

Working with Files and Data, including geospatial methods and plotting

Andrew D. Wickert

March 1, 2017

The geosciences are full of a range of data types—from mapping and surveys to chemical analyses to “data” from computer model outputs (spanning a whole range of complexities) to remotely-sensed images and more. Working with all of these data types in an efficient way is one of the primary uses of computers in the geosciences. In this section, you will learn:

1. How computers store data
2. How to work with basic ASCII and binary data sets, using Numpy tools and basic file-reading operations.
3. How to work with spreadsheets (e.g., Excel) when programming
4. How to work with geospatial data both within Python on its own and within a GIS framework

1 Data storage and retrieval, data plotting, GIS, and Python modules

Data can be stored as text (typically ASCII) or binary values. Often, one file contains one data set. However, more advanced data storage formats like NetCDF and HDF are also widely-used—these act as containers for multiple related pieces of multidimensional data. NetCDF is more common and is often used for atmospheric, oceanographic, glaciologic—generally climate-related data, atmospheric science data and models. It is now also being used by the Community Surface Dynamics Modeling System (CSDMS) as a standard format for geological model input and output. HDF was developed by NASA, and therefore is often used for satellite remotely-sensed data. Unless you do quite a lot of work with these kinds of data and models, you may not come across these—but it is important to know at least a little about them.

Along with the basics of data storage and retrieval, we are going to have to learn how to display the data. This will likely be one of the more useful skills taught in this course, because you will learn how to write reusable computer code to generate publication-quality figures. Learning the plotting commands takes a bit of time, as does writing the code—but in my experience, this saves hours in the long run. Having the ability to generate reusable plots also makes one more willing to revise data analyses, knowing that time will not be wasted on laboriously recreating plots by hand. It also lets one make a whole set of similar plots for displaying large amounts of data.

Oftentimes in the geosciences, we want to plot geospatial data. I will show how you can install mapping packages and use GIS tools to do just this. I will also point out and (in some cases) very briefly discuss the Python interfaces to GRASS GIS, Arc GIS, and QGIS, the three leading scientific geospatial platforms. These come with a set of premade tools to make it possible to run very involved analyses in just a few lines of code. We will spend some more significant time going over GDAL and OGR: these are the libraries that underlie all GIS engines, and can be accessed directly from Python.

pip and easy_install

Using some of data management (and other!) tools involves downloading and installing additional Python modules generated by the community of Python users. This may be done in an automated way by using `pip` and `easy_install`. You may need them for some of the packages to work with specific data sets in this guide. Indeed, there are so many useful such packages to do many of the activities that you might want to do that it is often good to do an Internet search for “python <task you want to do>” before starting any work.

2 ASCII

2.1 Theory

In the Introductory notes, you learned about the ASCII table for encoding text. Each ASCII character requires 7 bits to be encoded. This increases to 8 bits, or 1 byte, as a “stop” bit is added—this was historically used to note if there was an error in the transmission.

In ASCII, the number “45” would require 2 bytes of storage. “61.083” would require 6 bytes. “-15E-3”, where “E” denotes that 15 is multiplied by 10 to the following power (in this case, -3), would also require 6 bytes. In most situations, ASCII data take more storage space than binary data. They have the advantage, though, in that they are fully human-readable.

The great line-ending schism

A very important historical note is the controversy about line-endings. It may sound trivial, but it can turn a successful data import into a failure! Here it is:

- UNIX-based computers end lines of text with a newline or “line feed” character, `\n` (ASCII 10)
- DOS/Windows-based computers end lines of text with a carriage return, followed by a newline character `\r\n` (ASCII 13, then ASCII 10)

This seemingly small detail has its basis in the mechanics of the transition from typewriters to computers: with typewriters, one must advance to a new line (“line feed” or “newline”, ASCII 10) and push the roll of paper all the way back to the start of the line (“carriage return”, ASCII 13). Computers can instead produce a new line all at once, because they are not limited to moving paper in physical space. This led to ASCII 10 being used on UNIX systems, and the typewriter-looking [ASCII 13, ASCII 10] being used on Windows systems, and a potential whole world of trouble! It can mean that, if you have a file with delimited data (e.g., csv – comma-separated values), that you cannot tell where one line ends and another begins if your computer looks for the wrong newline character! This is a problem that can often happen if, like most computer users, one is working on Windows but then needs to do a project on a supercomputer (mostly UNIX). Many pieces of software are becoming smarter about this difference, but this is one of the great and unfortunate schisms in computing that happened in the early days, when we were really just learning what we are doing, and has been carried forward to the present.

2.2 Practice

OK, enough with the theory. Let’s get our hands dirty! Figuratively, of course, in case you have been eating in a particularly messy way over your computer keyboard. (Aside: have you ever shaken out an old keyboard? Don’t do it over your face. Especially not with your mouth open. No, I wouldn’t know.)

As mentioned in the introductory section, there is a Matplotlib tutorial at http://matplotlib.org/users/pyplot_tutorial.html. And for those of you who like to learn about these sorts of things by following examples, Matplotlib has an excellent plotting gallery—a place where it shows you a whole range of graphics and how to create them, at <http://matplotlib.org/gallery.html>.

You should have already installed `numpy`. This is the numerical Python package and is important for doing work with all sorts of arrays. We are going to start by importing some data from an ASCII text file. In this case, it will be a simple transect of elevations across the park by where I grew up. We will use both the Numpy `genfromtxt`

feature to create a numpy array and the `pandas read_csv` feature to generate a Pandas DataFrame. This example already includes plotting with matplotlib as does a section in the introductory notes. This, I hope, will be fairly self-explanatory.

