🔒 **sandeepsuryaprasad** / **python_tutorials**   Private

<> Code   ⊙ Issues   ⭠⭡ Pull requests   ▶ Actions   ⊞ Projects   📖 Wiki   ⊘ Securit

⑂ master ⌄          **python_tutorials** / 13_regex /          Go to file          ···

_regular_expressions.py / <> Jump to ⌄

Sandeep Suryaprasad cleanup          Latest commit b4042cb 5 days ago          🕒 History

👥 **0** contributors

605 lines (495 sloc)  │  24.7 KB          Raw    Blame    ✏    ⌄    ⧉    🗑

```python
 1   import re
 2   # ============ Characters ===================
 3   # . - Matches any character except new line
 4   # \. - Mathes a dot.
 5   # \\ - Matches backslash
 6   # \* - Matches astrick
 7   # ============ Character set ==================
 8   # [abcd] - any character which matches either 'a' or 'b' or 'c' or 'd'
 9   # [^abcd] - any character but not 'a' or 'b' or 'c' or 'd'
10   # [a-z] - any character between 'a' through 'z'
11   # =========== Special Sequences ===============
12   # \w - Word character. Same as [a-zA-Z0-9_]. Matches alphanumeric and undersc
13   # \W - Non-Word Character. Same as [^a-zA-Z0-9_]. Matches anything but word c
14   # \d - Matches a digit. Same as [0-9]
15   # \D - Matches a Non-Digit. Same as [^0-9]
16   # \s - Matches only whitespace.
17   # \S - Matches only Non-Whitespace.
18   # ========== Anchors =====================
19   # ^ - Start of String
20   # $ - End of String
21   # \b - Word boudary [a-zA-Z0-9_]
22   # \B - Not a word Boundry
23   # [ ] - Matches characters in square brackets
24   # [^ ] - Matches characters Not in square brackets
25
26   # Meta Characters that needs to be Escaped
27   # . ^ $ * + ? { } [ ] \ | ( )
```

```python
28
29   # Quantifiers
30   """
31   1. The standard quantifiers (?, +, *, and {min,max}) are greedy.
32   2. When one of these governs a subexpression, such as a?, (expr )+, [0-9]+,
33      there is a minimum number of matches that are required before it can be co
34      attempt it will ever attempt to match,
35   3. They always attempt to match as many times as they can, up to that maximum
36   4. The only time they settle for anything less than their maximum allowed is
37      ends up causing some later part of the regex to fail.
38   5. plus, question mark and star are called quantifiers, because they influenc
39   """
40   # * — Match expression 0 or more times
41   # + — Match expression 1 or more times
42   # ? — Match expression 0 or 1 times
43   # {min,max} — Matches expression exactly 3 times
44
45   # =========== Grouping ====================
46   # ("A"| "B" | "C") — Either "A" or "B" or "C"
47
48   # re.findall()   # returns a list of all the matches
49   # re.sub()  # replaces one pattern with other
50   # re.finditer()  # returns an iterator object
51   # re.search()  # stops at the first match
52
53   # Word Boundary (\b)
54   # "start of word" boudary is simply the position where a sequence of alphanum
55   # "end of word" is the position where a sequence of alphanumeric characters e
56   # ----------------------------------------------------------------------------
57   # Rule: The Match That Begins Earliest (from left to right) Wins
58   # ----------------------------------------------------------------------------
59   re.findall(r"the", "the theory of relativity")
60
61   re.findall(r"cat", "The dragging belly indicates your cat is too fat")
62
63   re.findall(r'python', 'python and java are object oriented')
64
65   re.findall(r'aeiou', 'hello how are you doing anna')
66
67   re.findall(r'aeiou', 'hello how are you doing anna, aeiou')
68   # ----------------------------------------------------------------------------
69   # Character class or set
70   # ----------------------------------------------------------------------------
71   # Matches with both "Smith" and "smith"
72   re.findall(r'[sS]mith', 'smith')
```

```python
 73    re.findall(r'[sS]mith', 'Smith')
