

Exercise: Getting data from files

One of the best ways to increase your effectiveness as a GIS programmer is to learn how to manipulate text-based information. GIS data is often collected and shared in text-based formats such as an Excel spreadsheet in CSV (comma-separated value) format, a list of coordinates in a text file, or an XML response received through a Web service.

When faced with these files, you should first understand if your GIS software already comes with a tool or script that can read or convert the data to a format it can use. If no tool or script exists, you'll need to do some programmatic work to read the file and separate out the pieces of text that you really need. This is called parsing the text. When you parse, you cycle through lines of text, treating them as strings, and pull out the useful information from those strings. In an XML file, for example, you may know that the information you want falls inside a particular tag, such as `<AvgTemp>46</AvgTemp>`. One approach to getting the value 46 might be to search for the line containing the substring "AvgTemp", then take all the characters that fall between the first > coming from the left of the string and the first < coming from the right.

In another case, you might know that the values you want fall after the second and third commas in a line of comma-separated values. You can split up the line based on comma locations and put all the segments of the string in a Python list. You can then take the third and fourth items in the list to get the values you want. (Remember the third and fourth items would come after the second and third commas, respectively.) In both cases, the keys to effective parsing are to know how to read lines in a file and know your string manipulation methods. It's helpful to know how to read a string, search for values in a string, split up a string based on some delimiter, and extract particular segments of a string.

Sometimes you can import helper modules, or libraries, into your code to make it easier to parse certain types of text. In this exercise, you will read a text file collected from a GPS unit. The lines in the file represent readings taken from the GPS unit as the user traveled along a path. The method used in this exercise is one of many that can be used to parse out the coordinates from each reading.

