



6. THE GREEDY METHOD

Raveen de Silva, r.desilva@unsw.edu.au

office: K17 202

Course Admin: Anahita Namvar, cs3121@cse.unsw.edu.au

School of Computer Science and Engineering
UNSW Sydney

Term 3, 2021

Table of Contents

1. Introduction

2. Example Problems

3. Applications to Graphs

3.1 Single Source Shortest Paths

3.2 Minimum Spanning Trees

4. Puzzle

The Greedy Method

Question

What is a greedy algorithm?

Answer

A greedy algorithm is one that solves a problem by dividing it into stages, and rather than exhaustively searching all the ways to get from one stage to the next, instead only considers the choice that appears best.

This obviously reduces the search space, but it is not always clear whether the locally optimal choice leads to the globally optimal outcome.

The Greedy Method

Question

Does the greedy method always work?

Answer

No!

Suppose you are searching for the highest point in a mountain range. If you always climb upwards from the current point in the steepest possible direction, you will find a peak, but not necessarily the highest point overall.

The Greedy Method

Question

Is there a framework to decide whether a problem can be solved using a greedy algorithm?

Answer

Yes, but we won't use it.

The study of *matroids* is covered in CLRS. We will instead prove the correctness of greedy algorithms on a problem-by-problem basis. With experience, you will develop an intuition for whether the greedy method is useful for a particular problem.

The Greedy Method

Question

How do we prove that a greedy algorithm is correct?

Answer

There are two main methods of proof:

1. *Greedy stays ahead*: prove that at every stage, no other algorithm could do better than our proposed algorithm.
2. *Exchange argument*: consider an optimal solution, and gradually transform it to the solution found by our proposed algorithm without making it any worse.

Again, experience is most important!

Table of Contents

1. Introduction

2. Example Problems

3. Applications to Graphs

3.1 Single Source Shortest Paths

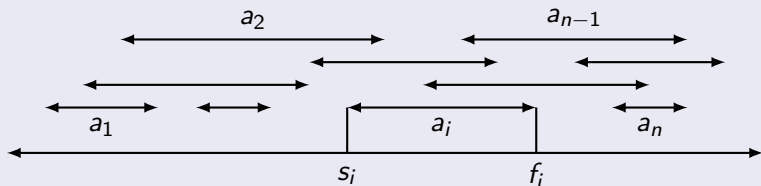
3.2 Minimum Spanning Trees

4. Puzzle

Activity Selection

Problem

Instance: A list of n activities a_i , with starting times s_i and finishing times f_i . No two activities can take place simultaneously.



Task: Find a *maximum size* subset of compatible activities.

Activity Selection

Attempt 1

Always choose the shortest activity which does not conflict with the previously chosen activities, then remove the conflicting activities and repeat.

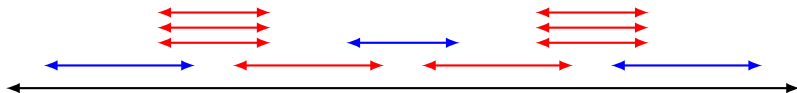


In the above example, our proposed algorithm chooses the activities in blue, then has to discard all the red activities, so clearly this does not work.

Activity Selection

Attempt 2

Maybe we should always choose an activity which conflicts with the fewest possible number of the remaining activities? It may appear that in this way we minimally restrict our next choice ...

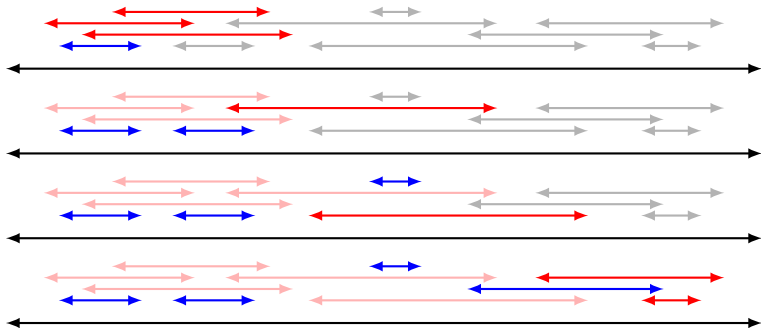


As appealing this idea is, the above figure shows this again does not work!

Activity Selection

Solution

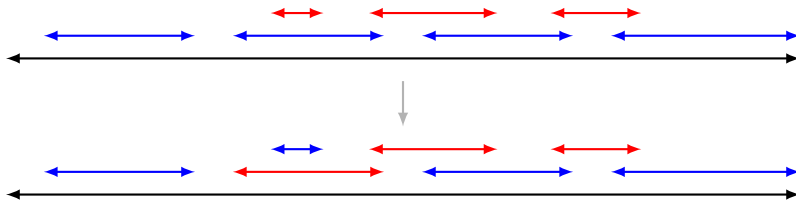
Among those activities which do not conflict with the previously chosen activities, always choose the activity with the earliest end time (breaking ties arbitrarily).



Activity Selection

To prove the correctness of our algorithm, we will use an *exchange argument*. We will show that any optimal solution can be transformed into a solution obtained by our greedy algorithm with the same number of activities.

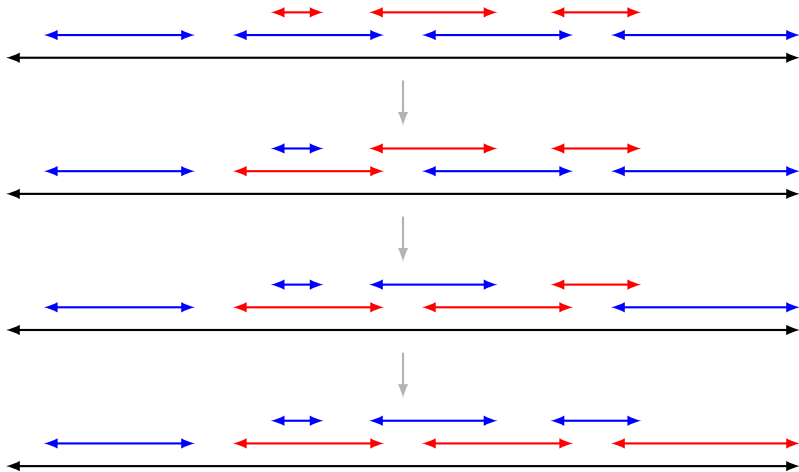
- Find the first place where the chosen activity violates the greedy choice.
- What if we replace that activity with the greedy choice?



Activity Selection

- Does the new selection have any conflicts? No!
- Does the new selection have the same number of activities? Yes!
- So the greedy choice is actually just as good as the choice used in the optimal solution!
- We replace it and repeat.
- Continuing in this manner, we eventually “morph” the optimal solution into the greedy solution, thus proving the greedy solution is also optimal.

Activity Selection



Activity Selection

- What is the time complexity of the algorithm?
- We represent activities by ordered pairs of their starting and their finishing times and sort them in increasing order of their finishing time (the second coordinate), in $O(n \log n)$ time.
- We go through this sorted list in order. How do we tell whether an activity conflicts with the already chosen activities?

