

12.09.2023

Регулярные выражения в реальной задаче

Ошуркова Татьяна,
Росбанк. Старший ИТ-менеджер



настоящие
возможности



О себе

Работаю в банковской сфере



Росбанк, Кредит Урал Банк (группа Газпромбанка)

Разработчик баз данных



А также web-разработчик, аналитик

Преподаватель программирования



Преподавала курсы для детей и взрослых



О чем будем говорить



Что такое регулярные выражения



Наиболее частые области применения регулярных выражений



Синтаксис регулярных выражений



Примеры использования регулярных выражений

Регулярные выражения

Регулярные выражения - это механизм для поиска и замены текста



Задачи регулярок



- Определить наличие последовательности символов
- Найти последовательность символов
- Заменить последовательность символов



Примеры реальных задач регулярных выражений



Выполнить валидацию данных

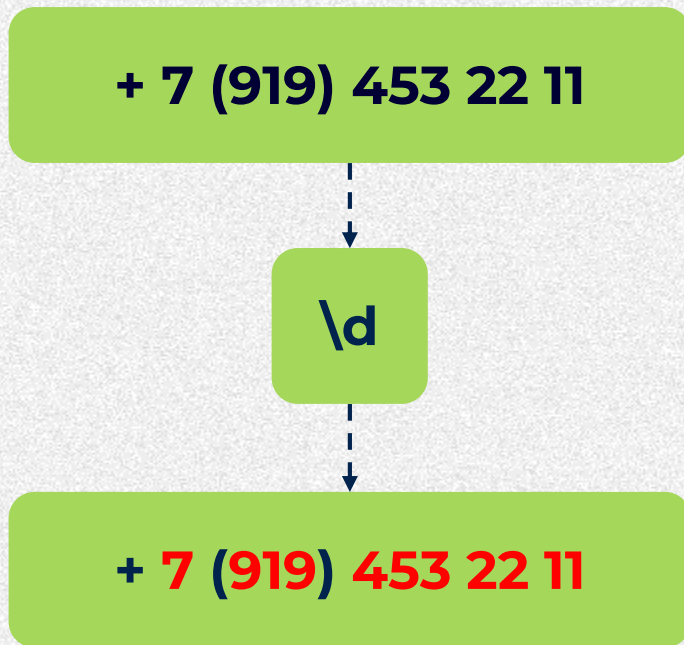
Проверить тип файла

Найти данные заданного формата

Удалить некорректные символы

Как составлять регулярки

Символьные классы (character classes)



<code>\d</code>	digit - цифра
<code>\s</code>	space - пробел
<code>\w</code>	word - слово
<code>\D</code>	не цифра (все кроме <code>\d</code>)
<code>\S</code>	не пробел (все кроме <code>\s</code>)
<code>\W</code>	не латиница, не знак подчёркивания и не цифра
<code>.</code>	любой символ кроме новой строки

Символьные классы

Примеры

REGULAR EXPRESSION

:/ \w

TEST STRING

abcd_ef • (145) • (14.4) • (15/5)

REGULAR EXPRESSION

:/ \s

TEST STRING

abcd_ef • (145) • (14.4) • (15/5)

REGULAR EXPRESSION

:/ \W

TEST STRING

abcd_ef • (145) • (14.4) • (15/5)

REGULAR EXPRESSION

:/ \S

TEST STRING

abcd_ef • (145) • (14.4) • (15/5)

Как составлять регулярки

Начало и конец строки

Начало строки	Конец строки	Совпадение строки
<p>test request bg request start end request</p> <p>↓</p> <p>^request</p> <p>↓</p> <p>test request bg request start end request</p>	<p>end of req test end start end bg</p> <p>↓</p> <p>end\$</p> <p>↓</p> <p>end of req test end start end bg</p>	<p>11:12 max call 16:45</p> <p>↓</p> <p>^\d\d:\d\d\$</p> <p>↓</p> <p>11:12 max call 16:45</p>

