
Problem Set 6

Both theory and programming questions are due on **Tuesday, November 15 at 11:59PM**. Please download the .zip archive for this problem set, and refer to the `README.txt` file for instructions on preparing your solutions.

We will provide the solutions to the problem set 10 hours after the problem set is due. You will have to read the solutions, and write a brief **grading explanation** to help your grader understand your write-up. You will need to submit the grading guide by **Tuesday, November 22, 11:59PM**. Your grade will be based on both your solutions and the grading explanation.

Problem 6-1. [30 points] I Can Haz Moar Frenzd?

Alyssa P. Hacker is interning at RenBook (人书 / 人書 in Chinese), a burgeoning social network website. She needs to implement a new friend suggestion feature. For two friends u and v (friendship is undirected), the **EdgeRank** $ER(u, v)$ can be computed in constant time based on the interest u shows in v (for example, how frequently u views v 's profile or comments on v 's wall). Assume that EdgeRank is directional and asymmetric, and that its value falls in the range $(0, 1)$. A user u is **connected** to a user v if they are connected through some mutual friends, i.e., $u = u_0$ has a friend u_1 , who has a friend u_2, \dots , who has a friend $u_k = v$. The integer k is the **vagueness** of the connection. Define the **strength** of such a connection to be

$$S(p) = \prod_{i=1}^k ER(u_{i-1}, u_i).$$

For a given user s , Alyssa wants to have a way to rank potential friend suggestions according to the strength of the connections s has with them. In addition, the vagueness of those connections should not be more than k , a value Alyssa will decide later.

Help Alyssa by designing an algorithm that computes the **strength** of the strongest connection between a given user s and **every other user** v to whom s is connected with vagueness at most k , in $O(kE + V)$ time (i.e., for every pair of s and v for $v \in V \setminus \{s\}$, compute the strength of the strongest connection between them with vagueness at most k). Assume that the network has $|V|$ users and $|E|$ friend pairs. Analyze the running time of your algorithm. For partial credit, give a slower but correct algorithm. You can describe your algorithm in words, pseudocode or both.

Solution: For each v , we want to maximize $S(p)$ for all p from s to v . This is equivalent to maximizing $\log S(p) = \sum \log ER(u_{i-1}, u_i)$, which is in turn equivalent to minimizing $w(p) = -\sum \log ER(u_{i-1}, u_i) = \sum \log \frac{1}{ER(u_{i-1}, u_i)}$. So for each edge (u, v) , we can assign its weight as $\log \frac{1}{ER(u, v)}$ which is nonnegative. Then the problem is converted to finding the shortest path from s to v for all $v \in V$ for paths with length less than or equal to k .

To do this, we use the modified Bellman-Ford algorithm with k iterations (see the pseudocode KLENGTH-BELLMAN-FORD below). During each iteration, we first check which edges need to be relaxed (line 6 - 8), but only update the distances after checking all the edges (line 9 - 10). This ensures that relaxation is only based on the distances from the previous iteration. This means after iteration i , $v.d = \min\{w(p) : |p| \leq i\}$.

KLENGTH-BELLMAN-FORD(G, s, k)

```

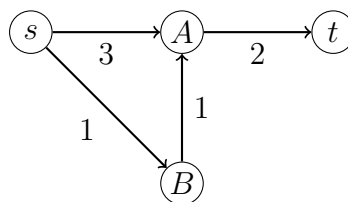
1  for each node  $v \in G.V$ 
2       $v.d = \infty$ 
3   $s.d = 0$ 
4  for  $i = 1$  to  $k$ 
5       $T = \{\}$ 
6      for each edge  $(u, v) \in G.E$ 
7          if  $v.d > u.d + w(u, v)$ 
8               $T[v] = u.d + w(u, v)$ 
9      for each  $v \in T$ 
10          $v.d = T[v]$ 

```

Time analysis: Initialization takes $O(V)$ time. Both checking and updating take $O(E)$ time, and there are k iterations. So total time is $O(kE + V)$.

Note that simply using unmodified BELLMAN-FORD algorithm with k iterations does not work because certain edge relaxation orders will find shortest paths with length greater than k in fewer than k iterations or even in 1 iteration. For example, in the graph below, suppose we want to find shortest paths with length at most 2. If we relax the edges in this order: (s, B) , (B, A) , (s, A) , (A, t) , after the 1st iteration, the shortest path to t is $\langle s, B, A, t \rangle$ which has length 3.

Similarly, using BFS with k iterations does not work either. Considering the same example above, in the second iteration, B and A are in the frontier. If we check B 's neighbors first and update A with weight 2, when we check A 's neighbors, we will update t with weight 4 but the length of this path is 3. The key is we need to use the distance values from the previous iteration instead of the new values in the current iteration for update.



