

Reg(ular)?Ex(pressions)?

A Multi-Tool for Working with Text

James Truitt

BitCurator Forum Philadelphia Satellite

2024-03-19

2/14/1865

5/22/1890

7/4/1878

9/30/1899

3/15/1885

1865-02-14

1890-05-22

1878-07-04

1899-09-30

1885-03-15

Keeran, Vincent What Is
Causing the Current Rise
In Prices 10/12/47

Keeran, Vincent.
What Is Causing the Current
Rise In Prices?
1947-10-12

personal|primary|naf|http://id.loc.gov/
authorities/names/|http://id.loc.gov/au
thorities/names/n85329150|Lockwoo
d, Belva Ann,
1830-1917|text|marcrelator|creator

relators:cre:person:Lockwood, Belva
Ann, 1830-1917

What are Regular Expressions?

A language
for pattern matching

Why Use Regular Expressions?

- Flexible
- Powerful
- Available
- Output-focused

Where to Use Regular Expressions?


```

1 import re
2 dates = [
3     '02/14/1865', '05/22/1890', '07/04/1878', '09/30/1899', '03/15/1885', '10/11/1871',
4     '08/29/1869', '01/05/1880', '11/19/1894', '12/12/1875', '04/26/1887', '06/07/1892',
5     '02/08/1867', '05/31/1895', '08/16/1873', '07/20/1882', '10/04/1897', '01/18/1868',
6     '03/23/1884', '09/05/1891', '12/29/1876', '11/10/1889', '04/13/1870', '06/21/1893',
7     '02/27/1879', '07/02/1886', '10/28/1898', '01/09/1872', '08/12/1881', '12/03/1866']
8
9 pattern2 = r'(\d\d)\./(\d\d)\./(\d{4})'
10 repl2 = r'\3-\1-\2'
11 for date in dates:
12     print(date, '→', re.sub(r'(\d\d)\./(\d\d)\./(\d{4})', r'\3-\1-\2', date), sep='\t')

```

OUTPUT DEBUG CONSOLE TERMINAL PORTS

✓ TERMINAL

● 14:55 ~/Downloads ✗ python3 regextest.py

```

02/14/1865 → 1865-02-14
05/22/1890 → 1890-05-22
07/04/1878 → 1878-07-04
09/30/1899 → 1899-09-30
03/15/1885 → 1885-03-15
10/11/1871 → 1871-10-11
08/29/1869 → 1869-08-29
01/05/1880 → 1880-01-05
11/19/1894 → 1894-11-19
12/12/1875 → 1875-12-12
04/26/1887 → 1887-04-26
06/07/1892 → 1892-06-07
02/08/1867 → 1867-02-08
05/31/1895 → 1895-05-31
08/16/1873 → 1873-08-16
07/20/1882 → 1882-07-20
10/04/1897 → 1897-10-04
01/18/1868 → 1868-01-18
03/23/1884 → 1884-03-23
09/05/1891 → 1891-09-05
12/29/1876 → 1876-12-29
11/10/1889 → 1889-11-10
04/13/1870 → 1870-04-13
06/21/1893 → 1893-06-21
02/27/1879 → 1879-02-27
07/02/1886 → 1886-07-02
10/28/1898 → 1898-10-28
01/09/1872 → 1872-01-09
08/12/1881 → 1881-08-12
12/03/1866 → 1866-12-03

```

○ 14:57 ~/Downloads ✗ □

```

1 import re
2 dates = [
3     '02/14/1865', '05/22/1890', '07/04/1878', '09/30/1899', '03/15/1885', '10/11/1871',
4     '08/29/1869', '01/05/1880', '11/19/1894', '12/12/1875', '04/26/1887', '06/07/1892',
5     '02/08/1867', '05/31/1895', '08/16/1873', '07/20/1882', '10/04/1897', '01/18/1868',
6     '03/23/1884', '09/05/1891', '12/29/1876', '11/10/1889', '04/13/1870', '06/21/1893',
7     '02/27/1879', '07/02/1886', '10/28/1898', '01/09/1872', '08/12/1881', '12/03/1866']
8
9 pattern2 = r'(\d\d)\.(\d\d)\.(\d{4})'
10 repl2 = r'\3-\1-\2'
11
12 for date in dates:
13     print(date, '→', re.sub(r'(\d\d)\.(\d\d)\.(\d{4})', r'\3-\1-\2', date), sep='\t')