Как составлять регулярки

Наборы (sets) и диапазоны (ranges)

Набор	Диапазон	Исключающий диапазон
<p data-bbox="186 529 825 682">srer strat seer</p> <p data-bbox="326 815 685 968">s[are]er</p> <p data-bbox="186 1100 825 1253">srer strat seer</p>	<p data-bbox="958 529 1602 682">reql reqt req0 req_</p> <p data-bbox="1059 815 1500 968">req[a-zl-9_]</p> <p data-bbox="958 1100 1602 1253">reql reqt req0 req_</p>	<p data-bbox="1730 529 2374 682">reql reqt req0 req_</p> <p data-bbox="1832 815 2272 968">req[^l-3a-c]</p> <p data-bbox="1730 1100 2374 1253">reql reqt req0 req_</p>

Немного практики

Кредитные договоры хранятся в формате трёхзначный идентификатор-год. Какие кредитные договоры **подходят под написанное регулярное выражение?**

1. 25c-2022
2. 31f-2022
3. 32b-2021
4. 72g-2020
5. 55a-2024

`\d[2,5][a-e]-202[0-2]`



Немного практики

Кредитные договоры хранятся в формате трёхзначный идентификатор-год. Какие кредитные договоры **подходят под написанное регулярное выражение?**

1. 25c-2022
2. 31f-2022
3. 32b-2021
4. 72g-2020
5. 55a-2024

`\d[2,5][a-e]-202[0-2]`

1. **25c-2022**
2. 31f-2022
3. **32b-2021**
4. 72g-2020
5. 55a-2024

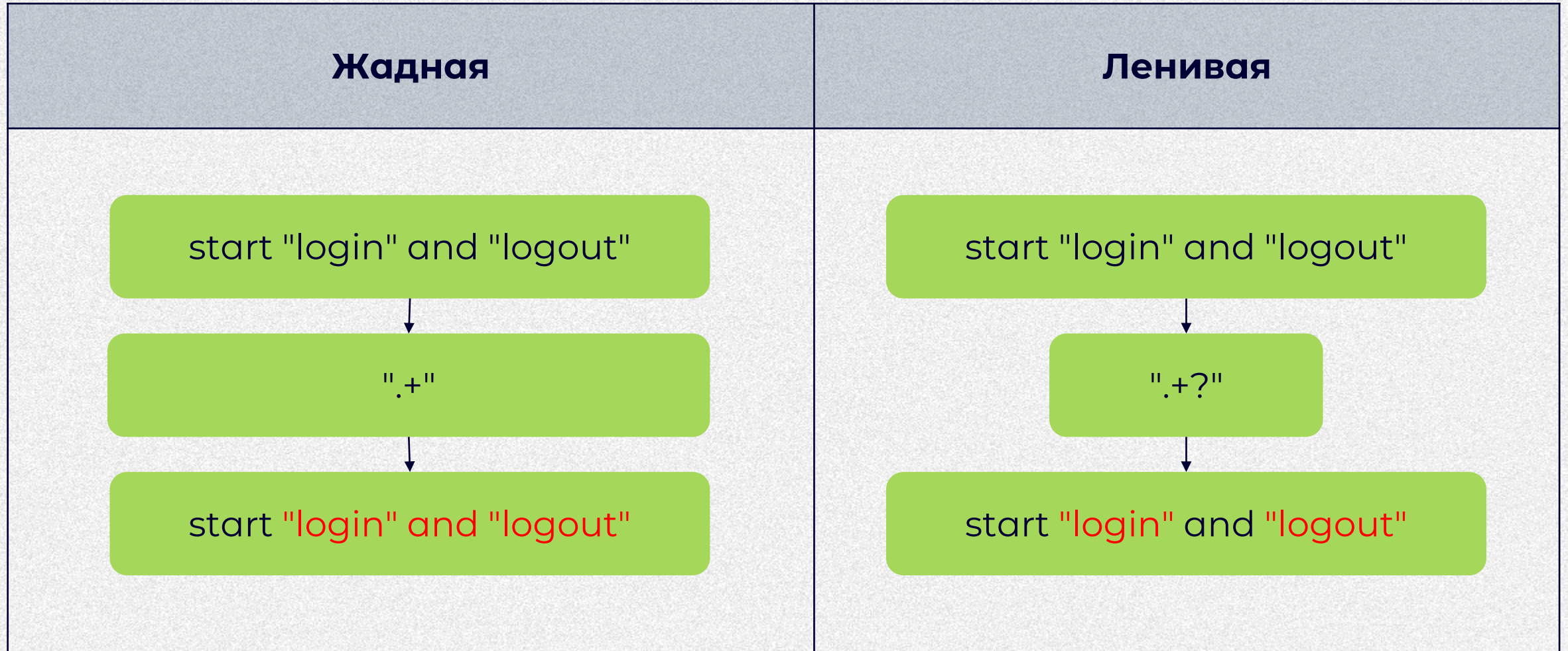
Как составлять регулярки

Квантификаторы (quantifiers)

Количество	Обозначения	
<p>start 123wrw 12tyws</p> <p>↓</p> <p><code>\d{3,5}\w{3,}</code></p> <p>↓</p> <p>start 123wrw 12tyws</p>	?	{0,1}
	+	{1,}
	*	{0,}
		<p>100 15 10.24 45.782 145.3 78.675</p> <p>↓</p> <p><code>[\s^\d+\.\d{1,2}[\s\$]</code></p> <p>↓</p> <p>100 15 10.24 45.782 145.3 78.675</p>

