

MIT OpenCourseWare
<http://ocw.mit.edu>

6.047 / 6.878 Computational Biology: Genomes, Networks, Evolution
Fall 2008

For information about citing these materials or our Terms of Use, visit: <http://ocw.mit.edu/terms>.

6.047/6.878 Fall 2008 Recitation 1 Notes

Pouya Kheradpour, Matt Rasmussen

September 5, 2008

1 Python

Python is popular a programming language that is frequently used in computational biology. Its main features are simple syntax, dynamic typing, and a large number of supporting libraries. Python is available within all MIT Athena accounts and is also available for download at <http://python.org> for all platforms (Windows, Linux, Mac OS X).

There are several tutorials and documentation sites for Python:

- Official tutorial, <http://docs.python.org/tut/tut.html>
- Library reference, <http://docs.python.org/lib/lib.html>
- WikiBooks, <http://en.wikibooks.org/wiki/Programming:Python>

We won't be using any complex or exotic features of Python, so it is probably not necessary to buy a Python book just for this course. If you would like one for the future, however, *Learning Python* and *Programming Python* by Mark Lutz are both excellent.

1.1 Brief summary of Python commands to learn

Here is a brief tour of basic Python language features, including print, variables, functions, lists, list comprehensions, loops, tuples, dictionaries, import, dir, and help. You should type these commands into the interactive interpreter to learn what each command does.

```
1 # hello world
2 print 'hello, world!'
3 print "hello, world!"
4
5 # functions & variables
6 def fact(n):
7     if n == 0 or n == 1:
8         return 1
9     else:
10        return n*fact(n-1)
11
12 print fact(8)
13
14 x = fact(8)
15
16
17 # print with formatting syntax
18 # In the C programming language this would be:
19 # printf("x = %d", x);
20 print 'x = %d' % x
21
22 # lists & loops
23 lst = [1, 2, 3, 4]
24 print lst
25 print lst[2]
26 lst[2] = 0
```

```
27 lst
28 lst.append(5)
29 lst
30 del lst[0]
31 lst
32
33 lst = range(1,5)
34 lst
35
36 print len(lst)
37
38 for i in range(1,5):
39     print i
40
41 # list comprehension (advanced feature)
42 print [x*x for x in [1,2,3,4]]
43
44 # equivalent to
45 lst = []
46 for x in [1,2,3,4]:
47     lst.append(x*x)
48 print lst
49
50 # filtering with list comprehension (advanced feature)
51 print [x*x for x in range(1,11) if x % 2 == 0]
52
53 # equivalent to
54 lst = []
55 for x in range(1,11):
56     if x % 2 == 0: # only append even numbers, x / 2 has remainder 0
57         lst.append(x*x)
58 print lst
59
60
61
62 # tuples
63 from math import sqrt
64 def csqrt(n):
65     if n >= 0:
66         return (sqrt(n),0)
67     else:
68         return (0,sqrt(-n))
69
70 real, imag = csqrt(-16)
71 print '%d+%di' % (real,imag)
72
73 # dictionaries (hash tables)
74 profs = {}
75 profs['6.047'] = 'kellis'
76 profs['7.012'] = 'lander'
77
78 profs['6.047']
79 profs['bogus']
80
81 # import, dir, help
82 from math import cos
```

```

83 cos(0)
84 import math
85 dir(math)
86 help(math.hypot)
87
88 # matrices (lists of lists)
89 m = [[0 for j in range(10)] for i in range(10)]
90 m
91 m[3][4]
92 m[4][5] = 6
93 m

```

2 Probability

1. We will quickly review some basic probability by considering an alternate way to represent motifs: a *position weight matrix* (PWM). We would like to model the fact that proteins may bind to motifs that are not fully specified. That is, some positions may require a certain nucleotide (e.g. A), while others positions are free to be a subset of the 4 nucleotides (e.g. A or C). A PWM represents the set of all DNA sequences that belong to the motif by using a matrix that stores the probability of finding each of the 4 nucleotides in each position in the motif. For example, consider the following PWM for a motif with length 4:

	1	2	3	4
A	0.6	0.25	0.10	1.0
G	0.4	0.25	0.10	0.0
T	0.0	0.25	0.40	0.0
C	0.0	0.25	0.40	0.0

We say that this motif can generate sequences of length 4. PWMs typically assume that the distribution of one position is not influenced by the base of another position. Notice that each position is associated with a probability distribution over the nucleotides (they sum to 1 and are nonnegative).

2. We can also model the *background distribution* of nucleotides (the distribution found across the genome):

A	0.1
G	0.4
T	0.1
C	0.4

Notice how the probabilities for A and T are the same and the probabilities of G and C are the same. This is a consequence of the complementarity DNA which ensures that the overall composition of A and T, G and C is the same overall in the genome.

3. Consider the sequence $S = GCAA$.

The probability of the motif generating this sequence is $P(S|M) = 0.4 \times 0.25 \times 0.1 \times 1.0 = 0.01$.

The probability of the background generating this sequence $P(S|B) = 0.4 \times 0.4 \times 0.1 \times 0.1 = 0.0016$.

4. Alone this isn't particularly interesting. However, given fraction of sequences that are generated by the motif, e.g. $P(M) = 0.1$, and assuming all other sequences are generated by the background ($P(B) = 0.9$) we can compute the probability that the motif generated the sequence using Bayes' Rule:

$$\begin{aligned}P(M|S) &= \frac{P(S|M)P(M)}{P(S)} \\ &= \frac{P(S|M)P(M)}{P(S|B)P(B) + P(S|M)P(M)} \\ &= \frac{0.01 \times 0.1}{0.0016 \times 0.9 + 0.01 \times 0.1} = 0.40984\end{aligned}$$

3 Basic definitions in molecular biology

1. The fundamental building blocks of DNA are A, T, G, C. RNA has the same *nucleotides* except for T which is replaced by U.
2. The central dogma of molecular biology states that DNA is *transcribed* to mRNA which is *translated* to proteins. Notice that because the nucleotide difference between DNA and mRNA is minimal, it is called transcription whereas the reading of mRNA to construct proteins is called translation.
3. Genes in DNA are **interrupted** by **introns** that do not code for proteins but often play an important role in regulation. mRNA has these introns stripped away and only contains **exons** or regions that are **expressed**.
4. Many organisms have their DNA broken into several chromosomes. Each chromosome contains two strands of DNA, which are complementary to each other but are read in opposite directions. Genes can occur on either strand of DNA. The DNA before a gene (in the 5' region) is considered "upstream" whereas the DNA after a gene (in the 3' region) is considered "downstream".