

Assignment 3: Finding Genes

Gene finder: All of your code for this problem should be in a file called `GeneFinder.py`. At the end you will also make a file with comments called `GeneFinder.txt`. This assignment is taken (with minor adaptations) from the companion material for CFB. Find the originals at <https://www.cs.hmc.edu/twiki/bin/view/CFB/GeneFinder>.

Introduction

In the text of CFB we discussed some investigations of *Salmonella* pathogenesis. Researchers identified genomic regions present in *Salmonella* but not in closely related species such as *E. coli*. Further study revealed that some of these regions are essential to *Salmonella*-specific pathogenic processes. One particular region was found to be needed for the entry of the bacterium into the epithelial cells of the gut. The researchers studying it wished to identify the protein-coding genes in it, so they sequenced it and analyzed that sequence for genes. In this homework, you will develop a simple gene finder, and use it to examine this sequence. You will also examine a related "mystery" sequence, and try to determine its function.

You can begin by creating a homework file called `GeneFinder.py`. You will make use of a number of functions we have already written. You can give yourself access to the functions in your `load.py` and `dna.py` files by making sure they are present in the same directory as `GeneFinder.py`. Then include this at the top of `GeneFinder.py`:

`GeneFinder.py`

```
1 from load import *  
2 from dna import *
```

You should find a file called `X73525.fa` in the scripts folder for this assignment. Make sure your `GeneFinder.py` file is in the same directory.

A simple gene-finding strategy

A simple gene-finding strategy is to look for large open reading frames. As we've seen, an open reading frame (ORF) is the stretch of sequence between a start codon and the next in frame stop codon. This strategy depends on the assumption that protein-coding sequence has larger ORFs than noncoding sequence. In noncoding regions "Start" and "Stop" codons appear by random chance, and the ORFs between them are not long. In many circumstances protein-coding genes have ORFs that are longer than this.

To look for genes in a particular sequence, you can identify its longest ORFs. But how do you know if these are genes? To help decide, you can look at many known noncoding regions to obtain a distribution of open reading frame lengths for noncoding DNA. You can then compare the longest ORFs in our test sequence with this distribution. If the test sequence ORFs are longer, they are likely to be genes.

In this problem the approach you will take is the following. You will compare open reading frames in our *Salmonella* sequence with open reading frame lengths from known noncoding sequences. How do you get known noncoding sequences? One way is to make them randomly. Here you will do this by randomly shuffling (reordering) our genomic sequences. More on that shortly!

Part 1: Retrofitting functions from `orf.py`

For your `orf.py` homework you created three functions. The first of these, `restOfORF(DNA)`, can be used as written. Copy and paste this into your present homework file.

Next you will write a modified version of `oneFrame` called `oneFrameV2(DNA)`. In our previous version of `oneFrame`, we searched the given reading frame for every start codon, and returned the ORF corresponding to it. This sometimes results in nested ORFs, and can be seen in the following example:

example

```
1 >>> oneFrame("ATGCCCATGGGAAATTTTGACCC")
2 ['ATGCCCATGGGAAATTT', 'ATGGGAAATTT']
```

Both ORFs end in the same stop codon. The second ORF is shorter, because it begins with an ATG which is internal to the first ORF. For our gene-finder application it is desirable to instead only return the largest ORF from a set of nested ORFs like this.

GeneFinder.py

```
1 >>> oneFrameV2("ATGCCCATGGGAAATTTTGACCC")
2 ['ATGCCCATGGGAAATTT']
```

A natural approach to writing `oneFrameV2` is to search through the string till you find the first ATG. Then call `restOfORF` to get the ORF. But when you resume searching for the next ATG, you skip ahead and start looking after the end of the ORF you just found. A while loop will be very convenient here!

Finally, create a `longestORFV2(DNA)` function which calls `oneFrameV2`.

Part 2: longestORFBothStrands(DNA)

We are given a DNA sequence in 5' to 3' order. A gene might appear on this strand or its reverse complement. Thus, our next function is a very short one called `longestORFBothStrands(DNA)`. This function takes a DNA string as input and finds the longest ORF on that DNA string or its reverse complement. This won't be hard! You can use the `longestORFV2` function you have already written. First ask it for the longest ORF in the given DNA and then ask it for the longest ORF on its reverse complement (use your `reverseComplement` function to find the reverse complement). The longer of those two is the longest ORF possible (break ties arbitrarily).

For example,

example

```
1 >>> longestORFBothStrands('CTATTTTCATG')
2 'ATGAAA'
```

Pause here to make sure you understand how we got that answer.

Part 3: longestORFNoncoding(DNA, numReps)

To assess whether long ORFs are genes, a researcher would ask the question, “is this sequence length indicative of a coding region or would I expect to see sequences this long in garbage?” By “garbage”, of course, we mean just random sequences of nucleotides.

