

# 2015 Computational Social Science Workshop

## Day 1 - Introduction to python - Part 2 / 3 (cont'd.)

with **Jongbin Jung** (jongbin at stanford.edu)

- PhD student. Decision analysis, MS&E

All material for days 1 (intro to python) and 2 (web scraping with python) publicly available [here \(https://github.com/jongbinjung/css-python-workshop\)](https://github.com/jongbinjung/css-python-workshop)

## 2. python development (with spyder) - 2 of 3

Now that we're quite familiar with some of python's basic concepts, let's touch some 'advanced' topics.

### Writing functions

Let's create a function to count the number of vowels in a given string

```
In [2]: def count_vowels(s):  
        """Count the number of vowels in a string."""  
        vowels = 'aeiouAEIOU'  
        nvowels = [s.count(v) for v in vowels] # count the number of each vowel in s  
        return sum(nvowels) # return the sum of elements in nvowel  
  
        # use the new function  
        count_vowels('Eels are delicious animals')
```

Out[2]: 12

- the **def** keyword declares a function **definition**, followed by a function name and the parenthesized list of formal parameters
- the statements that form the body of the function start at the next line, and must be indented
- the first statement of the function body can optionally be a string, also known as the **docstring** (<https://docs.python.org/2/tutorial/controlflow.html#tut-docstrings>)
- many tools (such as `spyder`) use the docstring to give users meaningful information - so help yourself, make a habit of writing meaningful docstrings
- functions that don't finish with a `return` statement return `None` (a special python object)

Functions can also return more a tuple of values. For example, let's modify our `count_vowels` function to return the number of vowel along with the list specifying the number of *each* vowel.

```
In [3]: def count_vowels(s):
        """
        Count the number of vowels in a string.

        returns: number of vowels, list containing number of appearance for each vowel
        """
        vowels = 'aeiouAEIOU'
        nvowels = [s.count(v) for v in vowels] # count the number of each vowel in s
        return sum(nvowels), list(zip(vowels, nvowels)) # return the sum and a zipped
        list

count_vowels('Eels are delicious animals')
```

```
Out[3]: (12,
        [('a', 3),
         ('e', 3),
         ('i', 3),
         ('o', 1),
         ('u', 1),
         ('A', 0),
         ('E', 1),
         ('I', 0),
         ('O', 0),
         ('U', 0)])
```

A returned tuple can also be 'unpacked' into multiple variables.

```
In [4]: total_count, individual_count = count_vowels('Eels are delicious animals')
        print 'Found total', total_count, 'vowels, each vowel as follows:'
        print individual_count

Found total 12 vowels, each vowel as follows:
[('a', 3), ('e', 3), ('i', 3), ('o', 1), ('u', 1), ('A', 0), ('E', 1), ('I', 0),
 ('O', 0), ('U', 0)]
```

## Functions with optional arguments

Let's further enhance the `count_vowels` function by letting the user specify

- which vowels to count ('aeiouAEIOU' by default)
- whether to return a single sum or a tuple of the sum and list (single sum by default)

This can be achieved by specifying default values in the function declaration.

```
In [5]: def count_vowels(s, vowels = 'aeiouAEIOU', returnAll = False):
        """
        Count the number of vowels in a string.

        returns: number of vowels, list containing number of appearance for each vowel
        """
        nvowels = [s.count(v) for v in vowels] # count the number of each vowel in s
        if returnAll:
            return sum(nvowels), list(zip(vowels, nvowels)) # return the sum and a zip
            ped list
        else:
            return sum(nvowels) # return just the sum

count_vowels('Eels are delicious animals')
```

```
Out[5]: 12
```

```
In [6]: count_vowels('Eels are delicious animals', vowels = 'aeiou') # no caps
```

```
Out[6]: 11
```

```
In [7]: count_vowels('Eels are delicious animals', returnAll = True) # give me EVERYTHING
```

```
Out[7]: (12,
        [('a', 3),
         ('e', 3),
         ('i', 3),
         ('o', 1),
         ('u', 1),
         ('A', 0),
         ('E', 1),
         ('I', 0),
         ('O', 0),
         ('U', 0)])
```

Be careful with having mutable defaults, though. Default values of a function's argument are shared between subsequent calls, and this might cause problems if you're manipulating the argument's value within the function. For example,

```
In [8]: def fun(n, stuff=[]):
        """Issues with mutable defaults."""
        stuff.append(n)
        return stuff

        print fun(1) # stuff is empty by default
        print fun(2) # stuff was manipulated, and is now [1] from the previous call!
        print fun(3) # even worse, stuff is now [1, 2] !!!
```

```
[1]
[1, 2]
[1, 2, 3]
```