```

OUTPUT DEBUG CONSOLE **TERMINAL** PORTS

▼ TERMINAL

● 14:55 ~/Downloads ✗ python3 regextest.py

```

02/14/1865 → 1865-02-14
05/22/1890 → 1890-05-22
07/04/1878 → 1878-07-04
09/30/1899 → 1899-09-30
03/15/1885 → 1885-03-15
10/11/1871 → 1871-10-11
08/29/1869 → 1869-08-29
01/05/1880 → 1880-01-05
11/19/1894 → 1894-11-19
12/12/1875 → 1875-12-12
04/26/1887 → 1887-04-26
06/07/1892 → 1892-06-07
02/08/1867 → 1867-02-08
05/31/1895 → 1895-05-31
08/16/1873 → 1873-08-16
07/20/1882 → 1882-07-20
10/04/1897 → 1897-10-04
01/18/1868 → 1868-01-18
03/23/1884 → 1884-03-23
09/05/1891 → 1891-09-05
12/29/1876 → 1876-12-29
11/10/1889 → 1889-11-10
04/13/1870 → 1870-04-13
06/21/1893 → 1893-06-21
02/27/1879 → 1879-02-27
07/02/1886 → 1886-07-02
10/28/1898 → 1898-10-28
01/09/1872 → 1872-01-09
08/12/1881 → 1881-08-12
12/03/1866 → 1866-12-03

```

○ 14:57 ~/Downloads ✗

An important aspect of identity investigations (as well as other types) is the ability to search the data for a list of keywords. *bulk_extractor* provides the capability to do that through two different means. First, the *find* scanner is a simple **regular expression** finder that uses **regular expressions**. The *find* scanner looks through the data for anything listed in the global find list. The format of the find list should be rows of **regular expressions** while any line beginning with a # is considered a comment. The following is an excerpt from a sample find list file:

```
# This is a comment line
\b\d{1,3}\.\d{1,3}\.\d{1,3}\.\d{1,3}\b
# another comment line
/^[a-z0-9_-]{3,16}$/
```

The first **regular expression** from the above example, beginning with \b, looks for the following in order: a word boundary followed a digit repeated between 1-3 times, a digit repeated between 1-3 times, a digit repeated 1-3 times, a '.', a digit repeated 1-3 times,

OpenRefine

allHazardNYYM

Permalink

Facet / Filter

Undo / Redo 1 / 1

118063 rows

Refresh

Reset all

Remove all

First name

Invert reset

☐ case sensitive

☒ regular expression

Show as: rows records

Show: 5 10 25 50 100 500 1000 rows

All

ID

First name

Last name

Volume

Page

Note

1.

0

Enos

Chase

1.1

4

Cornwal MM to m

2.

1

Joshua

White

1.1

7

Coemans MM with w Deborah and 4 ch: Michael, Job, John & Phebe. See also pg 4 and C 733 Vol 3.7 pg 5

Replace

Find

☐ case insensitive

☐ whole word

☒ regular expression

Leave blank to add the replacement string after each character.
Check "regular expression" to find special characters (new lines, tabulations...) or complex patterns.

Replace with

☒ use \n for new lines, \t for tabulation, \\n for \n, \\t for \t.

If "regular expression" option is checked and finding pattern contains groups delimited with parentheses, \$0 will return the complete string matching the pattern, and \$1, \$2... the 1st, 2nd... group.

OK

Cancel

OpenRefine

allHazardNYYM

Permalink

Facet / Filter

Undo / Redo 1 / 1

118063 rows

Refresh

Reset all

Remove all

First

Invert reset

☐ case s

☒ regular expression

Show as: rows records

Show: 5 10 25 50 100 500 1000 rows

All

ID

First name

Last name

Volume

Page

Note

1.

0

Enos

Chase

1.1

4

Cornwal MM to m

2.

1

Joshua

White

1.1

7

Coemans MM with w Deborah and 4 ch: Michael, Job, John & Phebe. See also pg 4 and C 733 Vol 3.7 pg 5

Custom text transform on column Note

Expression

Language

General Refine Expression Language (GREL) ▾

value

.replace(/REGEX/, "string")

.split(/REGEX/

.find(/REGEX/|

.contains(/REGEX/)

No syntax error.

Preview

History

Starred

Help

row

value

value .replace(/M\?M\?/, "Mo ...

1.

Cornwal MM to m

Error:
java.lang.ArrayIndexOutOfBoundsException
Index 0 out of bounds for length 0

2.

Coemans MM with w Deborah and 4 ch: Michael, Job, John & Phebe. See also pg 4 and C 733 Vol 3.7 pg 5

false

3.

