

# CMSC 451

## Design and Analysis of Computer Algorithms<sup>1</sup>

David M. Mount  
Department of Computer Science  
University of Maryland  
Fall 2017

---

<sup>1</sup> Copyright, David M. Mount, 2017, Dept. of Computer Science, University of Maryland, College Park, MD, 20742. These lecture notes were prepared by David Mount for the course CMSC 451, Design and Analysis of Computer Algorithms, at the University of Maryland. Permission to use, copy, modify, and distribute these notes for educational purposes and without fee is hereby granted, provided that this copyright notice appear in all copies.

# Lecture 1: Introduction to Algorithm Design

**What is an algorithm?** This course will focus on the study of the design and analysis of algorithms for discrete (as opposed to numerical) problems. We can define *algorithm* to be:

Any well-defined computational procedure that takes some values as *input* and produces some values as *output*.

The concept of a “well-defined computational procedure” dates back to ancient times. In fact, the word “algorithm” is derived from the Latin form of the Persian scholar Muhammad ibn Musa al-Khwarizmi, who lived in ninth century A.D. Al-Khwarizmi codified procedures for numerous arithmetic operations (such as addition, multiplication, and division with Arabic numerals) and algebraic and trigonometric operations (such as computing square roots and computing the digits of  $\pi$ ).

**Why study algorithm design?** While the study of algorithms predates digital computers, the field really took off with the advent of computers. The use of asymptotic (big-Oh) notation became popular in the 1960’s and 1970’s as a means to provide a rigorous mathematical measure of an algorithm’s running time. This evolved into the field of *computational complexity*, which seeks to categorize computational problems according to their complexity. This gave rise to the study of NP-Hard problems.

The field has also led to the development of general techniques for the design of efficient algorithms, such as divide-and-conquer, greedy algorithms, dynamic programming, and so on.

From a more practical perspective, algorithm design and analysis is often the first step the development of software for tricky combinatorial problems. Asymptotic analysis is used to identify promising solutions, which can then be prototyped in order to determine which methods perform best.

**Course Overview:** This course will consist of a number of major sections. The first will be a short review of some preliminary material, including asymptotics, summations and recurrences, sorting, and basic graph algorithms. These have been covered in earlier courses, and so we will breeze through them pretty quickly. Next, we will consider a number of common algorithm design techniques, including greedy algorithms, dynamic programming, and augmentation-based methods (particularly for network flow problems).

Most of the emphasis of the first portion of the course will be on problems that can be solved efficiently, in the latter portion we will discuss intractability and NP-hard problems. These are problems for which no efficient solution is known. Finally, we will discuss methods to approximate NP-hard problems, and how to prove how close these approximations are to the optimal solutions.

**Issues in Algorithm Design:** Algorithms are mathematical objects (in contrast to the much more concrete notion of a computer program implemented in some programming language and executing on some machine). As such, we can reason about the properties of algorithms mathematically. When designing an algorithm we need to be concerned both with its *correctness* and *efficiency*.

Intuitively, an algorithm’s efficiency is a function of the amount of computational resources it requires, measured typically as execution time and the amount of space, or memory, that the algorithm uses. The amount of computational resources can be a complex function of the size and structure of the input set. In order to reduce matters to their simplest form, it is common to consider efficiency as a function of *input size*, which is usually represented by the symbol  $n$ .

**Worst-case complexity:** Among all inputs of the same size, what is the *maximum* running time?

**Average-case complexity:** Among all inputs of the same size, what is the *expected* running time? This expectation is computed assuming that the inputs are drawn from some given *probability distribution*. The choice of distribution can have a significant impact on the final conclusions.

**Asymptotic Notation:** Asymptotic O-notation (“big-O”) provides us with a way to simplify the messy functions that often arise in analyzing the running times of algorithms. Suppose that we analyze two algorithms, and find that they have running times of:

$$T_1(n) = 3.9n + 4.17 \log n + 3.5n^2 \quad \text{and} \quad T_2(n) = \max(4.6n(\log n)^4, 6.4n^3 - 3n \log_3 n).$$

Which of these algorithms is better? Asymptotic analysis is based on (1) focusing on the growth rate by considering the performance as the value of  $n$  increases to infinity and (2) ignoring constant factors, which tend to rely on secondary issues such as programming style and machine architecture. We’ll give the formal definition later, but intuitively we can say that  $T_1$  grows on the order of  $n^2$  and  $T_2$  grows on the order of  $n^3$ , that is

$$T_1(n) = O(n^2) \quad \text{and} \quad T_2(n) = O(n^3).$$

Formally, we say that  $f(n)$  is  $O(g(n))$  if there exist constants  $c > 0$  and  $n_0 \geq 0$  such that,  $f(n) \leq c \cdot g(n)$ , for all  $n \geq n_0$ . Thus, big-O notation can be thought of as a way of expressing a sort of *fuzzy* “ $\leq$ ” relation between functions, where by fuzzy, we mean that constant factors are ignored and we are only interested in what happens as  $n$  tends to infinity.

Another way to think about asymptotics is in terms of limits. An alternative definition is that  $f(n)$  is  $O(g(n))$  if

$$\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} \leq c, \quad \text{for some constant } c \geq 0.$$

For example, we can say that  $T_1(n)$  is  $O(n^2)$  since

$$\lim_{n \rightarrow \infty} \frac{3.9n + 4.17 \log n + 3.5n^2}{n^2} = \lim_{n \rightarrow \infty} \left( 3.9 \frac{1}{n} + 4.17 \frac{\log n}{n^2} + 3.5 \right) = 3.5,$$

since in the limit  $1/n$  and  $\log n/n^2$  both tend to zero in the limit.

Big-O notation has a number of relatives, which are useful for expressing other sorts of relations. These include  $\Omega$  (“big-omega”),  $\Theta$  (“theta”),  $o$  (“little-oh”),  $\omega$  (“little-omega”). Let  $c$  denote an arbitrary positive *constant* (not 0, not  $\infty$ , and not depending on  $n$ ). Intuitively, each represents a form of “asymptotic relational operator”:

Notation	Relational Form	Limit Definition
$f(n)$ is $o(g(n))$	$f(n) \prec g(n)$	$\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = 0$
$f(n)$ is $O(g(n))$	$f(n) \preceq g(n)$	$\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = c$ or $0$
$f(n)$ is $\Theta(g(n))$	$f(n) \approx g(n)$	$\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = c$
$f(n)$ is $\Omega(g(n))$	$f(n) \succeq g(n)$	$\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = c$ or $\infty$
$f(n)$ is $\omega(g(n))$	$f(n) \succ g(n)$	$\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = \infty.$

By far, the most commonly arising functions in algorithm analysis are of one of three forms. Given a constant  $a > 0$ , they are  $(\log n)^a$  (called *polylogarithmic*),  $n^a$  (called *polynomial*), and  $a^n$  (called *exponential*). For any  $a, b, c$ , such that  $a, b > 0$  and  $c > 1$  we have the following relative order:

$$(\log n)^a \prec n^b \prec c^n.$$

To keep matters simple, we will focus almost exclusively on **worst-case analysis** measured using **asymptotic analysis**. in this course. You should be mindful, however, that worst-case analysis is not always the best way to analyze an algorithm's performance. For example, some algorithms have the property that they run very fast on typical inputs but might run extremely slowly (perhaps hundreds to thousands of times slower) on a very small fraction of *pathological* inputs. For such algorithms, an average case analysis may be a much more accurate reflection of the algorithm's true performance.

**Describing Algorithms:** Throughout out this course, when you are asked to present an algorithm, this means that you need to do three things:

**Present the Algorithm:** Give a clear, simple, and unambiguous description of the algorithm (in plain English prose or pseudo-code, for example). A guiding principal here is to remember that your description will be read by a human, and **not** a compiler. Obvious technical details should be kept to a minimum so that the key computational issues stand out.

**Prove its Correctness:** Present a justification (that is, an informal proof) of the algorithm's correctness. This justification may assume that the reader is familiar with the basic background material presented in class. Try to avoid rambling about obvious or trivial elements and focus on the key elements. A good proof provides a high-level overview of what the algorithm does, and then focuses on any tricky elements that may not be obvious.

**Analyze its Efficiency:** Present a worst-case analysis of the algorithms efficiency, typically it running time (but also its space, if space is an issue). Sometimes this is straightforward and other times it might involve setting up and solving a complex recurrence or a summation. When possible, try to reason based on algorithms that you have seen. For example, the recurrence  $T(n) = 2T(n/2) + n$  is common in divide-and-conquer algorithms (like Mergesort) and it is well known that it solves to  $O(n \log n)$ .

Note that your presentation does not need to be in this order. Often it is good to begin with an explanation of how you derived the algorithm, emphasizing particular elements of the design that establish its correctness and efficiency. Then, once this groundwork has been laid down, present the algorithm itself. If this seems to be a bit abstract now, don't worry. We will see many examples of this process throughout the semester.

**Background Information:** I will assume that you have familiarity with the information from a basic algorithms course, such as CMSC 351. As is indicated in the syllabus, it is expected that you have knowledge of:

- Basic programming skills (programming with loops, pointers, structures, recursion)
- Discrete mathematics (proof by induction, sets, permutations, combinations, and probability)
- Understanding of basic data structures (lists, stacks, queues, trees, graphs, and heaps)
- Knowledge of sorting algorithms (MergeSort, QuickSort, HeapSort) and basic graph algorithms (minimum spanning trees and Dijkstra's algorithm)
- Basic calculus (manipulation of exponentials, logarithms, differentiation, and integration)

## Lecture 2: Graph Basics

**Graphs and Digraphs:** A graph  $G = (V, E)$  is a structure that represents a discrete set  $V$  objects, called *vertices* or *nodes*, and a set of pairwise relations  $E$  between these objects, called *edges*. Edges may be *directed* from one vertex to another or may be *undirected*. The term "graph" means an undirected graph, and directed graphs are often called *digraphs* (see Fig. 1). Graphs and digraphs provide a flexible mathematical model for numerous application problems involving binary relationships between a

discrete collection of object. Examples of graph applications include *communication* and *transportation networks*, *social networks*, *logic circuits*, *surface meshes* used for shape description in computer-aided design and geographic information systems, *precedence constraints* in scheduling systems.

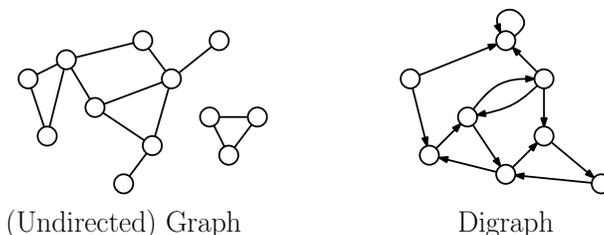


Fig. 1: Graphs and digraphs.

**Definition:** An *undirected graph* (or simply *graph*)  $G = (V, E)$  consists of a finite set  $V$  and a set  $E$  of *unordered pairs* of distinct vertices.

**Definition:** A *directed graph* (or *digraph*)  $G = (V, E)$  consists of a finite set  $V$  and a set  $E$  of *ordered pairs* of vertices.

Observe that multiple edges between the same two vertices are not allowed, but in a directed graph, it is possible to have two oppositely directed edges between the same pair of vertices. For undirected graphs, *self-loop* edges are not allowed, but they are allowed for directed graphs. Directed graphs and undirected graphs are different objects mathematically. Certain notions (such as path) are defined for both, but other notions (such as connectivity and spanning trees) may be defined only for one.

**Graph and Digraph Terminology:** Given an edge  $e = (u, v)$  in a digraph, we say that  $u$  is the *origin* of  $e$  and  $v$  is the *destination* of  $e$ . Given an edge  $e = \{u, v\}$  in an undirected graph,  $u$  and  $v$  are called the *endpoints* of  $e$ . The edge  $e$  is *incident* on (meaning that it touches) both  $u$  and  $v$ . Given two vertices in a graph or digraph, we say that vertex  $v$  is *adjacent* to vertex  $u$  if there is an edge  $\{u, v\}$  (for graphs) or  $(u, v)$  (for digraphs).

In a digraph, the number of edges coming out of  $v$  is called its *out-degree*, denoted  $\text{out-deg}(v)$ , and the number of edges coming in is called its *in-degree*, denoted  $\text{in-deg}(v)$ . In an undirected graph we just talk about the *degree* of a vertex as the number of incident edges, denoted  $\text{deg}(v)$ .

When discussing the size of a graph, we typically consider both the number of vertices and the number of edges. The number of vertices is typically written as  $n$ , and the number of edges is written as  $m$ . (**Beware:** There are many different conventions. The number of vertices may be expressed as  $n$ ,  $v$ ,  $V$ , or  $|V|$ , and the number of edges may be expressed as  $m$ ,  $e$ ,  $E$ , or  $|E|$ ).

Here are some basic combinatorial facts about graphs and digraphs. We will leave the proofs to you. Given a graph with  $n$  vertices and  $m$  edges then:

**In a graph:**

**Number of edges:**  $0 \leq m \leq \binom{n}{2} = n(n-1)/2 = O(n^2)$ .

**Sum of degrees:**  $\sum_{v \in V} \text{deg}(v) = 2m$ .

**In a digraph:**

**Number of edges:**  $0 \leq m \leq n^2$ .

**Sum of degrees:**  $\sum_{v \in V} \text{in-deg}(v) = \sum_{v \in V} \text{out-deg}(v) = m$ .

Notice that generally the number of edges in a graph may be as large as quadratic in the number of vertices. However, the large graphs that arise in practice typically have much fewer edges. A graph is said to be *sparse* if  $m$  is  $O(n)$ , and *dense*, otherwise. When giving the running times of algorithms, we will usually express it as a function of both  $n$  and  $m$ , so that the performance on sparse and dense graphs will be apparent.

**Paths and Cycles:** A *path* in a graph or digraph is a sequence of vertices  $\langle v_0, \dots, v_k \rangle$  such that  $(v_{i-1}, v_i)$  is an edge for  $i = 1, \dots, k$ . The *length* of the path is the number of edges,  $k$ . A path is *simple* if all vertices and all the edges are distinct. A *cycle* is a path containing at least one edge and for which  $v_0 = v_k$ . A cycle is *simple* if its vertices (except  $v_0$  and  $v_k$ ) are distinct, and all its edges are distinct.

A graph or digraph is said to be *acyclic* if it contains no simple cycles. An acyclic connected graph is called a *free tree* or simply *tree* for short (see Fig. 2). (The term “free” is intended to emphasize the fact that the tree has no root, in contrast to a *rooted tree*, as is usually seen in data structures.) An acyclic undirected graph (which need not be connected) is a collection of free trees, and is called a *forest*. An acyclic digraph is called a *directed acyclic graph*, or *DAG* for short (see Fig. 2).

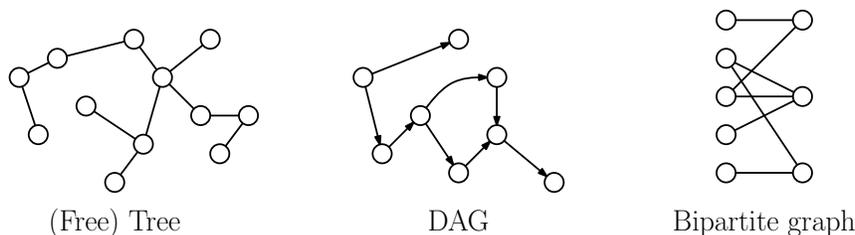


Fig. 2: Illustration of common graph terms.

A *bipartite graph* is one in which the vertices of a graph can be partitioned into two disjoint subsets, denoted  $V_1$  and  $V_2$ , such that all the edges have one endpoint in  $V_1$  and one in  $V_2$  (see Fig. 2). Note that every cycle in a bipartite graph contains an even number of edges.

We say that  $w$  is *reachable* from  $u$  if there is a path from  $u$  to  $w$ . Note that every vertex is reachable from itself by a trivial path that uses zero edges. An undirected graph is *connected* if every vertex can reach every other vertex. (Connectivity is a bit messier for digraphs, and we will define it later.) The subsets of mutually reachable vertices partition the vertices of the graph into disjoint subsets, called the *connected components* of the graph. In digraphs the notion of reachability is a bit different, because it is possible for  $u$  to reach  $w$  but not vice versa. A digraph is said to be *strongly connected* if for each  $u$  and  $w$ , there is a path from  $u$  to  $w$  and a path from  $w$  to  $u$ .

**Representations of Graphs and Digraphs:** There are two common ways of representing graphs and digraphs. First we show how to represent digraphs. Let  $G = (V, E)$  be a digraph with  $n = |V|$  and let  $m = |E|$ . We will assume that the vertices of  $G$  are indexed  $\{1, 2, \dots, n\}$ .

**Adjacency Matrix:** An  $n \times n$  matrix defined for  $1 \leq v, w \leq n$ .

$$A[v, w] = \begin{cases} 1 & \text{if } (v, w) \in E \\ 0 & \text{otherwise.} \end{cases}$$

(See Fig. 3.) If the digraph has weights we can store the weights in the matrix. For example if  $(v, w) \in E$  then  $A[v, w] = W(v, w)$  (the weight on edge  $(v, w)$ ). If  $(v, w) \notin E$  then generally  $W(v, w)$  need not be defined, but often we set it to some “special” value, e.g.  $A(v, w) = -1$ , or  $\infty$ . (By  $\infty$  we mean some number which is larger than any allowable weight.)

It might come as a surprise, but there are a number of interesting relationships between the use of matrices to represent graphs and the matrices that arise in linear algebra to represent linear

transformations. For example, the eigenvalues of the adjacency matrix of a graph provide a lot of information about the structure of the graph.

**Adjacency List:** An array  $Adj[1 \dots n]$  of pointers where for  $1 \leq v \leq n$ ,  $Adj[v]$  points to a list (e.g., a singly or doubly linked list) containing the vertices that are adjacent to  $v$  (i.e., the vertices that can be reached from  $v$  by a single edge). If the edges have weights then these weights may also be stored in the linked list elements (see Fig. 3).

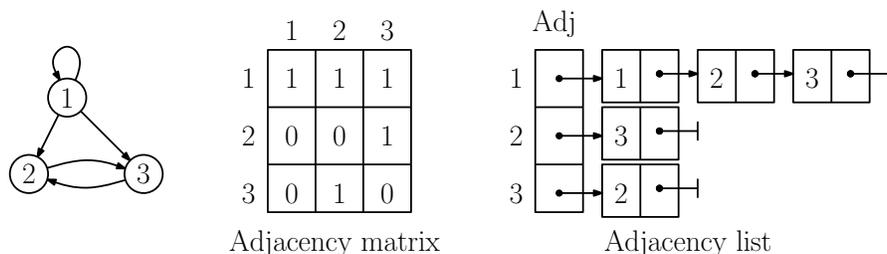


Fig. 3: Adjacency matrix and adjacency list for digraphs.

We can represent undirected graphs using exactly the same representation, but we will store each edge twice. In particular, we representing the undirected edge  $\{v, w\}$  by the two oppositely directed edges  $(v, w)$  and  $(w, v)$  (see Fig. 4). Notice that even though we represent undirected graphs in the same way that we represent digraphs, it is important to remember that these two classes of objects are mathematically distinct from one another.

This can cause some complications. For example, suppose you write an algorithm that operates by marking edges of a graph. You need to be careful when you mark edge  $(v, w)$  in the representation that you also mark  $(w, v)$ , since they are both the same edge in reality. When dealing with adjacency lists, it may not be convenient to walk down the entire linked list, so it is common to include *cross links* between corresponding edges.

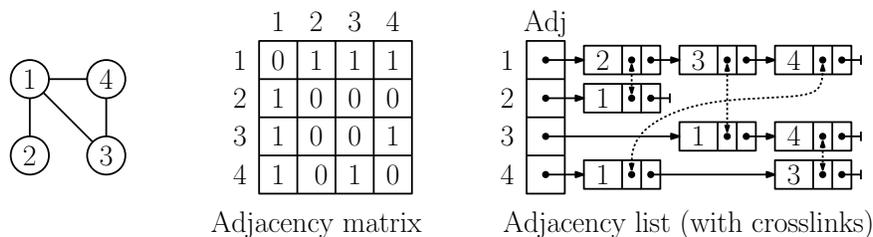


Fig. 4: Adjacency matrix and adjacency list for graphs.

An adjacency matrix requires  $\Theta(n^2)$  storage, and an adjacency list requires  $\Theta(n + m)$  storage. The  $n$  arises because there is one entry for each vertex in  $Adj$ . Since each list has  $\text{out-deg}(v)$  entries, when this is summed over all vertices, the total number of adjacency list records is  $\Theta(m)$ . For most applications, the adjacency list representation is standard.

**Depth-First Search:** One of the most important basic operations on a graph is to systematically visit all its vertices. These traversals naturally impose a type of tree structure (or generally a forest) on the graph, and trees are usually much easier to reason about than general graphs.

We are given a graph  $G = (V, E)$ , which may be directed or undirected. We employ four auxiliary arrays. To avoid revisiting the same vertex, we maintain a mark for each vertex: undiscovered, discovered, finished. Additional information can be stored as part of the traversal process:

**Discovery time:**  $d[u]$  indicates the time when vertex  $u$  was discovered, which coincides with the moment that the DFS process is started at this vertex.

**Finish time:**  $f[u]$  indicates the time when vertex  $u$  is finished processing. At this point, all of  $u$ 's neighboring nodes have been visited, and indeed, everything reachable from  $u$  has been discovered and possibly finished.

**Predecessor pointer:**  $p[u]$  indicates the vertex that discovered  $u$ . Each edge of the form  $(p[u], u)$  is a tree edge in the DFS recursion tree.

DFS induces a tree structure. In order to handle instances where not all vertices are reachable from the starting vertex, we include a main program that invokes DFS whenever an undiscovered vertex is encountered. The main program is shown in code block below and the recursive DFSvisit function is shown in the next code block. (Fig. 5 illustrates the execution on an undirected graph, and Fig. 6 shows an example on a directed graph.)

---

Depth-First Search (Main Program)

```

DFS(G) {                                     // main program
    time = 0
    for each (u in V)                         // initialization
        mark[u] = undiscovered

    for each (u in V)
        if (mark[u] == undiscovered)         // undiscovered vertex?
            DFSVisit(u)                       // ...start a new search here
}

```

---

DFS Visit (Process a single node)

```

DFSVisit(u) {                                // perform a DFS search at u
    mark[u] = discovered                     // u has been discovered
    d[u] = ++time
    for each (v in Adj(u)) {
        if (mark[v] == undiscovered) {      // undiscovered neighbor?
            pred[v] = u
            DFSVisit(v)                       // ...visit it
        }
    }
    mark[u] = finished                       // we're done with u
    f[u] = ++time
}

```

---

**Analysis:** The running time of DFS is  $O(n + m)$ . We'll do the analysis for undirected graphs. First observe that if we ignore the time spent in the recursive calls, the main DFS procedure runs in  $O(n)$  time. Each vertex is visited exactly once in the search, and hence the call `DFSVisit()` is made exactly once for each vertex. We can just analyze each one individually and add up their running times. Ignoring the time spent in the recursive calls, we can see that each vertex  $u$  can be processed in  $O(1 + \deg(u))$  time (the "+1" is needed in case the degree is 0). Thus the total time used in the procedure is

$$T(n) = n + \sum_{u \in V} (1 + \deg(u)) = n + \left( \sum_{u \in V} \deg(u) \right) + n = 2n + m = O(n + m).$$

A similar analysis holds if we consider DFS for digraphs.

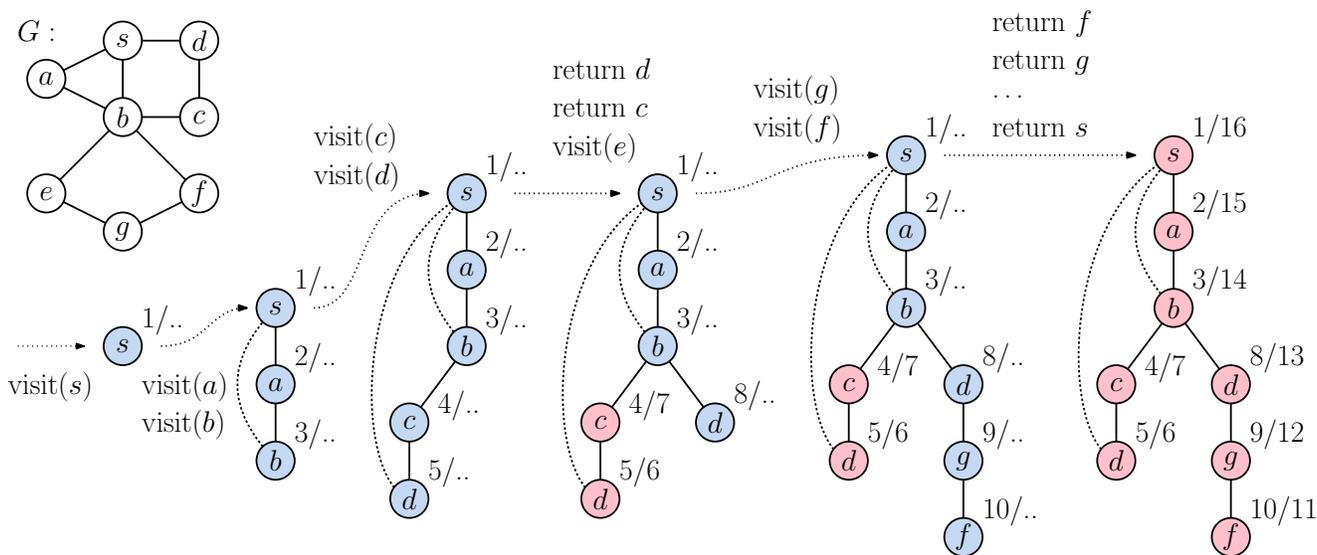


Fig. 5: Depth-first search on an undirected graph. (Blue nodes are discovered, and pink nodes are finished. Each node  $u$  is labeled with the values  $d[u]/f[u]$ .)

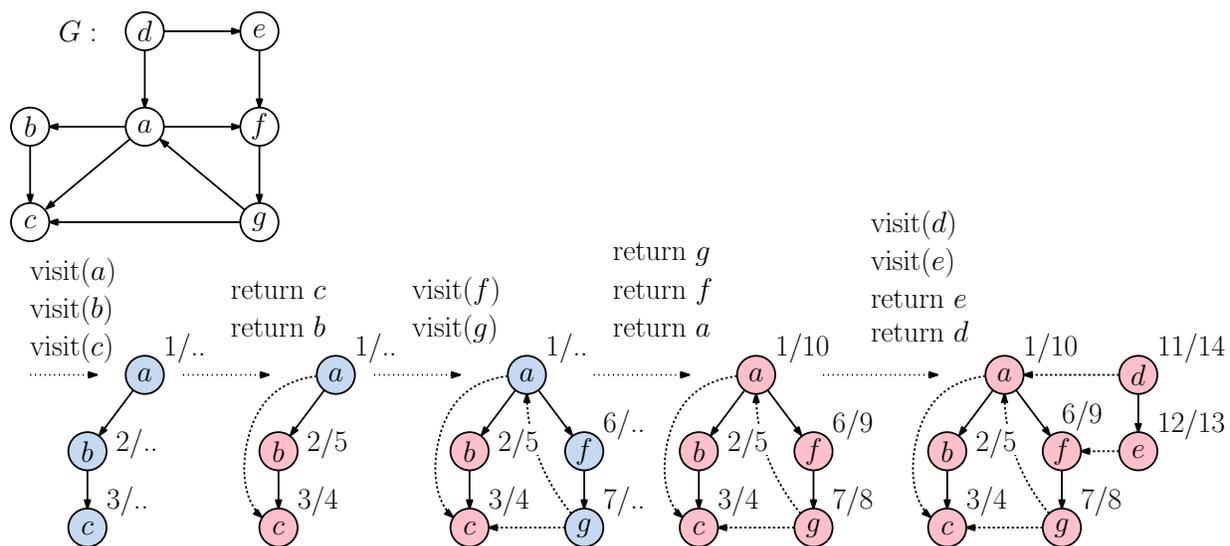


Fig. 6: Depth-first search on a directed graph. (Blue nodes are discovered, and pink nodes are finished. Each node  $u$  is labeled with the values  $d[u]/f[u]$ .)

**Parenthesis Lemma and Edge Types:** DFS naturally imposes a tree structure (actually a collection of trees, or a forest) on the structure of the graph. This is just the recursion tree, where the edge  $(u, v)$  arises when processing vertex  $u$  we call  $\text{DFSVisit}(v)$  for some neighbor  $v$ . The hierarchical structure naturally imposes a nesting structure on the discovery-finish time intervals. This is described in the following lemma (and illustrated in Fig. 7(a)).

**Lemma:** (Parenthesis Lemma) Given a graph  $G = (V, E)$  (directed or undirected), and any DFS tree for  $G$  and any two vertices  $u, v \in V$ :

- $u$  is a descendant of  $v$  iff  $[d[u], f[u]] \subseteq [d[v], f[v]]$ .
- $u$  is an ancestor of  $v$  iff  $[d[u], f[u]] \supseteq [d[v], f[v]]$ .
- $u$  and  $v$  are unrelated (in terms of ancestor/descendant) iff  $[d[u], f[u]]$  and  $[d[v], f[v]]$  are disjoint.

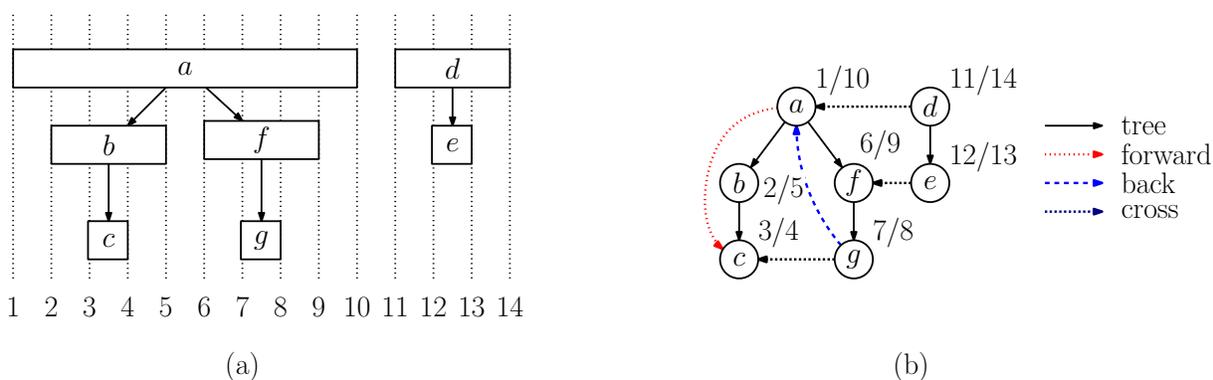


Fig. 7: (a) the Parenthesis Lemma and (b) the DFS edge types.

The structure of the remaining (non-tree) edges of the graph depend on the type of graph involved. For **undirected graphs**, the remaining edges are called *back edges*. An important observation is that for each back edge  $(u, v)$ ,  $u$  is either a proper ancestor or a proper descendant of  $v$ . To see why, consider any non-tree edge  $(u, v)$ . Since the graph is undirected, we may assume without loss of generality that  $u$  was discovered before  $v$ . By the parenthesis lemma, this means either that  $u$  is an ancestor of  $v$  (and we are done) or that their discovery-finish intervals are disjoint. If they are disjoint,  $u$  must finish before  $v$  is discovered. However, this is impossible, because as we are processing  $u$ , we will see the edge  $(u, v)$  and thus discover  $v$ .

For **directed graphs** the non-tree edges of the graph can be classified as follows (See Fig. 7(b)):

**Back edges:**  $(u, v)$  where  $v$  is a (not necessarily proper) ancestor of  $u$  in the tree. (Thus, a self-loop edge is considered to be a back edge.)

**Forward edges:**  $(u, v)$  where  $v$  is a proper descendant of  $u$  in the tree.

**Cross edges:**  $(u, v)$  where  $u$  and  $v$  are not ancestors or descendants of one another (in fact, the edge may go between different trees of the forest).

It is not difficult to classify the edges of a DFS tree on-the-fly by analyzing the vertex status (undiscovered, discovered, finished) and/or considering the time stamps. (This is left as an exercise.)<sup>2</sup>

<sup>2</sup>Be careful, however. Remember that in an undirected graph, every edge is represented twice. When classifying back edges, you should be sure that you are not seeing the other half of a tree edge.

## Lecture 3: Cycles and Strong Components

**Applications of DFS:** Last time we introduced depth-first search (DFS). Today, we discuss some applications of this powerful and efficient method of graph traversal.

**Directed Acyclic Graphs (Optional):** A *directed acyclic graph*, or *DAG*, is a directed graph that has no cycles. DAGs arise in many applications where there are precedence or ordering constraints. For example, if there are a series of tasks to be performed, and certain tasks must precede other tasks (e.g., in construction you have to build the walls before you install the windows). In general a *precedence constraint graph* is a DAG in which vertices are tasks and the edge  $(u, v)$  means that task  $u$  must be completed before task  $v$  begins.

It is easy to see that every DAG must have at least one vertex with no incoming edges, and at least one vertex with no outgoing edges. A vertex with no incoming edges (only outgoing) is called a *source* and a vertex with no outgoing edges (only incoming) is called a *sink*.

**Acyclicity Testing (Optional):** Let us consider the problem of determining whether a digraph is acyclic. We are given a directed graph  $G = (V, E)$ , and we wish to determine whether  $G$  contains a cycle. If so,  $G$  is not a DAG.

We will present a simple algorithm based on DFS. Recall that in addition to tree edges, a DFS forest contains three other types of edges, back edges (which go to a vertex's ancestor), forward edges (which go to a vertex's descendant), and cross edges (everything else). Observe that if the DFS forest of  $G$  contains at least one back edge, then  $G$  has a cycle. This is easy to see. If  $(u, v)$  is a back edge, then there is a path in the tree from the ancestor  $v$  to the descendant  $u$ , and the back edge from  $u$  to  $v$  completes the cycle. The following lemma shows that this condition is not only sufficient, but necessary.

**Claim:** If a digraph  $G$  has a cycle, then *any* DFS forest of  $G$  (i.e., no matter what order the vertices are visited) has a *back edge*.

**Proof:** The proof is based on a very simple observation about the various edge types and finish times. Recall from the Parenthesis Lemma (from the previous lecture) that if  $u$  is an ancestor of  $v$  then we have  $[d[v], f[v]] \subset [d[u], f[u]]$ . It follows that if  $(u, v)$  is a tree edge or forward edge then  $f[u] > f[v]$ . Also, observe that  $(u, v)$  is a cross edge, it must be  $u$  was discovered after  $v$  was finished (for otherwise,  $u$  would have made a DFSvisit call on  $v$ , implying that this would be a tree edge). Therefore  $d[u] > f[v]$ . Since a vertex cannot finish until after it was discovered, we have  $f[u] > f[v]$ .

In summary, for all these three edge types (tree, forward, and cross), the finish time of the origin is strictly greater than the finish time of the destination. It follows directly that it is impossible to complete a cycle from any combination of just these three edge types. Only for back edges do we have  $f[u] < f[v]$ , and therefore in order to form a cycle we need to have *at least one* back edge. Therefore, if a graph  $G$  has a cycle, in any DFS forest of  $G$  there must be at least one back edge.

The above theorem implies that in order to determine whether a graph  $G$  has a cycle, it suffices to test whether it has a back edge. How do we know whether an edge is a back edge. The proof of the above theorem provides an easy way. We can first apply DFS to  $G$ , and we then run through the edges, checking whether  $d[u] > d[v]$ . Can we do this on the fly as DFS is running? The answer is yes. Observe that a back edge goes from a vertex  $u$  to an ancestor  $v$ . Such an ancestor must have been discovered, but not yet finished. For the other non-tree edge types, the destination  $v$  will have already finished. The main DFS function is the same, only DFSvisit needs to be updated.