If you are on Windows or Mac, your Python distribution should have come with Pandas. If you are on Linux (or Mac and are using a package manager) and do not yet have Pandas, you can type `sudo apt-get install python-pandas` or a similar command.

```
2  #! /usr/bin/env python
3
4  # 1 --- numpy
5  import numpy as np
6  from matplotlib import pyplot as plt
7
8  # Strip out header information
9  # Just remember that columns are: x, z, lat, lon
10 data = np.genfromtxt(' ../../data/BattleCreekProfile.txt', skip_header=1, \
11                     delimiter=',')
12
13 # Plot as distance profile
14 plt.plot(data[:,0]/1000., data[:,1], 'k', linewidth=2)
15 plt.xlabel('Distance [km]')
16 plt.ylabel('Elevation [m]')
17 plt.title('NUMPY!')
18 plt.show()
19
20 # 2 --- pandas
21 import pandas as pd
22
23 # Import it as a data frame --- keep information in a coordinate system
24 data = pd.read_csv(' ../../data/BattleCreekProfile.txt')
25
26 # Plot --- hey look, distance is inferred!
27 plt.plot(data['Distance (m)']/1000., data['Elevation (m)'], 'k', linewidth=2)
28 plt.xlabel('Distance [km]')
29 plt.ylabel('Elevation [m]')
30 plt.title('PANDAS!')
31 plt.show()
32
33 # Think about how useful this could be for time-series data!
34 # Indeed, this is the main use of Pandas, which is designed as a
35 # data-management library
```

code/FilesData/CSVexample.py

Time series and the datetime library

If you want to work with time-series data, should look into the `datetime` library for powerful functions to deal with date/time objects.

You can also save ASCII data using the `np.savetxt` feature. I mostly use it as follows:

```
np.savetxt('OUTPUT_FILE_NAME.TXT', InputVariableName, fmt='%FORMAT_STRING')
```

Here, it is important to mention **formatting strings**. These follow the conventions from the C programming language. I will just introduce my most-used subset of them here.

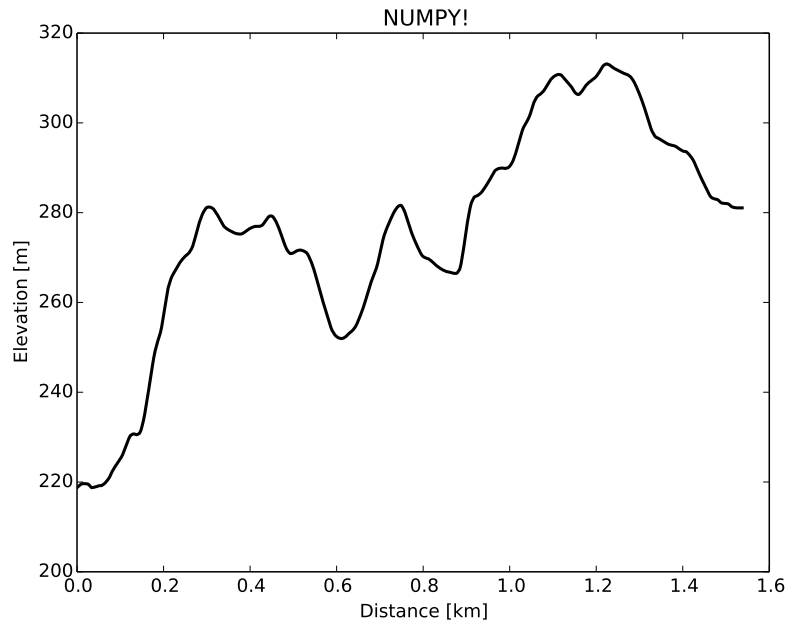


Figure 1: Topographic profile from the Mississippi River valley (left) through tributary valleys and an end moraine complex (right).

Formatting strings contain a character. This determines how the incoming value is treated, but all output is as a string:

- d: integer
- f: floating point (“float”)
- s: string

These are then modified by numbers placed in front of the digits. Both the formatting string and the value that one would like to format as such must be preceded by a %, and the formatting string must be inside quotes to denote that it is a string. I will explain through a set of examples:

```

2  #! /usr/bin/env python
3
4  # Formatting strings
5
6  # d — integer
7  # f — floating point
8  # s — plain string
9
10 # Let's see how these change the formatting of a number.

```

```

10 num = 14.81
12 print "Original number", num
13 print ""
14 print '%d:' ,
15 print ( '%d' %num)
16 print ""
17 print '%f:' ,
18 print ( '%f' %num)
19 print ""
20 print '%s:' ,
21 print ( '%s' %num)
22 print ""

24 # Note that the '%d' option just cuts off the decimal without rounding the
26 # number! So it is doing floor division .

28 # The '%f' option has a pre-set decimal precision .

30 # The '%s' option just converts the number into a string , in the same way that
32 # str(num) would do so .

34 # Now let's start having a bit more control .

36 # %d
37 # The number of digits in a decimal number may be specified
38 # This can create spaces before that number
39 print '%5d %14.81 :'
40 print ( '%5d' %14.81)
41 print ""
42 # But will not truncate the number: it simply sets the minimum number of
43 # digits
44 # that must be displayed .
45 print '%5d %6132614.81 :'
46 print ( '%5d' %6132614.81)
47 print ""
48 # It can also be used to generate zeros before a number instead of the spaces .
49 # This is known as zero-padding (0-padding)
50 print '%05d %14.81 :'
51 print ( '%05d' %14.81)
52 print ""

54 # %f
55 # Floating point numbers can be formatted as follows
56 # [TOTAL NUMBER OF CHARACTERS, INCLUDING DECIMAL].[NUMBER OF DIGITS AFTER
57 # DECIMAL POINT]
58 # — OR —
59 # .[NUMBER OF DIGITS AFTER THE DECIMAL POINT]
60 # (with the number of digits before the decimal point being fit to the number
61 # that you are formatting as a string)

62 # If there is not enough precision , it still is truncated . But in the
63 # floating-point case , it is rounded .
64 print '%4.1f %14.89 :'
65 print ( '%4.1f' %14.89)
66 print ""
67 # This is the perfect size for this
68 print '%5.2f %14.89 :'
69 print ( '%5.2f' %14.89)
70 print ""
71 # This automatically sets the proper space to the left of the decimal point
72 print '%.2f %14.89 :'
73 print ( '%.2f' %14.89)
74 print ""
75 # Extra space and decimal places

