

Design and Analysis of Algorithms

EKESH KUMAR*

April 7, 2020

These are my course notes for CMSC 451: Design and Analysis of Algorithms, taught by Professor Clyde Kruskal. Gaps in lecture material are filled in with CLRS and Kleinberg & Tardos. Please send corrections to ekumar1@terpmail.umd.edu.

Contents

1	Tuesday, January 28, 2020	3
1.1	Introduction	3
1.2	Stable Marriage Problem	3
2	Thursday, January 30, 2020	5
2.1	Optimality and Correctness of Gale-Shapley	5
3	Tuesday, February 4, 2020	6
3.1	Graph Terminology	6
3.2	Graph Representations	6
3.3	Graph Traversal	7
4	Thursday, February 6, 2020	9
4.1	Articulation Points	9
5	Tuesday, February 11, 2020	12
5.1	Articulation Point Algorithm Implementation	12
5.2	Strongly Connected Components	12
5.3	Classifying Edges in a DFS Tree	14
6	Thursday, February 13, 2019	15
6.1	Kosaraju's Algorithm	15
6.2	Topological Sorting	16
6.3	Bipartite Graphs	19
7	Tuesday, February 18, 2020	21
7.1	The Union-Find Data Structure	21
7.1.1	Motivating the Union-Find Data Structure	21

*Email: ekumar1@terpmail.umd.edu

7.1.2	Implementation of the Union-Find Data Structure	23
7.1.3	Analysis of Union-Find Operations	25
7.2	The Minimum Spanning Tree Problem	26
7.2.1	Problem Statement	26
7.2.2	Kruskal's Algorithm	26
7.2.3	Prim's Algorithm	28
8	Thursday, February 20, 2020	30
8.1	Interval Scheduling	30
8.1.1	Extensions: Minimizing Lateness	32
8.2	Caching	32
8.3	Farthest in Future Algorithm	33
9	Tuesday, February 25, 2020	34
9.1	Prefix Codes	34
9.1.1	Constructing a Huffman code	35
9.2	Matrix Multiplication	37
10	Thursday, February 27, 2020	38
10.1	Strassen's Algorithm	38
10.2	Closest Pair of Points	39
11	Tuesday, March 3, 2020	40
11.1	Closest Pair of Points	40
11.2	Counting Inversions	41
12	Thursday, March 5, 2020	44
12.1	Convolutions	44
12.2	The Fast Fourier Transform	45
12.2.1	Polynomial Evaluation	46
12.2.2	Polynomial Interpolation	47
13	Tuesday, April 7, 2020	49
13.1	Subset Sum Problem	49

§1 Tuesday, January 28, 2020

§1.1 Introduction

This is CMSC 451: Design and Analysis of Algorithms. We will cover graphs, greedy algorithms, divide and conquer algorithms, dynamic programming, network flows, NP-completeness, and approximation algorithms.

- Homeworks are due every other Friday or so; NP-homeworks are typically due every other Wednesday.
- There is a 25% penalty on late homeworks, and there's one get-out-of-jail free card for each type of homework.

§1.2 Stable Marriage Problem

As an introduction to this course, we'll discuss the [stable marriage problem](#), which is stated as follows:

Given a set of n men and n women, match each man with a woman in such a way that the matching is *stable*.

What do we mean when we call a matching is “stable”? We call a matching *unstable* if there exists some man M who prefers a woman W over the woman he is married to, and W also prefers M over the man she is currently married to.

In order to better understand the problem, let's look at the $n = 2$ case. Call the two men M_1 and M_2 , and call the two women W_1 and W_2 .

- First suppose M_1 prefers W_1 over W_2 and W_1 prefers M_1 over M_2 . Also, suppose that M_2 prefers W_2 over W_1 and W_2 prefers M_2 , then
- If both W_1 and W_2 prefer M_1 over M_2 , and both M_1 and M_2 prefer W_1 over W_2 , then it's still easy to see what will happen: M_i will always match with W_i .
- Now let's say M_1 prefers W_1 to W_2 , M_2 prefers W_2 to W_1 , W_1 prefers M_2 to M_1 , and W_2 prefers M_1 to M_2 . In this case, the two men rank different women first, and the two women rank different men first. However, the men's preferences “clash” with the women's preferences. One solution to this problem is to match M_1 with W_1 and M_2 with W_2 . This is stable since both men get their top preference even though the two women are unhappy.

The solution to the problem starts to get a lot more complicated when the people's preferences do not exhibit any pattern. So how do we solve this problem in the general case? We can use the [Gale-Shapley algorithm](#). Before discussing this algorithm, however, we can make the following observations about this problem:

- Each of the n men and M woman are initially unmarried. If an unmarried man M chooses the woman W who is ranked highest on their list, then we cannot immediately conclude whether we can match M and w in our final matching. This is clearly the case since if we later find out about some other man M_2 who prefers W over any other woman, W may choose M_2 if she likes him more than M . However, we cannot immediately rule out M being matched to W either since a man like M_2 may not ever come.
- Just because everyone isn't happy doesn't mean a matching isn't stable. Some people might be unhappy, but there might not be anything they can do about it (if nobody wants to switch).

Moreover, we introduce the notion of a man *proposing* to a woman, which a woman can either accept or reject. If she is already engaged and accepts a proposal, then her existing engagement breaks off (the previous man becomes unengaged).

Now that we've introduced these basic ideas, we can now present the algorithm:

Input: A list of n men and n women to be matched.

Output: A valid stable matching.

```
stable_matching {
    set each man and each woman to "free"
    while there exists a man m who still has a woman w to propose to {
        let w be the highest ranked woman m hasn't proposed to.

        if w is free {
            (m, w) become engaged
        } else {
            let m' be the man w is currently engaged to.
            if w prefers m' to m {
                (m', w) remain engaged.
            } else {
                (m, w) become engaged and m' loses his partner.
            }
        }
    }
}
```

Proposition 1.1

The Gale-Shapley algorithm terminates in $\mathcal{O}(n^2)$ time.

Proof. In the worst case, n men end up proposing to n women. The act of proposing to another person is a constant-time operation. Thus, the $\mathcal{O}(n^2)$ runtime is clear. \square

§2 Thursday, January 30, 2020

§2.1 Optimality and Correctness of Gale-Shapley

Last time, we introduced the Gale-Shapley algorithm to find a stable matching. Today, we'll prove that the algorithm is correct (i.e. it never produces an unstable matching), and it is optimal for men (i.e. the men always end up for their preferred choice).

First, we'll show that the algorithm is correct:

Proposition 2.1

The matching generated by the Gale-Shapley algorithm is never an unstable matching.

Proof. Suppose, for the sake of contradiction, that m and w prefer each other over their current partner in the matching generated by the Gale-Shapley algorithm. This can happen either if m never proposed to w , or if m proposed to w and w rejected m . In the former case, m must prefer his partner to w , which implies that m and w do not form an unstable pair. In the latter case, w prefers her partner to m , which also implies m and w don't form an unstable pair. Thus, we arrive at a contradiction. \square

Next, we'll prove that the algorithm is optimal for men. However, before presenting the proof, observe that it is not too hard to see intuitively that the algorithm “favors” the men. Since the men are doing all of the proposing and the women can only do the deciding, it turns out that the men always ends up with their most preferred choice (as long as the matching remains stable).

Proposition 2.2

The matching generated by the Gale-Shapley algorithm gives men their most preferred woman possible without contradicting stability.

Proof. To see why this is true, let A be the matching generated by the men-proposing algorithm, and suppose there exists some other matching B that is better for at least one man, say m_0 . If m_0 is matched in B to w_1 which he prefers to his match in A , then in A , m_0 must have proposed to w_1 and w_1 must have rejected him. This can only happen if w_1 rejected him in favor of some other man — call him m_2 . This means that in B , w_1 is matched to m_0 but she prefers m_2 to m_0 . Since B is stable, m_2 must be matched to some woman that he prefers to w_1 ; say w_3 . This means that in A , m_2 proposed to w_3 before proposing to w_1 , and this means that w_3 rejected him. Since we can perform similar considerations, we end up tracing a “cycle of rejections” due to the finiteness of the sets A and B . \square

§3 Tuesday, February 4, 2020

Today, we'll recap graph terminology and elementary graph algorithms.

§3.1 Graph Terminology

Definition 3.1. A **graph** $G = (V, E)$ is defined by a set of vertices V and a set of edges E .

The number of vertices in the graph, $|V|$, is the **order** of the graph, and the number of edges in the graph, $|E|$, is the **size** of the graph. Typically, we reserve the letter n for the order of a graph, and we reserve m for the size of a graph.

Definition 3.2. We say a graph is **directed** if its edges can only be traversed in one direction. Otherwise, we say the graph is **undirected**.

Definition 3.3. A graph is called **simple** if it's an undirected graph without any loops (edges that start and end at the same vertex).

Definition 3.4. A graph is **connected** if for every pair of vertices u, v , there exists a path between u and v .

§3.2 Graph Representations

There are two primary ways in which we can represent graphs: **adjacency matrices** and **adjacency lists**.

An adjacency matrix is an $n \times n$ matrix A in which $A[u][v]$ is equal to 1 if the edge (u, v) exists in the graph; otherwise, $A[u][v]$ is equal to 0. Note that the adjacency matrix is symmetric if and only if the graph is undirected.

On the other hand, an adjacency list is a list of $|V|$ lists, one for each vertex. For each vertex $u \in V$, the adjacency list $\text{Adj}[u]$ contains all vertices v for which there exists an edge (u, v) in E . In other words, $\text{Adj}[u]$ contains all of the vertices adjacent to u in G .

Each graph representation has its advantages and disadvantages in terms of runtime. This is summarized by the table below.

	ADJACENCY LIST	ADJACENCY MATRIX
Storage	$\mathcal{O}(n + m)$	$\mathcal{O}(n^2)$
Add vertex	$\mathcal{O}(1)$	$\mathcal{O}(n^2)$
Add edge	$\mathcal{O}(1)$	$\mathcal{O}(1)$
Remove vertex	$\mathcal{O}(n + m)$	$\mathcal{O}(n^2)$
Remove edge	$\mathcal{O}(m)$	$\mathcal{O}(1)$

Figure 1: Adjacency Matrix vs Adjacency List

An explanation of these runtimes are provided below:

- An adjacency list requires $\mathcal{O}(n+m)$ since there are n lists inside of the adjacency list. Now for each vertex v_i , there are $\deg(v_i)$ vertices in the i^{th} adjacency list. Since $\sum_i \deg(v_i) = \mathcal{O}(m)$, we conclude that the adjacency list representation of a graph requires $\mathcal{O}(n+m)$ space. On the other hand, the adjacency matrix representation of a graph requires $\mathcal{O}(n^2)$ space since we are storing an $n \times n$ matrix.
- We can add a vertex in constant time in an adjacency list by simply inserting a new list into the adjacency list. On the other hand, to insert a new vertex in an adjacency matrix, we need to increase the dimensions of the adjacency matrix from $n \times n$ to $(n+1) \times (n+1)$. This requires $\mathcal{O}(n^2)$ time since we need to copy over the old matrix to a new matrix.
- We can insert an edge (u, v) into an adjacency list in constant time by simply appending v to the end of u 's adjacency list (and u to the end of v 's adjacency list if the graph is undirected). Similarly, we can insert an edge in an adjacency matrix in constant time by setting $A[u][v]$ to 1 (and also setting $A[v][u]$ to 1 if the graph is undirected).
- Removing a vertex requires $\mathcal{O}(n+m)$ time in an adjacency list since we need to traverse the entire adjacency list and remove any incoming or outgoing edges to the vertex being removed. Similarly, this operation takes $\mathcal{O}(n^2)$ time in an adjacency matrix since we need to traverse the entire matrix to remove incoming and outgoing edges.
- Removing an edge (u, v) requires $\mathcal{O}(m)$ time in an adjacency matrix since we only need to search the adjacency lists of u and v (in the worst case, these vertices have all m edges in their adjacency list). On the other hand, this operation takes constant time in an adjacency matrix since we're just setting $A[u][v]$ to 0.