---

```

DFSvisit(u) {                                     // perform a DFS search at u
  mark[u] = discovered                           // u has been discovered
  d[u] = ++time
  for each (v in Adj(u)) {
    if (mark[v] == undiscovered) {              // undiscovered neighbor?
      pred[v] = u
      DFSvisit(v)                               // ...visit it
    }
    else if (mark[v] != finished) {            // found a back edge?
      output "G has a cycle!"
    }
  }
  mark[u] = finished                             // we're done with u
  f[u] = ++time
}

```

---

**Strong Components:** (The following material applies *only* to directed graphs!)

A digraph  $G = (V, E)$  is said to be *strongly connected* if for every vertex  $u$  and  $v$  there is a path from  $u$  to  $v$  and from  $v$  to  $u$ . It is easy that this *mutual reachability* relation between vertices is an equivalence relation. This implies that it partitions  $V$  into equivalence class, called the *strong components* (or *strongly-connected components*) of  $G$  (see Fig. 8(a) and (b)).

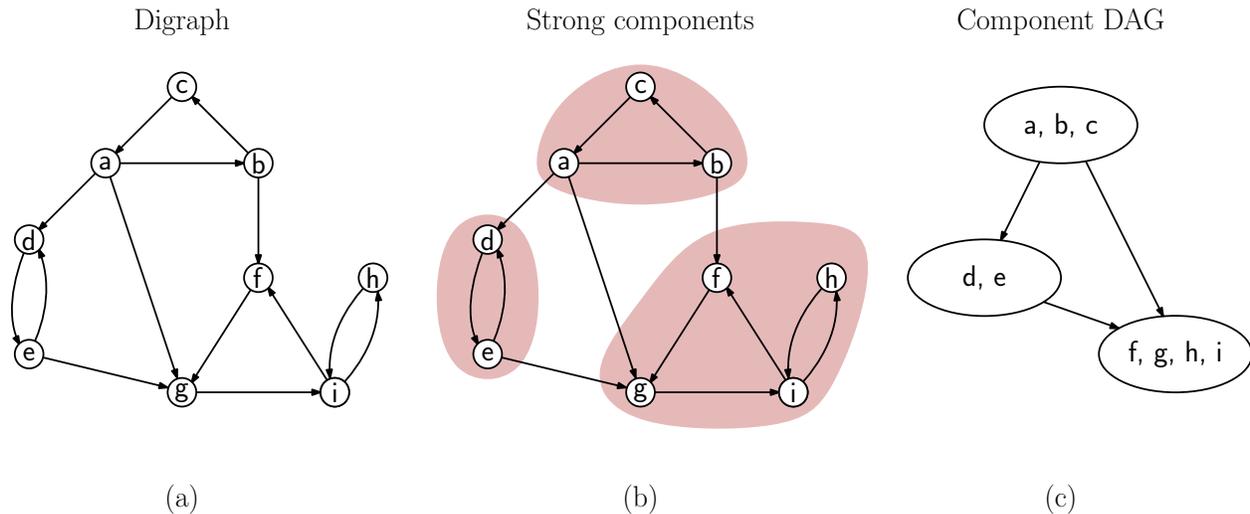


Fig. 8: Strong components and the component DAG.

If the vertices within each strong component are collapsed into a single vertex, the resulting digraph is called the *component digraph* (see Fig. 8(c)). It is easy to see that the component digraph must be acyclic (since if a number of components could be joined in a cycle, they would collapse into a single larger strong component). Therefore, this graph is usually called the *component DAG*.

There exists an  $O(n + m)$ -time DFS algorithm for computing strong components. It is based on the following lemma.

**Claim 1:** If DFSvisit is started at a vertex  $u$ , it will terminate precisely when all the vertices reachable from  $u$  have been visited.

**Proof:** This follows from the exhaustive nature of DFS. (Note that some of these vertices may have been reached by earlier calls to DFSvisit.)

**Claim 2:** If  $C$  and  $C'$  are two strong components, and there is an edge from a vertex in  $C$  to a vertex in  $C'$ , then the highest finish time in  $C$  is bigger than the highest finish time in  $C'$ .

**Proof:** There are two cases depending on whether the DFS first encounters a vertex from  $C$  or  $C'$ . If it first encounters a vertex  $u$  in  $C$ , then by Claim 1 the DFS will visit all the vertices of both  $C$  and  $C'$  before returning to  $u$ . Therefore,  $u$  will have the highest finish time of every vertex in  $C \cup C'$ . If it first visits a vertex in  $C'$ , then the DFS will get stuck in  $C'$  (since it is not possible to reach anything in  $C$ ). It follows that all the vertices of  $C$  will have higher discovery times than those of  $C'$ , which further implies that they will have higher finish times as well.

**Claim 3:** The vertex that receives the highest finish time in a DFS must lie in a source vertex of the component DAG. (Recall that a vertex in a DAG is called a source if it has no incoming edges.)

Claim 3 is equivalent to saying that the strong components can be linearly arranged in decreasing order of their highest finish times. By doing so, every edge in the component DAG will go from an earlier component in the linear order to a later one. How can we exploit this to obtain an efficient algorithm to find the strong components.

Claim 3 allows us to identify a vertex in some *source* of the component DAG. Unfortunately, this is not all that useful. What *would* be useful is to identify a vertex in a *sink* of the component DAG. If we could do this, we could start a DFS at this vertex, with the knowledge (by Claim 2) that no other strong components would be visited. We could then delete all these vertices (or equivalently, mark them as visited), and repeat the process. Eventually, all the strong components will be identified, each one arising as a separate subtree of the DFS forest.

So how do we convert an algorithm that identifies sources to one that identifies sinks in the component DAG? The trick is to reverse all the edges of  $G$ . Let  $G^R$  denote the directed graph that has the same vertex set as  $G$ , but every edge  $(u, v)$  is replaced by its reverse  $(v, u)$ . Note that the strong components of  $G^R$  are the same as  $G$ , but the direction of edges in the component DAG have all be *reversed*. Thus, the sources in the component DAG of  $G^R$  are sinks in the component DAG of  $G$ . This leads to the following (insanely clever) algorithm for computing strong components.

- (1) Given  $G$ , compute  $G^R$  (see Fig. 9(a)). (Note: This is a small programming exercise, which involves a simple traversal of  $G$ 's adjacency list. It can be done in  $O(n + m)$  time.)
- (2) Run DFS( $G^R$ ) and label each vertex of  $G$  with the finish time of the corresponding vertex of  $G^R$  (see Fig. 9(b)).
- (3) Sort the vertices of  $G$  in decreasing order of finish times. (Note: Since finish times are integers in the range  $[1, 2n]$ , this can be done in  $O(n)$  time through Bucket Sort.)
- (4) Run DFS( $G$ ), but in the outermost loop, whenever we need to find a new vertex to start DFSvisit, take the vertex with the highest finish time (using the above sorted order).
- (5) Each subtree of the DFS forest will be a strong component (see Fig. 9(c)).

The correctness of the above algorithm follows from the remarks made earlier. The running time is  $O(n + m)$ , dominated by the time to compute  $G^R$  and the times for the two DFS's.

## Lecture 4: Bridges and 2-Edge Connectivity

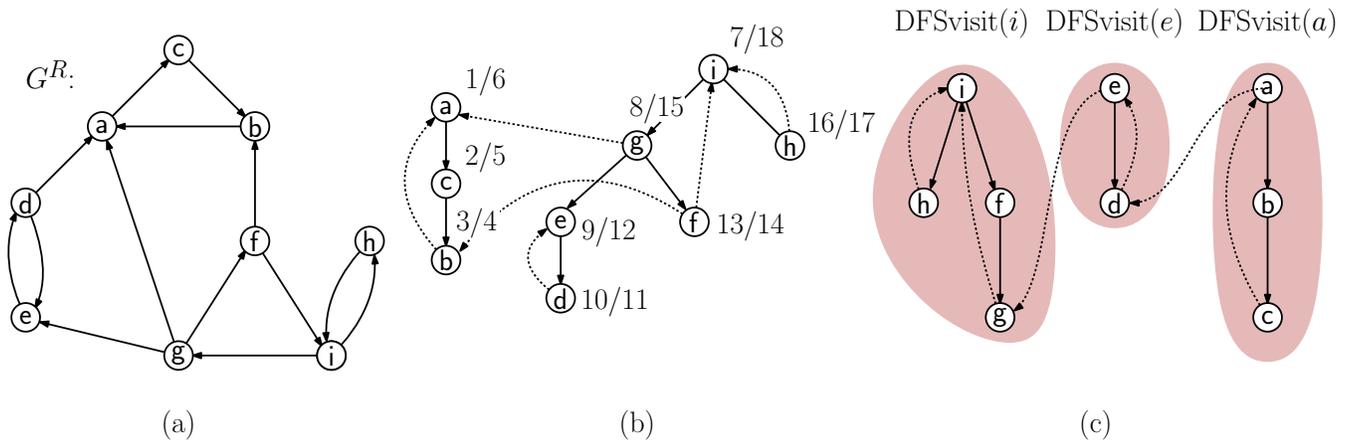


Fig. 9: Strong components and DFS.

**Higher-Order Graph Connectivity:** (The following material applies *only* to undirected graphs!)

Let  $G = (V, E)$  be an *connected* undirected graph. We often assume that our graphs are connected, but sometimes it is desirable to have a higher degree of connectivity. For example, if a graph can be disconnected through the removal of a single edge or vertex, the connectivity is rather “fragile.” Here are some definitions:

**Bridge:** Any edge whose removal results in a disconnected graph (see Fig. 10(a)).

**2-Edge Connected:** A graph is *2-edge connected* if it contains no bridges (see Fig. 10(b)). In general a graph is *k-edge connected* if the removal of any  $k - 1$  edges results in a connected graph.

Here are also vertex-based equivalents:

**Cut Vertex:** Any vertex whose removal (together with the removal of any incident edges) results in a disconnected graph (see Fig. 10(a)).

**Biconnected:** A graph is *biconnected* if it contains no cut vertices (see Fig. 10(c)). In general a graph is *k-connected* if the removal of any  $k - 1$  vertices results in a connected graph.

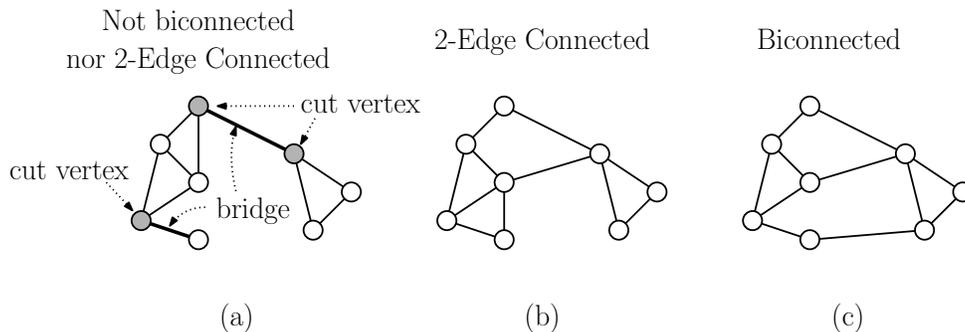


Fig. 10: Higher-order connectivity in graphs

We will present an  $O(n + m)$ -time algorithm for computing all the bridges of an undirected graph. (The algorithm can be modified to compute the cut vertices as well.)

Although we will not prove it, a useful fact about 2-edge connected graphs is that between every pair of distinct vertices, there are at least two edge-disjoint paths between them. (This is a consequence of a more general result called Menger's Theorem.) If a graph is biconnected, then for every pair of edges there is a simple cycle (that is a cycle that does not repeat any vertices) that contains both edges.

**Finding Bridges through DFS:** An obvious, but slow, way for computing bridges would be to delete each edge and then apply DFS to determine whether the resulting graph remains connected. However, this would take  $O(m(n + m))$  time. We will see that it is possible to identify all the bridges with a single application of DFS in  $O(n + m)$  time.

We assume that  $G$  is connected, which implies that there is a single DFS tree. Recall that the DFS tree consists of two types of edges: *tree edges*, which connect a parent with its child in the DFS tree, and *back edges*, which connect a (non-parent) ancestor with a (non-child) descendant.

Suppose that we are currently processing a vertex  $u$  in  $\text{DFSvisit}$ , and we see an edge  $(u, v)$  going to a neighbor  $v$  of  $u$ . If this edge is a back edge (that is, if  $v$  is an ancestor of  $u$ ) then  $(u, v)$  cannot be a bridge, because the tree edges between  $u$  and  $v$  provide a second way to connect these vertices. Therefore, we may limit consideration to when  $(u, v)$  is a tree edge, that is,  $v$  has not yet been discovered, and so we will invoke  $\text{DFSvisit}(v)$ . While we are doing this, we will keep track of the back edges in the subtree rooted at  $v$ . Observe that all these back edges remain entirely within this subtree (see Fig. 11(a)) then  $(u, v)$  is a bridge, since its removal completely disconnects this subtree from the rest of the tree. On the other hand, if there is even a single back edge leading out from this subtree, then  $(u, v)$  is *not* a bridge. Such a back edge must go from within the subtree to a proper ancestor of  $v$ . By the Parenthesis Lemma, this means that it leads to a vertex whose discovery time is strictly smaller than  $v$ 's discovery time. (Recall a vertex's ancestors are discovered before it.) In summary, we have established the following claim.

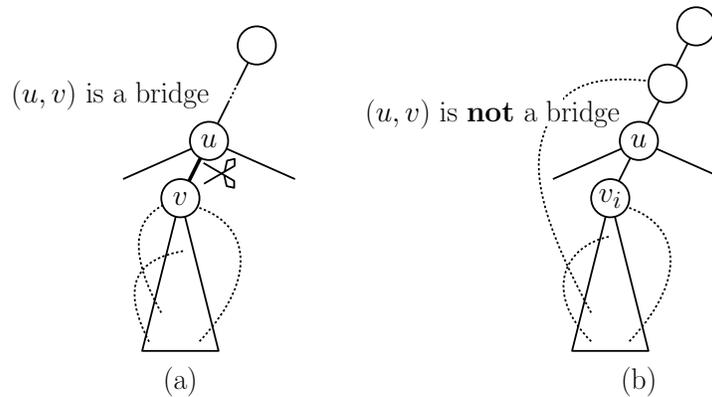


Fig. 11: Conditions for a vertex to be a cut vertex.

**Claim:** An edge  $(u, v)$  is a bridge if and only if it is a tree edge and (assuming that  $u$  is the parent of  $v$ ) there is no back edge within  $v$ 's subtree that leads to a vertex whose discovery time is strictly smaller than  $v$ 's discovery time.

**Tracking Back Edges:** The above claim provides us with a structural characterization of bridges. How can we design an algorithm that tests this condition? To do this, we will introduce an auxiliary quantity, which will be computed as the DFS runs. We define

$$\text{Low}[u] = \min(d[u], \min\{d[w] : \exists \text{ back edge } (v, w) \text{ where } v \text{ is a descendant of } u\}).$$

Note that we use the term "descendant" in the nonstrict sense, that is,  $v$  may be  $u$  itself.

Intuitively,  $\text{Low}[u]$  is the closest to the root that you can get in the tree by taking any one back edge from either  $u$  or any of its descendants. (Beware of this notation: “Low” means low discovery time, not “low” in our drawing of the DFS tree. In fact  $\text{Low}[u]$  tends to be “high” in the tree, in the sense of being close to the root.) Also note that you may consider *any* descendant of  $u$ , but you may only follow *one* back edge.

To compute  $\text{Low}[u]$  we use the following simple rules: Suppose that we are performing DFS on the vertex  $u$ .

**Initialization:**  $\text{Low}[u] = d[u]$ .

**Back edge** ( $u, v$ ):  $\text{Low}[u] = \min(\text{Low}[u], d[v])$ . Explanation: We have detected a new back edge coming out of  $u$ . If this goes to a lower  $d$ -value than the previous back edge then make this the new Low (see Fig. 12(a)).

**Tree edge** ( $u, v$ ):  $\text{Low}[u] = \min(\text{Low}[u], \text{Low}[v])$ . Explanation: Since  $v$  is in the subtree rooted at  $u$  any single back edge leaving the tree rooted at  $v$  is a single back edge for the tree rooted at  $u$  (see Fig. 12(b)).

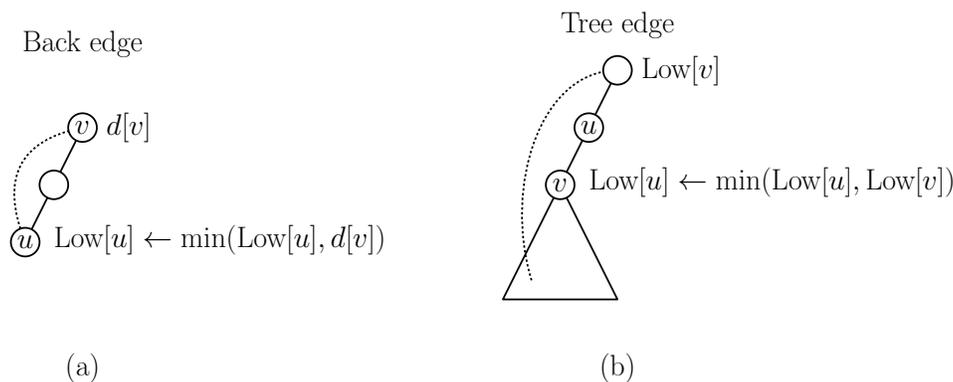


Fig. 12: Cut Vertices and the definition of  $\text{Low}[u]$ .

The code block below shows how to compute  $\text{Low}[u]$  for all vertices. We do not both computing finish times, since they are not needed for our purposes. Note that there is a subtlety in determining whether an edge is a back edge. Clearly, such an edge must go to a previously discovered vertex (which is why it is in the else-clause), but we also need to check that this vertex is not  $u$ 's parent. Recall that every edge of an undirected graph is reflected twice in the adjacency list (with  $v$  as a neighbor of  $u$  and  $u$  as a neighbor of  $v$ ). We need to check that we are not simply seeing the tree edge again, but from the child back to the parent. To do this, we check  $v$  is *not*  $u$ 's parent, that is  $v \neq \text{pred}[u]$ .

Observe that once  $\text{Low}[u]$  is computed for all vertices  $u$ , we can test whether a given tree edge  $(\text{pred}[v], v)$  is a bridge by testing whether there is no back edge in  $v$ 's subtree going to an earlier discovered vertex, that is,  $d[v] = \text{Low}[v]$ . (Actually, the test is more naturally stated as  $d[v] \geq \text{Low}[v]$ , but by definition  $\text{Low}[v] \leq d[v]$ , so testing for equality is equivalent.) The final code is shown in the code block below.

**Wrapup:** We have shown how to compute bridges in an undirected graph. There are a number of interesting problems that we still have not discussed. First, we claim that it is possible to adapt this algorithm to compute cut vertices as well. (The computation of Low is the same, but a different condition is applied to determine which vertices are cut vertices.) Second, if a graph fails to be 2-edge connected, it may be desirable to partition the vertices of the graph into 2-edge connected components. For example, in Fig. 13(c) the components consist of  $\{d, g, h\}$ ,  $\{b, c, f\}$ ,  $\{a\}$ , and  $\{e, i, j\}$ . This can be done by simple

```

DFSvisit(u) {
  mark[u] = discovered
  Low[u] = d[u] = ++time           // set discovery time and init Low
  for each (v in Adj(u)) {
    if (mark[v] == undiscovered) { // (u,v) is a tree edge
      pred[v] = u                 // v's parent is u
      DFSvisit(v)
      Low[u] = min(Low[u], Low[v]) // update Low[u]
    }
    else if (v != pred[u]) {      // (u,v) is a back edge
      Low[u] = min(Low[u], d[v]) // update Low[u]
    }
  }
}
}

```

```

findAllBridges(G) {
  time = 0
  for each (u in V) // initialize
    mark[u] = undiscovered

  for each (u in V) // undiscovered vertex?
    if (mark[u] == undiscovered) // ...start a new search here
      DFSvisit(u) // check for the bridges
  for each (v in V) {
    u = pred[v]
    if (u != null and d[v] == Low[v])
      output (u, v) as an bridge
  }
}

```

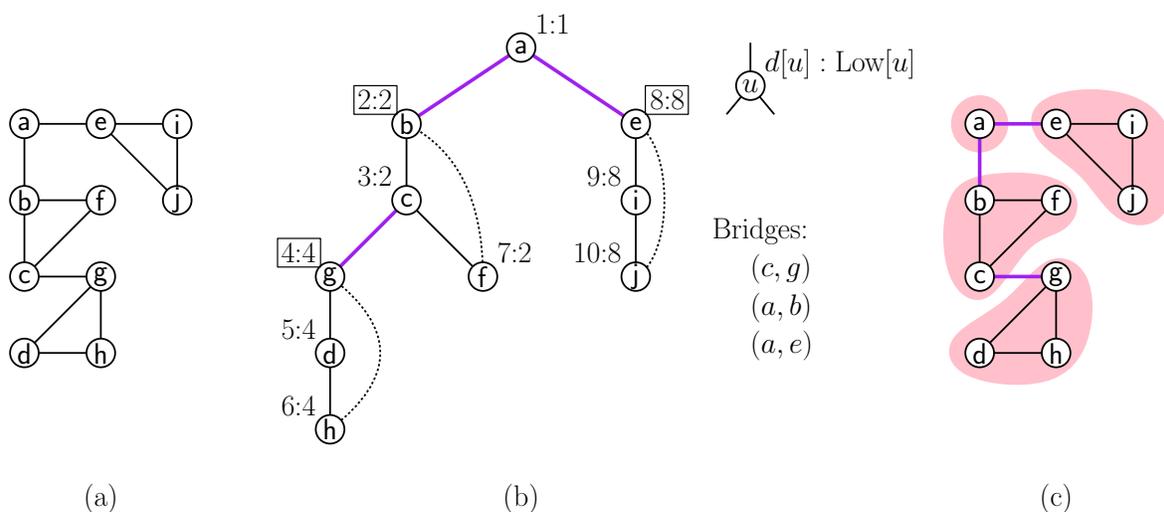


Fig. 13: Computing bridges via DFS.

extension as well. (The vertices are stored on a stack, and whenever a bridge is detected, we pop off an appropriate subset of the stack. We will leave the details as an exercise.) A similar approach can be applied to computing the biconnected components of a graph, which is a partition of the edge set of the graph.

## Lecture 5: Graph Shortest Paths: Dijkstra and Bellman-Ford

**Shortest Paths:** Today we consider the problem of computing shortest paths in a directed graph. We are given a digraph  $G = (V, E)$  and a source vertex  $s \in V$ , and we want to compute the shortest path from  $s$  to every other vertex in  $G$ . This is called the *single source shortest path problem*. The algorithms we will present work for undirected graphs as well, by simply assuming that each undirected edge consists of two directed edges going in opposite directions.

In general we assume that the graph is *weighted*, meaning that each edge  $(u, v) \in E$  has a numeric edge weight  $w(u, v)$ . The *cost* of a path is the sum of edge weights along the path. Define the *distance* from any vertex  $u$  to any vertex  $v$  to be the minimum cost over all the paths from  $u$  to  $v$ . We denote this by  $\delta(u, v)$ . Thus, we are interested in computing  $\delta(s, v)$  for all  $v \in V$ . We assume that every vertex has a trivial path of cost zero to itself, and hence  $\delta(v, v) = 0$ , for all  $v \in V$ .

**Background:** In earlier classes you may have learned about *breadth-first search*. This algorithm runs in  $O(n+m)$  algorithm for finding shortest paths from a single source vertex to all other vertices, assuming that the graph has no edge weights, or equivalently, that the weights are equal to unity. Thus, distance is just the number of edges on a path. This algorithm uses a first-in, first-out queue to process the vertices, visiting  $s$  first, then all the vertices accessible by a single edge, then all the vertices accessible by two edges, and so on, until all the vertices reachable from  $s$  have been processed. Today we will consider two algorithms that work for weighted graphs.

Since edge weights usually correspond to distances or travel times, it is typically the case that edge weights are positive, or at least, they are nonnegative. However, there are applications where negative edge weights make sense. For example, if an edge denotes a financial trade (buying or selling commodities), the cost may either be positive (a loss) or negative (a gain). In this context, the shortest path corresponds to a sequence of transactions that minimizes loss (or equivalently, maximizes gain).

In general it is possible to define shortest paths even for graphs with negative edge weights, but it should be noted that shortest paths are not well defined if the graph has *negative-cost cycles*. The reason is that the cost could be made arbitrarily small by traversing the cycle over and over. We will assume that the graphs we will be working with have no negative-cost cycles. Since allowing negative edge weights makes the problem more difficult to solve, we will consider these two variants separately.

**Dijkstra's Algorithm:** We first present a simple greedy algorithm for the single-source problem, which assumes that the edge weights are nonnegative. The algorithm, called *Dijkstra's algorithm*, was invented by the famous Dutch computer scientist, Edsger Dijkstra in 1956 (and published later in 1959). It is among the most famous algorithms in Computer Science.

In our presentation of the algorithm, we will stress the task of computing just the distance from the source to each vertex (not the path itself). We will store a *predecessor link* with each vertex, which points to way back to the source. Thus, the actual path can be found by traversing the predecessor links and reversing the resulting path. Since we store one predecessor link per vertex, the total space needed to store all the shortest paths is just  $O(n)$ .

**Shortest Paths and Relaxation:** The basic structure of Dijkstra's algorithm is to maintain an *estimate* of the shortest path for each vertex, call this  $d[v]$ . Intuitively  $d[v]$  stores the length of the shortest path from  $s$  to  $v$  that the algorithm currently knows of. Indeed, there will always exist a path of length  $d[v]$ , but it might not be the ultimate shortest path. Initially, we know of no paths, so  $d[v] = \infty$ , and

$d[s] = 0$ . As the algorithm proceeds and sees more and more vertices, it updates  $d[v]$  for each vertex in the graph, until all the  $d[v]$  values “converge” to the true shortest distances.

The process by which an estimate is updated is sometimes called *relaxation*. Here is how relaxation works. Intuitively, if you can see that your solution is not yet reached an optimum value, then push it a little closer to the optimum. In particular, if you discover a path from  $s$  to  $v$  shorter than  $d[v]$ , then you need to update  $d[v]$ . This notion is common to many optimization algorithms.

Consider an edge from a vertex  $u$  to  $v$  whose weight is  $w(u, v)$ . Suppose that we have already computed current estimates on  $d[u]$  and  $d[v]$ . We know that there is a path from  $s$  to  $u$  of weight  $d[u]$ . By taking this path and following it with the edge  $(u, v)$  we obtain a path to  $v$  of length  $d[u] + w(u, v)$ . If this path is better than the existing path of length  $d[v]$  to  $v$ , we should update  $d[v]$  to the value  $d[u] + w(u, v)$  (see Fig. 16.) We should also remember that the shortest path to  $v$  passes through  $u$ , which we do by setting  $\text{pred}[v]$  to  $u$  (see the code block below).

```

relax(u, v) {
  if (d[u] + w(u, v) < d[v]) { // is the path through u shorter?
    d[v] = d[u] + w(u, v) // yes, then take it
    pred[v] = u // record that we go through u
  }
}

```

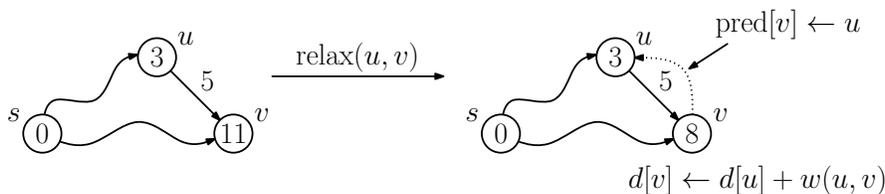


Fig. 14: Relaxation.

Observe that whenever we set  $d[v]$  to a finite value, there is always evidence of a path of that length. Therefore  $d[v] \geq \delta(s, v)$ . If  $d[v] = \delta(s, v)$ , then further relaxations cannot change its value.

It is not hard to see that if we perform  $\text{relax}(u, v)$  repeatedly over all edges of the graph, the  $d[v]$  values will eventually converge to the final true distance value from  $s$ . (This remark will reemerge when we discuss the Bellman-Ford algorithm later.) The cleverness of any shortest path algorithm is to perform the updates in a judicious manner, so the convergence is as fast as possible. In particular, the best possible would be to order relaxation operations in such a way that each edge is relaxed exactly once. Assuming that the edge weights are nonnegative, Dijkstra’s algorithm achieves this objective.<sup>3</sup>

**Dijkstra’s Algorithm:** Dijkstra’s algorithm operates by maintaining a subset of vertices,  $S \subseteq V$ , for which we claim we “know” the true distance, that is  $d[v] = \delta(s, v)$ . Initially  $S = \emptyset$ , the empty set, and we set  $d[s] = 0$  and all others to  $+\infty$ . One by one, we select vertices from  $V \setminus S$  to add to  $S$ . (If you haven’t seen it before, the notation “ $A \setminus B$ ” means the set  $A$  excluding the elements of set  $B$ . Thus  $V \setminus S$  consists of the vertices that are not in  $S$ .)

The set  $S$  can be implemented using an array of vertex marks. Initially all vertices are marked as “undiscovered,” and we set  $\text{mark}[v] = \text{finished}$  to indicate that  $v \in S$ .

<sup>3</sup>Note, by the way that while this objective is optimal in the worst case, there are instances where you might hope for much better performance. For example, given a cartographic road map of the entire United States, computing the shortest path between two locations near Washington DC should not require relaxing every edge of this road map. A better approach to this problem is provided by another greedy algorithm, called  $A^*$ -search.

How do we select which vertex among the vertices of  $V \setminus S$  to add next to  $S$ ? Here is where greedy selection comes in. Dijkstra recognized that the best (greediest) way is to process the vertex with the smallest  $d$ -value. That is, we take the unprocessed vertex that is closest (by our estimate) to  $s$ . This way, whenever a relaxation is being performed, it is possible to infer that result of the relaxation yields the final distance value. Later we will justify why this is the proper choice.

In order to perform this selection efficiently, we store the vertices of  $V \setminus S$  in a *priority queue* (e.g. a heap), where the key value of each vertex  $u$  is  $d[u]$ . We will need to make use of three basic operations that are provided by the priority queue:

**Build:** Create a priority queue from a list of  $n$  elements, each with an associated key value.

**Extract min:** Remove (and return a reference to) the element with the smallest key value.

**Decrease key:** Given a reference to an element in the priority queue, decrease its key value to a specified value, and reorganize if needed.

For example, using a standard binary heap (as in heapsort) the first operation can be done in  $O(n)$  time, and the other two can be done in  $O(\log n)$  time each. Dijkstra's algorithm is given in the code block below, and see Fig. 14 for an example.

---

Dijkstra's Algorithm

```
dijkstra(G,w,s) {
  for each (u in V) {                                // initialization
    d[u] = +infinity
    mark[u] = undiscovered
    pred[u] = null
  }
  d[s] = 0                                           // distance to source is 0
  Q = a priority queue of all vertices u sorted by d[u]
  while (Q is nonEmpty) {                            // until all vertices processed
    u = extract vertex with minimum d[u] from Q
    for each (v in Adj[u]) {                          // relax all outgoing edges from u
      if (d[u] + w(u,v) < d[v]) {
        d[v] = d[u] + w(u,v)
        decrease v's key in Q to d[v]
        pred[v] = u
      }
    }
    mark[u] = finished
  }
  [The pred pointers define an "inverted" shortest path tree]
}
```

---

Notice that the marking is not really used by the algorithm, but it has been included to make the connection with the correctness proof a little clearer.

To analyze Dijkstra's algorithm, recall that  $n = |V|$  and  $m = |E|$ . We account for the time spent on each vertex after it is extracted from the priority queue. It takes  $O(\log n)$  to extract this vertex from the queue. For each incident edge, we spend potentially  $O(\log n)$  time if we need to decrease the key of the neighboring vertex. Thus the time is  $O(\log n + \deg(u) \cdot \log n)$  time. The other steps of the update run in constant time. Recalling that the sum of degrees of the vertices in a graph is  $O(m)$ , the overall

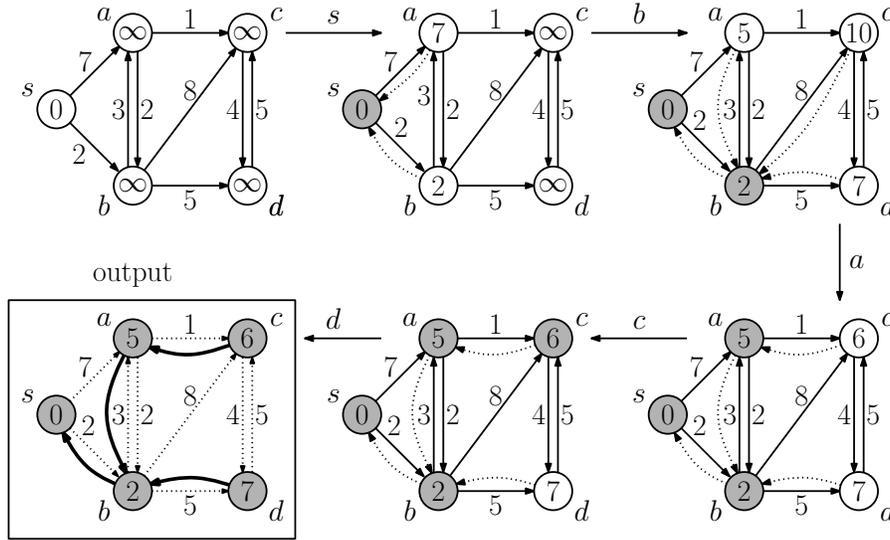


Fig. 15: Dijkstra's Algorithm example.

running time is given by  $T(n, m)$ , where

$$\begin{aligned}
 T(n, m) &= \sum_{u \in V} (\log n + \deg(u) \cdot \log n) = \sum_{u \in V} (1 + \deg(u)) \log n \\
 &= \log n \sum_{u \in V} (1 + \deg(u)) = (\log n)(n + 2m) = \Theta((n + m) \log n).
 \end{aligned}$$

Since  $G$  is connected,  $n$  is asymptotically no greater than  $m$ , so this is  $O(m \log n)$ . If you use a “smarter” heap (in particular, a Fibonacci heap, which supports very fast decrease-key operations, the running time is only  $O(n \log n + m)$ .

**Correctness:** Recall that  $d[v]$  is the distance value assigned to vertex  $v$  by Dijkstra's algorithm, and let  $\delta(s, v)$  denote the length of the true shortest path from  $s$  to  $v$ . To establish correctness, we need to show that  $d[v] = \delta(s, v)$  on termination. This is a consequence of the following lemma, which states that once a vertex  $u$  has been added to  $S$  (i.e., has been marked “finished”),  $d[u]$  is the true shortest distance from  $s$  to  $u$ .

**Lemma:** When a vertex  $u$  is added to  $S$ ,  $d[u] = \delta(s, u)$ .

**Proof:** Suppose to the contrary that at some point Dijkstra's algorithm *first* attempts to add a vertex  $u$  to  $S$  for which  $d[u] \neq \delta(s, u)$ . By our observations about relaxation,  $d[u]$  is never less than  $\delta(s, u)$ , thus we have  $d[u] > \delta(s, u)$ . Consider the situation just prior to the insertion of  $u$  into  $S$ , and consider the true shortest path from  $s$  to  $u$ . Because  $s \in S$  and  $u \in V \setminus S$ , at some point this path must first jump out of  $S$ . Let  $(x, y)$  be the first edge taken by the shortest path, where  $x \in S$  and  $y \in V \setminus S$  (see Fig. 16). (Note that it may be that  $x = s$  and/or  $y = u$ ).