```

```

76 print '%9.4f %14.81 : '
    print ( '%9.4f' %14.89)
    print ""
78 # 0-padding with extra decimal places
    print '%09.4f %14.81 : '
80 print ( '%09.4f' %14.89)

```

code/FilesData/formattingStrings.py

Saving multiple files with a numbering scheme

In many cases, you may have data that belong in a particular order—whether they are a set of time-steps, a series of different analyses, or a number of different sample ID's that you want to automatically generate. For the case with real dates and times, you may use that date and time as part of the file name. And if you use it as **yyyymmdd**, it will alphabetically sort correctly! This becomes harder though if it is an arbitrary time step – like say, millions of years ago, seconds, or just an arbitrary set of **[0, 1, 2, ..., 51, etc.]**. If you want the data to be sorted in order when you do a simple alphabetical sorting – great for viewing in file browsers, loading into data analyses, or just for general orderly storage, you can run into the problem in which your data show up as **[0, 1, 10, 100, 2, ..., 51, etc.]**. This is certainly not what we want? So how do we fix it?

We use **zero-padding** (see code above on string formatting). So in our above example, if we can safely say that all the numbers are integers, and nothing will go over 9999, we can write an output file name as:

```
fnpadded = 'some_descriptive_text_' + '%04d' %time_step + '.ext'
```

And then you can use this to save text output as a floating point with 2 decimal places of precision, called for example `plate_reconstruction`, as follows:

```
np.savetxt(fnpadded, plate_reconstruction, fmt='%0.2f')
```

Sometimes, ASCII data are not so easy to work with. This can happen when they are not in a simple grid. In that case, we have to use lower-level Python commands. To learn basic file handling, you may see the help at <http://www.pythonforbeginners.com/files/reading-and-writing-files-in-python>.

One example of this sort of data is a set of GPS tracks provided by Ben P. These look something like:

```

0
2 728609.215175176 7272865.670114477 3850.11220000000032
  728608.813141581 7272866.856349733 3850.17690000000058
4 ...
  728588.585519101 7272909.968807690 3851.28429999999993
6 END
  1
8 729772.085292112 7284330.246414313 3857.98080000000048
  ...
10 729789.783887492 7284223.112814812 3858.06919999999981
   END
12 ...

```

The numbers and “END” lines separate individual tracks. However, these are not nice 3-column entries. How do we parse a file like this?