§3.3 Graph Traversal

Before discussing recapping the two primary types of graph traversal, we will introduce some more terminology.

Definition 3.5. A **connected component** of a graph is a maximally connected subgraph of G . Each vertex belongs to one connected component as does each edge.

There are two primary ways in which we can traverse graphs: using **breadth-first search** or **depth-first search**. These two methods of graph traversal are very similar, and they allow us to explore every vertex in a connected components of a graph.

1. Breadth-first search starts at some source vertex v and all vertices with distance k away from v before visiting vertices with distance $k+1$ from v . This algorithm is typically implemented using a queue, and it can be used to find the shortest path (measured by the number of edges) from the source vertex.

2. Depth-first search starts from a source vertex and keeps on going outward until we cannot proceed any further. We must subsequently backtrack and begin performing the depth-first search algorithm again. This algorithm is typically implemented using a stack, whether it be the data structure or the function call stack.

Both of these algorithms run in $\mathcal{O}(n^2)$ time on an adjacency matrix and $\mathcal{O}(n + m)$ time on an adjacency list.

Since breadth-first search and depth-first search are guaranteed to visit all of the vertices in the same connected component as the starting vertex, we can easily write an algorithm that counts the number of connected components in a graph.

Some C++ code is provided below.

```
/* visited[] is a global Boolean array. */
/* AdjList is a global vector of vectors. */
void dfs(int v) {
    visited[v] = true;
    for (int i = 0; i < AdjList[v].size(); i++) {
        int u = AdjList[v][i];
        if (!visited[u]) {
            dfs(u);
        }
    }
}

int main(void) {
    /* Assume AdjList and other variables have been declared. */
    int numCC = 0;
    for (int i = 0; i < num_vertices; i++) {
        if (!visited[i]) {
            numCC = numCC + 1;
            dfs(i);
        }
    }
}
```

§4 Thursday, February 6, 2020

Today, we'll discuss algorithms to find articulation points and biconnected components.

§4.1 Articulation Points

Definition 4.1. An **articulation point** or **cut vertex** is a vertex in a graph $G = (V, E)$ whose removal (along with any incident edges) would disconnect G .

Definition 4.2. A graph is said to be **biconnected** if the graph not have any articulation points.

For example, consider the following graph:

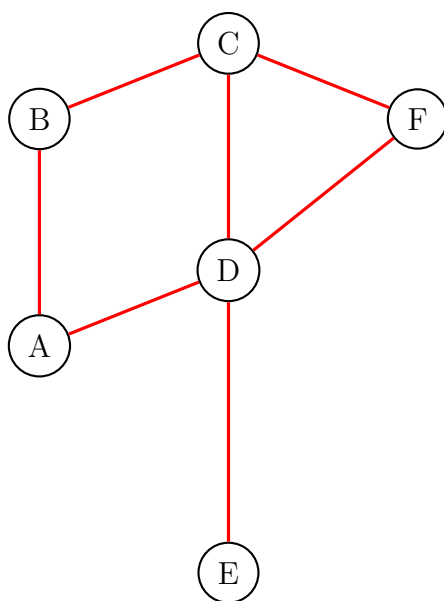


Figure 2: A Graph with an Articulation Point

In the diagram above, Vertex D is an articulation point. To see why, note that if we were to remove Vertex D (and any incident edges to D) from the graph, then we would end up with two connected components: the first component would contain the vertices A, B, C , and F , whereas the second component would only contain the vertex E .

Why are articulation points important? One example in which searching for articulation points is important is in the study of networks. In a network modeled by a graph, an articulation point represents a vulnerability: it is a single point whose failure would split the network into two or more components (preventing communication between the nodes in different networks).

How do we find an articulation points? The brute force algorithm is as follows:

1. Run an $\mathcal{O}(V + E)$ depth-first search or breadth-first search to count the number of connected components in the original graph $G = (V, E)$.

2. For each vertex $v \in V$, remove v from G , and remove any of v 's incident edges. Run an $\mathcal{O}(V + E)$ depth-first search or breadth-first search again, and check if the number of connected components increases. If so, then v is an articulation point. Restore v and any of its incident edges.

This naive algorithm calls the depth-first search or breadth-first search algorithm $\mathcal{O}(V)$ times. Hence, it runs in $\mathcal{O}(V \times (V + E)) = \mathcal{O}(V^2 + VE)$ time.

While this algorithm *works*, it is not as efficient as we can get. We will now describe a linear-time algorithm that runs the depth-first search algorithm just *once* to identify all articulation points and bridges. This algorithm is often accredited to Hopcraft and Tarjan.

In this modified depth-first search, we will now maintain two numbers for each vertex v : **dfs_num(v)** and **dfs_low(v)**. The quantity **dfs_num(v)** represents a label that we will assign to nodes in an increasing fashion. For instance, the vertex from which we call depth-first search would have a **dfs_num** of 0. The subsequent vertex we visit would be assigned a **dfs_num** of 1, and so on.

On the other hand, the quantity **dfs_low(v)**, also known as the **low-link value** of the vertex v , represents the smallest **dfs_num** reachable from that node while performing a depth-first (including itself).

Here's an example. Consider the following directed graph:

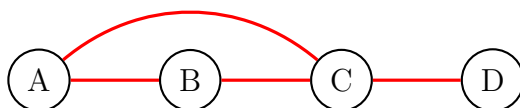


Figure 3: Articulation Point Example

Suppose we perform a depth-first search starting at Vertex A .

- Vertex A will be assigned a **dfs_num** of 0 since this is the first vertex that we're visiting. Moreover, 0 is the smallest **dfs_num** that is reachable from A (all other vertices have their **dfs_num** set to `nil` or `INFINITY`). Hence, we set **dfs_num(A) = 0** and **dfs_low(A) = 0**.
- Next, we visit vertex B . Vertex B is assigned a **dfs_num** of 1 since it's the second vertex we're visiting. Moreover, Vertex B has a **dfs_low** value of 0 since we can reach a vertex with a **dfs_num** value of 0 through the path $B \rightarrow C \rightarrow A$. Note that it would be invalid to say that the path $B \rightarrow A$ causes **dfs_low(B)** to equal 0 since we cannot go backwards in the depth-first search traversal.
- Applying similar reasoning, we find that vertex C ends up with a **dfs_num** value of 2, and it also has a **dfs_low** value of 0 (we can reach vertex A).

- Finally, vertex B ends up with a `dfs_num` value of 3; however, no vertices with a lower `dfs_num` value are reachable from D . Hence, the `dfs_low` value of D is also equal to 3. Note that it is incorrect to say that D has a `dfs_low` value of 0 through the path $D \rightarrow C \rightarrow A$ since we cannot revisit vertices while performing the depth-first search algorithm.

Why do we care about these `dfs_num` and `dfs_low` values? It becomes more clear when we consider the depth-first search tree produced by calling the depth-first search algorithm. The quantity `dfs_low(v)` represents the smallest `dfs_num` value reachable from the current depth-first search spanning subtree rooted at the vertex v . The value `dfs_low(v)` can only be made smaller if there's a back edge (an edge from a vertex v to an ancestor of v) in the depth-first search tree.

This leads us to make the following observation: If there's a vertex u with neighbor v satisfying `dfs_low(v) >= dfs_num(u)`, then we can conclude that vertex u is an articulation point. Note that this makes sense intuitively since it means that the *smallest* numbered vertex that we can ever reach starting from vertex v is greater than or equal to the number we assigned to u . Hence, removing u would disconnect v from any vertex with smaller `dfs_num` than `dfs_num(u)`.

Going back to the previous graph figure, we can note that the following:

$$3 = \text{dfs_num}(D) \geq \text{dfs_low}(C) = 0$$

As stated previously, this implies that Vertex C is an articulation point. Note that removing Vertex C would disconnect the vertices A and B from Vertex D .

Now, there's one special case to this algorithm. The root of the depth-first search spanning tree (the vertex that we choose as the source in the first depth-first search call) is an articulation point only if it has more than one children. This one case is not detected by the algorithm; however, it is easy to check in implementation.

§5 Tuesday, February 11, 2020

§5.1 Articulation Point Algorithm Implementation

Last time, we introduced the algorithm to find articulation points. Recall that if there's a vertex u with neighbor v satisfying $\text{dfs_low}(v) \geq \text{dfs_num}(u)$, then vertex u is an articulation point.

In terms of the depth-first search tree, the quantity $\text{dfs_low}(v)$ is the lowest value that you can reach by going down the depth-first search tree rooted at v and possibly taking a back edge up (we can't visit the immediate parent of v). The inequality $\text{dfs_low}(v) \geq \text{dfs_num}(u)$ implies that we cannot visit any vertex with dfs_num less than $\text{dfs_num}(u)$ when we start a depth-first search from v (there aren't any back edges that go to a vertex visited before vertex u).

Furthermore, recall that the root of the depth-first search tree is an exception — this vertex is an articulation point only if it has more than one child.

When actually implementing this algorithm, we need to be clever in order to maintain a linear time complexity. A pseudocode implementation is provided at <http://www.cs.umd.edu/class/spring2020/cmsc451/biconnected.pdf>.

§5.2 Strongly Connected Components

Recall that an undirected graph $G = (V, E)$ is called **connected** provided that for any pair of vertices $u, v \in V$, there exists a path between u and v .

The corresponding analogue for connectivity in a directed graph is presented below:

Definition 5.1. We call a *directed* graph **strongly connected** if, for every pair of vertices $u, v \in V$, there exists a directed path $u \rightsquigarrow v$.

We're often interested in checking whether or not a graph is strongly connected (e.g. starting from *anywhere* in a directed graph, is it possible to reach *everywhere* else?).

Like connected components in an undirected graph, strongly connected components in a directed graph form a partition of the set of vertices. This is formalized through the following result:

Lemma 5.2 (Klekleinberg and Tardos, 3.17)

For any two nodes s and t in a directed graph, their strong components are either identical or disjoint.

Proof. Consider any two nodes s and t that are mutually reachable. We claim that the strong components containing s and t are identical. This is clearly true due to the definition of a strongly connected component — for any node v , if s and v are

mutually reachable, then t and v are mutually reachable as well (we can always go $s \rightsquigarrow t \rightsquigarrow v$). Similarly, if t and v are mutually reachable, then s and v must be mutually reachable as well.

Conversely, suppose s and t are *not* mutually reachable. Then there cannot be a node v in the strong component of both s and t . Suppose such a node v existed. Then s and v would be mutually reachable, and v and t would be mutually reachable. But this would imply that s and t are mutually reachable, which is a contradiction. \square

A brute force algorithm to check whether a graph is strongly connected is presented below:

1. For each vertex $v \in V$ in our input graph $G = (V, E)$, perform a depth-first search starting with vertex v .
2. If there exists some vertex u that we cannot reach from a vertex v , then we can conclude that G is not strongly connected.
3. If we finish iterating over all vertices with no issues, we can conclude that our graph is strongly connected.

