КАЛУЖСКИЙ ФИЛИАЛ ФЕДЕРАЛЬНОГО ГОСУДАРСТВЕННОГО БЮДЖЕТНОГО ОБРАЗОВАТЕЛЬНОГО УЧРЕЖДЕНИЯ ВЫСШЕГО ОБРАЗОВАНИЯ «МОСКОВСКИЙ ГОСУДАРСТВЕННЫЙ ТЕХНИЧЕСКИЙ УНИВЕРСИТЕТ ИМЕНИ Н.Э. БАУМАНА (национальный исследовательский университет)»



Факультет «Информатика и управление»

Кафедра «Программное обеспечение ЭВМ, информационные технологии»

Высокоуровневое программирование

Лекция №13. «Обработка текстов. Регулярные выражения»

Регулярные выражения

- Регулярные выражения (regular expressions) описывают множество строк, используя специальный язык. (Строка, в которой задано регулярное выражение, будет называться шаблоном.)
- Для работы с регулярными выражениями в Python используется модуль **re**.

```
import re
pattern = r"[0-9]+"
number_re = re.compile(pattern)
number_re.findall("122 234 65435")

['122', '234', '65435']
```

Регулярные выражения

- В этом примере шаблон *pattern* описывает множество строк, которые состоят из одного или более символов из набора "0", "1", ..., "9".
- Функция re.compile() компилирует шаблон в специальный Regex-объект, который имеет несколько методов, в том числе метод findall() для получения списка всех непересекающихся вхождений строк, удовлетворяющих шаблону, в заданную строку.
- То же самое можно было сделать и так:

```
1 import re
2 re.findall(r"[0-9]+", "122 234 65435")
['122', '234', '65435']
```

• Предварительная компиляция шаблона предпочтительнее при его частом использовании, особенно внутри цикла.

- Синтаксис регулярных выражений в **Python** почти такой же, как в **Perl**, **grep** и некоторых других инструментах.
- Часть символов (в основном буквы и цифры) обозначают сами себя.
- Строка удовлетворяет (соответствует) шаблону, если она входит во множество строк, которые этот шаблон описывает.
- Здесь стоит также отметить, что различные операции используют шаблон по-разному. Так, search() ищет первое вхождение строки, удовлетворяющей шаблону, в заданной строке, а match() требует, чтобы строка удовлетворяла шаблону с самого начала.

Символ	Что обозначает в регулярном выражении
"."	Любой символ
11 ^ 11	Начало строки
"\$"	Конец строки
11 * 11	Повторение фрагмента нуль или более раз (жадное)
"+"	Повторение фрагмента один или более раз (жадное)
"?"	Предыдущий фрагмент либо присутствует, либо отсутствует
"{m,n}"	Повторение предыдущего фрагмента от m до n раз включительно (жадное)
"[]"	Любой символ из набора в скобках. Можно задавать диапазоны символов с идущими подряд кодами, например: $a-z$
"[^]"	Любой символ не из набора в скобках
"\"	Обратная косая черта отменяет специальное значение следующего за ней символа
" "	Фрагмент справа или фрагмент слева
"*?"	Повторение фрагмента нуль или более раз (не жадное)
"+?"	Повторение фрагмента один или более раз (не жадное)
"{m,n}?"	Повторение предыдущего фрагмента от m до n раз включительно (не жадное)

- Если A и В регулярные выражения, то их конкатенация AB является новым регулярным выражением.
- Можно считать, что конкатенация основной способ составления регулярных выражений.
- Скобки, описанные ниже, применяются для задания приоритетов и выделения групп (фрагментов текста, которые потом можно получить по номеру или из словаря, и даже сослаться в том же регулярном выражении).
- Алгоритм, который сопоставляет строки с регулярным выражением, проверяет соответствие того или иного фрагмента строки регулярному выражению. Например, строка "a" соответствует регулярному выражению "[a-z]", строка "fruit" соответствует "fruit| vegetable", а вот строка "apple" не соответствует шаблону "pineapple".