Because  $u$  is the first vertex where we made a mistake and since  $x$  was already processed, we have  $d[x] = \delta(s, x)$ . Since we applied relaxation to  $x$  when it was processed, we must have

$$d[y] = d[x] + w(x, y) = \delta(s, x) + w(x, y) = \delta(s, y).$$

Since  $y$  appears before  $u$  along the shortest path and edge weights are nonnegative, we have  $\delta(s, y) \leq \delta(s, u)$ . Also, because  $u$  (not  $y$ ) was chosen next for processing, we know that  $d[u] \leq d[y]$ .

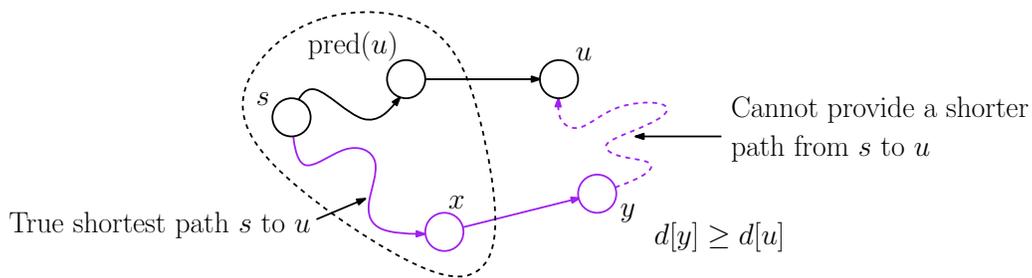


Fig. 16: Correctness of Dijkstra's Algorithm.

Putting this together, we have

$$\delta(s, u) < d[u] \leq d[y] = \delta(s, y) \leq \delta(s, u).$$

Clearly we cannot have  $\delta(s, u) < \delta(s, u)$ , which establishes the desired contradiction.

**Bellman-Ford Algorithm:** Let us now consider the question of how to solve the single-source shortest path problem when negative edge costs are allowed. Recall that we need to assume that the graph has no negative-cost cycles. (As an exercise, you are encouraged to think about how to detect whether this assumption is violated.)

We shall present the *Bellman-Ford algorithm*, which solves this problem. This algorithm slightly predates Dijkstra's algorithm, which was discovered in 1956. The algorithm was originally due to Alfonso Shimbel in 1955. It was rediscovered by Ford in 1956 and then by Bellman in 1958. This algorithm is slower than the best implementations of Dijkstra's algorithm in the worst case. Dijkstra's algorithm runs in time  $O(m \log n)$ , whereas Bellman-Ford runs in time  $O(nm)$ .

The idea behind the Bellman-Ford algorithm is quite simple. We initialize  $d[s] \leftarrow 0$  and  $d[v] \leftarrow \infty$  for all other vertices. We know that for each edge  $(u, v) \in E$ , the operation  $\text{relax}(u, v)$  propagates shortest-path information outwards from  $s$ . So, let's simply apply this operation repeatedly along each edge of  $E$  until the  $d$ -values converge.

Bellman-Ford Algorithm

```

bellman-ford(G, w, s) {
  for each (u in V) {                                // initialization
    d[u] = +infinity
    pred[u] = null
  }
  d[s] = 0
  repeat {                                          // repeat until convergence
    converged = true
    for each ((u, v) in E) {                        // relax along each edge
      if (d[u] + w(u, v) < d[v]) {
        d[v] = d[u] + w(u, v)
        pred[v] = u
        converged = false
      }
    }
  } until (converged)
  [The pred pointers define an "inverted" shortest path tree]
}

```

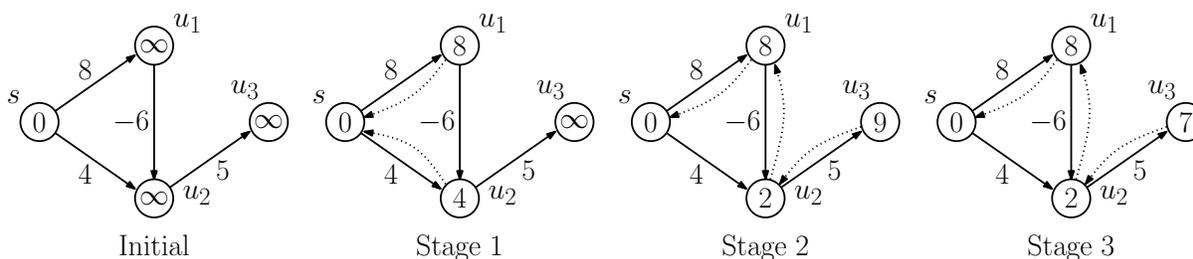


Fig. 17: Bellman-Ford Algorithm. (The algorithm will run for one more stage, but nothing will change, and it will terminate.)

**Correctness:** The following lemma establishes the correctness of the algorithm.

**Lemma:** On termination of the Bellman-Ford algorithm, for all  $v \in V$ ,  $d[v]$  contains the correct distance from  $s$  to  $v$ .

**Proof:** We assert first that *if* the algorithm converges, then  $d[v] = \delta(s, v)$  for all  $v \in V$ . As observed earlier,  $d[v]$  contains the cost of *some* path from  $s$  to  $v$ , so we have  $d[v] \geq \delta(s, v)$ . We will prove that  $d[v] \leq \delta(s, v)$  by induction on the length of the shortest path from  $s$  to  $v$ . In particular, if the shortest path from  $s$  to  $v$  consists of  $i$  edges, then after the  $i$ th iteration of the repeat-loop,  $d[v] = \delta(s, v)$ .

For the basis case ( $i = 0$ ) the only vertex whose shortest path is of length zero is  $s$  itself. By our initialization code,  $d[s] = 0$ , and by definition of shortest paths  $\delta(s, s) = 0$ , so we're done.

For the general case ( $i \geq 1$ ), consider any vertex  $v$  be any vertex such that the length of the shortest path from  $s$  to  $v$  consists of  $i$  edges. Let  $u$  be the vertex that immediately precedes  $v$  along this shortest path. It follows that (1) the shortest path from  $s$  to  $u$  is of length  $i - 1$ , and (2)  $\delta(s, v) = \delta(s, u) + w(u, v)$ . By the induction hypothesis, we know that after  $i - 1$  iterations of the repeat-loop, we have  $d[u] = \delta(s, u)$ . After the  $i$ th iteration of the repeat-loop, when we consider the edge  $(u, v)$  we will update the value of  $d[v]$  to

$$d[v] \leq d[u] + w(u, v) = \delta(s, u) + w(u, v) = \delta(s, v),$$

In conclusion, we have shown that  $\delta(s, v) \leq d[v] \leq \delta(s, v)$ , which implies that  $d[v] = \delta(s, v)$ . Which completes the proof.

**Running time:** Of course, all of the above is based on the assumption that the algorithm does indeed terminate. Next, we show that this will happen after at most  $n$  iterations of the repeat loop. We will make use of the following observation.

**Lemma:** If  $G$  has no negative cost cycles, then there is a shortest path from  $s$  to any vertex  $v$  that does not repeat any vertex.

**Proof:** Suppose to the contrary that the shortest path from  $s$  to  $v$  did repeat a vertex. Then there would be a cycle in the path. Since  $G$  has no negative cost cycles, we can remove this cycle from the path without increasing the path's total cost. If we repeat this for every repeated vertex, we will eventually have a path that contains no repetitions.

**Corollary:** For each  $v \in V$ , there is a shortest path from  $s$  to  $v$  consisting of at most  $n - 1$  edges (where  $n = |V|$ ).

It follows immediately from the corollary that after  $n - 1$  iterations, the  $d$ -values of all the vertices will have achieved their final values. After one more iteration, we will discover this fact, and the algorithm will terminate.

Since each relax operation required  $O(1)$  time, each iteration of the repeat-loop takes time  $O(m)$ . Since the algorithm terminates after  $n$  iterations, the total running time is  $O(nm)$ . Note that the algorithm may actually terminate earlier.

Now, looping back to the question we asked early, can we modify this algorithm to detect whether the graph has a negative cost cycle? Hopefully the easy answer will occur to you. (If not, think about it a bit more.)

## Lecture 6: Greedy Algorithms: Huffman Coding

**Greedy Algorithms:** In an *optimization problem*, we are given an input and asked to compute a structure, subject to various constraints, in a manner that either minimizes cost or maximizes profit. Such problems arise in many applications of science and engineering. Given an optimization problem, we are often faced with the question of whether the problem can be solved efficiently (as opposed to a brute-force enumeration of all possible solutions), and if so, what approach should be used to compute the optimal solution?

In many optimization algorithms a series of selections need to be made. A simple design technique for optimization problems is based on a *greedy* approach, that builds up a solution by selecting the best alternative in each step, until the entire solution is constructed. When applicable, this method can lead to very simple and efficient algorithms. (Unfortunately, it does not always lead to optimal solutions.)

Today, we will consider one of the most well-known examples of a greedy algorithm, the construction of Huffman codes.

**Huffman Codes:** Huffman codes provide a method of encoding data efficiently. Normally when characters are coded using standard codes like ASCII or the Unicode, each character is represented by a fixed-length *codeword* of bits (e.g., 8 or 16 bits per character). Fixed-length codes are popular, because its is very easy to break a string up into its individual characters, and to access individual characters and substrings by direct indexing. However, fixed-length codes may not be the most efficient from the perspective of minimizing the total quantity of data.

Consider the following example. Suppose that we want to encode strings over the (rather limited) 4-character alphabet  $C = \{a, b, c, d\}$ . We could use the following fixed-length code:

Character	a	b	c	d
Fixed-Length Codeword	00	01	10	11

A string such as “abacdaacac” would be encoded by replacing each of its characters by the corresponding binary codeword.

a	b	a	c	d	a	a	c	a	c
00	01	00	10	11	00	00	10	00	10

The final 20-character binary string would be “00010010110000100010”.

Now, suppose that you knew the relative probabilities of characters in advance. (This might happen by analyzing many strings over a long period of time. In applications like data compression, where you want to encode one file, you can just scan the file and determine the exact frequencies of all the characters.) You can use this knowledge to encode strings differently. Frequently occurring characters are encoded using fewer bits and less frequent characters are encoded using more bits. For example, suppose that characters are expected to occur with the following probabilities. We could design a *variable-length code* which would do a better job.

Character	a	b	c	d
Probability	0.60	0.05	0.30	0.05
Variable-Length Codeword	0	110	10	111

Notice that there is no requirement that the alphabetical order of character correspond to any sort of ordering applied to the codewords. Now, the same string would be encoded as follows.

a    b    a    c    d    a    a    c    a    c  
0   110   0   10   111   0   0   10   0   10

Thus, the resulting 17-character string would be “01100101110010010”. Thus, we have achieved a savings of 3 characters, by using this alternative code. More generally, what would be the expected savings for a string of length  $n$ ? For the 2-bit fixed-length code, the length of the encoded string is just  $2n$  bits. For the variable-length code, the expected length of a single encoded character is equal to the sum of code lengths times the respective probabilities of their occurrences. The expected encoded string length is just  $n$  times the expected encoded character length.

$$n(0.60 \cdot 1 + 0.05 \cdot 3 + 0.30 \cdot 2 + 0.05 \cdot 3) = n(0.60 + 0.15 + 0.60 + 0.15) = 1.5n.$$

Thus, this would represent a 25% savings in expected encoding length. (Of course, we would also need to consider the cost of transmitting the code book itself, but typically the code book is much smaller than the text being transmitted.) The question that we will consider today is how to form the *best code*, assuming that the probabilities of character occurrences are known.

**Prefix Codes:** One issue that we didn’t consider in the example above is whether we will be able to *decode* the string, once encoded. In fact, this code was chosen quite carefully. Suppose that instead of coding the character “a” as 0, we had encoded it as 1. Now, the encoded string “111” is ambiguous. It might be “d” and it might be “aaa”. How can we avoid this sort of ambiguity? You might suggest that we add separation markers between the encoded characters, but this will tend to lengthen the encoding, which is undesirable. Instead, we would like the code to have the property that it can be uniquely decoded.

Note that in both the variable-length codes given in the example above no codeword is a *prefix* of another. This turns out to be critical. Observe that if two codewords did share a common prefix, e.g.  $a \rightarrow 001$  and  $b \rightarrow 00101$ , then when we see  $00101\dots$  how do we know whether the first character of the encoded message is “a” or “b”. Conversely, if no codeword is a prefix of any other, then as soon as we see a codeword appearing as a prefix in the encoded text, then we know that we may decode this without fear of it matching some longer codeword. Thus we have the following definition.

**Prefix Code:** Mapping of codewords to characters so that no codeword is a prefix of another.

Observe that any binary prefix coding can be described by a binary tree in which the codewords are the leaves of the tree, and where a left branch means “0” and a right branch means “1”. The length of a codeword is just its depth in the tree. The code given earlier is a prefix code, and its corresponding tree is shown in Fig. 18.

Decoding a prefix code is simple. We just traverse the tree from root to leaf, letting the input character tell us which branch to take. On reaching a leaf, we output the corresponding character, and return to the root to continue the process.

**Expected encoding length:** Once we know the probabilities of the various characters, we can determine the total length of the encoded text. Let  $p(x)$  denote the probability of seeing character  $x$ , and let  $d_T(x)$

Character	a	b	c	d
Codeword	0	110	10	111

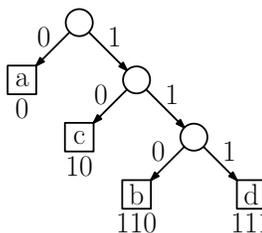


Fig. 18: A tree-representation of a prefix code.

denote the length of the codeword (depth in the tree) relative to some prefix tree  $T$ . The expected number of bits needed to encode a single character is given in the following formula:

$$B(T) = \sum_{x \in C} p(x)d_T(x).$$

This suggests the following problem:

**Optimal Code Generation:** Given an alphabet  $C$  and the probabilities  $p(x)$  of occurrence for each character  $x \in C$ , compute a prefix code  $T$  that minimizes the expected length of the encoded bit-string,  $B(T)$ .

There is an elegant greedy algorithm for finding such a code. It was invented in the 1950's by David Huffman, and is called a *Huffman code*. (While the algorithm is simple, it was not obvious. Huffman was a student at the time, and his professors, Robert Fano and Claude Shannon, two very eminent researchers, had developed their own algorithm, which as suboptimal.)

By the way, Huffman coding was used for many years by the Unix utility `pack` for file compression. Later it was discovered that there are better compression methods. For example, `gzip` is based on a more sophisticated method called the *Lempel-Ziv coding* (in the form of an algorithm called *LZ77*), and `bzip2` is based on combining the *Burrows-Wheeler transformation* (an extremely cool invention!) with run-length encoding, and Huffman coding.

**Huffman's Algorithm:** Here is the intuition behind the algorithm. Recall that we are given the occurrence probabilities for the characters. We are going to build the tree up from the leaf level. We will take two characters  $x$  and  $y$ , and “merge” them into a single *meta-character* called  $z$ , which then replaces  $x$  and  $y$  in the alphabet. The character  $z$  will have a probability equal to the sum of  $x$  and  $y$ 's probabilities. Then we continue recursively building the code on the new alphabet, which has one fewer character. When the process is completed, we know the code for  $z$ , say 010. Then, we append a 0 and 1 to this codeword, given 0100 for  $x$  and 0101 for  $y$ .

Another way to think of this, is that we merge  $x$  and  $y$  as the left and right children of a root node called  $z$ . Then the subtree for  $z$  replaces  $x$  and  $y$  in the list of characters. We repeat this process until only one meta-character remains. The resulting tree is the final prefix tree. Since  $x$  and  $y$  will appear at the bottom of the tree, it seem most logical to select the two characters with the smallest probabilities to perform the operation on. The result is Huffman's algorithm. It is illustrated in Fig. 19.

The pseudocode for Huffman's algorithm is given below. Let  $C$  denote the set of characters, and let  $n = |C|$ . Each character  $x \in C$  is associated with an occurrence probability  $\text{prob}[x]$ . Initially, the characters are all stored in a *priority queue*  $Q$ . Recall that this data structure can be built initially in  $O(n)$  time, and we can extract the element with the smallest key in  $O(\log n)$  time and insert a new element in  $O(\log n)$  time. The objects in  $Q$  are sorted by probability. Note that with each execution of the for-loop, the number of items in the queue decreases by one. So, after  $n - 1$  iterations, there is exactly one element left in the queue, and this is the root of the final prefix code tree.

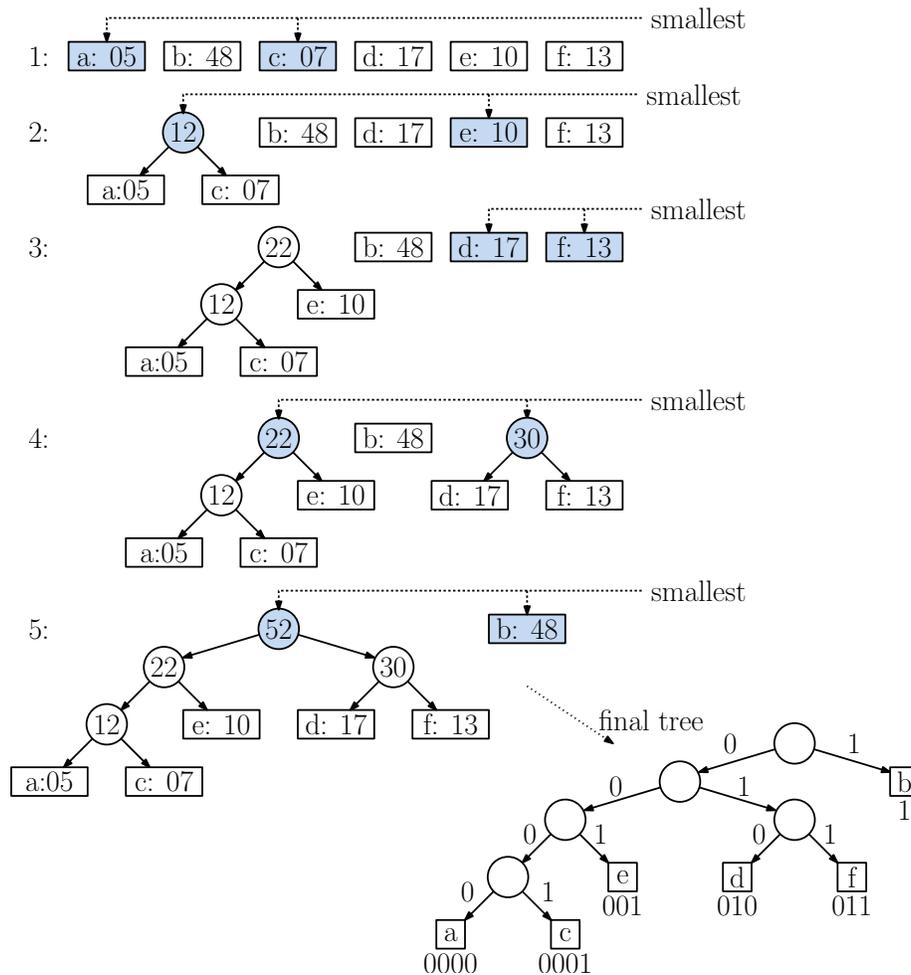


Fig. 19: Huffman's Algorithm.

---

Huffman's Algorithm

```

huffman(C, prob) {
    for each (x in C) {
        add x to Q sorted by prob[x]
    }
    for (i = 1 to |C| - 1) {
        z = new internal tree node
        left[z] = x = extract-min from Q // extract min probabilities
        right[z] = y = extract-min from Q
        prob[z] = prob[x] + prob[y] // z's probability is their sum
        insert z into Q // z replaces x and y
    }
    return the last element left in Q as the root
}
  
```

---

**Correctness:** The big question that remains is why is this algorithm correct, that is, does it compute the tree that minimizes the expected encoding length? Recall that the cost of any encoding tree  $T$  is  $B(T) = \sum_x p(x)d_T(x)$ . Our approach will be to show that any tree that differs from the one constructed by Huffman's algorithm can be converted into one that is equal to Huffman's tree without increasing its cost. This is done by identifying an appropriate place where the two solutions differ, modify the non-greedy solution so that it is a bit closer to the greedy solution, and showing that this modification can be done so that the cost does not increase. By repeating this, we will eventually modify any solution into the greedy solution in a manner that does not increase this cost. This implies that the greedy solution is has the minimum cost.

Our approach is based a few observations. First, observe that the Huffman tree is a *full binary tree*, meaning that every internal node has exactly two children. (It would never pay to have an internal node with only one child, since we could replace this node with its child without increasing the tree's cost.) So we may safely limit consideration to full binary trees. Our next observation (proved below) is that in any optimal code tree, the two characters with the lowest probabilities will be siblings at the maximum depth in the tree. Once we have this fact, we will merge these two characters into a single meta-character whose probability is the sum of their individual probabilities. As a result, we will now have one less character in our alphabet. This will allow us to apply induction to the remaining  $n - 1$  characters.

Let's first prove the above assertion that the two characters of lowest probability may be assumed to be siblings at the lowest level of the tree.

**Claim 1:** Consider the two characters,  $x$  and  $y$  with the smallest probabilities. Then there is an optimal code tree in which these two characters are siblings at the maximum depth in the tree.

**Proof:** Let  $T$  be any optimal prefix code tree, and let  $b$  and  $c$  be two siblings at the maximum depth of the tree. (There may be many such siblings, and if so pick any such pair.) If  $\{x, y\} = \{b, c\}$  we are done. Otherwise, from the fact that  $x$  and  $y$  have the lowest probabilities, we may label the nodes such that  $p(b) \leq p(c)$  and  $p(x) \leq p(y)$ .

Now, since  $x$  and  $y$  have the two smallest probabilities it follows that  $p(x) \leq p(b)$  and  $p(y) \leq p(c)$ . (In both cases they may be equal.) Because  $b$  and  $c$  are at the deepest level of the tree we know that  $d_T(b) \geq d_T(x)$  and  $d_T(c) \geq d_T(y)$ . (Again, they may be equal.) Thus, we have  $p(b) - p(x) \geq 0$  and  $d_T(b) - d_T(x) \geq 0$ , and hence their product is nonnegative. Now, suppose that we switch the positions of  $x$  and  $b$  in the tree, resulting in a new tree  $T'$  (see Fig. 20).

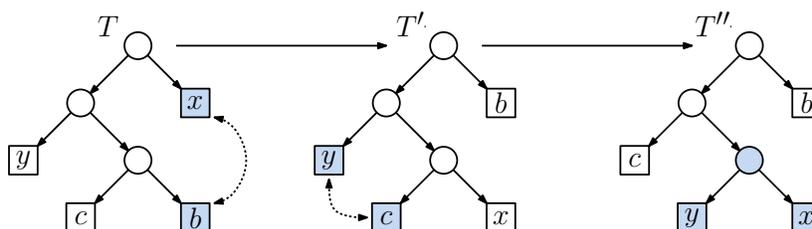


Fig. 20: Showing that the lowest probability nodes are siblings at the tree's lowest level.

Next let us see how the cost changes as we go from  $T$  to  $T'$ . Almost all the nodes contribute the same to the expected cost in both trees. The only exceptions are nodes  $x$  and  $b$ . By subtracting the old contributions of these nodes and adding in the new contributions we have

$$\begin{aligned} B(T') &= B(T) - (\text{old cost for } b \text{ and } x) + (\text{new cost for } b \text{ and } x) \\ &= B(T) - (p(x)d_T(x) + p(b)d_T(b)) + (p(x)d_T(b) + p(b)d_T(x)). \end{aligned}$$

With a little algebraic manipulation we obtain

$$\begin{aligned} B(T') &= B(T) + p(x)(d_T(b) - d_T(x)) - p(b)(d_T(b) - d_T(x)) \\ &= B(T) - (p(b) - p(x))(d_T(b) - d_T(x)) \\ &\leq B(T), \end{aligned}$$

where the last step follows because  $(p(b) - p(x))(d_T(b) - d_T(x)) \geq 0$ . Thus the cost does not increase. (Given our assumption that  $T$  was already optimal, it certainly cannot decrease either, since otherwise we would have a contradiction.) Since  $T$  was an optimal tree,  $T'$  is also an optimal tree.

By a similar argument, we can switch  $y$  with  $c$  to obtain a new tree  $T''$ . Again, the same sort of argument implies that  $T''$  is also optimal. The final tree  $T''$  satisfies the statement of the claim.

The above claim applies to just one pair of nodes, those with the lowest probabilities. To show that the *entire* Huffman tree is optimal, we need to extend this argument. We will do this by induction. In order to reduce from  $n$  characters to  $n - 1$ , we will do the same reduction that Huffman's algorithm does; namely we will *merge* characters  $x$  and  $y$  into a new meta-character  $z$ , whose probability is the sum of the probabilities of  $x$  and  $y$ .

**Claim 2:** Let  $T_n$  be any prefix-code tree that satisfies the property of Claim 1 (lowest probability symbols  $x$  and  $y$  are siblings at the deepest level). Let  $T_{n-1}$  be the tree that results by replacing these two nodes and their parent with a single leaf node  $z$  of probability  $p(z) = p(x) + p(y)$ . Then  $B(T_n) = B(T_{n-1}) + p(z)$ .

**Proof:** Let  $d$  denote the depths of  $x$  and  $y$  in  $T_n$ . Clearly,  $z$  is at depth  $d - 1$  in  $T_{n-1}$  (see Fig. 21).

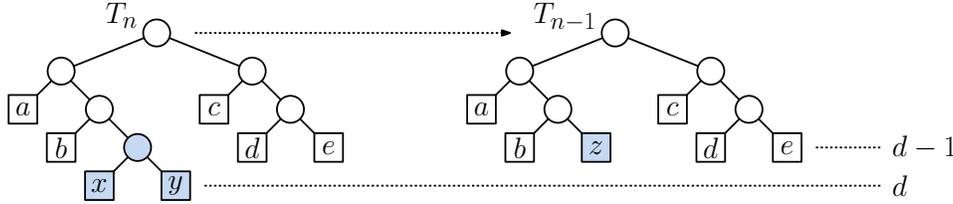


Fig. 21: Proving the correctness of Huffman's algorithm.

Because  $z$  replaces  $x$  and  $y$  the costs of the two trees satisfies

$$\begin{aligned} B(T_n) &= B(T_{n-1}) - (z\text{'s cost in } B(T_{n-1})) + (x \text{ and } y\text{'s costs in } B(T_n)) \\ &= B(T_{n-1}) - p(z)(d - 1) + (p(x)d + p(y)d) \\ &= B(T_{n-1}) - p(z)(d - 1) + p(z)d \\ &= B(T_{n-1}) + p(z). \end{aligned}$$

Note that the cost of trees  $T_n$  and  $T_{n-1}$  differ only by the fixed term  $p(z)$ , which does not depend on the tree's structure. Therefore (subject to this replacement), minimizing the cost of  $T_n$  is equivalent to minimizing the cost of  $T_{n-1}$ . This allows us to prove our main result.

**Claim 3:** Huffman's algorithm produces an optimal prefix code tree.

**Proof:** The proof is by induction on  $n$ , the number of characters. The basis case ( $n = 1$ ) is trivial, since there is only one tree possible. If  $n \geq 2$ , then by Claim 1, we know that the two characters  $x$  and  $y$  of lowest probability are siblings at the deepest level of an optimal tree. Huffman's algorithm replaces these nodes by a character  $z$  whose probability is the sum of their probabilities. By

induction, Huffman's algorithm computes the optimum tree over the resulting alphabet of  $n - 1$  symbols. Call it  $T_{n-1}$ . Replacing  $z$  with nodes  $x$  and  $y$  results in a tree  $T_n$  whose cost is higher by the fixed amount  $p(z) = p(x) + p(y)$ . Since  $T_{n-1}$  is optimal, and the cost of replacement does not depend on the tree's structure,  $T_n$  is also optimal.

## Lecture 7: Greedy Algorithms for Scheduling

**Interval Scheduling:** We continue our discussion of greedy algorithms with a number of problems motivated by applications in resource scheduling. Our first problem is called *interval scheduling*. We are given a set  $R$  of  $n$  activity requests that are to be scheduled to use some resource. Each activity has a given start time  $s_i$  and a given finish time  $f_i$ . For example, these may represent bids to use a picnic area in a neighborhood park. The Department of Parks and Recreation would like to grant as many of the bids as possible, but only one group can use the picnic area at a time.

We say that two requests  $i$  and  $j$  *conflict* if their start-finish intervals overlap, that is,

$$[s_i, f_i] \cap [s_j, f_j] \neq \emptyset.$$

(We do not allow finish time of one request to overlap the start time of another one, but this is easily remedied in practice.) Here is a formal problem definition.

**Interval scheduling problem:** Given a set  $R$  of  $n$  activity requests with start-finish times  $[s_i, f_i]$  for  $1 \leq i \leq n$ , determine a subset of  $R$  of maximum cardinality consisting of requests that are mutually non-conflicting.

An example of an input and two possible (optimal) solutions is given in Fig. 22. Notice that goal here is to maximize the *number* of activities that are granted (as opposed, say to some other criterion, like maximizing the total time that the resource is being used).

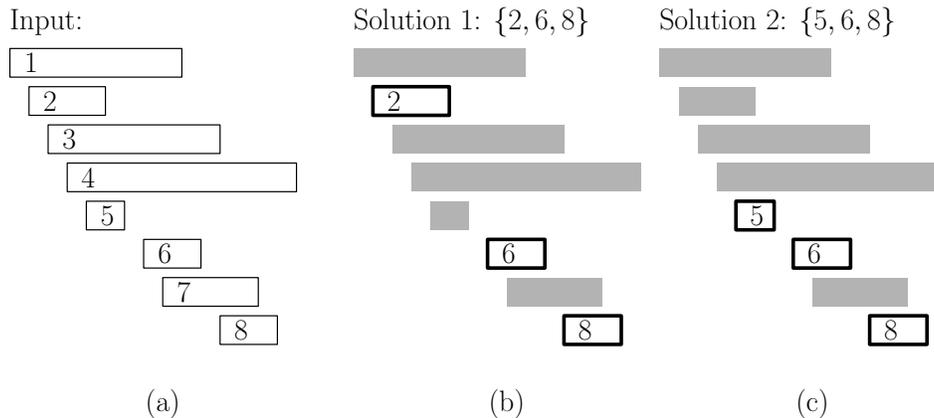


Fig. 22: An input and two possible solutions to the interval scheduling problem.

How do we schedule the largest number of activities on the resource? There are a number of ideas on how to proceed. As we shall see, there are a number of seemingly reasonable approaches that *do not* guarantee an optimal solution.

**Earliest Activity First:** Repeatedly select the activity with the earliest start time, provided that it does not overlap any of the previously scheduled activities.

**Shortest Activity First:** Repeatedly select the activity with the smallest duration ( $f_i - s_i$ ), provided that it does not conflict with any previously scheduled activities.

**Lowest Conflict Activity First:** Repeatedly select the activity that conflicts with the smallest number of remaining activities, provided that it does not conflict with of the previously scheduled activities. (Note that once an activity is selected, all the conflicting activities can be effectively deleted, and this affects the conflict counts for the remaining activities.)

As an exercise, show (by producing a counterexample) that each of the above strategies may not generate an optimal solution.

If at first you don't succeed, keep trying. Here, finally, is a greedy strategy that does work. Since we do not like activities that take a long time, let us select the activity that finishes first and schedule it. Then, we skip all activities that conflict with this one, and schedule the next one that has the earliest finish time, and so on. Call this strategy *Earliest Finish First* (EFF). The pseudo-code is presented in the code-block below. It returns the set  $S$  of scheduled activities.

---

Greedy Interval Scheduling

```
greedyIntervalSchedule(s, f) { // schedule tasks with given start/finish times
  sort tasks by increasing order of finish times
  S = empty // S holds the sequence of scheduled activities
  prev_finish = -infinity // finish time of previous task
  for (i = 1 to n) {
    if (s[i] > prev_finish) { // task i doesn't conflict with previous?
      append task i to S // ...add it to the schedule
      prev_finish = f[i] // ...and update the previous finish time
    }
  }
  return S
}
```

---

An example is given in Fig. 23. The start-finish intervals are given in increasing order of finish time. Activity 1 is scheduled first. It conflicts with activities 2 and 3. Then activity 4 is scheduled. It conflicts with activities 5 and 6. Finally, activity 7 is scheduled, and it interferes with the remaining activity. The final output is  $\{1, 4, 7\}$ . Note that this is not the only optimal schedule.  $\{2, 4, 7\}$  is also optimal.

The algorithm's correctness will be shown below. The running time is dominated by the  $O(n \log n)$  time needed to sort the jobs by their finish times. After sorting, the remaining steps can be performed in  $O(n)$  time.

**Correctness:** Let us consider the algorithm's correctness. First, observe that the output is a valid schedule in the sense that no two conflicting tasks appear in the final schedule. This is because we only add a task if its start time exceeds the previous finish time, and the previous finish time increases monotonically as the algorithm runs.

Second, we consider optimality. The proof's structure is worth noting, because it is common to many correctness proofs for greedy algorithms. It begins by considering an arbitrary solution, which may assume to be an optimal solution. If it is equal to the greedy solution, then the greedy solution is optimal. Otherwise, we consider the first instance where these two solutions differ. We replace the alternate choice with the greedy choice and show that things can only get better. Thus, by applying this argument inductively, it follows that the greedy solution is as good as an optimal solution, thus it is optimal.

**Claim:** The EFF strategy provides an optimal solution to interval scheduling.

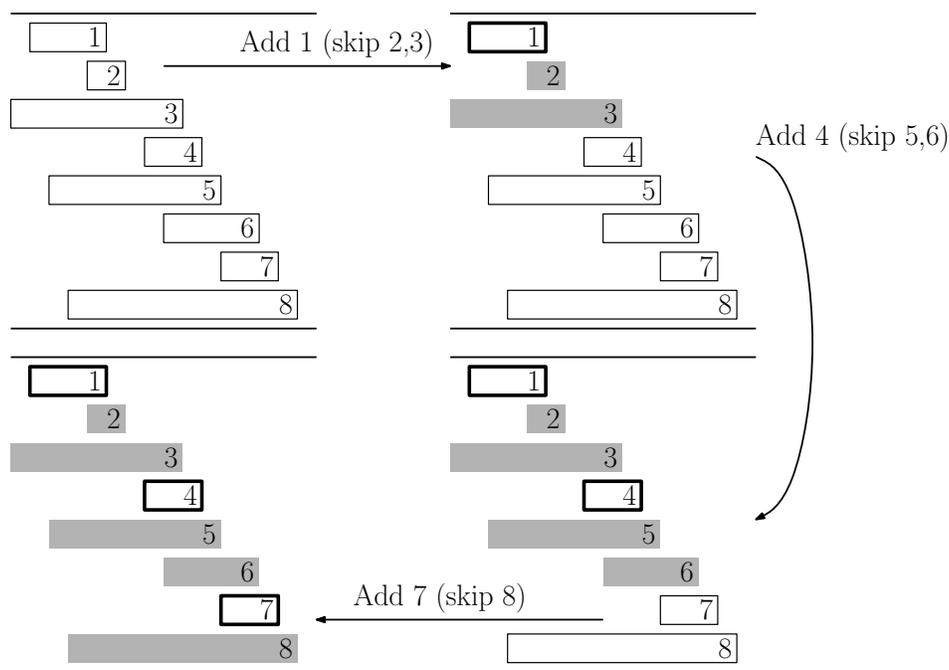


Fig. 23: An example of the greedy algorithm for interval scheduling. The final schedule is  $\{1, 4, 7\}$ .

**Proof:** Let  $O = \langle x_1, \dots, x_k \rangle$  be the activities of an *optimal solution* listed in increasing order of finish time, and let  $G = \langle g_1, \dots, g_{k'} \rangle$  be the activities of the EFF solution similarly sorted. If  $G = O$ , then we are done. Otherwise, observe that since  $O$  is optimal, it must contain at least as many activities as the greedy schedule, and hence there is a first index  $j$  where these two schedules differ. That is, we have:

$$\begin{aligned} O &= \langle x_1, \dots, x_{j-1}, x_j, \dots \rangle \\ G &= \langle x_1, \dots, x_{j-1}, g_j, \dots \rangle, \end{aligned}$$

where  $g_j \neq x_j$ . (Note that  $k \geq j$ , since otherwise  $G$  would have more activities than  $O$ , which would contradict  $O$ 's optimality.) The greedy algorithm selects the activity with the earliest finish time that does not conflict with any earlier activity. Thus, we know that  $g_j$  does not conflict with any earlier activity, and it finishes no later than  $x_j$  finishes.