 74
 75    # Matches separate or saperate
 76    re.findall(r's[ea]p[ae]rate', 'seperate')
 77    re.findall(r's[ea]p[ae]rate', 'saparate')
 78
 79    # Match any one character in the character set (either a, e, i, o, u)
 80    re.findall(r'[aeiou]', 'hello how are you doing anna')
 81
 82    # Match either a, b, c, d
 83    re.findall(r'[abcd]', 'hello world')
 84    re.findall(r'[abcd]', 'abcdefghijk')
 85
 86    # Matching any number between 0-9
 87    re.findall(r'[0123456789]', 'The cost of the book is Rs.100')
 88
 89    # Matching HTML headers
 90    re.findall(r'<h[123456]>', "<h1>")
 91    re.findall(r'<h[123456]>', "<h2>")
 92    re.findall(r'<h[123456]>', "<h3>")
 93    re.findall(r'<h[123456]>', "<h4>")
 94    re.findall(r'<h[123456]>', "<h5>")
 95    re.findall(r'<h[123456]>', "<h6>")
 96    # ----------------------------------------------------------------------
 97    # Range "-"
 98    # ----------------------------------------------------------------------
 99    # Matches any number between 0-9
100    re.findall(r'[0-9]', 'The cost of the book is Rs.100')
101
102    # Matches only lower case letters
103    re.findall(r"[a-z]", 'hello HOW ARE YOU')
104
105    # Matches only upper case letters
106    re.findall(r"[A-Z]", 'hello HOW ARE YOU')
107
108    # Matches all upper case and lower case characters
109    re.findall(r"[a-zA-Z]", 'hello HOW ARE YOU')
110
111    # Matches any number between 1-6
112    re.findall(r"<h[1-6]>", "<h1>")
113    re.findall(r"<h[1-6]>", "<h2>")
114    re.findall(r"<h[1-6]>", "<h3>")
115    re.findall(r"<h[1-6]>", "<h4>")
116    re.findall(r"<h[1-6]>", "<h5>")
117    re.findall(r"<h[1-6]>", "<h6>")
```

```python
118
119    # Count total number of Upper case and Lower case letters
120    sentence = "Hello World Welcome To Python"
121    upper_case = re.findall(r'[A-Z]', sentence)
122    lower_case = re.findall(r'[a-z]', sentence)
123
124    print(f'Total No of upper case letters {len(upper_case)}')
125    print(f'Total No of lower case letters {len(lower_case)}')
126
127    # Write a program to count the number of white spaces in a given string
128    sentence = "Hello world welcome to Python Hi  How are you. Hi how are you"
129    spaces = re.findall(r' ', sentence)
130
131    # Write a program to count the number of occurrences of each lower case and u
132    sentence = 'hello@world! welcome!!! Python$ hi how are you & where are you?'
133    chrs = re.findall(r'[a-zA-Z]', sentence)
134    d = {chr: chrs.count(chr) for chr in chrs}
135
136    # ----------------------------------------------------------------------
137    # Meta Character "+" (matches 1 or more occurances of previous expression)
138    # ----------------------------------------------------------------------
139    re.findall(r'[0-9]+', 'The cost of the book is Rs.100')
140
141    re.findall(r'[abcd]+', 'abcdefg hijkab')
142
143    re.findall(r'an+a', 'annnnnnnnnnna')
144
145    # Matches each word in the string
146    re.findall(r"[a-zA-Z]+", "Hello World Welcome To Python Programming Pyt123on"
147
148    # Count the characters in each word. Please ignore special characters if ther
149    sentence = "Hi there! How are you:) How are you doing today!"
150    words = re.findall(r'[a-zA-Z]+', sentence)
151    word_len = { word: len(word) for word in words}
152
153    # Sum all the numbers in the below string.
154    word = "Pytho12nReg567exp2" # 1 + 2 + 5 + 6 + 7 + 2
155    total = 0
156    numbers = re.findall(r'[0-9]', word)
157    for number in numbers:
158        total += int(number)
159
160    # Adding 12 + 567 + 2
161    word = "Pytho12nReg567exp2"
162    total = 0
```

```python
163    numbers = re.findall(r'[0-9]+', word)
164    for number in numbers:
165        total += int(number)
166
167    # Match file names and extensions
168    message = "Downloading file archive.zip to downloads folder..."
