непосредственно перед или после шаблона, но при этом не включать найденное в `match`-объект.

Шаблон	Описание	Пример	Применяем к тексту
<code>(?=...)</code>	<i>lookahead assertion</i> , соответствует каждой позиции, сразу после которой начинается соответствие шаблону ...	Isaac <code>(?=Asimov)</code>	Isaac Asimov, Isaac other
<code>(?!...)</code>	<i>negative lookahead assertion</i> , соответствует каждой позиции, сразу после которой НЕ может начинаться шаблон ...	Isaac <code>(?!Asimov)</code>	Isaac Asimov, Isaac other
<code>(?<=...)</code>	<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.	<code>(?<=abc)def</code>	abc def , bcdef
<code>(?<!=...)</code>	<i>negative lookbehind assertion</i> , соответствует каждой позиции, которой НЕ может заканчиваться шаблон ...	<code>(?<!=abc)def</code>	abcdef, bc def



На всякий случай ещё раз. Каждый из этих шаблонов проверяет лишь то, что идёт непосредственно перед позицией или непосредственно после позиции. Если пару таких шаблонов написать рядом, то проверки будут независимы (то есть будут соответствовать AND в каком-то смысле).

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

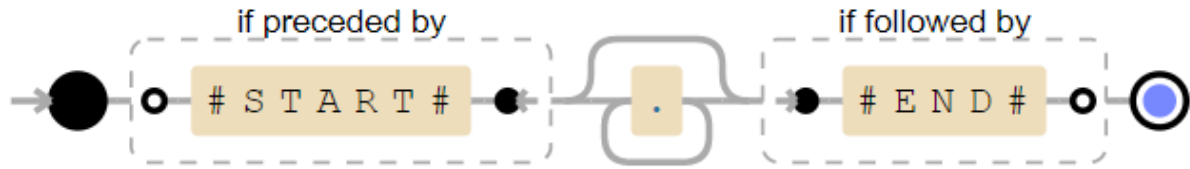
(?<=Людовик)VI — «шестой», но только если Людовик

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

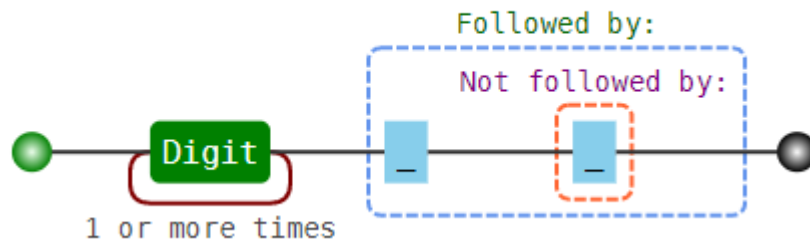
(?<!Людовик)VI — «шестой», но только если не Людовик

КарлIV, КарлIX, КарлV, КарлVI, КарлVII, КарлVIII,
 ЛюдовикIX, ЛюдовикVI, ЛюдовикVII, ЛюдовикVIII, ЛюдовикX, ...,
 ЛюдовикXVII,
 ФилиппI, ФилиппII, ФилиппIII, ФилиппIV, ФилиппV, ФилиппVI

Шаблон	Комментарий	Применяем к тексту
(?<!\d)\d(?!\d)	Цифра, окружённая не-цифрами	Text ABC 123 A <u>1</u> B <u>2</u> C <u>3</u> !
(?<=#START#).*?(?=#END#)	Текст от #START# до #END#	text from #START# <u>till</u> #END#
\d+(?=_(!_))	Цифра, после которой идёт ровно одно подчёркивание	<u>12</u> _34__56
^(?:(!boo).)*?\$	Строка, в которой нет boo (то есть нет такого символа, перед которым есть boo)	<u>a foo and</u> boo and zoo <u>and others</u>
^(?:(!boo)(!foo).)*?\$	Строка, в которой нет ни boo, ни foo	a foo and boo and zoo <u>and others</u>



RegExp: `/\d+(?=_(!_))/?`



Прочие фичи

Конечно, здесь описано не всё, что умеют регулярные выражения, и даже не всё, что умеют регулярные выражения в питоне. За дальнейшим можно обращаться к [этому разделу](#). Из полезного за кадром осталась компиляция регулярных выражений для ускорения многократного использования одного шаблона, использование именованных групп и разные хитрые трюки.

А уж какие извращения можно делать с регулярными выражениями в языке Perl — поручик Ржевский просто отдыхает 😊

Задачи — 4

Задача 10. CamelCase -> under_score

Владимир написал свой открытый проект, именуя переменные в стиле «ВерблюжийРегистр».

И только после того, как написал о нём статью, он узнал, что в питоне для имён переменных принято использовать подчёркивания для разделения слов (under_score). Нужно срочно всё исправить, пока его не «закидали тапками».

Задача могла бы оказаться достаточно сложной, но, к счастью, Владимир совсем не использовал строковых констант и классов.

Поэтому любая последовательность букв и цифр, внутри которой есть заглавные, — это имя переменной, которое нужно поправить.

Ввод		Вывод	
1.	MyVar17 = OtherVar + YetAnother2Var	1.	my_var17 = other_var + yet_another2_var
2.	TheAnswerToLifeTheUniverseAndEverything = 42	2.	the_answer_to_life_the_universe_and_everything = 42

Задача 11. Удаление повторов

Довольно распространённая ошибка — это повтор слова.

Вот в предыдущем предложении такая допущена. Необходимо исправить каждый такой повтор (слово, один или несколько пробельных символов, и снова то же слово).