Как составлять регулярки

Жадная (greedy) и ленивая (lazy) квантификация



Как составлять регулярки

Группы (groups)

my@mail.com
my?@mail
his@site.com.uk
Mike.Tyson@yandex

`[-.\w]+@[([\w-]+\.)+\w+`

my@mail.com
my?@mail
his@site.com.uk
Mike.Tyson@yandex

20.02.2022
2023-11-22
10-10-2012

`(\d{1,2})([-.V])(\d{1,2})([-.V])(\d{4})`

20.02.2022
2023-11-22
10-10-2012

Как составлять регулярки

Группы с захватом (capturing) и без (non-capturing)

20.02.2022
2023-11-22
10-10-2012
10-10.2012

`(\d{1,2})([-.V])(\d{1,2})\2(\d{4})$`

20.02.2022
2023-11-22
10-10-2012
10-10.2012

20.02.2022
2023-11-22
10-10-2012
10-10.2012

`(?:\d{1,2})([-.V])(\d{1,2})\1(\d{4})$`

20.02.2022
2023-11-22
10-10-2012
10-10.2012

Немного практики

Необходимо выполнить соответствие логина, который создает пользователь, определенным требованиям. **Какие из логинов корректны в соответствии с написанным для валидации регулярным выражением?**

1. kate23@mail
2. nikita.8.9
3. maria_petrova
4. evgenia1995
5. oleg2-23
6. Dima01

`[a-z0-9_.]{3,11}`



Немного практики

Необходимо выполнить соответствие логина, который создает пользователь, определенным требованиям. **Какие из логинов корректны в соответствии с написанным для валидации регулярным выражением?**

1. kate23@mail
2. nikita.8.9
3. maria_petrova
4. evgenia1995
5. oleg2-23
6. Dima01

[a-z0-9_.]{3,11}

1. kate23@mail
2. **nikita.8.9**
3. maria_petrova
4. **evgenia1995**
5. oleg2-23
6. Dima01