Синта	ксис регулярных выражений
Обозначение	Описание
"(регвыр)"	Обособляет регулярное выражение в скобках и выделяет группу
"(?:регвыр)"	Обособляет регулярное выражение в скобках без выделения группы
"(?=регвыр)"	Взгляд вперед: строка должна соответствовать заданному регулярному выражению, но дальнейшее сопоставление с шаблоном начнется с того же места
"(?!регвыр)"	То же, но с отрицанием соответствия
"(?<=регвыр)"	Взгляд назад: строка должна соответствовать, если до этого момента соответствует регулярному выражению. Не занимает места в строке, к которой применяется шаблон. Параметр регвыр должен быть фиксированной длины (то есть, без "+" и "*")
"(? регвыр)"</th <th>То же, но с отрицанием соответствия</th>	То же, но с отрицанием соответствия
"(?Р<имя>регвыр)"	Выделяет именованную группу с именем имя
" (?Р=имя) "	Точно соответствует выделенной ранее именованной группе с именем имя
"(?#регвыр)"	Комментарий (игнорируется)
"(?(имя)рв1 рв2)"	Если группа с номером или именем имя оказалась определена, результатом будет сопоставление с рв1, иначе - с рв2. Часть рв2 может отсутствовать

может отсутствовать "(?флаг)" Задает флаг для всего данного регулярного выражения. Флаги

необходимо задавать в начале шаблона

Последовательность	Чему соответствует
"\1" - "\9"	Группа с указанным номером. Группы нумеруются, начиная с 1
"\A"	Промежуток перед началом всей строки (почти аналогично "^")
"\Z"	Промежуток перед концом всей строки (почти аналогично "\$")
"\b"	Промежуток между символами перед словом или после него
"\B"	Наоборот, не соответствует промежутку между символами на границе слова
"\d"	Цифра. Аналогично "[0-9]"
"\s"	Любой пробельный символ. Аналогично "[\t\n\r\f\v]"
"\S"	Любой непробельный символ. Аналогично "[^\t\n\r\f\v]"
"\W"	Любая цифра или буква (зависит от флага LOCALE)
"\W"	Любой символ, не являющийся цифрой или буквой (зависит от флага LOCALE)

Флаги, используемые с регулярными выражениями

• "(?i)", re.I, re.IGNORECASE

Сопоставление проводится без учета регистра букв.

• "(?L)", re.L, re.LOCALE

Влияет на определение буквы в "\w", "\W", "\b", "\В" в зависимости от текущей культурной среды (locale).

• "(?m)", re.M, re.MULTILINE

Если этот флаг задан, "^" и "\$" соответствуют началу и концу любой строки.

•"(?s)", re.S, re.DOTALL

Если задан, "." соответствует также и символу конца строки "\n".

- "(?x)", re.X, re.VERBOSE
- Если задан, пробельные символы, не экранированные в шаблоне обратной косой чертой, являются незначащими, а все, что расположено после символа "#", -- комментарии. Позволяет записывать регулярное выражение в несколько строк для улучшения его читаемости и записи комментариев.
- "(?u)", re.U, re.UNICODE

В шаблоне и в строке использован Unicode.

Методы объекта-шаблона

- В результате успешной компиляции шаблона функцией re.compile() получается шаблонобъект (он именуется SRE_Pattern).
- match(s)

Сопоставляет строку в с шаблоном, возвращая в случае удачного сопоставления объект с результатом сравнения (объект SRE_Match). В случае неудачи возвращает None. Сопоставление начинается от начала строки.

• search(s)

Аналогичен match (s), но ищет подходящую подстроку по всей строке s.

Методы объекта-шаблона

• split(s[, maxsplit=0])

Разбивает строку на подстроки, разделенные подстроками, заданными шаблоном. Если в шаблоне выделены группы, они попадут в результирующий список, перемежаясь с подстроками между разделителями. Если указан maxsplit, будет произведено не более maxsplit разбиений.

