

Lecture 7

Numerical Computing

Prof. Steven Ludtke
N410.07, sludtke@bcm.edu

Lab #2

- Download the starter program from the class website. The starter program queries PubMed for articles written by a specific author, and prints the number of publications and number of coauthors. Try the program and understand what it's doing Use this program as a starting point for the following exercise. Each person in the group should select and complete one task (any of the 4), then one person should combine all of the work into a single program, which should be emailed to the TA before the end of class. You should coordinate your data representations so your code will all work together. Make sure each person in the group understands how all of the components work before leaving. As with lab 1, the tasks are ordered from most difficult to easiest:
 1. Rather than simply counting the number authors in the retrieved records, tabulate the number of times each author name appears, then print a) the number of authors and b) the top 10 authors with their publication count.
 2. Similar to 1), tabulate all of the words present in the title of all of the retrieved records. Eliminate common words like "the" and "and", and print an ordered list of the top 20 occurring words.
 3. Very often it is impossible to identify authors uniquely by their PubMed names. "Ludtke SJ" could be "Ludtke S". "Chen F" could represent dozens of different authors. While there is no perfect solution to this problem, to help assess the impact of the first issue, count the number of unique AU values, and also count the number of unique family names. Print the difference between these counts as the number of possible overlaps.
 4. Modify the program to query based on a provided keyword rather than an author's name. Pubmed does not like overly large queries, so be sure to restrict the number of retrieved records to a reasonable number (no more than 1000). Be sure to print a warning if this limit was reached.

Lab 2 Review

```
from Bio import Entrez
from Bio import Medline

author=raw_input("Enter author name: ")

Entrez.email = "user@bcm.edu"

handle = Entrez.esearch(db="pubmed", term="{}[Author]".format(author), retmax=500)
record = Entrez.read(handle)
handle.close()

ids=record["IdList"]

handle=Entrez.efetch(db="pubmed",id=tuple(ids),rettype="medline")
records = list(Medline.parse(handle))
handle.close()

coauth=set()
for r in records:
    coauth.update(r["AU"])

print "Coauthors:",coauth
print len(ids)," matching records"
print len(coauth)," coauthors"
```

Medline Terms

<u>Affiliation [AD]</u>	<u>Investigator [IR]</u>	<u>Pharmacological Action [PA]</u>
<u>Article Identifier [AID]</u>	<u>ISBN [ISBN]</u>	<u>Place of Publication [PL]</u>
<u>All Fields [ALL]</u>	<u>Issue [IP]</u>	<u>PMID [PMID]</u>
<u>Author [AU]</u>	<u>Journal [TA]</u>	<u>Publisher [PUBN]</u>
<u>Author Identifier [AUID]</u>	<u>Language [LA]</u>	<u>Publication Date [DP]</u>
<u>Book [book]</u>	<u>Last Author [LASTAU]</u>	<u>Publication Type [PT]</u>
<u>Comment Corrections</u>	<u>Location ID [LID]</u>	<u>Secondary Source ID [SI]</u>
<u>Corporate Author [CN]</u>	<u>MeSH Date [MHDA]</u>	<u>Subset [SB]</u>
<u>Create Date [CRDT]</u>	<u>MeSH Major Topic [MAJR]</u>	<u>Supplementary Concept [NM]</u>
<u>Completion Date [DCOM]</u>	<u>MeSH Subheadings [SH]</u>	<u>Text Words [TW]</u>
<u>EC/RN Number [RN]</u>	<u>MeSH Terms [MH]</u>	<u>Title [TI]</u>
<u>Editor [ED]</u>	<u>Modification Date [LR]</u>	<u>Title/Abstract [TIAB]</u>
<u>Entrez Date [EDAT]</u>	<u>NLM Unique ID [JID]</u>	<u>Transliterated Title [TT]</u>
<u>Filter [FILTER]</u>	<u>Other Term [OT]</u>	<u>UID [PMID]</u>
<u>First Author Name [1AU]</u>	<u>Owner</u>	<u>Version</u>
<u>Full Author Name [FAU]</u>	<u>Pagination [PG]</u>	<u>Volume [VI]</u>
<u>Full Investigator Name [FIR]</u>	<u>Personal Name as Subject [PS]</u>	
<u>Grant Number [GR]</u>		