Какие задачи мы рассмотрим

Валидация данных и определение направления платежа в платежном поручении

Разбор логов. Найти нужный запрос, достать строку с ключевыми словами

Задача с платежным поручением

Определение правил
в реквизитах

Формализация
средствами регулярных
выражений

Создание sql запроса

Определение правил в реквизитах

Корреспондентский счет

- Начало счета
- Соответствие счету получателя

Счет получателя платежа

- Начало счета (план)
- Сочетание цифр в счете

Формализация правил

Корр. счет	Счёта получателя
(*)	^475(*)
(*)	^40[392]09.....4(*)
(*)	^60323(*)
<div data-bbox="122 729 489 863" style="background-color: #002060; color: white; padding: 5px;">Начало строки</div> ^40102(*)	<div data-bbox="1082 722 1454 856" style="background-color: #002060; color: white; padding: 5px;">Список совпадения</div> ^3100(*)
^40102(*)	^3212(*)
^40102(*)	^3222(*)
^40102(*)	<div data-bbox="2015 936 2440 1071" style="background-color: #002060; color: white; padding: 5px;">Любые символы и количество</div> ^41[892]07.....5(*)
^40102(*)	^03242(*)

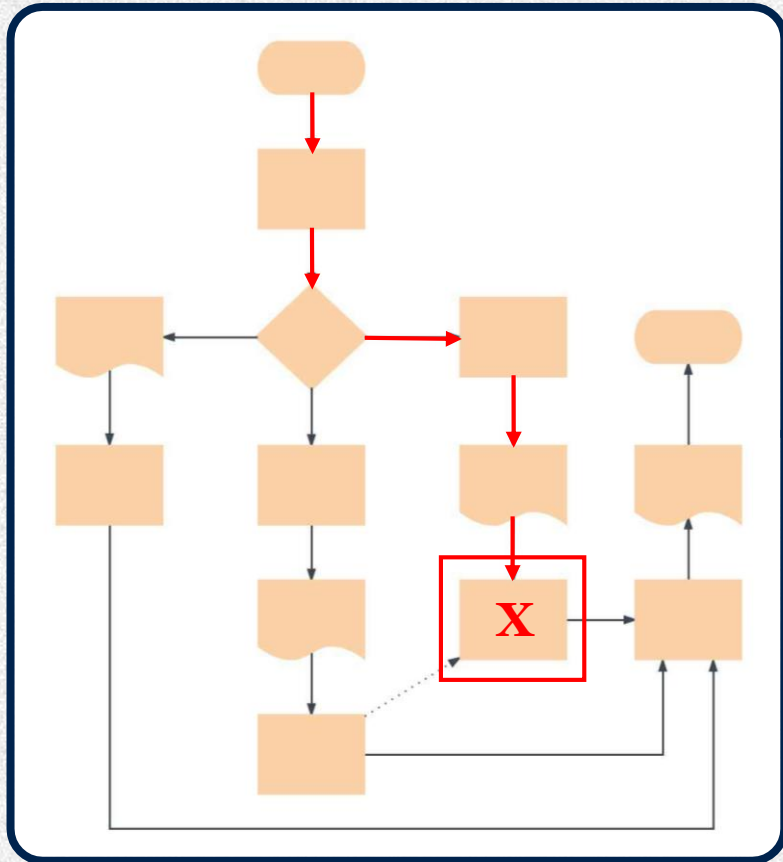
Регулярные выражения в запросе

Таблица payment_rules

BIC	CORR_ACCOUNT	PAY_ACCOUNT	KBK	DIRECTION
(*)	(*)	^68097(*)	(*)	Перевод внутри банка
(*)	^40307(*)	^3541(*)	(*)	Бюджетный платеж
(*)	^40307(*)	^3921(*)	(*)	Бюджетный платеж

```
select pr.DIRECTION
from payment_rules pr
where REGEXP_LIKE(v_Bank_Bik, pr.BIC)
and REGEXP_LIKE(v_Corr_Account, pr.CORR_ACCOUNT)
and REGEXP_LIKE(v_Pay_Account, pr.PAY_ACCOUNT)
and REGEXP_LIKE(v_KBK, pr.KBK)
```