Activity Selection

- Suppose we are up to activity i , starting at s_i and finishing at f_i , with all earlier finishing activities already processed.
- If all previously chosen activities finished before s_i , activity i can be chosen without a conflict. Otherwise, there will be a clash, so we discard activity i .
- We would prefer not to go through all previously chosen activities each time.

Activity Selection

- We need only keep track of the latest finishing time among chosen activities.
- Since we process activities in increasing order of finishing time, this is just the finishing time of the last activity to be chosen.
- Every activity is therefore either chosen (and the last finishing time updated) or discarded in constant time, so this part of the algorithm takes $O(n)$ time.
- Thus, the algorithm runs in total time $O(n \log n)$, dominated by sorting.

Activity Selection

A related problem

Instance: A list of n activities a_i with starting times s_i and finishing times $f_i = s_i + d$; thus, all activities are of the same duration. No two activities can take place simultaneously.

Task: Find a subset of compatible activities of *maximal total duration*.

Solution

Since all activities are of the same duration, this is equivalent to finding a selection with a largest number of non conflicting activities, i.e., the previous problem.

Activity Selection

Question

What happens if the activities are not all of the same duration?

Solution

The greedy strategy no longer works - we will need a more sophisticated technique.

Cell Towers

Problem

Instance: Along the long, straight road from Loololong (in the West) to Goolagong (in the East), houses are scattered quite sparsely, sometimes with long gaps between two consecutive houses. Telstra must provide mobile phone service to people who live alongside the road, and the range of Telstra's cell tower is 5km.



Task: Design an algorithm for placing the minimal number of cell towers alongside the road, that is sufficient to cover all houses.

Cell Towers

- Let's attempt a greedy algorithm, processing the houses west to east.
- The first house must be covered by some tower, which we place 5km to the east of this house.
- This tower may cover some other houses, but eventually we should reach a house that is out of range of this tower. We then place a second tower 5km to the east of that house.
- Continue in this way until all houses are covered.

Cell Towers

- At each house, we need to decide whether to place a new tower. This can be done in constant time by referring to the most recently created tower, which can itself be updated in constant time if necessary.
- Therefore this algorithm runs in $O(n)$ time if the houses are provided in order, and $O(n \log n)$ time otherwise.
- We can prove the correctness of this algorithm using an exchange argument; this is left as an exercise.

Cell Towers

- One of Telstra's engineers started with the house closest to Loololong and put a tower 5km away to the east. He then found the westmost house not already in the range of the tower and placed another tower 5 km to the east of it and continued in this way till he reached Goolagong.
- His junior associate did exactly the same but starting from Goolagong and moving westwards and claimed that his method required fewer towers.
- Is there a placement of houses for which the associate is right?

Minimising Job Lateness

Problem

Instance: A start time T_0 and a list of n jobs a_i , with duration times t_i and deadlines d_i . Only one job can be performed at any time; all jobs have to be completed. If a job a_i is completed at a finishing time $f_i > d_i$ then we say that it has incurred lateness $l_i = f_i - d_i$.

Task: Schedule all the jobs so that the lateness of the job with the largest lateness is minimised.

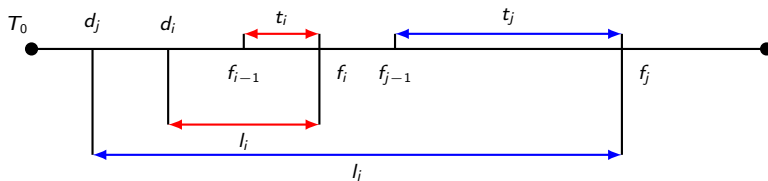
Minimising Job Lateness

Solution

Ignore job durations and schedule jobs in the increasing order of deadlines.

Proof of optimality

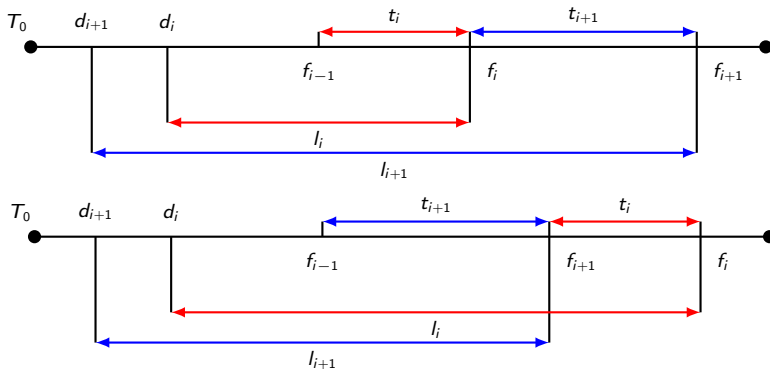
Consider any optimal solution. We say that jobs a_i and jobs a_j form an inversion if job a_i is scheduled before job a_j but $d_j < d_i$.



Minimising Job Lateness

- We will show that there exists a scheduling without inversions which is also optimal.
- Recall the BUBBLESORT algorithm: if we manage to eliminate all inversions between adjacent jobs, eventually all the inversions will be eliminated.

Minimising Job Lateness



Note that swapping adjacent inverted jobs reduces the larger lateness!

Tape Storage

Problem

Instance: A list of n files f_i of lengths l_i which have to be stored on a tape. Each file is equally likely to be needed. To retrieve a file, one must start from the beginning of the tape and scan it until the file is found and read.

Task: Order the files on the tape so that the average (expected) retrieval time is minimised.

Tape Storage

- If the files are stored in order l_1, l_2, \dots, l_n , then the expected time is proportional to

$$\begin{aligned} & l_1 + (l_1 + l_2) + (l_1 + l_2 + l_3) + \dots + (l_1 + l_2 + l_3 + \dots + l_n) \\ &= nl_1 + (n-1)l_2 + (n-2)l_3 + \dots + 2l_{n-1} + l_n. \end{aligned}$$

- This is minimised if $l_1 \leq l_2 \leq l_3 \leq \dots \leq l_n$, so we simply sort the files by increasing order of length for an $O(n \log n)$ solution.

Tape Storage II

Problem

Instance: A list of n files of lengths l_i and probabilities to be needed p_i , $\sum_{i=1}^n p_i = 1$, which have to be stored on a tape. To retrieve a file, one must start from the beginning of the tape and scan it until the file is found and read.

Task: Order the files on the tape so that the **expected** retrieval time is minimised.

Tape Storage II

- If the files are stored in order l_1, l_2, \dots, l_n , then the expected time is proportional to

$$p_1 l_1 + (l_1 + l_2) p_2 + (l_1 + l_2 + l_3) p_3 \\ + \dots + (l_1 + l_2 + l_3 + \dots + l_n) p_n$$

- We now show that this is minimised if the files are ordered in a decreasing order of values of the ratio p_i/l_i .