Since we perform $\mathcal{O}(V)$ depth-first search calls in the algorithm above, the runtime of this algorithm runs in $\mathcal{O}(V \times (V + E)) = \mathcal{O}(V^2 + VE)$ time on an Adjacency List. However, this is not as efficient as we can get.

It turns out that we can solve the problem of determining whether a graph is strongly connected in linear time using two depth-first search calls. Before presenting this algorithm, we'll need the following terminology:

Definition 5.3. Given a directed graph $G = (V, E)$, the **transpose graph** of G is the directed graph G^T obtained by reversing the orientation of each edge from (u, v) to (v, u) .

A summary of Kosaraju's algorithm is presented below:

1. Pick an arbitrary vertex $v \in V$ in our initial graph $G = (V, E)$.
2. Perform a depth-first search from v and verify that every other vertex in the graph can be reached from v . If there exists some vertex u that cannot be reached from v , then we can immediately conclude that G is not strongly connected.
3. Compute G^T , the transpose graph of G . Perform a depth-first search on G^T with the same source vertex v . If we can reach every vertex from v in G^T as well, then we can conclude that G is strongly connected.

Why does this work? Because a graph and its transpose always have the same connected components (for each directed $u \rightsquigarrow v$ path, we can just go in the reverse direction).

Now, this algorithm tells us *if* a graph is strongly connected; however, it doesn't tell us *what* the strongly connected components are (i.e. if a graph has many strongly connected components, which component does an arbitrary vertex v belong in?). To answer this question, we'll first present a way to classify the edges in a depth-first search tree.

We will see that this edge-classification system is very closely related to finding strongly connected components in a graph.

§5.3 Classifying Edges in a DFS Tree

While performing a depth-first search traversal, we generate a depth-first search spanning tree. In particular, this DFS tree's root is the source vertex from which we started the DFS traversal, and we add the edge (u, v) if we traverse the edge (u, v) during the DFS procedure.

Within the depth-first search tree, we can classify each edge into exactly one of four disjoint categories:

1. **Tree edges** are edges traversed by the depth-first search traversal (i.e. they are neighbors in the original graph, and we go from one of the vertices to the other). These are the only type of edges that are actually explored.
2. **Back edges** are edges that are part of a cycle in the original graph. In particular, a back edge is an edge (u, v) that we discover when we have started (but not finished) a DFS traversal from v and we're exploring the neighbors of vertex u .
3. **Forward edges** and **cross edges** are edges of the form (u, v) where we have started (but not finished) the depth-first search traversal from u , and we find a vertex v that has already been fully explored.

§6 Thursday, February 13, 2019

Last time, we started discussing strongly connected components, and we presented an edge-classification system. Today, we'll show how we can use our edge-classification system to identify what vertices lie in strongly connected components.

§6.1 Kosaraju's Algorithm

Now, we'll show how we can identify strongly connected components in linear time. The algorithm that we will describe is [Kosaraju's algorithm](#).

The pseudocode corresponding to the algorithm is presented below:

```

procedure kosarajuSCC(graph G) {

    for each node v in G:
        color v gray.

    let L be an empty list.
    for each node v in G:
        if v is gray:
            run DFS starting at v, appending each node to list L when it
            is we've finished processing that node.

    let G' be the transpose graph of G

    for each node v in G':
        color v gray.

    let SCC be a new array of length n.
    let index = 0

    for each node v in L, in reverse order:
        if v is gray:
            run DFS on v in G' and set scc[u] = index
            for each node u visited during the traversal.
            index = index + 1

    return scc
}

```

How is this working?

1. Firstly, we look at the original graph $G = (V, E)$, and we perform a depth-first search on the components of G . Once we've finished visiting each node v , we append v to the end of a list L (we are placing the vertices into L in [reverse-topological order](#)). The list L ends up being sorted in reverse-order of

finishing time. The entire purpose of this first depth-first search traversal is to be able to number the vertices according to their finish time.

2. Next, we'll construct the transpose graph G^T , and we'll iterate over L in reverse-order. Recall that the strongly connected components in G^T are exactly the same as those in G . Also, we mark each
3. For each vertex v we visit in L , if we haven't already call DFS on while iterating over L , any set of vertices that we visit forms a strongly connected component.

Some more intuition is provided below.

Note that, when performing a depth-first search in G^T in post-order from a node v , the depth-first search first visits nodes that can reach v followed by v itself, and finally followed by nodes that cannot reach v . On the other hand, when we perform a depth-first search in pre-order on the original graph G from a node v , the depth-first search first visits v , followed by any nodes reachable from v , and finally the nodes that are not reachable from v .

§6.2 Topological Sorting

Next, we'll begin discussing our next problem. First, we'll present a couple of definitions.

Definition 6.1. A **directed acyclic graph**, also known as a “DAG,” is (as its name suggests), a directed graph that doesn't have any cycles.

Definition 6.2. A **topological sort** of a directed acyclic graph $G = (V, E)$ is a linear ordering of all its vertices such that if G contains an edge (u, v) , then u precedes v in the ordering.

Clearly, a graph with a cycle cannot be topologically sorted — we wouldn't be able to order the vertices that form the cycle.

It's important to remember that, unlike number sorting algorithms, topological sorts are not unique. Each graph G can have multiple valid topological sorts.

Topological sorts are really helpful when we're considering a graph that represents precedences among events or objects. Here are a few examples:

Example 6.3 (Figure 22.7, CLRS)

Professor Bumstead gets dressed in the morning. The professor must wear certain garments before others (e.g. socks before shoes), whereas other pairs of items can be put on in any order (e.g. socks and pants). We can represent this situation with a directed acyclic graph $G = (V, E)$ in which a directed edge (u, v) indicates that garment u must be donned before garment v . The professor can topologically sort this graph in order to get a valid order for getting dressed.

Here's another example.

Example 6.4 (Pick-up Sticks)

The game of *pick-up sticks* involves two players. The game consists of dropping a bundle of sticks. Subsequently, players take turns trying to remove sticks without disturbing any of the others. In order to model this game, we can use a directed graph $G = (V, E)$ in which each vertex represents a stick. We place a directed edge (u, v) between sticks u and v if stick u is on top of stick v . By topologically sorting the graph, we can find a valid way to pick up the sticks on top first.

Now, we've seen a couple of examples in which topological sorts can be useful, but how do we perform a topological sort?

It turns out we can topologically sort a graph in linear time. We will present two algorithms.

Firstly, we present **Kahn's algorithm**, which relies on the following fact:

Proposition 6.5

Every directed acyclic graph has at least one vertex with in-degree 0.

Proof. Suppose not. For each vertex v , we can move backwards through an incoming edge. But due to the finiteness of the graph G and absence of a cycle, this process must eventually terminate. The vertex we terminate must have in-degree 0. \square

Now that we've established this fact, a summary of Kahn's algorithm is presented below:

1. Enqueue all vertices with in-degree 0 into a priority queue Q . At least one such vertex must exist due to Proposition 6.5.
2. Let L be an empty list. This will store our vertices in topologically sorted order.
3. While the Q isn't empty, extract the next vertex u from Q . Remove the vertex u from the original graph G along with any incident edges, and add u to L . If this removal causes another vertex v to have in-degree 0, then enqueue v into Q .
4. Once the while-loop terminates, L will contain every vertex in topologically sorted order.

While we won't prove correctness for this algorithm, it should be a little clear as to why it works. Since we're always choosing vertices with in-degree 0, we know that there is no other vertex that should come before the vertex we're choosing. Hence, the vertices we pick are always "safe." This is pretty similar to the selection sort algorithm used to sort numbers in which we repeatedly pick the minimum element in an array to place at the front of the array. This algorithm runs in $\mathcal{O}(V + E)$ time

on an adjacency list.

Here's a second algorithm that correctly performs a topological sort. This is just a slight modification to the DFS algorithm.

1. Let $G = (V, E)$ be our original graph. Mark each vertex $v \in V$ as “unvisited.”
2. For each unvisited vertex, call `DFS(v)`, and prepend v into an array A once we've finished visiting all of its neighbors.
3. Once we've finished visiting every vertex in G , the array A will be in reverse-topological order. We can reverse the array in linear time, and we're done.

This algorithm runs in $\mathcal{O}(V + E)$ time as the runtime is dominated by our depth-first search calls.

Once again, we won't prove correctness of this algorithm, but it should be clear why this algorithm works. Our call to depth-first search will end pushing vertices with out-degree 0 onto the stack first (because they won't have any more neighbors to visit), which are always safe to place at the end of the topological ordering since no vertex is “greater” than them. This is followed by other vertices in ascending order of out-degree.

A C++ implementation of this algorithm is presented below:

```
vector<vector<int>>> AdjList; /* Our graph. */
vector<int> toposort; /* Global array to store topological sort. */
bool visited[10000];

void dfs(int u) {
    visited[u] = true;
    for (int i = 0; i < AdjList[u].size(); i++) {
        int v = AdjList[u][i];
        if (!visited[v]) {
            dfs(v);
        }
    }
    toposort.push_back(u);
}

int main(void) {
    memset(visited, false, sizeof(visited));
    for (int i = 0; i < V; i++) {
        if (!visited[i]) {
            dfs(i);
        }
    }
    reverse(toposort.begin(), toposort.end());
    /* Topological sort is complete. */
}
```

§6.3 Bipartite Graphs

Finally, we'll discuss bipartite graphs.

Definition 6.6. A graph $G = (V, E)$ is called **bipartite** if we can partition its vertex set V into two disjoint sets U and V such that each edge $(u, v) \in E$ has one endpoint in U and the other endpoint in V .

Here's an equivalent definition that we sometimes like to use:

Definition 6.7. A graph $G = (V, E)$ is said to be **bipartite** if we can color each vertex either black or white such that no two adjacent vertices have the same color.

In order to test whether a graph is bipartite, we can perform a graph search in which we color vertices as we go along. Although we can use either breadth-first search or depth-first search for this check, breadth-first search is often the more natural approach. Pretty much, we start by coloring the source vertex with value 0, color the direct neighbors of the source vertex with 1, the neighbors of the neighbors of the source vertex with color 0, and so on. If we encounter any violations (i.e. two adjacent vertices with the same color) as we go along, then we can conclude that the given graph is not bipartite.

A C++ implementation is provided below:

```
vector<vector<int>> AdjList; /* Our graph. */

bool isBipartite(int src) {
    queue<int> q;
    q.push(src);
    vector<int> color(V, INFINITY);
    color[src] = 0;
    bool isBipartite = true;

    while (!q.empty() && isBipartite) {
        int u = q.front(); q.pop();

        for (int i = 0; i < AdjList[u].size(); i++) {
            int v = AdjList[u][i];
            if (color[v] == INFINITY) {
                /* We haven't colored v yet. */
                color[v] = 1 - color[u];
                q.push(v);
            } else if (color[v] == color[u]) {
                /* We've found a violation. */
                isBipartite = false;
                break;
            }
        }
    }
    return isBipartite;
}
```

The runtime of this algorithm is dominated is $\mathcal{O}(V + E)$ on an adjacency list since we're just performing a breadth-first search.

Another useful fact regarding bipartite graphs is the following:

Fact 6.8. A graph is bipartite if and only if it has no odd cycles (i.e. cycles of length 3, 5, 7, etc).

