

functions

Genome 559: Introduction to Statistical
and Computational Genomics

Prof. James H. Thomas

Take a deep breath and think how much you've learned.

4 weeks ago, this would have been gibberish:

```
import sys
matrixFile = open(sys.argv[1], "r")
matrix = []
line = matrixFile.readline().strip()
while len(line) > 0:
    fields = line.split("\t")
    intList = []
    for field in fields:
        intList.append(int(field))
    matrix.append(intList)
    line = matrixFile.readline().strip()
matrixFile.close()
for row in matrix:
    for val in row:
        print val,
    print ""
```

```
# initialize empty matrix
# read first line
# until end of file
# split line on tabs, giving a list of strings
# create an int list to fill
# for each field in current line
# append the int value of field to intList
# after intList is filled, append it to matrix
# read next line and repeat loop
# go through the matrix row by row
# go through each value in the row
# print each value without line break
# add a line break after each row is done
```

Problem Set 3 - code clarity

Pick names that make sense:

`count` for counting something

`index` for an index in a list or string

`xFileName` for a file name

`xFile` for a file

Use a counter in a loop only if you need one:

```
for line in lineList:  
    line.do-something
```

rather than

```
for index in range(0, len(lineList)):  
    line[index].do-something
```

Once you have a program that is bug-free and works, look for ways to make it simpler:

Instead of:

```
myText = myFile.read()
myList = myText.split("\n")
finalText = myList[0]
for index in range(1:len(myList)):
    finalText += myList[index]
```

(Hmm, all this does is replace new lines in original text...)

How about:

```
myText = myFile.read()
finalText = myText.replace("\n", "")
```

```
(or just finalText = myFile.read().replace("\n", ""))
```

These loops "work" to find every instance of foo in seq, but what is wrong with them?

```
query = "foo"  
for index in len(seq):  
    pos = seq.find(foo, index)
```

or

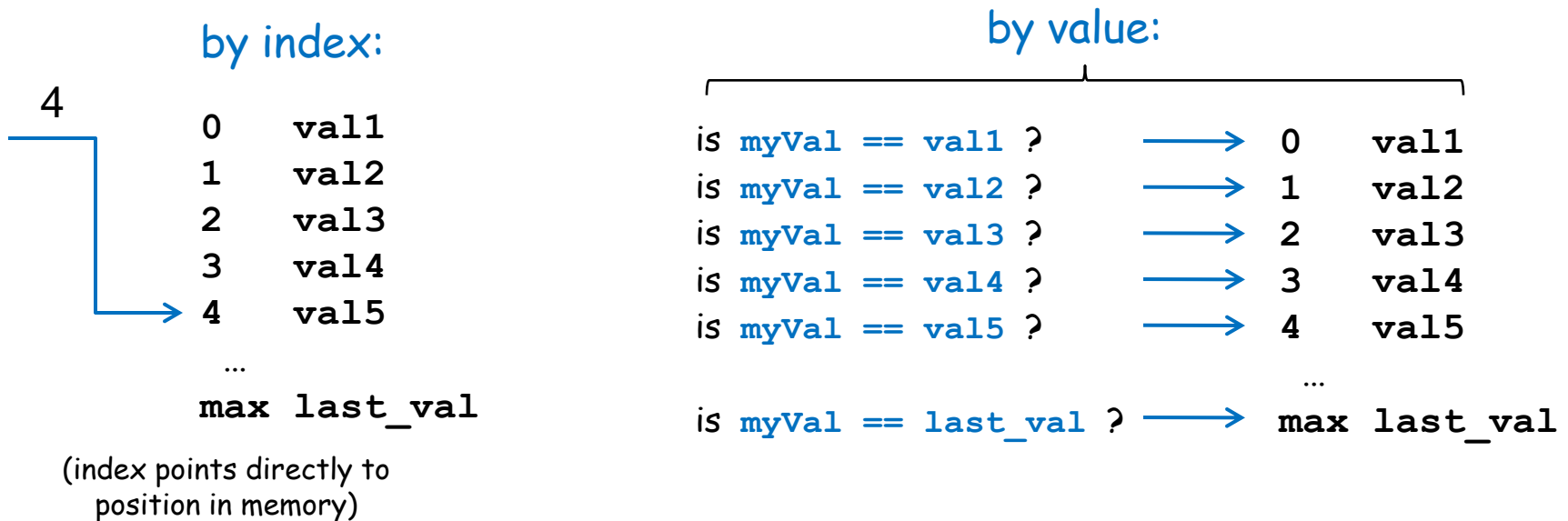
```
query = "foo"  
while pos >= 0:  
    pos = seq.find(foo)  
    seq = seq(1:)
```

Review

- Start paying attention to program robustness and speed.
- During program development, use `print` liberally to see intermediate values. (then remove them)
- Dictionaries - key : value pairs.
- Dictionaries are useful when you want to look up some data (value) based on a key.

