

Python Regex

Steven Wang

August 17, 2016

```
1 # Author : Steven Wang  Date: 20160816
2
3 import re
4
5 # r"string" : "\" means just a backslash, no "escape"
6
7 if re.search("ape","The ape was at the apex"):
8     print("There is an ape")
9
10 allApes = re.findall("ape.", "The ape was at the apex")
11 # '.' will replace any single character or single space
12 for i in allApes:
13     print(i)
14
15 theStr = "The ape was at the apex"
16
17 for i in re.finditer("ape.", theStr):
18     locTuple = i.span()
19     # return a Tuple that contain its location in the string
20     print(locTuple)
21     # >>> (4, 8)
22     # >>> (19, 23)
23     print(theStr[locTuple[0]:locTuple[1]])
24     # >>> ape
25     # >>> apex
26
27 animalStr = "Cat rat mat pat"
28 allAnimals = re.findall("[crmf]at", animalStr)
29 # match certain character
30 for i in allAnimals:
31     print(i)
32
33 # >>>rat
34 # >>>mat
35 # >>>pat
36
```



```

83 # >>>This is a long string that goes on for many lines
84
85 # \b : Backspace
86 # \f : Form Feed
87 # \r : Carriage Return
88 # \t : Tab
89 # \v : Vertical Tab
90
91 # \r\n : This would show up a lot in windows and html
92
93 # Matching numbers
94 # \d : [0-9]      These 2 are the same, will match any number
95 # \D : [^0-9]    These 2 are the same, will match any character
                   but a number
96
97 randStr = "12345"
98 print("Matches :", len(re.findall("\d", randStr)))
99 # >>>Matches : 5
100 print("Matches :", len(re.findall("\d{5}", randStr)))
101 # match a 5 digit number
102 # >>>Matches : 1
103 print("Matches :", len(re.findall("\d{2}", randStr)))
104 # >>>Matches : 2
105 # 1st : 12 ; 2nd : 34
106
107 # Matching within a range between 5 and 5 digits
108 numStr = "123 12345 123456 1234567"
109 print("Matches :", len(re.findall("\d{5,7}", numStr)))
110 # either 5, 6 or 7 digits
111 # >>>Matches : 3
112
113 # Matching any single letter or number
114 # \w : [a-zA-Z0-9_]
115 # \W : [^a-zA-Z0-9_]
116
117 phNum = "857-555-3158"
118 if re.search("\w{3}-\w{3}-\w{4}", phNum):
119     print("It is a phone number")
120 # >>>It is a phone number
121
122 # \s : [\f\n\r\t\v]
123 # will match all kinds of space
124 # \S : [^\f\n\r\t\v]
125
126 if re.search("\w{2,20}\s\w{2,20}", "Steven Wang"):
127     print("It is valid")
128

```

```

129 # + : will match one or more characters
130 print("Matches :", len(re.findall("a+", "a as ape bug")))
131 # >>>Matches : 3
132
133 # Matching an email addresss
134 email_list = "ap@apple.com, mail@gmail.com, w@b.com ce@.com"
135 print("Email Matches :", len(re.findall("[\w._%+-]{2,20}@[ \w
136 .-]{2,20}\.[A-Za-z]{2,3}", email_list)))
137
138 # summary
139 # Did you find a match
140 # if re.search("REGEX", yourString)
141
142 # Get list of matches
143 # print("Matches :", len(re.findall("REGEX", yourString)))
144
145 # Get a pattern object
146 # regex = re.compile("REGEX")
147
148 # Substitute the match
149 # yourString = regex.sub("substitution", yourString)
150
151 # [ ] : Match what is in the brackets
152 # [^ ] : Match anything not in the brackets
153 # . : Match any 1 character or space
154 # + : Match 1 or more of what proceeds
155 # \n : Newline
156 # \d : Any 1 number
157 # \D : Anything but a number
158 # \w : Same as [a-zA-Z0-9_]
159 # \W : Same as [^a-zA-Z0-9_]
160 # \s : Same as [\f\n\r\t\v]
161 # \S : Same as [^\f\n\r\t\v]
162 # {5} : Match 5 of what proceeds the curly brackets
163 # {5,7} : Match values that are between 5 and 7 in length
164
165 randStr = "cat cats"
166 regex = re.compile("[cat]+s?")
167 # "+" will match one or more trailing characters
168 matches = re.findall(regex, randStr)
169 # The question mark is called a quantifier
170 # colour?r matches both colour and color
171 # Nov(ember)? matches Nov and November
172 for i in matches:
173     print(i)
174