• findall(s)

Ищет все неперекрывающиеся подстроки s, удовлетворяющие шаблону.

• finditer(s)

Возвращает итератор по объектам с результатами сравнения для всех неперекрывающихся подстрок, удовлетворяющих шаблону.

Методы объекта-шаблона

• sub(repl, s)

Заменяет в строке s все (или только count, если он задан) вхождения неперекрывающихся подстрок, удовлетворяющих шаблону, на строку, заданную с помощью repl. В качестве repl может выступать строка или функция. Возвращает строку с выполненными заменами. В первом случае строка тер1 подставляется не просто так, а интерпретируется с заменой вхождений "\номер" на группу с соответствующим номером и вхождений "\q<имя>" на группу с номером или именем имя. В случае, когда repl - функция, ей передается объект с результатом каждого успешного сопоставления, а из нее возвращается строка для замены.

• subn(repl, s)

Аналогичен sub (), но возвращает кортеж из строки с выполненными заменами и числа замен.

```
match = re.search(r'\d\d\D\d\d', r'Телефон 123-12-12')
 2 print(match)
 3 print(match[0] if match else 'Not found')
<re.Match object; span=(9, 14), match='23-12'>
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')
```

```
match = re.fullmatch(r'\d\d\D\d\d', r'T. 12-12')
 2 print('YES' if match else 'NO')
NO
 1 print(re.split(r'\W+', 'Скажите, дядя, ведь не даром...'))
['Скажите', 'дядя', 'ведь', 'не', 'даром', '']
   print(re.findall(r'\d\d\.\d\d\.\d{4}',
                     r'Эта строка написана 24.11.2020, а могла бы и 18.11.2020'))
['24.11.2020', '18.11.2020']
```

Эта строка написана DD.MM.YYYY, а могла бы и DD.MM.YYYY

```
1 import re
 2 delim_re = re.compile(r"[:,;]")
 3 text = "This,is;example"
 4 print(delim re.split(text))
['This', 'is', 'example']
 1 | delim_re = re.compile(r"([:,;])")
 2 print(delim re.split(text))
['This', ',', 'is', ';', 'example']
```

• r"\b\w+\b"

Соответствует слову из букв и знаков подчеркивания.

Соответствует целому числу. Возможно, со знаком.

Число, стоящее в скобках. Скобки используются в самих регулярных выражениях, поэтому они экранируются "\".

• r"[a-cA-C]{2}"

Соответствует строке из двух букв "a", "b" или "c". Например, "Ac", "CC", "bc".

• r"aa|bb|cc|AA|BB|CC"

Строка из двух одинаковых букв.

• r"([a-cA-C])\1"

Строка из двух одинаковых букв, но шаблон задан с использованием групп

•r"aa|bb"

Соответствует "aa" или "bb"

• r"a(a|b)b"

Соответствует "aab" или "abb"

Соответствует строке, которая начинается с набора из восьми или четырех цифр и двоеточия. Все, что идет после двоеточия и после следующих за ним пробелов, выделяется в группу с номером 1, тогда как набор цифр в группу не выделен.

•
$$r''(\w+) = .*\b\1\b''$$

Слова слева и справа от знака равенства присутствуют. Операнд "\1" соответствует группе с номером 1, выделенной с помощью скобок.

• r"(?P<var>\w+)=.*\b(?P=var)\b"

То же самое, но теперь используется именованная группа var.