Dictionary and List access times

- Accessing a list by index is very fast.
- Accessing a dictionary by key is very fast.
- Accessing a list by value (e.g. `list.index(myVal)` or `list.count(myVal)`) can be **SLOW**.

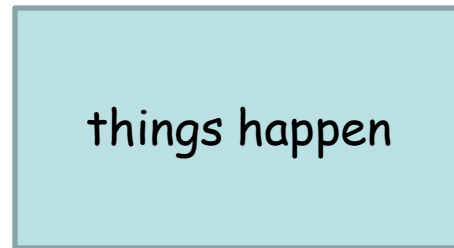


What is a function?

- Reusable piece of code
 - write once, use many times
- Takes defined inputs (**arguments**) and may produce (**return**) a defined output
- Helps simplify and organize your program
- Helps avoid duplication of code

What a function does

stuff goes in (arguments)



other stuff comes out (return)

Other than the arguments and the return, everything else inside the function is invisible outside the function (variables assigned, etc.).

The function doesn't need to have a return - if it does something to one of the arguments, this may be visible outside the function (for example: if the argument is a list, the function could sort the list).

What is a function?

```
import math
```

```
def jc_dist(rawdist):  
    if rawdist < 0.75 and rawdist > 0.0:  
        newdist = (-3.0/4.0) * math.log(1.0 - (4.0/3.0)* rawdist)  
        return newdist  
    elif rawdist >= 0.75:  
        return 1000.0  
    else:  
        return 0.0
```

define the function and
argument(s) names

do something

return a computed
value

```
def <function_name>(<arguments>):  
    <function code block>  
    <usually return something>
```

Using a function

<function defined here>

```
import sys
rawDist = sys.argv[1]
correctedDist = jc_dist(rawDist)
```

Building a function

Jukes-Cantor distance correction written directly in program:

```
import sys
import math

rawdist = float(sys.argv[1])
if rawdist < 0.75 and rawdist > 0.0:
    newdist = (-3.0/4.0) * math.log(1.0 - (4.0/3.0)* rawdist)
    print newdist
elif rawdist >= 0.75:
    print 1000.0
else:
    print 0.0
```

Building a function - step 1

add a function
definition



```
import sys
import math

def jc_dist(rawdist):
    rawdist = float(sys.argv[1])
    if rawdist < 0.75 and rawdist > 0.0:
        newdist = (-3.0/4.0) * math.log(1.0 - (4.0/3.0)* rawdist)
        print newdist
    elif rawdist >= 0.75:
        print 1000.0
    else:
        print 0.0
```

Building a function - step 2

```
import sys
import math
```

```
def jc_dist(rawdist):
    rawdist = float(sys.argv[1])
    if rawdist < 0.75 and rawdist > 0.0:
        newdist = (-3.0/4.0) * math.log(1.0 - (4.0/3.0)* rawdist)
        print newdist
    elif rawdist >= 0.75:
        print 1000.0
    else:
        print 0.0
```

add a function
definition

delete line - use function
argument instead of argv

Building a function - step 3

```
import sys
import math
```

```
def jc_dist(rawdist):
    if rawdist < 0.75 and rawdist > 0.0:
        newdist = (-3.0/4.0) * math.log(1.0 - (4.0/3.0)* rawdist)
        return newdist
    elif rawdist >= 0.75:
        return 1000.0
    else:
        return 0.0
```

add a function
definition

deleted line - use function
argument instead of argv

return value rather
than printing it

Use the function

```
raw = 0.23  
corrected = jc_dist(raw)  
print corrected
```

Once you've written the function, you can forget about it and just use it!

We've used lots of functions before

```
log()  
readline(), readlines(), read()  
sort()  
split(), replace(), lower()
```

These functions are part of the Python programming environment (in other words they are already written for you).

Note - some of these are functions attached to objects (called object "methods") rather than stand-alone functions. We'll cover this soon.