We'll test this by generating a bunch of “garbage” sequences of the same length as our test DNA sequence, and measuring the maximum ORF length in each. Then, we'll ask the following question. Is the very longest ORF among these still shorter than some ORFs we observe in our real DNA? If the real DNA ORFs are significantly longer than what we see in the garbage sequence, that is a very strong indicator that we did in fact find genes in our original DNA!

Our first task is to write a function `longestORFNoncoding(DNA, numReps)` that makes a bunch of garbage sequences, finds the very longest ORF in all of these, and returns its length. Note: this function returns a number rather than a DNA string.

OK, so now it's time to generate garbage sequences. We could generate totally random strings of the same length as our DNA string, but that might not be a very accurate test since our DNA string might have more nucleotides of one type and fewer of another. To be fair, our garbage strings should have the same nucleotides but just reordered or “shuffled” randomly.

To do that, you'll first need to take your DNA string and turn it into a list of its constituent symbols. There is a built-in function called `list` that takes as input a string and returns the list of symbols in that string.

example

```
1 >>> myList = list("hello")
2 >>> myList
3 ['h', 'e', 'l', 'l', 'o']
```

Now, in the random package (get it by including `import random` at the top of your file) there is a function called `shuffle` that scrambles the list. Unlike other functions that we've used until now, this `shuffle` function doesn't return a new list, but rather it actually changes the list that you gave it. Here's an example:

example

```
1 >>> import random
2 >>> myList = list("hello")
3 >>> myList
4 ['h', 'e', 'l', 'l', 'o']
5 >>> random.shuffle(myList)
6 >>> myList           # <-- myList will be changed now!
7 ['o', 'l', 'h', 'e', 'l']
```

In other words, do not do this in your code...

example

```
1 myList = random.shuffle(myList)
```

but instead do just this...

example

```
1 random.shuffle(myList)  # this will actually change myList by scrambling its contents
```

Oh, wait! But now we want that shuffled list to be glued back together as a string, since after all we're dealing with DNA strings. That's easy, here's a function that takes as input a list of symbols and returns back the string that we get by gluing those symbols together. Take a minute to see how it works. It's really short and sweet.

Then, copy it into your `GeneFinder.py` file and try it out!

example

```
1 def collapse(L):
2     output = ""          # This is our initial output string
3     for s in L:          # for each string in the list...
4         output = output + s    ... construct a new output string
5     return output        # ... and return the final output string
```

Here's the function in action:

example

```
1 >>> collapse(['o', 'l', 'h', 'e', 'l'])
2 'olhel'
```

For each garbage sequence you make, calculate the longest ORF using `longestORFBothStrands`. You should repeat this process `numReps` times, and return a number indicating the length of the longest ORF you see in all those repetitions.

You can test your `longestORFNoncoding` function on our real *Salmonella* sequence.

example

```
1 >>> X73525=loadSeq("X73525.fa")
```

Now run `longestORFNoncoding` a few times. Note that because it makes use of randomness, it will not give you exactly the same number each time. It will, however, be consistent enough for our purposes.

example

```
1 >>> longestORFNoncoding(X73525,50)
2 624
3 >>> longestORFNoncoding(X73525,50)
4 693)
```

Part 4: findORFs

Our next step is to write a function `findORFs(DNA)` that will identify all the ORFs in the real (unshuffled) DNA and return them as a list. If there are none, it should return an empty list.

Once again, this task is easy because we can make use of functions we have already written. `findORFs` should call `oneFrameV2` in each of the three possible reading frames of the sequence. It should then combine all of the ORFs found in each frame and return them. Note that this strategy makes it possible to find overlapping reading frames. `findORFs` will likely be very similar to the `longestORF` that you wrote above.

example

```
1 >>> findORFs("ATGGGATGAATTAACCATGCCCTAA")
2 ['ATGGGA', 'ATGCCC', 'ATGAAT']
3 >>> findORFs("GGAGTAAGGGGG")
4 []
```

Part 5: findORFsBothStrands

Next write a function called `findORFsBothStrands(DNA)` that searches both the forward and reverse complement strands for ORFs and returns a list with all the ORFs found. For example:

example

```
1 >>> findORFsBothStrands('ATGAAACAT')
2 ['ATGAAACAT', 'ATGTTTCAT']
```

Part 6: getCoordinates

The functions we've written so far return the sequence of an ORF. However another bit of information we might like to know is its coordinates in the original DNA sequence. We'll now write a function `getCoordinates(orf,DNA)` which returns the beginning and end coordinates of an ORF in DNA. We will follow the convention of reporting everything in forward strand coordinates (even if the ORF is actually on the reverse complemented strand).