Using the more basic Python file read/write commands (see links above), one can input the lines of the file, parse them into lists, and then turn them into arrays. This is as follows:

```
#!/usr/bin/env python
2
import numpy as np
4 from matplotlib import pyplot as plt

6 # Open text file for reading
f = open('../data/PG_astgtm2_dgps_tracks.txt', 'r')
8
# Prepare a list to contain each individual track
10 tracks = []

12 # Now there are two ways to do this. The first is easier, and requires more
# memory on your computer. In this, we will read in the whole file at once
14 # and then parse it.

16 # This reads all lines in the file, and creates a list in which
# each entry is a string that is that line.
18 text = f.readlines()

20
# Go back to the beginning of the file
22 f.seek(0)

24 # This list will contain the points in each individual track.
track = []
26 # and this will contain the track number as recorded by the GPS
track_numbers = []

28 # This next statement is *asking* for an infinite loop. That is why
# I am introducing a new flow-control statement, break, that will get
30 # us out of it.
32 while True:
    # readline() reads the next line; strip() removes newline characters
    34 line = f.readline().strip()
    if line:
        36 try:
            # if the line is an integer, it is a track number
            38 track_numbers.append(int(line))
        except:
            40 # Check if the line is ending a section; if it is, package the section
            # and ship it off to the "tracks" list in its own numpy array.
            42 if line == 'END':
                # Only do this if we need another track to be entered; there are
                # two "END"s at the end of the file, so this will prevent it from
                # adding an empty track there that does not correspond to the numbers
                46 if len(track_numbers) > len(tracks):
                    tracks.append(np.array(track))
                # reset individual track list for the next one
                48 track = []
            else:
                50 # I will here in two steps split the line with data into a list and
                # turn it into a numpy array of floating point values.
                52 tmp0 = line.split(' ') # Split it at the spaces
                # Everything is still a string, so need to tell numpy to make the
                # array of floating-point values
                54 tmp1 = np.array(tmp0, dtype=float)
                # Now append it to track — we will change this to an array in the
                # step before we append it to the "tracks" master list (above)
                56 track.append(tmp1)
```



```

60     else:
61         # This is how we get out of the potentially infinite loop: if the
62         # line is empty
63         break
64
65 # Now let's plot all of these tracks' x and y components
66 # with the default set of different colors per line
67 # Remember, [:,i] means ALL ROWS IN iTH COLUMN
68 fig = plt.figure()
69 for line in tracks:
70     plt.plot(line[:,0], line[:,1]) # Easting, Northing
71     plt.title('GPS tracks', fontsize=20, fontweight='bold')
72     plt.ylabel('Northing', fontsize=20)
73     plt.xlabel('Easting', fontsize=20)
74     plt.tight_layout() # Formatting helper
75     plt.show()
76
77 # And let's now combine all of the tracks together into a single numpy array
78 # to create a large set of points
79 # Numpy arrays play nicely with numpy lists, so we just need to use the
80 # concatenate command to do this in one step!
81 # Other ways to combine numpy arrays include using "np.vstack" and "np.hstack"
82 alltracks = np.concatenate(tracks)
83
84 # And let's see if there is any statistical clustering of elevations
85 plt.hist(alltracks[:, -1], bins=100)
86 plt.title('GPS track hypsometry', fontsize=20, fontweight='bold')
87 plt.ylabel('Number of measurements', fontsize=20)
88 plt.xlabel('Elevation [m]', fontsize=20)
89 plt.tight_layout() # Formatting helper
90 plt.show()
91 # Yep, there are some distinct hypsometric peaks!
92 # But look at the x-axis and how crowded those labels are. We'll fix that
93 # in the next plot.
94
95 # Now, let's get some practice plotting both together:
96 plt.figure(figsize=(14,6)) # Wide figure
97 # 1 row, 2 columns, panel 1
98 ax1 = plt.subplot(121)
99 for line in tracks:
100     ax1.plot(line[:,0], line[:,1]) # Easting, Northing
101 # You have to set the axis object properties here
102 ax1.set_title('GPS tracks', fontsize=20, fontweight='bold')
103 # Rotate x-ticks
104 # Can also use strings like "vertical"
105 # And can use plt.xticks (or plt.yticks) to set many other tick-related
106 # parameters.
107 # See, for a straightforward example:
108 # http://matplotlib.org/examples/ticks\_and\_spines/ticklabels\_demo\_rotation.html
109 plt.xticks(rotation=60) # use of "plt" explained below for ax2
110 ax1.set_ylabel('Northing', fontsize=20)
111 ax1.set_xlabel('Easting', fontsize=20)
112 # 1 row, 2 columns, panel 2
113 ax2 = plt.subplot(122)
114 ax2.hist(alltracks[:, -1], bins=100)
115 # Check this out — "plt" applies to the currently-selected axis. So it works
116 # here just as well as "ax2."
117 plt.title('GPS track hypsometry', fontsize=20, fontweight='bold')
118 plt.xticks(rotation=60)
119 ax2.set_ylabel('Number of measurements', fontsize=20)
120 ax2.set_xlabel('Elevation [m]', fontsize=20)
121 # Now autoformat a nice layout and show the figure
122 plt.tight_layout() # Formatting helper
123 plt.show()

```

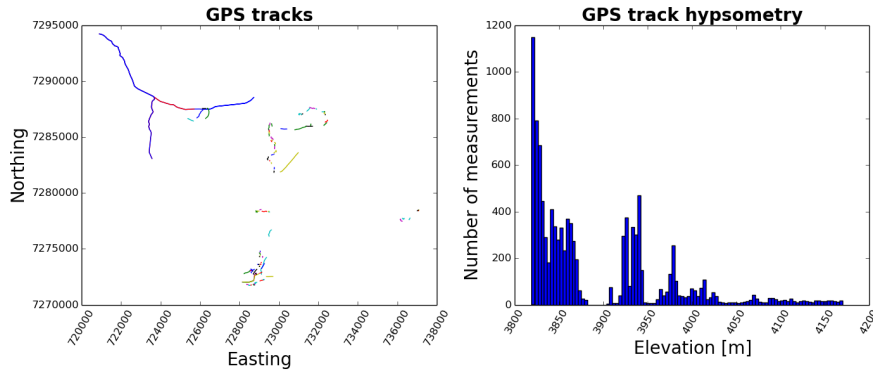


Figure 2: GPS tracks (left) and hypsometry (right).

3 Binary

3.1 Theory – and how to get the most out of storage space

Binary data are represented as a set of ones and zeros. In the Introductory notes, you learned how binary works, and a bit about how numbers may be generated.

Binary data are usually more compact than ASCII data, especially if used correctly. For example, a number between 0 and 65535 may be represented by two bytes (16 bits); this can be seen in binary because $2^{16} = 65536$, and we want to include 0, so have to shift the maximum value down one. This is called a 16-bit unsigned integer, because it has 16 bits of data, and does not include a + or - sign (and hence is always positive). A number like this would require 5 bytes in ASCII to represent.

If you are representing a large number of binary values (0 or 1) as 16-bit unsigned integers, you would be using 2 bytes of data per value, while ASCII would require only 1 byte. So in this case, ASCII would be better! But this is where it becomes important to *intelligently choose binary representations of data*. If we represented each of these values as binary logical values, each would only require 1 bit of storage space—an $8\times$ improvement over ASCII and a $16\times$ improvement over the 16-bit unsigned integer, which is really just too much storage space for binary data.

We often denote these values as follows:

- unsigned integer: `uint`
- signed integer (so can be + or -): `int`

- floating point: `float`

You might also see terms like `single`, `double`, `char`, `word`, etc.; these are less-descriptive terms that also relate to number of bytes in data. I use the three above to be more clear to those who have not memorized what all of these are.