Ввод	Вывод
Довольно распространённая ошибка ошибка — это лишний повтор слова слова. Смешно, не правда ли? Не нужно портить хор хоровод.	Довольно распространённая ошибка — это лишний повтор слова. Смешно, не правда ли? Не нужно портить хор хоровод.

Задача 12. Близкие слова

Для простоты будем считать словом любую последовательность букв, цифр и знаков `_` (то есть символов `\w`).

Дан текст. Необходимо найти в нём любой фрагмент, где сначала идёт слово «олень», затем не более 5 слов, и после этого идёт слово «заяц».

Ввод	Вывод
1. Да он олень, а не заяц!	1. олень, а не заяц

Задача 13. Форматирование больших чисел

Большие целые числа удобно читать, когда цифры в них разделены на тройки запятыми.

Переформатируйте целые числа в тексте.

Ввод	Вывод
1. 12 мало	1. 12 мало
2. лучше 123	2. лучше 123
3. 1234 почти	3. 1,234 почти
4. 12354 хорошо	4. 12,354 хорошо
5. стало 123456	5. стало 123,456
6. супер 1234567	6. супер 1,234,567

Задача 14. Разделить текст на предложения

Для простоты будем считать, что:

- каждое предложение начинается с заглавной русской или латинской буквы;
- каждое предложение заканчивается одним из знаков препинания **. ; ! ?**;
- между предложениями может быть любой непустой набор пробельных символов;
- внутри предложений нет заглавных и точек (нет пакостей в духе «Мы любим творчество А. С. Пушкина»).

Разделите текст на предложения так, чтобы каждое предложение занимало одну строку.

Пустых строк в выводе быть не должно. Любые наборы из более одного пробельного символа замените на один пробел.

Ввод		Вывод	
1.	В этом	1.	В этом предложении разрывы строки...
2.	предложении разрывы строки... Но это	2.	Но это не так важно!
3.	не так важно! Совсем? Да, совсем!	3.	Совсем?
4.	И это	4.	Да, совсем!
5.	не должно мешать.	5.	И это не должно мешать.

Задача 15. Форматирование номера телефона

Если вы когда-нибудь пытались собирать номера мобильных телефонов, то наверняка знаете, что почти любые 10 человек используют как минимум пять различных способов записать номер телефона. Кто-то начинает с **+7**, кто-то просто с **7** или **8**, а некоторые вообще не пишут префикс. Трёхзначный код кто-то отделяет пробелами, кто-то при помощи дефиса, кто-то скобками (и после скобки ещё пробел некоторые добавляют). После следующих трёх цифр кто-то ставит пробел, кто-то дефис, кто-то ничего не ставит. И после следующих двух цифр — тоже. А некоторые начинают за здоровье, а заканчивают... В общем очень неудобно!

На вход даётся номер телефона, как его мог бы ввести человек. Необходимо его переформатировать в формат **+7 123 456-78-90**. Если с номером что-то не так, то нужно вывести строку **Fail!**.

Ввод	Вывод
1. +7 123 456-78-90	1. +7 123 456-78-90
1. 8(123)456-78-90	1. +7 123 456-78-90
1. 7(123) 456-78-90	1. +7 123 456-78-90
1. 1234567890	1. +7 123 456-78-90
1. 123456789	1. Fail!
1. +9 123 456-78-90	1. Fail!
1. +7 123 456+78=90	1. Fail!
1. +7(123 45678-90	1. +7 123 456-78-90
1. 8(123 456-78-90	1. Fail!

Задача 16. Поиск e-mail'ов — 2

В предыдущей задаче мы немного схалтурили.
Однако к этому моменту задача должна стать посильной!

На вход даётся текст. Необходимо вывести все e-mail адреса, которые в нём встречаются. При этом e-mail не может быть частью слова, то есть слева и справа от e-mail'а должен быть либо конец строки, либо не-буква и при этом не один из символов ' . _ + - ', допустимых в адресе.

Ввод		Вывод	
1.	Иван Иванович!	1.	ivanoff@ivan-chai.ru
2.	Нужен ответ на письмо от ivanoff@ivan-chai.ru.	2.	serge'o-lupin@mail.ru
3.	Не забудьте поставить в копию		
4.	serge'o-lupin@mail.ru- это важно.		
1.	NO: foo@ya.ru, foo@ya.ru, foo@foo@foo	1.	boo1@ya.ru
2.	NO: +foo@ya.ru, foo@ya.ru	2.	boo2@ya.ru
3.	NO: foo@ya_ru, -foo@ya.ru-, foo@ya.ru+	3.	boo3@ya.ru
4.	NO: foo..foo@ya.ru		
5.	YES: (boo1@ya.ru), boo2@ya.ru!, boo3@ya.ru		

Post scriptum

PS. Текст длинный, в нём наверняка есть опечатки и ошибки. Пишите о них скорее в личку, я тут же исправлю.

PSS. Ух и намаялся я нормальный html в хабра-html перегонять. Кажется, парсер хабра писан на регулярках, иначе как объяснить все те странности, которые приходилось вылавливать бинпоиском? 😊

[Регулярные выражения в Python от простого к сложному. Подробности, примеры, картинки, упражнения](#)

[DebuggexBeta](#)

Опубликовано **Вадим В. Костерин** ст. преп. кафедры ЦЭиИТ. Автор более 130 научных и учебно-методических работ. Лауреат ВДНХ (серебряная медаль). [Посмотреть больше записей](#)

Оставьте комментарий

Ваш адрес email не будет опубликован. Обязательные поля помечены *

Комментарий *

Имя *

Email *

Сайт

☐ Сохранить моё имя, email и адрес сайта в этом браузере для последующих моих комментариев.

восемь - один =



Отправить комментарий

PYTHON. Все права защищены. 2019-2022 © В.В. Костерин, Челябинск