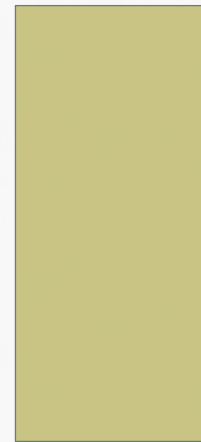


Даты

- 02.12.2014 – лекция
- 03.12.2014 – практика
- **09.12.2014 – (ВТОРНИК!) практика**
- **10.12.2014 – (СРЕДА!) теория. контрольная работа**
- 16.12.2014 – ...
- **17.12.2014 – практика. КОНТРОЛЬНАЯ РАБОТА**
- 23.12.2014 – зачет
- 24.12.2014 – практика. зачет

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ





ИСТОРИЯ

- Кен Томпсон добавил поиск по регулярному выражению в редактор QED в конце 1960-х, позаимствовав нотацию из теоретической статьи Клини
- Из QED регулярные выражения переключались в ed – стандартный текстовый редактор системы

СТАНДАРТЫ

- Стандарт POSIX:
 - Basic Regular Expressions (BRE)
 - Extended Regular Expressions (ERE)

Поддерживаются в утилитах Unix

- Perl Compatible Regular Expressions

Из Perl синтаксис заимствован в Java, .NET, Python, Ruby, JavaScript, и т.д.

ПРИМЕР

- `ls *.txt`

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ

- Формальный язык поиска и осуществления манипуляций с подстроками в тексте, основанный на использовании метасимволов (wildcard)
- По сути это «шаблон», состоящий из символов и метасимволов и задающий правило поиска.

ШАБЛОНЫ

Символы в шаблонах делятся на два типа:

- Литералы – обычные символы
- Метасимволы – символы, которые используются для замены других символов или их последовательностей

ЛИТЕРАЛЫ

Литералы:

- все символы за исключением специальных
 - `[] \ ^ $. | ? * + () { }`
- специальные символы предваренные \
 - например: `\[` или `\$`

МЕТАСИМВОЛ .

Обозначает один любой символ

Пример:

– st..d – регулярное выражение

– под его описание подходит:

standard

stand

astddd

СИМВОЛЬНЫЕ КЛАССЫ

Позволяет указать, что на данном месте в строке может стоять один из перечисленных символов.

- [A-Z] – любая заглавная латинская буква
- [a-d] – строчная буква от a до d
- [A-Za-z0-9] – Латинская буква или цифра
- [А-Яа-яЁё] – любая русская буква

СИМВОЛЬНЫЕ КЛАССЫ

- Спецсимвол отрицания в символьных классах:

^ (крышка)

- [^abc] - все символы (не буквы, а именно символы) кроме букв латинского алфавита a, b, c.

ПЕРЕЧИСЛЕНИЕ

- Вертикальная черта разделяет допустимые варианты. Например, `gray | grey` соответствует `gray` или `grey`
- `gr(a | e)y` описывают строку `gray` или `grey`

ПОЗИЦИЯ ВНУТРИ СТРОКИ

Представление	Позиция	Пример	Соответствие
<code>^</code>	Начало строки	<code>^a</code>	aaa aaa
<code>\$</code>	Конец строки	<code>a\$</code>	aaa aaa
<code>\b</code>	Граница слова	<code>a\b</code>	aaa aaa
		<code>\ba</code>	aaa aaa
<code>\B</code>	Не граница слова	<code>\Ba\b</code>	aaa aaa

КВАНТИФИКАЦИЯ

Представление	Число повторений	Пример	Соответствие
{n}	Ровно n раз	colou{3}r	colouuur
{m,n}	От m до n включительно	colou{2,4}r	colouur, colouuur, colouuuur
{m,}	Не менее m	colou{2,}r	colouur, colouuur, colouuuur и т. д.
{,n}	Не более n	colou{,3}r	color, colour, colouur, colouuur
* {0,}	Ноль или более	colou*r	color, colour, colouur и т. д.
+ {1,}	Одно или более	colou+r	colour, colouur и т. д. (но не color)
? {0,1}	Ноль или одно	colou?r	color, colour