§7 Tuesday, February 18, 2020

Last time, we finished graph algorithms. Today, we'll begin **greedy algorithms**, which are a class of algorithms that repeatedly make “locally optimal” decisions in an attempt to find a globally optimal solution.

§7.1 The Union-Find Data Structure

§7.1.1 Motivating the Union-Find Data Structure

Before we introduce Kruskal's algorithm, we'll need to first introduce a data structure known as the **union-find** or **disjoint-set** data structure. Why? Because this data structure is used in the implementation of Kruskal's algorithm, which is one of the two minimum spanning tree algorithms we will be talking about.

The union-find data structure consists of a collection of disjoint sets (i.e. a set of sets). Each disjoint set is uniquely determined by a **set representative**, which is some member of the set. In most applications, it doesn't actually matter which member of the set is used as the representative; all we care is that, if we ask for the representative of a set twice without making any modifications, we should get the same answer both times.

The union-find data structure supports the following operations:

1. The **MAKE-SET**(x) operation creates a new set whose only member is x . Since x is the only member of this newly created set, x must also be the representative of this set. Moreover, since we require the sets to be disjoint, we require that x not already be in some other set.
2. The **UNION**(x , y) operation unites the two sets that contain the elements x and y . More precisely, if S_x and S_y are the sets containing x and y , then we remove both of these sets from our collection of sets, and we form a new set $S \stackrel{\text{def}}{=} S_x \cup S_y$, which is subsequently added to the collection of sets. What element becomes the representative of the new set? Typically, if S_x was originally larger than S_y , then we make the representative of S_x the representative of S . Otherwise, we make the representative of S_y the representative of S .
3. The **FIND-SET**(x) method takes in an element x and returns the representative of the set containing x . Note that this means that **FIND-SET**(x) might return x itself (if x is the representative of its set).

There are several applications of the union-find data structure. One of the many applications arises when we are trying to determine the connected components in an undirected graph. In particular, we can answer queries of the form “Are vertices u and v in the same connected component?” with a quick running time by using this data structure.

Consider the following pseudocode:

Input: A graph G .

Output: Nothing. This function is called as a preprocessing step
in order to use the function SAME-COMPONENT(u, v).

```
CONNECTED-COMPONENTS( $G$ ) {  
  for each vertex  $v \in G$  {  
    MAKE-SET( $v$ )  
  }  
  for each edge  $(u, v) \in G$  {  
    if FIND-SET( $u$ )  $\neq$  FIND-SET( $v$ ) {  
      UNION( $u, v$ )  
    }  
  }  
}
```

Input: Two vertices u and v . CONNECTED-COMPONENTS(G) must be
called prior to using this function.

Output: True if u and v are in the same connected component;
otherwise false.

```
SAME-COMPONENT( $u, v$ ) {  
  return (FIND-SET( $u$ ) == FIND-SET( $v$ ))  
}
```

How does these functions work?

- We use a single union-find data structure that is initially empty. At first, we create a new disjoint set for each vertex. Each disjoint set in our union-find data structure will represent a connected component in our graph.
- Next, we traverse every edge in our graph G . For each edge (u, v) , we merge the two disjoint sets containing u and v (since they must be in the same component).
- Finally, we can call the SAME-COMPONENT function with two vertices u and v which simply compare the representatives of the sets u and v are in to determine whether the two vertices are in the same component.

§7.1.2 Implementation of the Union-Find Data Structure

In our connected components example, we use a union-find data structure, but we never explain how the functions **MAKE-SET**, **UNION**, or **FIND-SET** are implemented. In this section, we'll discuss how to implement these three methods.

Union-find data structures are typically implemented as a **disjoint-set forest** in which each member only points to its parent (the root of each tree is the representative of the disjoint set, and it is its own parent). The following figure from CLRS illustrates this idea:

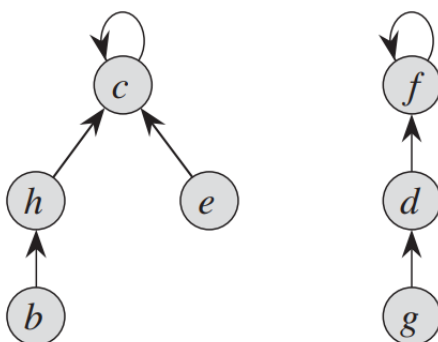


Figure 4: A Disjoint Forest

The disjoint-set forest above represents the two sets $\{c, h, b, e\}$ with representative c and $\{f, d, g\}$ with representative f . Note that the parent of any representative is itself.

How do we keep track of the parent of each vertex? This is easy — we can just include an array called **parent** as a part of our data structure implementation. For any vertex v , we can store the parent of v in **parent**[v].

Now, we will discuss two heuristics to improve the running-time of various union-find operations. The first heuristic, known as the **union by rank heuristic**, is a heuristic that is applied when performing the **UNION** operation. In particular, this heuristic specifies to make the root of the tree with fewer nodes to point to the root of the tree with more nodes. Why? Because following the union by rank heuristic minimizes the overall depth of the resulting tree.

The following diagram illustrates the resulting tree that comes from performing the **UNION** operation on two elements in the disjoint sets from the previous figure:

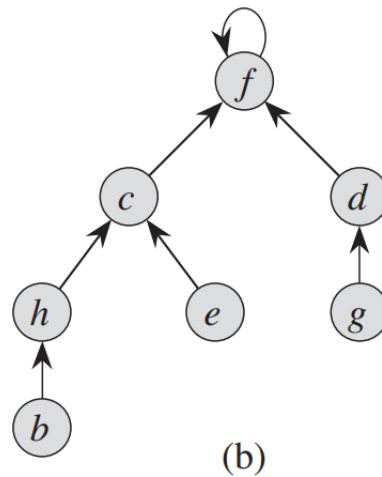


Figure 5: Our Disjoint Forest after performing UNION

Note that f is the representative of the resulting tree since we make the tree with fewer nodes point to the root of the tree with more nodes.

It would be computationally expensive to keep on recomputing the number of roots in each tree whenever we perform a UNION operation. Thus, we can instead just maintain an array **rank** which stores an upper bound on the height of each node. During a UNION operation, we simply make the root with a smaller rank point to the root with the larger rank.

The second heuristic, known as **path compression** is a heuristic that is used during FIND-SET operations to make each node on the find path point directly to the root. This technique is fairly easy to implement, and its purpose is to keep the depth of the tree small.

A C++ implementation of the union-find data structure is presented below:

```

/* An implementation of the union-find data structure. */
class UnionFind {
private:
    vector<int> parent;
    vector<int> rank;
public:
    /* A constructor to initialize a union-find data structure with
       capacity N. */
    UnionFind(int N) {
        parent.assign(N, 0);
        rank.assign(N, 0);

        /* Each vertex is initially its own parent. */
        for (int i = 0; i < N; i++) {
            parent[i] = i;
        }
    }
}
  
```



```

    }

    /* findSet(u) returns the representative of the set that u belongs
       to. */
    int findSet(int u) {
        if (parent[u] == u) {
            /* u is the representative of its set. */
            return u;
        }
        /* Path compression heuristic. */
        return parent[u] = findSet(parent[u]);
    }

    /* inSameSet(u, v) returns true if u and v are in the same set; false
       otherwise. */
    bool inSameSet(int u, int v) {
        /* We compare the set representatives. */
        return findSet(u) == findSet(v);
    }

    /* Union the sets that u and v belong in. */
    void unionSet(int u, int v) {
        if (!inSameSet(u, v)) {
            int rep1 = findSet(u);
            int rep2 = findSet(v);
            /* Union by rank heuristic. */
            if (rank[rep1] > rank[rep2]) {
                parent[rep2] = rep1;
            } else {
                parent[rep1] = rep2;
                if (rank[rep1] == rank[rep2]) {
                    rank[rep2]++;
                }
            }
        }
    }
};

```

§7.1.3 Analysis of Union-Find Operations

In order to discuss the running time of each of the union-find operations, we will need to use the **inverse Ackermann function**, denoted $\alpha(n)$. For our purposes, all we need to know is that this is an *extremely* slowly growing function (for all practical purposes, its value never exceeds 5).

While we won't derive the bound, we will take it for granted that the UNION and FIND-SET operations run in $\mathcal{O}(\alpha(n))$ time (approximately constant time). A full derivation is provided in CLRS 21.4.

§7.2 The Minimum Spanning Tree Problem

§7.2.1 Problem Statement

The **minimum spanning tree** problem is stated as follows:

“Given a graph $G = (V_1, E_1)$, find a connected subgraph $H = (V_2, E_2)$ such that $V_1 = V_2$ and the quantity

$$\sum_{(u,v) \in E_2} \text{weight}(u, v)$$

is as minimal as possible.”

The following proposition shows that H will always be a tree:

Proposition 7.1

If $H = (V_2, E_2)$ is a connected subgraph of $G = (V_1, E_1)$ with the properties described above, then H is a tree.

Proof. By definition, H must be connected. Thus, it suffices to show that H doesn’t have any cycles. Suppose H contained a cycle C . Let e be an edge on C , and consider the graph $H \setminus \{e\}$. This graph is still connected since removing an edge in a cycle can’t disconnect a graph, but this graph is also “cheaper” than H ; this is a contradiction. \square

A simple brute force algorithm to find the minimum spanning tree would work by generating each possible spanning tree and storing the generated tree if its cost is less than our previously stored minimum. Unfortunately, this algorithm is not feasible since graph has exponentially many different spanning trees. Thus, we are compelled to look for more efficient solutions.

Today, we will discuss **Kruskal’s algorithm** and **Prim’s algorithm**, both of which are used to find the minimum spanning tree of a graph.

Both of these algorithms are classified as **greedy algorithms** — they repeatedly make locally optimal choices in an attempt to find a globally optimal solution.

§7.2.2 Kruskal’s Algorithm

Let S be an initially empty set, and let $G = (V, E)$ be our graph. Kruskal’s algorithm works by iteratively adding edges the least weight to S as long as (u, v) does not form a cycle with any of the other edges in S . The algorithm terminates when adding any edge in $E \setminus S$ to S would result in a cycle.

How do we quickly check if adding an edge (u, v) to S will result in a cycle? This can be done quite easily using the union-find data structure.

The pseudocode for Kruskal’s algorithm is below:

Input: A graph G .

Output: A set of edges that form a minimum spanning tree of G .

```

KRUSKAL( $G$ ) {
    let  $S$  be an empty set.

    for each vertex  $v \in G$  {
        MAKE-SET( $v$ )
    }

    sort the edges in  $G$ .Edges into nondecreasing order by weight.

    for each edge  $(u, v) \in G$ .Edges taken in sorted order {
        if FIND-SET( $u$ ) != FIND-SET( $v$ ) {
            #  $(u, v)$  won't form a cycle.
            Add the edge  $(u, v)$  to  $S$ .
            UNION( $u, v$ )
        }
    }
    return  $S$ 
}

```

As mentioned earlier, there isn't too much to this algorithm:

1. First, we sort the edges in non-decreasing order by weight so that we can traverse the list of edges from lowest weight to highest weight.
2. For each weight we look at, we check whether we can add the weight without adding a cycle. This is done by maintaining a union-find data structure.
3. Finally, we return the set of edges that form our minimum spanning tree.