Lab 2 Review

```
from Bio import Entrez
from Bio import Medline

author=raw_input("Enter search word: ")

Entrez.email = "user@bcm.edu"

handle = Entrez.esearch(db="pubmed", term="{}[TIAB]".format(author), retmax=500)
record = Entrez.read(handle)
handle.close()

ids=record["IdList"]

handle=Entrez.efetch(db="pubmed",id=tuple(ids),rettype="medline")
records = list(Medline.parse(handle))
handle.close()

coauth=set()
for r in records:
    coauth.update(r["AU"])

print "Coauthors:",coauth
print len(ids)," matching records"
print len(coauth)," coauthors"
```

Lab 2 Review

```
from Bio import Entrez
from Bio import Medline

author=raw_input("Enter search word: ")

Entrez.email = "user@bcm.edu"

handle = Entrez.esearch(db="pubmed", term="{}[TIAB]".format(author), retmax=500)
record = Entrez.read(handle)
handle.close()

ids=record["IdList"]

handle=Entrez.efetch(db="pubmed",id=tuple(ids),rettype="medline")
records = list(Medline.parse(handle))
handle.close()

coauth=set()
coauthlast=set()
for r in records:
    coauth.update(r["AU"])
    coauthlast.update(r["AU"].split()[0]) # <-- nope... won't work right

print "Coauthors:",coauth
print len(ids)," matching records"
print len(coauth)," coauthors with ",len(coauth)-len(coauthlast),"possible overlaps"
```

Lab 2 Review

```
from Bio import Entrez
from Bio import Medline

author=raw_input("Enter search word: ")

Entrez.email = "user@bcm.edu"

handle = Entrez.esearch(db="pubmed", term="{}[TIAB]".format(author), retmax=500)
record = Entrez.read(handle)
handle.close()

ids=record["IdList"]

handle=Entrez.efetch(db="pubmed",id=tuple(ids),rettype="medline")
records = list(Medline.parse(handle))
handle.close()

coauth=set()
coauthlast=set()
for r in records:
    for a in r["AU"]:
        coauth.add(a)
        coauthlast.add(a.split()[0])

print "Coauthors:",coauth
print len(ids)," matching records"
print len(coauth)," coauthors with ",len(coauth)-len(coauthlast),"possible overlaps"
```

Lab 2 Review

```
handle=Entrez.efetch(db="pubmed",id=tuple(ids
),rettype="medline")
records = list(Medline.parse(handle))
handle.close()
```

```
coauth=set()
for r in records:
    coauth.update(r["AU"])

print "Coauthors:",coauth
```


Lab 2 Review

```
handle=Entrez.efetch(db="pubmed",id=tuple(ids
),rettype="medline")
records = list(Medline.parse(handle))
handle.close()
```

```
from collections import Counter
c=Counter()
for r in records:
    c.update(r["AU"])

print "Coauthors:",c.most_common(10)
```

Lab 2 Review

```
handle=Entrez.efetch(db="pubmed",id=tuple(ids),ret
type="medline")
records = list(Medline.parse(handle))
handle.close()
```

```
auth={}
for r in records:
    for a in r["AU"]:
        try: auth[a]+=1
        except: auth[a]=1
```

```
authl=sorted(auth.items(),lambda a,b: b[1]-a[1])
```

```
print "Coauthors:",authl[:10]
```

Lab 2 Review

```
handle=Entrez.efetch(db="pubmed",id=tuple(ids),rettype="medline")
records = list(Medline.parse(handle))
handle.close()
```

```
allauth=[]
for r in records: allauth.extend(r["AU"])
```

```
uniqauth=set(allauth)
authcount=[]
for n in uniqauth: authcount.append((allauth.count(n),n))
```

```
authcount.sort(reverse=True)
```

```
print "Coauthors:",authcount[:10]
```

Homework Review

(print a power table)

```
print " ",
for e in xrange(2,11):
    print "{:>10d}".format(e),
print ""

for v in xrange(2,11):
    print "{:3d}".format(v),
    for e in xrange(2,11):
        print "{:10.5g}".format(v**e),

print ""
```