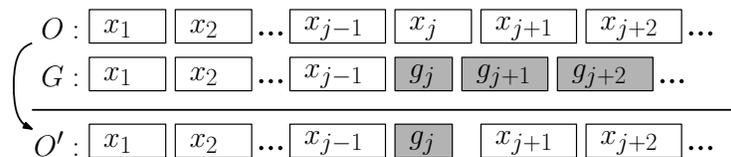


Fig. 24: Proof of optimality for the greedy schedule.

Consider the modified “greedier” schedule  $O'$  that results by replacing  $x_j$  with  $g_j$  in the schedule  $O$  (see Fig. 24). That is,  $O' = \langle x_1, \dots, x_{j-1}, g_j, x_{j+1}, \dots, x_k \rangle$ . Clearly,  $O'$  is a valid schedule, because  $g_j$  finishes no later than  $x_j$ , and therefore it cannot create any new conflicts. This new schedule has the same number of activities as  $O$ , and so it is at least as good with respect to our optimization criterion.

By repeating this process, we will eventually convert  $O$  into  $G$  without ever decreasing the number of activities. It follows that  $G$  is optimal.

**Interval Partitioning:** Next, let us consider a variant of the above problem. In interval scheduling, we assumed that there was a single exclusive resource, and our objective was to schedule as many nonconflicting activities as possible on this resource. Let us consider a different formulation, where instead we have an infinite number of possible exclusive resources to use, and we want to schedule *all* the activities using the smallest number resources. (The Department of Parks and Recreation can truck in as many picnic tables as it likes, but there is a cost, so it wants to keep the number small.)

As before, we are given a collection  $R$  of  $n$  activity requests, each with a start and finish time  $[s_i, f_i]$ . The objective is to find the smallest number  $d$ , such that it is possible to partition  $R$  into  $d$  disjoint subsets  $R_1, \dots, R_d$ , such that the events of  $R_j$  are mutually nonconflicting, for each  $j$ ,  $1 \leq j \leq d$ .

We can view this as a *coloring problem*. In particular, we want to assign colors to the activities such that two conflicting activities must have different colors. (In our example, the colors are rooms, and two lectures at the same time must be assigned to different class rooms.) Our objective is to find the minimum number  $d$ , such that it is possible to color each of the activities in this manner.

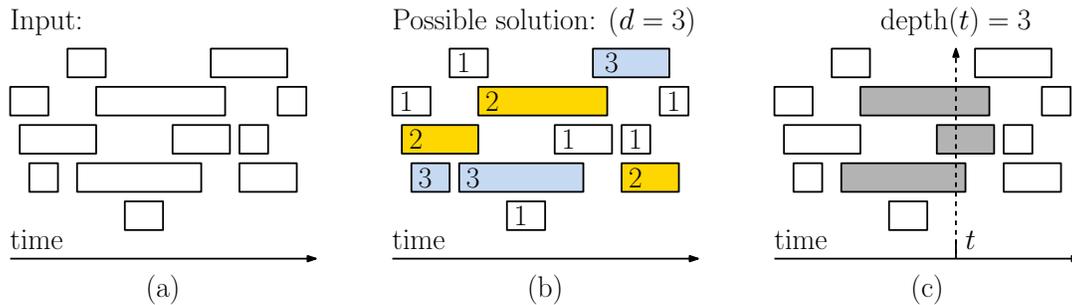


Fig. 25: Interval partitioning: (a) input, (b) possible solution, and (c) depth.

We refer to the subset of activities that share the same color as a *color class*. The activities of each color class are assigned to the same room. (For example, in Fig. 25(a) we give an example with  $n = 12$  activities and in (b) show an assignment involving  $d = 3$  colors. Thus, the six activities labeled 1 can be scheduled in one room, the three activities labeled 2 can be put in a second room, and the three activities labeled 3 can be put in a third room.)

In general, coloring problems are hard to solve efficiently (in the sense of being NP-hard). However, due to the simple nature of intervals, it is possible to solve the interval partitioning problem quite efficiently by a simple greedy approach. First, we sort the requests by increasing order of start times. We then assign each request the smallest color (possibly a new color) such that it conflicts with no other requests of this color class. The algorithm is presented in the following code block.

(The solution given in Fig. 25(b) comes about by running the above algorithm.) With its two nested loops, it is easy to see that the algorithm's running time is  $O(n^2)$ . If we relax the requirement that the color be the smallest available color (instead allowing any available color), it is possible to reduce this to  $O(n \log n)$  time with a bit of added cleverness.<sup>4</sup>

<sup>4</sup>Rather than have the for-loop iterate through just the start times, sort both the start times and the finish times into one large list of size  $2n$ . Each entry in this sorted list stores a record consisting of the type of event (start or finish), the index of the activity (a number  $1 \leq i \leq n$ ), and the time of the event (either  $s_i$  or  $f_i$ ). The algorithm visits each time instance from left to right, and while doing this, it maintains a stack containing the collection of *available colors*. It is not hard to show that each of the  $2n$  events can be processed in  $O(1)$  time. We leave the implementation details as an exercise. The total running time to sort the records is  $O((2n) \log(2n)) = O(n \log n)$ , and the total processing time is  $2n \cdot O(1) = O(n)$ . Thus, the overall running time is  $O(n \log n)$ .

---

```

greedyIntervalPartition(s, f) {    // schedule tasks with given start/finish times
    sort requests by increasing start times
    for (i = 1 to n) do {          // classify the ith request
        E = emptyset              // E stores excluded colors for activity i
        for (j = 1 to i-1) do {
            if ([s[j],f[j]] overlaps [s[i],f[i]]) add color[j] to E
        }
        Let c be the smallest color NOT in E
        color[i] = c
    }
    return color[1...n]
}

```

---

**Correctness:** Let us now establish the correctness of the greedy interval partitioning algorithm. We first observe that the algorithm never assigns the same color to two conflicting activities. This is due to the fact that the inner for-loop eliminates the colors of all preceding conflicting tasks from consideration. Thus, the algorithm produces a valid coloring. The question is whether it produces an optimal coloring, that is, one having the minimum number of distinct colors.

To establish this, we will introduce a helpful quantity. Let  $t$  be any time instant. Define  $\text{depth}(t)$  to be the number of activities whose start-finish interval contains  $t$  (see Fig. 25(c)). Given an set  $R$  of activity requests, define  $\text{depth}(R)$  to be the maximum depth over all possible values of  $t$ . Since the activities that contribute to  $\text{depth}(t)$  conflict with one another, clearly we need at least this many resources to schedule these activities. Therefore, we have the following:

**Claim:** Given any instance  $R$  of the interval partitioning problem, the number of resources needed is at least  $\text{depth}(R)$ .

This claim states that, if  $d$  denotes the minimum number of colors in any schedule, we have  $d \geq \text{depth}(R)$ . This does not imply, however, that this bound is necessarily achievable. But, in the case of interval partitioning, we can show that the depth bound is achievable, and indeed, the greedy algorithm achieves this bound.

**Claim:** Given any instance  $R$  of the interval partitioning problem, the number of resources produced by the greedy partitioning algorithm is at most  $\text{depth}(R)$ .

**Proof:** It will simplify the proof to assume that all start and finish times are distinct. (Let's assume that we have perturbed them infinitesimally to guarantee this.) We will prove a stronger result, namely that at any time  $t$ , the number of colors assigned to the activities that overlap time  $t$  is at most  $\text{depth}(t)$ . The result follows by taking the maximum over all times  $t$ .

To see why this is true, consider an arbitrary start time  $s_i$  during the execution of the algorithm. Let  $t^- = s_i - \varepsilon$  denote the time instant that is immediately prior to  $s_i$ . (That is, there are no events, start or finish, occurring between  $t^-$  and  $s_i$ .) Let  $d$  denote the depth at time  $t^-$ . By our hypothesis, just prior to time  $s_i$ , the number of colors being used is at most the current depth, which is  $d$ . Thus, when time  $s_i$  is considered, the depth increases to  $d + 1$ . Because at most  $d$  colors are in use prior to time  $s_i$ , there exists an unused color among the first  $d + 1$  colors. Therefore, the total number of colors used at time  $s_i$  is  $d + 1$ , which is not greater than the total depth.

To see whether you really understand the algorithm, ask yourself the following question. Is sorting of the activities essential to the algorithm's correctness? For example, can you answer the following questions?

- If the sorting step is eliminated, is the result necessarily optimal?
- If the tasks are sorted by some other criterion (e.g., finish time or duration) is the result necessarily optimal?

**Scheduling to Minimize Lateness:** Finally, let us discuss a problem of scheduling a set of  $n$  tasks where each task is associated with a *execution time*  $t_i$  and a *deadline*  $d_i$ . (Consider, for example, the assignments from your various classes and their due dates.) The objective is to schedule the tasks, no two overlapping in time, such that they are all completed before their deadline. If this is not possible, define the *lateness* of the  $i$ th task to be amount by which its finish time exceeds its deadline. The objective is to minimize the maximum lateness over all the tasks.

More formally, given the execution times  $t_i$  and deadlines  $d_i$ , the output is a set of  $n$  starting times,  $S = \{s_1, \dots, s_n\}$ , for the various tasks. Define the the finish time of the  $i$ th task to be  $f_i = s_i + t_i$  (its start time plus its execution time). The intervals  $[s_i, f_i]$  must be pairwise disjoint. The *lateness* of the  $i$ th task is the amount of time by which it exceeds its deadline, that is,  $\ell_i = \max(0, f_i - d_i)$ . The *maximum lateness* of  $S$  is defined to be

$$L(S) = \max_{1 \leq i \leq n} \max(0, f_i - d_i) = \max_{1 \leq i \leq n} \ell_i.$$

The overall objective is to compute  $S$  that minimizes  $L(S)$ .

An example is shown in Fig. 26. The input is given in Fig. 26(a), where the execution time is shown by the length of the rectangle and the deadline is indicated by an arrow pointing to a vertical line segment. A suboptimal solution is shown in Fig. 26(b), and the optimal solution is shown in Fig. 26(c). The width of each red shaded region indicates the amount by which the task exceeds its allowed deadline. The longest such region yields the maximum lateness.

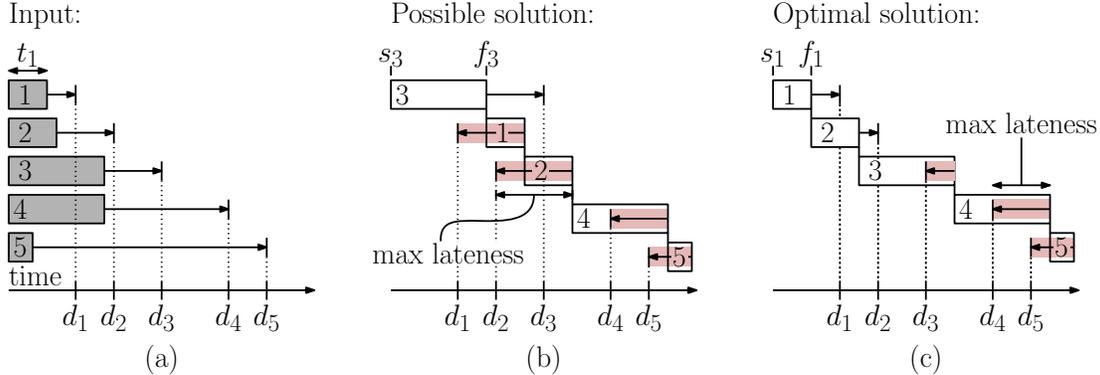


Fig. 26: Scheduling to minimize lateness.

Let us present a greedy algorithm for computing a schedule that minimizes maximum lateness. As before, we need to find a quantity upon which to base our greedy choices. Here are some ideas that *do not* guarantee an optimal solution.

**Smallest duration first:** Sort tasks by increasing order of execution times  $t_i$  and schedule them in this order.

**Smallest slack-time first:** Define the *slack time* of task  $x_i$  as  $d_i - t_i$ . This statistic indicates how long we can safely wait before starting a task. Schedule the tasks in increasing order of slack-time.

As before, see if you can generate a counterexample showing that each of the above strategies may fail to give the optimal solution.

So what is the right solution? The best strategy turns out to find the task that needs to finish first and get it out of the way. Define the *Earliest Deadline First* (EDF) strategy work by sorting the tasks by their deadline, and then schedule them in this order. (This is counterintuitive, because it completely ignores part of the input, namely the running times.) Nonetheless, we will show that this is the best possible. The pseudo-code is presented in the following code block.

---

```

greedySchedule(t, d) {
    // schedule given execution times and deadlines
    sort tasks by increasing deadline (d[1] <= ... <= d[n])
    f_prev = 0 // f is the finish time of previous task
    for (i = 1 to n) do {
        assign task i to start at s[i] = f_prev // start next task
        f_prev = f[i] = s[i] + t[i] // its finish time
        lateness[i] = max(0, f[i] - d[i]) // its lateness
    }
    return array s // return array of start times
}

```

---

The solution shown in Fig. 26(c) is the result of this algorithm. Observe that the algorithm's running time is  $O(n \log n)$ , which is dominated by the time to sort the tasks by their deadline. After this, the algorithm runs in  $O(n)$  time.

**Correctness:** It is easy to see that this algorithm produces a valid schedule, since we never start a new job until the previous job has been completed. We will show that this greedy algorithm produces an optimal schedule, that is, one that minimizes the maximum lateness. As with the interval scheduling problem, our approach will be to show that is it possible to “morph” any optimal schedule to look like our greedy schedule. In the morphing process, we will show that schedule remains valid, and the maximum lateness can never increase, it can only remain the same or decrease.

To begin, we observe that our algorithm has no *idle time* in the sense that the resource never sits idle during the running of the algorithm. It is easy to see that by moving tasks up to fill in any idle times, we can only reduce lateness. Henceforth, let us consider schedules that are idle-free. Let  $G$  be the schedule produced by the greedy algorithm, and let  $O$  be any optimal idle-free schedule. If  $G = O$ , then greedy is optimal, and we are done. Otherwise,  $O$  must contain at least one *inversion*, that is, at least one pair of tasks that have not been scheduled in increasing order of deadline. Let us consider the first instance of such an inversion. That is, let  $x_i$  and  $x_j$  be the first two consecutive tasks in the schedule  $O$  such that  $d_j < d_i$ . We have:

- (a) The schedules  $O$  and  $G$  are identical up to these two tasks
- (b)  $d_j < d_i$  (and therefore  $x_j$  is scheduled before  $x_i$  in schedule  $G$ )
- (c)  $x_i$  is scheduled before  $x_j$  in schedule  $O$

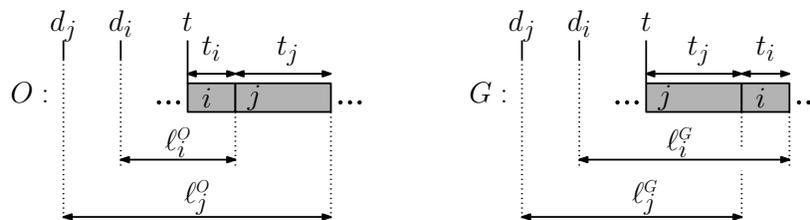


Fig. 27: Optimality of the greedy scheduling algorithm for minimizing lateness.

We will show that by swapping  $x_i$  and  $x_j$  in  $O$ , the maximum lateness cannot increase. The reason that swapping  $x_i$  and  $x_j$  in  $O$  does not increase lateness can be seen intuitively from Fig. 27. The lateness is reflected in the length of the horizontal arrowed line segments in the figure. It is evident that the worst lateness involves  $x_j$  in schedule  $O$  (labeled  $\ell_j^O$ ). Unfortunately, a picture is not a formal argument. So, let us see if we put this intuition on a solid foundation.

First, let us define some notation. The lateness of task  $i$  in schedule  $O$  will be denoted by  $\ell_i^O$  and the lateness of task  $j$  in  $O$  will be denoted by  $\ell_j^O$ . Similarly, let  $\ell_i^G$  and  $\ell_j^G$  denote the respective latenesses of tasks  $i$  and  $j$  in schedule  $G$ . Because the two schedules are identical up to these two tasks, and because there is no slack time in either, the first of the two tasks starts at the same time in both schedules. Let  $t$  denote this time (see Fig. 27). In schedule  $O$ , task  $i$  finishes at time  $t + t_i$  and (because it needs to wait for task  $i$  to finish) task  $j$  finishes as time  $t + (t_i + t_j)$ . The lateness of each of these tasks is the maximum of 0 and the difference between the finish time and the deadline. Therefore, we have

$$\ell_i^O = \max(0, t + t_i - d_i) \quad \text{and} \quad \ell_j^O = \max(0, t + (t_i + t_j) - d_j).$$

Applying a similar analysis to  $G$ , we can define the latenesses of tasks  $i$  and  $j$  in  $G$  as

$$\ell_i^G = \max(0, t + (t_i + t_j) - d_i) \quad \text{and} \quad \ell_j^G = \max(0, t + t_j - d_j).$$

The “max” will be a pain to carry around, so to simplify our formulas we will exclude reference to it. (You are encouraged to work through the proof with the full and proper definitions.)

Given the individual latenesses, we can define the maximum lateness contribution from these two tasks for each schedule as

$$L^O = \max(\ell_i^O, \ell_j^O) \quad \text{and} \quad L^G = \max(\ell_i^G, \ell_j^G).$$

Our objective is to show that by swapping these two tasks, we do not increase the overall lateness. Since this is the only change, it suffices to show that  $L^G \leq L^O$ . To prove this, first observe that,  $t_i$  and  $t_j$  are nonnegative and  $d_j < d_i$  (and therefore  $-d_j > -d_i$ ). Recalling that we are dropping the “max”, we have

$$\ell_j^O = t + (t_i + t_j) - d_j > t + t_i - d_i = \ell_i^O.$$

Therefore,  $L^O = \max(\ell_i^O, \ell_j^O) = \ell_j^O$ . Since  $L^G = \max(\ell_i^G, \ell_j^G)$ , in order to show that  $L^G \leq L^O$ , it suffices to show that  $\ell_i^G \leq L^O$  and  $\ell_j^G \leq L^O$ . By definition we have

$$\ell_i^G = t + (t_i + t_j) - d_i < t + (t_i + t_j) - d_j = \ell_j^O = L^O,$$

and

$$\ell_j^G = t + t_j - d_j \leq t + (t_i + t_j) - d_j = \ell_j^O = L^O.$$

Therefore, we have  $L^G = \max(\ell_i^G, \ell_j^G) \leq L^O$ , as desired. In conclusion, we have the following.

**Claim:** The greedy scheduling algorithm minimizes maximum lateness.

## Lecture 8: Greedy Approximation: The $k$ -Center Problem

**Greedy Approximation for NP-Hard Problems:** One of the common applications of greedy algorithms is for producing approximation solutions to NP-hard problems.

As we shall see later this semester, NP-hard optimization problems represent very challenging computational problems in the sense that there is no known exact solution that has worst-case polynomial-time running time. Given an NP-hard problem, there are no ideal algorithmic solutions. One has to compromise between optimality or running time. Nonetheless, there are a number of examples of NP-hard problems where simple greedy heuristics produce solutions that are not far from optimal.

**Clustering and Center-Based Clustering:** Clustering is a widely studied problem with applications in statistics, pattern recognition, and machine learning. In a nutshell, it involves partitioning a large set of objects, called *points*, into a small number of subsets whose members are all mutually close to one another. In machine learning, clustering is often performed in high-dimensional spaces. Each data object is associated with a *property vector*, a numeric vector whose length may range from tens to thousands that describes the salient properties of this object. Then clustering is performed to group similar objects together.

In this geometric view of clustering, we think of the data set as a set  $P$  of  $n$  points in a multi-dimensional space (see Fig. 28(a)). The output of the clustering algorithm is the set  $C = \{c_1, \dots, c_k\}$  of points called *cluster centers* or simply *centers* (see Fig. 28(b)). Depending on the method being employed, it may be required that the center points are drawn from  $P$  itself, or they may just be arbitrary points in space, sometimes called *Steiner points*. In our example, the cluster centers have been chosen from  $P$ .

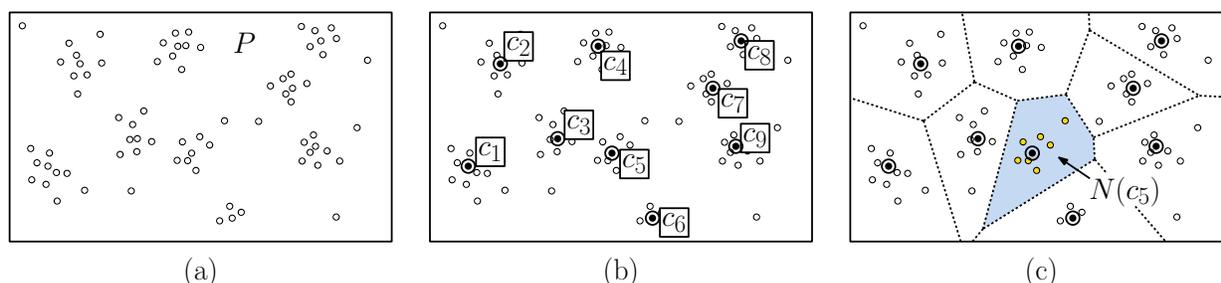


Fig. 28: Center-based geometry clustering.

Given a center point  $c_i$ , the associated cluster are the points of  $P$  that are closer to  $c_i$  center than to any other cluster center. We refer to this as  $c_i$ 's *neighborhood*, denoted  $N(c_i)$ . In Fig. 28(c), we show a subdivision of space into nine regions according to which cluster center is closest, and highlight the neighborhood  $N(c_5)$  in particular. (By the way, this subdivision is a famous structure in geometry, called the *Voronoi diagram* of the cluster centers. Computing Voronoi diagrams is covered in a course on Computational Geometry, but it is not necessary for our purposes. In  $O(nk)$  time we can compute the distance between each of the  $n$  points and each of the  $k$  centers, and assign each point to its closest center.)

**The  $k$ -Center Problem:** There are many ways in which to define clustering as an optimization problem. We will consider one of the simplest. We are given a set  $P$  of  $n$  points in space. For each pair of points  $u, v \in P$ , let  $\delta(u, v)$  denote the distance between  $u$  and  $v$ . We will assume that the distance satisfies the typical properties of any natural distance function:

**Positivity:**  $\delta(u, v) \geq 0$  and  $\delta(u, v) = 0$  if and only if  $u = v$ .

**Symmetry:**  $\delta(u, v) = \delta(v, u)$

**Triangle inequality:**  $\delta(u, w) \leq \delta(u, v) + \delta(v, w)$

For the rest of the lecture, we can think of  $\delta$  as the Euclidean distance, but the algorithm that we will present can be applied to any function satisfying the above properties.

**$k$ -center problem:** Given a set  $P$  of  $n$  points in space and an integer  $k \leq n$ , find a set  $C \subseteq P$  of  $k$  points in order to minimize the maximum distance of any point of  $P$  to its closest center in  $C$ .

Why maximum distance and not, say, sum of distances? There are many formulations of clustering, and this is the one that we have chosen for now.

We can view the  $k$ -center problem as a covering problem by balls. Given a point  $x$  in space and radius  $r$ , define the *ball*  $B(x, r)$  to be the (closed) ball of radius  $r$  centered at  $x$ . Given any solution  $C$  to the  $k$ -center problem, let  $\Delta(C)$  denote the maximum distance from any point of  $P$  to its closest center. If we now place balls of radius  $\Delta(C)$  about each point in  $C$ , it is easy to see that every point of  $P$  lies within the union of these balls. By definition of  $\Delta(C)$ , one of the points of  $P$  will lie on the boundary of one of these balls (for otherwise we could make  $\Delta(C)$  smaller). The neighborhood of each cluster will lie within its associated ball (see Fig. 29(b)).

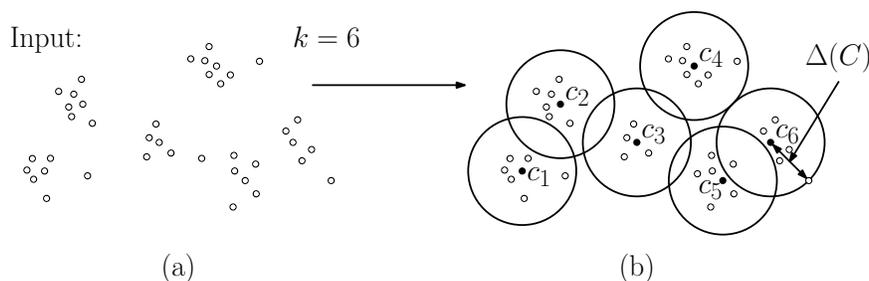


Fig. 29: The  $k$ -center problem in the Euclidean plane.

Given this perspective, we can see that the  $k$ -center problem is equivalent to the following problem:

**$k$ -center problem (equivalent form):** Given a set  $P$  of  $n$  points in space and an integer  $k \leq n$ , find the minimum radius  $\Delta$  and a set of balls of radius  $\Delta$  centered at  $k$  points of  $P$  such that  $P$  lies within the union of these balls.

**Greedy Approximation Algorithm:** Like many clustering problems, the  $k$ -center problem is known to be NP-hard, and so we will not be able to solve it exactly. (We will show this later this semester for a graph-based variant of the  $k$ -center problem.) Today, we will present a simple greedy algorithm that does not produce the optimum value of  $\Delta$ , but only an approximation to its value. Our algorithm will produce an estimate that is at most twice as large as the optimum value of  $\Delta$ .

Let us start with a couple of useful definitions. Given a set  $C = \{c_1, \dots, c_k\}$  of centers and for any  $c_i \in C$ , define the *bottleneck distance* of  $c_i$  to be the distance to its farthest point in  $N(c_i)$ , that is,

$$\Delta(c_i) = \max_{v \in N(c_i)} \delta(v, c_i).$$

Clearly, the maximum *bottleneck distance* over all the centers is just  $\Delta(C)$ .

Intuitively, if we think of the cluster centers as the locations of Starbucks (or your favorite retailer), then each customer (point in  $P$ ) is associated with the closest Starbucks. The set  $N(c_i)$  are the customers that go to the  $i$ th Starbucks location, and  $\Delta(c_i)$  is the maximum distance any of these customers needs to travel to get to  $c_i$ .  $\Delta(C)$  is the maximum distance that *anyone* needs to travel to their nearest Starbucks.

The greedy algorithm begins by selecting any point of  $P$  to be the initial center  $g_1$ . We then repeat the following process until we have  $k$  centers. Let  $G_i = \{g_1, \dots, g_i\}$  denote the current set of centers. Recall that  $\Delta(G_i)$  is the maximum distance of any point of  $P$  from its nearest center. Let  $u$  be the point achieving this distance. Intuitively,  $u$  is the *most dissatisfied customer*, since he/she has to drive the farthest to get to the nearest Starbucks. The greediest way to satisfy  $u$  is to put the next center directly at  $u$ . (Thus plopping the next Starbucks right on top of  $u$ 's house. Are you satisfied now?) In other words, set

$$g_{i+1} \leftarrow u \quad \text{and} \quad G_{i+1} \leftarrow G_i \cup \{g_{i+1}\}.$$

The pseudocode is presented in the code block below. The value  $d[u]$  denotes the distance from  $u$  to its closest center. (We make one simplification, by starting with  $G$  being empty. When we select the first center, all the points of  $P$  have infinite distances, so the initial choice is arbitrary.)

Greedy Approximation for  $k$ -center

```

GreedyKCenter(P, k) {
  G = empty
  for each u in P do                // initialize distances
    d[u] = INFINITY

  for (i = 1 to k) {
    Let u be the point of P such that d[u] is maximum
    Add u to G                       // u is the next cluster center
    for (each v in P) {              // update distance to nearest center
      d[v] = min(d[v], distance(v,u))
    }
    Delta = max_{v in P} d[v]        // update the bottleneck distance
  }
  return (G, Delta)                 // final centers and max distance
}

```

It is easy to see that the algorithm's running time is  $O(kn)$ . In general  $k \leq n$ , so this is  $O(n^2)$  in the worst case. One step of the algorithm is illustrated in Fig. 30. Assuming that we have three centers  $G = \{g_1, g_2, g_3\}$ , let  $g_4$  be the point that is farthest from its nearest center ( $g_1$  in this case). In each step we create a center at  $g_4$ , so now  $G = \{g_1, \dots, g_4\}$ . In anticipation of the next step, we find the point that maximizes the distance to its closest center ( $g_5$  in this case), and if the algorithm continues, it will be the location of the next center.

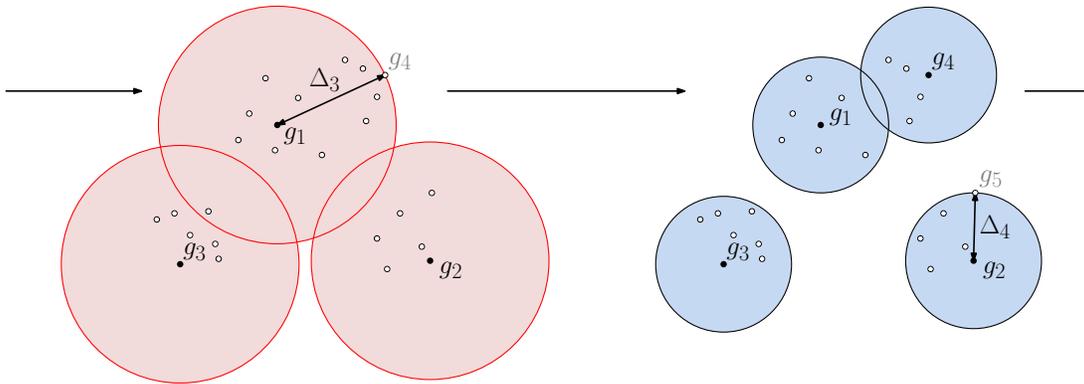


Fig. 30: Greedy approximation to  $k$ -center (from stage 3 to 4).

**Approximation Bound:** Now, let us show that this algorithm is at most a factor of 2 from the optimum. Let  $G = \{g_1, \dots, g_k\}$  denote the set of centers computed by the greedy algorithm, and let  $\Delta(G)$  denote its bottleneck distance. Let  $O = \{o_1, \dots, o_k\}$  denote the optimum set of centers, that is the set of  $k$  centers such that  $\Delta(O)$  is the smallest possible. We will show that greedy is at most twice as bad as optimal, that is,

$$\Delta(G) \leq 2\Delta(O).$$

Of course, we don't know what  $O$  is or even what  $\Delta(O)$  is. Our approach will be to determine a lower bound  $\Delta_{\min}$  on the optimum bottleneck distance ( $\Delta_{\min} \leq \Delta(O)$ ). We will then show that our

algorithm produces a value that is at most twice this lower bound value ( $\Delta(G) \leq 2\Delta_{\min}$ ). It will follow therefore that  $\Delta(G) \leq 2\Delta(O)$ .

The analysis is based on the following three claims, each of which is quite straightforward to prove. Define  $G_i$  to be the set of greedy centers after the  $i$ th execution of the algorithm, and let  $\Delta_i = \Delta(G_i)$  denote its overall bottleneck distance (the farthest any point is from its closest center in  $G_i$ ). The greedy algorithm stops with  $g_k$ , but for the sake of the analysis it is convenient to consider the next center to be added if we ran it for one more iteration. That is, define  $g_{k+1}$  to be the point of  $P$  that maximizes the distance to its closest center in  $G_k$ . (This distance is  $\Delta(G_k)$ .) Also, define  $G_{k+1} = \{g_1, \dots, g_{k+1}\}$ .

**Claim 1:** For  $1 \leq i \leq k+1$ ,  $\Delta_{i+1} \leq \Delta_i$ . That is, the sequence of bottleneck distances is monotonically nonincreasing. (In Fig. 30 this is represented by the fact that the radii of the covering disks decrease with each stage.)

**Proof:** Whenever we add a new center, the distance to each point's closest center will either be the same or will decrease. Therefore, the maximum of such a set can never increase.

**Claim 2:** For  $1 \leq i \leq k+1$ , every pair of greedy centers in  $G_i$  is separated by a distance of at least  $\Delta_{i-1}$ .

**Proof:** Consider the  $i$ th stage. By the induction hypothesis, the first  $i-1$  centers are separated from each other by distance  $\Delta_{i-2} \geq \Delta_{i-1}$ . The  $i$ th center is, by definition, at distance  $\Delta_{i-1}$  from its closest center, and therefore it is at distance at least  $\Delta_{i-1}$  from all the other centers.

**Claim 3:** Let  $\Delta_{\min} = \Delta(G)/2$ . Then for any set  $C$  of  $k$  cluster centers,  $\Delta(C) \geq \Delta_{\min}$ .

**Proof:** By definition of  $\Delta(C)$  we know that every point of  $P$  lies within distance  $\Delta(C)$  of some point of  $C$ , and since  $G \subseteq P$ , this is true for  $G$  as well. Since  $|G_{k+1}| = k+1$ , by the pigeonhole principle, there exists at least two centers  $g, g' \in G_{k+1}$  that are in the same neighborhood of some center  $c \in C$ , that is,  $\max(\delta(g, c), \delta(g', c)) \leq \Delta(c)$ . Since  $g, g' \in G_{k+1}$ , by Claim 2,  $\delta(g, g') \geq \Delta_k = \Delta(G)$ . By applying the triangle inequality to the triple  $(g, c, g')$  and symmetry of the distance function, we have

$$\begin{aligned} \Delta(G) &\leq \delta(g, g') \leq \delta(g, c) + \delta(c, g') \quad (\text{by triangle inequality}) \\ &\leq \delta(g, c) + \delta(g', c) \quad (\text{by distance symmetry}) \\ &\leq \Delta(c) + \Delta(c) \leq \Delta(C) + \Delta(C) = 2\Delta(C). \end{aligned}$$

Therefore,  $\Delta(C) \geq \Delta(G)/2 = \Delta_{\min}$ , as desired.

Now, by applying Claim 3 to  $O$ , we have  $\Delta(O) \geq \Delta_{\min}$ . By definition,  $\Delta_{\min} = \Delta(G)/2$ , and so  $\Delta(G) \leq 2\Delta(O)$ , completing the analysis.

You might wonder whether this bound is tight. We will leave it as an exercise to prove that there is a input such that (if you are unlucky about how you choose the first point) the greedy algorithm returns a value that is arbitrarily close to twice that of the optimum.

## Lecture 9: Greedy Approximation: Set Cover

**Set Cover:** An important class of optimization problems involves covering a certain domain, with sets of a certain characteristics. Many of these problems can be expressed abstractly as the *set cover problem*. We are given a pair  $\Sigma = (X, S)$ , called a *set system*, where  $X = \{x_1, \dots, x_m\}$  is a finite set of objects, called the *universe*, and  $S = \{s_1, \dots, s_n\}$  is a collection of subsets of  $X$ , such that every element of  $X$  belongs to at least one set of  $S$ . Set systems arise in many applications of science and engineering.

**Undirected Graph:** An undirected graph  $G = (V, E)$  is a set system where  $V$  constitutes the universe, and the edges  $E$  are subsets of cardinality two.

**Geometric set systems:** The universe consists of  $n$  points in space, and the sets are the subsets of points that are contained within some specified geometric shape (balls, cubes, rectangles, triangles, etc.)

**Wireless Coverage:** The universe consists all the locations on a college campus, and for each possible location of a wireless transmitter there is an associated region of the campus that is covered by placing a wireless transmitter at this location.

