

Summer Research School in Computational Biology and Bioinformatics - coding challenge.

This test is to assess student applicants for the Summer Research School in Computational Biology and Bioinformatics. The test questions and examples are framed in the Julia language (<https://julialang.org/>), and we prefer solutions that use Julia. However, we understand that some of you may be too busy to learn a new language, so we will also accept solutions in Python, C++, or Rust. The preference for Julia may be a factor when we're making the assessment, however.

Please submit your work through the github classroom system. If you found this document, it means you already have a git repository with these questions and a few other files. We can access what you upload to this repository, which is how we will assess your solutions.

Note: by submitting, your solutions (including your name and email) will be stored on GitHub, and you consent to us accessing and processing your submission for the purpose of selecting applicants for the Summer Research School. If you object to uploading your solutions, name, or email address to GitHub, please contact the organizers and we can discuss this, and possibly make an alternative arrangement for your submission.

General rules:

- Answer these questions without help.
- Answer as many questions as you can. You do not have to answer them all to submit. Some of these questions are difficult, and we do not necessarily expect anyone to solve them all, so please submit whatever you do manage to solve!
- You are allowed to use or adapt code that you find on the internet, but please do not ask for help online, or post these questions online in any form.

Steps:

1. Install the latest version of Julia (<https://julialang.org/>), along with IJulia and Jupyter Notebooks (instructions here: <https://github.com/JuliaLang/IJulia.jl>). Check that you can open a Jupyter notebook and run some Julia commands. If you get stuck on any of these installation steps, Google will prove very helpful.
2. Clone or Download the initial challenge repository from GitHub *to your computer*.
3. Open the “FunctionTesting.ipynb” Jupyter notebook *on your computer*, and run the code that will test your solutions. This tests each of the functions in the “functions.jl” file. Initially, you should get messages that all of them fail.
4. For each question, fill in the code for the corresponding function in “functions.jl”. You can re-run the tests in “FunctionTesting.ipynb” to see how they perform on a single test case.
5. **Critical:** When you’re content with your solutions, or feel like you won’t be able to get any further, please update the “functions.jl” in your **online** GitHub repository so that it contains your solutions (don’t just change the file locally on your computer!). If you are familiar with Git, you can use a “push” command, but you can also simply upload your modified “functions.jl” file using the “Upload files” button through the web browser interface while viewing your repository on GitHub. Note, your repository URL should look like this: <https://github.com/MurrellGroup/ki-comp-bio-coding-challenge-2023-YourGitHubUsername> Please check that “functions.jl” at this address (but modified for your username) contains the version with your solutions.

Other considerations:

- Avoid putting any print/println statements in the final functions you submit, as these will make the test output more difficult to read.
- For questions you do not solve, just leave the relevant function as is - no need to delete it.
- When we assess you, we will download your repository and run your "functions.jl" against a more comprehensive set of tests, on many different cases, sometimes including timing of the solutions.
- Do not modify "test.jl".
- For the trickier algorithmic problems, optimal solutions are nice, but don't let that stop you from submitting. We aren't quite sure how hard these are, nor whether anyone will provide optimal solutions to all of the problems, so please just submit the best solution you can find for each of the problems, or leave that problem blank.

Questions:

Q0) Literally just your name and email (rating: starter):

Make the Q0 function in “functions.jl” return your name and email address, so we can contact you if you are selected. **Note:** *If you get this one wrong, the rest won't matter.*

Q1) GC content (rating: starter):

Count the proportion (as a percent) of bases in a DNA sequence that are either G or C.
Example input:

"GCATGCACATAGCAGCGAGCTACTACATCGCGGCTAGACTACTGAGCGA"

Q2) Translation (rating: starter):

In DNA, 3 successive nucleotides forms a “codon” (https://en.wikipedia.org/wiki/Genetic_code), where each codon codes for an Amino Acid. Write a function that takes a DNA string, and returns the corresponding amino acid string.

Q3) A magically annoying coin (rating: intermediate):

You are given a coin that can be flipped up to K times. On the n^{th} flip, the probability of landing heads (H) is n/K , and tails (T) is $1 - (n/K)$. Write a function that computes the logarithm of the probability of observing a sequence of flips, such as:

TTTTTTTTTTTTTTTTHTTTTTHHTTTTTTHTTHTTTHNTTTNHTHTTNNHHHHHHHHHHHNTTTHNHTHTNHNHTHTTTNHHHHHTTNNHNHTNN

(What we might ask if you get called for an interview: Do you think that this example sequence above was generated by the magical coin? Why?)

[Continued on next page]

Q4) The biggest product (rating: tricky):

Consider an array of strictly positive numbers, each greater than zero. For any contiguous subsequence of such numbers, you can compute their product. Write a function that, in the shortest compute time, calculates the starting point and ending point of a contiguous subsequence of a strictly positive array that has the largest product.

For example:

```
A = [1.6, 0.56, 1.3, 1.5, 0.9, 1.5, 2.5, 1.1, 0.46, 0.65]
```

The contiguous subsequence with the greatest product is from index 3 to index 8, with a product of 7.239375. On this example, the function should return: (3, 8)

Note: The speed of your implementation, and how it scales to large input arrays, will be considered in the assessment. Will your function be able to run on arrays with 10,000 elements in a reasonable amount of time? How about 1 million or 1 billion elements?

Q5) Data dumpster diving (rating: wat):

The following .csv file contains a symmetric 2D array of pairwise distances between 1722 elements:

<https://drive.google.com/open?id=1MDv1UViwPUoeLJfQ1XH1jvAr-nAzDTTX>

What sort of thing are these elements? What does the element at position 430 (ie. the 430th row or column of the distance matrix) correspond to? How did you figure this out?