Tape Storage II

- Let us see what happens if we swap two adjacent files, say files k and $k + 1$.
- The expected time before the swap and after the swap are, respectively,

$$\begin{aligned} E = & l_1 p_1 + (l_1 + l_2) p_2 + (l_1 + l_2 + l_3) p_3 + \dots \\ & + (l_1 + l_2 + l_3 + \dots + l_{k-1} + l_k) p_k \\ & + (l_1 + l_2 + l_3 + \dots + l_{k-1} + l_k + l_{k+1}) p_{k+1} \\ & + \dots + (l_1 + l_2 + l_3 + \dots + l_n) p_n \end{aligned}$$

and

$$\begin{aligned} E' = & l_1 p_1 + (l_1 + l_2) p_2 + (l_1 + l_2 + l_3) p_3 + \dots \\ & + (l_1 + l_2 + l_3 + \dots + l_{k-1} + l_{k+1}) p_{k+1} \\ & + (l_1 + l_2 + l_3 + \dots + l_{k-1} + l_{k+1} + l_k) p_k \\ & + \dots + (l_1 + l_2 + l_3 + \dots + l_n) p_n. \end{aligned}$$

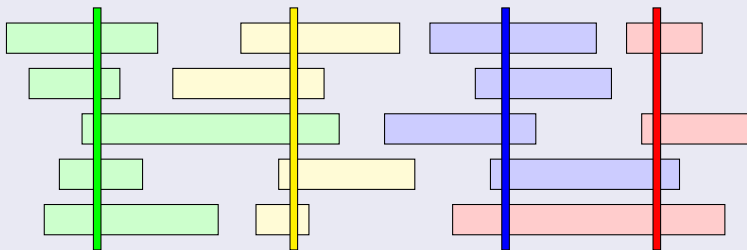
Tape Storage II

- Thus, $E - E' = l_k p_{k+1} - l_{k+1} p_k$, which is positive whenever $l_k p_{k+1} > l_{k+1} p_k$, i.e. when $p_k/l_k < p_{k+1}/l_{k+1}$.
- Consequently, $E > E'$ if and only if $p_k/l_k < p_{k+1}/l_{k+1}$, which means that the swap decreases the expected time whenever $p_k/l_k < p_{k+1}/l_{k+1}$, i.e., if there is an inversion: file $k + 1$ with a larger ratio p_{k+1}/l_{k+1} has been put after file k with a smaller ratio p_k/l_k .
- For as long as there are inversions there will be inversions of consecutive files and swapping will reduce the expected time. Consequently, the optimal solution is the one with no inversions.

Interval Stabbing

Problem

Instance: Let X be a set of n intervals on the real line, described by two arrays $X_L[1..n]$ and $X_R[1..n]$, representing their left and right endpoints. We say that a set P of points stabs X if every interval in X contains at least one point in P .

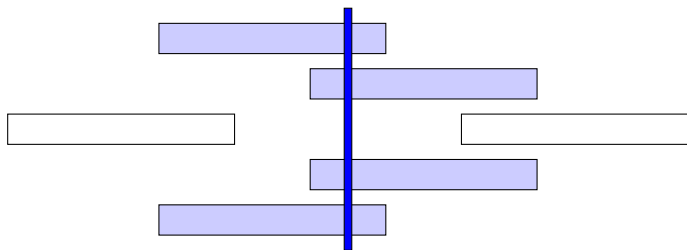


Task: Describe and analyse an efficient algorithm to compute the smallest set of points that stabs X .

Interval Stabbing

Attempt 1

Is it a good idea to stab the largest possible number of intervals?



No! In the above example, this strategy needs three stabbing points rather than two.

Interval Stabbing

Hint

The interval which ends the earliest has to be stabbed somewhere.

What is the best place to stab it?

Fractional Knapsack

Problem

Instance: A list of n items described by their weights w_i and values v_i , and a maximal weight limit W of your knapsack. You can take each item any number of times (not necessarily integer).

Task: Select a non-negative quantity of each item, with total weight not exceeding W and maximal total value.

Fractional Knapsack

Solution

Fill the entire knapsack with the item with highest value per unit weight!

0-1 Knapsack

Problem

Instance: A list of n *discrete items* described by their weights w_i and values v_i , and a maximal weight limit W of your knapsack.

Task: Find a subset S of the items with total weight not exceeding W and maximal total value.

0-1 Knapsack

- Can we always choose the item with the highest value per unit weight?
- Assume there are just three items with weights and values:

$A(10 \text{ kg}, \$60), B(20 \text{ kg}, \$100), C(30 \text{ kg}, \$120)$

and a knapsack of capacity $W = 50 \text{ kg}$.

- The greedy strategy would choose items A and B , while the optimal solution is to take items B and C !
- So when does the Greedy Strategy work??
- Unfortunately there is no easy rule ...

Array Merging

- Assume you are given n sorted arrays of different sizes.
- You are allowed to merge any two arrays into a single new sorted array and proceed in this manner until only one array is left.
- Design an algorithm that achieves this task and uses minimal total number of moves of elements of the arrays and give an informal justification why your algorithm is optimal.
- This problem is somewhat related to the next problem, which is arguably among the most important applications of the greedy method!

The Huffman Code

- Assume you are given a set of symbols, for example the English alphabet plus punctuation marks and a blank space (to be used between words).
- You want to encode these symbols using binary strings, so that sequences of such symbols can be decoded in an unambiguous way.

The Huffman Code

- One way of doing so is to reserve bit strings of equal and sufficient length, given the number of distinct symbols to be encoded. For example, if you have 26 letters and up to 6 punctuation symbols, you could use strings of 5 bits, as $2^5 = 32$.
- To decode a piece of text you would partition the bit stream into groups of 5 bits and use a lookup table to decode the text.

The Huffman Code

- However this might not be the most economical way: all the symbols have codes of equal length but the symbols are not equally frequent.
- One would prefer an encoding in which frequent symbols such as 'a', 'e', 'i' or 't' have short codes while infrequent ones, such as 'q', 'x' and 'z' can have longer codes.

The Huffman Code

- However, if the codes are of variable length, then how can we partition a bitstream *uniquely* into segments each corresponding to a code?
- One way of ensuring unique readability of codes from a single bitstream is to ensure that no code of a symbol is a prefix of a code for another symbol.
- Codes with such property are called *prefix codes*.

The Huffman Code

- We can now formulate the problem as follows:
Given the frequencies (probabilities of occurrences) of each symbol, design an optimal prefix code, i.e. a prefix code which minimises the expected length of an encoded text.
- Note that this amounts to saying that the *average* number of bits per symbol in an “average” text is as small as possible.
- We now sketch the algorithm informally; please see the textbook for details and the proof of optimality.
- MATH3411 Information, Codes & Ciphers covers this and much more!

