# PYTHON LAB BOOK

## Python For Programmers
## *UCSC Extension Online*

### Lab 19  re Module

Topics

- re - Regular Expressions (Optional)

- Search and replace (Optional)

- Named groups (Optional)

v.4.0

lab18_1.py

```
 1 #!/usr/bin/env python
 2 """lab18_1.py
 3 Make an UpIt(str) function that returns the input string, but
 4 with all caps.  Your UpIt function will be different from
 5 str.upper() in that, if any of the characters in the input string
 6 are already uppercase, it raises an exception.  Invent your own
 7 exception and put it in a reasonable spot in the exception
 8 hierarchy.
 9 """
10 import sys
11
12 class CaseError(ValueError):
13     pass
14
15 def UpIt(s):
16     if not s.islower():
17         raise CaseError
18     return s.upper()
19
20 if __name__ == '__main__':
21     try:
22         print UpIt(sys.argv[1])
23     except CaseError:
24         print sys.argv[1], 'has at least one upper case letter.'
25     except IndexError:
26         print 'Usage:', sys.argv[0], 'string_of_lower_case_letters'
27
28 """
29 $ lab18_1.py HI
30 HI has at least one upper case letter.
31 $ lab18_1.py hi
32 HI
33 $ lab18_1.py Hi
34 Hi has at least one upper case letter.
35 $
36 """
```

```
lab18_2.py
 1 #!/usr/bin/env python
 2 """lab18_2.py -- collects 2 sides of a right triangle and prints
 3 the hypotenuse."""
 4
 5 import math
 6
 7 def GetXY(prompt):
 8     """Returns a tuple of 2 floats, or None if the user just hits the
 9     enter key or enters q."""
10
11     while True:
12         try:
13             incoming = raw_input(prompt)
14         except (KeyboardInterrupt, EOFError):
15             return None
16         if incoming == '' or incoming[0].lower() == 'q':
17             return None
18         try:
19             x, y = [float(x) for x in incoming.split(',')]
20         except (TypeError, NameError, ValueError, SyntaxError):
21             print "Two numbers are required."
22         else:
23             return x,y
24
25 def Hypot(x, y):
26     """Returns sqrt(x**2 + y**2)"""
27     return math.sqrt(x * x + y * y)
28
29 def main():
30     while True:
31         xy = GetXY("Give me 2 sides of a right triangle: (x, y) ")
32         if not xy:
33             print
34             break
35         answer = Hypot(*xy)
36         if answer < 0:
37             continue
38         print 'Hypotenuse is %.2f' % (answer)
39
40 if __name__ == '__main__':
41     main()
42
43 """
44 $ lab18_2.py
```

```
45 Give me 2 sides of a right triangle: (x, y) 4, 5
46 Hypotenuse is 6.40
47 Give me 2 sides of a right triangle: (x, y) a
48 Two numbers are required.
49 Give me 2 sides of a right triangle: (x, y) 1
50 Two numbers are required.
51 Give me 2 sides of a right triangle: (x, y) 1,a
52 Two numbers are required.
53 Give me 2 sides of a right triangle: (x, y) 1,2,3
54 Two numbers are required.
55 Give me 2 sides of a right triangle: (x, y) 2,,3
56 Two numbers are required.
57 Give me 2 sides of a right triangle: (x, y) 1,2
58 Hypoteneuse is 2.24
59 Give me 2 sides of a right triangle: (x, y)
60 $ lab18_2.py
61 Give me 2 sides of a right triangle: (x, y)
62 $ lab18_2.py
63 Give me 2 sides of a right triangle: (x, y) Q
64 $"""
```

Answers for Lab 18_3

1.
```
>>> a_list = ['loop']
>>> a_list += a_list
>>> a_list
['loop', 'loop']
>>> a_list.append(a_list)
>>> a_list
['loop', 'loop', [...]]
```

So a_list has itself in it, that a_list inside has a_list inside, has a_list inside, ... Python uses **...** to indicate this infinite pattern.

I can't explain the difference between append and +=.

2.
```
>>> bottles = 100
>>> def HowMany():
...     print bottles
...
>>> HowMany()
100
>>> def HowMany():
...     bottles -= 1
...     print bottles
...
>>> HowMany()
Traceback (most recent call last):
  File "<stdin>", line 1, in ?
  File "<stdin>", line 2, in HowMany
UnboundLocalError: local variable 'bottles' referenced before
assignment
```

While interpreting, Python sees the bottles -= 1 as an assignment and decides that bottles is a local name and sets it up that way. But at run-time, the interpreter tries to access bottles and doesn't find the local variable.

How can you fix this so that the HowMany's bottles is really the global bottles?

```
>>> def HowMany():
...     global bottles
...     bottles -= 1
...     print bottles
...
>>> HowMany()
99
```

3.
```
>>> def Snake(rattle=[]):
...     rattle += ["hiss"]
...     print rattle
...
>>> Snake([100])
[100, 'hiss']
```

No surprises there but:

```
>>> Snake()
['hiss']
```

As expected.

```
>>> Snake()
['hiss', 'hiss']
>>> Snake()
['hiss', 'hiss', 'hiss']
```

Woah, since it used the mutable default in the last call, and then it altered it, it altered that default!

