

CS 124 Programming Assignment 1: Spring 2021

Your name(s) (up to two):

Collaborators: (You shouldn't have any collaborators but the up-to-two of you, but tell us if you did.)

No. of late days used on previous psets:

No. of late days used after including this pset:

Homework is due Wednesday 2021-02-24 at 11:59pm ET. You are allowed up to **twelve** (college)/**forty** (extension school) late days through the semester, but the number of late days you take on each assignment must be a nonnegative integer at most **two** (college)/**four** (extension school).

Overview: The purpose of this assignment is to experience some of the problems involved with implementing an algorithm (in this case, a minimum spanning tree algorithm) in practice. As an added benefit, we will explore how minimum spanning trees behave in random graphs.

Assignment: You may work in groups of two, or by yourself. Both partners will receive the same grade and turn in a single joint report. (You should each submit a copy of the report in Gradescope.)

I recommend using a common programming language such as Java, C, or C++. If you use a more obscure language, that is fine, but if there are errors that are correspondingly harder to find it may cause your grade to be lower. I might advise you not to use a scripting language like Python, although some students have successfully done the assignment in Python.

We will be considering *complete*, *undirected* graphs. A graph with n vertices is complete if all possible $\binom{n}{2}$ edges are present in the graph.

Consider the following types of graphs:

- Complete graphs on n vertices, where the weight of each edge is a real number chosen uniformly at random on $[0, 1]$.
- Complete graphs on n vertices, where the vertices are points chosen uniformly at random inside the unit square. (That is, the points are (x, y) , with x and y each a real number chosen uniformly at random from $[0, 1]$.) The weight of an edge is just the Euclidean distance between its endpoints.
- Complete graphs on n vertices, where the vertices are points chosen uniformly at random inside the unit cube (3 dimensions) and hypercube (4 dimensions). As with the unit square case above, the weight of an edge is just the Euclidean distance between its endpoints.

Your first goal is to determine in each case how the expected (average) weight of the minimum spanning tree (not an edge, the whole MST) grows as a function of n . This will require implementing an MST algorithm, as well as procedures that generate the appropriate random graphs. (You should check to see what sorts of random number generators are available on your system, and determine how to seed them, say with a value from the machine's clock.) You may implement any MST algorithm (or algorithms!) you wish; however, I suggest you choose carefully.

For each type of graph, you must choose several values of n to test. For each value of n , you must run your code on several randomly chosen instances of the same size n , and compute the average value for your runs. Plot your values vs. n , and interpret your results by giving a simple function $f(n)$ that describes

your plot. For example, your answer might be $f(n) = \log n$, $f(n) = 1.5\sqrt{n}$, or $f(n) = \frac{2n}{\log n}$. Try to make your answer as accurate as possible; this includes determining the constant factors as well as you can. On the other hand, please try to make sure your answer seems reasonable.

Code setup:

So that we may test your code ourselves as necessary, please make sure your code accepts the following command line form:

```
./randmst 0 numpoints numtrials dimension
```

The flag 0 is meant to provide you some flexibility; you may use other values for your own testing, debugging, or extensions. Note: you should test your program – design tests to make sure your program is working, for example by checking it on some small examples (or find other tests of your choosing). You don’t need to put your tests in your writeup, that is for you.

The value numpoints is n , the number of points; the value numtrials is the number of runs to be done; the value dimension gives the dimension. (Use dimension = 2 for the square, and 3 and 4 for cube and hypercube, respectively; use dimension = 0 for the case where weights are assigned randomly. Notice that dimension 1 is just not that interesting.) The output for the above command line should be the following:

```
average numpoints numtrials dimension
```

where average is the average minimum spanning tree weight over the trials.

Please pay attention to the following requirements. In order to grade appropriately, our objective is to ensure that we can run the programs without any special per-student attention. Hence we require:

- If possible, for compatibility reasons, the code should run on the Harvard system (nice.fas.harvard.edu cluster of Linux machines – you can use an ssh client to log in if you have a Harvard account), even if you code on another system.
- We expect an executable with the code as described above (input as described, answers should go to standard output as described).
- The code should compile with make; no instructions for humans. That is, the command “make randmst” should produce an executable from your directory. You may need to read up on makefiles to make this happen.

What to hand in: Besides *submitting a copy of the code you created*, your group should hand in a single well organized and clearly written report describing your results. For the first part of the assignment, this report must contain the following quantitative results (for each graph type):

- A table listing the average tree size for several values of n . (A *graph* is insufficient, although you can have that too; we need to see the actual numbers. If you have a graph but no table of actual numbers, I will take points off.)
- A description of your guess for the function $f(n)$.

Run your program for $n = 128; 256; 512; 1024; 2048; 4096; 8192; 16384; 32768; 65536; 131072; 262144$; and larger values, if your program runs fast enough. (Having your code handle up to at least $n = 262144$ vertices is one of the assignment requirements; however, handling n up to 131072 will result in only losing

1-2 points. Providing results only for smaller n will hurt your score on the assignment.) Run each value of n at least five times and take the average. (Make sure you re-seed the random number generator appropriately!)

In addition, you are expected to discuss your experiments in more depth. This discussion should reflect what you have learned from this assignment; the actual issues you choose to discuss are up to you. Here are some possible suggestions for the second part:

- Which algorithm did you use, and why?
- Are the growth rates (the $f(n)$) surprising? Can you come up with an explanation for them?
- How long does it take your algorithm to run? Does this make sense? Do you notice things like the cache size of your computer having an effect?
- Did you have any interesting experiences with the random number generator? Do you trust it?

Your grade will be based primarily on the correctness of your program and your discussion of the experiments. Other considerations will include the size of n your program can handle. Please do a careful job of solid writing in your writeup. Length will not earn you a higher grade, but clear descriptions of what you did, why you did it, and what you learned by doing it will go far.

Hints:

To handle large n , you may want to consider simplifying the graph. For example, for the graphs in this assignment, the minimum spanning tree is extremely unlikely to use any edge of weight greater than $k(n)$, for some function $k(n)$. We can estimate $k(n)$ using small values of n , and then try to throw away edges of weight larger than $k(n)$ as we increase the input size. Notice that throwing away too many edges may cause problems. Why will throwing away edges in this manner never lead to a situation where the program returns the wrong tree?

You may invent any other techniques you like, as long as they give the same results as a non-optimized program. Be sure to explain any techniques you use as part of your discussion and attempt to justify why they should give the same results as a non-optimized program!