

## Regular Expressions

A regular expression is a special sequence of characters that helps you match or find other strings or sets of strings, using a specialized syntax held in a pattern. Regular expressions are widely used in UNIX world.

The module `re` provides full support for Perl-like regular expressions in Python. The `re` module raises the exception `re.error` if an error occurs while compiling or using a regular expression.

### Match, Search, Findall, Substitute

`re.match(pattern, string, flags=0)` This function attempts to match RE `pattern` to `string` with optional `flags`.

`re.search(pattern, string, flags=0)` Scan through `string` looking for the first location where the regular expression `pattern` produces a match, and return a corresponding *match object*. Return `None` if no position in the string matches the pattern; note that this is different from finding a zero-length match at some point in the string.

`re.findall(pattern, string, flags=0)` Return all non-overlapping matches of `pattern` in `string`, as a list of strings. The `string` is scanned left-to-right, and matches are returned in the order found.

`re.sub(pattern, repl, string, count=0, flags=0)` Return the string obtained by replacing the leftmost non-overlapping occurrences of `pattern` in `string` by the replacement `repl`. If the pattern isn't found, `string` is returned unchanged. `repl` can be a string or a function.

### Exercise 1

```
import re
line = "My Dogs are smarter than your Cats"
results = re.findall('[dc]\S+s', line, re.M|re.I)
if results:
    print ("results: ", results)
else:
    print ("Nothing found")
```

### Exercise 2

```
import re
phone = "(631) 615-4215 # School Nova phone number"
num = re.sub('#.*$', '', phone)
print ("phone num: ", num)
num = re.sub('\D', '', phone)
print ("phone num: ", num)
```

Anchors	Quantifiers	Groups and Ranges
<code>^</code> Start of string <code>\A</code> Start of string <code>\$</code> End of string <code>\Z</code> End of string <code>\b</code> Word boundary <code>\B</code> Not word boundary <code>\&lt;</code> Start of word <code>\&gt;</code> End of word	<code>*</code> 0 or more <code>+</code> 1 or more <code>?</code> 0 or 1 <code>{3}</code> Exactly 3 <code>{3,}</code> 3 or more <code>{3,5}</code> 3, 4 or 5	<code>.</code> Any character except new line ( <code>\n</code> ) <code>(a b)</code> a or b <code>(...)</code> Group <code>(?:...)</code> Passive Group <code>[abc]</code> Range (a or b or c) <code>[^abc]</code> Not a or b or c <code>[a-q]</code> Letter between a and q <code>[A-Q]</code> Upper case letter between A and Q <code>[0-7]</code> Digit between 0 and 7 <code>\n</code> nth group/subpattern <i>Note: Ranges are inclusive.</i>
Character Classes	Quantifier Modifiers	Pattern Modifiers
<code>\c</code> Control character <code>\s</code> White space <code>\S</code> Not white space <code>\d</code> Digit <code>\D</code> Not digit <code>\w</code> Word <code>\W</code> Not word <code>\x</code> Hexadecimal digit <code>\O</code> Octal digit	<i>"x" below represents a quantifier</i> <code>x?</code> Ungreedy version of "x"	<code>g</code> Global match <code>i</code> Case-insensitive <code>m</code> Multiple lines <code>s</code> Treat string as single line <code>x</code> Allow comments and white space in pattern <code>e</code> Evaluate replacement <code>U</code> Ungreedy pattern
POSIX	Escape Character	Metacharacters (must be escaped)
<code>[[:upper:]]</code> Upper case letters <code>[[:lower:]]</code> Lower case letters <code>[[:alpha:]]</code> All letters <code>[[:alnum:]]</code> Digits and letters <code>[[:digit:]]</code> Digits <code>[[:xdigit:]]</code> Hexadecimal digits <code>[[:punct:]]</code> Punctuation <code>[[:blank:]]</code> Space and tab <code>[[:space:]]</code> Blank characters <code>[[:cntrl:]]</code> Control characters <code>[[:graph:]]</code> Printed characters <code>[[:print:]]</code> Printed characters and spaces <code>[[:word:]]</code> Digits, letters and underscore	<code>\</code> Escape Character	<code>^</code> [ . * <code>\$</code> { + <code>(</code> \   ? <code>)</code>   >
Assertions	Special Characters	String Replacement (Backreferences)
<code>?=</code> Lookahead assertion <code>?!</code> Negative lookahead <code>?&lt;=</code> Lookbehind assertion <code>?!= or ?&lt;!</code> Negative lookbehind <code>?&gt;</code> Once-only Subexpression <code>?()</code> Condition [if then] <code>?()</code> Condition [if then else] <code>?#</code> Comment	<code>\n</code> New line <code>\r</code> Carriage return <code>\t</code> Tab <code>\v</code> Vertical tab <code>\f</code> Form feed <code>\xxx</code> Octal character xxx <code>\xhh</code> Hex character hh	<code>\$n</code> nth non-passive group <code>\$2</code> "xyz" in <code>/^(abc(xyz))\$/</code> <code>\$1</code> "xyz" in <code>/^(?:abc)(xyz)\$/</code> <code>\$`</code> Before matched string <code>\$'</code> After matched string <code>\$+</code> Last matched string <code>\$&amp;</code> Entire matched string
Sample Patterns		
<i>Pattern</i> <code>{[A-Za-z0-9-]+}</code> <code>(\d{1,2}\V\d{1,2}\V\d{4})</code> <code>([^\s]+(?:=.(\.jpg gif png))\.\d{2})</code> <code>(^[1-9]{1}\$ ^1-4{1}[0-9]{1}\$ ^50\$)</code> <code>(#?([A-Fa-f0-9]){3}([A-Fa-f0-9]){3})?</code> <code>((?=.*)d)(?=.*[a-z])(?=.*[A-Z]).{8,15}</code>	<i>Will Match</i> Letters, numbers and hyphens Date (e.g. 21/3/2006) jpg, gif or png image Any number from 1 to 50 inclusive Valid hexadecimal colour code String with at least one upper case letter, one lower case letter, and one digit (useful for passwords).	
<code>(\w+@[a-zA-Z_]+?\.[a-zA-Z]{2,6})</code> <code>(\&lt;/?[\^&gt;]+\&gt;)</code>	Email addresses HTML Tags	
<i>Note: These patterns are intended for reference purposes and have not been extensively tested. Please use with caution and test thoroughly before use.</i>		

## Homework

Search through your hard drive for email addresses using the sample program below. Adjust the program to fit your environment (root directory, file extension etc.):

```
import os, fnmatch, re, sys
found = set()

pattern = re.compile(r'\b[A-Z0-9._%+-]+@[A-Z0-9.-]+\.[A-Z]{2,4}\b',re.I)

for root, dir, files in os.walk("/Users/serge/projects"):
    #print (root)

    for file in fnmatch.filter(files, "*.txt"):
        file = os.path.join(root, file)
        #print(file)
        if os.path.isfile(file) and os.access(file, os.R_OK):

            try:

                for line in open(file, 'r'):

                    found.update(pattern.findall(line))

            except UnicodeDecodeError:

                print ("Program could not open file ", file)

for myMatch in found:

    print (myMatch)
```

Additional Documentation can be found at <https://docs.python.org/3/library/re.html>