Table of Contents

- 1. Introduction
- 2. Example Problems
- 3. Applications to Graphs
 - 3.1 Single Source Shortest Paths
 - 3.2 Minimum Spanning Trees
- 4. Puzzle

Tsunami Warning

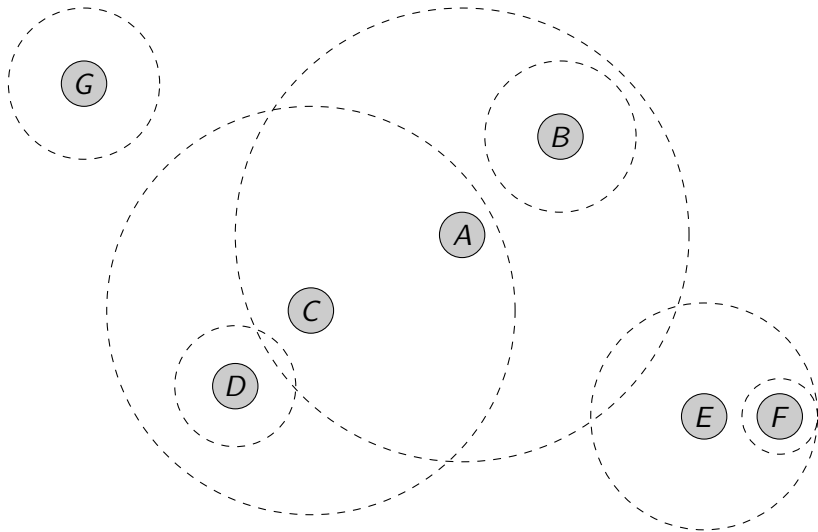
Problem

Instance: There are n radio towers for broadcasting tsunami warnings. You are given the (x, y) coordinates of each tower and its radius of range. When a tower is activated, all towers within the radius of range of the tower will also activate, and those can cause other towers to activate and so on.

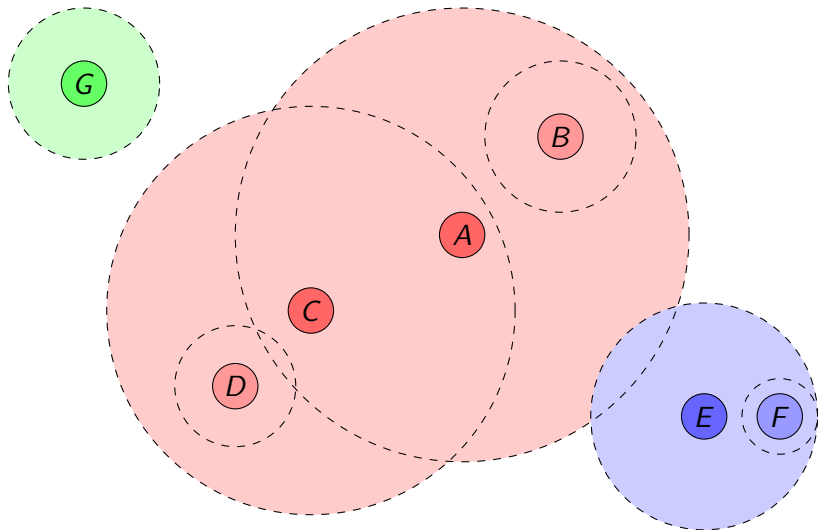
You need to equip some of these towers with seismic sensors so that when these sensors activate the towers where these sensors are located all towers will eventually get activated and send a tsunami warning.

Task: Design an algorithm which finds the fewest number of towers you must equip with seismic sensors.

Tsunami Warning



Tsunami Warning



Tsunami Warning

Attempt 1

Find the unactivated tower with the largest radius (if multiple with the same radius, pick the any of them). Place a sensor at this tower. Find and remove all activated towers. Repeat.

Attempt 2

Find the unactivated tower with the largest number of towers within its range. If there is none, place a sensor at the leftmost tower. Repeat.

Exercise

Give examples which show that neither Algorithm 1 nor Algorithm 2 solve the problem correctly.

Tsunami Warning

- It is useful to consider the towers as vertices of a directed graph, where an edge from tower A to tower B indicates that the activation of A *directly* causes the activation of B , that is, B is within the radius of A .

Observation

Suppose that activating tower A causes tower B to also be activated, and vice versa. Then we never want to place sensors at both towers; indeed, placing a sensor at A *is equivalent to* placing a sensor at B .

Tsunami Warning

Observation

We can extend this notion to a larger number of towers.

Suppose S is a subset of the towers and that activating any tower in S causes the activation of all towers in S .

We never want to place more than one sensor in S , and if we place one, then it doesn't matter where we put it.

In this way, we can treat all of S as a unit; a *super-tower*.

Strongly Connected Components

Definition

Given a directed graph $G = (V, E)$ and a vertex v , the *strongly connected component* of G containing v consists of all vertices $u \in V$ such that there is a path in G from v to u and a path from u to v . We will denote it by C_v .

In the terms of our problem, strongly connected components are *maximal* super-towers.

Strongly Connected Components

- How do we find the strongly connected component $C_v \subseteq V$ containing v ?
- Construct another graph $G_{rev} = (V, E_{rev})$ consisting of the same set of vertices V but with the set of edges E_{rev} obtained by reversing the direction of all edges E of G .

Claim

u is in C_v if and only if u is reachable from v and v is reachable from u .

Equivalently, u is reachable from v in both G and G_{rev} .

Strongly Connected Components

- Use BFS to find the set $R \subseteq V$ of all vertices in V which are reachable in G from v .
- Similarly find the set $R' \subseteq V$ of all vertices which are reachable in G_{rev} from v .
- The strongly connected component of G containing v is given by $C_v = R \cap R'$.
- Finding all strongly connected components in this way could require $O(|V|)$ traversals of the graph.

Strongly Connected Components

- Faster algorithms exist! Famously, Kosaraju's algorithm and Tarjan's algorithm find all strongly connected components of a directed graph in $O(|V| + |E|)$.
- We won't cover them in this course, but we will in COMP4128.
- These faster algorithms are *linear time* (by which we mean linear in the size of the input, which consists of $|V|$ vertices and $|E|$ edges).
- A linear time algorithm is asymptotically “no slower than” reading the input, so we can run these algorithms “for free”, i.e. without worsening the time complexity of our solution to a problem.

The Condensation Graph

- It should be clear that distinct strongly connected components are disjoint sets, so the strongly connected components form a partition of V .
- Let C_G be the set of all strongly connected components of a graph G .

Definition

Define the *condensation graph* $\Sigma_G = (C_G, E^*)$, where

$$E^* = \{(C_{u_1}, C_{u_2}) \mid (u_1, u_2) \in E, C_{u_1} \neq C_{u_2}\}.$$

The vertices of Σ_G are the strongly connected components of G , and the edges of Σ_G correspond to those edges of G that are not within a strongly connected component, with duplicates ignored.

The Condensation Graph

- We begin our solution to the tsunami warning problem by finding the condensation graph.
- Now we have the set of super-towers, and we know for each super-tower which others it can activate.
- Our task is to decide which super-towers need a sensor installed in order to activate all the super-towers.
- We need to know one more property about the condensation graph.

The Condensation Graph

Claim

The condensation graph Σ_G is a directed acyclic graph.

Proof Outline

Suppose there is a cycle in Σ_G . Then the vertices on this cycle are not *maximal* strongly connected sets, as they can be merged into an even larger strongly connected set.