169    # image.jpeg
170    # index.xhtml
171    # python.py
172    re.findall(r'[a-z]+\.[a-z]+', message)
173    # -----------------------------------------------------------------------
174    # Meta Character "?" (matches 0 or 1 occurance of previous expression)
175    # -----------------------------------------------------------------------
176    re.findall(r'colou?r', "what color do you like")
177
178    re.findall(r'https?://', 'https://www.google.com')
179
180    re.findall(r'https?://', 'http://www.google.com')
181
182    re.findall(r'July?', "Jul the 26th day")
183
184    re.findall(r'an?a', "ana")
185
186    re.findall(r'an?a', "anna")
187    # -----------------------------------------------------------------------
188    # Meta Character "*" (matches 0 or more occurances of previous expression)
189    # -----------------------------------------------------------------------
190    re.findall(r"an*a", "hello ana")
191
192    re.findall(r"an*a", "hello aa")
193
194    re.findall(r"an*a", "hello annna")
195
196    # Regular Expression for Matching Inbox, Inbox(1), .... Inbox(N)
197    re.findall(r"Inbox\(?\d*\)?", "Inbox(10)")
198    re.findall(r"Inbox\(?\d*\)?", "Inbox")
199    # -----------------------------------------------------------------------
200    # Negation "^"
201    # -----------------------------------------------------------------------
202    # Matches everything apart from numbers between 0-9
203    re.findall(r'[^0-9]', 'The cost of the book is Rs.100')
204
205    # Matches everything apart from alphabets 'a', 'b', 'c' and 'd'
206    re.findall(r'[^abcd]', 'abcdefg hijkab')
207
```

```python
208    # Matches everything apart from numbers between 0-9
209    re.findall(r'[^0-9]+', 'The cost of the book is Rs.100')
210
211    re.findall(r'[^abcd]+', 'abcdefg hijkab')
212
213    # Match only those characters excepts digits
214    word = '@hello12world34welcome!123'
215    re.findall(r'[^0-9]', word)
216
217    # Count the number of special characters in the below string
218    sentence = 'hello@world! welcome!!! Python$ hi26 how are you & where are you?
219    re.findall(r"[^a-zA-Z0-9_\s]", sentence)
220    # ----------------------------------------------------------------
221    # Starts with "^" and ends with "$"
222    # ----------------------------------------------------------------
223    re.findall(r"^hello", "hello world")
224
225    re.findall(r"^hello", "world hello")
226
227    re.findall(r"hello$", "world hello")
228
229    re.findall(r"hello$", "hello world")
230
231    re.findall(r'hello$', 'hello world welcome to python')
232
233    # Matching the only those lines which ends with "UDP"
234    with open("./data_files/sample.log") as f:
235        for line in f:
236            match = re.findall(r".*UDP$", line)
237            if match:
238                print("".join(match))
239
240    # string starts with "hello" and ends with "hello" (meaning exactly one word
241    re.findall(r"^hello$", "hello")
242
243    # Phone Number pattern (4DIGITS-3DIGITS-3DIGITS)
244    re.findall(r'\d{3}-\d{3}-\d{4}', '456-9832-098')
245
246    # matching only 800 and 900 numbers
247    re.findall(r"^[89]00-\d{3}-\d{4}", '800-123-123')
248    # ----------------------------------------------------------------
249    # Word Boundary (\b) The expression should be a word boundry
250    # (Transition between non-word character to word character or word character
251    # ----------------------------------------------------------------
252    # starts with word boundry
```

```python
253    re.findall(r"\bday", "what a beautiful day today is")
254
255    # ends with word boundry
256    re.findall(r"day\b", "what a beautiful day today is")
257
258    # starts and ends with word boundry
259    re.findall(r"\bday\b", "what a beautiful day today is")
260
261    re.findall(r"\b[0-9]{6}\b", 'Pincode of Bangalore is 560001 and the number is
262
263    # Regular expression which matches words that starts with "h"
264    re.findall(r"\bh[a-zA-Z0-9_]+", 'hello world hi hello universe how are you ha
265
266    # Regular expression which matches words that starts with "P or J"
267    re.findall(r"\b[PJ][a-zA-Z0-9_]+", 'Python is a programming language. Python
268
269    # Regular expression which matches words that ends with "y"
270    re.findall(r"[a-zA-Z0-9_]+y\b", 'hello world hi hello universe how are you ha
271
272    # print only those words starting with vowel character
273    sentence = "hello hi american engieers and indian writers officers united sta
274    words = re.findall(r"\b[aeiou][a-zA-Z0-9_]+", sentence)
275
276    # Matches only Capital Letter words
277    re.findall(r"\b[A-Z]+\b", "This is PYTHON programming LANGUAGE and REGEX")
278
279    # Matches only lower case words
280    re.findall(r"\b[a-z]+\b", "This is PYTHON programming LANGUAGE and REGEX")
281
282    # Matching only pdf files
283    re.findall(r"[a-zA-Z0-9]+\.pdf\b", "downloading apple.pdf to downloads folder