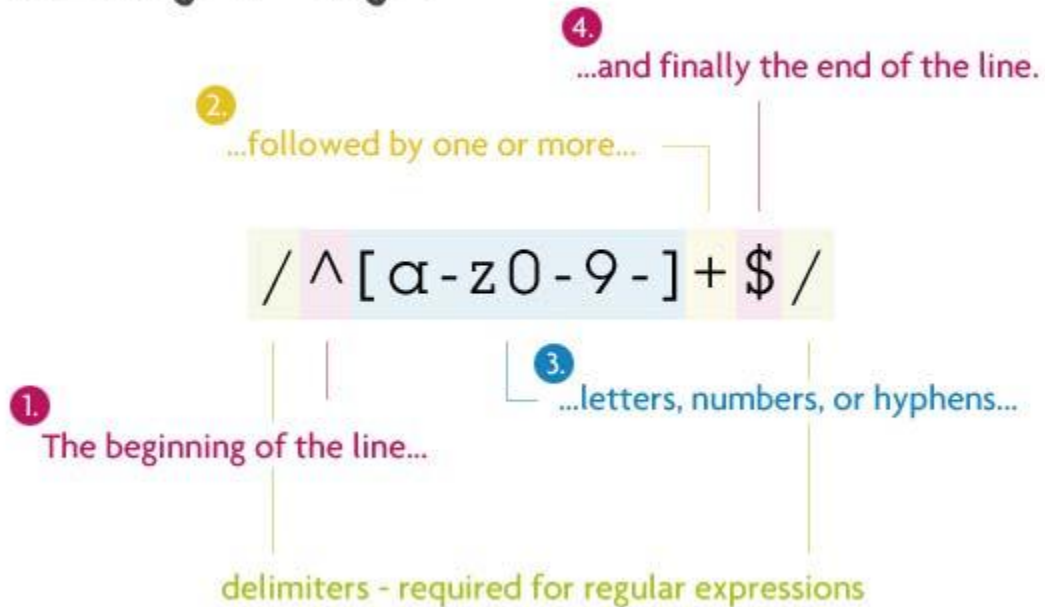
ПРИМЕР

Проверка MAC-адреса

- `/^(?:[0-9A-Fa-f]{2}:){5}[0-9A-Fa-f]{2}$/`

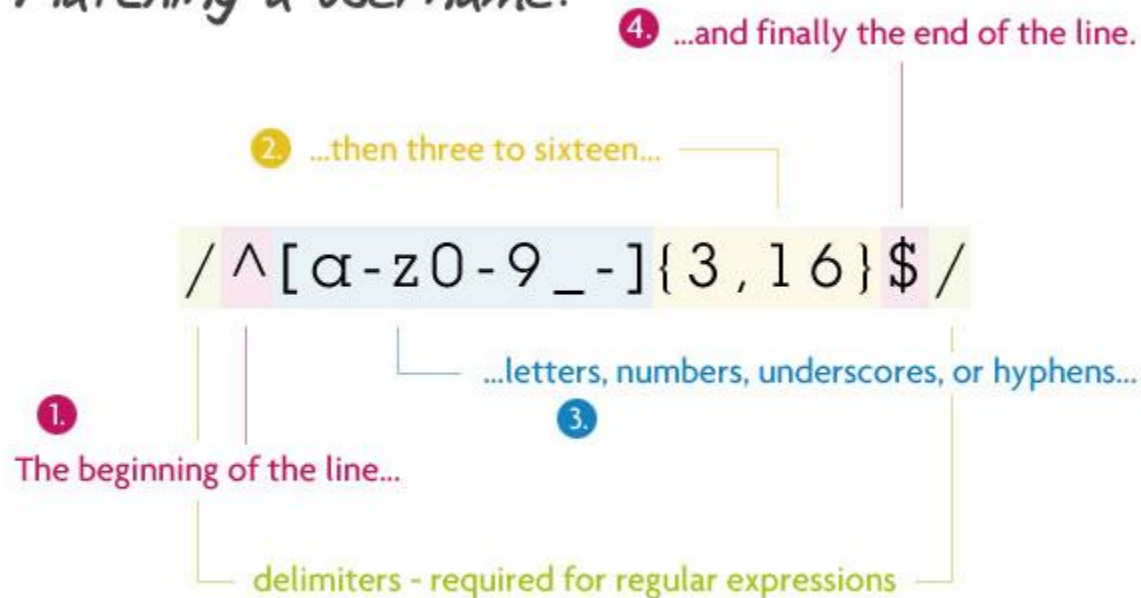
ПРИМЕР

Matching a "slug":



ПРИМЕР

Matching a username:



КВАНТИФИКАЦИЯ

- «Ленивые» выражения
- «Жадные» выражения
- «Ревнивые» (сверхжадные) выражения

Жадный	Ленивый	Ревнивый
*	*?	*+
?	??	?+
+	+	++
{n,}	{n,}?	{n,}+

ЖАДНАЯ КВАНТИФИКАЦИЯ

- Выражение (`<.*>`) соответствует строке, содержащей несколько тегов HTML-разметки, целиком.
- `<p>Википедия` — свободная энциклопедия, в которой `<i>каждый</i>` может изменить или дополнить любую статью`</p>`

ЛЕНИВАЯ КВАНТИФИКАЦИЯ

- Чтобы выделить отдельные теги, можно применить ленивую версию этого выражения: (`<.*?>`) Ей соответствует не вся показанная выше строка, а отдельные теги (выделены цветом):
- `<p>Википедия` — свободная энциклопедия, в которой `<i>каждый</i>` может изменить или дополнить любую статью `</p>`.