```

```

175 randStr = "cat cats cat's"
176 regex = re.compile("[cat]+['s]*")
177 # will match all three
178
179 # eg.
180 regex = re.compile("[\w\s]+[\r]?\\n")
181 # \s space
182 # find a line
183
184 # greedy and lazy matching
185 randStr = "<name>Steven</name><name>Wang</name>"
186 # we want to get things between the tags
187 'regex = re.compile("<name>.*</name>")
188 # However, "*" is what we called "greedy" and it search for the
    largest scope; It grab the biggest match possible
189
190 'regex = re.compile("<name>.*?</name>")
191 # add "?" to make it lazy
192 # +?
193 # {5}?
194
195 # greedy : grab the biggest match possible
196 # lazy : grab the smallest match possible
197
198 # Word boundaries
199 # \b
200 randStr = "ape at the apex"
201 # we want only the first match
202 regex = re.compile(r"\bape\b")
203
204 # String boundaries
205 # ^ : Beginning of the string
206 # $ : End of the string
207 randStr = "Match everying up to @"
208 regex = re.compile(r"^.*[^\@]")
209 matches = re.findall(regex, randStr)
210
211 randStr = "@ Get this string"
212 regex = re.compile(e"[^\s].*$")
213
214 randStr = '''Ape is big
215 Turtle is slow
216 Steven is good
217 '''
218 # multiline string
219
220 regex = re.compile(r"(?m)^.*?\s")

```

```

221 # (?m) target each line in the multiline strings
222
223 # Substrings
224 randStr = "My anumer is 888-000-1111"
225 regex = re.compile((r"888-(.*)")
226 # only get 000-1111
227
228 regex = re.compile(r"888-(.*)-(.*)")
229 # 555 1111
230
231 # Summary
232 # Did you find a match
233 # if re.search("REGEX", yourString)
234
235 # Get list of matches
236 # print("Matches :", len(re.findall("REGEX", yourString)))
237
238 # Get a pattern object
239 # regex = re.compile("REGEX")
240
241 # Substitute the match
242 # yourString = regex.sub("substitution", yourString)
243
244 # [ ]      : Match what is in the brackets
245 # [^ ]     : Match anything not in the brackets
246 # ( )      : Return surrounded submatch
247 # .        : Match any 1 character or space
248 # +        : Match 1 or more of what proceeds
249 # ?        : Match 0 or 1
250 # *        : Match 0 or More
251 # *?       : Lazy match the smallest match
252 # \b       : Word boundary
253 # ^        : Beginning of String
254 # $        : End of String
255 # \n       : Newline
256 # \d       : Any 1 number
257 # \D       : Anything but a number
258 # \w       : Same as [a-zA-Z0-9_]
259 # \W       : Same as [^a-zA-Z0-9_]
260 # \s       : Same as [\f\n\r\t\v]
261 # \S       : Same as [^\f\n\r\t\v]
262 # {5}      : Match 5 of what proceeds the curly brackets
263 # {5,7}    : Match values that are between 5 and 7 in length
264 # ($m)     : Allow ^ on multiline string
265
266 # ----- Back References -----
267 # A back reference allows you to to reuse the expression