Tsunami Warning

Solution

The correct greedy strategy is to only place a sensor in each super-tower without incoming edges in the condensation graph.

Proof

These super-towers cannot be activated by another super-tower, so they each require a sensor. This shows that there is no solution using fewer sensors.

Tsunami Warning

Proof (continued)

We still have to prove that this solution activates all super-towers.

Consider a super-tower with one or more incoming edges. Follow any of these edges backwards, and continue backtracking in this way.

Since the condensation graph is acyclic, this path must end at some super-tower without incoming edges. The sensor placed here will then activate all super-towers along our path.

Therefore, all super-towers are activated as required.

Topological Sorting

Definition

Let $G = (V, E)$ be a directed graph, and let $n = |V|$. A *topological sort* of G is a linear ordering (enumeration) of its vertices $\sigma : V \rightarrow \{1, \dots, n\}$ such that if there exists an edge $(v, w) \in E$ then v precedes w in the ordering, i.e., $\sigma(i) < \sigma(j)$.

Property

A directed acyclic graph permits a topological sort of its vertices.

Note that the topological sort is not necessarily unique, i.e., there may be more than one valid topological ordering of the vertices.

Topological Sorting

```
1: function TOPOLOGICAL-SORT( $G$ )
2:    $L \leftarrow$  Empty list that will contain ordered elements
3:    $S \leftarrow$  Set of all nodes with no incoming edge
4:   while  $S$  is non-empty do
5:     remove a node  $u$  from  $S$ ;
6:     add  $u$  to tail of  $L$ ;
7:     for each node  $v$  with an edge  $e = (u, v)$  do
8:       remove edge  $e$  from the graph;
9:       if  $v$  has no other incoming edges then
10:        insert  $v$  into  $S$ ;
11:       end if
12:     end for
13:   end while
```


Topological Sorting

```
14:   if the graph has edges left, then  
15:       return error (graph has at least one cycle)  
16:   else  
17:       return  $L$  (a topologically sorted order)  
18:   end if  
19: end function
```

Topological Sorting

- This algorithm runs in $O(|V| + |E|)$, that is, *linear time*.
- Once again, we can run this algorithm “for free” as it is asymptotically no slower than reading the graph.
- In problems involving directed acyclic graphs, it is often useful to start with a topological sort and then think about the actual problem!

Table of Contents

- 1. Introduction
- 2. Example Problems
- 3. Applications to Graphs
 - 3.1 Single Source Shortest Paths
 - 3.2 Minimum Spanning Trees
- 4. Puzzle

Single Source Shortest Paths

- Suppose we are given a directed graph $G = (V, E)$ with *non-negative* weight $w(e) \geq 0$ assigned to each edge $e \in E$, and a designated vertex $v \in V$.
- For simplicity, we will assume that for every $u \in V$ there is a path from v to u .
- The task is to find for every $u \in V$ the shortest path from v to u .
- This is accomplished by a very elegant greedy algorithm developed by Edsger Dijkstra in 1959.

Single Source Shortest Paths

We first prove a simple fact about shortest paths, which will form the basis of the algorithm.

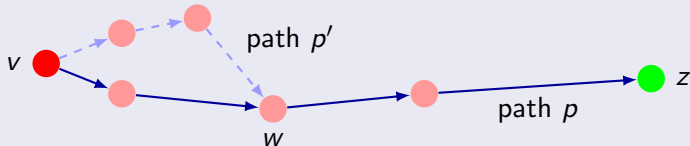
Claim

For every vertex w on a shortest path from v to z , the shortest path from v to w is just the truncation of that path ending at w .

Single Source Shortest Paths

Proof

Consider a shortest path p in G from a vertex v to a vertex z (shown in dark blue). Assume the claim is false, i.e., there is a shorter path from v to w (shown in dashed light blue) which is not an initial segment of the shortest path from v to z .

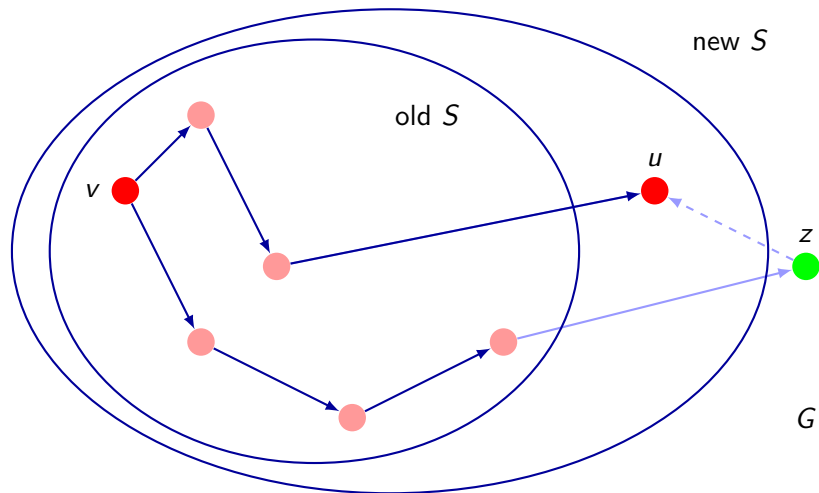


However, in this case we could remove the part of path p up to w and replace it with the path p' , thus obtaining a shorter path from v to z . This is a contradiction!

Dijkstra's Shortest Paths Algorithm

- The algorithm builds a set S of vertices for which the shortest path has been already established, starting with an empty set and adding one vertex at a time.
- At each stage of the construction, we add the vertex $u \in V \setminus S$ which has the shortest path from v to u with all intermediate vertices already in S .
- Note that the first vertex added to S is v itself, with path length zero.

Dijkstra's Shortest Paths Algorithm



Dijkstra's Shortest Paths Algorithm

- Why does this produce a shortest path p from v to u in G ?
- Assume the opposite, that there exists a shorter path p' from v to u .
- Throughout the proof, we will denote by S the *old* set, as it was before adding u .
- The path p' cannot be entirely in S , so let z be the first vertex on p' which is outside S .
- The part of p' up to z is at least as long as p , and since there are no negative weight edges, the full path p' is also at least as long as p . This is a contradiction, completing the proof.

Efficient Implementation of Dijkstra's Algorithm

Let $n = |V|$ and $m = |E|$, and suppose the vertices of the graph are labelled $1, \dots, n$, where vertex v is the source.

Attempt 1

Maintain an array $d[1, \dots, n]$ where $d[i]$ stores the length of the shortest path from vertex v to vertex i with all intermediate vertices in S . Initially $d[v] = 0$ and $d[i] = \infty$ otherwise.

At each stage of the algorithm, perform a linear search of the array d , ignoring those vertices already in S , and select the vertex u with smallest $d[u]$ to be added to S . For each outgoing edge from vertex u to another vertex $z \in V \setminus S$ with weight $w(u, z)$, we check whether $d[z] > d[u] + w(u, z)$, and if so we update $d[z]$ to the value $d[u] + w(u, z)$.