How fast is Kruskal's algorithm? Firstly, note that the first for-loop performs V MAKE-SET operations. Subsequently, we sort the list of edges; doing so requires $\mathcal{O}(E \log(E))$ time. Finally, the second for-loop performs $\mathcal{O}(E)$ FIND-SET and UNION operations. Putting everything together, we have a runtime of $\mathcal{O}((V + E)\alpha(V))$ time. But since $\alpha(|V|) = \mathcal{O}(\log(V)) = \mathcal{O}(\log(E))$ (where the second equality follows due to the fact that $|E| \geq |V| - 1$ in a connected graph), the total running time of Kruskal's algorithm is $\mathcal{O}(E \log(E))$.

§7.2.3 Prim's Algorithm

Prim's algorithm works by starting with an empty set S and iteratively adding edges to S until our minimum spanning tree is complete. We start by adding an arbitrary vertex to S , and at each step we add a vertex that is connected to some other vertex in S .

Some pseudocode illustrating how Prim's algorithm works is shown below:

```
# Input:  A graph G and a source vertex v.

# Output: A set S containing the edges that represent a minimum
# spanning tree.

PRIM(G, v) {
    let key[1...V] be an array.
    let Q be an empty minimum priority queue

    for each vertex u ∈ G {
        key[u] = ∞
        parent[u] = NIL
        enqueue u into Q.
    }
    key[v] = 0

    # This is the main loop.
    while Q isn't empty {
        let u = EXTRACT-MIN(Q)
        for each vertex v in Adj[u] {
            if v ∈ Q and weight(u, v) < key[v] {
                key[v] = weight(u, v)
                parent[v] = u
            }
        }
    }
}
```

How does this algorithm work?

- We start building our minimum spanning tree from an arbitrary vertex v . This vertex is passed in as a parameter to our function.

- Next, we process enqueue all of our vertices into a minimum priority queue Q which allows us to extract elements with the minimum **key** value in logarithmic time.
- While Q isn't empty, we take the vertex with the smallest **key** value from Q ; denote this vertex by u . Note that, on the first iteration, the vertex we grab is always v .
- For each neighbor v of u , we check whether the edge (u, v) is cheaper than the stored **key** value of v . If so, we update the key value of v to the weight of edge (u, v) . We additionally store the vertex u from which we took v .

The purpose of the minimum priority queue is to iteratively identify the cheapest edge that we can add to our minimum spanning tree. The **key** value of a vertex v represents the “cheapest” amount that we can pay in order to add that vertex to our spanning tree.

Prim's algorithm greedily selects the pair (u, v) in front of the priority queue—which has the minimum weight w —if the end point of this edge, namely v , has not been taken before. When the **while** loop terminates, the minimum spanning tree consists of the set of edges

$$A = \{(v, \text{parent}[v]) \mid v \in V - \{r\} - Q\}.$$

§8 Thursday, February 20, 2020

Today, we'll discuss two more algorithmic problems, both of which have greedy optimal solutions.

§8.1 Interval Scheduling

The first problem we'll discuss is known as the **interval scheduling problem**, which is stated as follows:

Given a pair of parallel arrays `start[1...N]` and `finish[1...N]`, call a set of indices S **compatible** if, for any pair of indices $i, j \in S$, the intervals $(\text{start}[i], \text{finish}[i])$ and $(\text{start}[j], \text{finish}[j])$ are disjoint. Moreover, call a compatible set S **optimal** if its cardinality is maximal. The goal is to find an optimal set.

Why do we care about this problem? For each index $1 \leq k \leq N$, we can interpret the quantity `start[k]` and `finish[k]` as the starting time and ending time of an event. Under this interpretation, our task is to fit as many events as possible into our calendar.

There are several algorithms that we can implement that following the greedy heuristic:

1. One approach is to always select the next available event that always starts the earliest (i.e. keep on picking $\operatorname{argmin}_{k \in \{1, 2, \dots, N\}} \text{start}[k]$), and remove k from our set afterwards. This method, however, is not optimal. A counterexample can be generated by considering the case in which the event with the earliest start time is very very long. By accepting this request, we'll miss out on many other events.
2. A second approach is to keep on picking $\operatorname{argmin}_{k \in \{1, 2, \dots, N\}} \text{finish}[k] - \text{start}[k]$ and remove the index we picked from our set. While this is better than the other approach, this isn't optimal either.
3. A third approach is to pick the next request that finishes first (that is, pick $k = \operatorname{argmin}_{k \in \{1, 2, \dots, N\}} \text{finish}[k]$) over and over again. This algorithm seems similar to our first idea. Surprisingly, however, this is the optimal solution.

Some pseudocode illustrating how this procedure works is presented below:

```
# Input:  A set  $S$  representing the
# Output: An optimal solution  $A$ .
SCHEDULING( $S$ ) {
    let  $A$  be the empty set.
```

```

while S isn't empty {
    let e be the event in S with the smallest finishing time.
    add request e to set A.
    remove any events that aren't compatible with e from S.
}
return A
}

```

Next, we'll prove that the set A returned by this algorithm is an optimal solution.

Proposition 8.1

The set A returned by our algorithm is a compatible set of events.

Proof. On each iteration, we add an event, and we remove any events that *aren't* compatible with the event we just added. Since the compatibility relationship between events is symmetric, we know that we'll never have a pair of incompatible events in A . \square

Now we need to show that the set A produced by this algorithm has maximal cardinality. In order to do so, let \mathcal{O} be an optimal set of intervals. We want to show $|\mathcal{O}| = |A|$.

In other words, if $A = \{i_1, i_2, \dots, i_k\}$ and $\mathcal{O} = \{j_1, j_2, \dots, j_m\}$, then our goal is to show $k = m$.

In order to show that this is true, we need to make use of the following lemma:

Lemma 8.2

For any indices $r \leq k$, we have $\text{finish}[i_r] \leq \text{finish}[j_r]$.

Proof. For brevity, this proof writes $f(k)$ and $s(k)$ represent $\text{finish}[k]$ and $\text{start}[k]$, respectively.

When $r = 1$, the statement holds since our algorithm always picks the index i_1 corresponding to the event with the minimum finish time. Now suppose $f(i_{r-1}) \leq f(j_{r-1})$. We want to show $f(i_r) \leq f(j_r)$. But this is clearly true since $f(j_{r-1}) \leq s(j_r)$ implies $f(i_{r-1}) \leq s(j_r)$. This means that j_r is in the set S of compatible events at the time when the greedy algorithm chooses i_r . Since the greedy algorithm always picks the event with the minimum finish time, we must have $f(i_r) \leq f(j_r)$. \square

Lemma 8.2 means precisely that our greedy algorithm's intervals are finished at least as soon as the corresponding intervals in \mathcal{O} .

We can now prove our original claim:

Proposition 8.3

The set A returned by our greedy algorithm has maximal cardinality.

Proof. If A doesn't have maximal cardinality, then an optimal set \mathcal{O} must have more requests. In other words, we require $m > k$. Applying our previous lemma with $r = k$, we find $f(i_k) \leq f(j - k)$. But since $m > k$, there exists some request j_{k+1} in \mathcal{O} . Since this request starts after the event corresponding to j_k ends, deleting all of the requests that aren't compatible with i_1, \dots, i_k will still contain j_{k+1} . However, this means that the greedy algorithm stops with a request present in set, when it's actually only supposed to stop when S is empty. \square

§8.1.1 Extensions: Minimizing Lateness

Once again, consider the situation in which we have a set of n events that we want to schedule in an interval of time. But now, instead of a start time and a finish time, each event has a *deadline*. We say an event k is **late** if our finish time is greater than its deadline. Moreover, we define the *lateness* of a late event as the difference between the time at which it was finished and the time of the deadline. The objective of this problem is to minimize the number of late events.

The greedy algorithm in this problem is to sort the jobs in increasing order of their deadlines, and schedule them in this order (i.e. we process the events with the earliest deadline first). We will not prove the correctness of this algorithm.

§8.2 Caching

A **cache** is a piece of hardware or software that stores data in a special location so that future requests for that data can be served in a high-speed manner. The idea of caching is to store frequently-used values in a special area so that we can access the values in a quick manner. If a value is *not* cached, then we say that the value is stored in **main memory**.

In order to have an effective cache, it should usually be the case that when we're trying to access a piece of data, it's already present in the cache. Today, we'll talk about a cache maintenance algorithm that determines what to keep in the cache and what to toss out of the cache when new data is brought in.

Our problem is stated as follows:

Let U be a set containing n pieces of data stored in main memory, and let C denote our cache that can hold $k < n$ pieces of memory. Given a sequence of data items d_1, d_2, \dots, d_m drawn from U , we must process them in order and determine which of the k items to keep in the cache. When an item d_i is presented that isn't in C , we say a **cache miss** occurs (we want to minimize these), and we have the option to evict some other

data element in C in exchange for d_i . Thus, our problem consists of computing the minimum number of cache misses necessary to process our data sequence.

Example 8.4 (Caching Example)

Suppose $U = \{a, b, c\}$, and our cache size is $k = 2$. Moreover, suppose we are presented with the sequence

$$a, b, c, b, c, a, b.$$

If the cache initially contains items a and b , then on the third item in the sequence, we can evict a to bring in c , and on the sixth item, we could evict c to bring in a . This results in two total cache misses. It can be shown that no solution can have fewer than two cache misses.

§8.3 Farthest in Future Algorithm

Surprisingly, the solution to the caching problem is fairly short. When data element d_i needs to be brought into the cache, we should always evict the item that is needed the farthest into the future. This is known as the **Farthest-in-Future algorithm**, and it was discovered by Belady.

We won't prove optimality; however, it's important to take note that that greedy algorithms might not always be obvious (why do we evict the element farthest in the future as opposed to the least frequent element?)

§9 Tuesday, February 25, 2020

§9.1 Prefix Codes

One particular class of encoding schemes are **prefix codes**. A prefix code for a set S of letters is a function γ that maps each letter $x \in S$ to some sequence of zeros and ones in such a way that for any $x, y \in S$ with $x \neq y$, the sequence $\gamma(x)$ is not a prefix of the sequence $\gamma(y)$. Why do many encoding schemes fall into this class? Because it removes ambiguity: — if there exists a pair of letters where the bit string that encodes one letter is a prefix of the bit string that encodes the other, then there might be multiple interpretations of the same string.

The ambiguity of encoding schemes that aren't prefix codes is demonstrated through the following example:

Example 9.1 (Ambiguity Morse Code)

In Morse code, we typically encode letters with dashes and dots. For our purpose, we can think of dots and dashes as zeros and ones. Suppose e maps to 0 (a single dot), t maps to 1, and a maps to 01. Then the string 0101 can have several interpretations: it can mean *eta*, *aa*, *etet*, or *aet*. If the morse code were a prefix code, then this problem wouldn't be present.

Now, here's an example illustrating the ease of using a prefix code:

Example 9.2 (Prefix Code Example)

Suppose we have a set $S = \{a, b, c, d, e\}$ with the encoding $\gamma(a) = 11, \gamma(b) = 01, \gamma(c) = 001, \gamma(d) = 10, \gamma(e) = 000$. This defines a prefix code since no encoding is a prefix of any other. The string *cecab* is encoded as 0010000011101, and a recipient of this message can decipher this message to our single unique message.

In order to efficiently decipher a prefix code, we need an effective way to represent the prefix code so that we can easily pick off the codeword. This is typically done with a binary tree in which the leaves of the tree store the characters of our alphabet. How does this work? We interpret the binary codeword for a character as a simple path from the root to that character; the bit 0 tells us to go to the left child, whereas the bit 1 tells us to go to the right child.