Mary Jacacks dt of Benjamin Jacacks, dec & Sarah. S of Tiddeman,dec & Elizabeth.

true

Boring Google Sheet ☆ 📁 ☁

File Edit View Insert Format Data Tools Extensions Help

🔍 ↶ ↷ 🖨 📋 100% | \$ % .0 .00 123 | Defaul... ▾ | - 10 +

A13 ▾ | fx =regex

13 =regex

REGEXMATCH
Whether a piece of text matches regular expression.
REGEXEXTRACT
REGEXREPLACE
Tab to accept. ⬆ ⬇ to navigate

Boring Google Doc ☆ 📁 ☁

File Edit View Insert Format Data Tools Extensions Help

🔍 ↶ ↷ 🖨 📋 100% | \$ % .0 .00 123 | Defaul... ▾ | - 10 +

E16 ▾

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

Find and replace

×

Find

Replace with

Search

All sheets ▾

☐ Match case

☐ Match entire cell contents

☒ Search using regular expressions [Help](#)

☐ Also search within formulas

☐ Also search within links

Find

Replace






Replace all

Done


Learning the Syntax


regex101.com

bit.ly/bcc-regex


    


SAVE & SHARE


 Save new Regex %+s


 Add to Community Li...


FLAVOR


 PCRE2 (PHP >=7.3)


 **PCRE (PHP <7.3)** ✓


 ECMAScript (JavaScri...

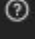
 Python

 Golang


 Java 8

 .NET 7.0 (C#)

 Rust

 [Regex Flavor Guide](#)

FUNCTION


 Match

PLEASE SUPPORT REGEX101

If you're running an ad blocker, consider whitelisting regex101 to support the website. [Read more.](#)

REGULAR EXPRESSION

no match (2 steps, 0.0ms) 

 / REGEX-PATTERN

/ gmi



TEST STRING

TEXT-STRING

SUBSTITUTION

no result

insert your replacement value here

TEXT-STRING

Most Characters Match Themselves

Whitespace

Pattern	Meaning
<code>\t</code>	Tab
<code>\n</code>	Newline
<code>\r</code>	Carriage return

Boundaries

Pattern	Meaning
<code>^</code>	Start of line
<code>\$</code>	End of line
<code>\b</code>	Word boundary

Character Classes

Pattern	Matches
<code>[abcde]</code>	catherine's 3 dogs
<code>^[abcde]</code>	catherine's 3 dogs
<code>[3te]</code>	catherine's 3 dogs
<code>^[3te]</code>	catherine's 3 dogs
<code>[0-9]</code>	catherine's 3 dogs
<code>^[0-9]</code>	catherine's 3 dogs
<code>[a-e3]</code>	catherine's 3 dogs

Character Class Shorthands

Shorthand	Equivalent to
.	<code>[^\n\r]</code>
<code>\w</code>	<code>[a-zA-Z0-9_]</code>
<code>\W</code>	<code>[^\w]</code>
<code>\s</code>	<code>[\n\r\t]</code>
<code>\S</code>	<code>[^\s]</code>
<code>\d</code>	<code>[0-9]</code>
<code>\D</code>	<code>[^\d]</code>
<code>[^\w\s]</code>	Punctuation
<code>[\s\S]</code>	Any character

Quantifiers

Pattern	Meaning
$X?$	0 or 1 instances of X
X^*	0 or more instances of X
X^+	1 or more instances of X

Quantifiers

Greedy	Lazy	Meaning
<code>X?</code>	<code>X??</code>	0 or 1 instances of X
<code>X*</code>	<code>X*?</code>	0 or more instances of X
<code>X+</code>	<code>X+?</code>	1 or more instances of X

(Capture Groups)

Pattern	Matches
<code>ab+</code>	<code>ab</code> , <code>abb</code> , <code>abbb</code>
<code>(ab)+</code>	<code>ab</code> , <code>abab</code> , <code>ababab</code>
<code>(ab)\1</code>	<code>abab</code>
<code>(..)\1</code>	<code>abab</code> , <code>cdcd</code> , <code>adad</code>
<code>(\d\d)\1</code>	<code>03/03</code> , <code>07/07</code> , <code>12/12</code>

(Replacement Patterns)

\1, \$1

Resources for learning more

- Tutorials
 - Library Carpentry:
<https://librarycarpentry.org/lc-data-intro/index.html>
 - RegexOne: <https://regexone.com/>
 - The Programming Historian:
<https://programminghistorian.org/en/lessons/understanding-regular-expressions>
 - Regular-Expressions.info:
<https://www.regular-expressions.info/tutorial.html>
- Cheatsheets
 - Core RegEx:
<https://web.mit.edu/hackl/www/lab/turkshop/slides/regex-cheatsheet.pdf>
 - Python's RegEx flavor:
<https://www.debuggex.com/cheatsheet/regex/python>
- Tools
 - <https://regexr.com/>
 - <https://regex101.com>
- Exercises: <http://alf.nu/RegexGolf>