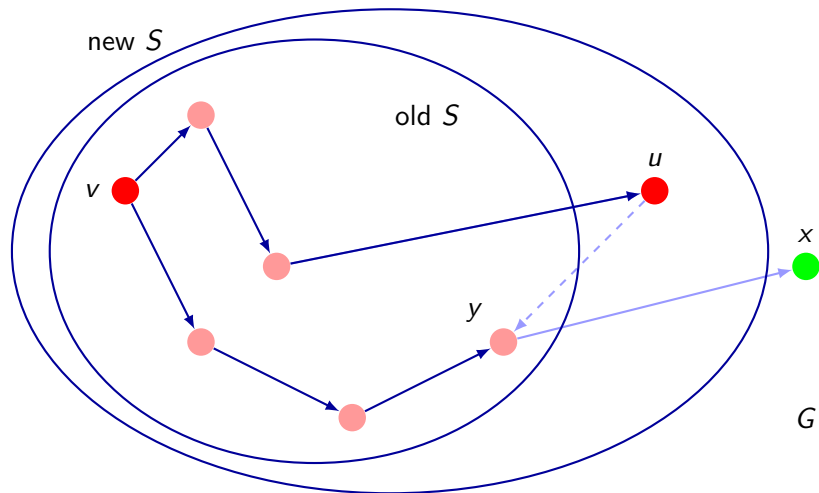
Dijkstra's Shortest Paths Algorithm

- In this algorithm, we have made use of a claim that we haven't yet proven.
- Recall that $d[i]$ stores the length of the shortest path from vertex v to vertex i with all intermediate vertices in S .
- When we add u to S , we only update these values for those paths where u is the penultimate vertex.

Question

Is it possible that the inclusion of u in S allows for a new shortest path through S from v to some vertex x without the last edge going from u to x ?

Dijkstra's Shortest Paths Algorithm



Dijkstra's Shortest Paths Algorithm

Answer

No!

Suppose that the inclusion of u in S allows for a new shortest path through S from v to x with penultimate vertex y .

Such a path must include u , or else it would not be new. Thus the path is $p = v \rightarrow \cdots \rightarrow u \rightarrow \cdots \rightarrow y \rightarrow x$.

Since y was added to S before u was, we know that the shortest path p' from v to y through S does not pass through u .

Finally, appending the edge from y to x to p' produces a path through S from v to x which is no longer than p , which is a contradiction.

Efficient Implementation of Dijkstra's Algorithm

- What is the time complexity of this algorithm?
- At each of n steps, we scan an array of length n . We also run the constant time update procedure at most once for each edge. The algorithm therefore runs in $O(n^2 + m)$.
- In a simple graph (no self loops or parallel edges) we have $m \leq n(n - 1)$, so we can simplify to $O(n^2)$.
- If the graph is *dense*, this is fine. But this is not guaranteed!

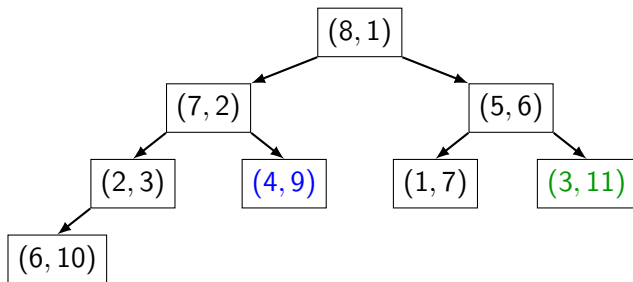
Efficient Implementation of Dijkstra's Algorithm

- Can we do better when $m \ll n^2$?
- The key is to avoid the linear search at each stage. We would like to pick out the vertex $u \in V \setminus S$ with smallest $d[u]$ in faster than linear time.
- What data structure enables this?
- We would like to use a heap, but the standard heap doesn't allow us to update arbitrary elements.

Augmented Heaps

- We will use a heap represented by an array $A[1, \dots, n]$; the left child of $A[i]$ is stored in $A[2i]$ and the right child in $A[2i + 1]$.
- Every element of A is of the form $(i, d(i))$, where the key $d(i)$ is the distance of a path from vertex v to vertex i with all intermediate vertices in S . The min-heap property is maintained with respect to the keys only.
- We will also maintain another array P of the same length which stores the position of elements in the heap. If $A[j]$ refers to vertex i , then we set $P[i] = j$, so that $A[P[i]] = (i, d(i))$.
- Changing the key of an element i is now an $O(\log n)$ operation. We look up its position $P[i]$ in A , change the key of the element located at $A[P[i]]$, and then perform the Heapify operation to ensure that the heap property is preserved.

Augmented Heaps



i	1	2	3	4	5	6	7	8
$A(i)$	8	7	5	2	4	1	3	6
	1	2	6	3	9	7	11	10
$P(i)$	6	4	7	5	3	8	2	1

Efficient Implementation of Dijkstra's Algorithm

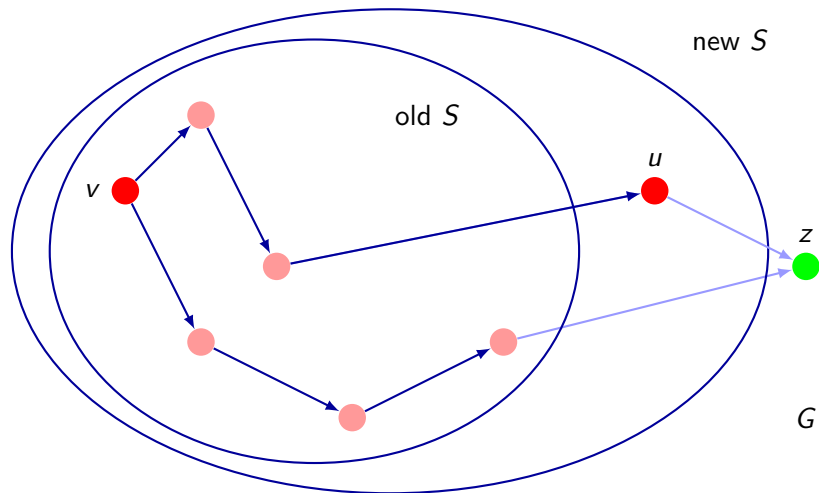
Algorithm

Initially we create a heap array A and a position array P , both of length n , with keys $d(v) = 0$ and $d(i) = \infty$ otherwise.

At each stage, we pop from the top of the heap, to obtain the vertex u with the smallest key and add it to set S .

We then go over all outgoing edges from u to a vertex $z \in V \setminus S$ with weight $w(u, z)$. If $d[z] > d[u] + w(u, z)$ then we update $d[z]$ to the value $d[u] + w(u, z)$, using the update procedure discussed earlier.

Dijkstra's Shortest Paths Algorithm



Dijkstra's Shortest Paths Algorithm

- What is the time complexity of our algorithm?
- Each of n stages requires a deletion from the heap (when an element is added to S), which takes $O(\log n)$ many steps.
- Each edge causes at most one update of a key in the heap, also taking $O(\log n)$ many steps.
- Thus, in total, the algorithm runs in time $O((n + m) \log n)$.
- Assuming the graph is connected, $m \geq n - 1$, so we can simplify to $O(m \log n)$.