This behavior isn't necessarily a problem, and it might even make sense in some contexts. However, it's definitely worth keeping in mind to avoid being surprised. If you want to prevent such behavior, one simple work-around is to set the default to `None`, and check if it is indeed `None`, before assigning the 'true' default, such as:

```
In [9]: def fun(n, stuff=None):
        """Fix for mutable defaults."""
        if stuff is None:
            stuff = []
        stuff.append(n)
        return stuff

        print fun(1) # unspecified argument stuff is None, then set to []
        print fun(2) # unspecified argument stuff is None, then set to []
        print fun(3) # unspecified argument stuff is None, then set to []
        print fun(3, [1,2]) # and we can always specify stuff if we need to!
```

```
[1]
[2]
[3]
[1, 2, 3]
```

Finally, to capture an arbitrary number of arguments in a function, you can use the `*name` and `**name` parameters. Note that, if both are present, `*name` **must** occur before `**name`, and both must occur after all the formal parameters. When present, the `*name` parameter receives a tuple containing the positional arguments beyond the formal parameter list, and `**name` receives a dictionary containing the key-value pair of the named arguments, except for those corresponding to a formal parameter. For example:

```
In [10]: def fun(n, name='Jongbin', *arguments, **keywords):
        """Demo of *name and **name parameters."""
        print '\n' + '=' * 79
        print 'Function called with n=', n, # note that a leading ',' in the print sta
        tement will avoid new lines
        print 'Name=', name
        print 'Arguments received:'
        print '\t', # a tab character to print appropriate indents
        for arg in arguments:
            print arg, '|',
        print '\nNamed arguments received:'
        print '\t', # a tab character to print appropriate indents
        for key, value in keywords.iteritems():
            print key, '=', value, '|',

fun(1) # supply minimal arguments
fun(2, 'Padme', 'Amidala', 'Princess', 'testing additional arguments') # some addi
tional arguments
fun(2, 'Luke', gender='male', affiliation='Rebel Alliance', text='testing named arg
uments') # named arguments
fun(3, 'Anakin', 'Skywalker', 'Jedi', 2015, weapon='Lightsaber', skill='force') #
both

=====
Function called with n= 1 Name= Jongbin
Arguments received:

Named arguments received:

=====
Function called with n= 2 Name= Padme
Arguments received:
    Amidala | Princess | testing additional arguments |
Named arguments received:

=====
Function called with n= 2 Name= Luke
Arguments received:

Named arguments received:
    gender = male | text = testing named arguments | affiliation = Rebel Alli
    ance |
=====
Function called with n= 3 Name= Anakin
Arguments received:
    Skywalker | Jedi | 2015 |
Named arguments received:
    weapon = Lightsaber | skill = force |
```

Sometimes, an opposite situation may occur, where the required arguments are in a list/tuple or keyword arguments are in a dictionary, and you would like to unpack them programatically in the function call. In such cases, you can use the `*name` and `**name` conventions introduced above in the function call. For example:

```
In [11]: print 'regular call:', range(1, 10, 2) # the range function takes arguments (start
        , stop[, step])
        args = [1, 10, 2] # pack the arguments (equivalent to above) into a list
        print 'unpack from list:', range(*args) # all the function by unpacking the list

regular call: [1, 3, 5, 7, 9]
unpack from list: [1, 3, 5, 7, 9]
```

```
In [12]: def print_info(name, email, phone):
         """Quick demo of keyword argument unpacking."""
         print 'Name:', name
         print 'email:', email
         print 'phone:', phone

         kwargs = {'name': 'Jongbin Jung', 'email': 'jongbin at stanford.edu', 'phone': '650-123-4567'}
         print_info(**kwargs)

Name: Jongbin Jung
email: jongbin at stanford.edu
phone: 650-123-4567
```

## Lambda expressions

Anonymous one-liner functions can be created with the `lambda` keyword wherever function objects are required, but you don't want or need to define a full function.

```
In [13]: # sort the list of words by the number of vowels in each word,
         # but without defining a separate count_vowel function
         words.sort(key=lambda word: sum([word.count(v) for v in 'aeiouAEIOU']))
         print words

-----
NameError                                Traceback (most recent call last)
<ipython-input-13-09e3e0c6e5fc> in <module>()
      1 # sort the list of words by the number of vowels in each word,
      2 # but without defining a separate count_vowel function
----> 3 words.sort(key=lambda word: sum([word.count(v) for v in 'aeiouAEIOU']))
      4 print words

NameError: name 'words' is not defined
```

## Modules

Once you start building functions, you might want to collect certain functions as a general 'toolbox' to be used across multiple projects. In python, you can put definitions in a file with a `.py` extension. Such a file is called a module. Once you save your functions into a module, you can import them. Let's practice with some examples.