Поиск ошибки



BIC	CORR_ACCOUNT	PAY_ACCOUNT	KBK	DIRECTION
(*)	(*)	^68097(*)	(*)	Перевод внутри банка
(*)	^40307(*)	^3541(*)	(*)	Бюджетный платеж
(*)	^40307(*)	^3921(*)	(*)	Бюджетный платеж

```
select pr.DIRECTION
from payment_rules pr
where REGEXP_LIKE(v_Bank_Bik, pr.BIC)
and REGEXP_LIKE(v_Corr_Account, pr.CORR_ACCOUNT)
and REGEXP_LIKE(v_Pay_Account, pr.PAY_ACCOUNT)
and REGEXP_LIKE(v_KBK, pr.KBK)
```

Пример поиска ошибки

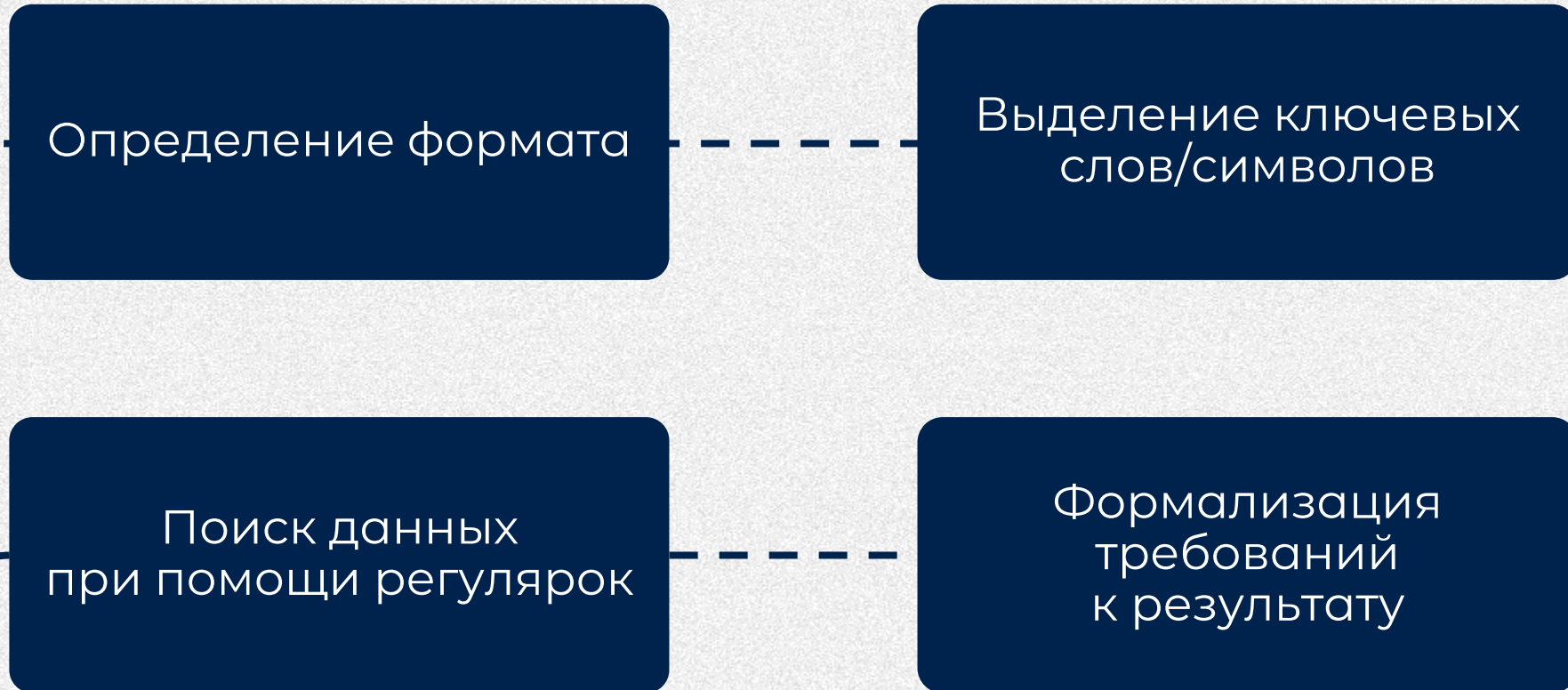


BIC	CORR_ACCOUNT	PAY_ACCOUNT	КБК	DIRECTION
(*)	^40307(*)	^35 <u>4</u> 1(*)	(*)	Бюджетный платеж
(*)	^40307(*)	^35 <u>2</u> 1(*)	(*)	Перевод внутри банка

Настройки с использованием регулярных выражений

BIC	CORR_ACCOUNT	PAY_ACCOUNT	KBK	DIRECTION
(*)	(*)	^68097(*)	(*)	Перевод внутри банка
(*)	^40307(*)	^40[392]09.....4(*)	(*)	Бюджетный платеж
(*)	^40308(*)	^40[392]07.....5(*)	(*)	Платеж в другой банк

Задача с логами



Файл лога

```
03-29-13:17:38.600-14828_14828:I*open_connection-...some_text.␣
03-29-13:17:38.651-14828-14828:D*MainModule:␣
03-29-13:17:40.655_14828-14828:D*InputProcessor:8␣
03-29-13:17:40.657-14828-14828:D*MainModule:восемь␣
03-29_13:17:40.658-14828_15420:D*InputCommand:8␣
...some_text...␣
...some_text...␣
...some_text...␣
03-29_13:17:41.239-14828-15420:D*WmsApi:Request:Login_request('03-429-
2022_13:17:40','arg_1','arg_2','arg_31','arg_A','arg_51','arg_6','arg_7','arg_8','arg_9',
'arg_101','arg_11','arg_12','[CDATA[0.....]]')␣
...some_text...␣
...some_text...␣
...some_text...
```

Большой
объем

Нерелевантные
данные