A fundamental question involving set systems is determining the smallest number of sets needed to cover the entire universe. A *cover* of  $S$  is defined to be a subcollection of sets whose union covers  $X$ . For example, in Fig. 31(a), we elements of  $X$  are the black circles, and the sets  $s_1, \dots, s_6$  are indicated by rectangles. In this case there exists a cover of size three, consisting of  $s_3, s_4,$  and  $s_5$  (see Fig. 31(b)). (The sets of this cover do not overlap, which is sometimes called an *exact cover*. In general, the sets of the cover are allowed to overlap.)

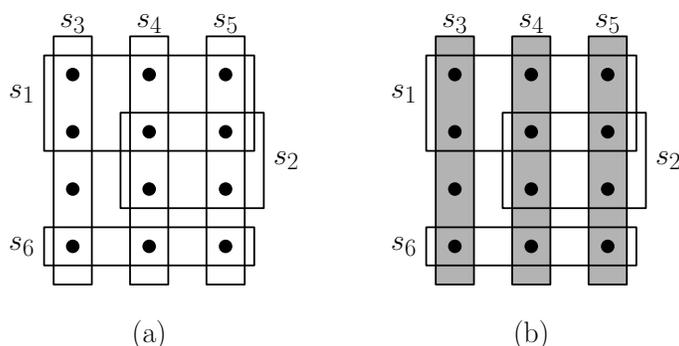


Fig. 31: Set cover. The optimum cover consists of the three sets  $\{s_3, s_4, s_5\}$ .

Notice that the output of set cover is not a set, but rather a set of sets. If we think of the sets of  $S$  as being indexed by the integers from 1 to  $n$ , then we can think of a cover  $C$  more conveniently as a subset of  $\{1, \dots, n\}$ . This suggests the following definition.

**Set Cover Problem:** Given a set system  $\Sigma = (X, S)$ , where  $S = \{s_1, \dots, s_n\}$ , compute a set  $C \subseteq \{1, \dots, n\}$  of minimum cardinality such that

$$X = \bigcup_{i \in C} s_i$$

The set cover problem is a very important and powerful optimization problem. It arises in a vast number of applications. Determining the fewest locations to place wireless transmitters to cover the entire campus is an example. Unfortunately, the set cover problem is known to be NP-hard. We will present a simple *greedy heuristic* for this problem.

How do we determine how good our approximation is? Given an input instance  $\Sigma = (X, S)$ , let  $O(\Sigma)$  denote an optimum cover (of minimum cardinality) and let  $G(\Sigma)$  denote the cover produced by our greedy heuristic. Clearly, greedy cannot have fewer sets than the optimum, and so we have  $|G(\Sigma)| \geq |O(\Sigma)|$ . We say that  $G$  achieves an *approximation ratio* of  $\rho$  if  $|G(\Sigma)| \leq \rho |O(\Sigma)|$ , for any input  $\Sigma$ . Ideally, we would like  $\rho$  to be as small as possible, say, a small constant. Unfortunately, the best that we can show for set cover is that  $\rho$  is a slowly growing function of  $m = |X|$ , and in particular  $\rho = \ln m$ . (This might strike you as being rather weak, but there are compelling reasons from the theory of computational complexity that logarithmic approximation ratio is the best that we might hope for assuming that  $P \neq NP$ . With a bit more work, it is possible to improve this slightly to an approximation ratio of  $\rho = (\ln m')$ , where  $m'$  is the maximum cardinality of any set of  $S$ .)

**Greedy Set Cover:** A simple greedy approach to set cover works by at each stage selecting the set that covers the greatest number of uncovered elements. The algorithm is presented in the code block below. The set  $C$  contains the indices of the sets of the cover, and the set  $U$  stores the elements of  $X$  that are still uncovered. Initially,  $C$  is empty and  $U \leftarrow X$ . We repeatedly select the set of  $S$  that covers the greatest number of elements of  $U$  and add it to the cover.

Greedy Set Cover

---

```

Greedy-Set-Cover( $X, S$ ) {
   $U = X$  //  $U$  stores the uncovered items
   $C = \text{empty}$  //  $C$  stores the sets of the cover
  while ( $U$  is nonempty) {
    select  $s[i]$  in  $S$  that covers the most elements of  $U$ 
    add  $i$  to  $C$ 
    remove the elements of  $s[i]$  from  $U$ 
  }
  return  $C$ 
}

```

---

We will not worry about implementing this algorithm in the most efficient manner. If we assume that  $U$  and the sets  $s_i$  are each represented as a simple list of elements of  $X$  (each of length at most  $m$ ), then we can perform each iteration of the main while loop in time  $O(mn)$ , for a total running time of  $O(mn^2)$ .

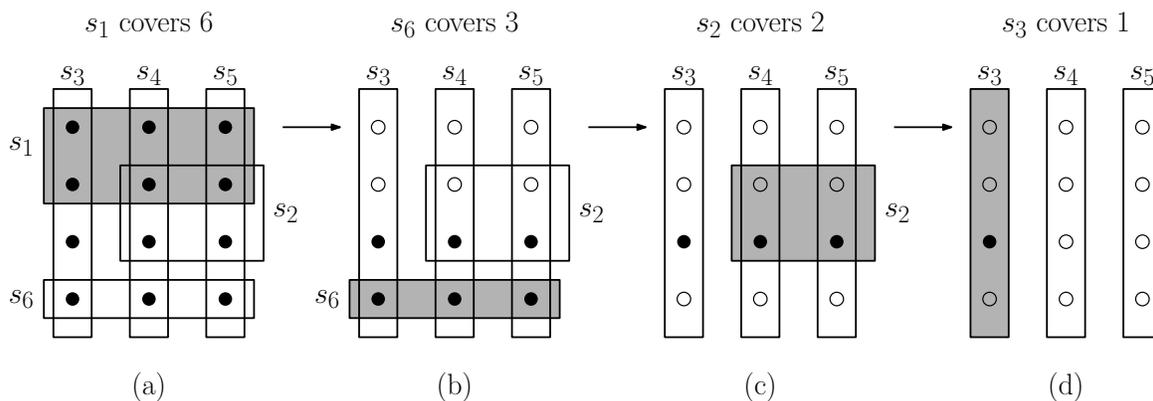


Fig. 32: The greedy heuristic. Final cover is  $\{s_1, s_6, s_2, s_3\}$ .

For the example given earlier the greedy-set cover algorithm would select  $s_1$  (see Fig. 32(a)), then  $s_6$  (see Fig. 32(b)), then  $s_2$  (see Fig. 32(c)) and finally  $s_3$  (see Fig. 32(d)). Thus, it would return a set cover of size four, whereas the optimal set cover has size three.

**What is the approximation factor?** The problem with the greedy heuristic is that it can be “fooled” into picking the wrong set, over and over again. Consider the example shown in Fig. 33 involving a universe of 32 elements. The optimal set cover consists of sets  $s_7$  and  $s_8$ , each of size 16. Initially all three sets  $s_1$ ,  $s_7$ , and  $s_8$  have 16 elements. If ties are broken in the worst possible way, the greedy algorithm will first select the set  $s_1$ . We remove all the covered elements. Now  $s_2$ ,  $s_7$  and  $s_8$  all cover eight of the remaining elements. Again, if we choose poorly,  $s_2$  is chosen. The pattern repeats, choosing  $s_3$  (covering four of the remainder),  $s_4$  (covering two) and finally  $s_5$  and  $s_6$  (each covering one). Although there are ties for the greedy choice in this example, it is easy to modify the example so that the greedy choice is unique.

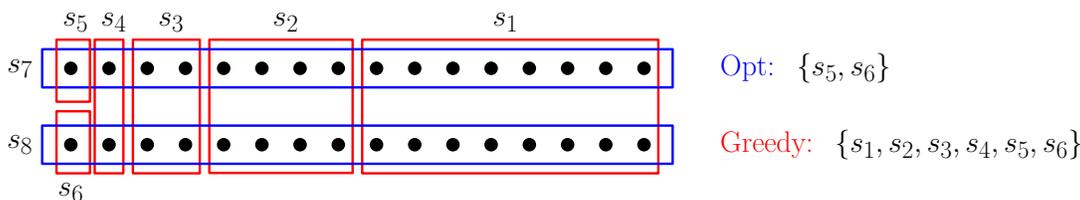


Fig. 33: Repeatedly fooling the greedy heuristic.

From the pattern, you can see that we can generalize this to any number of elements that is a power of 2. While there is an optimal solution with 2 sets, the greedy algorithm will select roughly  $\lg m$  sets, where  $m = |X|$ . (Recall that “lg” denotes logarithm base 2, and “ln” denotes the natural logarithm.) Thus, on this example the greedy heuristic achieves an approximation factor of roughly  $(\lg m)/2$ . There were many cases where ties were broken badly here, but it is possible to redesign the example such that there are no ties, and yet the algorithm has essentially the same ratio bound.

We will show that the greedy set cover heuristic never performs worse than a factor of  $\ln m$ . Before giving the proof, we need one useful technical inequality.

**Lemma:** For all  $c > 0$ ,

$$\left(1 - \frac{1}{c}\right)^c \leq \frac{1}{e}.$$

where  $e$  is the base of the natural logarithm.

**Proof:** We use the fact that for any real  $z$  (positive, zero, or negative),  $1 + z \leq e^z$ . (This follows from the Taylor’s expansion  $e^z = 1 + z + \frac{z^2}{2!} + \frac{z^3}{3!} + \dots \geq 1 + z$ .) Now, if we substitute  $-1/c$  for  $z$  we have  $(1 - 1/c) \leq e^{-1/c}$ . By raising both sides to the  $c$ th power, we have the desired result.

We now prove the approximation bound.

**Theorem:** Given any set system  $\Sigma = (X, S)$ , let  $G$  be the output of the greedy heuristic and let  $O$  be an optimum cover. Then  $|G| \leq |O| \cdot \ln m$ , where  $m = |X|$ .

**Proof:** We will cheat a bit. Let  $c$  denote the size of the optimum set cover, and let  $g$  denote the size of the greedy set cover minus 1. We will show that  $g \leq c \cdot \ln m$ . (Note that we should really show that  $g + 1 \leq c \cdot \ln m$ , but this is close enough and saves us some messy details.)

Let’s consider how many new elements we cover with each round of the algorithm. Initially, there are  $m_0 = m$  elements to be covered. Let  $m_i$  denote the number of elements remaining to be covered after  $i$  iterations of the greedy algorithm. After  $i - 1$  rounds there are  $m_{i-1}$  elements that remain to be covered. We know that there is a cover of size  $c$  for these elements (namely, the optimal cover), and so by the pigeonhole principle there exists some set that covers at least  $m_{i-1}/c$  elements. Since the greedy algorithm selects the set covering the largest number of remaining elements, it must select a set that covers at least this many elements. The number of elements that remain to be covered is at most

$$m_i \leq m_{i-1} - \frac{m_{i-1}}{c} = m_{i-1} \left(1 - \frac{1}{c}\right).$$

Thus, with each iteration the number of remaining elements decreases by a factor of at least  $(1 - 1/c)$ . If we repeat this  $i$  times, we have

$$m_i \leq m_0 \left(1 - \frac{1}{c}\right)^i = m \left(1 - \frac{1}{c}\right)^i.$$

How long can this go on? Since the greedy heuristic ran for  $g + 1$  iterations, we know that just prior to the last iteration we must have had at least one remaining uncovered element, and so we have

$$1 \leq m_g \leq m \left(1 - \frac{1}{c}\right)^g = m \left(\left(1 - \frac{1}{c}\right)^c\right)^{g/c}.$$

(In the last step, we just rewrote the expression in a manner that makes it easier to apply the above technical lemma.) By the above lemma we have

$$1 \leq m \left(\frac{1}{e}\right)^{g/c}.$$

Now, if we multiply by  $e^{g/c}$  on both sides and take natural logs we find that  $g$  satisfies:

$$e^{g/c} \leq m \quad \Rightarrow \quad \frac{g}{c} \leq \ln m \quad \Rightarrow \quad g \leq c \cdot \ln m.$$

Therefore (ignoring the missing “+1” as mentioned above) the greedy set cover is larger than the optimum set cover by a factor of at most  $\ln m$ .

## Lecture 10: Dynamic Programming: Weighted Interval Scheduling

**Dynamic Programming:** In this lecture we begin our coverage of an important algorithm design technique, called *dynamic programming* (or *DP* for short). The technique is among the most powerful for designing algorithms for optimization problems. Dynamic programming is a powerful technique for solving optimization problems that have certain well-defined clean structural properties. (The meaning of this will become clearer once we have seen a few examples.)

There is a superficial resemblance to divide-and-conquer, in the sense that it breaks problems down into smaller subproblems, which can be solved recursively. However, unlike divide-and-conquer problems, in which the subproblems are disjoint, in dynamic programming the subproblems typically overlap each other, and this renders straightforward recursive solutions inefficient.

Dynamic programming solutions rely on two important structural qualities, *optimal substructure* and *overlapping subproblems*.

**Optimal substructure:** This property (sometimes called the *principle of optimality*) states that for the global problem to be solved optimally, each subproblem should be solved optimally. While this might seem intuitively obvious, not all optimization problems satisfy this property. For example, it may be advantageous to solve one subproblem suboptimally in order to conserve resources so that another, more critical, subproblem can be solved optimally.

**Overlapping Subproblems:** While it may be possible subdivide a problem into subproblems in exponentially many different ways, these subproblems overlap each other in such a way that the number of distinct subproblems is reasonably small, ideally *polynomial* in the input size.

An important issue is how to generate the solutions to these subproblems. There are two complementary (but essentially equivalent) ways of viewing how a solution is constructed:

**Top-Down:** A solution to a DP problem is expressed recursively. This approach applies recursion directly to solve the problem. However, due to the overlapping nature of the subproblems, the same recursive call is often made many times. An approach, called *memoization*, records the results of recursive calls, so that subsequent calls to a previously solved subproblem are handled by table look-up.

**Bottom-up:** Although the problem is formulated recursively, the solution is built iteratively by combining the solutions to small subproblems to obtain the solution to larger subproblems. The results are stored in a table.

In the next few lectures, we will consider a number of examples, which will help make these concepts more concrete.

**Weighted Interval Scheduling:** Let us consider a variant of a problem that we have seen before, the Interval Scheduling Problem. Recall that in the original (unweighted) version we are given a set  $S = \{1, \dots, n\}$  of  $n$  activity requests, where each activity is expressed as an interval  $[s_i, f_i]$  from a given start time  $s_i$  to a given finish time  $f_i$ . The objective is to compute any maximum sized subset of non-overlapping intervals (see Fig. 34(a)).

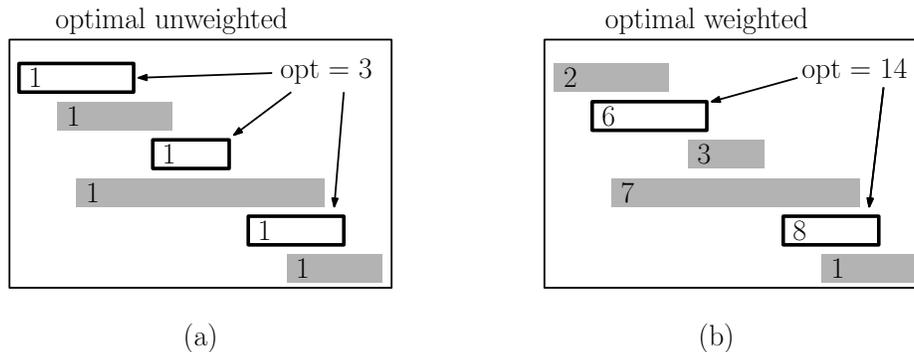


Fig. 34: Weighted and unweighted interval scheduling.

In *weighted interval scheduling*, we assume that in addition to the start and finish times, each request is associated with a numeric *weight* or *value*, call it  $v_i$ , and the objective is to find a set of non-overlapping requests such that sum of values of the scheduled requests is maximum (see Fig. 34(b)). The unweighted version can be thought of as a special case in which all weights are equal to 1. Although a greedy approach works fine for the unweighted problem, no greedy solution is known for the weighted version. We will demonstrate a method based on dynamic programming.

**Recursive Formulation:** Dynamic programming solutions are based on a decomposition of a problem into smaller subproblems. Let us consider how to do this for the weighted interval scheduling problem. As we did in the greedy algorithm, it will be convenient to sort the requests in nondecreasing order of finish time, so that  $f_1 \leq \dots \leq f_n$ .

Here is the idea behind DP in a nutshell. Consider the *last* request  $[s_n, f_n]$ . There are two possibilities. If this request *is not* in the optimum schedule, then we can safely ignore it, and recursively compute the optimum solution of the first  $n - 1$  requests. Otherwise, this request *is* in the optimal solution. We will schedule this request (receiving the profit of  $v_n$ ) and then we must eliminate all the requests whose intervals overlap this one. Because requests have been sorted by finish time, this involves finding the largest index  $p$  such that  $f_p < s_n$ . Thus, we solve the problem recursively on the first  $p$  requests.

But we don't know the optimum solution, so how can we select among these two options? The answer is that we will compute the cost of both of them recursively, and take the better of the two.

Let's now implement this idea. Recall that the requests are sorted by finish times. For the sake of generality, let's assume that we want to solve the problem on requests 1 through  $j$ , where  $0 \leq j \leq n$ . If  $j = 0$ , there is nothing to do. Otherwise, given any request  $j$ , define  $p(j)$  to be the largest integer such that  $f_{p(j)} < s_j$ . Thus, request 1 through  $p(j)$  do not overlap request  $j$ . If no such  $i$  exists, let  $p(j) = 0$  (see Fig. 35).

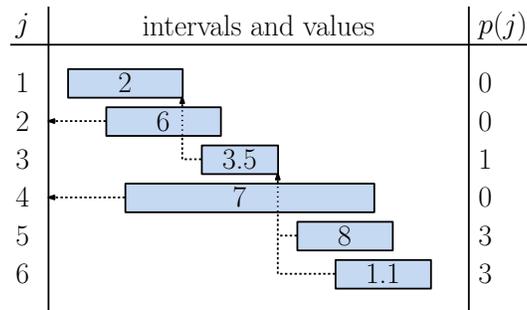


Fig. 35: Weighted interval scheduling input and  $p$ -values.

For now, let's just concentrate on computing the *optimum total value*. Later we will consider how to determine which *requests* produce this value. A natural idea would be to define a function that gives the optimum value for just the first  $i$  requests.

**Definition:** For  $0 \leq j \leq n$ ,  $\text{opt}(j)$  denotes the maximum possible value achievable if we consider just tasks  $\{1, \dots, j\}$  (assuming that tasks are given in order of finish time).

As a basis case, observe that if we have no tasks, then we have no value. Therefore,  $\text{opt}(0) = 0$ . If we can compute  $\text{opt}(j)$  for each value of  $j$ , then clearly, the final desired result will be the maximum value using *all* the requests, that is,  $\text{opt}(n)$ .

Summarizing our earlier observations, in order to compute  $\text{opt}(j)$  for an arbitrary  $j$ , we observe that there are two possibilities:

**Request  $j$  is not in the optimal schedule:** If  $j$  is not included in the schedule, then we should do the best we can with the remaining  $j - 1$  requests. Therefore,  $\text{opt}(j) = \text{opt}(j - 1)$ .

**Request  $j$  is in the optimal schedule:** If we add request  $j$  to the schedule, then we gain  $v_j$  units of value, but we are now limited as to which other requests we can take. We cannot take any of the requests following  $p(j)$ . Thus we have  $\text{opt}(j) = v_j + \text{opt}(p(j))$ .

How do we know which of these two options to select? The danger in trying to be too smart (e.g., trying a greedy choice) is that the choice may not be correct. Instead, the best approach is often simple brute force:

**DP Selection Principle:**  
When given a set of feasible options to choose from, try them *all* and take the *best*.

In this cases there are two options, which suggests the following recursive rule:

$$\text{opt}(j) = \max \left\{ \begin{array}{ll} \text{opt}(j - 1) & \text{(request } j \text{ is not in opt)} \\ v_j + \text{opt}(p(j)) & \text{(request } j \text{ is in opt)} \end{array} \right\}.$$

We could express this in pseudocode as shown in the following code block:

I have left it as self-evident that this simple recursive procedure is correct. Indeed the only subtlety is the inductive observation that, in order to compute  $\text{opt}(j)$  optimally, the two subproblems that are used to make the final result  $\text{opt}(j - 1)$  and  $\text{opt}(p(j))$  should also be computed optimally. This is an example of the principle of optimality, which in this case is clear.<sup>5</sup>

<sup>5</sup>Why would you want to solve a subproblem suboptimally? Suppose, that you had the additional constraint that the final schedule can only contain 23 intervals. Now, it might be advantageous to solve a subproblem suboptimally, so that you have a few extra intervals left over to use at a later time.

---

```

recursive-WIS(j) {
  if (j == 0) return 0
  else return max( recursive-WIS(j-1), v[j] + recursive-WIS(p[j]) )
}

```

---

**Memoized Formulation:** The principal problem with this elegant and simple recursive procedure is that it has a *horrendous* running time. To make this concrete, let us suppose that  $p(j) = j - 2$  for all  $j$ . Let  $T(j)$  be the number of recursive function calls to  $\text{opt}(0)$  that result from a single call to  $\text{opt}(j)$ . Clearly,  $T(0) = 1$ ,  $T(1) = T(0) + T(0)$ , and for  $j \geq 2$ ,  $T(j) = T(j - 1) + T(j - 2)$ . If you start expanding this recurrence, you will find that the resulting series is essentially a Fibonacci series:

$j$	0	1	2	3	4	5	6	7	8	...	20	30	50
$T(j)$	1	2	3	5	8	13	21	34	55	...	17,711	2,178,309	32,951,280,099

It is well known that the Fibonacci series  $F(j)$  grows exponentially as a function of  $j$ .<sup>6</sup> This may seem ludicrous. (And it is!) Why should it take 32 billion recursive calls to fill in a table with just 50 entries? If you look at the recursion tree, the reason jumps out immediately (see Fig 36). The problem is that the same recursive calls are being generated over and over again. But there is no reason to make even a second call this, since they all return exactly the same value.

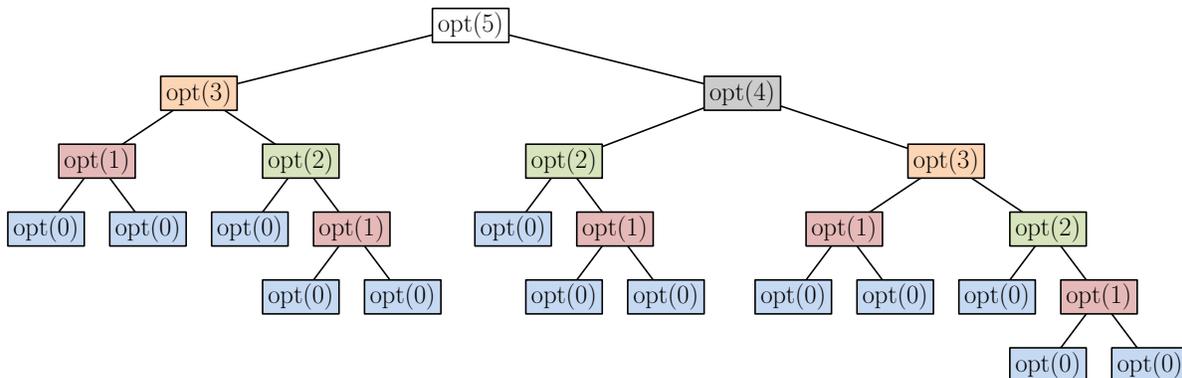


Fig. 36: The exponential nature of recursive-WIS.

This suggests a smarter version of the algorithm. Once a value has been computed for  $\text{recursive-WIS}(j)$  we store the value in a global array  $M[1..n]$ , and all future attempts to compute this value will simply access the array, rather than making a recursive call. This technique is called *memoizing* (I guess because we are making a memo to ourselves of what the value is for future reference). You might imagine that we initialize all the  $M[j]$  entries to  $-1$  initially, and use this special value to determine whether an entry has already been computed.

We will add one additional piece of information, which will help in reconstructing the final schedule. Whenever a choice is made between two options, we'll save a *predecessor pointer*,  $\text{pred}[j]$ , which reminds of which choice we made. (Its value is either  $j - 1$  or  $p(j)$ , depending on which recursive call was made). The resulting algorithm is presented in the following code block.

The memoized version runs in  $O(n)$  time. To see this, observe that each invocation of memoized-WIS either returns in  $O(1)$  time (with no recursive calls), or it computes one new entry of  $M$ . Since there are  $n$  entries in the table, the latter can occur at most  $n$  times.

---

<sup>6</sup>The  $n$ th Fibonacci number is roughly grows as  $\Theta(\phi^n)$ , where  $\phi \approx 1.618$  is the celebrated Golden Ratio.

---

```

memoized-WIS(j) {
  if (j == 0) return 0 // basis case - no requests
  else if (M[j] has been computed) return M[j]
  else {
    leaveWeight = memoized-WIS(j-1) // total weight if we leave j
    takeWeight = v[j] + memoized-WIS(p[j]) // total weight if we take j
    if ( leaveWeight > takeWeight ) {
      M[j] = leaveWeight // better to leave j
      pred[j] = j-1 // previous is j-1
    } else {
      M[j] = takeWeight // better to take j
      pred[j] = p[j] // previous is p[j]
    }
    return M[j] // return final weight
  }
}

```

---

**Bottom-up Construction:** (Optional) Yet another method for computing the values of the array, is to dispense with the recursion altogether, and simply fill the table up, one entry at a time. We need to be careful that this is done in such an order that each time we access the array, the entry being accessed is already defined. This is easy here, because we can just compute values in increasing order of  $j$ . As before, we include the computation of the predecessor values.

---

```

bottom-up-WIS() {
  M[0] = 0
  for (i = 1 to n) {
    leaveWeight = M[j-1] // total weight if we leave j
    takeWeight = v[j] + M[p[j]] // total weight if we take j
    if ( leaveWeight > takeWeight ) {
      M[j] = leaveWeight // better to leave j
      pred[j] = j-1 // previous is j-1
    }
    else {
      M[j] = takeWeight // better to take j
      pred[j] = p[j] // previous is p[j]
    }
  }
}

```

---

Do you think that you understand the algorithm now? If so, answer the following question. Would the algorithm be correct if, rather than sorting the requests by finish time, we had instead sorted them by start time? How about if we didn't sort them at all?

**Computing the Final Schedule:** So far we have seen how to compute the value of the optimal schedule, but how do we compute the schedule itself? This is a common problem that arises in many DP problems, since most DP formulations focus on computing the numeric optimal value, without consideration of the object that gives rise to this value. The solution is to leave ourselves a few hints in order to reconstruct the final result.

In `bottom-up-WIS()` we did exactly this. We know that value of  $M[j]$  arose from two distinct possibilities, either (1) we didn't take  $j$  and just used the result of  $M[j - 1]$ , or (2) we did take  $j$ , added its

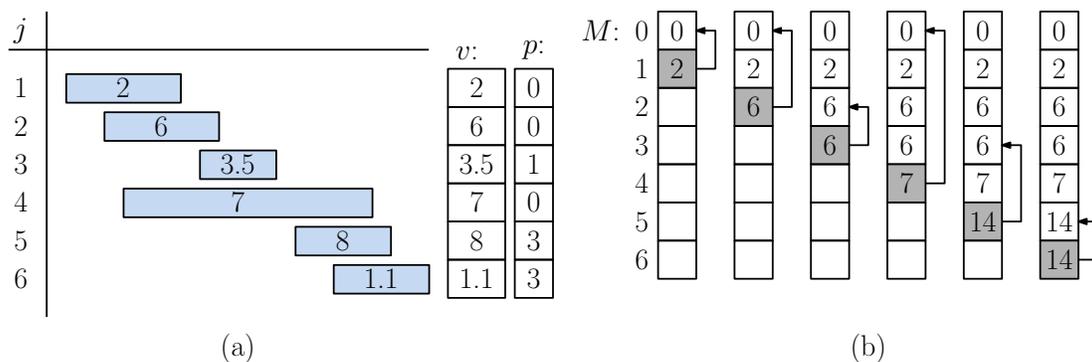


Fig. 37: (a) The input intervals and  $p$  values, (b) the bottom-up construction of the table and predecessor values, (c) the final predecessor values.

value  $v_j$ , and used  $M[p(j)]$  to complete the result. To remind us of how we obtained the best choice for  $M[j]$  was to store a predecessor pointer  $\text{pred}[j]$ .

In order to generate the final schedule, we start with  $M[n]$  and work backwards. In general, if we arrive at  $M[j]$ , we check whether  $\text{pred}[j] = p[j]$ . If so, we can surmise that we did use the  $j$ th request, and we continue with  $\text{pred}[j] = p[j]$ . If not, then we know that we did not include request  $j$  in the schedule, and we then follow the predecessor link to continue with  $\text{pred}[j] = j - 1$ . The algorithm for generating the schedule is given in the code block below.

---

Computing Weighted Interval Scheduling Schedule

```

get-schedule() {
    j = n // start with the last request
    sched = empty
    while (j > 0) {
        if (pred[j] == p[j]) { // take request j
            prepend j to the front of sched
        }
        j = pred[j] // continue with j's predecessor
    }
    return sched // return the final schedule
}

```

---

The computation of the final schedule is illustrated in Fig. 38 (where predecessor values are shown as arrows).

- Since  $\text{pred}[6] = 5 \neq p[6]$ , we do *not* use request 6, and we continued with  $\text{pred}[6] = 5$ .
- Since  $\text{pred}[5] = 3 = p[5]$ , we *use* request 5, and we continue with  $\text{pred}[5] = 3$ .
- Since  $\text{pred}[3] = 2 \neq p[3]$ , we do *not* use request 3, and we continued with  $\text{pred}[3] = 2$ .
- Since  $\text{pred}[2] = 0 = p[2]$ , we *use* request 2, and we continue with  $\text{pred}[2] = 0$ , and terminate.

We obtain the final list  $\langle 5, 2 \rangle$ .

## Lecture 11: Dynamic Programming: Longest Common Subsequence

**Strings:** One important area of algorithm design is the study of algorithms for character strings. Finding patterns or similarities within strings is fundamental to various applications, ranging from document

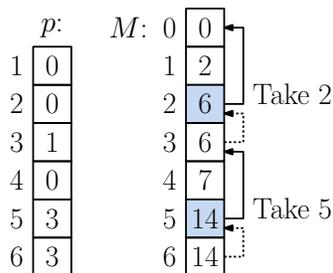


Fig. 38: Following the predecessor links to compute the final schedule.

analysis to computational biology. One common measure of similarity between two strings is the lengths of their longest common subsequence. Today, we will consider an efficient solution to this problem based on dynamic programming.

**Longest Common Subsequence:** Let us think of character strings as sequences of characters. Given two sequences  $X = \langle x_1, x_2, \dots, x_m \rangle$  and  $Z = \langle z_1, z_2, \dots, z_k \rangle$ , we say that  $Z$  is a *subsequence* of  $X$  if there is a strictly increasing sequence of  $k$  indices  $\langle i_1, i_2, \dots, i_k \rangle$  ( $1 \leq i_1 < i_2 < \dots < i_k \leq m$ ) such that  $Z = \langle x_{i_1}, x_{i_2}, \dots, x_{i_k} \rangle$ . For example, let  $X = \langle \text{ABRACADABRA} \rangle$  and let  $Z = \langle \text{AADAA} \rangle$ , then  $Z$  is a subsequence of  $X$ .

Given two strings  $X$  and  $Y$ , the *longest common subsequence* of  $X$  and  $Y$  is a longest sequence  $Z$  that is a subsequence of both  $X$  and  $Y$ . For example, let  $X = \langle \text{ABRACADABRA} \rangle$  and let  $Y = \langle \text{YABBADABBADOO} \rangle$ . Then the longest common subsequence is  $Z = \langle \text{ABADABA} \rangle$  (see Fig. 39).

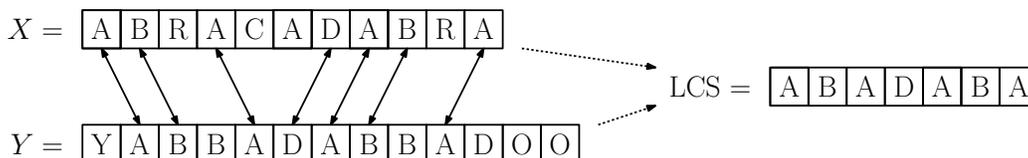


Fig. 39: An example of the LCS of two strings  $X$  and  $Y$ .

The *Longest Common Subsequence Problem* (LCS) is the following. Given two sequences  $X = \langle x_1, \dots, x_m \rangle$  and  $Y = \langle y_1, \dots, y_n \rangle$  determine the length of their longest common subsequence, and more generally the sequence itself. Note that the subsequence is not necessarily unique. For example the LCS of  $\langle \text{ABC} \rangle$  and  $\langle \text{BAC} \rangle$  is either  $\langle \text{AC} \rangle$  or  $\langle \text{BC} \rangle$ .

**DP Formulation for LCS:** The simple brute-force solution to the problem would be to try all possible subsequences from one string, and search for matches in the other string, but this is hopelessly inefficient, since there are an exponential number of possible subsequences.

Instead, we will derive a dynamic programming solution. In typical DP fashion, we need to break the problem into smaller pieces. There are many ways to do this for strings, but it turns out for this problem that considering all pairs of *prefixes* will suffice for us. A *prefix* of a sequence is just an initial string of values,  $X_i = \langle x_1, \dots, x_i \rangle$ .  $X_0$  is the empty sequence.

The idea will be to compute the longest common subsequence for every possible pair of prefixes. Let  $\text{lcs}(i, j)$  denote the length of the longest common subsequence of  $X_i$  and  $Y_j$ . For example, in the above case we have  $X_5 = \langle \text{ABRAC} \rangle$  and  $Y_6 = \langle \text{YABBAD} \rangle$ . Their longest common subsequence is  $\langle \text{ABA} \rangle$ . Thus,  $\text{lcs}(5, 6) = 3$ .

Let us start by deriving a recursive formulation for computing  $\text{lcs}(i, j)$ . As we have seen with other DP problems, a naive implementation of this recursive rule will lead to a very inefficient algorithm. Rather

than implementing it directly, we will use one of the other techniques (memoization or bottom-up computation) to produce a more efficient algorithm.

**Basis:** If either sequence is empty, then the longest common subsequence is empty. Therefore,  $\text{lcs}(i, 0) = \text{lcs}(j, 0) = 0$ .

**Last characters match:** Suppose  $x_i = y_j$ . For example: Let  $X_i = \langle ABCA \rangle$  and let  $Y_j = \langle DACA \rangle$ . Since both end in ‘A’, it is easy to see that the LCS *must* also end in ‘A’.<sup>7</sup> Also, there is no harm in assuming that the last two characters of both strings will be matched to each other, since matching the last ‘A’ of one string to an earlier instance of ‘A’ of the other can only limit our future options.

Since the ‘A’ is the last character of the LCS, we may find the overall LCS by (1) removing ‘A’ from both sequences, (2) taking the LCS of  $X_{i-1} = \langle ABC \rangle$  and  $Y_{j-1} = \langle DAC \rangle$  which is  $\langle AC \rangle$ , and (3) adding ‘A’ to the end. This yields  $\langle ACA \rangle$  as the LCS. Therefore, the length of the final LCS is the length of  $\text{lcs}(X_{i-1}, Y_{j-1}) + 1$  (see Fig. 40), which provides us with the following rule:

$$\text{if } (x_i = y_j) \text{ then } \text{lcs}(i, j) = \text{lcs}(i - 1, j - 1) + 1$$

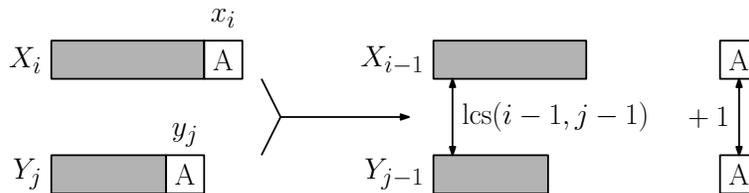


Fig. 40: LCS of two strings whose last characters are equal.