The file for this example is called `gps_track.txt`. The figure below shows a snapshot of what the contents of this file looks like. Open the file and examine its contents.

```
gps_track.txt - Notepad
File Edit Format View Help

type,ident,lat,long,y_proj,x_proj,new_seg,display,color,altitude,depth,temp,time,model,filename,ltime
TRACK,ACTIVE LOG,40.78966141,-77.85948515,4627251.76270444,1779451.21349775,True,False,255,358.228393554688,0,0,2008/06/11-14:08:30,etrex Venture, ,2008/06/11 09:08:30
TRACK,ACTIVE LOG,40.78963995,-77.85945952,4627248.40489401,1779446.18060893,False,False,255,358.228393554688,0,0,2008/06/11-14:09:43,etrex Venture, ,2008/06/11 09:09:43
TRACK,ACTIVE LOG,40.78961849,-77.85957098,4627245.69008772,1779444.78476531,False,False,255,357.747802734375,0,0,2008/06/11-14:09:44,etrex Venture, ,2008/06/11 09:09:44
TRACK,ACTIVE LOG,40.78953266,-77.85965681,4627234.83213242,1779439.20202706,False,False,255,353.421875,0,0,2008/06/11-14:10:18,etrex Venture, ,2008/06/11 09:10:18
TRACK,ACTIVE LOG,40.78957558,-77.85972118,4627238.65402635,1779432.89982442,False,False,255,356.786376953125,0,0,2008/06/11-14:11:57,etrex Venture, ,2008/06/11 09:11:57
TRACK,ACTIVE LOG,40.78968287,-77.85976410,4627249.79592111,1779427.14663093,False,False,255,354.383178710938,0,0,2008/06/11-14:12:18,etrex Venture, ,2008/06/11 09:12:18
TRACK,ACTIVE LOG,40.78979015,-77.85961390,4627264.19055204,1779437.76243578,False,False,255,351.499145507813,0,0,2008/06/11-14:12:50,etrex Venture, ,2008/06/11 09:12:50
TRACK,ACTIVE LOG,40.78983507,-77.85961390,4627268.97701256,1779436.91622722,False,False,255,357.747802734375,0,0,2008/06/11-14:14:24,etrex Venture, ,2008/06/11 09:14:24
TRACK,ACTIVE LOG,40.78987598,-77.85950661,4627275.37008691,1779445.16407113,False,False,255,360.151123046875,0,0,2008/06/11-14:14:35,etrex Venture, ,2008/06/11 09:14:35
TRACK,ACTIVE LOG,40.78996181,-77.85946369,4627285.58504744,1779447.10971155,False,False,255,357.747802734375,0,0,2008/06/11-14:14:47,etrex Venture, ,2008/06/11 09:14:47
TRACK,ACTIVE LOG,40.79011202,-77.85957098,4627300.72880585,1779435.05431356,False,False,255,350.057250976563,0,0,2008/06/11-14:15:10,etrex Venture, ,2008/06/11 09:15:10
TRACK,ACTIVE LOG,40.79015493,-77.85945952,4627305.83572582,1779436.02726889,False,False,255,353.902587890625,0,0,2008/06/11-14:15:45,etrex Venture, ,2008/06/11 09:15:45
TRACK,ACTIVE LOG,40.79015493,-77.85950661,4627306.47872977,1779439.66424433,False,False,255,357.747802734375,0,0,2008/06/11-14:16:21,etrex Venture, ,2008/06/11 09:16:21
TRACK,ACTIVE LOG,40.79021931,-77.85929203,4627316.87391712,1779456.58252508,False,False,255,352.460571289063,0,0,2008/06/11-14:16:47,etrex Venture, ,2008/06/11 09:16:47
TRACK,ACTIVE LOG,40.79024076,-77.85927057,4627319.58761408,1779457.97853529,False,False,255,352.941162109375,0,0,2008/06/11-14:16:48,etrex Venture, ,2008/06/11 09:16:48
TRACK,ACTIVE LOG,40.79036951,-77.85907745,4627336.83984196,1779471.8086655,False,False,255,352.460571289063,0,0,2008/06/11-14:17:17,etrex Venture, ,2008/06/11 09:17:17
TRACK,ACTIVE LOG,40.79036951,-77.85909891,4627336.51825492,1779469.98973889,False,False,255,343.327880859375,0,0,2008/06/11-14:18:13,etrex Venture, ,2008/06/11 09:18:13
TRACK,ACTIVE LOG,40.79047680,-77.85899162,4627350.09107746,1779476.96808367,False,False,255,342.847290039063,0,0,2008/06/11-14:18:32,etrex Venture, ,2008/06/11 09:18:32
TRACK,ACTIVE LOG,40.79064846,-77.85920620,4627366.01910842,1779455.39595965,False,False,255,346.212036132813,0,0,2008/06/11-14:18:59,etrex Venture, ,2008/06/11 09:18:59
TRACK,ACTIVE LOG,40.79069138,-77.85922766,4627370.48398108,1779452.73079612,False,False,255,345.731323242188,0,0,2008/06/11-14:19:42,etrex Venture, ,2008/06/11 09:19:42
TRACK,ACTIVE LOG,40.79079866,-77.85939932,4627379.87553205,1779436.06595964,False,False,255,344.77001953125,0,0,2008/06/11-14:20:14,etrex Venture, ,2008/06/11 09:20:14
TRACK,ACTIVE LOG,40.79011324,-77.85967827,4627399.62551803,1779408.19190874,False,False,255,346.692626953125,0,0,2008/06/11-14:20:42,etrex Venture, ,2008/06/11 09:20:42
TRACK,ACTIVE LOG,40.79101324,-77.85965681,4627399.9470929,1779410.01081578,False,False,255,346.692626953125,0,0,2008/06/11-14:21:16,etrex Venture, ,2008/06/11 09:21:16
TRACK,ACTIVE LOG,40.79118490,-77.85987139,4627415.87524858,1779388.43897893,False,False,255,347.173217773438,0,0,2008/06/11-14:21:47,etrex Venture, ,2008/06/11 09:21:47