4.
```
>>> class X:
...     pass
...
>>> class Y:
...     pass
...
>>> X.a = 1
>>> X.b = 2
>>> X.c = 3
>>> Y.a = X.a + X.b + X.c
>>> for X.i in range(Y.a):
...     print X.i
...
0
1
2
3
4
5
>>> dir (X)
['__doc__', '__module__', 'a', 'b', 'c', 'i']
```

Handling Regular Expressions in Python (Labs 19-20 are optional)

```
>>> import re

>>> mo = re.search("needle", "Is there a needle in this haystack?")
>>> print mo                              mo == None if there's no match
<_sre.SRE_Match object at 0xb7ea5f00>

>>> mo.span()
(11, 17)
>>> mo.start()
11
>>> mo.end()
17

>>> re.findall("needle", "Is there a needle in this haystack of needles?")
['needle', 'needle']

>>> re.sub("needle", "thread", "Is there a needle in this haystack of needles?")
'Is there a thread in this haystack of threads?'

>>> re.subn("needle", "thread", "Is there a needle in this haystack of needles?")
('Is there a thread in this haystack of threads?', 2)

>>> c = re.compile("needle")

>>> c.findall("Is there a needle in this haystack of needles?")
['needle', 'needle']
```

Raw Strings!

```
"(\d+\.\d*|\d*\.\d+|\d+)"
```

is a regular expression that matches numbers.

All those '\' need escaping:

```
"(\\d+\\.\\d*|\\d*\\.\\d+|\\d+)" TOO HARD!
```

```
r"(\d+\.\d*|\d*\.\d+|\d+)" is the "raw string" notation.
```

```
>>> c = re.compile(r"(\d+\.\d*|\d*\.\d+|\d+)")
>>> c.findall("3.2, a .4 and a 5. x.x .y z. . .")
['3.2', '.4', '5.']
```

However, your pattern might contain special re characters that you are looking for as literals.

```
>>> re.findall("cat?", "Where's the cat?  Is the cat in the car?")
['cat', 'cat', 'ca']
```

It didn't find that question mark but thought we were looking for 0 or 1 ts, so it just matched the *cat*. In this case, you need to put the escape character before the '?'. `re.escape` does it for you:

```
>>> re.findall(re.escape("cat?"), "Where's the cat?  Is the cat in the car?")
['cat?']
```

```
>>> import re

Setting up some capturing groups:

>>> p = re.compile('(a(b)c)d')    >>> p = re.compile(r'\b((\w+) ?(\w?\.?)) Jones')
>>> m = p.search('abcd')          >>> m = p.search("John J. Jones")
>>> m.group()                     >>> m.group()
'abcd'                            'John J. Jones'
>>> m.groups()                    >>> m.groups()
('abc', 'b')                      ('John J.', 'John', 'J.')
>>> m.group(1)                    >>> m.group(2)
'abc'                             'John'
                                  >>> p.findall("John J. Jones, Mary Jones")
                                  [('John J.', 'John', 'J.'), ('Mary', 'Mary', '')]
Names Groups:

>>> r'\b((?P<first>\w+) ?(?P<middle>\w?\.?)) Jones')
>>> m = p.search("John J. Jones")
>>> m.groups(2)
'John'
>>> m.group('first')
'John'
>>> m.group('middle')
'J.'


Works nicely with sub:

>>> p = re.compile('section{(?P<which>[^}]*)}')
>>> p.sub(r'subsection{\g<which>}', 'section{First}')
'subsection{First}'
>>> p.sub(r'subsection{\1}', 'section{First}')
'subsection{First}'
>>> p.sub(r'subsection{\g<1>}', 'section{First}')
'subsection{First}'
```

With p.sub or re.sub, instead of a replacement string, you can give a
function name!