These can also be used to denote how many bits are involved in each data type. For example:

- unsigned 16-bit integer: `uint16`
- signed 32-bit integer (so can be + or -): `int32`
- 64-bit floating point: `float64`

3.2 Practice – raw binary files

```

1  #!/usr/bin/env python
2  import numpy as np

4  # Create a 2x3 numpy array
a = np.array([[1,2,3],[4,5,6]])

6  # Save it as a straight binary output with different precisions.
8  # Check the filesize after each of these
filename = 'testout.bin'

10 # 8-bit unsigned integer
12 a.astype('uint8').tofile(filename)

14 # 64-bit floating point
a.astype('float64').tofile(filename)

16 # Now load the saved file
18 b = np.fromfile(filename, dtype='float64')
# Hm, we lost the information about the line ending! This is because it
20 # is just a string of binary without any shape data.

22 # Let's see what happens if we try to load it as a 8-bit integer (signed)
c = np.fromfile(filename, dtype='int8')
24 # WOW! It worked. But what do the values look like? Hm... so you
# could combine those sets of binary values together in clumps of 8, and
26 # convert those 64-bit clumps of binary data into the original set of
# values. Not bad, huh? That's binary for ya!

```

code/FilesData/BinaryIO.py

As a quick mental exercise, imagine that you have the numbers -15, 35, 119, and 43. What is the ideal number of bits with which to represent each number? How about 0, 50612, 151, 10512, 85, 3160? *Answers: 8 bits (int8), 16 bits (uint16).*

3.3 Binary containers

3.3.1 Numpy files

Standard binary data formats are all right. But there are formats that can remember rows and columns of data, as well as “container” formats that can contain multiple arrays. So let’s look at the native ones in numpy first.

```

1  #!/usr/bin/env python
   import numpy as np
3
4  # Create two 2x3 numpy arrays
5  # int8 is more than enough to represent these data
   a = np.array([[1,2,3],[4,5,6]], dtype='int8')
7  # float16 is a really rarely-used format, but I am using it here just to
   # illustrate how it stores these data
8  b = np.array([[1.6,0.1,3.512],[-27,5,6.5109]], dtype='float16')
9  # look at how imprecise b is!
10 print b
11
12 # Let's look at a magic trick that numpy does. You know already that int8 can
   # only take numbers from -128 to +127. Well, what happens if we add 500 to a?
13 c = a + 500
14 print c
15 print c.dtype
   # Hey — it made it be a int16! That's pretty cool
16 # (Other languages like C would instead wrap around to -128 and continue
   # forward, causing potential big problems)
17
18 # Now let's save an array into an *.npy file
19 np.save('testout.npy', a)
   # And load it — it keeps the structure of the rows and columns, great!
20 d = np.load('testout.npy')
21
22 # How about saving multiple arrays? we can use a compressed *.npz file
23 np.savez('testout.npz', a=a, b=b, c=c)
24 # The reason for the i=i format is because this is telling us to take array
   # "c", for example, and to also call it "c" in the *.npz file.
25 # We could likewise change the names:
26 np.savez('testout.npz', x=a, y=b, z=c)
27
28 # Now let's load it
29 zfile = np.load('testout.npz')
30 print zfile['x'] == a
31 print zfile['y'] == b
32 print zfile['z'] == c

```

code/FilesData/npIO.py

3.3.2 Raster data – a geospatial example

Instead of just showing any standard raster data set, I am going to use a GeoTIFF that covers Berlin and the edge of Potsdam. In order to work with this, you are going to have to install the Python GDAL library.

GDAL and OGR GDAL (raster) and OGR (vector) are the key software libraries behind most GIS applications. You can access them directly in Python, in addition to being able to access more general GIS libraries.

Installing these libraries can be done with either your package manager (`sudo apt-get install python-gdal`), or something like this), or with pip (`sudo pip install gdal`). Seriously, do a good Google search about pip and how it works. It is super useful, with tools for remote sensing, dealing with multiple data types, geological modeling, and of course lots of non-geological work.

```

2  #!/usr/bin/env python
3
4  import gdal
5  from matplotlib import pyplot as plt
6  import numpy as np
7
8  # GDAL is the GIS library that runs inside of standard GIS software.
9  # It is used for raster data (OGR is used for vector.)
10 # You can use it directly with Python.
11
12 # Here we are loading the 1 arcsecond SRTM DEM of Berlin (extending to around
13 # Potsdam Hbf) as a GeoTIFF
14 ds = gdal.Open('../data/n52_e013_1arc_v3.tif')
15
16 # Try typing "ds." (without the quotes) and then two tabs in iPython. Look at
17 # how much extra information is packaged with this GeoTIFF!
18 # Try these:
19 ds.GetProjectionRef()
20 ds.GetGeoTransform()
21
22 band = ds.GetRasterBand(1) # only one band — elevation
23 elevation = band.ReadAsArray() # So read it as a numpy array.
24
25 # OK — let's plot it!
26 plt.imshow(elevation)
27 plt.colorbar()
28 plt.title('Berlin and surrounding Brandenburg')
29 plt.show()
30
31 # All right, now let's plot the proper coordinate system.
32 loc = ds.GetGeoTransform()
33 x = np.linspace(loc[0], \
34                 loc[0] + loc[1]*elevation.shape[1], \
35                 elevation.shape[1])
36 y = np.linspace(loc[3] + loc[5]*elevation.shape[0], \
37                 loc[3], \
38                 elevation.shape[0])
39
40 # Now plot
41 plt.imshow(elevation, extent=[x.min(), x.max(), y.min(), y.max()])
42 cbar = plt.colorbar() # Instantiate a class so I can modify its
43 # characteristics
44 cbar.set_label('Elevation [m]', fontsize=16, fontweight='bold')
45 plt.xlabel('Longitude E', fontsize=16, fontweight='bold')
46 plt.ylabel('Latitude N', fontsize=16, fontweight='bold')
47 plt.title('Berlin and surrounding Brandenburg\nSRTM 1-arcsecond data', \
48          fontsize=20)
49 plt.tight_layout()
50 plt.show()

```

code/FilesData/GDAL_DEM.py

3.3.3 Working with NetCDF and HDF files – in brief!

These “container” file formats require special libraries to work with them. This is further complicated by the fact that there are different standards: NetCDF3 is being supplanted by NetCDF4, and HDF4 is still used for many satellite data sets but HDF5 being the better-supported current “standard” in Python. Did this last sentence just confuse you with too many acronyms? Me too! So let’s try to simplify this problem.