The following binary search tree illustrates a prefix code representation:

If we had the sequence 001011101, then we could start at the root, and we'd scan our sequence from left to right. First, we counter two zeros, so we go to the left twice. At this point, we'd be at the vertex labelled 58. Next, we encounter a 1, so we go to the right. Thus, we obtain the character b . Next, we start at the root again, and we follow our procedure again. The next character that we decipher is d . This process continues until there are no more bits to process.

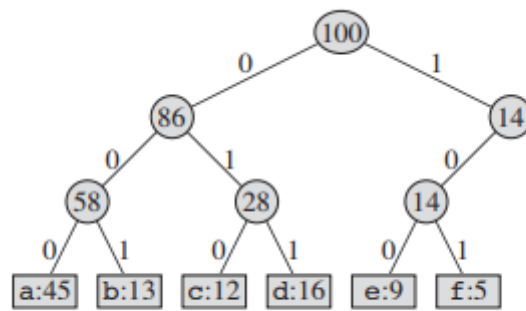


Figure 6: Prefix Code Representation

Given a tree T corresponding to a prefix code, we can now easily compute the number of bits required to encode a file. In particular, for each character c in our alphabet S , we can let $\text{freq}[c]$ denote the frequency of c in our file. Moreover, we can let $d_T(c)$ denote the depth of c 's leaf in the tree. With this notation, the number of bits required to encode a file is given by

$$\sum_{c \in S} \text{freq}[c] \cdot d_T(c).$$

We call this the **cost** of the tree T .

§9.1.1 Constructing a Huffman code

Now that we've introduced prefix codes, we'll talk about an optimal prefix code known as a **Huffman code**, whose tree representation has minimum cost. The algorithm constructing the tree is presented below:

Input: A set C representing the set of all possible characters that
 # might appear in our text, and an array $\text{freq}[]$ in which $\text{freq}[k]$ represents
 # the frequency of the character k in our text.
 # Output: The root of a binary representing our encoding minimum cost.

```
HUFFMAN(C, freq) {
    let Q be a minimum priority queue
    for each element c in C { enqueue c into Q }

    for i = 1 to n - 1 {
        let z be a new node
        x = EXTRACT-MIN(Q)
        y = EXTRACT-MIN(Q)
        z.left = x
        z.right = y
        freq[z] = freq[x] + freq[y]
        insert z into Q.
    }
    return EXTRACT-MIN(Q) /* Return the root of the tree. */
}
```

How does this algorithm work?

1. Firstly, we enqueue all of the characters in C into our minimum priority queue Q .
2. The for-loop repeatedly extracts the two vertices with the lowest frequency and replaces them in the queue with a new node representing their “merger” (parent). The frequency of z is the sum of the frequencies of x and y .
3. After $n - 1$ merges, there’s only one node left in the queue, which is the root of the code tree.

If we’re using a binary heap, then the algorithm runs in $\mathcal{O}(n \log(n))$ time since we perform $\mathcal{O}(n)$ calls to **EXTRACT-MIN**, which is an $\mathcal{O}(\log(n))$ operation.

While we won’t show it, it can be shown that this construction of a tree is optimal. This procedure counts as a greedy algorithm since, at each step, we greedily extract the characters with the lowest frequency.

§9.2 Matrix Multiplication

The next problem we'll discuss is stated as follows:

Given two $n \times n$ matrices A and B , compute the $n \times n$ matrix C whose $(i, j)^{\text{th}}$ entry is defined by $c_{ij} = \sum_{k=1}^n a_{ik}b_{kj}$. In other words, we want to compute the product $C = AB$.

The brute force solution is $\mathcal{O}(n^3)$. In this algorithm, we just use three loops, and we compute each value c_{ij} in C as the summation provided in the problem statement. Of course, we want to do better.

Another idea is to perform a divide-and-conquer technique on the matrix. In particular, we can divide the matrix into four submatrices (top left corner, top right corner, bottom left corner, bottom right corner), and we can calculate the products recursively. The time complexity of this algorithm is given by the recurrence $T(n) = 8T(n/2) + \mathcal{O}(n^2)$. Unfortunately, by Master's Theorem, we know that the solution to this recurrence will be $\mathcal{O}(n^3)$, which isn't any better.

§10 Thursday, February 27, 2020

Last time, we introduced the matrix multiplication problem whose brute force solution runs in $\mathcal{O}(n^3)$ time. Today, we'll introduce **Strassen's algorithm**, a quicker solution to the matrix multiplication problem.

§10.1 Strassen's Algorithm

Suppose we want to compute the matrix product $\mathbf{C} = \mathbf{AB}$. Strassen's algorithm is a divide-and-conquer algorithm that works by partitioning the three matrices \mathbf{A} , \mathbf{B} , and \mathbf{C} into equally sized block matrices as follows:

$$\mathbf{A} = \begin{pmatrix} \mathbf{A}_{1,1} & \mathbf{A}_{1,2} \\ \mathbf{A}_{2,1} & \mathbf{A}_{2,2} \end{pmatrix} \quad \mathbf{B} = \begin{pmatrix} \mathbf{B}_{1,1} & \mathbf{B}_{1,2} \\ \mathbf{B}_{2,1} & \mathbf{B}_{2,2} \end{pmatrix} \quad \mathbf{C} = \begin{pmatrix} \mathbf{C}_{1,1} & \mathbf{C}_{1,2} \\ \mathbf{C}_{2,1} & \mathbf{C}_{2,2} \end{pmatrix}$$

Our naive algorithm would compute the following quantities:

1. $\mathbf{C}_{1,1} = \mathbf{A}_{1,1}\mathbf{B}_{1,1} + \mathbf{A}_{1,2}\mathbf{B}_{2,1}$,
2. $\mathbf{C}_{1,2} = \mathbf{A}_{1,1}\mathbf{B}_{1,2} + \mathbf{A}_{1,2}\mathbf{B}_{2,2}$,
3. $\mathbf{C}_{2,1} = \mathbf{A}_{2,1}\mathbf{B}_{1,1} + \mathbf{A}_{2,2}\mathbf{B}_{2,1}$,
4. $\mathbf{C}_{2,2} = \mathbf{A}_{2,1}\mathbf{B}_{1,2} + \mathbf{A}_{2,2}\mathbf{B}_{2,2}$,

With this construction, however, we require 8 total multiplications to calculate our matrix. Strassen's algorithm works by cleverly rewriting some of these expressions so that we only require 7 multiplications (similar to how Karatsuba's algorithm for large-integer multiplication). More precisely, we define the following matrices:

1. $\mathbf{M}_1 \stackrel{\text{def}}{=} (\mathbf{A}_{1,1} + \mathbf{A}_{2,2})(\mathbf{B}_{1,1} + \mathbf{B}_{2,2})$,
2. $\mathbf{M}_2 \stackrel{\text{def}}{=} (\mathbf{A}_{2,1} + \mathbf{A}_{2,2})\mathbf{B}_{1,1}$,
3. $\mathbf{M}_3 \stackrel{\text{def}}{=} \mathbf{A}_{1,1}(\mathbf{B}_{1,2} - \mathbf{B}_{2,2})$,
4. $\mathbf{M}_4 \stackrel{\text{def}}{=} \mathbf{A}_{2,2}(\mathbf{B}_{2,1} - \mathbf{B}_{1,1})$,
5. $\mathbf{M}_5 \stackrel{\text{def}}{=} (\mathbf{A}_{1,1} + \mathbf{A}_{1,2})\mathbf{B}_{2,2}$
6. $\mathbf{M}_6 \stackrel{\text{def}}{=} (\mathbf{A}_{2,1} - \mathbf{A}_{1,1})(\mathbf{B}_{1,1} + \mathbf{B}_{1,2})$,
7. $\mathbf{M}_7 \stackrel{\text{def}}{=} (\mathbf{A}_{1,2} - \mathbf{A}_{2,2})(\mathbf{B}_{2,1} + \mathbf{B}_{2,2})$.

Note that computing the values of these matrices requires only 7 multiplications (one for each \mathbf{M}_k) instead of the usual 8. We can now express our block matrices in terms of the \mathbf{M}_k matrices as follows:

1. $\mathbf{C}_{1,1} = \mathbf{M}_1 + \mathbf{M}_4 - \mathbf{M}_5 + \mathbf{M}_7$,

2. $C_{1,2} = M_3 + M_5,$
3. $C_{2,1} = M_2 + M_4$
4. $C_{2,2} = M_1 - M_2 + M_3 + M_6.$

We can iterate the procedure of dividing our matrices into blocks recursively until the submatrices are just numbers.

How fast does Strassen's algorithm run? Let $f(n)$ denote the number of multiplication operations we perform on a $2^n \times 2^n$ matrix. By the recursive application of Strassen's algorithm, we find $f(n) = 7f(n-1) + c4^n$, where c is some positive constant that depends on the number of additions performed at each step of the operation. Thus, we find $f(n) = (7 + o(1))^n$. Letting $N = 2^n$, we conclude that Strassen's algorithm runs in $\mathcal{O}(N^{\log_2(7)+o(1)}) \approx \mathcal{O}(N^{2.8074})$ time.

§10.2 Closest Pair of Points

The closest pair of points problem is stated as follows:

Given n points in the plane $P = \{p_1, p_2, p_3, \dots, p_n\}$, find two points p_i and p_j such that the Euclidean distance between p_i and p_j is minimal.

A simple brute force solution is to consider all $\binom{n}{2}$ pairs of points, and keep track of the minimum distance value seen so far (this `minimum` variable would initially be set to ∞). The runtime of this algorithm is $O(n^2)$ since computing the distance between two points is a constant-time operation.

Next class, we'll present a more efficient solution.

§11 Tuesday, March 3, 2020

Recall the closest pair of points problem:

Given n points in the plane $P = \{p_1, p_2, p_3, \dots, p_n\}$, find two points p_i and p_j such that the Euclidean distance between p_i and p_j is minimal.

Last time, we discussed a brute force $\mathcal{O}(n^2)$ solution to this problem. Today, we'll see a more efficient solution. We'll also introduce two new problems.

§11.1 Closest Pair of Points

Our plan is to use a divide and conquer approach. We'll find the closest pair among the points in the “left half” of our plane, and we'll find the closest pair of points in the “right half” of our plane. Using this information, we'll construct our final answer in linear time. If we develop an algorithm with this structure, then our recurrence will have an $\mathcal{O}(n \log(n))$ solution. Note that our “combining” phase is more tricky than it seems: we haven't considered the case in which one point is in the left half of the plane and another point is in the right half of the plane.

Before any recursion begins, we sort all of our points in P by x -coordinate and again by y -coordinate, producing two lists P_x and P_y . Moreover, we define Q to be the set of points in the first $\lceil n/2 \rceil$ positions of P_x (i.e. the “left half” of the plane), and we let R be the set of points in the final $\lfloor n/2 \rfloor$ positions of P_x (i.e. the “right half” of the plane).

With a single for-loop, we can iterate over P_x and P_y and create the lists Q_x , consisting of the points in Q sorted by increasing x -coordinate, Q_y , the set of points in Q sorted by increasing y -coordinate, and analogous lists for R_x and R_y .

Next, we'll discuss how to combine the solutions.

Suppose q_0 and q_1 are returned as a closest pair of points in Q . Moreover, suppose r_0 and r_1 are returned as a closest pair of points in R . How do we combine our solution to get the closest pair of points in the plane? First, we need to introduce some more notation.