РЕВНИВАЯ (СВЕРХЖАДНАЯ) КВАНТИФИКАЦИЯ

- Захватывает самое большое вхождение. В каком-то смысле, ещё «жаднее» жадных и идет дальше них: *один раз что-то «схватив», они никогда не откатываются назад, они не «отдают» кусочки схваченного ими следующим частям регекспа.*

ПРИМЕР

"one" "two" "three" "test" "me"

- ".*"
- ".*?"
- ".*+"

ПРИМЕР

"one" "two" "three" "test" "me"

- ".*" "one" "two" "three" "test" "me"

"one" "two" "three" "test" "me"

- ".*?" "one" "two" "three" "test" "me"

"one"

- ".*+" "one" "two" "three" "test" "me"

Ничего!

ГРУППИРОВКА

- Круглые скобки используются для определения области действия и приоритета операций.
- Например, выражение $(\text{тр}[\text{ау}]\text{м-?})^*$ найдёт последовательность вида трам-трам-трумтрам-трум-трамтрум.

ГРУППИРОВКА С ОБРАТНОЙ СВЯЗЬЮ

- При обработке выражения подстро́ки, найденные по шаблону внутри группы, сохраняются в отдельной области памяти и получают номер начиная с единицы.
- Каждой подстроке соответствует пара скобок в регулярном выражении.
- Квантификация группы не влияет на сохранённый результат, то есть сохраняется лишь первое вхождение.

ПРИМЕР

- Пример:

(та | ту)-\1

- Найдет:

та-та или ту-ту, но не та-ту

ГРУППИРОВКА БЕЗ ОБРАТНОЙ СВЯЗИ

(?:шаблон)

- Под результат такой группировки не выделяется отдельная область памяти и, соответственно, ей не назначается номер.
- Это положительно влияет на скорость выполнения выражения .

АТОМАРНАЯ ГРУППИРОВКА

(?>шаблон)

- Не создает обратных связей.
- Такая группировка запрещает возвращаться назад по строке, если часть шаблона уже найдена.

<code>a(?>bc b x)cc</code>	<code>abccaxcc</code> , но не <code>abccaxcc</code>
<code>a(?>x*)xa</code>	не найдётся <code>axxxa</code>

НАПОМИНАНИЕ

- Существуют три типа регулярных выражений:
 - BRE
 - ERE
 - PCRE

BASIC REGULAR EXPRESSIONS

- Все символы трактуются буквально, исключая перечень в таблице
- $\backslash(.^*[\.,]\backslash)^*$
 - Точка в $[\]$ и вне трактуется по-разному
 - $()$ и $\{ \}$ в качестве синтаксического элемента необходимо предварять “\”
- $([0-9]^*\backslash.[0-9]^*\backslash\$)$
 - Чтобы искать собственно точку, доллар и пр. метасимволы, их нужно предварять “\”
 - $()$ и $\{ \}$ без “\” – ищет сами символы скобок!

EXTENDED REGULAR EXPRESSIONS

- **Добавлено:**

- `? + |`

- **Исключено:**

- `\n` – из-за высокой вычислительной стоимости

- **Изменено:**

- Символы скобок `() { }` как синтаксические элементы не требуют “\” перед собой, для поиска самих этих символов “\” теперь нужен.

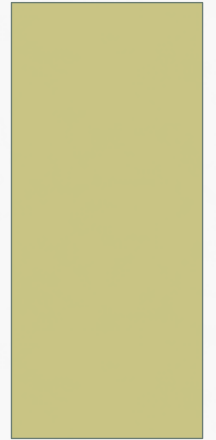
ВЫРАЖЕНИЯ В СТИЛЕ PERL

- Ленивые квантификаторы: `*?`, `+?`, `??`
- Сверхжадные квантификаторы: `*+`, `++`, `?+`
- Сокращенные записи символьных классов:
`\w`, `\W`, `\s`, `\S`, ...
- Lookaheads и lookbehinds – подсказки алгоритму поиска
- Именованные группы связывания (named capture groups)
- Рекурсивные шаблоны.



Python

REGEXP



REGEXP

```
>>> import re
>>> regexp = "{{(. *?)}}"
>>> str = "{{this}} is {{strange}} string"
>>> for match in re.findall(regexp, str):
...     print("FIND: "+match)
...
FIND: this
FIND: strange
```

REGEXP - COMPILED

```
>>> import re
>>> regexp = re.compile("{{(. *?)}}")
>>> str = "{{this}} is {{strange}} string"
>>> for match in regexp.findall(str):
...     print("FIND: "+match)
...
FIND: this
FIND: strange
```

REGEXP

- finditer
- match
- search