First, install packages in the terminal (use “sudo” if you are on a UNIX-like system):

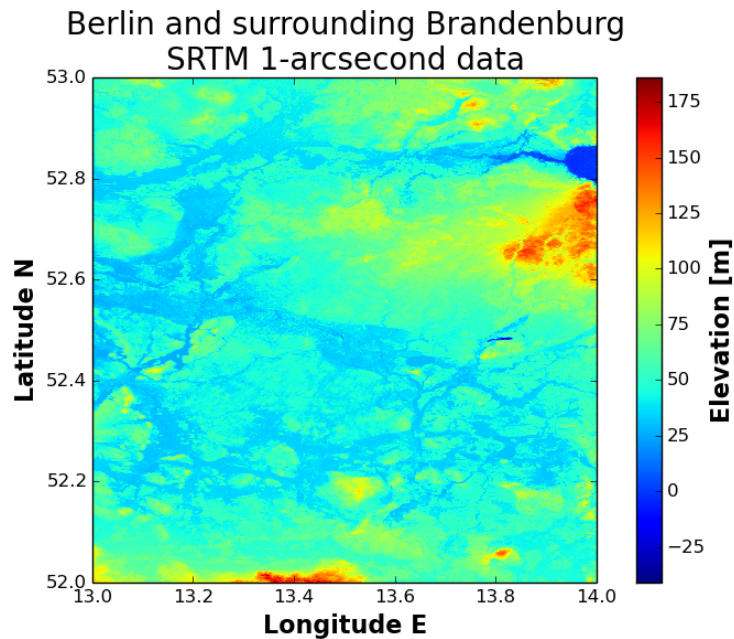


Figure 3: The DEM of the Berlin area does not show a “well-ordered” erosional fluvial landscape. Instead, it shows a landscape that has been greatly altered by glacial and fluvial processes around the ice margin.

```
2 pip install h4py
   pip install netcdf4
```

In addition to these, there exists a NetCDF-3 library that is built-in to SciPy. Each of these has good help (and example) files. See:

- h4py: <http://download.nexusformat.org/sphinx/examples/h5py/index.html> and http://www.icare.univ-lille1.fr/howto_hdf/#data_extraction_Python
- netCDF4 (works with NetCDF-4 and NetCDF-3 files): <https://github.com/erdc-cm/netCDF4-Python/blob/master/examples/tutorial.py>
- SciPy NetCDF: http://docs.scipy.org/doc/scipy-0.14.0/reference/generated/scipy.io.netcdf.netcdf_file.html (there used to be a more full tutorial but I cannot find it any longer)

In addition, this webpage has a lot of information on NetCDF (in general) and on the

Basemap mapping toolkit, which will be discussed below: <http://www-pord.ucsd.edu/~cjiang/python.html>.

So it seems that in this case, the only thing that is left out is HDF4. What do we do with that? We convert it! Tools exist to convert HDF4 into either NetCDF or HDF5.

- HDF4 → HDF5: <http://www.hdfgroup.org/h4toh5/> (main page); <http://www.hdfgroup.org/h4toh5/download.html> (Windows and Linux download); <http://hdfeos.org/software/h4toh5/bin/mac/> (Mac download)
- HDF4 → NetCDF4: <http://hdfeos.org/software/h4cflib.php>

Below, I have converted an HDF4 file from the TRMM satellite (Tropical Rainfall Measurement Mission) into a NetCDF-4 file, and import and plot it using Python. I first plot it using the standard matplotlib libraries, but later in the notes, use CartoPy to project it properly.

```
#!/usr/bin/env python
2
3 from netCDF4 import Dataset
4 from matplotlib import pyplot as plt
5 import numpy as np
6
7 # Open NetCDF file for reading
8 ncf = Dataset('../data/3B43.20040601.7A.nc', 'r', format='NETCDF4')
9
10 # Check out the variables
11 print ncf.variables.keys()
12 # Let's see a bit more about precipitation
13 print ncf.variables['precipitation']
14 # And let's get the lat and lon values
15 lat = ncf.variables['nlat'][:]
16 lon = ncf.variables['nlon'][:]
17
18 # Let's get the dates of the observation
19 dates = ncf.FileHeader.split('\n')[4:6]
20 print dates
21 start_date = dates[0].split('=')[-1]
22 end_date = dates[1].split('=')[-1]
23
24 # Now let's plot precipitation
25 plt.figure(figsize=(14,6))
26 # Show image
27 plt.imshow(np.flipud(ncf.variables['precipitation'][:].transpose()), \
28            extent=[lon.min(), lon.max(), lat.min(), lat.max()])
29 plt.colorbar()
30 # Label it with the contained information
31 plt.title('June 2004: ' + ncf.variables['precipitation'].long_name + ', [ '
32          + ncf.variables['precipitation'].units + ' ]')
33 plt.xlabel('Longitude E')
34 plt.ylabel('Latitude N')
35 plt.tight_layout()
36 plt.show()
```

code/FilesData/netCDF4basic.py

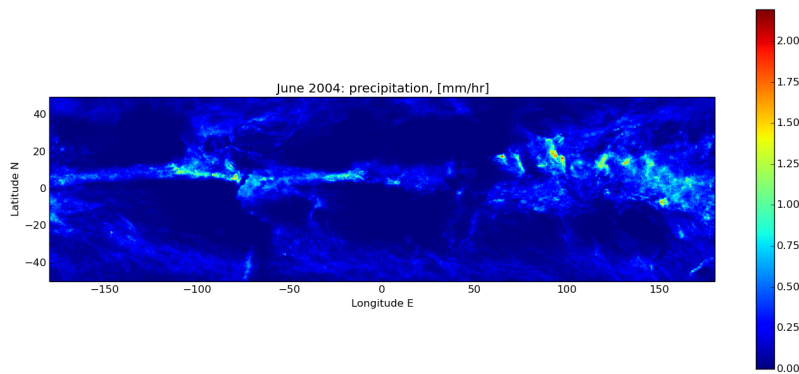


Figure 4: TRMM data from June 2004 (which is when the writer of these course notes graduated from high school).