Homework Review

(center of mass from x/y file)

```
from sys import argv

lines=file(argv[1],"rU").readlines()
values=[]
for l in lines:
    v=l.split()
    values.append((float(v[0]),float(v[1])))

xsum,ysum=0,0
for x,y in values:
    xsum+=x
    ysum+=y

print xsum/len(values),ysum/len(values)
```

Homework Review

```
from sys import argv

lines=file(argv[1],"rU").readlines()
values=[]
for l in lines:
    v=l.split()
    values.append((float(v[0]),float(v[1])))

x,y=zip(*values)

print sum(x)/len(values),sum(y)/len(values)
```

Homework Review

```
from sys import argv
```

```
lines=file(argv[1],"rU").readlines()
```

```
values=[(float(v.split()[0]),float(v.split()[1])) for v in lines]
```

```
x,y=zip(*values)
```

```
print sum(x)/len(values),sum(y)/len(values)
```

Homework Review

```
import numpy as np
from sys import argv

xy=np.loadtxt(argv[1]).transpose()

print np.mean(xy[0]),np.mean(xy[1])
```


IEEE Floating Point

```
a=10000.0
```

```
b=0
```

```
for i in xrange(10000): b+=1
```

```
a==b
```

?

IEEE Floating Point

```
a=10000.0
```

```
b=0
```

```
for i in xrange(80000): b+=0.125
```

```
a==b
```

```
?
```

IEEE Floating Point

```
a=10000.0
```

```
b=0
```

```
for i in xrange(100000): b+=0.1
```

```
a==b
```

?

Decimal Numbers

```
1
0 1
0 0 1
0 0 0 1
X X X X
```

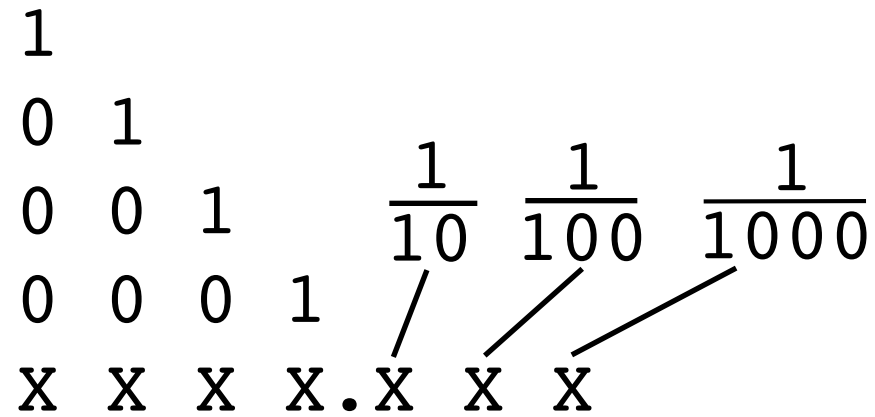
Binary Numbers

```
1
2 6 3 1
8 4 2 6 8 4 2 1
X X X X X X X X
```

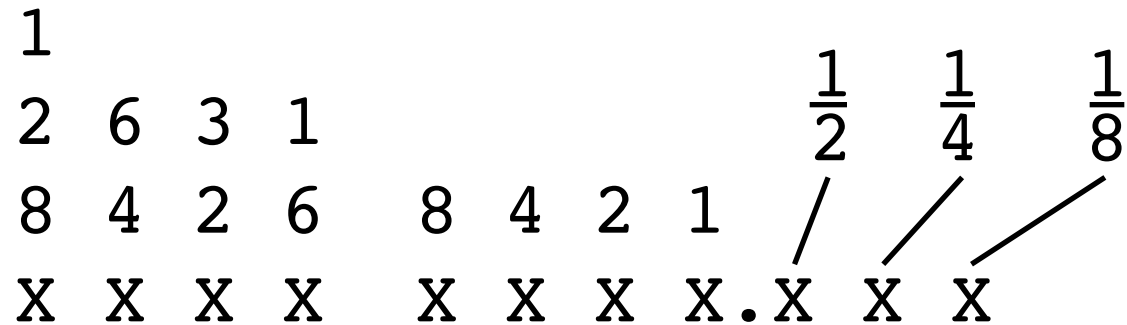
dec	binary
1	1
2	10
4	100
8	1000
15	1111
212	1101 0100

Decimal Numbers

1
0 1
0 0 1 $\frac{1}{10}$ $\frac{1}{100}$ $\frac{1}{1000}$
0 0 0 1
X X X X . X X X



Binary Numbers



0.25	0000 0000.010
0.625	0000 0000.101
0.3	0000 0000.0100 1100 1100 ...