Сложные
правила отбора

Определение формата

Цифровые
данные

REGULAR EXPRESSION 1 match (518 steps, 0.4ms)

```
:/ ^.+open_connection.+$/ gm
```

TEST STRING

```
03-29-13:17:38.600-14828_14828:I•open_connection-...some_text.↵
03-29-13:17:38.651-14828-14828:D•MainModule:↵
03-29-13:17:40.655_14828-14828:D•InputProcessor:8↵
03-29-13:17:40.657-14828-14828:D•MainModule:восемь↵
03-29_13:17:40.658-14828_15420:D•InputCommand:8•↵
...some_text...↵
...some_text...↵
...some_text...↵
03-29_13:17:41.239-14828-15420:D•WmsApi:Request:Login_request('03-429-
2022_13:17:40','arg_1','arg_2','arg_31','arg_A','arg_51','arg_6','arg_7'
,'arg_8','arg_9','arg_101','arg_11','arg_12','[CDATA[0.....]]')↵
...some_text...↵
...some_text...↵
...some_text...
```

Ключевые
слова

Выделение ключевых слов/символов

Нужны данные
из Login_request

```
REGULAR EXPRESSION 1 match (18 steps, 0.1ms)  
:/ Login_request\(.*\)/ gm  
TEST STRING  
03-29-13:17:38.600-14828_14828:I•open_connection-...some_text.␣  
03-29-13:17:38.651-14828-14828:D•MainModule:␣  
03-29-13:17:40.655_14828-14828:D•InputProcessor:8␣  
03-29-13:17:40.657-14828-14828:D•MainModule:восемь␣  
03-29_13:17:40.658-14828_15420:D•InputCommand:8␣  
...some_text...␣  
...some_text...␣  
...some_text...␣  
03-29_13:17:41.239-14828-15420:D•WmsApi:Request:Login_request('03-429-  
2022_13:17:40','arg_1','arg_2','arg_31','arg_A','arg_51','arg_6','arg_7'  
, 'arg_8', 'arg_9', 'arg_101', 'arg_11', 'arg_12', '[CDATA[0.....]]')␣  
...some_text...␣  
...some_text...␣  
...some_text...
```