**Last characters do not match:** Suppose that  $x_i \neq y_j$ . In this case  $x_i$  and  $y_j$  cannot both be in the LCS (since they would have to be the last character of the LCS). Thus either  $x_i$  is *not* part of the LCS, or  $y_j$  is *not* part of the LCS (and possibly *both* are not part of the LCS).

At this point it may be tempting to try to make a “smart” choice. By analyzing the last few characters of  $X_i$  and  $Y_j$ , perhaps we can figure out which character is best to discard. However, this approach is doomed to failure (and you are strongly encouraged to think about this, since it is a common point of confusion). Remember the DP selection principle: *When given a set of feasible options to choose from, try them all and take the best.* Let’s consider both options, and see which one provides the better result.

**Option 1:** ( $x_i$  is not in the LCS) Since we know that  $x_i$  is out, we can infer that the LCS of  $X_i$  and  $Y_j$  is the LCS of  $X_{i-1}$  and  $Y_j$ , which is given by  $\text{lcs}(i - 1, j)$ .

**Option 2:** ( $y_j$  is not in the LCS) Since  $y_j$  is out, we can infer that the LCS of  $X_i$  and  $Y_j$  is the LCS of  $X_i$  and  $Y_{j-1}$ , which is given by  $\text{lcs}(i, j - 1)$ .

We compute both options and take the one that gives us the longer LCS (see Fig. 41).

$$\text{if } (x_i \neq y_j) \text{ then } \text{lcs}(i, j) = \max(\text{lcs}(i - 1, j), \text{lcs}(i, j - 1))$$

Combining these observations we have the following recursive formulation:

$$\text{lcs}(i, j) = \begin{cases} 0 & \text{if } i = 0 \text{ or } j = 0, \\ \text{lcs}(i - 1, j - 1) + 1 & \text{if } i, j > 0 \text{ and } x_i = y_j, \\ \max(\text{lcs}(i - 1, j), \text{lcs}(i, j - 1)) & \text{if } i, j > 0 \text{ and } x_i \neq y_j. \end{cases}$$

As mentioned earlier, a direct recursive implementation of this rule will be very inefficient. Let’s consider two alternative approaches to computing it.

<sup>7</sup>We will leave the formal proof as an exercise, but intuitively this is proved by contradiction. If the LCS did not end in ‘A’, then we could make it longer by adding ‘A’ to its end.

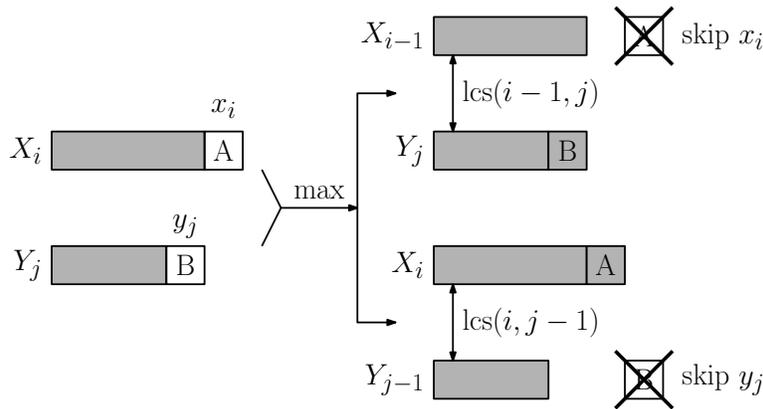


Fig. 41: The possible cases in the DP formulation of LCS.

**Memoized implementation:** The principal source of the inefficiency in a naive implementation of the recursive rule is that it makes repeated calls to  $\text{lcs}(i, j)$  for the same values of  $i$  and  $j$ . To avoid this, it creates a 2-dimensional array  $\text{lcs}[0..m, 0..n]$ , where  $m = |X|$  and  $n = |Y|$ . The memoized version first checks whether the requested value has already been computed, and if so, it just returns the stored value. Otherwise, it invokes the recursive rule to compute it. See the code block below. The initial call is  $\text{memoized-lcs}(m, n)$ .

---

Memoized Longest Common Subsequence

```

memoized-lcs(i,j) {
  if (lcs[i,j] has not yet been computed) {
    if (i == 0 || j == 0) // basis case
      lcs[i,j] = 0
    else if (x[i] == y[j]) // last characters match
      lcs[i,j] = memoized-lcs(i-1, j-1) + 1
    else // last chars don't match
      lcs[i,j] = max(memoized-lcs(i-1, j), memoized-lcs(i, j-1))
  }
  return lcs[i,j] // return stored value
}

```

---

The running time of the memoized version is  $O(mn)$ . To see this, observe that there are  $m + 1$  possible values for  $i$ , and  $n + 1$  possible values for  $j$ . Each time we call  $\text{memoized-lcs}(i, j)$ , if it has already been computed then it returns in  $O(1)$  time. Each call to  $\text{memoized-lcs}(i, j)$  generates a constant number of additional calls. Therefore, the time needed to compute the initial value of any entry is  $O(1)$ , and all subsequent calls with the same arguments is  $O(1)$ . Thus, the total running time is equal to the number of entries computed, which is  $O((m + 1)(n + 1)) = O(mn)$ .

**Bottom-up implementation:** The alternative to memoization is to just create the lcs table in a bottom-up manner, working from smaller entries to larger entries. By the recursive rules, in order to compute  $\text{lcs}[i, j]$ , we need to have already computed  $\text{lcs}[i - 1, j - 1]$ ,  $\text{lcs}[i - 1, j]$ , and  $\text{lcs}[i, j - 1]$ . Thus, we can compute the entries row-by-row or column-by-column in increasing order. See the code block below and Fig. 42(a). The running time and space used by the algorithm are both clearly  $O(mn)$ .

**Extracting the LCS:** The algorithms given so far compute only the length of the LCS, not the actual sequence. The remedy is common to many other DP algorithms. Whenever we make a decision, we save some information to help us recover the decisions that were made. We then work backwards,

```

bottom-up-lcs() {
  lcs = new array [0..m, 0..n]
  for (i = 0 to m) lcs[i,0] = 0           // basis cases
  for (j = 0 to n) lcs[0,j] = 0
  for (i = 1 to m) {                     // fill rest of table
    for (j = 1 to n) {
      if (x[i] == y[j])                 // take x[i] (= y[j]) for LCS
        lcs[i,j] = lcs[i-1, j-1] + 1
      else
        lcs[i,j] = max(lcs[i-1, j], lcs[i, j-1])
    }
  }
  return lcs[m, n]                       // final lcs length
}

```

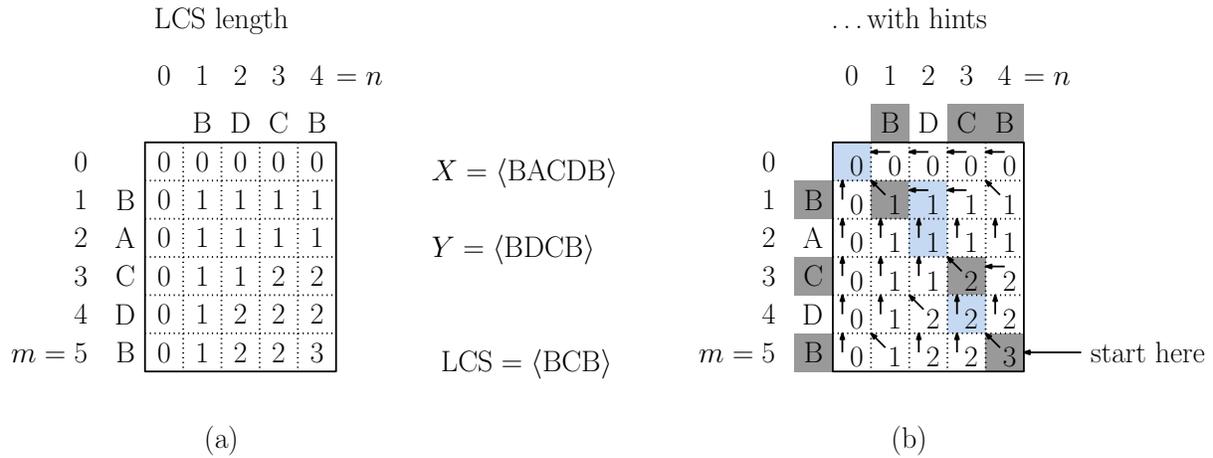


Fig. 42: Contents of the lcs array for the input sequences  $X = \langle BACDB \rangle$  and  $Y = \langle BCDB \rangle$ . The numeric table entries are the values of  $lcs[i, j]$  and the arrow entries are used in the extraction of the sequence.

unraveling these decisions to determine all the decisions that led to the optimal solution. In particular, the algorithm performs three possible actions:

**add<sub>XY</sub>**: Add  $x_i(=y_j)$  to the LCS (‘↖’ in Fig. 42(b)) and continue with  $\text{lcs}[i-1, j-1]$

**skip<sub>X</sub>**: Do not include  $x_i$  to the LCS (‘↑’ in Fig. 42(b)) and continue with  $\text{lcs}[i-1, j]$

**skip<sub>Y</sub>**: Do not include  $y_j$  to the LCS (‘←’ in Fig. 42(b)) and continue with  $\text{lcs}[i, j-1]$

An updated version of the bottom-up computation with these added hints is shown in the code block below and Fig. 42(b).

---

Bottom-up Longest Common Subsequence with Hints

```
bottom-up-lcs-with-hints() {
  lcs = new array [0..m, 0..n]           // stores lcs lengths
  h = new array [0..m, 0..n]           // stores hints
  for (i = 0 to m) { lcs[i,0] = 0; h[i,0] = skipX }
  for (j = 0 to n) { lcs[0,j] = 0; h[0,j] = skipY }
  for (i = 1 to m) {
    for (j = 1 to n) {
      if (x[i] == y[j])
        { lcs[i,j] = lcs[i-1, j-1] + 1; h[i,j] = addXY }
      else if (lcs[i-1, j] >= lcs[i, j-1])
        { lcs[i,j] = lcs[i-1, j]; h[i,j] = skipX }
      else
        { lcs[i,j] = lcs[i, j-1]; h[i,j] = skipY }
    }
  }
  return lcs[m, n]                       // final lcs length
}
```

---

How do we use the hints to reconstruct the answer? We start at the the last entry of the table, which corresponds to  $\text{lcs}(m, n)$ . In general, suppose that we are visiting the entry corresponding to  $\text{lcs}(m, n)$ . If  $h[i, j] = \text{add}_{XY}$ , we know that  $x_i(=y_j)$  is appended to the LCS sequence, and we continue with entry  $[i-1, j-1]$ . If  $h[i, j] = \text{skip}_X$  we know that  $x_i$  is not in the LCS sequence, and we continue with entry  $[i-1, j]$ . If  $h[i, j] = \text{skip}_Y$  we know that  $y_j$  is not in the LCS sequence, and we continue with entry  $[i, j-1]$ . Because the characters of the LCS are generated in reverse order, we *prepend* each one to a sequence, so that when we are done, the sequence is in proper order.

---

Extracting the LCS using the Hints

```
get-lcs-sequence() {
  LCS = new empty character sequence
  i = m; j = n                          // start at lower right
  while(i != 0 or j != 0)               // loop until upper left
    switch h[i,j]
      case addXY:                        // add x[i] (= y[j])
        prepend x[i] (or equivalently y[j]) to front of LCS
        i--; j--; break
      case skipX: i--; break             // skip x[i]
      case skipY: j--; break             // skip y[j]
  return LCS
}
```

---

## Lecture 12: Dynamic Programming: Chain Matrix Multiplication

**Chain matrix multiplication:** This problem involves the question of determining the optimal sequence for performing a series of operations. This general class of problem is important in compiler design for code optimization and in databases for query optimization. We will study the problem in a very restricted instance, where the dynamic programming issues are easiest to see.

Suppose that we wish to multiply a series of matrices

$$C = A_1 \cdot A_2 \cdots A_n$$

Matrix multiplication is an associative but not a commutative operation. This means that we are free to parenthesize the above multiplication however we like, but we are not free to rearrange the order of the matrices. Also recall that when two (nonsquare) matrices are being multiplied, there are restrictions on the dimensions. A  $p \times q$  matrix has  $p$  rows and  $q$  columns. You can multiply a  $p \times q$  matrix  $A$  times a  $q \times r$  matrix  $B$ , and the result will be a  $p \times r$  matrix  $C$  (see Fig. 43). The number of columns of  $A$  must equal the number of rows of  $B$ . In particular for  $1 \leq i \leq p$  and  $1 \leq j \leq r$ , we have

$$C[i, j] = \sum_{k=1}^q A[i, k] \cdot B[k, j].$$

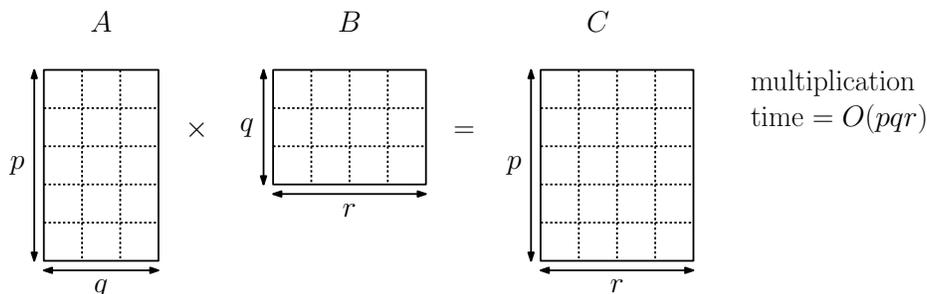


Fig. 43: Matrix Multiplication.

This corresponds to the (hopefully familiar) rule that the  $[i, j]$  entry of  $C$  is the dot product of the  $i$ th (horizontal) row of  $A$  and the  $j$ th (vertical) column of  $B$ . Observe that there are  $pr$  total entries in  $C$  and each takes  $O(q)$  time to compute, thus the total time to multiply these two matrices is proportional to the product of the dimensions,  $pqr$ .

Note that although any legal “parenthesization” will lead to a valid result, not all involve the same number of operations. Consider the case of 3 matrices:  $A_1$  be  $5 \times 4$ ,  $A_2$  be  $4 \times 6$  and  $A_3$  be  $6 \times 2$ .

$$\begin{aligned} \text{cost}[(A_1 A_2) A_3] &= (5 \cdot 4 \cdot 6) + (5 \cdot 6 \cdot 2) = 180, \\ \text{cost}[A_1 (A_2 A_3)] &= (4 \cdot 6 \cdot 2) + (5 \cdot 4 \cdot 2) = 88. \end{aligned}$$

Even for this small example, considerable savings can be achieved by reordering the evaluation sequence.

**Chain Matrix Multiplication Problem:** Given a sequence of matrices  $A_1, \dots, A_n$  and dimensions  $p_0, \dots, p_n$  where  $A_i$  is of dimension  $p_{i-1} \times p_i$ , determine the order of multiplication (represented, say, as a binary tree) that minimizes the number of operations.

**Important Note:** This algorithm *does not* perform the multiplications, it just determines the best order in which to perform the multiplications and the total number of operations.

**Dynamic programming approach:** A naive approach to this problem, namely that of trying all valid ways of parenthesizing the expression, will lead to an exponential running time. We will solve it through dynamic programming.

This problem, like other dynamic programming problems involves determining a structure (in this case, a parenthesization). We want to break the problem into subproblems, whose solutions can be combined to solve the global problem. As is common to any DP solution, we need to find some way to break the problem into smaller subproblems, and we need to determine a recursive formulation, which represents the optimum solution to each problem in terms of solutions to the subproblems. Let us think of how we can do this.

Since matrices cannot be reordered, it makes sense to think about sequences of matrices. Let  $A_{i..j}$  denote the result of multiplying matrices  $i$  through  $j$ . It is easy to see that  $A_{i..j}$  is a  $p_{i-1} \times p_j$  matrix. (Think about this for a second to be sure you see why.) Now, in order to determine how to perform this multiplication optimally, we need to make many decisions. What we want to do is to break the problem into problems of a similar structure. In parenthesizing the expression, we can consider the highest level of parenthesization. At this level we are simply multiplying two matrices together. That is,

$$A_{1..n} = A_{1..k} \cdot A_{k+1..n} \quad \text{for } 1 \leq k \leq n-1.$$

Thus the problem of determining the optimal sequence reduces to two decisions:

- What is the best place to split the chain? (what is  $k$ ?)
- How do we parenthesize each of the subsequences  $A_{1..k}$  and  $A_{k+1..n}$ ?

Clearly, the problems of computing the two subsequences can be solved recursively, by applying the same scheme. (This is an instance of the principle of optimality. There is no advantage to be gained by solving a subproblem suboptimally.) So, let us think about the problem of determining the best value of  $k$ . At this point, you may be tempted to consider some clever ideas. For example, since we want matrices with small dimensions, pick the value of  $k$  that minimizes  $p_k$ . Although this is not a bad idea, in principle. (After all it might work. It just turns out that it doesn't in this case. This takes a bit of thinking, which you should try.)

Instead, as is the case in the other dynamic programming solutions we have seen, we will try *all possible* choices of  $k$  and take the best of them. This is not as inefficient as it might sound, since there are only  $O(n^2)$  different sequences of matrices. (There are  $\binom{n}{2} = n(n-1)/2$  ways of choosing  $i$  and  $j$  to form  $A_{i..j}$  to be precise.) Thus, we do not encounter exponential growth in our algorithm's complexity, only polynomial growth.

Notice that our chain matrix multiplication problem satisfies the principle of optimality. In particular, once we decide to break the sequence into the product  $A_{1..k} \cdot A_{k+1..n}$ , it is in our best interest to compute each subsequence optimally. That is, for the global problem to be solved optimally, the subproblems should be solved optimally as well.

**Recursive formulation:** Let's explore how to express the optimum cost of multiplication in a recursive form. Later we will consider how to efficiently implement this recursive rule. We will subdivide the problem into subproblems by considering subsequences of matrices. In particular, for  $1 \leq i \leq j \leq n$ , let  $m(i, j)$  denote the minimum number of multiplications needed to compute  $A_{i..j}$ . The desired total cost of multiplying all the matrices is that of computing the entire chain  $A_{1..n}$ , which is given by  $m(1, n)$ . The optimum cost can be described by the following recursive formulation.

**Basis:** Observe that if  $i = j$  then the sequence contains only one matrix, and so the cost is 0. (There is nothing to multiply.) Thus,  $m(i, i) = 0$ .

**Step:** If  $i < j$ , then we are asking about the product  $A_{i..j}$ . This can be split into two groups  $A_{i..k}$  times  $A_{k+1..j}$ , by considering each  $k$ ,  $i \leq k < j$  (see Fig. 44).

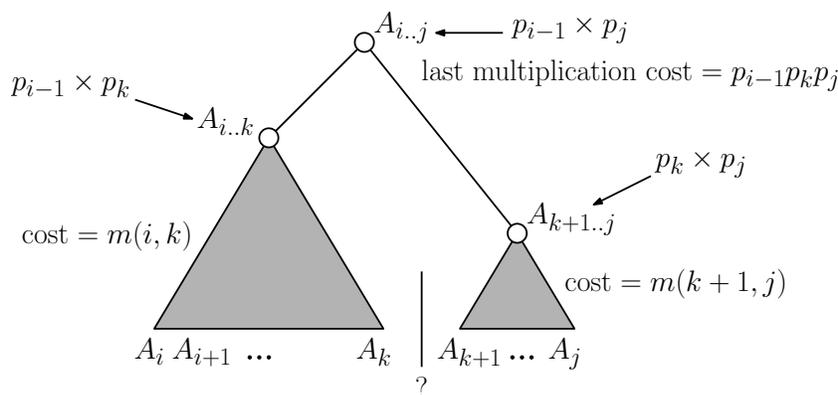


Fig. 44: Dynamic programming decision.

The optimum times to compute  $A_{i..k}$  and  $A_{k+1..j}$  are, by definition,  $m(i, k)$  and  $m(k+1, j)$ , respectively. Let us assume inductively that we can compute these values. (Note that they each involve a strictly smaller number of matrices, so there is no possibility of circularity.) Since  $A_{i..k}$  is a  $p_{i-1} \times p_k$  matrix, and  $A_{k+1..j}$  is a  $p_k \times p_j$  matrix, the time to multiply them is  $p_{i-1} p_k p_j$ . This suggests the following recursive rule for  $m(i, j)$ .

$$\begin{aligned}
 m(i, i) &= 0 \\
 m(i, j) &= \min_{i \leq k < j} (m(i, k) + m(k+1, j) + p_{i-1} p_k p_j) \quad \text{for } i < j.
 \end{aligned}$$

**Bottom-up implementation:** As with other DP problems, there are two natural implementations of the recursive rule that will lead to an efficient algorithm. One is memoization (which we will leave as an exercise), and the other is bottom-up calculation. We will consider just the latter.

To do this, we will store the values of  $m(i, j)$  in a 2-dimensional array  $m[1..n, 1..n]$ . The trickiest part of the process is arranging the order in which to compute the values. In the process of computing  $m(i, j)$  we need to access values  $m(i, k)$  and  $m(k+1, j)$  for  $k$  lying between  $i$  and  $j$ . Note that we cannot just compute the matrix in the simple row-by-row order that we used for the longest common subsequence problem. To see why, suppose that we are computing the values in row 3. When computing  $m[3, 5]$ , we would need to access both  $m[3, 4]$  and  $m[4, 5]$ , but  $m[4, 5]$  is in row 4, which has not yet been computed.

Instead, the trick is to compute *diagonal-by-diagonal* working out from the middle of the array. In particular, we organize our computation according to the number of matrices in the subsequence. For example,  $m[3, 5]$  represents a chain of  $5 - 3 + 1 = 3$  matrices, whereas  $m[3, 4]$  and  $m[4, 5]$  each represent chains of only two matrices. We first solve the problem for chains of length 1 (which is trivial), then chains of length 2, and so on, until we come to  $m[1, n]$ , which is the total chain of length  $n$ .

To do this, for  $1 \leq i \leq j \leq n$ , let  $L = j - i + 1$  denote the length of the subchain being multiplied. How shall we set up the loops to do this? The case  $L = 1$  is trivial, since there is only one matrix, and nothing needs to be multiplied, so we have  $m[i, i] = 0$ . Otherwise, our outer loop runs from  $L = 2, \dots, n$ . If a subchain of length  $L$  starts at position  $i$ , then  $j = i + L - 1$ . Since  $j \leq n$ , we have  $i + L - 1 \leq n$ , or in other words,  $i \leq n - L + 1$ . So our inner loop will be based on  $i$  running from 1 up to  $n - L + 1$ . The code is presented in the code block below. (Also, see Fig. 45 for an example.) We will explain below the purpose of the  $s$  array.

The array  $s[i, j]$  will be explained below. It will be used to extract the actual multiplication sequence. The running time of the procedure is  $O(n^3)$ . This is because we have three nested loops, and each can iterate at most  $n$  times. (A more careful analysis would show that the total number of iterations grows roughly as  $n^3/6$ .)

```

Matrix-Chain(p[0..n]) {
  s = array[1..n-1, 2..n]
  for (i = 1 to n) m[i, i] = 0           // initialize
  for (L = 2 to n) {                   // L = length of subchain
    for (i = 1 to n - L + 1) {
      j = i + L - 1
      m[i,j] = INFINITY
      for (k = i to j - 1) {           // check all splits
        cost = m[i, k] + m[k+1, j] + p[i-1]*p[k]*p[j]
        if (cost < m[i, j]) {         // found a new optimum?
          m[i, j] = cost              // ...save its cost
          s[i, j] = k                 // ...and the split marker
        }
      }
    }
  }
  return m[1, n] (final cost) and s (splitting markers)
}

```

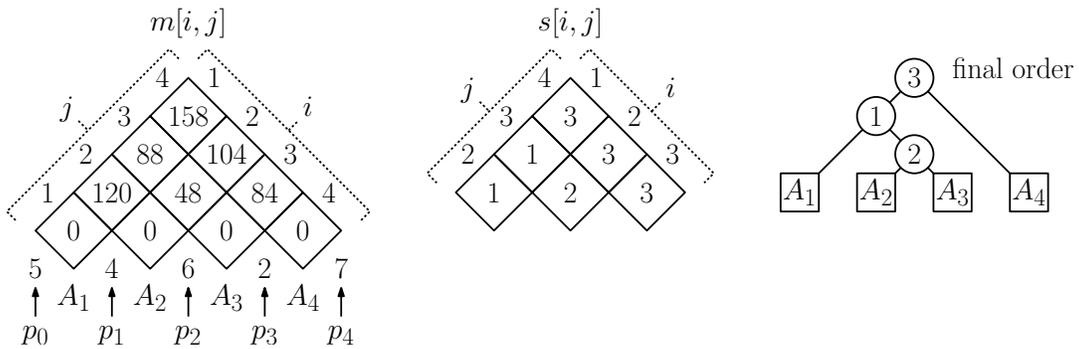


Fig. 45: Chain matrix multiplication for the product  $A_1 \cdots A_4$ , where  $A_i$  is of dimension  $p_{i-1} \times p_i$ .

**Extracting the final Sequence:** Extracting the actual multiplication sequence is a fairly easy extension.

The basic idea is to leave a *split marker* indicating what the best split is, that is, the value of  $k$  that leads to the minimum value of  $m[i, j]$ . We can maintain a parallel array  $s[i, j]$  in which we will store the value of  $k$  providing the optimal split. For example, suppose that  $s[i, j] = k$ . This tells us that the best way to multiply the subchain  $A_{i..j}$  is to first multiply the subchain  $A_{i..k}$  and then multiply the subchain  $A_{k+1..j}$ , and finally multiply these together. Intuitively,  $s[i, j]$  tells us what multiplication to perform *last*. Note that we only need to store  $s[i, j]$  when we have at least two matrices, that is, if  $j > i$ .

The actual multiplication algorithm uses the  $s[i, j]$  value to determine how to split the current sequence. Assume that the matrices are stored in an array of matrices  $A[1..n]$ , and that  $s[i, j]$  is global to this recursive procedure. The recursive procedure `do-mult` does this computation and below returns a matrix (see Fig. 45).

---

Extracting Optimum Sequence

```
do-mult(i, j) {
    if (i == j)                // basis case
        return A[i]
    else {
        k = s[i, j]
        X = do-mult(i, k)      // X = A[i]...A[k]
        Y = do-mult(k+1, j)   // Y = A[k+1]...A[j]
        return X * Y         // multiply matrices X and Y
    }
}
```

---

It's a good idea to trace through this example to be sure you understand it.

## Lecture 13: All-Pairs Shortest Paths and the Floyd-Warshall Algorithm

**All-Pairs Shortest Paths:** Earlier, we saw that Dijkstra's algorithm and the Bellman-Ford algorithm both solved the problem of computing shortest paths in graphs from a single source vertex. Suppose that we want instead to compute shortest paths between *all pairs* of vertices. We could do this applying either Dijkstra or Bellman-Ford using every vertex as a source, but today we will consider an alternative approach, which is based on dynamic programming.

Let  $G = (V, E)$  be a directed graph with edge weights. For each edge  $(u, v) \in E$ , let  $w(u, v)$  denote its weight. The *cost* of a path is the sum of its edge weights, and the *distance* between two vertices, denoted  $\delta(u, v)$ , minimum cost of any path between  $u$  and  $v$ . As in Bellman-Ford, we allow negative cost edges, but no negative-cost cycles. (Recall that negative-cost cycles imply that shortest paths are not defined.)

The algorithm that we will present is called the *Floyd-Warshall algorithm*. It runs in  $O(n^3)$  time, where  $n = |V|$ . It dates back to the early 1960's. The algorithm can be adapted for use in a number of related applications as well.

**Transitive Closure:** You are given a *binary relation*  $R$  on a set  $X$ , by which we mean that  $R$  is a subset of ordered pairs  $(x, y) \subseteq X \times X$ . A relation is said to be *transitive* if for any  $x, y, z \in X$ , if  $(x, y) \in R$  and  $(y, z) \in R$  then  $(x, z) \in R$ . The *transitive closure* of  $R$ , denoted  $R^*$  is the smallest extension of  $R$  that is transitive.

We can think of  $(X, R)$  as a directed graph, where  $X$  are the vertices and the pairs of  $R$  are edges. The transitive closure of  $R$  is effectively the same as the *reachability* relation in this graph, that

is,  $(x, y) \in R^*$  if and only if there exists a path from  $x$  to  $y$  in  $R$ . The Floyd-Warshall algorithm can be modified to compute the transitive closure in time  $O(n^3)$ , where  $n = |X|$ .

**All-Pairs Max-Capacity Paths:** Let  $G = (V, E)$  be a directed graph with positive edge capacities  $c(u, v)$ . Think of each edge as a pipe, and the capacity  $c(u, v)$  as the amount of flow that can be pushed through this pipe per unit interval. The *capacity* of a path is the minimum capacity of any edge along the path. (Intuitively, the minimum capacity edge forms a bottleneck, which limits the total amount of flow along the path. By the way, the capacity of the trivial path from  $u$  to  $u$  is  $+\infty$ .)

Given any two nodes  $u$  and  $v$ , we would like to know the maximum capacity path between them. It is easy to modify the Floyd-Warshall algorithm to compute the maximum capacity path between every pair of vertices in  $O(n^3)$  time, where  $n = |V|$ .

**Input/Output Representation:** We assume that the digraph is represented as an adjacency matrix, rather than the more common adjacency list. (Adjacency lists are generally more efficient for sparse graphs, but storing all the inter-vertex distances will require  $\Omega(n^2)$  storage anyway.) Because the algorithm is matrix-based, we will employ common matrix notation, using  $i, j$  and  $k$  to denote vertices rather than  $u, v$ , and  $w$  as we usually do.

The input is an  $n \times n$  matrix  $w$  of edge weights, which are based on the edge weights in the digraph. We let  $w_{ij}$  denote the entry in row  $i$  and column  $j$  of  $w$ .

$$w_{ij} = \begin{cases} 0 & \text{if } i = j, \\ w(i, j) & \text{if } i \neq j \text{ and } (i, j) \in E, \\ +\infty & \text{if } i \neq j \text{ and } (i, j) \notin E. \end{cases}$$

Setting  $w_{ij} = \infty$  if there is no edge, intuitively means that there is no direct link between these two nodes, and hence the direct cost is infinite. The reason for setting  $w_{ii} = 0$  is that there is always a trivial path of length 0 (using no edges) from any vertex to itself.

The output will be an  $n \times n$  distance matrix  $D = d_{ij}$  where  $d_{ij} = \delta(i, j)$ , the shortest path cost from vertex  $i$  to  $j$ . In order to recover the shortest path, we will also compute a helper matrix  $\text{helper}[i, j]$ . The value of  $\text{helper}[i, j]$  will be any vertex that is midway along the shortest path from  $i$  to  $j$ . If the shortest path travels directly from  $i$  to  $j$  without passing through any other vertices, then  $\text{helper}[i, j]$  will be set to *null*. (Later we will see how to use these values to compute the final path.)

**Floyd-Warshall Algorithm:** As with any DP algorithm, the key is reducing a large problem to smaller problems. A natural way of doing this is to follow the approach of Bellman-Ford of constructing shortest paths based on the number of edges in the path. The first iteration finds all shortest paths of length one, the next finds shortest paths of length two, and so on. It turns out that this does not lead to the fastest algorithm, however. (A more efficient variant finds shortest paths by doubling, first paths of length one, then two, then four, etc. However, this is still not the fastest.)

Rather than restricting the number of edges on the path, the trick is to restrict the set of vertices through which the path is allowed to pass. Given a path  $p = \langle v_1, v_2, \dots, v_\ell \rangle$  we refer to the vertices  $v_2, \dots, v_{\ell-1}$  as the *intermediate vertices* of this path. Note that a path consisting of a single edge has no intermediate vertices.

- Define  $d_{ij}^{(k)}$  to be the shortest path from  $i$  to  $j$  such that any intermediate vertices on the path are chosen from the set  $\{1, \dots, k\}$ .

In other words, we consider a path from  $i$  to  $j$  which either consists of the single edge  $(i, j)$ , or it visits some intermediate vertices along the way, but these intermediate can only be chosen from among  $\{1, \dots, k\}$ . (It does not need to visit all of these vertices, and it may visit none of them.) The path is free to visit any subset of these vertices, and to do so in any order. For example, in the digraph shown in the Fig. 46(a), notice how the value of  $d_{5,6}^{(k)}$  changes as  $k$  varies.

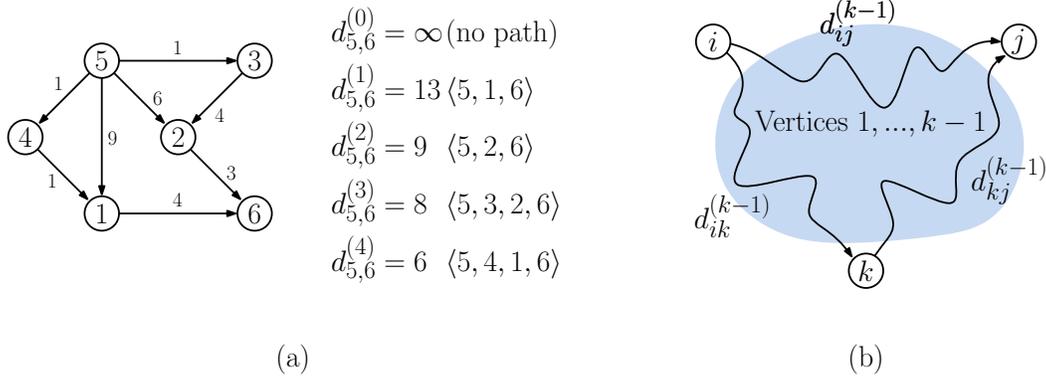


Fig. 46: Limiting intermediate vertices. For example  $d_{5,6}^{(3)}$  can go through any combination of the intermediate vertices  $\{1, 2, 3\}$ , of which  $\langle 5, 3, 2, 6 \rangle$  has the lowest cost of 8.

How do we compute  $d_{ij}^{(k)}$  assuming that we have already computed the previous matrix  $d^{(k-1)}$ ? For the basis case ( $k = 0$ ) the path cannot go through any intermediate vertices, and so  $d_{ij}^{(0)} = w(i, j)$  for all  $i, j$ . For the induction step ( $k \geq 1$ ), there are two cases, depending on the ways that we might get from vertex  $i$  to vertex  $j$ , assuming that the intermediate vertices are chosen from  $\{1, 2, \dots, k\}$ :

**Don't go through  $k$  at all:** The shortest path from node  $i$  to node  $j$  uses intermediate vertices  $\{1, \dots, k - 1\}$ , and hence the length of the shortest path is  $d_{ij}^{(k-1)}$ .

**Go through  $k$ :** First observe that a shortest path does not pass through the same vertex twice, so we can assume that we pass through  $k$  exactly once. (The assumption that there are no negative cost cycles is being used here.) That is, we go from  $i$  to  $k$ , and then from  $k$  to  $j$ . In order for the overall path to be as short as possible we should take the shortest path from  $i$  to  $k$ , and the shortest path from  $k$  to  $j$ . Since both of these paths use intermediate vertices only in  $\{1, \dots, k - 1\}$ , the length of the path is  $d_{ik}^{(k-1)} + d_{kj}^{(k-1)}$ .<sup>8</sup>

From the above discussion, we have the following recursive rule (the DP formulation) for computing  $d^{(k)}$ , which is illustrated in Fig. 46(b).

$$\begin{aligned}
 d_{ij}^{(0)} &= w_{ij}, \\
 d_{ij}^{(k)} &= \min \left( d_{ij}^{(k-1)}, d_{ik}^{(k-1)} + d_{kj}^{(k-1)} \right) \quad \text{for } k \geq 1.
 \end{aligned}$$

The final answer is  $d_{ij}^{(n)}$ , as this allows all possible vertices as intermediate vertices.