Simply modifying BELLMAN-FORD to keep track of the length of the current path from s to v and avoid relaxation if the length is greater than k does not work either. Consider the graph above again, and we want to find shortest paths with length at most 2. If we relax the edges in this order: (s, B) , (B, A) , (s, A) , (A, t) and keep track of the path lengths, we will not find a path to

t . When we try to relax edge (A, t) , the shortest path to A is $\langle s, B, A \rangle$ which has length 2. The algorithm will avoid update $t.d$ because its shortest path would have length 3. So $t.d$ remains to be ∞ . However, the correct algorithm should return $t.d = 5$ through path $\langle s, A, t \rangle$ with length 2.

Problem 6-2. [30 points] **RenBook Competitor**

Having experienced RenBook and its alternatives, you've come to the conclusion that the social networking websites currently out there lack focus. You've decided to create a social network that caters specifically to the algorithms-in-Python crowd. Excited about this idea, you've pitched to various investors, a few of whom have traded you a check in return for convertible debt. You're ready to go! Now you have to implement your idea.

The first step is to rent a machine on the web, a task which you've easily accomplished using your newly-acquired funding. The next step, however, is harder to complete: you must install a web server library and all of its dependencies. Each library that you wish to install can depend on a number of other libraries, which you will have to install first. Each of those libraries can in turn have its own dependencies.

You can try to install the libraries one-by-one and resolve the dependencies by hand, but you'd like to write a script that will do the work for you. Fortunately, having studied graphs in 6.006, you're confident that you will be able to accomplish this task. You will need to determine which libraries need to be installed and then generate the order in which the libraries will be installed so that there will be no dependency problems.

Examining the software library repository, you see that there are V total libraries, which together have a total of E dependencies. The repository enforces the rule that dependencies cannot be cyclic. Libraries rarely all depend on one another, so you can safely assume that $E \ll V^2$.

- (a) An **installation order** is an ordering of all the libraries such that each library's dependencies appear prior to it in the sequence. If we install each library in this sequence in order, we are guaranteed to avoid dependency problems. Describe in detail how to generate an installation order for the entire repository in $O(V + E)$ time.

Solution: Each library is represented as a vertex. For each pair of dependency such that library u depends on library v , add a directed edge (u, v) in G . As G is directed and acyclic, we can topological sort on G (i.e., return the vertices in increasing order of their finishing times in DFS) to produce the installation order. Constructing G takes $O(V + E)$ time and topological sort (DFS) takes $O(V + E)$ time. So overall the running time is $O(V + E)$.

We wish to install a web server library along with its dependencies. Suppose that some libraries are already installed on your system, and that only P libraries remain to be installed (you can determine whether a library has already been installed by performing a dictionary lookup in $O(1)$ time). Assume that the maximum number of dependencies for any given library is D .

- (b) Give pseudocode for an algorithm that generates an installation order for the non-installed libraries that are needed for installing the web server library in $O(P + PD)$

time. Describe your algorithm. You may use any routine given in CLRS as a subroutine in your pseudocode, and you can use a textual description, a clarifying example, or a correctness proof for the description.

Solution: Suppose that you have an unweighted graph of $G = (V, E)$ that represents the dependencies. Each vertex $v \in V$ represents a software library, and each edge $(u, v) \in E$ represents a dependency (v must be installed before u). Let the web server library be called s . Perform a breadth-first search starting at s through the dependency graph. At each iteration of BFS, perform a lookup in constant time to check whether the nodes currently under consideration have already been installed. If you discover that a node has been installed, cut the edge leading to it. You can safely do this because you know that since a library has been installed, all of its dependencies have also been installed. This takes $O(PD)$ time because we only visit vertices that are not installed. We can now perform a topological sort on this new, smaller, graph. The only vertices that remain (V') in this graph are the packages that need to be installed (P), and the number of edges (E') at each node is at most D , so the running time is $O(V' + E') = O(P + PD)$.

You can also omit the BFS and do a modified topological sort that does not consider the installed vertices. Below is the pseudocode.

INSTALL-ORDER(G, s)

```

1  order = []
2  visited = {s}
3  DFS-VISIT( $G, s, order, visited$ )
4  return order
```

DFS-VISIT($G, v, order, visited$)

```

1  for  $n \in G.adj[v]$ 
2      if  $n$  is not installed and not in  $visited$ 
3          Add  $n$  to  $visited$ 
4          DFS-VISIT( $G, n, order, visited$ )
5  append  $v$  to  $order$ 
```

Problem 6-3. [30 points] Rubik's Cube

In this problem, you will develop algorithms for solving the $2 \times 2 \times 2$ Rubik's Cube, known as the *Pocket Cube*. Call a configuration of the cube " k levels from the solved position" if it can reach the solved configuration in exactly k twists, but cannot reach the solved configuration in any fewer twists.

The `rubik` directory in the problem set package contains the Rubik's Cube library and a graphical user interface to visualize your algorithm.

We will solve the Rubik's Cube puzzle by finding the shortest path between two configurations (the start and goal) using BFS.