284
285    # Regular expression for matching only 3 letter words in the given string
286    sentence = "hello hi how are you what is your name he is older than me how ol
287    re.findall(r'\b[a-zA-Z0-9_]{3}\b', sentence)
288    # o/p ['how', 'are', 'you', 'how', 'old', 'are', 'you']
289
290    # Extract only 4 digit numbers from the string
291    re.findall(r"\b\d{4}\b", "Copyright 1998. All rights reserved")
292
293    # Regular expression for matching the words which starts with "he"
294    sentence = "he helps the community and he is the hero of the day"
295    re.findall(r"\bhe[a-zA-Z0-9_]*", sentence)
296
297    # Regular expression for matching the words which either starts with "he" or
```

```python
298    sentence = "he helps the community and he is the hero of the day she sells se
299    re.findall(r"\b(?:he|se)[a-zA-Z0-9_]*", sentence)
300
301    # Regular Expression - PAN Number
302    sentence = "my pan number is ABCDE1234X and the other number is XYZTR3104J id
303    re.findall(r'\b[A-Z]{5}[0-9]{4}[A-Z]\b', sentence)
304
305    # Different Combintations
306    line = "03/22 08:51:06 WARNING :.....mailslot_create: setsockopt(MCAST_ADD) f
307    re.findall(r"[A-Z]+", line)
308    re.findall(r"\b[A-Z]+", line)
309    re.findall(r"[A-Z]+\b", line)
310    re.findall(r"\b[A-Z]+\b", line)
311
312    # Matches all digits
313    re.findall(r"\d", "654 this string is starting with 12 and ending with number
314
315    # Matches whole numbers
316    re.findall(r"\d+", "654 this string is starting with 12 and ending with numbe
317
318    # Matches only 3 Digit numbers
319    re.findall(r"\d{3}", "654 this string is starting with 12 and ending with num
320
321    # Matches all digit numbers
322    re.findall(r"\b\d{3}\b", "654 this string is starting with 12 and ending with
323
324    # Matches the string that ends with 3 digit number
325    re.findall(r"\b\d{3}\b$", "4632746327 this string is ending with 235")
326
327    # Matches the string that starts with 3 digit number
328    re.findall(r"^\b\d{3}\b", "654 this string is starting with and ending with n
329    # ------------------------------------------------------------------------
330    # Count the number of white spaces in the file
331    def count_spaces():
332        with open('./data_files/sample.log') as f:
333            white_spaces = 0
334            for line in f:
335                count = len(re.findall(r"\s", line))
336                white_spaces += count
337        return white_spaces
338    # ------------------------------------------------------------------------
339    # Count the number of Capital Letter words in the file
340    def count_caps():
341        with open('./data_files/sample.log') as f:
342            capital_words = 0
```

```python
343            for line in f:
344                count = len(re.findall(r"\b[A-Z]+\b", line))
345                capital_words += count
346        return capital_words
347    # ------------------------------------------------------------------------
348    # Count the number of Capital Letters in the file
349    def count_cap_letters():
350        with open('./data_files/sample.log') as f:
351            capital_letters = 0
352            for line in f:
353                count = len(re.findall(r"[A-Z]", line))
354                capital_letters += count
355        return capital_letters
356    # ------------------------------------------------------------------------
357    # Count the number of INFO, TRACE, WARNING, EVENT messages in the file
358    def count_messages(message_name):
359        with open('./data_files/sample.log') as f:
360            message_count = 0
361            for line in f:
362                count = len(re.findall(message_name, line))
363                message_count += count
364        return message_count
365    # ------------------------------------------------------------------------
366    def get_all_ip():
367        """Extract all the ip addresses in the sample log file"""
368        ips = [ ]
369        with open('./data_files/sample.log') as f:
370            for line in f:
371                ip = re.findall(r'\d{1,3}\.\d{1,3}\.\d{1,3}\.\d{1,3}', line)
372                if ip:
373                    for item in ip:
374                        ips.append(item)
375        return ips
376    # ------------------------------------------------------------------------
377    def count_words():
378        """Counts the number of words in the file"""
379        word_count = 0
380        with open('./data_files/sample.log') as f:
381            for line in f:
382                words = re.findall(r"\b[A-Za-z]+\b", line)
383                word_count += len(words)
384        return word_count
385    # ------------------------------------------------------------------------
386    def count_total_letters():
387        """Counts number of lowecase and upper case characters in the file"""
```