Note: This should be answered by returning the answer as text in the relevant function. This will not be included in the automatic assessment, but will be checked by a human.

Hint: https://en.wikipedia.org/wiki/Multidimensional_scaling

Q6) Array yeet (rating: difficult):

You have two arrays of numeric values, not necessarily of the same length. Consider a set of deletions in each array, resulting in the arrays being the same length. For every possible set of deletions, you can compute a score, which is the sum of the element-wise squared differences, plus the number of elements you had to delete (across both arrays). Write a function that, in the shortest possible compute time, decides which elements to delete in both arrays to minimize this score. For example, given:

```
A = [-5, 9, -3, -2, -9, -2, 2, 11, -1, -4, -10, 21]
```

```
B = [-4, -4, -2, -9, -2, 8, 3, 10, -2, -3, -9]
```

Deleting the three elements in bold gives you:

```
A = [-5, -3, -2, -9, -2, 2, 11, -1, -4, -10]
```

```
B = [-4, -4, -2, -9, -2, 3, 10, -2, -3, -9]
```

And taking the total element-wise squared difference (in Julia syntax): `sum((A .- B).^2)`

Gives you 7, plus 3 for the 3 deleted elements, gives you a total score of 10. Your task is to find the deletions, if any, that would minimize this score, for any two numeric arrays. Return a tuple of the optimal score, the position of the deletions in the first input array (an array of integers), and the position of the deletions in the second array (also an array of integers). Return value for the above example:

```
(10.0, [2,12], [6])
```

Note: The speed of your implementation, and how it scales to large input arrays, will be considered in the assessment.

Q7) Visualization

A collaborating biologist gives you a file they found in Newick format (https://en.wikipedia.org/wiki/Newick_format) and asks you to make them a slide to "visualize the data". You ask them more questions, but they don't answer with any additional clarity. They mumble something about "maybe the format isn't completely standard - there might be extra stuff in there".

Produce a visualization for them that will fit on a single slide (PDF preferred), providing as much insight into the data as possible. You will be rated on the informativeness and the aesthetics of the visualization.

The Newick file can be found here:

<https://drive.google.com/file/d/1BMjtNSXA39rgHcOEs3LPt9u2c8yw-0FC>

Instead of an answer in the code, just commit the PDF directly to the classroom as Q7.pdf.

Note: one starting point for producing figures is Julia Plots (<https://docs.juliaplots.org/stable/>)

Q8) Markov's path (rating: difficult)

You have a set of N reference sequences, indexed from 1 to N , with each sequence being a string consisting of the characters A, T, G, and C. Each sequence has a fixed length of L . A path P is a list of length L with integers representing which reference is chosen for each index in the path. The path undergoes a "switch" at index t whenever $P_t \neq P_{t+1}$. The following is an example set of reference sequences where N is equal to 3 and L is equal to 5:

1. [A, C, T, C, A]
2. [A, C, G, T, A]
3. [C, G, C, G, C]

Further, any path emits a sequence corresponding to the characters in the position of the reference sequence indicated by the path. For example, the path $[2, 2, 2, 2, 2]$ would mean that reference 2 is chosen for all indices. This path does not have any switches, and emits the sequence [A, C, G, T, A]. The path $[3, 3, 1, 1, 1]$, would mean that you are first at reference 3, but you swap, once, to reference 1. So this emits the sequence [C, G, T, C, A], with one switch from 3 to 1.

You have an N -by- N switching cost matrix C , where C_{ij} is the cost of switching from reference i to reference j . The switching cost of an entire path is the sum of the costs of all of the individual switches. The "total score" for a path is the number of mismatches between a query sequence Q and the sequence emitted by the path, plus the switching cost of the path. The "minimal switches", which is what you will need to calculate, is the number of switches that minimizes the total score. If there are different numbers of switches that tie for the minimal total score, return the lower of these.

Your friend Markov gives you a query sequence Q and a switching cost matrix of the form:

- If $i = j$ (ie. no switch), then $C_{ij} = 0$
- If $i \neq j$, then $C_{ij} = K$.

So Markov provides you with Q and the value of K . You calculate the minimal switches.

Note: The speed of your implementation, and how it scales to large input arrays, will be considered in the assessment.

Q9) Barcodes (rating: unknown):

You are asked to design a set of DNA barcodes. A DNA barcode is made up of a sequence of the letters A,C,G and T, of a particular length. For example, a DNA barcode of length 7 could be "TAGCTAG". *All barcodes in a set must have the same length.* You will have to generate a set of barcodes with the following constraints:

- 1) In a single barcode, no letter can be immediately repeated (ie. "TAGGTAG" is not allowed).
- 2) The minimum hamming distance (https://en.wikipedia.org/wiki/Hamming_distance) between any pair of barcodes must be greater than or equal to some threshold, D. For example, here is a set of 6 barcodes of length 4, where D = 3:

CGTA
TCTC
TACG
GCAT
AGAG
ATGC

Your task is to write a function that takes, as input, the number of barcodes required, and the minimum distance threshold D, and returns a set of barcodes, *as short as you can find while satisfying the requirements*, using a reasonable amount of compute time. As many as 100 barcodes might be requested, with typical values of D ranging from 1 to 6. *Barcodes should be returned as an array of strings.*

Note: The speed of your implementation, and how it scales to large input arrays, will be considered in the assessment.

Hint: Especially here, do not let perfect be the enemy of good.