For illustration purposes, create let's create two modules that contain one function of the same name each:

```
In [14]: # save this function to a file named module1.py
         def speak():
             """Make module 1 say something"""
             print 'Module 1 speaking ...'
```

```
In [15]: # save this function to a file named module2.py
         def speak():
             """Make module 2 say something"""
             print 'Hi, this is module 2 speaking!'
```

You can import each module (and the functions in them) using the `import` statement as follows:

```
In [16]: import module1
         import module2
```

Note that the name you use in the `import` statement is just the file name of the module, without the `.py` extension.

When you `import` a module, python creates an isolated 'space' for each module. This allows different modules to have functions of the same name, without causing confusion. But because of this, when ever you want to use a function from a certain module, you have to specify the module name before calling the function. Compare:

```
In [17]: module1.speak()  
         module2.speak()  
  
Module 1 speaking ...  
Hi, this is module 2 speaking!
```

This can be a bit painful (and messy) if your module names get longer. There are typically two ways to work around this:

1. `import` with the keyword `as` to assign your own name to a model
2. assign your own function name to a module's function

Each approach is illustrated below, which to use should depend on the context and personal style:

```
In [18]: import module1 as m1  
         import module2 as m2  
         m1.speak()  
         m2.speak()  
  
Module 1 speaking ...  
Hi, this is module 2 speaking!
```

```
In [19]: import module1  
         import module2  
  
         speak1 = module1.speak  # note the lack of parentheses  
         speak2 = module2.speak  # when assigning functions to a new name  
  
         speak1()  
         speak2()  
  
Module 1 speaking ...  
Hi, this is module 2 speaking!
```

Finally, modules can also be executed as standalone scripts. However, to do this, the module must know when it's been imported or executed. This is done in python by specifying a `__name__` variable within each module's 'space'. When a module is imported, it's `__name__` variable is set to the filename it was imported from:

```
In [20]: module1.__name__  
  
Out[20]: 'module1'
```

However, if a module is executed, for example from the terminal with the command,

```
python module_name.py
```

then the `__name__` variable is set to `__main__`.

To illustrate this, let's create a new module, `module3.py`:

```
In [ ]: # save this code to a file named module3.py
def speak():
    """Make module 3 say something"""
    print 'My __name__ is', __name__

    if __name__ == '__main__':
        speak()
        print 'You\'ve executed me!'
```

```
In [21]: import module3

module3.speak()

My __name__ is module3
```

Now, instead of importing module3, execute it from a command prompt with the command:

```
python module3.py
```

(you can open a command prompt within spyder)

The output should look like:

```
My __name__ is __main__
You've executed me!
```

## A little more on module execution ...

When executing a module from the command prompt, you can also pass arguments to the module in the form of

```
python filename.py arguments
```

The arguments are passed to the module via a list in the `sys` standard module, and can be accessed by calling `sys.argv`. (*standard* modules are modules that are built-in to python). The first (position 0) element of `sys.argv` contains the execution call of the module, so arguments that are passed through the command prompt start from position 1.

For example, we can write a module that takes a single argument from the command prompt as follows:

```
In [ ]: # save this code to a file named module4.py
import sys # import the standard module sys

if __name__ == '__main__':
    print 'The first element of sys.argv is', sys.argv[0]
    print 'The argument passed was:', sys.argv[1]
```

Then, execute from the command prompt with an argument, for example:

```
python module4.py hello
```

This should print to the screen:

```
The first element of sys.argv is module4.py  
The argument passed was: hello
```

Note that all arguments are passed as a string by default. If you want to use a different type, you will have to convert it within python, e.g., `int(sys.argv[1])` to convert the first argument into an integer).

(if you want to do some serious argument parsing, you should take a look at the [argparse module \(https://docs.python.org/2/howto/argparse.html\)](https://docs.python.org/2/howto/argparse.html))

### Exercise 3.

Continuing from the previous exercise ...

1. Write a function `top_n(d, n=5)`, which takes a dictionary of word counts (such as that created in the previous exercise) and an optional argument `n`, and prints words that have the top `n` count, along with the actual count.

- to sort a dictionary by its values, use the built-in function `sorted(iterable, cmp=None, key=None, reverse=False)`; you can set the sorting `key` to the dictionary's value by setting `key=d.get`, and sort in descending order by setting `reverse=True`.
- you might want to create a simple word count dictionary to test your function

2. Expand your function from the previous question into a module that can be executed with a target filename and integer `n` as an argument, i.e.,

```
$ python module_name.py target_file.txt 5
```

which,

- A. reads the contents of the target file
- B. generates a word occurrence count dictionary from the text
- C. prints word/count of the words with top `n` occurrences