```python
388        with open('./data_files/sample.log') as f:
389            total_letters = 0
390            for line in f:
391                count = len(re.findall(r"[a-zA-Z]", line))
392                total_letters += count
393        return total_letters
394    # --------------------------------------------------------------------
395    def get_log_lines():
396        """Prints only those log lines which are logged in the year 98"""
397        with open("../data_files/sample.log") as f:
398            for line in f:
399                matches = findall(r"(?:0[1-9]|1[0-2])/98", line)
400                if matches:
401                    print(line, end="")
402    # --------------------------------------------------------------------
403    def regex_strip(line):
404        """Strips leading and trailing whitespaces and returns a new string"""
405        stripped_line = sub(r"(^\s*|\s*$)", "", line)
406        return stripped_line
407    # --------------------------------------------------------------------
408    def count_leading_trailing_spaces(line):
409        """Returns a tuple of leading and trailing whitespaces"""
410        leading_spaces = findall(r"\s+$", line)
411        trailing_spaces = findall(r"^\s+",line)
412        # join the list
413        count_leading_spaces = len("".join(leading_spaces))
414        count_trailing_spaces = len("".join(trailing_spaces))
415        return (count_leading_spaces, count_trailing_spaces)
416    # --------------------------------------------------------------------
417    # Flitering only floating point values from file
418    with open("./data_files/points.txt") as f:
419        floats = [ ]
420        for line in f:
421            match = re.findall(r"-?[0-9]+\.[0-9]+", line)
422            # match = re.findall(r"-?[0-9]+\.[0-9]{3}", line) # matches 3 digits
423            for item in match:
424                floats.append(item)
425    # --------------------------------------------------------------------
426    # Groups ( )
427    # --------------------------------------------------------------------
428    # Matches either "python" or "java"
429    re.findall('(python|java)', 'python and java are object oriented')
430
431    # Matches 09:00 am/pm or 9:00 am/pm
432    sentence = 'The meeting is between 9:00 am and 12:30 pm'
```

```python
433     re.findall(r'[0-9]?[0-9]:[0-9][0-9]\s(?:am|pm)', sentence)
434
435     # Regular Expression - YYYY-MM-DD date format
436     _dates = ['2019-01-02', '2019-13-02', '2019-12-26', '26-08-2019', '20-19-20',
437     re.findall(r'\d{4}-(?:0[1-9]|[12][0-9]|3[01])-(?:0[1-9]|[12][0-9]|3[01])', '2
438
439     # Regular Expression - 24hr time format
440     _formats = ['00:00:00', '23:59:59', '24:00:00', '1:59:20', '12:9:10', '10:20:
441     re.findall(r"(?:[01]\d|2[0-3]):[0-5]\d:[0-5]\d",'23:59:59')
442     # ------------------------------------------------------------------------
443     # Replacing patterns
444     # ------------------------------------------------------------------------
445     # Replace whitespaces with newline character in the below string
446     sentence = "Hello world welcome to python"
447     words = re.sub(r'\s', '\n', sentence)
448     print(words)
449
450     # Replace all vowels with "*"
451     sentence = "hello world welcome to python"
452     words = re.sub(r'[aeiou]', '*', sentence)
453     print(words)
454
455     # Replace all occurances of digits with "*"
456     sentence = 'hello123world welcome456to python012'
457     words = re.sub(r'\d', '*', sentence)
458
459     # Replace all occurances of special characters with "*"
460     sentence = 'hello#$%world welcome@!#$%to python*&^%'
461     words = re.sub(r'[^a-zA-Z\s]', "*", sentence)
462
463     # Replace "And" with "&"
464     sentence = "Java and Python are programming languages"
465     _sentence = re.sub(r"\sAnd\s", " & ", sentence)
466
467     # Replace all occurances of "Java" with "Python" in a file
468     with open('java.txt', 'r') as jf:
469         with open('python.txt', 'a') as pf:
470             for line in jf:
471                 new_line = re.sub('Java', 'Python', line)
472                 pf.write(new_line)
473     # ------------------------------------------------------------------------
474     # re.finditer()
475     # ------------------------------------------------------------------------
476     matches = re.finditer(r"hello", 'hello world welcome to python hello world')
477     for match in matches:
```

```python
478            print(match.group())
479
480    # Write a program to find the index of nth occurrence of a sub-string in a st
481    sentence = "hello world welcome to python hello hi how are you hello there"
482    matches = re.finditer(r'hello', sentence)
483    positions = [ (match.start(), match.end()) for match in matches]
484
485    # -----------------------------------------------------------------------
486    # re.search()