Function names and access

Giving a function an informative name is very important!
Long names are fine if needed.

```
def makeDictFromTwoLists(keyList, valueList):  
def translateDNA(dna_seq):  
def getFastaSequences(fileName):
```

- For now, your function will have to be defined within your program and before you use it.
- Later you'll learn how to save a function in a module so that you can load your module and use the function just the way we do for Python modules.
- Usually, potentially reusable parts of your code should be written as functions.
- Your program (outside of functions) will often be very short - largely reading arguments and making output.

Sample problem #1

Below is part of the program from a sample problem last class. It reads key - value pairs from a tab-delimited file and makes them into a dictionary. Rewrite it so that there is a function called `makeDict` that takes a file name as an argument and returns the dictionary.

Use:

```
scoreDict = makeDict(myFileName)
```

```
import sys
myFile = open(sys.argv[1], "r")
# make an empty dictionary
scoreDict = {}
for line in myFile:
    fields = line.strip().split("\t")
    # record each value with name as key
    scoreDict[fields[0]] = float(fields[1])
myFile.close()
```

Here's what the file contents look like:

```
seq00036<tab>784
seq57157<tab>523
seq58039<tab>517
seq67160<tab>641
seq76732<tab>44
seq83199<tab>440
seq92309<tab>446
etc.
```

Solution #1

```
import sys
```

```
def makeDict(fileName):  
    myFile = open(fileName, "r")  
    myDict = {}  
    for line in myFile:  
        fields = line.strip().split("\t")  
        myDict[fields[0]] = float(fields[1])  
    myFile.close()  
    return myDict
```

```
myFileName = sys.argv[1]  
scoreDict = makeDict(myFileName)
```

name used
inside function

name used to
call function

Two things to notice here:

- you can use any file name (string) when you call the function
- you can assign any name to the function return

(in programming jargon, the function lives in its own namespace)

Sample problem #2

Write a function that mimics the `<file>.readlines()` method. Your function will have a `file` object as the argument and will return a `list` of strings (in exactly the format of `readlines()`). Use your new function in a program that reads the contents of a file and prints it to the screen.

You can use other file methods within your function - just don't use the `<file>.readlines()` method directly.

This isn't a useful function, since Python developers already did it for you, but the point is that the functions you write are just like the ones we've already been using. BTW you will learn how to attach functions to objects a bit later (things like the `split()` function of strings, as in `myString.split()`).

Solution #2

```
import sys

def readlines(file):
    text = file.read()
    tempLines = text.split("\n")
    lines = []
    for tempLine in tempLines:
        lines.append(tempLine + "\n")
    return lines

myFile = open(sys.argv[1], "r")
lines = readlines(myFile)
for line in lines:
    print line.strip()
```

Challenge problem

Write a program that reads a file containing a tab-delimited matrix of pairwise distances and puts them into a 2-dimensional list of distances (floats). Have the program accept two additional arguments, which are the names of 2 sequences from the matrix, and print their distance.

Make the matrix reading a function.

Be sure it works with ANY matrix file with this format! (the file will always be a square matrix of size $N+1 \times N+1$ (N for each distance and 1 row and column for names)).

Here's what the file contents look like:

```
names<tab>seq1<tab>seq2<tab>seq3
seq1<tab>0<tab>0.1<tab>0.2
seq2<tab>0.1<tab>0<tab>0.3
seq3<tab>0.2<tab>0.3<tab>0
```

```
>python dist.py matrixFile seq2 seq3
0.3
```

Hints - use the first line to make a dictionary of names to list indices; your function should return a 2-dimensional list of floats.

Challenge solution

```
import sys

def makeMatrix(fileName):
    myFile = open(fileName, "r")
    myMatrix = []
    lines = myFile.readlines()
    for rowIndex in range(1, len(lines)):
        fields = lines[rowIndex].strip().split("\t")
        matRow = []
        for colIndex in range(1, len(fields)):
            matRow.append(float(fields[colIndex]))
        myMatrix.append(matRow)
    myFile.close()
    return myMatrix

def makeNameMap(line):
    nameMap = {}
    fields = line.strip().split("\t")
    for index in range(1, len(fields)):
        nameMap[fields[index]] = index - 1
    return nameMap

distMatrix = makeMatrix(sys.argv[1])
nameMap = makeNameMap(open(sys.argv[1], "r").readline())
print distMatrix[nameMap[sys.argv[2]]][nameMap[sys.argv[3]]]
```

I wrote both complex parts as functions; this makes the point that once these are written and debugged, the program is simple and easy to read (the last three lines).

(this could be done more efficiently - this way you open the file twice)

looks up the argument string as the key in nameMap, which returns the index of the name in the 2-dimensional list of distance values