As with other DP problems, a recursive implementation of this rule would be prohibitively slow because the same values may be reevaluated many times. Instead, we store the computed values in a table. An honest implementation of the above rule would involve a 3-dimensional array,  $d[i, j, k]$ . However, a careful analysis of the algorithm reveals that the third dimension does not need to be explicitly stored. We will present the simpler version of the algorithm (which omits the  $k$ th dimension) and leave it as an exercise that the algorithm is still a correct implementation of the above rule.

The complete algorithm is presented in the code fragment below. We have also included helper values,  $\text{helper}[i, j]$ , that will be useful for extracting the final shortest paths. Recall that it stores any vertex

<sup>8</sup>Although the figure suggests that  $i \neq j \neq k$ , you should convince yourself that this holds even if some combination of  $i, j$ , and  $k$  are equal to each other. For example, if  $j = k$ , then  $d_{kj}^{(k-1)} = d_{jj}^{(k-1)} = w(j, j) = 0$ , and so the formula degenerates to  $d_{ik}^{(k-1)} + d_{kj}^{(k-1)} = d_{ik}^{(k-1)} = d_{ij}^{(k-1)}$ .

along the shortest path from  $i$  to  $j$ . If the path goes through  $k$ , then we can store  $k$  as the helper. An example of the algorithm's execution is shown in Fig. 47.

---

Floyd-Warshall Shortest-Path Algorithm

```
floyd-warshall(n, w) {
    d = array [1..n, 1..n]           // distance matrix
    for (i = 1 to n) {               // initialize
        for (j = 1 to n) {
            d[i, j] = w[i, j]
            helper[i, j] = null
        }
    }
    for (k = 1 to n) {               // use intermediates {1..k}
        for (i = 1 to n) {           // ...from i
            for (j = 1 to n) {       // ...to j
                if (d[i, k] + d[k, j]) < d[i, j]) {
                    d[i, j] = d[i, k] + d[k, j] // new shorter path length
                    helper[i, j] = k           // new path is through k
                }
            }
        }
    }
    return d                          // d[i,j] holds the distance from i to j
}
```

---

Clearly the algorithm's running time is  $O(n^3)$ . The space used by the algorithm is  $O(n^2)$ . Observe that we deleted all references to the superscript ( $k$ ) in the code. It is left as an exercise that this does not affect the correctness of the algorithm. (Hint: The danger is that values may be overwritten and then used later in the same phase. Consider which entries might be overwritten and then reused, they occur in row  $k$  and column  $k$ . It can be shown that the overwritten values are equal to their original values.)

**Extracting the Shortest Path:** Let's next see how to use the helper values  $\text{helper}[i, j]$  to extract the shortest path. Recall that whenever we discover that the shortest path from  $i$  to  $j$  passes through an intermediate vertex  $k$ , we set  $\text{helper}[i, j] = k$ . If the shortest path does not pass through any intermediate vertex, then  $\text{helper}[i, j] = \text{null}$ . To find the shortest path from  $i$  to  $j$ , we consult  $\text{helper}[i, j]$ . If it is  $\text{null}$ , then the shortest path is just the edge  $(i, j)$ . Otherwise, we recursively compute the shortest path from  $i$  to  $\text{helper}[i, j]$  and concatenate this with the shortest path from  $\text{helper}[i, j]$  to  $j$ .

## Lecture 14: Network Flows: Basic Definitions

**Network Flow:** "Network flow" is the name of a variety of related graph optimization problems, which are of fundamental value. We are given a *flow network*, which is essentially a directed graph with nonnegative edge weights. We think of the edges as "pipes" that are capable of carrying some sort of "stuff." In applications, this stuff can be any measurable quantity, such as fluid, megabytes of network traffic, commodities, currency, and so on. Each edge of the network has a given *capacity*, that limits the amount of stuff it is able to carry. The idea is to find out how much flow we can push from a designated source node to a designated sink node.

Although the network flow problem is defined in terms of the metaphor of pushing fluids, this problem and its many variations find remarkably diverse applications. These are often studied in the area of operations research. The network flow problem is also of interest because it is a restricted version of

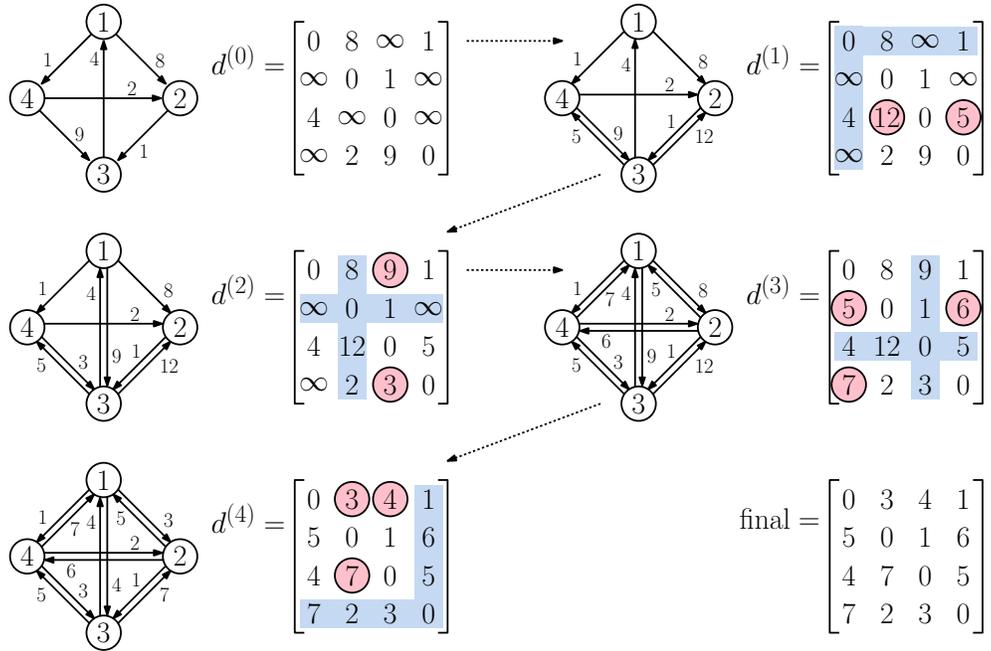


Fig. 47: Floyd-Warshall example. Newly updates entries are circled.

Printing the Shortest Path

```

get-path(i, j) {
    if (helper[i, j] == null)                // path is a single edge
        output(i, j)
    else {
        get-path(i, helper[i, j])           // path goes through helper
        get-path(helper[i, j], j)          // output the path from i to helper
    }
}

```

a more general optimization problem, called *linear programming*. A good understanding of network flows is helpful in obtaining a deeper understanding of linear programming.

**Flow Networks:** A *flow network* is a directed graph  $G = (V, E)$  in which each edge  $(u, v) \in E$  has a nonnegative *capacity*  $c(u, v) \geq 0$ . (In our book, the capacity of edge  $e$  is denoted by  $c_e$ .) If  $(u, v) \notin E$  we model this by setting  $c(u, v) = 0$ . There are two special vertices: a *source*  $s$ , and a *sink*  $t$  (see Fig. 48).

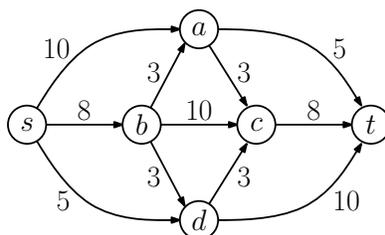


Fig. 48: A flow network.

We assume that there is no edge entering  $s$  and no edge leaving  $t$ . Such a network is sometimes called an *s-t network*. We also assume that every vertex lies on some path from the source to the sink.<sup>9</sup> This implies that  $m \geq n - 1$ , where  $n = |V|$  and  $m = |E|$ . It will also be convenient to assume that all capacities are integers. (We can assume more generally that the capacities are rational numbers, since we can convert them to integers by multiplying them by the least common multiple of the denominators.)

**Flows, Capacities, and Conservation:** Given an *s-t network*, a *flow* (also called an *s-t flow*) is a function  $f$  that maps each edge to a nonnegative real number and satisfies the following properties:

**Capacity Constraint:** For all  $(u, v) \in E$ ,  $f(u, v) \leq c(u, v)$ .

**Flow conservation (or flow balance):** For all  $v \in V \setminus \{s, t\}$ , the sum of flow along edges into  $v$  equals the sum of flows along edges out of  $v$ .

We can state flow conservation more formally as follows. First off, let us make the assumption that if  $(u, v)$  is *not* an edge of  $E$ , then  $f(u, v) = 0$ . We then define the total flow into  $v$  and total flow out of  $v$  as:

$$f^{\text{in}}(v) = \sum_{(u,v) \in E} f(u, v) \quad \text{and} \quad f^{\text{out}}(v) = \sum_{(v,w) \in E} f(v, w).$$

Then flow conservation states that  $f^{\text{in}}(v) = f^{\text{out}}(v)$ , for all  $v \in V \setminus \{s, t\}$ . Note that flow conservation *does not* apply to the source and sink, since we think of ourselves as pumping flow from  $s$  to  $t$ .

Two examples are shown in Fig. 49, where we use the notation  $f/c$  on each edge to denote the flow  $f$  and capacity  $c$  for this edge.

The quantity  $f(u, v)$  is called the *flow* along edge  $(u, v)$ . We are interested in defining the total flow, that is, the total amount of fluid flowing from  $s$  to  $t$ . The *value* of a flow  $f$ , denoted  $|f|$ , is defined as the sum of flows out of  $s$ , that is,

$$|f| = f^{\text{out}}(s) = \sum_{w \in V} f(s, w),$$

(For example, the value of the flow shown in Fig. 49(a) is  $5 + 8 + 5 = 18$ .) From flow conservation, it follows easily that this is also equal to the flow into  $t$ , that is,  $f^{\text{in}}(t)$ . We will prove this later.

<sup>9</sup>Neither of these is an essential requirement. Given a network that fails to satisfy these assumptions, we can easily generate an equivalent one that satisfies both.

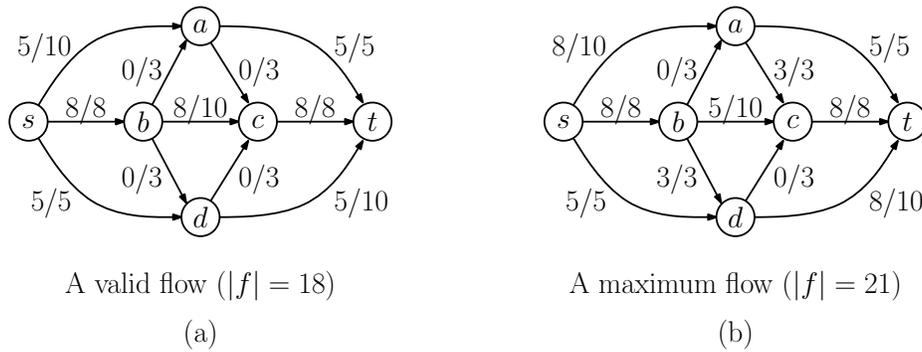


Fig. 49: A valid flow and a maximum flow.

**Maximum Flow:** Given an  $s$ - $t$  network, an obvious optimization problem is to determine a flow of maximum value. More formally, the *maximum-flow problem* is, given a flow network  $G = (V, E)$ , and source and sink vertices  $s$  and  $t$ , find the flow of maximum value from  $s$  to  $t$ . (For example, in Fig. 49(b) we show flow of value  $8 + 8 + 5 = 21$ , which can be shown to be the maximum flow for this network.) Note that, although the value of the maximum flow is unique, there may generally be many different flow functions that achieve this value.

**Path-Based Flows:** The definition of flow we gave above is sometimes called the *edge-based* definition of flows. An alternative, but mathematically equivalent, definition is called the *path-based* definition of flows. Define an  $s$ - $t$  path to be any simple path from  $s$  to  $t$ . For example, in Fig. 48,  $\langle s, a, t \rangle$ ,  $\langle s, b, a, c, t \rangle$  and  $\langle s, d, c, t \rangle$  are all examples of  $s$ - $t$  paths. There may generally be an exponential number of such paths (but that is alright, since this is just a mathematical definition).

A *path-based flow* is a function that assigns each  $s$ - $t$  path a nonnegative real number such that, for every edge  $(u, v) \in E$ , the sum of the flows on all the paths containing this edge is at most  $c(u, v)$ . Note that there is no need to provide a flow conservation constraint, because each path that carries a flow into a vertex (excluding  $s$  and  $t$ ), carries an equivalent amount of flow out of that vertex. For example, in Fig. 50(b) we show a path-based flow that is equivalent to the edge-based flow of Fig. 50(a). The paths carrying zero flow are not shown.

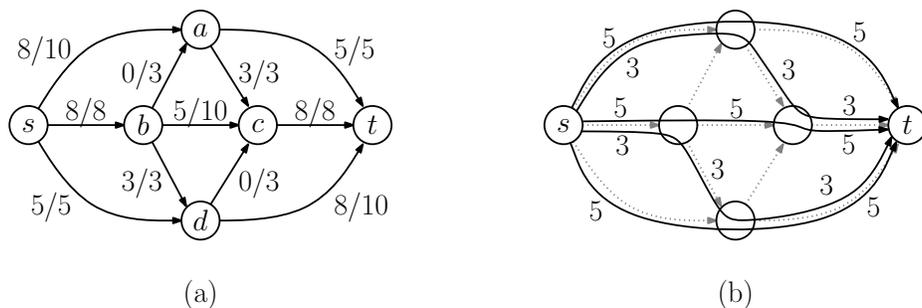


Fig. 50: (a) An edge-based flow and (b) its path-based equivalent.

The *value* of a path-based flow is defined to be the total sum of all the flows on all the  $s$ - $t$  paths of the network. Although we will not prove it, the following claim is an easy consequence of the above definitions.

**Claim:** Given an  $s$ - $t$  network  $G$ , under the assumption that there are no edges entering  $s$  or leaving  $t$ ,  $G$  has an edge-based flow of value  $x$  if and only if  $G$  has a path-based flow of value  $x$ .

**Multi-source, multi-sink networks:** It may seem overly restrictive to require that there is only a single source and a single sink vertex. Many flow problems have situations in which many source vertices  $s_1, \dots, s_k$  and many sink vertices  $t_1, \dots, t_l$ . This can easily be modeled by just adding a special *super-source*  $s'$  and a *super-sink*  $t'$ , and attaching  $s'$  to all the  $s_i$  and attach all the  $t_j$  to  $t'$ . We let these edges have infinite capacity (see Fig. 51). Now by pushing the maximum flow from  $s'$  to  $t'$  we are effectively producing the maximum flow from all the  $s_i$ 's to all the  $t_j$ 's.

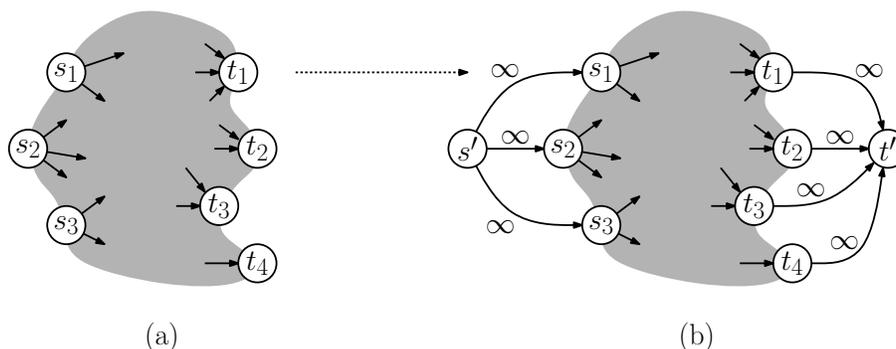


Fig. 51: Reduction from (a) multi-source/multi-sink to (b) single-source/single-sink.

Note that we don't assume any correspondence between flows leaving source  $s_i$  and entering  $t_j$ . Flows from one source may flow into *any* sink vertex. In some cases, you would like to specify the flow from a certain source must arrive at a designated sink vertex. For example, imagine that the sources are manufacturing production centers and sinks are retail outlets, and you are told the amount of commodity from  $s_i$  to arrive at  $t_j$ . This variant of the flow problem, called the *multi-commodity flow problem*, is a much harder problem to solve (in fact, some formulations are NP-hard).

## Lecture 15: Network Flows: The Ford-Fulkerson Algorithm

**Network Flow:** We continue discussion of the network flow problem. Last time, we introduced basic concepts, such as the concepts  $s$ - $t$  networks and flows. Today, we discuss the Ford-Fulkerson Max Flow algorithm, cuts, and the relationship between flows and cuts.

Recall that a *flow network* is a directed graph  $G = (V, E)$  in which each edge  $(u, v) \in E$  has a nonnegative *capacity*  $c(u, v) \geq 0$ , with a designated source vertex  $s$  and sink vertex  $t$ . We assume that there are no edges entering  $s$  or exiting  $t$ . A *flow* is a function  $f$  that maps each edge to a nonnegative real number that does not exceed the edge's capacity, and such that the total flow into any vertex other than  $s$  and  $t$  equals the total flow out of this vertex. The total *value* of a flow is equal to the sum of flows coming out of  $s$  (which, by flow conservation, is equal to the total flow entering  $t$ ). The objective of the *max flow problem* is to compute a flow of maximum value. Today we present an algorithm for this problem.

**Why Greedy Fails:** Before considering our algorithm, we start by considering why a simple greedy scheme for computing the maximum flow does not work. The idea behind the greedy algorithm is motivated by the path-based notion of flow. (Recall this from the previous lecture.) Initially the flow on each edge is set to zero. Next, find any path  $P$  from  $s$  to  $t$ , such that the edge capacities on this path are all strictly positive. Let  $c_{\min}$  be the minimum capacity of any edge on this path. This quantity is called the *bottleneck capacity* of the path. Push  $c_{\min}$  units through this path. For each edge  $(u, v) \in P$ , set  $f(u, v) \leftarrow c_{\min} + f(u, v)$ , and decrease the capacity of  $(u, v)$  by  $c_{\min}$ . Repeat this until no  $s$ - $t$  path (of positive capacity edges) remains in the network.

While this may seem to be a very reasonable algorithm, and will generally produce a valid flow, it may fail to compute the maximum flow. To see why, consider the network shown in Fig. 52(a). Suppose we push 5 units along the topmost path, 8 units along the middle path, and 5 units along the bottommost path. We have a flow of value 18. After adjusting the capacities (see Fig. 52(b)) we see that there is no path of positive capacity from  $s$  to  $t$ . Thus, greedy gets stuck.

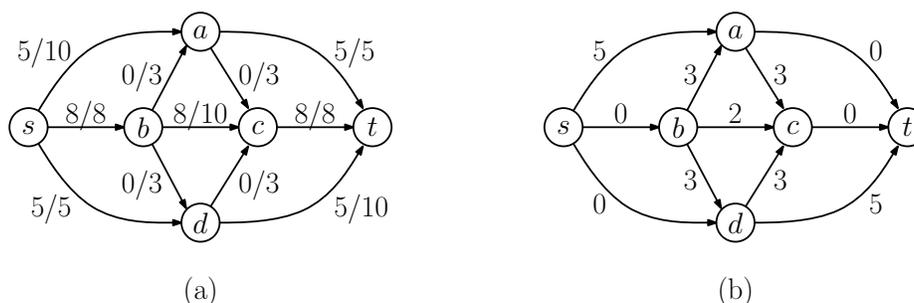


Fig. 52: The greedy flow algorithm can get stuck before finding the maximum flow.

**Residual Network:** The key insight to overcoming the problem with the greedy algorithm is to observe that, in addition to increasing flows on edges, it is possible to *decrease* flows on edges that already carry flow (as long as the flow never becomes negative). It may seem counterintuitive that this would help, but we shall see that it is exactly what is needed to obtain an optimal solution.

To make this idea clearer, we first need to define the notion of the residual network and augmenting paths. Given a flow network  $G$  and a flow  $f$ , define the *residual network*, denoted  $G_f$ , to be a network having the same vertex set and same source and sink, and whose edges are defined as follows:

**Forward edges:** For each edge  $(u, v)$  for which  $f(u, v) < c(u, v)$ , create an edge  $(u, v)$  in  $G_f$  and assign it the capacity  $c_f(u, v) = c(u, v) - f(u, v)$ . Intuitively, this edge signifies that we can *increase* the flow along this edge by up to  $c_f(u, v)$  units without violating the original capacity constraint.

**Backward edges:** For each edge  $(u, v)$  for which  $f(u, v) > 0$ , create an edge  $(v, u)$  in  $G_f$  and assign it a capacity of  $c_f(v, u) = f(u, v)$ . Intuitively, this edge signifies that we can *decrease* the existing flow by as much as  $c_f(v, u)$  units. Conceptually, by pushing positive flow along the reverse edge  $(v, u)$  we are decreasing the flow along the original edge  $(u, v)$ .

Observe that every edge of the residual network has *strictly positive* capacity. (This will be important later on.) Note that each edge in the original network may result in the generation of up to two new edges in the residual network. Thus, the residual network is of the same asymptotic size as the original network.

An example of a flow and the associated residual network are shown in Fig. 53(a) and (b), respectively. For example, the edge  $(b, c)$  of capacity 2 signifies that we can add up to 2 more units of flow to edge  $(b, c)$  and the edge  $(c, b)$  of capacity 8 signifies that we can cancel up to 8 units of flow from the edge  $(b, c)$ .

The capacity of each edge in the residual network is called its *residual capacity*. The key observation about the residual network is that if we can push flow through the residual network then we can push this additional amount of flow through the original network. This is formalized in the following lemma. Given two flows  $f$  and  $f'$ , we define their *sum*,  $f + f'$ , in the natural way, by summing the flows along each edge. If  $f'' = f + f'$ , then  $f''(u, v) = f(u, v) + f'(u, v)$ . Clearly, the value of  $f + f'$  is equal to  $|f| + |f'|$ .

**Lemma:** Let  $f$  be a flow in  $G$  and let  $f'$  be a flow in  $G_f$ . Then  $(f + f')$  is a flow in  $G$ .

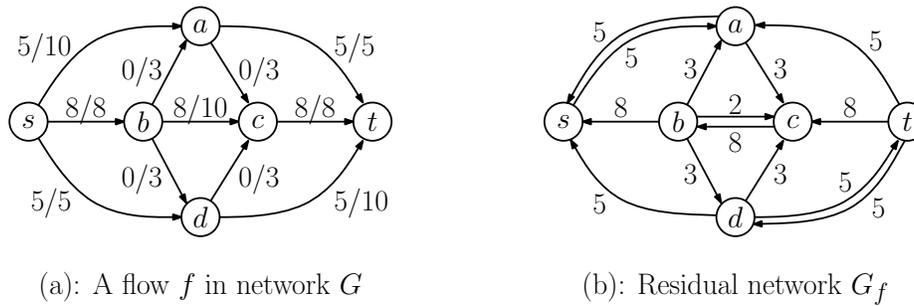


Fig. 53: A flow  $f$  and the residual network  $G_f$ .

**Proof:** (Sketch) To show that the resulting flow is valid, we need to show that it satisfies both the capacity constraints and flow conservation. It is easy to see that the capacities of  $G_f$  were exactly designed so that any flow along an edge of  $G_f$  when added to the flow  $f$  of  $G$  will satisfy  $G$ 's capacity constraints. Also, since both flows satisfy flow conservation, it is easy to see that their sum will as well. (More generally, any linear combination  $\alpha f + \beta f'$  will satisfy flow conservation.)

This lemma suggests that all we need to do to increase the flow is to find any flow in the residual network. This leads to the notion of an augmenting path.

**Augmenting Paths and Ford-Fulkerson:** Consider a network  $G$ , let  $f$  be a flow in  $G$ , and let  $G_f$  be the associated residual network. An *augmenting path* is a simple path  $P$  from  $s$  to  $t$  in  $G_f$ . The *residual capacity* (also called the *bottleneck capacity*) of the path is the minimum capacity of any edge on the path. It is denoted  $c_f(P)$ . (Recall that all the edges of  $G_f$  are of strictly positive capacity, so  $c_f(P) > 0$ .) By pushing  $c_f(P)$  units of flow along each edge of the path, we obtain a valid flow in  $G_f$ , and by the previous lemma, adding this to  $f$  results in a valid flow in  $G$  of strictly higher value.

For example, in Fig. 54(a) we show an augmenting path of capacity 3 in the residual network for the flow given earlier in Fig. 53. In (b), we show the result of adding this flow to every edge of the augmenting path. Observe that because of the backwards edge  $(c, b)$ , we have decreased the flow along edge  $(b, c)$  by 3, from 8 to 5.

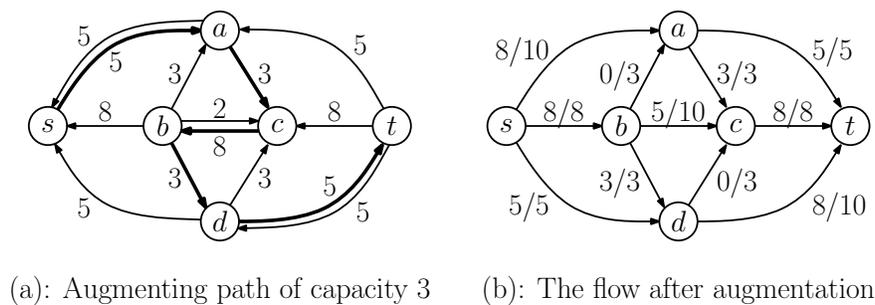


Fig. 54: Augmenting path and augmentation.

How is this different from the greedy algorithm? The greedy algorithm only increases flow on edges. Since an augmenting path may increase flow on a backwards edge, it may actually *decrease* the flow on some edge of the original network.

This observation naturally suggests an algorithm for computing flows of ever larger value. Start with a flow of weight 0, and then repeatedly find an augmenting path. Repeat this until no such path exists.

This, in a nutshell, is the simplest and best known algorithm for computing flows, called the *Ford-Fulkerson method*. (We do not call it an “algorithm,” since the method of selecting the augmenting path is not specified. We will discuss various strategies in future lectures.) It is summarized in the code fragment below.

---

Ford-Fulkerson Network Flow

```

ford-fulkerson-flow(G = (V, E, s, t)) {
  f = 0 (all edges carry zero flow)
  while (true) {
    G' = the residual-network of G for f
    if (G' has no s-t augmenting path)
      break // no augmentations left
    P = any-augmenting-path of G' // augmenting path
    c = minimum capacity edge of P // augmentation amount
    augment f by adding c to the flow on every edge of P
  }
  return f
}

```

---

There are three issues to consider before declaring this a reasonable algorithm.

- How efficiently can we perform augmentation?
- How many augmentations might be required until converging?
- If no more augmentations can be performed, have we found the max-flow?

Let us consider first the question of how to perform augmentation. First, given  $G$  and  $f$ , we need to compute the residual network,  $G_f$ . This is easy to do in  $O(n + m)$  time, where  $n = |V|$  and  $m = |E|$ . We assume that  $G_f$  contains only edges of strictly positive capacity. Next, we need to determine whether there exists an augmenting path from  $s$  to  $t$  in  $G_f$ . We can do this by performing either a DFS or BFS in the residual network starting at  $s$  and terminating as soon (if ever)  $t$  is reached. Let  $P$  be the resulting path. Clearly, this can be done in  $O(n + m)$  time as well. Finally, we compute the minimum cost edge along  $P$ , and increase the flow  $f$  by this amount for every edge of  $P$ .

Two questions remain: What is the best way to select the augmenting path, and is this correct in the sense of converging to the maximum flow? Next, we consider the issue of correctness. Before doing this, we will need to introduce the concept of a cut.

**Cuts:** In order to show that Ford-Fulkerson leads to the maximum flow, we need to formalize the notion of a “bottleneck” in the network. Intuitively, the flow cannot be increased forever, because there is some subset of edges, whose capacities eventually become saturated with flow. Every path from  $s$  to  $t$  must cross one of these saturated edges, and so the sum of capacities of these edges imposes an upper bound on size of the maximum flow. Thus, these edges form a bottleneck.

We want to make this concept mathematically formal. Since such a set of edges lie on every path from  $s$  to  $t$ , their removal defines a partition separating the vertices that  $s$  can reach from the vertices that  $s$  cannot reach. This suggests the following concept.

Given a network  $G$ , define a *cut* (also called an *s-t cut*) to be a partition of the vertex set into two disjoint subsets  $X \subseteq V$  and  $Y = V \setminus X$ , where  $s \in X$  and  $t \in Y$ . We define the *net flow* from  $X$  to  $Y$  to be the sum of flows from  $X$  to  $Y$  minus the sum of flows from  $Y$  to  $X$ , that is,

$$f(X, Y) = \sum_{x \in X} \sum_{y \in Y} f(x, y) - \sum_{y \in Y} \sum_{x \in X} f(y, x).$$

Observe that  $f(X, Y) = -f(Y, X)$ .

For example, Fig. 55 shows a flow of value 17. It also shows a cut  $(X, Y) = (\{s, a\}, \{b, c, d, t\})$ , where  $f(X, Y) = 17$ .

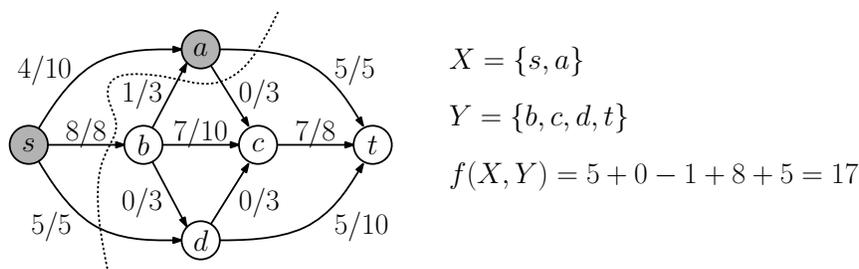


Fig. 55: Flow of value 17 across the cut  $(\{s, a\}, \{b, c, d, t\})$ .

**Lemma:** Let  $(X, Y)$  be any  $s$ - $t$  cut in a network. Given any flow  $f$ , the value of  $f$  is equal to the net flow across the cut, that is,  $f(X, Y) = |f|$ .

**Proof:** Recall that there are no edges leading into  $s$ , and so we have  $|f| = f^{\text{out}}(s) = f^{\text{out}}(s) - f^{\text{in}}(s)$ . Since all the other nodes of  $X$  must satisfy flow conservation it follows that

$$|f| = \sum_{x \in X} (f^{\text{out}}(x) - f^{\text{in}}(x))$$

Now, observe that every edge  $(u, v)$  where both  $u$  and  $v$  are in  $X$  contributes one positive term and one negative term of value  $f(u, v)$  to the above sum, and so all of these cancel out. The only terms that remain are the edges that either go from  $X$  to  $Y$  (which contribute positively) and those from  $Y$  to  $X$  (which contribute negatively). Thus, it follows that the value of the sum is exactly  $f(X, Y)$ , and therefore  $|f| = f(X, Y)$ .

Define the *capacity* of the cut  $(X, Y)$  to be the sum of the capacities of the edges leading from  $X$  to  $Y$ , that is,

$$c(X, Y) = \sum_{x \in X} \sum_{y \in Y} c(x, y).$$

There is a noteworthy asymmetry between the net flow of a cut and the capacity of a cut. In the net flow, we consider both the flows from  $X$  to  $Y$  and (negated) from  $Y$  to  $X$ . In contrast, in the definition of cut capacity, we only consider capacities from  $X$  to  $Y$ . To understand why this asymmetry exists, suppose that there are two edges  $(x, y)$  and  $(y, x)$ , where  $x \in X$  and  $y \in Y$ , where the edges have the same capacities and both carry the same flow. The flow from  $y$  to  $x$  effectively cancels out the flow from  $x$  to  $y$  (as if you have an “eddy” in the middle of a river, where the flow cycles around). Thus, this cycle does not contribute to the total flow value, and so it is important to consider both in the net flow. On the other hand, the capacity of the edge from  $y$  to  $x$  does not contribute to nor detract from the maximum amount of flow we can push from the  $X$ -side to the  $Y$ -side of the cut. Therefore, we do not need to consider it in computing the cut’s total capacity.

It is easy to see that it is not possible to push more flow through a cut than its capacity. Combining this with the above lemma we have:

**Lemma:** Given any  $s$ - $t$  cut  $(X, Y)$  and any flow  $f$  we have  $|f| \leq c(X, Y)$ .

The optimality of the Ford-Fulkerson method is based on the following famous theorem, called the *Max-Flow/Min-Cut Theorem*. Basically, it states that in any flow network the minimum capacity cut acts like a bottleneck to limit the maximum amount of flow. The Ford-Fulkerson method terminates when it finds this bottleneck, and hence on termination, it finds both the minimum cut and the maximum flow.

**Max-Flow/Min-Cut Theorem:** The following three conditions are equivalent.

- (i)  $f$  is a maximum flow in  $G$ ,
- (ii) The residual network  $G_f$  contains no augmenting paths,
- (iii)  $|f| = c(X, Y)$  for some cut  $(X, Y)$  of  $G$ .

**Proof:**

- (i)  $\Rightarrow$  (ii): (by contradiction) If  $f$  is a max flow and there were an augmenting path in  $G_f$ , then by pushing flow along this path we would have a larger flow, a contradiction.
- (ii)  $\Rightarrow$  (iii): If there are no augmenting paths then  $s$  and  $t$  are not connected in the residual network. Let  $X$  be those vertices reachable from  $s$  in the residual network, and let  $Y$  be the rest. Clearly,  $(X, Y)$  forms a cut. Clearly, each edge crossing the cut from  $X$  to  $Y$  must be saturated (or else we could reach one more vertex on the  $Y$ -side by adding flow along this edge), and each edge crossing the cut from  $Y$  to  $X$  must be carrying zero flow (or else we could reach one more vertex on the  $Y$ -side by reducing flow along this edge). It follows that the flow across the cut equals the capacity of the cut, thus  $|f| = c(X, Y)$ .
- (iii)  $\Rightarrow$  (i): This follows directly from the previous lemma. No flow value can exceed any cut capacity, and so if *any* flow's value matches *any* cut's capacity, this flow must be maximum.

We have established that, on termination, Ford-Fulkerson generates the maximum flow. But, is it guaranteed to terminate and, if so, how long does it take? We will consider this question next time.

## Lecture 16: Network Flow Algorithms

**Algorithmic Aspects of Network Flow:** In the previous lecture, we presented the Ford-Fulkerson algorithm. We showed that on termination this algorithm produces the maximum flow in an  $s$ - $t$  network. In this lecture we discuss the algorithm's running time, and discuss more efficient alternatives.

**Analysis of Ford-Fulkerson:** Before discussing the worst-case running time of the Ford-Fulkerson algorithm, let us first consider whether it is guaranteed to terminate. We assume that all edge capacities are integers.<sup>10</sup> Every augmentation by Ford-Fulkerson increases the flow by an integer amount. Thus, the resulting residual network also has integer capacities. Therefore, after a finite number of augmentations the algorithm must terminate.

**Lemma:** Given an  $s$ - $t$  network with integer capacities, the Ford-Fulkerson algorithm terminates. Furthermore, it produces an integer-valued flow function.

Recall our convention that  $n = |V|$  and  $m = |E|$ . Since we assume that every vertex is reachable from  $s$ , it follows that  $m \geq n - 1$ . Therefore,  $n = O(m)$ . Running times of the form  $O(n + m)$  can be expressed more simply as  $O(m)$ .

As we saw last time, the residual network can be computed in  $O(n + m) = O(m)$  time and an augmenting path can also be found in  $O(m)$  time. Therefore, the running time of each augmentation step is  $O(m)$ . How many augmentations are needed? Unfortunately, the number could be very large. To see this, consider the example shown in Fig. 56.