TRACK,ACTIVE LOG,40.79133511,-77.86006451,4627429.7329001,1779369.10903715,False,False,255,347.65393064063,0,0,2008/06/11-14:22:05,etrex Venture, ,2008/06/11 09:22:05
TRACK,ACTIVE LOG,40.79161406,-77.85987139,4627463.7335061,1779379.97757661,False,False,255,344.77001953125,0,0,2008/06/11-14:22:33,etrex Venture, ,2008/06/11 09:22:33
TRACK,ACTIVE LOG,40.79172134,-77.85991430,4627475.0562103,1779374.22548807,False,False,255,343.80859375,0,0,2008/06/11-14:23:03,etrex Venture, ,2008/06/11 09:23:03
TRACK,ACTIVE LOG,40.79172134,-77.85991430,4627475.0562103,1779374.22548807,False,False,255,340.924682617188,0,0,2008/06/11-14:23:45,etrex Venture, ,2008/06/11 09:23:45
TRACK,ACTIVE LOG,40.79167843,-77.86004305,4627468.34161271,1779364.15903736,False,False,255,344.289428710938,0,0,2008/06/11-14:24:01,etrex Venture, ,2008/06/11 09:24:01
TRACK,ACTIVE LOG,40.79159260,-77.86015034,4627457.16214745,1779356.75766578,False,False,255,345.731323242188,0,0,2008/06/11-14:24:17,etrex Venture, ,2008/06/11 09:24:17
TRACK,ACTIVE LOG,40.79154968,-77.86016491,4627451.16052359,1779339.41750767,False,False,255,348.134643554688,0,0,2008/06/11-14:24:31,etrex Venture, ,2008/06/11 09:24:31
TRACK,ACTIVE LOG,40.79169989,-77.86055803,4627463.01825969,1779320.08774845,False,False,255,349.095825195313,0,0,2008/06/11-14:24:48,etrex Venture, ,2008/06/11 09:24:48
TRACK,ACTIVE LOG,40.79185009,-77.86075115,4627476.87491346,1779300.75825936,False,False,255,351.018432617188,0,0,2008/06/11-14:25:03,etrex Venture, ,2008/06/11 09:25:03
TRACK,ACTIVE LOG,40.79194146,-77.86081553,4627483.08881601,1779294.03255604,False,False,255,351.979858398438,0,0,2008/06/11-14:25:26,etrex Venture, ,2008/06/11 09:25:26
TRACK,ACTIVE LOG,40.79189301,-77.86083698,4627480.37530763,1779292.63741228,False,False,255,352.941162109375,0,0,2008/06/11-14:26:20,etrex Venture, ,2008/06/11 09:26:20
TRACK,ACTIVE LOG,40.79169989,-77.86060095,4627462.37514912,1779316.4499767,False,False,255,350.537841796875,0,0,2008/06/11-14:26:43,etrex Venture, ,2008/06/11 09:26:43
TRACK,ACTIVE LOG,40.79161406,-77.86045074,4627455.05411118,1779330.87351158,False,False,255,350.057250976563,0,0,2008/06/11-14:27:04,etrex Venture, ,2008/06/11 09:27:04
TRACK,ACTIVE LOG,40.79159260,-77.86045074,4627452.66088714,1779331.29660567,False,False,255,350.057250976563,0,0,2008/06/11-14:27:06,etrex Venture, ,2008/06/11 09:27:06
TRACK,ACTIVE LOG,40.79159556,-77.86012888,4627431.16046351,1779363.23028903,False,False,255,347.173217773438,0,0,2008/06/11-14:27:32,etrex Venture, ,2008/06/11 09:27:32
TRACK,ACTIVE LOG,40.79112053,-77.85982847,4627409.33982599,1779393.34591179,False,False,255,346.212036132813,0,0,2008/06/11-14:27:57,etrex Venture, ,2008/06/11 09:27:57
TRACK,ACTIVE LOG,40.79086304,-77.85950661,4627385.44746832,1779425.70289218,False,False,255,345.250610351563,0,0,2008/06/11-14:28:24,etrex Venture, ,2008/06/11 09:28:24
TRACK,ACTIVE LOG,40.79075575,-77.85935640,4627375.73353514,1779440.5498292,False,False,255,345.731323242188,0,0,2008/06/11-14:28:44,etrex Venture, ,2008/06/11 09:28:44
TRACK,ACTIVE LOG,40.79051371,-77.85905600,4627353.9163485,1779470.6652577,False,False,255,346.212036132813,0,0,2008/06/11-14:29:10,etrex Venture, ,2008/06/11 09:29:10
TRACK,ACTIVE LOG,40.79045534,-77.85945314,4627348.34087781,1779481.02821242,False,False,255,348.134643554688,0,0,2008/06/11-14:29:35,etrex Venture, ,2008/06/11 09:29:35
TRACK,ACTIVE LOG,40.79045534,-77.85945314,4627348.34087781,1779481.02821242,False,False,255,348.134643554688,0,0,2008/06/11-14:29:35,etrex Venture, ,2008/06/11 09:29:35
TRACK,ACTIVE LOG,40.79032660,-77.85877705,4627336.55616698,1779498.11631935,False,False,255,352.941162109375,0,0,2008/06/11-14:30:29,etrex Venture, ,2008/06/11 09:30:29
TRACK,ACTIVE LOG,40.79017639,-77.85838393,4627322.69873444,1779517.44674923,False,False,255,350.057250976563,0,0,2008/06/11-14:30:55,etrex Venture, ,2008/06/11 09:30:55
TRACK,ACTIVE LOG,40.79019785,-77.85862684,4627324.44891961,1779513.3865948,False,False,255,352.941162109375,0,0,2008/06/11-14:31:34,etrex Venture, ,2008/06/11 09:31:34
TRACK,ACTIVE LOG,40.79034805,-77.85881996,4627338.3052446,1779494.05637582,False,False,255,351.979858398438,0,0,2008/06/11-14:31:55,etrex Venture, ,2008/06/11 09:31:55
TRACK,ACTIVE LOG,40.79024076,-77.85903454,4627323.12460431,1779477.98422633,False,False,255,350.537841796875,0,0,2008/06/11-14:32:17,etrex Venture, ,2008/06/11 09:32:17
```

