

Introduction to Scientific Computing for Biologists: Final Project

Julie Szymaszek

21 March 2015

1 Introduction to the Problem

In the age of cheaply available sequencing, there is a wealth of data available through databases like NCBI's SRA. This is useful for exploratory projects, like rotations, that want to look at data without spending too much money on generating data. For my rotation last quarter, I was interested in parsing through publicly available RNAseq data in NCBI's DRA database in order to test hypotheses about the evolution of cell types.

However, all of these RNAseq data files have a lot of information. Simply put, the problem is that a lot of this information was not useful to the project I was working on. I obtained a set of adipose tissue RNAseq files that were processed using an open, web-based platform called Galaxy. The formatted RNAseq files were of varying kinds of adipose tissue from a variety of animals, though mostly from mammals. One of the first things I wanted to do with the files was find out what genes were "expressed" and compile a list of all of those in a single file for downstream applications.

The easiest way to do this would be to use the column in the file called "TPM" as a measure for expression of the transcripts reported. As a general rule for my project, I was using any TPM value of 2 or higher as a criteria for expression. I wanted to pull out every transcript that had a TPM value of 2 or higher. From there, I needed to know what the associated gene name is for that transcript ID, because the only things in that column are the names for transcripts presented in how the Ensembl database reports it. After getting the gene name, I'd like to know what genes are commonly expressed in all of these tissues.

Doing all of this by hand would take a long time because I'd have to first sort through one file by TPM number and only pull out transcripts. Then I'd have to look up the transcript names in order to figure out what the associated gene name for each one is. After I did that for every single, I'd have to use a venn diagram calculator to find out what genes all of these samples had in common. The quarter is only 10 weeks long, and this was only the first step in a series of things I was supposed to accomplish. Once I know what genes I'm interested in - I can look up whether they are of interest to me, and start planning experiments with cell cultures in the lab.

2 One Possible Solution

Since I have a file that contains every single possible transcript and their associated names from every single species I am working with, I want to use this as a dictionary to lookup the names of genes that are expressed, and use that output to find all of the genes in common.

I wrote two separate scripts which are both fold in the main directory called Final Project. The first script entitled main.py extracts ENSEMBL Transcript ID for TPM values greater than 2 from all of the aligned RNAseq files present in the specified subdirectory called transcript reads. After extracting all of the transcript read file, the script identifies the "locus" in the files, which is the transcript ID. Those entries with a TPM value greater than 2 (an arbitrarily selected cutoff value that would mean this transcript is "expressed" in the tissue) are exported into a list. Should I decide that there is reason to use another TPM cutoff value, I can update the script (for example, Kin et al, 2015 used a TPM cutoff of 3 [1]). The locus is compared to the giant file I have of every possible transcript and its gene name, and pulls out the gene name. The gene name output is saved into individual files for each tissue sample as "output original file name.txt." These are all saved into the subdirectory entitled "associated genes."

The second script I wrote called associated genes.py can then be called from the command line. This script takes all of the files in the subdirectory associated genes and finds the intersection of its contents in sequence. It saves its output into the main folder in a file called intersection genes.txt.

3 Code

Before I can access the code, I have to clone my directory from github, using the following command in terminal:

```
1  
2 git clone https://github.com/jszymaszek/scicom-project.git
```

This directory contains the two necessary scripts, and the subdirectory that contains all of the transcript read data files that will be necessary to complete the analyses and run the two scripts successfully.

The first script, used to lookup the gene names of all of the transcripts that are expressed in the different samples, is presented below:

```
1 #!/usr/bin/python  
2  
3 import os  
4 #this imports a module that lets me read all of the files ↪  
   ↪ within a specific directory later in this code  
5  
6 def build_dict():  
7     dict = {}  
8     with open("associated_gene_names.txt") as file:  
9         #build a dictionary using the associated_gene_names.txt ↪  
           ↪ file as its input  
10         for line in file:  
11             arr = line.split("\t")  
12             if len(arr) > 1:  
13                 locus = arr[0].lower()  
14                 associated_gene = arr[1].lower()  
15                 dict[locus] = associated_gene  
16     return dict  
17 #whatever values i'm going to want to call from this ↪  
   ↪ dictionary will look up in the locus column  
18 #then it will just spit out whatever value is associated ↪  
   ↪ with that locus column entry from the associated_gene ↪
```

```

    ↪ column
19
20 for root, dirs, files in ↪
    ↪ os.walk("./transcript_reads",topdown=True):
21 #reading everything inside the folder transcript_reads
22 dict = build_dict()
23 #creates a dictionary that is comprised of the build_dict() ↪
    ↪ from above
24 with open("garbage.txt","w") as garbagefile:
25     garbagefile.truncate()
26 #creates a file called garbage.txt where some errors will go
27 #truncate clears out the file every time it's running so ↪
    ↪ that way there aren't duplicates
28 for name in files:
29     filename = os.path.join(root, name)
30     with open(filename) as file:
31         print "-----"+filename+"-----"
32         lines = file.readlines()
33         column_names = lines[0].strip().split("\t")
34         locus_index = column_names.index("locus")
35         TPM_index = column_names.index("TPM")
36         outfilename = ↪
            ↪ os.path.join("./associated_genes","output_"+name)
37 #within each file that is opened it takes the locus and the ↪
    ↪ TPM columns from the files within transcript_reads
38 #outfilename section is specifying that the output files ↪
    ↪ from this script will go here with each file being ↪
    ↪ called
39 #output_name of the original file.txt
40     with open(outfilename,"w") as outfile:
41         for i,line in enumerate(lines):
42             if i !=0 and len(lines[i].strip().split("\t")) ↪
                ↪ >= max(locus_index,TPM_index):
43 #ignore line 0 (column names), makes sure that the index ↪
    ↪ being accessed is within the range of the line
44                 locus = lines[i].strip() ↪
                    ↪ .split("\t")[locus_index].lower()
45                 TPM = float(lines[i].strip() ↪
                    ↪ .split("\t")[TPM_index])