Table of Contents

1. Introduction

2. Example Problems

3. Applications to Graphs

3.1 Single Source Shortest Paths

3.2 Minimum Spanning Trees

4. Puzzle

Minimum Spanning Trees

Definition

A minimum spanning tree T of a connected graph G is a subgraph of G (with the same set of vertices) which is a tree, and among all such trees it minimises the total length of all edges in T .

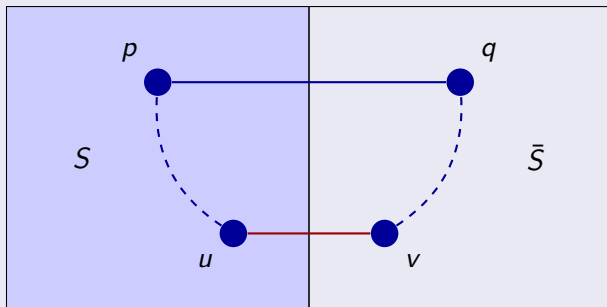
Lemma

Let G be a connected graph with all lengths of edges E of G distinct and S a non empty proper subset of the set of all vertices V of G . Assume that $e = (u, v)$ is an edge such that $u \in S$ and $v \notin S$ and is of minimal length among all the edges having this property. Then e must belong to every minimum spanning tree T of G .

Minimum Spanning Trees

Proof

Assume that there exists a minimum spanning tree T which does not contain such an edge $e = (u, v)$.



Minimum Spanning Trees

Proof (continued)

- Since T is a spanning tree, there exists a path from u to v within T , and this path must leave S by some edge, say (p, q) where $p \in S$ and $q \notin S$.
- However, (u, v) is shorter than any other edge with one end in S and one end outside S , including (p, q) .
- Replacing the edge (p, q) with the edge (u, v) produces a new tree T' with smaller total edge weight.
- This contradicts our assumption that T is a minimum spanning tree, completing the proof.

Minimum Spanning Trees

- There are two famous greedy algorithms for the minimum spanning tree problem.
- Both algorithms build up a forest, beginning with all n isolated vertices and adding edges one by one.
- *Prim's algorithm* uses one large component, adding one of the isolated vertices to it at each stage. This algorithm is very similar to Dijkstra's algorithm, but adds the vertex closest to S rather than the one closest to the starting vertex v .
- We will instead focus on *Kruskal's algorithm*.

Kruskal's Algorithm

- We order the edges E in a non-decreasing order with respect to their weights.
- An edge e is added if its inclusion does not introduce a cycle in the graph constructed thus far, or discarded otherwise.
- The process terminates when the forest is connected, i.e. when $n - 1$ edges have been added.

Kruskal's Algorithm

Claim

The Kruskal algorithm produces a minimal spanning tree, and if all weights are distinct, then such a Minimum Spanning Tree is unique.

Proof

We consider the case when all weights are distinct.

- Consider an edge $e = (u, v)$ added in the course of the Kruskal algorithm, and let F be the forest in its state *before* adding e .

Kruskal's Algorithm

Proof (continued)

- Let S be the set of vertices reachable from u in F . Then clearly $u \in S$ but $v \notin S$.
- The original graph does not contain any edges shorter than e with one end in S and the other outside S . If such an edge existed, it would have been considered before e and included in F , but then both its endpoints would be in S , contradicting the definition.
- Consequently, edge e is the shortest edge between a vertex of S and a vertex of \bar{S} and by the previous lemma it must belong to every spanning tree.

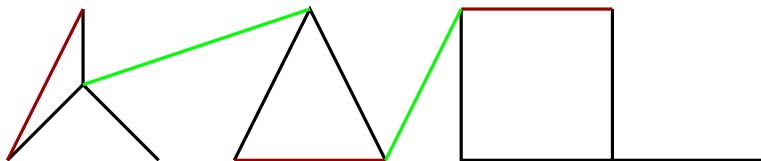
Kruskal's Algorithm

Proof (continued)

- Thus, the set of edges produced by the Kruskal algorithm is a subset of the set of edges of every spanning tree.
- But the graph produced by the Kruskal algorithm by definition has no cycles and is connected, so it is a tree.
- Therefore in the case where all edge weights are distinct, Kruskal's algorithm produces the unique minimum spanning tree.

Efficient Implementation of Kruskal's Algorithm

- To efficiently implement Kruskal's algorithm, we need to quickly determine whether a certain new edge will introduce a cycle.
- An edge $e = (u, v)$ will introduce a cycle in the forest F if and only if there is already a path between u and v , i.e., u and v are in the same connected component.



Union-Find

In our implementation of Kruskal's algorithm, we use a data structure known as *Union-Find*, which handles disjoint sets. This data structure supports three operations:

1. `MAKEUNIONFIND(S)`, which returns a structure in which all elements are placed into distinct singleton sets. This operation runs in time $O(n)$ where $n = |S|$.
2. `FIND(v)`, which returns the (label of the) set to which v belongs. This operation runs in time $O(1)$.
3. `UNION(A, B)`, which changes the data structure by replacing sets A and B with the set $A \cup B$. A sequence of k initial consecutive `UNION` operations runs in time $O(k \log k)$.

Union-Find

- Note that we do not give the run time of a single UNION operation but of a sequence of k consecutive such operations.
- Such time complexity analysis is called *amortized analysis*; it (essentially) estimates average cost of an operation in a sequence of operations, in this case $\log k$.
- When we refer to a sequence of k *initial consecutive* UNION operations, we refer to the first k UNION operations performed after the construction operation MAKEUNIONFIND. There may be FIND operations interspersed between these UNION operations.

Union-Find

- In this application, S is the vertex set V of a graph, of the form $\{1, 2, \dots, n\}$. We will label each set by one of its elements, called the *representative* of the set.
- The simplest implementation of the Union-Find data structure consists of:
 1. an array A such that $A[i] = j$ means that i belongs to the set with representative j ;
 2. an array B such that $B[i]$ contains the number of elements in the set with representative i ;
 3. an array L such that $L[i]$ contains pointers to the head and tail of a linked list containing the elements of the set with representative i .
- If i is not the representative of any set, then $B[i]$ is zero and the list $L[i]$ is empty.

Union-Find

Given two sets I and J with representatives i and j , $\text{UNION}(i, j)$ is defined as follows:

- assume $|I| \geq |J|$ (i.e. $B[i] \geq B[j]$); otherwise perform $\text{UNION}(j, i)$ instead;
- for each $m \in J$, update $A[m]$ from j to i ;
- update $B[i]$ to $B[i] + B[j]$ and $B[j]$ to zero;
- append the list $L[j]$ to the list $L[i]$ and replace $L[j]$ with an empty list.

Union-Find

Property

Observe that the new value of $B[i]$ is at least twice the old value of $B[j]$. Therefore if $\text{UNION}(i, j)$ changes the label of the set containing an element m (namely $A[m]$), then the number of elements in the set containing m (namely $B[A[m]]$) must have at least doubled.