```

```

268 # that proceeds it
269
270 randStr = "<a href='#><b>the ling</b></a>"
271
272 # eg.
273 regex = re.compile(r"<b>(.*?)</b>")
274 newStr = re.sub(regex, r"\1", randStr)
275 # "\1" back reference
276 print(newStr)
277
278 # eg 2.
279 randStr = "123-456-7890"
280 regex = re.compile(r"([\d]{3})-([\d]{3}-[\d]{4})")
281 newStr = re.sub(regex, r"(\1)\2", randStr)
282 # \1 : ([\d]{3})
283 # \2 : ([\d]{3}-[\d]{4})
284 print(randStr)
285
286 # eg 3.
287 randStr = "https://www.youtube.com http://www.google.com"
288 # goal
289 # <a href='https://www.youtube.com'>www.youtube.com</a>
290 # <a href='http://www.google.com'>www.google.com</a>
291 regex = re.compile(r"(https?://[\w.]+)")
292 # ? : "s" may or may not be there
293 newStr = re.sub(regex, r"<a href='\1'>\2</a>", randStr)
294
295 # Look ahead
296 # (?=expression)
297 randStr = "one two three four"
298 regex = re.compile(r"\w+(?=\b)")
299 # (?=\b) go not return word boundaries
300 # eg. (?=,); (?=.)
301
302 # Look behind
303 # (?<=expression)
304 randStr = "1. Bread 2. Apples 3. Lettuce"
305 regex = re.compile(r"(?<=\d.\s)\w+")
306 # (?<=\d.\s) : look for it but do not want return it
307
308 randStr"<h1>This is the title</h1> <h1>Another title</h1>"
309 regex = re.compile(r"(?<=<h1>).+?(?=</h1>)")
310 # .+? : lazy
311
312 # negative look head/behind
313 # look for text that does not match the pattern
314

```

```

315 # (?!expression) : Negative Look Ahead
316 # (?<!expression) : Negative Look Behind
317
318 randStr = "8 Apples $3, 1 Bread $1, 1 Orange $2"
319 regex = re.compile(r"(?<!\$)\d+")
320 matches = re.findall(regex, randStr)
321 print(len(matches))
322
323 #####
324 # summary
325 # [ ] : Match what is in the brackets
326 # [^ ] : Match anything not in the brackets
327 # ( ) : Return surrounded submatch
328 # . : Match any 1 character or space
329 # + : Match 1 or more of what proceeds
330 # ? : Match 0 or 1
331 # * : Match 0 or More
332 # *? : Lazy match the smallest match
333 # \b : Word boundary
334 # ^ : Beginning of String
335 # $ : End of String
336 # \n : Newline
337 # \d : Any 1 number
338 # \D : Anything but a number
339 # \w : Same as [a-zA-Z0-9_]
340 # \W : Same as [^a-zA-Z0-9_]
341 # \s : Same as [\f\n\r\t\v]
342 # \S : Same as [^\f\n\r\t\v]
343 # {5} : Match 5 of what proceeds the curly brackets
344 # {5,7} : Match values that are between 5 and 7 in length
345 # ($m) : Allow ^ on multiline string
346
347 # Use a back reference to substitute what is between the
348 # bold tags and eliminate the bold tags
349 # re.sub(r"<b>(.*?)</b>", r"\1", randStr)
350
351 # Use a look ahead to find all characters of 1 or more
352 # with a word boundary, but don't return the word
353 # boundary
354 # re.findall(r"\w+(?=\b)", randStr)
355
356 # Use a look behind to find words starting with a number,
357 # period and space, but only return the word that follows
358 # re.findall(r"(?<=\d.\s)\w+", randStr)
359
360 # Use a negative look behind to only return numbers without
361 # a $ in front of them

```



```

362 # re.findall(r"(?<!\$)\d+", randStr)
363
364 #####
365 # | : or
366 randStr = "1. Dog 2. Cat 3. Turtle"
367 regex = re.compile(r"\d\.\s(Dog|Cat)")
368 matches = re.findall(regex, randStr)
369 for i in matches:
370     print(i)
371
372 # Group
373
374 birthday = input("Enter your birthday(mm-dd-yyyy):")
375 bdRegex = re.search(r"(\d{1,2})-(\d{1,2})-(\d{4})", birthday)
376 print("You were born on", bdRegex.group())
377 print("Month", bdRegex.group(1))
378 print("Day", bdRegex.group(2))
379 print("Year", bdRegex.group(3))
380
381 # match object
382 match = re.search(r"\d{2}", "The chicken is 13 months old.")
383 print("Match : ", match.group())
384 print("Span : ", match.span())
385 print("Start : ", match.start())
386 print("End : ", match.end())
387
388 # Name groups
389 randStr = "December 21 1999"
390 regex = r"^(?P<month>\w+)\s(?P<day>\d+)\s(?P<year>\d+)"
391 print("Span : ", match.group('month'))
392 print("Start : ", match.group('day'))
393 print("End : ", match.group('year'))

```

python regex