```

```

46 #parse the locus and the TPM columns
47     if locus in dict:
48         if TPM > 2 and dict[locus].strip() != "":
49             associated_gene = dict[locus]
50             outfile.write(associated_gene)
51 #look at all of the corresponding locus entries and pull ↪
    ↪ out the locus name is the TPM value is greater than 2
52 #then use that locus name to lookup the associated gene ↪
    ↪ name in the dictionary function written above
53 #print out the gene name of that locus with a value into ↪
    ↪ the file output_name of the file.txt
54     else:
55         with open("garbage.txt","a") as garbagefile:
56             garbagefile.write("~~filename:")
57             garbagefile.write(filename)
58             garbagefile.write(" ~~locus:")
59             garbagefile.write(locus)
60             garbagefile.write("\n")
61 #however, sometimes the locus in the single files might not ↪
    ↪ have a corresponding gene name in the dictionary
62 #if this is the case, print them out into this file called ↪
    ↪ garbage.txt
63 #this file also includes the name of the file and the locus ↪
    ↪ name so i can later see what "failed"
64
65
66     '''print "---dirs---"
67     for name in dirs:
68         print root
69         print name
70         print(os.path.join(root, name))'''
71 print "END"

```

Now, it's time to use the second script to find what genes all of these samples have in common:

```

1  #!/usr/bin/python -tt
2  # -*- coding: utf-8 -*-
3
4  import os
5  #imports operating system module to be able to read the ↵
   ↵ list of files within a directory
6
7  def build_list(filename):
8      list = []
9      file = open(filename, 'r')
10     for line in file:
11         line = line.strip()
12         list.append(line)
13     file.close()
14     return list
15
16 #this builds a function called build_list that takes the ↵
   ↵ input from the directory associated_genes
17 #it builds a set with them
18
19 def get_intersect(listoffilenames):
20     flag = 1
21     for filename in listoffilenames:
22         print filename
23         if flag == 1:
24             s = set(build_list(filename))
25             flag = 0
26             s &= set(build_list(filename))
27         print '\n'.join(s)
28     return s
29
30 #this builds a function called get_intersect of all of the ↵
   ↵ filenames set
31 #starts with a file and it takes the intersection of it ↵
   ↵ with the next file and so forth
32 #basically takes the intersection of the files in sequence, ↵
   ↵ and joins them together
33

```

```

34 def build_output(listoffilenames,outputfilename):
35     file = open(outputfilename,'w')
36     intersect = get_intersect(listoffilenames)
37     for i in intersect:
38         file.write('%s\n' % i)
39     file.close
40 #taking the output of the get_intersect function and writes ↪
    ↪ it to file
41
42 if __name__ == "__main__":
43     listoffilenames = []
44     for root, dirs, files in ↪
        ↪ os.walk("./associated_genes",topdown=True):
45         for name in files:
46             filename = os.path.join(root, name)
47             listoffilenames.append(filename)
48     build_output(listoffilenames,'intersection_genes.txt')
49 #this uses the main function and os to take all of the ↪
    ↪ files in the folder "associated_genes"
50 #and uses the list of filenames as the input for this script
51 #then it builds the output of this script into a file ↪
    ↪ called "intersection of expressed genes

```

4 Future improvements

As I worked through this problem, I thought of a couple of ways in which this process could be streamlined further. For example, instead of executing two scripts, I could "join" them together using the main function in order to use execute the second script as a module of the first one. In the future, I may want to run both simultaneously, so this would be useful to do.

Another possibility area for improvement would be to use a package like Mechanize in order to obtain all of the transcript/associated gene name information I have contained in my dictionary file. I already had made this file a while back, so I decided to use it here, but all of the data comes from Ensembl's BioMart website. I could use a package like Mechanize to write a script that would access the information I needed, so that I wouldn't be limited to just the species I have in my dictionary.

References

- [1] Koryu Kin, Mauris C Nnamani, Vincent J Lynch, Elias Michaelides, and Gunter P Wagner. Cell-type phylogenetics and the origin of endometrial stromal cells. *Cell reports*, 10(8):1398–1409, 2015.