First consider the case where the ORF falls on the forward (non reverse complemented) strand. In that case, we can obtain its beginning coordinate with the following syntax:

example

```
1 >>> testDNA="ACGTTCTGA"
2 >>> testORF="GTT"
3 >>> testDNA.find(testORF)
4 2
```

We can then get the end coordinate by adding the length of the ORF to the start coordinate.

But what if the ORF is actually on the reverse complemented strand? Notice what happens when we use `.find` with a sequence that is not present:

```
example
1 >>> testDNA="ACGTTCTGA"
2 >>> testORF="GAA"
3 >>> testDNA.find(testORF)
4 -1
```

The method returns a -1. If we search the forward strand of DNA with an ORF and get a -1, then we should try searching with the reverse complement of orf. In our example, that would look like this:

```
example
1 >>> testDNA="ACGTTCTGA"
2 >>> revCompTestORF=reverseComplement("GAA")
3 >>> testDNA.find(revCompTestORF)
4 3
```

Here are some examples of `getCoordinates`:

```
example
1 >>> getCoordinates("GTT", "ACGTTCTGA")
2 [2, 5]
3 >>> getCoordinates("CGAA", "ACGTTCTGA")
4 [3, 7]
```

Part 7: Gene finding at last!

Write a function called `geneFinder(DNA, minLen)` that identifies ORFs longer than `minLen`, and returns a list with information about each.

`geneFinder` should first call `findORFsBothStrands` to obtain a list of ORFs in the input DNA. It should then run through this list, keeping only those ORFs which are longer than `minLen`.

For each ORF which is long enough, `geneFinder` should calculate:

- The beginning and end positions of the ORF in DNA using `getCoordinates`
- The protein sequence of the ORF using `codingStrandToAA`

These should then be placed in a list:

```
example
1 [beginningCoord, endCoord, proteinSequence]
```

There will be a list like this for every ORF that is long enough. You will collect these lists in another list (a list of lists). A final step is to sort this list of lists before returning it.

Say our final list of lists is called `finalOutputList`. Its elements are lists of this type: `[beginningCoord, endCoord, proteinSequence]`. We can sort `finalOutputList` by beginning coordinate, and return it like this:

```
example
1 finalOutputList.sort()
2 return finalOutputList
```

Part 8: Using our gene finder

Now its time to apply our method to data and see what we get.

1. Load X73525.fa into Python

example

```
1 >>> X73525 = loadSeq("X73525.fa")
```

and apply our gene-finding strategy to it. The first step is to decide on a threshold for `geneFinder`. Our strategy is to determine what the longest ORF we see in noncoding sequence is, and then to define ORFs longer than this as putative genes. Run `longestORFNoncoding(X73525,1500)`. This should be a relatively conservative way to pick a threshold. Now run `geneFinder` using the threshold value you get for `minLen`.

example

```
1 >>> geneList = geneFinder(X73525, put\_your\_minLen\_value\_here )
```

Next write a short function `printGenes(geneList)` which prints the output of `geneFinder` in a nice human-readable form. Use it to print your results. Note that our gene-finding strategy is very simple, and is also relatively conservative. As a result, we are likely to miss some true genes which are short. But we hope that the genes our method does find are real ones. Paste your `printGenes` output into a file called `GeneFinder.txt`.

2. Pick one of the genes from your output and blast its protein sequence. To start, open this link in a separate window: NCBI Blast, <https://blast.ncbi.nlm.nih.gov/Blast.cgi>. (On a mac, you can do that by holding down the “control” key while clicking on the link. You’ll then have the option to open the link in a new window.)

Go down to Basic Blast and follow the protein blast link. There should be a large box into which you can now input a protein sequence. Then scroll down to the bottom of the page and click the “Blast” button. When you blast a sequence, the application carries out a search through its databases for known sequences that either match your sequence or are close to it, and returns those to you in a list.

The BLAST output page will contain some graphics at the top. Scroll down past these to the section labelled “Descriptions”. The first links in this section will be the closest sequences BLAST could find to your original input. Clicking on one of these will take you to a page with information about that sequence, including the name and what organism it can be found in. The right pane on that page has many useful links. Based on what you learn from these top BLAST hits, briefly describe the likely function of your gene at the top of `GeneFinder.txt`.

3. The famous biologist, Professor P.I. Pette, has been studying *Salmonella* in her lab. She has obtained the following sequence from a novel pathogenic strain: `sa1DNA.fa` (this file is in your scripts directory). Prof. Pette obtained this from a region unique to Salmonella. She believes it is involved in pathogenesis, and may be from pathogenicity island SPI1. Download this file and place it in the same directory as `GeneFinder.py`.

Then use your code to help Prof Pette determine the function of this sequence. Identify any putative genes, and then blast them at NCBI BLAST. Note you will need to run `longestORFNoncoding` on this sequence to determine a good threshold for it. Using the links which NCBI BLAST gives you determine the name and function of your mystery gene(s). In `GeneFinder.txt` write a paragraph explaining what you have found.