IEEE Floating Point

- "Binary Scientific Notation"
- 127.25 1111111.01
- $1.2725 \times 10^2 \rightarrow 1.11111101 \times 10^{110}$

- Single:
- S EEEEEEEE EMMMMMMM MMMMMMMM MMMMMMMM
- S – Sign bit 0=+
- E – Exponent, bias 127
- M – Significand (Mantissa), implicit 1 when normalized

the 'decimal point' in binary

Hexadecimal

0xD2.A7

Binary

11010010.10100111

$\frac{1}{2}$ $\frac{1}{4}$ $\frac{1}{8}$ $\frac{1}{16}$ $\frac{1}{32}$ $\frac{1}{64}$ $\frac{1}{128}$ $\frac{1}{256}$

=210.65234375

=2.1065234375x10²

or

1.101001010100111x10¹¹¹

Digital Representation of Numbers

- Bit 0-1
- Nibble (4 bits) 0-15
- Byte (char) (8 bits) 0-255
- Word (short) (16 bits) 0-65,535
- Longword (long) (32 bits) 0-4,294,967,296
- Long Longword (64 bits) 0-1.844x10¹⁹
- Float (32 bits) 10³⁸ (7 digits)
- Double (64 bits) 10³⁰⁸ (15 digits)

normal
Python



Python “long integers” are a special object, not a standard number.

Numerical Processing

- Decimal - When you need accurate decimal values
- Fraction - Rational fractions
- NumPy - Fast arrays, linear algebra, etc.
- Matplotlib - Matlab-style plotting interface
- SciPy - Large library of numerical computing algorithms

Installing SciPy, etc.

- NumPy and matplotlib will be sufficient for the homework, SciPy is also a useful tool
- www.scipy.org
- Can consider: <http://matplotlib.sourceforge.net/users/installing.html>
- Linux
 - Use your package installer, should be available
- Mac (system python includes numpy)
 - WARNING - many installations may require that you have Xcode with command-line tools installed on your Mac ! It is free, but you have to install it !
 - If you are using 10.10 (and you should be) or 10.7-10.9, you may try:
 - <http://fonnesbeck.github.com/ScipySuperpack/>
 - For Snow Leopard or as an alternative for 10.7-10.9:
 - You can get Matplotlib easily as a side effect of installing (which we will use in a later lecture):
 - http://ncmi.bcm.edu/ncmi/software/counter_222/software_86
 - Otherwise, this tip may help:
http://blake.bcm.edu/emanwiki/EMAN2/COMPILE_EMAN2_MAC_OS_X#matplotlib_.281.0.29
- Windows
 - 32 bit - <http://sourceforge.net/projects/scipy/files/scipy/0.13.2>
 - 64 bit - You can try: <http://www.lfd.uci.edu/~gohlke/pythonlibs/>

What's wrong with this ?

```
from math import *
```

```
x=[i/10000000.0 for i in xrange(10000000)]  
y=[sin(i) for i in x]
```

How about this ?

```
from numpy import *  
  
x=arange(0,10.0,0.000001)  
y=sin(x)
```

NumPy

- <http://csc.ucdavis.edu/~chaos/courses/nlp/Software/NumPyBook.pdf> # numpy book

- `from numpy import *`
- `a=arange(60)`
- `b=a.reshape(10,6)` # make 2-D matrix
- `c=a.reshape(3,4,5)` # make 3-D (tensor)
- `b.shape` # current dimensions
- `b.size` # total number of elements
- `b.ndim` # dimensionality
- `b.dtype` # type of value stored
- `b.astype("")`