Let δ be the minimum distance between q_0 and q_1 and between r_0 and r_1 . That is, let $\delta \stackrel{\text{def}}{=} \min\{d(q_0, q_1), d(r_0, r_1)\}$. We want to figure out whether there exist points $q \in Q$ and $r \in R$ such that $d(q, r) < \delta$ (if no such points exists, then δ is our answer; otherwise, q and r are even closer points in our plane).

Let x^* denote the rightmost x -coordinate in Q , and let L denote the vertical line $x = x^*$. This line L separates the sets Q and R in the sense that any point in Q is either on the line or to the left of the line, and any point in R is strictly to the right of the line.

The following key observation is used to help us combine our solutions:

Proposition 11.1 (Existence of a Better Solution)

If there exists $q \in Q$ and $r \in R$ for which $d(q, r) < \delta$, then each of q and r lies within a distance of δ within the line L .

Proof. Suppose such q and r exist. The inequality $q_x \leq x^* \leq r_x$ implies

$$x^* - q_x \leq r_x - q_x \leq d(q, r) < \delta,$$

which yields

$$r_x - x^* \leq r_x - q_x \leq d(q, r) < \delta.$$

However, this means precisely that each q and r has an x -coordinate within δ of x^* . Hence, they lie within distance δ of L . \square

The immediate consequence of Proposition 11.1 is that, once we've found the two closest pairs of points from our recursive calls, we only need to search for a better solution within a strip of length δ of the line L . We now present a method to do this in linear time.

Let $S \subseteq P$ be the set of points in P within δ of L . Let S_y denote the list consisting of the points in S sorted by increasing y -coordinate. We can now restate Proposition 11.1 in terms of S as follows:

There exist $q \in Q$ and $r \in R$ for which $d(q, r) < \delta$ if and only if there exist $s, s' \in S$ for which $d(s, s') < \delta$.

Now, it can be show that if $s, s' \in S$ have the property that $d(s, s') < \delta$, then s and s' are within 15 positions of each other in the sorted list S_y (proof omitted). While this bound is not tight, the important note is that the distance between the points is is an absolute constant. We can now conclude the algorithm by making a single pass through S_y and, for each $s \in S_y$, computing the distance to the next 15 points in S_y . The runtime of this procedure is linear, so we've successfully figured out how to combine our solutions in linear time.

Thus, the recurrence for our algorithm takes the form

$$T(n) = 2T(n/2) + \mathcal{O}(n).$$

By Master's Theorem, we conclude $T(n) = \mathcal{O}(n \log(n))$.

§11.2 Counting Inversions

Definition 11.2. Given an array A , an **inversion** if a pair of indices (i, j) for which both $i < j$ and $A[i] > A[j]$ hold.

The inversion problem is stated as follows:

Given an array A , count the number of inversions in A .

We can think of the number of inversion in an array as the number of “bubble sort swaps” (swap between pairs of consecutive items) needed in order to sort the array.

Example 11.3 (Reverse-Sorted Array)

The array $A = [3, 2, 1]$ has exactly 3 inversions. Namely, $(1, 2)$, $(1, 3)$ and $(2, 3)$.

Example 11.4 (Unsorted Array)

The array $A = [3, 2, 1, 4]$ also has 3 inversions.

The most obvious solution is to simply use two nested for-loops and increment an `answer` variable every time we find an inversion. Here's an implementation of the brute force solution:

```
int main(void) {
    /* Assume we have initialized an array "A" */
    int answer = 0;
    for (int i = 0; i < N; i++) {
        for (int j = i + 1; j < N; j++) {
            /* The condition i < j is always true. */
            if (A[i] > A[j]) {
                /* (i, j) is an inversion. */
                answer = answer + 1;
            }
        }
    }
    cout << "Number of inversions: " << answer << endl;
}
```

The runtime of this algorithm is $\mathcal{O}(n^2)$, but of course, we want to do better.

Once again, we can take a divide and conquer approach for the inversion problem. More precisely, we can just modify the **MergeSort** algorithm. The key observation is that during the **merge** process of merge sort, if the front of the right sorted sublist is taken rather than the front of the left sorted sublist, then we can say that one or more inversions occur. We increment our inversion counter by the size of the current left sublist since all of those indices cause an inversion with the current element we are looking at in our right sublist.

The runtime of this algorithm is $\mathcal{O}(n \log(n))$ since we're only adding a few constant-time operations to the merge sort procedure.

A full implementation is provided on the next page.

```
int merge(vector<int>& A, int l, int m, int r) {
    vector<int> B(r - l + 1);
    int i = l, j = m + 1, k = 0;
    int inversions = 0;

    while (i <= m && j <= r) {
        if (A[i] <= A[j]) {
            B[k++] = A[i++];
        } else {
            B[k++] = A[j++];
            inversions += (m + 1 - i);
        }
    }

    /* Only one of the following two while-loops
    will be executed. */
    while (i <= m) B[k++] = A[i++];
    while (j <= r) B[k++] = A[j++];

    for (int i = l; i <= r; i++) {
        A[i] = B[i - l];
    }

    return inversions;
}

int mergesort(vector<int> &A, int l, int r) {
    int inversions = 0;
    if (r > l) {
        int m = l + (r - l)/2;
        inversions += mergesort(A, l, m);
        inversions += mergesort(A, m + 1, r);
        inversions += merge(A, l, m, r);
    }
    return inversions;
}

/* (i, j) is an inversion if A[i] > A[j] and i < j.
O(n*log(n)) inversion counting. */
int inversion_count(vector<int>& A) {
    return mergesort(A, 0, A.size() - 1);
}
```

§12 Thursday, March 5, 2020

Today, we'll begin discussing our last divide-and-conquer topic: the Fast Fourier ("four-ee-aye") Transform. The problem that we are trying to solve is stated as follows:

Given two vectors $a = (a_0, a_1, \dots, a_{n-1})$ and $b = (b_0, b_1, \dots, b_{n-1})$, compute the convolution $a \star b$ of a and b .

Before discussing the algorithm that lets us do this, let's first discuss convolutions and why they're important.

§12.1 Convolutions

A **convolution** of two vectors a and b is a method of "combining" the two vectors. More precisely, we define the convolution of the vectors $a = (a_0, a_1, \dots, a_{n-1})$ and $b = (b_0, b_1, \dots, b_{n-1})$ by the vector $c = (c_0, c_2, \dots, c_{2n-2})$ in which

$$c_k = \sum_{(i,j)|i+j=k} a_i b_j.$$

In other words, we have,

$$a \star b = (a_0 b_0, a_0 b_1 + a_1 b_0, \dots, a_{n-1} b_{n-1}).$$

Note that each summand in the k^{th} component of this vector exhausts all possible pairs of indices that sum to k . Moreover, note that the convolution of two n -dimensional vectors produces a $(2n - 1)$ -dimensional vector. However, unlike the vector sum and inner product, the convolution can easily be generalized to vectors of different lengths: if $a = (a_0, a_1, \dots, a_{m-1})$ and $b = (b_0, b_1, \dots, b_{n-1})$, then we define $a \star b$ to be a vector with $m + n - 1$ coordinates, where the k^{th} coordinate is equal to the sum over all $a_i b_j$ in which $i + j = k$, $i < m$ and $j < n$.

Why do we care about the convolution? Here are some examples in which convolutions are useful:

Example 12.1 (Polynomial Multiplication)

Suppose we have two polynomials $A(x) = a_0 + a_1x + a_2x^2 + \dots + a_{m-1}x^{m-1}$ and $B(x) = b_0 + b_1x + b_2x^2 + \dots + b_{n-1}x^{n-1}$ and we wish to compute the product $C(x) = A(x) \cdot B(x)$. In order to do so, we can define the vectors $a = (a_0, a_1, \dots, a_{m-1})$ and $b = (b_0, b_1, \dots, b_{n-1})$ and compute the convolution $c = a \star b$. In the polynomial $C(x)$, the coefficient of x^k is equal to the k^{th} component of c .

Example 12.2 (Combining Histograms)

Suppose we're studying a population of people, and we have two histograms. The first histogram shows the annual income of all the men in the population, and the other shows the annual income of all the women. We would like to produce a new histogram showing for each k the number of pairs (M, W) for which man M and woman W have a combined income of k . This problem can be restated as a convolution. More precisely, let $a = (a_0, \dots, a_{m-1})$ and $b = (b_0, \dots, b_{n-1})$ be our histograms, and let c_k denote the number of (m, w) pairs with combined income k . Observe that c_k is the number of ways to choose a man with income a_i and woman with income b_j with $i + j = k$. This quantity is given by a convolution.

Example 12.3 (Sum of Independent Random Variables)

If one is familiar with probability theory, then they may have encountered a theorem which tells us that the probability distribution function for the sum of two random variables is a convolution of the distributions of the summands.

Now that we've motivated the importance of convolutions, we'll now discuss how to compute convolutions efficiently. For simplicity, we consider the case in which our two vectors have equal length (i.e. $m = n$); however, our results hold for vectors of unequal length.

Computing a convolution efficiently is more difficult than it seems. If, for each k , we just calculate the sum $\sum_{(i,j)|i+j=k} a_i b_j$ and use it as the k^{th} coordinate in our convolution vector, we end up with an $\mathcal{O}(n^2)$ algorithm. Fortunately, we can do better — the **fast Fourier Transform** allows us to compute convolutions in $\mathcal{O}(n \log(n))$ time.

§12.2 The Fast Fourier Transform

In order to compute convolutions quickly, we will make use of the connection between the convolution and polynomial multiplication. However, rather than using convolutions to perform polynomial multiplication, we will exploit the connection in the opposite direction.

Given two vectors $a = (a_0, \dots, a_{n-1})$ and $b = (b_0, \dots, b_{n-1})$, we define $A(x)$ and $B(x)$ to be the polynomials $a_0 + a_1x + \dots + a_{n-1}x^{n-1}$ and $b_0 + b_1x + \dots + b_{n-1}x^{n-1}$, respectively. Under this interpretation, we wish to compute the product $C(x) = A(x)B(x)$ in $\mathcal{O}(n \log(n))$ time. From there, we can simply “read off” the convolution directly from the coefficients of $C(x)$.

Now, instead of multiplying A and B directly, we can treat them as functions of the variable x and multiply them with the following three steps:

1. Choose $2n$ values x_1, x_2, \dots, x_{2n} and evaluate $A(x_j)$ and $B(x_j)$ for each $j = 1, 2, \dots, 2n$.

2. Now for each index $1 \leq j \leq 2n$, we can easily compute $C(x_j)$. In particular, $C(x_j)$ is equal to the product of the two numbers $A(x_j)$ and $B(x_j)$.
3. Finally, we need to recover the polynomial C from its values on x_1, x_2, \dots, x_{2n} . Since any polynomial of degree d is fully determined by a set of $d + 1$ or more points, this is clearly possible. Since each A and B have degree at most $n - 1$, their product C has degree at most $2n - 2$. Thus, it can be reconstructed from the values $C(x_1), C(x_2), \dots, C(x_{2n})$ that we computed earlier.

This approach to multiplying polynomials sounds promising, but there are a couple of issues we need to address. Evaluating the polynomials A and B on a single point takes $\Omega(n)$ operations (using Horner's method¹, and our plan calls for performing $2n$ such evaluations. This brings us back up to quadratic time immediately. Moreover, we need a way to quickly reconstruct the polynomial C from the points $C(x_1), C(x_2), \dots, C(x_{2n})$.