re_swap.py
```
 1 #!/usr/bin/env python
 2 """Swapping cats and dogs again, but this time doing a better job with
 3 regular expressions.  The old, faulty solution is:
 4
 5 def DoSwap(text, apple, orange):
 6     """Swap apple for orange and orange for apple in the text.
 7
 8     Return the swapped text.
 9     """
10
11     dummy = 'wxyz'
12     while True:
13         if text.find(dummy) == -1:
14             break
15         dummy *= 2
16     text = text.replace(apple, dummy)
17     text = text.replace(orange, apple)
18     text = text.replace(dummy, orange)
19     return text
20
21 """
22
23 import re
24
25 def DoSwap(text, apple, orange):
26     """Swaps all occurances of apple for orange, and orange
27     for apple in the text."""
28     # Using VERBOSE and named groups for readability
29     compiled_re = re.compile(r"""
30     \b                  # matches a word boundary
31     (?P<found>%s|%s) # matches apple or orange and puts the match
32                         # in a group named "found" if the whole thing matches
33     (?P<rest>s?)\b   # matches a word boundary or 's' and a word boundary
34                         # and puts the 's', or not, into a group named "rest"
35     """ % (re.escape(apple), re.escape(orange)), re.IGNORECASE | re.VERBOSE)
36
37     def MatchCase(answer, like_string):
38         if like_string.isupper():
39             return answer.upper()
40         if like_string.islower():
41             return answer.lower()
42         if like_string.istitle():
43             return answer.title()
44         return answer
```

```
45
46     def SwapMatch(match_object):
47         found = match_object.group('found')
48         if found.lower() == apple.lower():
49             x = MatchCase(orange, found) + match_object.group('rest')
50             return x
51         x = MatchCase(apple, found) + match_object.group('rest')
52         return x
53
54     fixed = compiled_re.sub(SwapMatch, text)
55     return fixed
56
57 def main():
58     print DoSwap("DOGS: 14dogs are a lot of Dogs.", 'dog', 'elephant')
59     print DoSwap("3 Categories of CATS.", 'cat','dog')
60
61 if __name__ == '__main__':
62     main()
63 """
64 $ re_swap_caseless.py
65 ELEPHANTS: 14dogs are a lot of Elephants.
66 3 Categories of DOGS.
67 $
68 """
```

```
quiz.py
  1 #!/usr/bin/env python
  2 """Quiz answer."""
  3
  4 import sys
  5 sys.path.insert(0, '..')
  6 import pet_store.pet_def
  7
  8 class Rock(pet_store.pet_def.Pet):
  9     no_rocks = 0
 10
 11     def __init__(self, name, weight):
 12         super(Rock, self).__init__(name)
 13         pet_store.pet_def.Pet.__init__(self, name)
 14         Rock.no_rocks += 1
 15         self.__weight = weight
 16
 17 def main():
 18     rocky = Rock("Rocky", "3 oz")
 19     gravel = Rock("Gravel", "1/2 oz")
 20     diamond = Rock("Diamond", ".03 oz")
 21     print "Made %d pet rocks." % Rock.no_rocks
 22     print rocky.__weight
 23
 24 if __name__ == '__main__':
 25     main()
 26 """
 27 $ quiz.py
 28 Made 3 pet rocks.
 29 Traceback (most recent call last):
 30   File "./quiz.py", line 23, in <module>
 31     main()
 32   File "./quiz.py", line 20, in main
 33     print rocky.__weight
 34 AttributeError: 'Rock' object has no attribute '__weight'
 35 $ tree.py mall
 36 |-- /mall
 37 |   |-- /basket
 38 |   |   |-- quiz.py
 39 |   |-- /pet_store
 40 |   |   |-- pet_def.py
 41 |   |   |-- __init__.py
 42 $ """
```

Lab 19_re_Module (Optional)  Check out `http://regexlib.com` to find some helpful regular expressions and other tools.

1. Find `labs/lab_11_Packages/numbers.txt` again.  Use this regular expression to add up all the numbers:

   `"\d+\.\d*|\d*\.\d+|\d+"`

2. Say you have a dictionary of search and replace words:

   ```
   replace_d = {
   'orchestra leader':'singer',
   'china':'mexico',
   'zen':'mariachi',
   'master':'teacher',
   'sword':'baton',
   'through':'around'}
   ```

   Find file `labs/lab_19_re_Module/zen.story` and change all those words. Copy the dictionary from `zen.dictionary` in the same directory.

3. A big optional project! A few potentially useful, but imperfect, regular expressions:
   For email addresses:

   `r"[\s,:;(<\['\"]+([^@\s,;:\"]+)@([^@\-\s,:;>\)'\"]+\.[a-z,A-Z]{2,3})[\s,;:\"'>\)]"`

   For finding URL's:

   `r"http://[^\s\"<>()']*\b"`

   Copy and paste these from `labs/lab_19_re_Module/lab.data` Or, if you know better ones, use them, and please tell me.
   With the urllib module, you can:

   `stream = urllib.urlopen("http://www.pythontrainer.com")`

   and then handle `stream` as if it was an open file.

   If there is an error opening the page, the exception raised is `httplib.InvalidURL` so you need to `import httplib` if you want to catch such an exception. If the page asks for a password, it'll hang. So, you might want to catch that by:

   `import labs/lab_12_Function_Fancies.time_out_decorator as time_out`

Then decorate over any function that you expect will hang:

```
@time_out.TimeOut(3) # This make the following function
                     # timeout in 3 seconds.
def CollectPage(url):
    whatever that should do

try:
    CollectPage(some_url)
except RuntimeError:
    whatever that should do
```

Write an evil program that takes a starting url ('www.ngc.com' behaves pretty well), finds all the email addresses on that page, and all the links to urls on that page, and then reads all the links, and their links, etc., collecting email addresses.