• r"\bregular(?=\s+expression)"

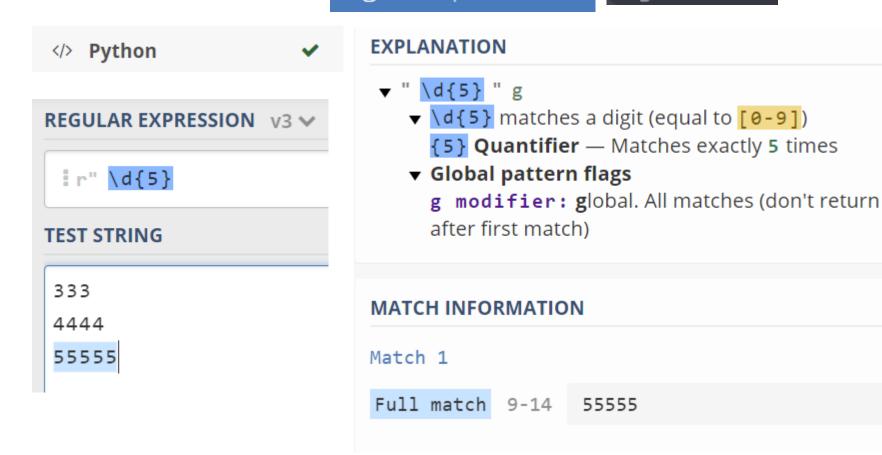
Соответствует слову "regular" только в том случае, если за ним после пробелов следует "expression"

• r"(?<=regular)expression"</pre>

Соответствует слову "expression", перед которым стоит "regular" и один пробел.

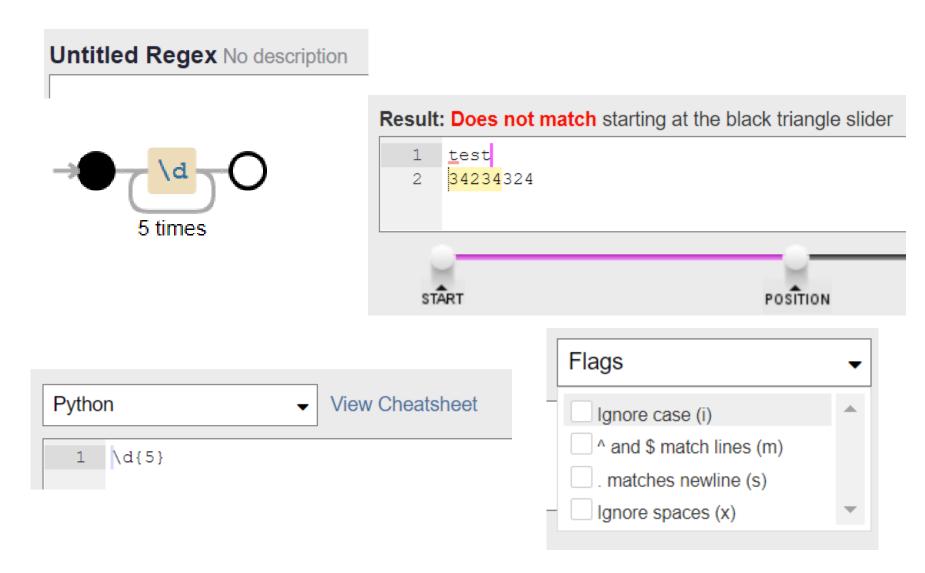
Он-лайн отладка пегуляпных выражений

regular expressions 101 regex 101.com



Он-лайн отладка регулярных выражений

DebuggexBeta www.debuggex.com



Задачи для самостоятельного решения

- Найдите все натуральные числа (возможно, окружённые буквами);
- Найдите все «слова», написанные капсом (то есть строго заглавными), возможно внутри настоящих слов (аааББввв);
- Найдите слова, в которых есть русская буква, а когда-нибудь за ней цифра;
- Найдите все слова, начинающиеся с русской или латинской большой буквы (\b граница слова);
- Найдите слова, которые начинаются на гласную (\b граница слова);

Задачи для самостоятельного решения

- Найдите все натуральные числа, не находящиеся внутри или на границе слова;
- Найдите строчки, в которых есть символ * (. это точно не конец строки!);
- Найдите строчки, в которых есть открывающая и когда-нибудь потом закрывающая скобки;
- Выделите одним махом весь кусок оглавления (в конце примера, вместе с тегами);
- Выделите одним махом только текстовую часть оглавления, без тегов;
- Найдите пустые строчки.