Notice that the file starts with a header line, explaining the meaning of the values contained in the readings from the GPS unit. Each subsequent line contains one reading. The goal in this exercise is to create a Python list containing the X,Y coordinates from each reading. Specifically, the script should be able to read the above file and print a text string like the one shown below.

```
[['-77.85948515', '40.78966141'], ['-77.85954952', '40.78963995'], ['-77.85957098', '40.78961849'], etc.]
```

Step 1: Plan approach for parsing the GPS track

Before you start parsing a file, it's helpful to outline what you're going to do and break up the task into manageable chunks. Here's some pseudocode for the approach we'll take in this example:

- 1. Open the file.
- 2. Read the header line.
- 3. Loop through the header line to find the index positions of the "lat" and "long" values.
- 4. Read the rest of the lines.
- 5. Split each line into a list of values, using the comma as a delimiter.
- 6. Find the values in the list that correspond to the lat and long coordinates and write them to a new list.

Step 2: Opening the file

The first thing the script needs to do is open the file. Python contains a built-in open() method for doing this. The parameters for this method are the path to the file and the mode in which you want to open the file (read, write, etc.). In this exercise you will only need to read the contents of the file and print them out to the user. "r" stands for read-only mode, "w" for write mode, etc.

```
# Reads a GPS-produced text file and writes the lat and long values
# to a list of coordinates
gpsTrack = open("C:\\Users\\Me\\Desktop\\GIS Programming\\Training\\Data\\gps_track.txt", "r")
```

Step 3: Reading the header line

Opening the file with the open() method gets you a file object (called gpsTrack in our case). You can read the first line by calling the file.readline() method, like this:

```
# Reads a GPS-produced text file and writes the lat and long values
# to a list of coordinates
gpsTrack = open("C:\\Users\\Me\\Desktop\\GIS Programming\\Training\\Data\\gps_track.txt", "r")

# Figure out position of lat and long in the header
headerLine = gpsTrack.readline()
```

This returns the string:

"type,ident,lat,long,y_proj,x_proj,new_seg,display,color,altitude,depth,temp,time,model,filename,ltime".

Step 3: Looping through the header line to find the index positions of the "lat" and "long" values

You need to search through this string and find the position of "lat" and "long". If you start counting comma-separated values in this string beginning from zero, it's easy to see that "lat" is at index position 2 and "long" is at index position 3. However, it's a good practice not to hard-code numbers like 2 and 3 into your script. Hard-coded numbers other than 0 or 1 are sometimes derided as *magic numbers*, suggesting that if you're not the programmer, you might have to use magic to know where the numbers came from!