A BFS that goes as deep as 14 levels (the diameter of the pocket cube) will take a few minutes (not to mention the memory needed). A few minutes is too slow for us: we want to solve the cube very quickly!

Instead, we will take advantage of a fact that we saw in lecture: the number of nodes at level 7 (half the diameter) is much smaller than half the total number of nodes.

With this in mind, we can instead do a two-way BFS, starting from each end at the same time, and meeting in the middle. At each step, expand one level from the start position, and one level from the end position, and then check to see whether any of the new nodes have been discovered in both searches. If there is such a node, we can read off parent pointers (in the correct order) to return the shortest path.

Write a function `shortest_path` in `solver.py` that takes two positions, and returns a list of moves that is a shortest path between the two positions.

Test your code using `test_solver.py`. Check that your code runs in less than 5 seconds.

Solution: The solution archive on the course Web site contains the staff's solution.

Problem 6-4. [30 points] **From Berklee to Berkeley**

Jack Florey and his fellow hackers are planning to put a TARDIS¹ on Berkeley's most symbolic tower. However, the company responsible for transporting the construction material mistook the destination as Berklee College of Music. In order to save the extra cost of transportation back to Berkeley, Jack wants to help them to find the fastest route from Berklee to Berkeley. He downloaded the data from the National Highway Planning Network (NHPN)². However, as he did not take 6.006 and does not know how to implement the Dijkstra's algorithm, he turns to you for help.

The partial code is in the `dijkstra` directory in the zip file for this problem set. It contains a `data` directory which includes `node` and `link` text files from the NHPN. Open `nhpn.nod` and `nhpn.lnk` in a text editor to get a sense of how the data is stored (`datadict.txt` and `format.txt` have more precise descriptions of the data fields and their meanings). The Python module `nhpn.py` provided contains code to load the text files into `Node` and `Link` objects. Read `nhpn.py` to understand the format of the `Node` and `Link` objects.

In `dijkstra.py`, the `PathFinder` object contains a source node, a destination node, and a `Network` object which represents the highway network with a list of `Node` objects. For each node, `node.adj` contains a list of all nodes adjacent to `node`.

Your task is implementing the method

```
PathFinder.dijkstra(weight, nodes, source, destination)
```

using Dijkstra's algorithm. It is given a function `weight(node1, node2)` that returns the weight of the link between `node1` and `node2`, a list of all the nodes, a source node and a destination node in the network. The method should return a tuple of the shortest path from

¹Doctor Who's Time And Relative Dimension In Space

²<http://www.fhwa.dot.gov/planning/nhpn/>

the source to the destination as a list of nodes, and the number of nodes visited during the execution of the algorithm. A node is visited if the shortest path of it from the source is computed. You should stop the search as soon as the shortest path to the destination is found.

Function `distance(node1, node2)` is a weight function used by the main program to call the `dijkstra` method. It returns the distance between two NHPN nodes. Nodes come with latitude and longitude (in millionths of degrees). For simplicity, we treat these as (x, y) coordinates on a flat surface, where the distance between two points can be calculated using the Pythagorean Theorem.

The `NodeDistancePair` object wraps a given node and its distance which can be used as a key in the priority queue.

The Python module `priority_queue.py` contains a min-heap based implementation of priority queue. It is augmented with a map of keys to their indices in the heap, so that the `decrease_key(key)` method takes $O(1)$ time to lookup the key in the priority queue.

After implementing `PathFinder.dijkstra()`, you can run

```
python dijkstra.py < tests/0boston_berkeley.in
```

to see the shortest distance from Boston, MA to Berkeley, CA.

To visualize the result, you can run the following command.

```
TRACE=kml python dijkstra.py < tests/0boston_berkeley.in
```

On Windows, use the following command instead:

```
dijkstra_kml.bat < tests/0boston_berkeley.in
```

This will create two files, `path_flat.kml` and `path_curved.kml`. Both should be paths from Boston, MA to Berkeley, CA. `path_flat.kml` is created using the distance function described earlier, and `path_curved.kml` is created using a distance function that does not assume that the Earth is flat. `.kml` files can be viewed using Google Maps, by putting the file in a web-accessible location (such as your Athena Public directory), going to <http://maps.google.com> and putting the URL (`http://...`) in the search box. Try asking Google Maps for driving directions from Berklee to Berkeley to get a sense of how similar their answer is. Two sample `.kml` files `path_flat_sol.kml` and `path_curved_sol.kml` are provided in the same directory and you can check your results against the samples.

You can use the following command to run all the tests on your Dijkstra's implementation:

```
python dijkstra_test.py
```

When your code passes all the tests and runs reasonably fast (the tests should complete in less than 40s when using CPython), upload your modified `dijkstra.py` to the course submission site. Our automated grading code will use our versions of `dijkstra_test.py`, `nHPN.py`, and `priority_queue.py`, so please do not modify these files.

Solution: The solution archive on the course Web site contains the staff's solution and secret test cases.