Type	Bit-Width	Character
<code>bool_</code>	<code>boolXX</code>	<code>'?'</code>
<code>byte</code>	<code>intXX</code>	<code>'b'</code>
<code>short</code>		<code>'h'</code>
<code>intc</code>		<code>'i'</code>
<code>int_</code>		<code>'l'</code>
<code>longlong</code>		<code>'q'</code>
<code>intp</code>		<code>'p'</code>
<code>ubyte</code>	<code>uintXX</code>	<code>'B'</code>
<code>ushort</code>		<code>'H'</code>
<code>uintc</code>		<code>'I'</code>
<code>uint</code>		<code>'L'</code>
<code>ulonglong</code>		<code>'Q'</code>
<code>uintp</code>		<code>'P'</code>
<code>single</code>	<code>floatXX</code>	<code>'f'</code>
<code>float_</code>		<code>'d'</code>
<code>longfloat</code>		<code>'g'</code>
<code>csingle</code>	<code>complexXX</code>	<code>'F'</code>
<code>complex_</code>		<code>'D'</code>
<code>clongfloat</code>		<code>'G'</code>
<code>object_</code>		<code>'O'</code>
<code>str_</code>		<code>'S#'</code>
<code>unicode_</code>		<code>'U#'</code>
<code>void</code>		<code>'V#'</code>

NumPy

- `a=zeros((nx,ny,...))`
- `a=fromfunction(lambda i,j:i+j,(4,5))`
- `a=loadtxt(filename)` # read multicolumn data from a text file
- `a=arange(0,20,.1)` # fractional version of `range()`
- `a*10` # multiply each element !
- `b=sin(a)` # `sin()` of each element
- `c=a[c>0]` # condition, returns elements `>0`
- `c.sort()` # sort values in-place
- `c.mean(),var(),std(),prod()` # average, variance, standard dev, product
- `inner(a,b), outer(a,b)` # inner and outer matrix products
- `dot(a,b), cross(a,b)` # dot and cross products (similar to above)
- `histogram(a,bins,range)` # compute a histogram of 'a'

Lambda

- Unnamed single-use functions:

```
lambda x,y: x*y*23
```

- equivalent to:

```
def f(x,y) : return x*y*23
```

SciPy

- Clustering package (scipy.cluster)
- Constants (scipy.constants)
- Fourier transforms (scipy.fftpack)
- Integration and ODEs (scipy.integrate)
- Interpolation (scipy.interpolate)
- Input and output (scipy.io)
- Linear algebra (scipy.linalg)
- Maximum entropy models (scipy.maxentropy)
- Miscellaneous routines (scipy.misc)
- Multi-dimensional image processing (scipy.ndimage)
- Orthogonal distance regression (scipy.odr)
- Optimization and root finding (scipy.optimize)
- Signal processing (scipy.signal)
- Sparse matrices (scipy.sparse)
- Sparse linear algebra (scipy.sparse.linalg)
- Spatial algorithms and data structures (scipy.spatial)
- Special functions (scipy.special)
- Statistical functions (scipy.stats)
- Image Array Manipulation and Convolution (scipy.stsci)

matplotlib (pylab)

- Matlab-like plotting library
- http://matplotlib.sourceforge.net/users/pyplot_tutorial.html

```
ipython --pylab          # special mode for interaction with pylab
x=arange(0,4*pi,0.05)    # from numpy
y=sin(x)                 # easy to apply a function to a list of values
plot(x,y)                # plot x,y and open a display window
```

OR

```
python
from pylab import *      # <-- only if you don't use ipython
x=arange(0,4*pi,0.05)
y=sin(x)                 # easy to apply a function to a list of values
plot(x,y)                # plot x,y and open a display window
show()                   # opens the plot window (blocks on some machines)
```

matplotlib (pylab)

```
ipython --pylab      # special mode for interaction with pylab

x=arange(0,4*pi,0.05) # from numpy

y=sin(x)             # easy to apply a function to a list of values

y2=cos(x)

figure(1)            # start a new figure

subplot(211)         # make a 1x2 set of plots and move to 1st

plot(x,y)            # plot x,y and open a display window

ylabel("sin(x)")

subplot(212)         # start on the 2nd subplot

plot(x,y2)          # second plot

ylabel("cos(x)")

xlabel("x")
```

Homework 4

- Install NumPy & Matplotlib on your computer before Monday lab. These packages come with SciPy, so you may optionally install that (or may have already)
- The rest of homework 4 will be tied into Lab 3, and will be assigned on Monday