Требования к результату

REGULAR EXPRESSION 1 match (845 steps, 0.6ms)

```
:/ (?<=Login_request\(\).*?(?!\)) /gm
```

TEST STRING

```
03-29-13:17:38.600-14828_14828:I•open_connection-...some_text.␣
03-29-13:17:38.651-14828-14828:D•MainModule:␣
03-29-13:17:40.655_14828-14828:D•InputProcessor:8␣
03-29-13:17:40.657-14828-14828:D•MainModule:восемь␣
03-29_13:17:40.658-14828_15420:D•InputCommand:8•␣
...some_text...␣
...some_text...␣
...some_text...␣
03-29_13:17:41.239-14828-15420:D•WmsApi:Request:Login_request('03-429-
2022_13:17:40','arg_1','arg_2','arg_31','arg_A','arg_51','arg_6','arg_7'
,'arg_8','arg_9','arg_101','arg_11','arg_12','[CDATA[0.....]]')␣
...some_text...␣
...some_text...␣
...some_text...
```

$(?<=Y)X$
ищет совпадение с X,
если перед ним есть Y

$(?<!Y)X$
ищет совпадение с X
при условии, что перед
ним нет Y

Опережающие
и ретроспективные проверки
(lookahead and lookbehind)

Проверка регулярного выражения

REGULAR EXPRESSION

```
:/ (?<=test)\.txt
```

TEST STRING

```
test.txt↵  
readme.txt↵  
mytest.txt↵  
main_file.txt
```

SUBSTITUTION

```
.log
```

```
test.log↵  
readme.txt↵  
mytest.log↵  
main_file.txt
```

```
1 import re  
2  
3 regex = r"(?<=test)\.txt"  
4  
5 test_str = ("test.txt\n"  
6            "readme.txt\n"  
7            "my_test.txt\n"  
8            "main_file.txt")  
9  
10 subst = ".log"  
11  
12 result = re.sub(regex, subst, test_str)  
13  
14 print(result)  
15  
16 # test.log  
17 # readme.txt  
18 # my_test.log  
19 # main_file.txt
```

Ограничения использования регулярных выражений



Ограниченность функционала

Зависимость от применяемой технологии

Сложность прочтения/понимания

Сложная адаптация одной реализации
к разным задачам

Итоги. Плюсы использования регулярных выражений

Возможность решения широкого класса задач

Необходимый инструмент при работе с текстовыми документами

Кроссплатформенность

Быстрый поиск и анализ текстовых данных



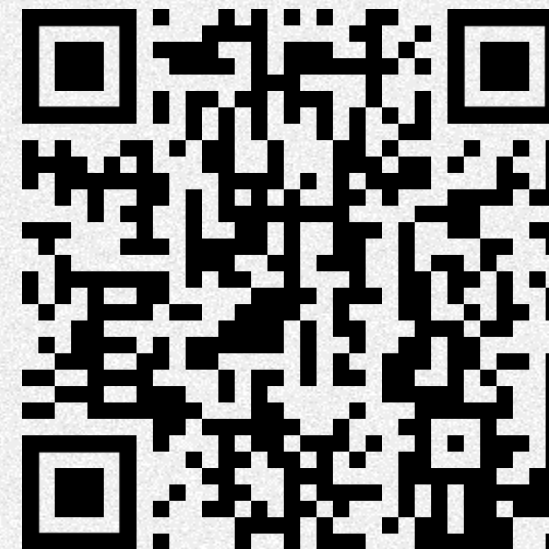
Будет полезно



regex101.com



regexcrossword.com



[справка на GitHub](https://docs.github.com)