We address these two issues separately.

§12.2.1 Polynomial Evaluation

We need to evaluate the polynomials A and B on $2n$ different points quickly. The key idea to doing this quickly is to find a set of $2n$ points x_1, \dots, x_{2n} that are related in some way so that the work in evaluating A and B on all of them can be shared across different evaluations. A set that works very well for us is the roots of unity.

Definition 12.4. An n^{th} root of unity is a number z satisfying the equation $z^n = 1$.

It can be shown that there are n n^{th} roots of unity. Moreover, these roots are given by $e^{2k\pi i/n}$ for $k = 0, 1, \dots, n - 1$. Clearly, each of these complex numbers satisfy our definition since

$$(e^{2k\pi i/n})^n = e^{2k\pi i} = (e^{2\pi i})^k = 1^k = 1.$$

For our numbers x_1, \dots, x_{2n} on which to evaluate A and B , we will choose the $(2n)^{\text{th}}$ roots of unity, and we propose a recursive procedure to compute A on each of the $(2n)^{\text{th}}$ roots of unity. For simplicity, we henceforth assume that n is a power of 2.

Let $A_{\text{even}}(x)$ and $A_{\text{odd}}(x)$ be two polynomials that consist of the even and odd coefficients of A , respectively. That is, we have,

$$A_{\text{even}}(x) = a_0 + a_2x + a_4x^2 + \dots + a_{n-2}x^{(n-2)/2},$$

and

$$A_{\text{odd}}(x) = a_1 + a_3x + a_5x^2 + \dots + a_{n-1}x^{(n-2)/2}.$$

By simple algebra, we can see that we can express $A(x)$ as

¹https://en.wikipedia.org/wiki/Horner%27s_method

$$A(x) = A_{\text{even}}(x^2) + xA_{\text{odd}}(x^2),$$

which demonstrates that we can compute $A(x)$ in a constant number of operations provided that we already have A_{even} and A_{odd} . Now suppose we evaluate both A_{even} and A_{odd} on the n^{th} roots of unity. This is an exact replica of the problem we face with A and the $(2n)^{\text{th}}$ roots of unity, except the input is half as large; the degrees of our two polynomials are $(n-2)/2$ rather than $n-1$. Moreover, we have n roots of unity rather than $2n$. Thus, we can perform these evaluations recursively in time $T(n/2)$ for each of A_{even} and A_{odd} , for a total of $2T(n/2)$ time.

But, how do we perform these evaluations? This can be done with $\mathcal{O}(n)$ additional operations given the results from the recursive calls on A_{even} and A_{odd} . Let $\omega = e^{2\pi ik/2n}$ be a $(2n)^{\text{th}}$ root of unity for some integer k . The quantity ω^2 is equal to $e^{2\pi ki/n}$, which is an n^{th} root of unity.

Thus, when we go to compute $A(\omega) = A_{\text{even}}(\omega^2) + \omega \cdot A_{\text{odd}}(\omega^2)$, we find that both evaluations on the right-hand side have been performed in a recursive step, which means that we can compute $A(\omega)$ in a constant number of operations. Repeating for each of the $2n$ roots of unity is therefore $\mathcal{O}(n)$ additional operations.

Therefore, our bound $T(n)$ on the number of operations satisfies $T(n) = 2T(n/2) + \mathcal{O}(n)$, which gives us the desired $\mathcal{O}(n \log(n))$ bound for the first step of our algorithm.

§12.2.2 Polynomial Interpolation

Next, we'll discuss how to Now, we've seen how to evaluate A and B on the set of all $(2n)^{\text{th}}$ roots of unity using $\mathcal{O}(n \log(n))$ operations. Also, we can clearly perform the second step of our algorithm naively in linear time. Thus, to conclude the algorithm for multiplying A and B , we need to reconstruct the polynomial C from its values on the $(2n)^{\text{th}}$ roots of unity in $\mathcal{O}(n \log(n))$ time.

The reconstruction of C can be achieved by defining an appropriate polynomial and evaluating it at the $(2n)^{\text{th}}$ roots of unity. This is exactly what we've just seen how to do using $\mathcal{O}(n \log(n))$ operations, so we'll do it here again. This requires an additional $\mathcal{O}(n \log(n))$ operations, and it concludes our algorithm.

Consider a polynomial $C(x) = \sum_{s=0}^{2n-1} c_s x^s$ that we want to reconstruct from its values at the $C(\omega_{s,2n})$ at the $(2n)^{\text{th}}$ roots of unity. Define a new polynomial $D(x) \stackrel{\text{def}}{=} \sum_{s=0}^{2n-1} d_s x^s$ where $d_s = C(\omega_{s,2n})$. We now consider the values of $D(x)$ at the $(2n)^{\text{th}}$ roots of unity:

$$\begin{aligned}
D(\omega_{j,2n}) &= \sum_{s=0}^{2n-1} C(\omega_{s,2n}) \omega_{j,2n}^s \\
&= \sum_{s=0}^{2n-1} \left(\sum_{t=0}^{2n-1} c_t \omega_{s,2n}^t \right) \omega_{j,2n}^s.
\end{aligned}$$

Now since $\omega_{s,2n} = (e^{2\pi i/2n})^s$, we get that

$$D(\omega_{j,2n}) = \sum_{t=0}^{2n-1} c_t \sum_{s=0}^{2n-1} \omega_{t+j,2n}^s.$$

However, note that for any $(2n)^{\text{th}}$ root of unity $\omega \neq 1$, we have $\sum_{s=0}^{2n-1} \omega^s = 0$. Thus, the only term that of the last line's outer sum that is not equal to 0 is for c_t such that $\omega_{t+j,2n} = 1$. This happens precisely when $t + j = 2n - j$.

It follows immediately that for any polynomial $C(x) = \sum_{s=0}^{2n-1} c_s x^s$ and corresponding polynomial $D(x) = \sum_{s=0}^{2n-1} C(\omega_{s,2n}) x^s$, we have $c_s = \frac{1}{2n} D(\omega_{2n-s,2n})$.

Thus, we can reconstruct the polynomial C from its values on the $(2n)^{\text{th}}$ roots of unity, and the coefficients of C are the coordinates in the convolution vector $a \star b$ that we were originally seeking. Therefore, we are done.

§13 Tuesday, April 7, 2020

§13.1 Subset Sum Problem

Today, we'll discuss another classical dynamic programming problem, known as the **subset sum problem**.² The subset problem is stated as follows:

Suppose we are given n items $\{1, 2, \dots, n\}$ each with nonnegative weight w_i . We are also given a bound W . How do we select a subset S of the items so that $\sum_{i \in S} w_i$ is maximized subject to $\sum_{i \in S} w_i \leq W$?

In other words, given n items each with nonnegative weights, what's the closest we can get to a weight of W without going over?

Example 13.1 (Subset Sum Example)

Suppose $n = 3$ with $w_1 = 2$, $w_2 = 3$, $w_3 = 4$, and $W = 5$. The solution to this instance of the subset problem is $\boxed{5}$ — it is optimal to choose w_1 and w_2 .

Does a greedy solution work? One greedy rule might be to sort the items in ascending order by weight and always pick the item with maximal weight that hasn't been taken yet. However, this greedy rule fails in [Example 13.1](#); we'll end up with a total weight of 4, which is sub-optimal.

We demonstrate how we can use dynamic programming to solve this problem. Recall that the main principles of dynamic programming are to come up with a recurrence so that we can relate the problem we want to solve to “smaller” subproblems. The tricky issue is determining what a good set of subproblems consists of.

One general strategy in dynamic programming is to consider subproblems consisting of only the first i “requests,” or items. We can use this strategy here. Formally, let $\text{OPT}(i)$ denote the best possible solution using only the subset $\{1, \dots, i\}$ of the original set of items. Now, the key to this problem is to concentrate on an optimal solution and consider two different cases, depending on whether or not the last item n we processed is part of this optimum solution or not. Let \mathcal{O} denote an optimal solution.

If $n \notin \mathcal{O}$, then $\text{OPT}(n) = \text{OPT}(n - 1)$. This is obvious — if the last item we processed isn't a part of our optimal solution, then the optimal solution using only $\{1, 2, \dots, n - 1\}$ shouldn't change when n is included (we now have the option to take n , but we don't want n anyways!).

²The subset sum problem is a special case of another classical dynamic programming problem, known as the **knapsack problem**. In the knapsack problem, each item $1 \leq i \leq n$ has a value v_i and weight w_i . For each item, we want to assign a number $x_i \in \{0, 1\}$ so that $\sum_i x_i v_i$ is maximized subject to $\sum_i x_i w_i \leq W$.

The only other case we have to consider is the case in which $n \in \mathcal{O}$. What we need to find is a simple recursion that tells us the best possible value we can obtain for solutions containing the last request n . Note that accepting request n does not immediately imply that we have to reject any other requests. Instead, it means that for the subset of requests $S \subseteq \{1, 2, \dots, n-1\}$ that we will accept, we have less available weight left. More precisely, we will have $W - w_n$ weight left for the remaining set of items we accept.

This suggests that we need more subproblems: we cannot just use the value $\text{OPT}(n-1)$ when we're including item n since the combined weight of the items in $\text{OPT}(n-1)$ and item n might exceed W . What we precisely need is the best solution using the first $n-1$ items when the total weight allowed is $W - w_n$. Thus, we require many more subproblems: one for each initial set $\{1, 2, \dots, i\}$ of the items, and each possible value for the remaining weight available w .

More precisely, if each of our items $1, 2, \dots, n$ have integer weights w_i and our maximum weight bound is W , then we have a subproblem for each $i = 0, 1, \dots, n$ and each integer $0 \leq w \leq W$. We henceforth use $\text{OPT}(i, w)$ to denote the value of the optimal solution using the subset of the items $\{1, 2, \dots, i\}$ with maximum allowed weight w . That is,

$$\text{OPT}(i, w) = \max_{S \subseteq \{1, 2, \dots, i\}} \sum_{j \in S} w_j \quad \text{subject to} \quad \sum_{j \in S} w_j \leq w.$$

Using this new set of subproblems, we can note that the final answer we want is $\text{OPT}(n, W)$. Moreover, we can express the value $\text{OPT}(i, w)$ as an expression from smaller problems. These results are summarized below:

1. If $n \notin \mathcal{O}$, then $\text{OPT}(n, W) = \text{OPT}(n-1, W)$ since we can ignore the item n .
2. If $n \in \mathcal{O}$, then $\text{OPT}(n, W) = w_n + \text{OPT}(n-1, W - w_n)$ since we now want to use the remaining capacity $W - w_n$ in an optimal way across the first $n-1$ items.

If $W < w_n$ for some item n , then we require $\text{OPT}(n, W) = \text{OPT}(n-1, W)$ since we aren't allowed to take item n due to our constraint. Now that we've considered both cases, we can get the optimum solution by simply taking the better of these two options. Therefore, we obtain the following recurrence:

$$\text{OPT}(i, W) = \begin{cases} \max(\text{OPT}(i-1, w), w_i + \text{OPT}(i-1, W - w_i)) & \text{if } w_i \leq W \\ \text{OPT}(i-1, w) & \text{otherwise.} \end{cases}$$

Also, we have the base cases $\text{OPT}(i, w) = 0$ provided that $i = 0$ since we aren't allowed to take any items when $i = 0$.

With our recurrence and base cases established, we are done.