487    # -----------------------------------------------------------------------
488    match = re.search(r"hello", 'hello world hello world')
489    print(match.group())
490
491    # -----------------------------------------------------------------------
492    # dot "." matches with everything
493    # -----------------------------------------------------------------------
494    re.findall(r'.', "hello world")
495
496    re.findall(r'h.', "hello")
497
498    re.findall(r'h.', "hello world hi how how are you")
499
500    re.findall(r'a.b', "acb")
501
502    re.findall(r'a.b', "a b")
503
504    re.findall(r'a.b', "ab")
505
506    re.findall(r'a.b', "a*b a?b")
507
508    re.findall(r'a.b', "abcde")
509
510    re.findall(r'a.a', "ana")
511
512    re.findall(r'a..a', "anna")
513
514    re.findall(r'a.*a', "annnnnna")
515
516    re.findall(r'a.*a', 'aa')
517
518    re.findall(r"^a.*a$", "anna")
519
520    re.findall(r"^a.*a$", "hello anna")
521
522    re.findall(r'a.*a', 'abcad')
```

```python
523
524   re.findall(r'a.*a$', 'abcad')
525
526   re.findall(r'a.*a$', 'abcada')
527
528   re.findall(r'a.+a', 'ana')
529
530   re.findall(r'a.+a', 'aa')
531   # ----------------------------------------------------------------------------
532   # ------------------------------- Back-referencing ---------------------------
533   """
534   1. Back-referencing is a regular-expression feature which allows you to match
535   matched earlier in the expression without specifically knowing the text when
536   2. Back-references provide a convenient way to identify a repeated character
537   For example, if the input string contains multiple occurrences of an arbitrar
538   first occurrence with a capturing group, and then use a backreference to matc
539   3. \1 will try to match what ever is matched in the first bracket
540   """
541   # Repeated word sequences
542   m = re.findall(r"(world)\1", "thethe python hello worldworld the")
543
544   # Repeated words
545   m = re.findall(r"([a-z]+\s)\1", "the the python hello world world the")
546
547   # Repeated words for 2 consecutive times
548   m = re.finditer(r"([a-z]+\s)\1{2}", "the the python hello world world the the
549
550   # Repeated characters
551   m = re.findall(r"([a-z])\1", "hello hurry programming")
552
553   # Repeated numbers
554   n = re.findall(r"([0-9])\1", "hello 123345, welcome to 001 98799")
555
556   # Repeated numbers pattern
557   n = re.findall(r"([0-9]+\s)\1", "hello 12345 12345 , welcome to 001 98799")
558
559   n = re.findall(r"([0-9])+", "hello 1234512345 , welcome to 00198799")
560
561   # finding the words that are repeated at the beginning and at the end of the
562   sentence = "hello world welcome to regex hello"
563   re.findall(r"^([a-z]+).*\1$", sentence)
564   # ----------------------------------------------------------------------------
565   # Regular Expression - IP Addresses
566   # ----------------------------------------------------------------------------
567   id_address_format = '\d{1,3}\.\d{1,3}\.\d{1,3}\.\d{1,3}'
```

```python
568    ips = ['10.1.2.3', '127.0.0.0', '199.99.9.9', '199.9.9999.9', '127-0-0-0']
569    # ----------------------------------------------------------------------
570    # Regular Expression - Email format
571    # ----------------------------------------------------------------------
572    email_pattern = r'[\w-]+\.?[\w-]@[\w]+\.(com|edu|in|gov)'
573    emails = ['test.user@company.com',
574              'test.user2@company.com',
575              'test_user@company.com',
576              'testing@company.com',
577              'test-T.user@company.com',
578              'testing@company',
579              'testingcompany.com'
580              ]
581    # ----------------------------------------------------------------------
582    # Regular Expression - URL Pattern
583    url_pattern = r'https?://[\w.]+'
584    urls = ['http://www.youtube.com',
585            'https://www.google.com',
586            'http://www.amazon.in',
587            'https://www.mail.yahoo.com',
588            'ftp://test.com'
589            'https://www.facebook.com/'
590            ]
591    # ----------------------------------------------------------------------
592
593    # ----------------------------------------------------------------------
594    # ------------------------------- LookAhead and LookBehind Anchors ------
595           # Under development ...
596    # ----------------------------------------------------------------------
597    # Misc
598    # ----------------------------------------------------------------------
599    # Finding the all the links in a html document using regex
600    from requests import request
601    response = request("GET", "http://demowebshop.tricentis.com/")
602    html_code = response.text
603    links = re.findall(r"\s*<a\s.+", html_code)
604    links = [ link.strip() for link in links ]
605    # ----------------------------------------------------------------------
```