4 Working with Spreadsheets (Excel or LibreOffice Calc)

So far, we've discussed all of these data formats that might be a bit more complex – and these are all important for work with computers and larger data sets. However, many of us in the Earth sciences deal with spreadsheet data. Up until now, most (all?) of you have worked only by hand with this data. But no longer!

There many packages to work with spreadsheets. I will show you first the builtin method with Pandas here, but you should know about at least the two following ones (available via `pip`):

- `openpyxl` (Excel)
- `odfpy` (LibreOffice)
- `oosheet` (LibreOffice calc – a spreadsheet-focused alternative to `odfpy`)

Here, we are going to plot an oxygen isotope history from a Greenland ice core:

```
import pandas as pd
2 from matplotlib import pyplot as plt

4 # Import data
  sheet_title = 'Renland d18O' # set outside for use in plotting
6 data = pd.read_excel('../data/vinther2008renland-agassiz.xlsx', \
                        sheetname=sheet_title, header=72)
8 # Note that they also include a % sign in front of all header lines.
  # How might you re-code this to dynamically learn how long the header is
10 # and skip over it?

12 # Give headers better names
  data.columns = ['years [b2k]', 'depth [m]', 'del18O [permil]']
14
```



```

# And plot
16 plt.plot(data['years [b2k]'], data['del18O [permil]'])
# plt.gca().invert_yaxis() # Grabs current axes and inverts x-axis
18 # plt.gca().invert_yaxis() # uncomment if you prefer to see time this way.
# Automatically use column labels in the plot
20 plt.title(sheet_title, fontsize=16, fontweight='bold')
plt.xlabel(data.columns[0], fontsize=16)
22 plt.ylabel(data.columns[-1], fontsize=16)
plt.show()
24
# If we want better axis labels, we can use the LaTeX formatting:
26 plt.figure(figsize=(12,5)) # make figure extra-wide (width, height)
plt.plot(data['years [b2k]'], data['del18O [permil]'])
28 plt.title('Greenland Ice Core: Renland', fontsize=16, fontweight='bold')
plt.xlabel('Years b2k', fontsize=16)
30 # Below, I use LaTeX formatting to write the Greek "delta" and the superscript
plt.ylabel('δ18O', fontsize=16)
32 plt.tight_layout()
plt.show()

```

code/FilesData/GreenlandXLSXpandas.py

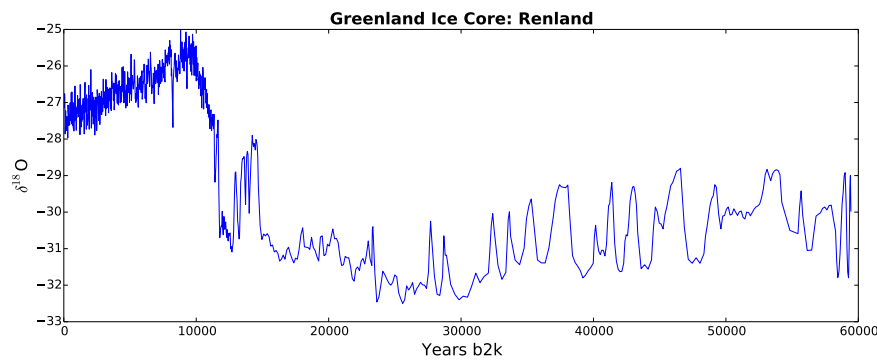


Figure 5: This is what the final graph of $\delta^{18}\text{O}$ should look like from the above code. Up is warm and down is cold, roughly-speaking: these values are also related to global ice volume. Note the Bølling-Allerød warm period ~ 14.5 ka, the Younger Dryas ~ 11.7 – 13 ka, and the 8.2 ka event (the little negative spike in the Holocene related to catastrophic meltwater release from North American proglacial lakes into the North Atlantic ocean).

Openpyxl is also easy-to-use. Here is a quick example that I wrote for some real research – to analyze sets of grain sizes of sediment from point counts of clasts. These data were transcribed from our field notebooks into a standard Excel format (thanks to Sara Savi), and I wrote the code that uses her standard Excel format to parse and compute statistics of different counts of grain size.