Union-Find

- Any sequence of k initial consecutive UNION operations can touch at most $2k$ elements of S (which happens if all UNION operations were applied to singleton sets).
- Thus, the set containing an element m after k initial consecutive UNION operations must have at most $2k$ elements.
- Since every UNION operation which changes $A[m]$ at least doubles $B[A[m]]$, we deduce that $A[m]$ has changed at most $\log 2k$ times.
- Thus, since we have at most $2k$ elements, any sequence of k initial consecutive UNION operations will have at most $2k \log 2k$ label changes in A .

Union-Find

- Each UNION operation requires only constant time to update the size array B and the list array L .
- Thus, every sequence of k initial consecutive UNION operations has time complexity of $O(k \log k)$.
- This Union-Find data structure is good enough to get the sharpest possible bound on the run time of the Kruskal algorithm.
- See the textbook for a Union-Find data structure based on pointers and path compression, which further reduces the amortised complexity of the UNION operation at the cost of increasing the complexity of the FIND operation from $O(1)$ to $O(\log n)$.

Efficient Implementation of Kruskal's Algorithm

- We now use the previously described Union-Find data structure to efficiently implement the Kruskal algorithm on a graph $G = (V, E)$ with n vertices and m edges.
- We first have to sort m edges of graph G which takes time $O(m \log m)$. Since $m < n^2$, we can rewrite this as $O(m \log n^2) = O(m \log n)$.
- As we progress through the Kruskal algorithm execution, we will start with n isolated vertices, which will be merged into connected components until all vertices belong to a single connected component. We use the Union-Find data structure to keep track of the connected components constructed at any stage.

Efficient Implementation of Kruskal's Algorithm

- For each edge $e = (u, v)$ on the sorted list of edges, we use two FIND operations to determine whether vertices u and v belong to the same component.
- If they do not belong to the same component, i.e., if $\text{FIND}(u) = i$ and $\text{FIND}(v) = j$ where $i \neq j$, we add edge $e = (u, v)$ to the spanning tree being constructed and perform $\text{UNION}(i, j)$ to merge the connected components containing u and v .
- We perform $2m$ FIND operations, each costing $O(1)$, as well as $n - 1$ UNION operations which in total cost $O(n \log n)$.
- The initial sorting of the edges in $O(m \log n)$ dominates, so this is the overall time complexity.

k -clustering of maximum spacing

Problem

Instance: A complete graph G with weighted edges representing distances between the two vertices.

Task: Partition the vertices of G into k disjoint subsets so that the minimal distance between two points belonging to different sets of the partition is as large as possible. Thus, we want a partition into k disjoint sets which are as far apart as possible.

k -clustering of maximum spacing

Solution

Sort the edges in increasing order and start performing the usual Kruskal's algorithm for building a minimal spanning tree, but stop when you obtain k connected components, rather than a single spanning tree.

k -clustering of maximum spacing

Proof of optimality

- Let d be the distance associated with the first edge of the minimal spanning tree which was not added to our k connected components.
- It is clear that d is the minimal distance between two vertices belonging to different connected components.
- All the edges included in the connected components produced by our algorithm are of length at most d .

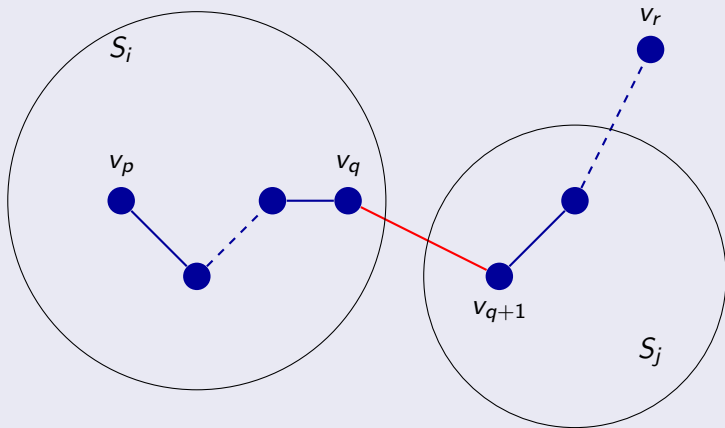
k -clustering of maximum spacing

Proof of Optimality (continued)

- Consider any partition \mathcal{S} into k subsets different from the one produced by our algorithm.
- This means that there is a connected component produced by our algorithm which contains vertices v_p and v_r such that $v_p \in S_i$ and $v_r \notin S_i$ for some $S_i \in \mathcal{S}$.

k -clustering of maximum spacing

Proof of optimality (continued)



k -clustering of maximum spacing

Proof of optimality (continued)

- Since v_p and v_r belong to the same connected component, there is a path in that component connecting v_p and v_r .
- Let v_q and v_{q+1} be two consecutive vertices on that path such that v_q belongs to S_i and $v_{q+1} \notin S_i$.
- Thus, $v_{q+1} \in S_j$ for some $j \neq i$.

k -clustering of maximum spacing

Proof of optimality (continued)

- Since (v_q, v_{q+1}) was an edge chosen by our proposed algorithm, we know that $d(v_q, v_{q+1}) \leq d$.
- It follows that the distance between these two clusters $S_i, S_j \in \mathcal{S}$ is at most d .
- Thus, such a partition cannot be a more optimal clustering than the one produced by our algorithm.

k -clustering of maximum spacing

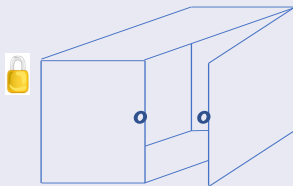
- What is the time complexity of this algorithm?
- We have $O(n^2)$ edges; thus sorting them by weight will take $O(n^2 \log n^2) = O(n^2 \log n)$.
- Running the (partial) Kruskal algorithm requires $O(m \log n) = O(n^2 \log n)$ steps, making use of the UNION-FIND data structure.
- So the grand total for the whole algorithm is $O(n^2 \log n)$ many steps.

Table of Contents

- 1. Introduction
- 2. Example Problems
- 3. Applications to Graphs
 - 3.1 Single Source Shortest Paths
 - 3.2 Minimum Spanning Trees
- 4. Puzzle

PUZZLE!!

Problem



Bob is visiting Elbonia and wishes to send his teddy bear to Alice, who is staying at a different hotel. Both Bob and Alice have boxes like the one shown on the picture as well as padlocks which can be used to lock the boxes.

PUZZLE!!

Problem (continued)

However, there is a problem. The Elbonian postal service mandates that when a nonempty box is sent, it must be locked. Also, they do not allow keys to be sent, so the key must remain with the sender. Finally, you can send padlocks only if they are locked. How can Bob safely send his teddy bear to Alice?

PUZZLE!!

Hint

The way in which the boxes are locked (via a padlock) is important. It is also crucial that *both* Bob and Alice have padlocks and boxes. They can also communicate over the phone to agree on the strategy.

There are two possible solutions; one can be called the “AND” solution, the other can be called the “OR” solution. The “AND” solution requires 4 mail one way services while the “OR” solution requires only 2.



That's All, Folks!!