Avoiding magic numbers gives you greater flexibility. If you wanted to re-use this script with a file in which "lat" and "long" were in different positions, you wouldn't have to modify your code. Even if "lat" and "long" went by some other name, it would be easier to find and change a string in your script instead of finding and changing a "magic number".

So how can you programmatically determine that "lat" is at index 2 and "long" is at index 3? Below is one way that uses the string.split() method. This method puts each "item" in the line into a list. The parameter you pass to the split method determines the delimiter, or the character that determines a new list item. In our case, it's the comma:

```
# Reads a GPS-produced text file and writes the lat and long values
# to a list of coordinates
gpsTrack = open("C:\\Users\\Me\\Desktop\\GIS Programming\\Training\\Data\\gps_track.txt", "r")

# Figure out position of lat and long in the header
headerLine = gpsTrack.readline()
valueList = headerLine.split(",")
```

The above method call returns:

"type,ident,lat,long,y_proj,x_proj,new_seg,display,color,altitude,depth,temp,time,model,filename,ltime".

The key is that now you can cycle through this list and discover the position of "lat" and "long". To do this, you could write a loop that searched through the list for "lat" and "long", but a quicker way is to use the helper method `index()` that gets you the index position of any item in the list:

```
# Reads a GPS-produced text file and writes the lat and long values
# to a list of coordinates
gpsTrack = open("C:\\Users\\Me\\Desktop\\GIS Programming\\Training\\Data\\gps_track.txt", "r")

# Figure out position of lat and long in the header
headerLine = gpsTrack.readline()
valueList = headerLine.split(",")

latValueIndex = valueList.index("lat")
lonValueIndex = valueList.index("long")
```

After running the above lines, `latValueIndex` is equal to 2 and `lonValueIndex` is equal to 3. With those variables set, you're now ready to start reading the rest of the lines in the file.

Step 4: Processing the rest of the lines in the file

When you have an open text file, you can always call `file.readline()` to go to the next line. In our case, we know we're going to use all the rest of the lines in the file, so it's more efficient to call `file.readlines()` to read them all at once. (This might not be efficient with an extremely long file.) The `readlines()` method returns a list of all the remaining lines in the file.

Now you can cycle through each GPS reading and split it up based on commas the same way you split up the header. You specifically need to pull out the values in index positions 2 and 3 of your list (represented by `latValueIndex` and `lonValueIndex`, respectively) and write those to a new list (`coordList`).

```
# Reads a GPS-produced text file and writes the lat and long values
# to a list of coordinates
gpsTrack = open("C:\\Users\\Me\\Desktop\\GIS Programming\\Training\\Data\\gps_track.txt", "r")

# Figure out position of lat and long in the header
headerLine = gpsTrack.readline()
valueList = headerLine.split(",")

latValueIndex = valueList.index("lat")
lonValueIndex = valueList.index("long")

# Read lines in the file and append to coordinate list
coordList = []

for line in gpsTrack.readlines():
    segmentedLine = line.split(",")
    coordList.append([segmentedLine[lonValueIndex], segmentedLine[latValueIndex]])
```

Notice that we first created the list and then placed values into it. “coordList” actually contains a bunch of small lists within a big list. Each small list is a coordinate pair representing the x (longitude) and y (latitude) location of one GPS reading. The list.append() method is used to add items to coordList. Notice again that you can append a list itself (representing the coordinate pair) using this method.

Step 5: Printing the coordinate values

```
# Reads a GPS-produced text file and writes the lat and long values
# to a list of coordinates
gpsTrack = open("C:\\Users\\Me\\Desktop\\GIS Programming\\Training\\Data\\gps_track.txt", "r")

# Figure out position of lat and long in the header
headerLine = gpsTrack.readline()
valueList = headerLine.split(",")

latValueIndex = valueList.index("lat")
lonValueIndex = valueList.index("long")

# Read lines in the file and append to coordinate list
coordList = []

for line in gpsTrack.readlines():
    segmentedLine = line.split(",")
    coordList.append([segmentedLine[lonValueIndex], segmentedLine[latValueIndex]])

print coordList
```

Step 6: Save and run script

Save the script as getGPScoords.py and run the script