Go ahead and brainstorm how you might use this in your own work; we may spend some time going over this in class, as spreadsheets are really common.

```

1 import pandas as pd
2 from matplotlib import pyplot as plt
3
4 # Import data
5 sheet_title = 'Renland d18O' # set outside for use in plotting
6 data = pd.read_excel('../data/vinther2008renland-agassiz.xlsx', \
7                      sheetname=sheet_title, header=72)
8 # Note that they also include a % sign in front of all header lines.
9 # How might you re-code this to dynamically learn how long the header is
10 # and skip over it?
11
12 # Give headers better names
13 data.columns = ['years [b2k]', 'depth [m]', 'del18O [permil]']
14
15 # And plot
16 plt.plot(data['years [b2k]'], data['del18O [permil]'])
17 #plt.gca().invert_yaxis() # Grabs current axes and inverts x-axis
18 # uncomment if you prefer to see time this way.
19 # Automatically use column labels in the plot
20 plt.title(sheet_title, fontsize=16, fontweight='bold')
21 plt.xlabel(data.columns[0], fontsize=16)
22 plt.ylabel(data.columns[-1], fontsize=16)
23 plt.show()
24
25 # If we want better axis labels, we can use the LaTeX formatting:
26 plt.figure(figsize=(12,5)) # make figure extra-wide (width, height)
27 plt.plot(data['years [b2k]'], data['del18O [permil]'])
28 plt.title('Greenland Ice Core: Renland', fontsize=16, fontweight='bold')
29 plt.xlabel('Years b2k', fontsize=16)
30 # Below, I use LaTeX formatting to write the Greek "delta" and the superscript
31 plt.ylabel('δ18O', fontsize=16)
32 plt.tight_layout()
33 plt.show()

```

code/FilesData/GreenlandXLSXpandas.py

5 Stream gauging: an example

@TODO: Write new example for Monday with code from laptop that I wrote for the Mississippi

6 Plotting geospatial data using Basemap

You might want to do some mapping without starting a full-fledged GIS. For that, I have often used the Python mapping toolkit called basemap. This one cannot be installed via pip or easy_install: you will have to download it and run `python setup.py install`. For instructions, see the Basemap toolkit webpage at <http://matplotlib.org/basemap/users/installing.html>. You can make really gorgeous maps with it; see <http://matplotlib.org/basemap/users/examples.html>.

There is also a newer mapping toolkit called CartoPy. It *is* available via pip, and you can read its help files at <http://scitools.org.uk/cartopy/docs/latest/index.html>. It is supposed to be better than Basemap in some ways (see <http://ocefpaf.github.io/python4oceanographers/blog/2013/09/23/cartopy/>), and has a really nice looking map gallery with good code examples provided as well <http://>

scitools.org.uk/cartopy/docs/latest/gallery.html though as with many newer things, it is still in development.

Additional help

Additional support for numpy, matplotlib, and CartoPy can all be found in the SciTools GitHub repository at <https://github.com/SciTools/courses>.

Here, I return to the TRMM dataset that we used earlier to show a global plot that is properly projected.

```
2  #! /usr/bin/env python
3
4  from netCDF4 import Dataset
5  import matplotlib as mpl
6  from matplotlib import pyplot as plt
7  import numpy as np
8  import cartopy
9  import cartopy.crs as ccrs
10
11 # Open NetCDF file for reading
12 ncfile = Dataset('../data/3B43.20040601.7A.nc', 'r', format='NETCDF4')
13 lats = ncfile.variables['nlat'][:]
14 lons = ncfile.variables['nlon'][:]
15 precip = np.flipud(ncfile.variables['precipitation'][:].transpose())
16
17 # After this, we will set up the plotting domain
18 fig = plt.figure()
19 ax = fig.add_axes()
20
21 # Next, we will set up the projection and background map features
22 ax = plt.axes(projection=ccrs.Orthographic(100, 15))
23 ax.coastlines()
24 ax.set_global()
25 ax.gridlines()
26
27 # And then plot precipitation
28 im = ax.contourf(lons, lats, precip,
29                 transform=ccrs.PlateCarree(),
30                 cmap='spectral')
31 fig.colorbar(im)
32 ax.set_title('June 2004: mean ' + ncfile.variables['precipitation'].long_name
33             + ' [' + ncfile.variables['precipitation'].units + ']', \
34             fontweight='bold', fontsize=16)
35 plt.tight_layout()
36 plt.show()
```

code/FilesData/netCDFglobe.py

7 Plotting geospatial data using CartoPy

8 Python–GIS integration

8.1 GDAL and OGR

<https://pcjericks.github.io/py-gdalogr-cookbook/>

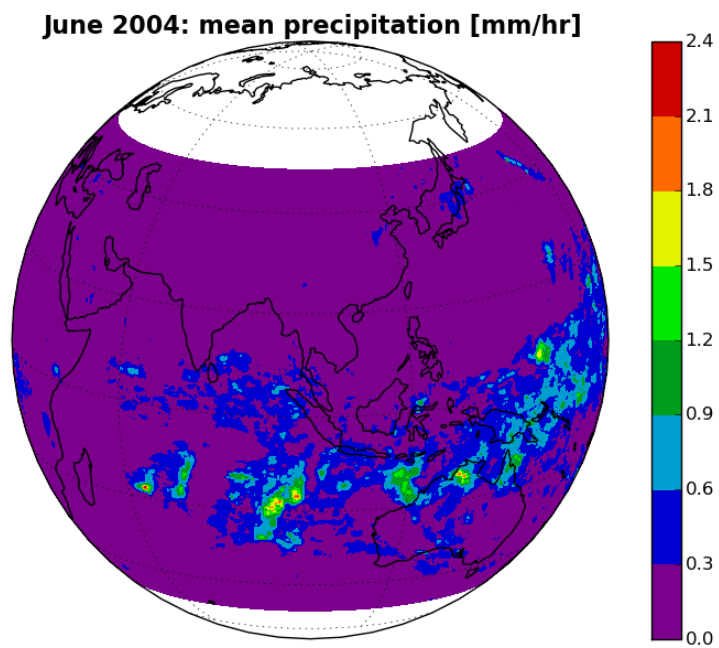


Figure 6: TRMM data from June 2004 (which is when the writer of these course notes graduated from high school), now nicely plotted on a real Earth. Yes, we're getting closer to being computational Earth scientists!