---

<sup>10</sup>This is not an unreasonable assumption. First, observe that if we multiply all the capacities by some positive constant  $c$ , all the flow values can be increased by this same factor. Assuming that the capacities are represented as fixed-point decimals with  $k$  digits to the right of decimal point, we can scale each capacity by  $c = 10^k$ , which converts all the capacities to integers. We solve this integer-capacity problem, and then divide the final result by  $c$ , thus mapping the solution back to the original instance.

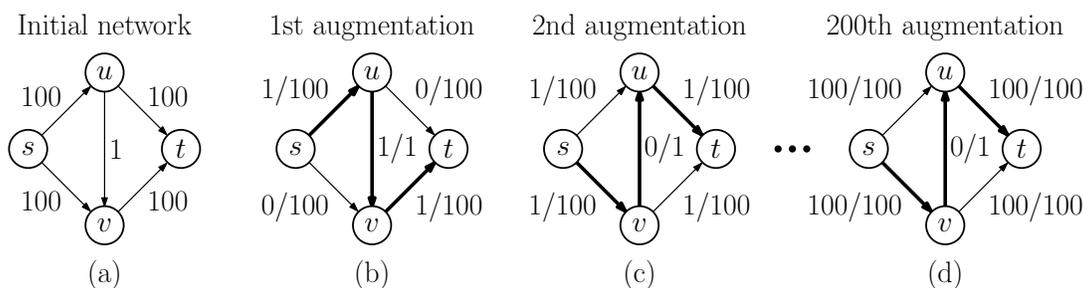


Fig. 56: Bad example for Ford-Fulkerson.

Suppose that we (foolishly) elect to augment using a path from  $s$  to  $t$  that uses the vertical edge in the middle, alternately increasing its flow to 1 and reducing it to 0. It would take 200 augmentation steps to converge. We could replace 100 with whatever huge value we want and make the running time arbitrarily high.

If we let  $|f|$  denote the final maximum flow value, the number of augmentation steps can be as high as  $|f|$ . If we make the reasonable assumption that each augmentation step takes at least  $\Omega(m)$  time, the total running time can be as high as  $\Omega(m|f|)$ . Since  $|f|$  may be arbitrarily high (it depends neither on  $n$  or  $m$ ), this running time could be arbitrarily high, as a function of  $n$  and  $m$ .

**Scaling Algorithm:** In the above example, we made the apparently foolish decision to augment on a path of very *low* capacity. Can we do better by selecting paths of high capacity? We will see an example of such an algorithm, called the *scaling algorithm*. For the sake of efficiency, the algorithm does not augment along the path of highest possible capacity, merely a path of relatively high capacity.

The idea is to start with an upper bound on the maximum possible flow. The sum of capacities of the edges leaving  $s$  suffices, that is,  $C = \sum_{(s,v) \in E} c(s,v)$ . Clearly, the maximum flow value cannot exceed  $C$ . We initialize the scaling parameter  $\Delta$  to be the largest power of 2, such that  $\Delta \leq C$ . Given any flow  $f$  (initially the flow of value 0), define  $G_f(\Delta)$  to be the residual network consisting *only of edges of residual capacity at least  $\Delta$* . Since we only deal with integer capacities, when  $\Delta = 1$ ,  $G_f(\Delta)$  is the true residual network. Intuitively, whenever we find an augmenting path in  $G_f(\Delta)$ , we are guaranteed to push at least  $\Delta$  units of flow, which means that we are guaranteed to make good progress. Next, find an  $s$ - $t$  path in  $G_f(\Delta)$ , augment the flow along this path, and update  $G_f(\Delta)$  accordingly. We repeat the process until no augmenting paths remain. We then set  $\Delta \leftarrow \Delta/2$  and repeat. When the value of  $\Delta$  falls below 1, we terminate the algorithm and return the final flow. See the code block below and Fig. 57.

**Analysis of the Scaling Algorithm:** We refer to the Kleinberg and Tardos book for a complete analysis of the scaling algorithm. Intuitively, the algorithm is efficient because each augmentation increases the flow by an amount of at least  $\Delta$ . The minimum cut can have at most  $m$  edges. So, after  $O(m)$  such augmentations, we will have effectively increased the flow along every edge of the minimum cut so much that its capacity in the residual graph falls below  $\Delta$ . When all the edges of the minimum cut disappear from  $G_f(\Delta)$ , it is not possible to augment further, and the algorithm goes on to the next smaller value of  $\Delta$ .

Since each augmentation involves running BFS (or DFS) in  $O(n + m) = O(m)$  time, it follows that after  $O(m^2)$  time, we will exhaust augmentations in  $G_f(\Delta)$ , and will proceed to the next smaller value of  $\Delta$ . After  $O(\log C)$  halvings of  $\Delta$ , we will have  $\Delta < 1$ , and the algorithm will terminate. Thus, the overall running time is  $O(m^2 \log C)$ . (It can be shown that a more efficient implementation runs in time  $O(nm \log C)$ .)

**Pseudo-Polynomial and Strong Polynomial Time:** Earlier, we complained that the Ford-Fulkerson

---

```

scaling-flow(G = (V, E, s, t)) {
    f = 0 // all edges carry zero flow
    D = largest power of 2 not larger than sum of capacities out of s
    while (D >= 1) { // when D = 1, Gf(D) has all edges
        Gf(D) = residual of G w.r.t. f, keeping only edges of capacity >= D
        while (there is an s-t path in Gf(D)) {
            P = any augmenting s-t path in Gf(D) // augment along the "fat" edges
            f' = augmenting flow for P
            f = f + f'
        }
        D = D/2 // shrink D's value
    }
    return f // return the final flow
}

```

---

algorithm is not really efficient, since its running time depends on the maximum flow value. The scaling algorithm also depends on the maximum flow value (albeit logarithmically rather than linearly). So, in what sense can we claim it is “efficient”?

Observe first that if the capacities are all small integers, then we can ignore the  $\log C$  component, and so the running time is just  $O(m^2)$  which is polynomial in the input size. (Efficient algorithms generally run in polynomial time, as opposed to exponential time.)

The algorithm is really *inefficient* when the capacities are *huge* numbers. Observe that if  $C$  is extremely large, then we require many bits to represent these numbers. Indeed, it takes  $\Theta(\log C)$  bits to represent a number of magnitude  $C$ . Thus, if we think of the input size from the perspective of *total number of bits* needed to represent the input graph, it can be shown that the scaling algorithm runs in time that is polynomial in this number of bits. (We will leave the details as an exercise.)

An algorithm whose running time is polynomial in the number of *bits* of input is called a *pseudopolynomial time algorithm*. In contrast, an algorithm that runs in time that is polynomial in the number of *words* of input (such as  $n$  and  $m$ ), is referred as running in *strongly polynomial time*.

**Edmonds-Karp Algorithm:** As mentioned above, neither of the algorithms we have seen runs in strongly polynomial time, that is, polynomial in  $n$  and  $m$ , irrespective of the magnitudes of the capacity. Edmonds and Karp developed such an algorithm in the 1970's (and it is claimed Dinic actually discovered this algorithm independently a couple years earlier).

This algorithm uses Ford-Fulkerson as its basis, but with the minor change that the  $s$ - $t$  path that is chosen in the residual network has the *smallest number of edges*. In particular, this just means that we run BFS in the residual graph from  $s$  to  $t$  to compute the augmenting path. It can be shown that the total number of augmenting steps using this method is  $O(nm)$ .<sup>11</sup> Since each augmentation takes  $O(m)$  time to run BFS, the overall running time is  $O(nm^2)$ .

**Even Faster Algorithms:** The max-flow problem is widely studied, and there are many different algorithms. No one knows that what the lowest possible running time is for network flow, but a running time of  $O(nm)$  has stood as an important milestone. See Table 1 for a summary of important results.

After a long sequence of improvements, in 2013 it was shown that the problem can be solved in this time. The final algorithm is not very elegant. It is a hybrid of two different algorithms, one by King,

---

<sup>11</sup>This is not trivial to prove. Neither of our textbooks gives a proof, but one can be found in the algorithms book by Cormen, Leiserson, Rivest, and Stein. Intuitively, it can be shown that after at most  $m$  augmentations, some vertex's distance from  $s$  increases by one, never to decrease again. Since a vertex's distance from  $s$  cannot exceed  $n$ , it follows that there are at most  $O(nm)$  augmentations.

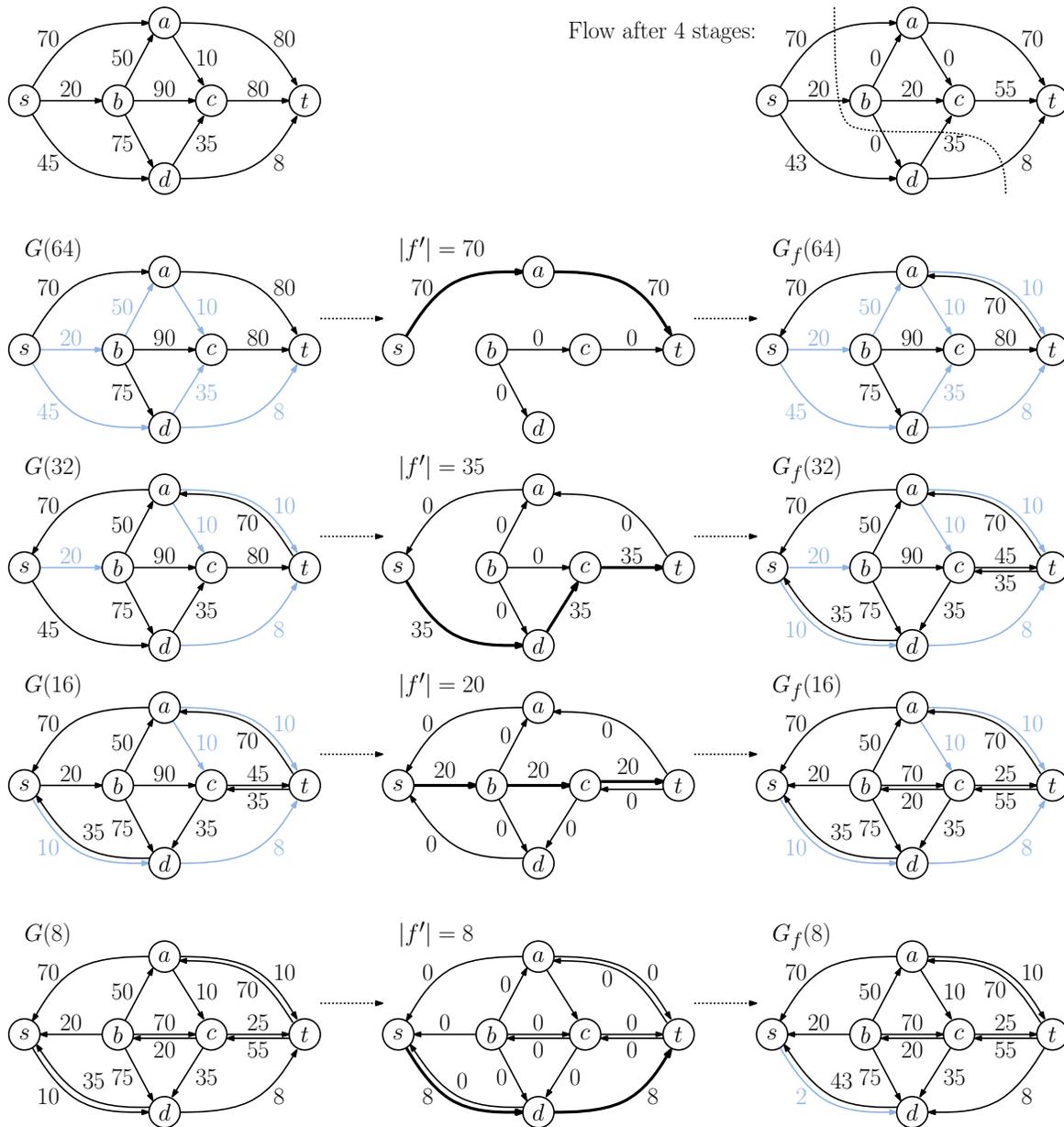


Fig. 57: The first four stages of the scaling algorithm for network flow. Note that at each stage, the value of the augmenting flow  $f'$  is at least  $\Delta$ . (The algorithm will run for  $G(4)$ ,  $G(2)$ , and  $G(1)$ , but no changes to the flow will occur.)

Rao, and Tarjan (KRT) that runs in  $O(nm)$  time for dense graphs and another by Orlin that runs in  $O(nm)$  time for sparse graphs. Whether there exist a unified algorithm that achieves this bound or even faster algorithms remain as open research questions.

Table 1: Running times of various network-flow algorithms ( $n = |V|$ ,  $m = |E|$ ,  $C$  is any upper bound on the maximum flow).

Algorithm	Year	Time	Notes
Ford-Fulkerson	1956	$O(mC)$	
Gabow	1985	$O(nm \log C)$	Scaling
Edmonds-Karp	1972	$O(nm^2)$	Ford-Fulkerson + augment shortest paths
Dinic	1970	$O(n^2m)$	Blocking flows in a layered graph
Dinic + Tarjan	1983	$O(nm \log n)$	Dinic + better data structures
Preflow push	1986	$O(nm \log(n^2/m))$	Goldberg and Tarjan
King, Rao, Tarjan	1994	$O(mn \log \frac{m}{n \log n} n)$	$O(nm)$ if $m = O(n^{1+\epsilon})$
Orlin + KRT	2013	$O(nm)$	Orlin: $O(nm)$ time for $m \leq O(n^{16/15-\epsilon})$ KRT: $O(nm)$ for $m > n^{1+\epsilon}$

**Applications of Max-Flow:** The network flow problem has a huge number of applications. Many of these applications do not appear at first to have anything to do with networks or flows. This is a testament to the power of this problem. In this lecture and the next, we will present a few applications from our book. (If you need more convincing of this, however, see the exercises in Chapter 7 of KL. There are over 40 problems, most of which involve reductions to network flow.)

**Maximum Matching in Bipartite Graphs:** There are many applications where it is desirable to compute a pairing between two sets of objects. We present such a problem, called *bipartite matching* in the form of a “dating game,” but the algorithm can be applied whenever it is desired to find pairing between objects of different classes subject to some compatibility criterion.

Suppose you are running a dating service, and there are a set of men  $X$  and a set of women  $Y$ . Using a questionnaire you establish which men are compatible which women. Your task is to pair up as many compatible pairs of men and women as possible, subject to the constraint that each man is paired with at most one woman, and vice versa. (It may be that some men are not paired with any woman.)

Recall that an undirected graph  $G = (V, E)$  is said to be *bipartite* if  $V$  can be partitioned into two sets  $X$  and  $Y$ , such that every edge has one endpoint in  $X$  and the other in  $Y$ . This problem can be modeled as an undirected, bipartite graph whose vertex set is  $V = X \cup Y$  and whose edge set consists of pairs  $(u, v)$ ,  $u \in X$ ,  $v \in Y$  such that  $u$  and  $v$  are compatible (see Fig. 58(a)). Given a graph, a *matching* is defined to be a subset of edges  $M \subseteq E$  such that for each  $v \in V$ , there is at most one edge of  $M$  incident to  $v$ . Clearly, the objective to the dating problem is to find a maximum matching in  $G$  that has the highest cardinality. Such a matching is called a *maximum matching* (see Fig. 58(b)).

The resulting undirected graph has the property that its vertex set can be divided into two groups such that all its edges go from one group to the other. This problem is called the *maximum bipartite matching problem*.

We will now show a reduction from maximum bipartite matching to network flow. In particular, we will show that, given any bipartite graph  $G$  (see Fig. 59(a)) for which we want to solve the maximum matching problem, we can convert it into an instance of network flow  $G'$ , such that the maximum matching on  $G$  can be extracted from the maximum flow on  $G'$ .

To do this, we construct a flow network  $G' = (V', E')$  as follows. Let  $s$  and  $t$  be two new vertices and

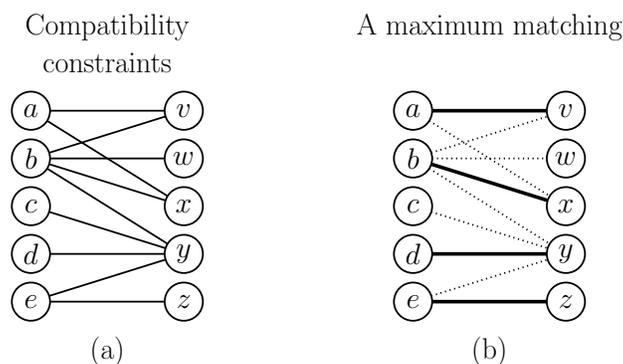


Fig. 58: A bipartite graph  $G$  and a maximum matching in  $G$ .

let  $V' = V \cup \{s, t\}$ .

$$E' = \begin{cases} \{(s, u) \mid u \in X\} \cup & \text{(connect source to left-side vertices)} \\ \{(v, t) \mid v \in Y\} \cup & \text{(connect right-side vertices to sink)} \\ \{(u, v) \mid (u, v) \in E\} & \text{(direct } G\text{'s edges from left to right).} \end{cases}$$

Set the capacity of all edges in this network to 1 (see Fig. 59(b)).

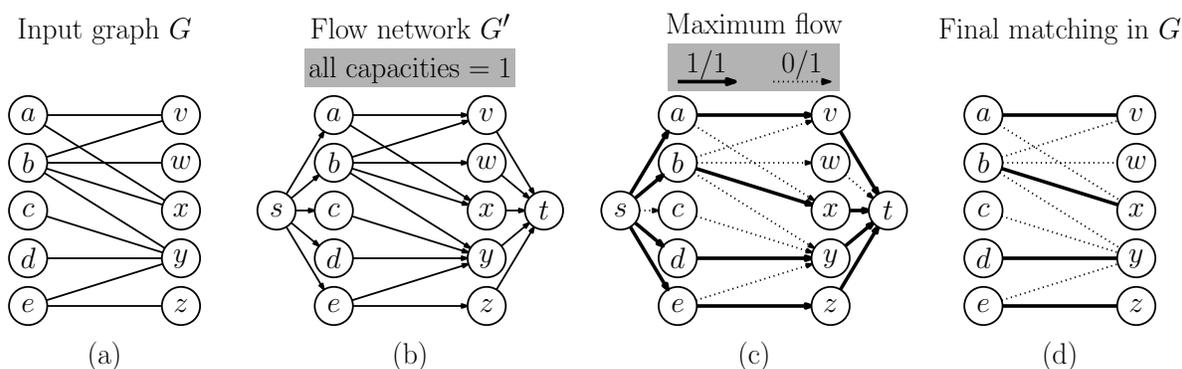


Fig. 59: Reducing bipartite matching to network flow.

Compute the maximum flow in  $G'$  (see Fig. 59(c)). The following lemma show that maximizing the flow in  $G'$  is equivalent to finding a maximum matching in  $G$ .

**Lemma:** Given a bipartite graph  $G$ ,  $G$  has a matching of size  $x$  if and only if  $G'$  (constructed above) has a flow of value  $x$ .

**Proof:** ( $\Rightarrow$ ) Let  $M$  denote any matching in  $G$ . We construct a flow in  $G'$  as follows. For each edge  $(x, y) \in M$ , set the flow along edges  $(s, x)$ ,  $(x, y)$ , and  $(y, t)$  to 1. All the edges remaining edges of  $G$  are assigned a flow of 0. We assert that  $f$  is a valid flow for  $G'$ . By our construction, each vertex  $x$  receives one unit of flow coming in from  $s$ , sends one unit to  $y$ , and  $y$  sends one unit to  $t$ . Therefore, we have flow conservation. Second, since  $M$  is a matching, each vertex of  $X$  or  $Y$  is incident to a single edge of  $M$ , which implies that it carries at most one unit of flow, which implies that the capacity constraints are all satisfied. Therefore,  $f$  is a valid flow in  $G'$ . By our construction,  $|f| = |M|$ .

( $\Leftarrow$ ) Suppose that  $G'$  has a flow  $f$ . We may assume that this is an integer flow. Since all edges have capacity 1, it follows that the flow value on each edge is either 0 or 1.

Let  $M$  denote the edges of  $X \times Y$  that are carrying unit flow in  $f$ . Observe that for every vertex of  $X$ , it has exactly one incoming edge (from  $s$ ) of capacity 1, and hence it can be incident to at most one edge of  $M$ . Symmetrically, every vertex of  $Y$  has exactly one outgoing edge (to  $t$ ) of capacity 1, and hence it also can be incident to at most one edge of  $M$ . Therefore,  $M$  is a matching in the original graph  $G$ . (An example is shown in Fig. 59(d)). Since each edge carries one unit of flow, the total value of the flow is the number of edges of  $M$ , that is,  $|f| = |M|$ .

Because the capacities are so low, we do not need to use a fancy implementation. Recall that Ford-Fulkerson runs in time  $O(m \cdot F^*)$ , where  $F^*$  is the final maximum flow. In our case  $F^*$  is at most  $n$  (the size of the largest possible matching). Therefore, the running time of Ford-Fulkerson on this instance is  $O(m \cdot F^*) = O(nm)$ .

There are other algorithms for maximum bipartite matching. The best is due to Hopcroft and Karp, and runs in  $O(\sqrt{n} \cdot m)$  time.

## Lecture 17: Network Flow: Extensions

**Extensions of Network Flow:** Network flow is an important problem because it is useful in a wide variety of applications. We will discuss two useful extensions to the network flow problem. We will show that these problems can be reduced to network flow, and thus a single algorithm can be used to solve both of them. Many computational problems that would seem to have little to do with flow of fluids through networks can be expressed as one of these two extended versions.

**Circulation with Demands:** There are many problems that are similar to network flow in which, rather than transporting flow from a single source to a single sink, we have a collection of *supply nodes* that want to ship flow (or products or goods) and a collection of *demand nodes* that want to receive flow. Each supply node is associated with the amount of product it wishes to ship and each demand node is associated with the amount that it wishes to receive. The question that arises is whether there is some way to get the products from the supply nodes to the demand nodes, subject to the capacity constraints. This is a *decision problem* (or *feasibility problem*), meaning that it has a yes-no answer, as opposed to maximum flow, which is an *optimization problem*.

We can model both supply and demand nodes elegantly by associating a single numeric value with each node, called its *demand*. If  $v \in V$  is a demand node, let  $d_v$  the amount of this demand. If  $v$  is a supply node, we model this by assigning it a negative demand, so that  $-d_v$  is its available supply. Intuitively, supplying  $x$  units of product is equivalent to demanding receipt of  $-x$  units.<sup>12</sup> If  $v$  is neither a supply or demand node, we let  $d_v = 0$ .

Suppose that we are given a directed graph  $G = (V, E)$  in which each edge  $(u, v)$  is associated with a positive capacity  $c(u, v)$  and each vertex  $v$  is associated with a supply/demand value  $d_v$ . Let  $S$  denote the set of *supply nodes* ( $d_v < 0$ ), and let  $T$  denote the set of *demand nodes* ( $d_v > 0$ ). Note that vertices of  $S$  may have incoming edges and vertices of  $T$  may have outgoing edges. (For example, in Fig. 60(a), we show a network in which each node is each labeled with its demand and each edge with its capacity.)

Recall that, given a flow  $f$  and a node  $v$ ,  $f^{\text{in}}(v)$  is the sum of flows along incoming edges to  $v$  and  $f^{\text{out}}(v)$  is the sum of flows along outgoing edges from  $v$ . We define a *circulation* in  $G$  to be a function  $f$  that assigns a nonnegative real number to each edge that satisfies the following two conditions.

<sup>12</sup>I would not advise applying this in real life. I doubt that the IRS would appreciate it if you paid your \$100 tax bill by demanding that they send you  $-\$100$  dollars.

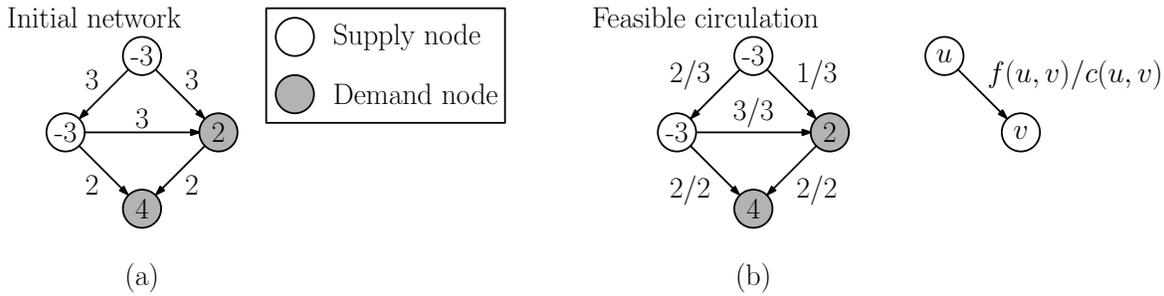


Fig. 60: (a) A circulation network and (b) a feasible circulation.

**Capacity constraints:** For each  $(u, v) \in E$ ,  $0 \leq f(u, v) \leq c(u, v)$ .

**Supply/Demand constraints:** For vertex  $v \in V$ ,  $f^{\text{in}}(v) - f^{\text{out}}(v) = d_v$ .

For example, in Fig. 60(b) we show a valid circulation for the network of part (a). Observe that demand constraints correspond to the flow-balance in the original max flow problem, since if a vertex is not in  $S$  or  $T$ , then  $d_v = 0$  and we have  $f^{\text{in}}(v) = f^{\text{out}}(v)$ . Also it is easy to see that the total demand must equal the total supply, otherwise we have no chance of finding a feasible circulation. That is, we require that

$$\sum_{v \in V} d_v = 0 \quad \text{or equivalently} \quad -\sum_{v \in S} d_v = \sum_{v \in T} d_v.$$

Define  $D = \sum_{v \in T} d_v$  denote the *total demand*. (Note that this is equal to the negation of the total supply,  $\sum_{v \in S} d_v$ .)

**Reducing Circulation to Max-Flow:** Rather than devise an algorithm for the circulation problem, we will show that we can reduce any instance  $G$  of the circulation problem into an equivalent network flow problem. The input to our reduction is a network  $G = (V, E)$ . For each vertex  $v$ , let  $d_v$  denote the demand value and for each edge  $(u, v)$ , let  $c(u, v)$  denote its capacity. First, observe that we may assume that sum of supplies equals total demand (since if not, we can simply answer “no” immediately.) Otherwise:

- Create a new network  $G' = (V', E')$  that has all the same vertices and edges as  $G$  (that is,  $V' \leftarrow V$  and  $E' \leftarrow E$ )
- Add to  $V'$  a *super-source*  $s^*$  and a *super-sink*  $t^*$
- For each supply node  $v \in S$ , we add a new edge  $(s^*, v)$  of capacity  $-d_v$
- For each demand node  $u \in T$ , we add a new edge  $(u, t^*)$  of capacity  $d_u$

(The output of the reduction is illustrated in Fig. 61(b).) We then invoke any max-flow algorithm on  $G'$ . Recalling that  $D$  denotes the total demand, we check whether the value of the maximum flow equals  $D$ . If so, we answer “yes,”  $G$  has a feasible circulation, and otherwise we answer “no,”  $G$  does not have a feasible circulation. (For example, in Fig. 61(c), there exists a flow of value  $D = 6$ , implying that the original network has a circulation.)

We prove correctness below, but intuitively, the newly created edges from  $s^*$  will be responsible for providing the necessary supply for the supply vertices of  $S$  and newly created edges into  $t^*$  are responsible for draining off the excess demand from the vertices of  $T$ . Suppose that we now compute the maximum flow in  $G'$  (by whichever maximum flow algorithm you like).

**Lemma:** There is a feasible circulation in  $G$  if and only if  $G'$  has an  $s^*$ - $t^*$  flow of value  $D$ .

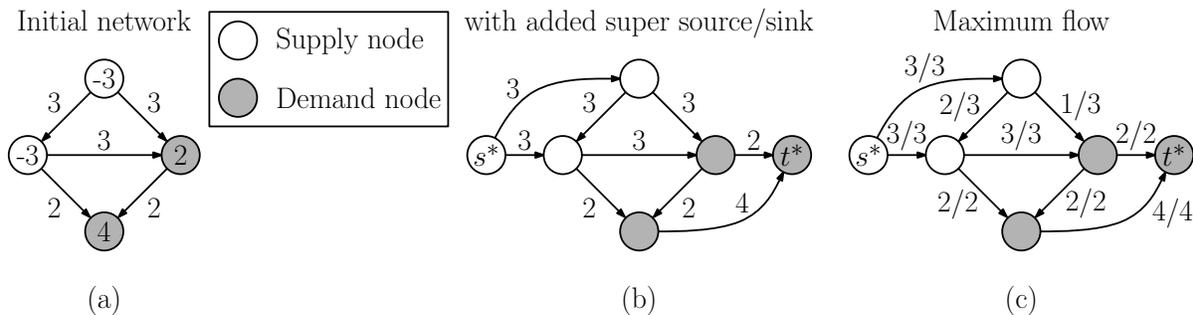


Fig. 61: Reducing the circulation problem to network flow.

**Proof:** ( $\Rightarrow$ ) Suppose that there is a feasible circulation  $f$  in  $G$ . The value of this circulation (the net flow coming out of all supply nodes) is clearly  $D$ . We can create a flow  $f'$  of value  $D$  in  $G'$ , by saturating all the edges coming out of  $s^*$  and all the edges coming into  $t^*$ . We claim that this is a valid flow for  $G'$ . Clearly it satisfies all the capacity constraints. To see that it satisfies the flow balance constraints observe that for each vertex  $v \in V$ , we have one of three cases:

- ( $v \in S$ ) The flow into  $v$  from  $s^*$  matches the supply coming out of  $v$  from the circulation.
- ( $v \in T$ ) The flow out of  $v$  to  $t^*$  matches the demand coming into  $v$  from the circulation.
- ( $v \in V \setminus (S \cup T)$ ) We have  $d_v = 0$ , which means that it satisfies flow conservation by the supply/demand constraints.

( $\Leftarrow$ ) Conversely, suppose that we have a flow  $f'$  of value  $D$  in  $G'$ . It must be that each edge leaving  $s^*$  and each edge entering  $t^*$  is saturated. Therefore, by the flow conservation of  $f'$ , all the supply nodes and all the demand nodes have achieved their desired supply/demand constraints. All the other nodes satisfy their supply/demand constraints because by the flow conservation of  $f'$  the incoming flow equals the outgoing flow. Therefore, the resulting flow is a circulation for  $G$ .

It is not hard to see that the reduction can be performed in  $O(n + m)$  time by a simple analysis of the network's structure. Thus, the overall running time is dominated by the time to compute the network flow (which is  $O(nm)$  according to the current state-of-art).

**Circulations with Upper and Lower Capacity Bounds:** Sometimes, in addition to having a certain maximum flow value, we would also like to impose minimum capacity constraints. That is, given a network  $G = (V, E)$ , for each edge  $(u, v) \in E$  we would like to specify two constraints  $\ell(u, v)$  and  $c(u, v)$ , where  $0 \leq \ell(u, v) \leq c(u, v)$ . A circulation function  $f$  must satisfy the same demand constraints as before, but must also satisfy both the upper and lower flow bounds:

**(New) Capacity Constraints:** For each  $(u, v) \in E$ ,  $\ell(u, v) \leq f(u, v) \leq c(u, v)$ .

**Demand Constraints:** For vertex  $v \in V$ ,  $f^{\text{in}}(v) - f^{\text{out}}(v) = d_v$ .

Henceforth, we will use the term *upper flow bound* in place of *capacity* (since it doesn't make sense to talk about a lower bound as a capacity constraint). An example of such a network is shown in Fig. 62(a), and a valid circulation is shown in Fig. 62(b).

We will reduce this problem to a standard circulation problem (with just the usual upper capacity bounds). To help motivate our reduction, suppose (for conceptual purposes) that we generate an initial (possibly invalid) circulation  $f_0$  that exactly satisfies all the lower flow bounds. In particular, we let  $f_0(u, v) = \ell(u, v)$  (see Fig. 63(a)). This circulation may be invalid because  $f_0$  need not satisfy the demand constraints (which, recall, provide for flow balance as well). We will modify the supply/demand

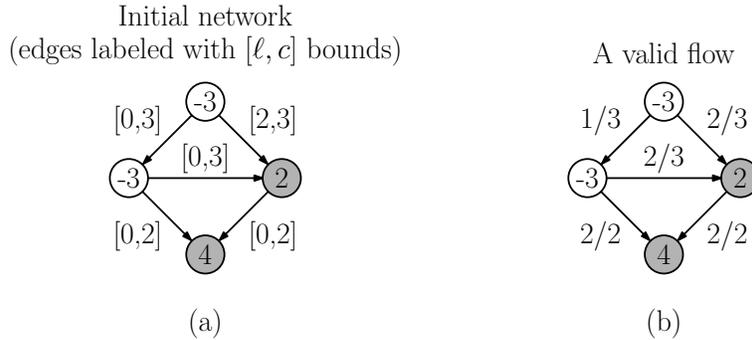


Fig. 62: (a) A network with both upper and lower flow bounds and (b) a valid circulation.

values to compensate for this imbalance. Since the lower-bound constraints are all satisfied, it will be possible to apply a standard circulation algorithm (without lower flow bounds) to solve the problem. To make this formal, for each  $v \in V$ , let  $L_v$  denote the *excess flow* coming into  $v$  in  $f_0$ , that is

$$L_v = f_0^{\text{in}}(v) - f_0^{\text{out}}(v) = \sum_{(u,v) \in E} \ell(u,v) - \sum_{(v,w) \in E} \ell(v,w).$$

(Note that this may be negative, which means that we have a flow deficit.) If we are lucky, then  $L_v = d_v$ , and  $v$ 's supply/demand needs are already met. Otherwise, we will adjust the supply and demand values so that by (1) computing any valid circulation  $f_1$  for the adjusted values and (2) combining this with  $f_0$ , we will obtain a flow that satisfies all the requirements.

How do we adjust the supply/demand values? We want to generate a net flow of  $d_v$  units coming into  $v$  and cancel out the excess  $L_v$  coming in. That is, we want  $f_1$  to satisfy:

$$f_1^{\text{in}}(v) - f_1^{\text{out}}(v) = d_v - L_v.$$

(In particular, this means that if we combine  $f_0$  and  $f_1$  by summing their flows for each edge, then the final flow into  $v$  will be  $d_v$ , as desired.)

How do we determine whether there exists such a circulation  $f_1$ ? Let's consider how to set the edge capacities. We have already sent  $\ell(u,v)$  units of flow through the edge  $(u,v)$ , which implies that we have  $c(u,v) - \ell(u,v)$  capacity remaining. (Note that unlike our definition of residual graphs, we do not want to allow for the possibility of "undoing" flow. Can you see why not?)

We are now ready to put the pieces together. Given the network  $G = (V, E)$  as input (with vertex demands  $d_v$  and lower and upper flow bounds  $\ell(u,v)$  and  $c(u,v)$ ):

1. Create an initial pseudo-circulation  $f_0$  by setting  $f(u,v) = \ell(u,v)$  (see Fig. 63(a))
2. Create a new network  $G' = (V', E')$  that has all the same vertices and edges as  $G$  (that is,  $V' \leftarrow V$  and  $E' \leftarrow E$ )
3. For each  $(u,v) \in E'$ , set its adjusted capacity to  $c'(u,v) \leftarrow c(u,v) - \ell(u,v)$
4. For each  $v \in V'$ , set its adjusted demand to  $d'_v \leftarrow d_v - L_v$ , where,  $L_v = f_0^{\text{in}}(v) - f_0^{\text{out}}(v)$  (see Fig. 63(b))

Now, invoke any (standard) circulation algorithm on  $G'$  (see Fig. 63(c)). Note that there is no need to consider lower flow bounds with  $G'$ , because  $f_0$  has already taken care of those. If the algorithm reports that there is no valid circulation for  $G'$ , then we declare that there is no valid circulation for the original network  $G$ . If the algorithm returns a valid circulation  $f_1$  for  $G'$ , then the final circulation for  $G$  is the combination  $f_0 + f_1$  (see Fig. 63(d)).

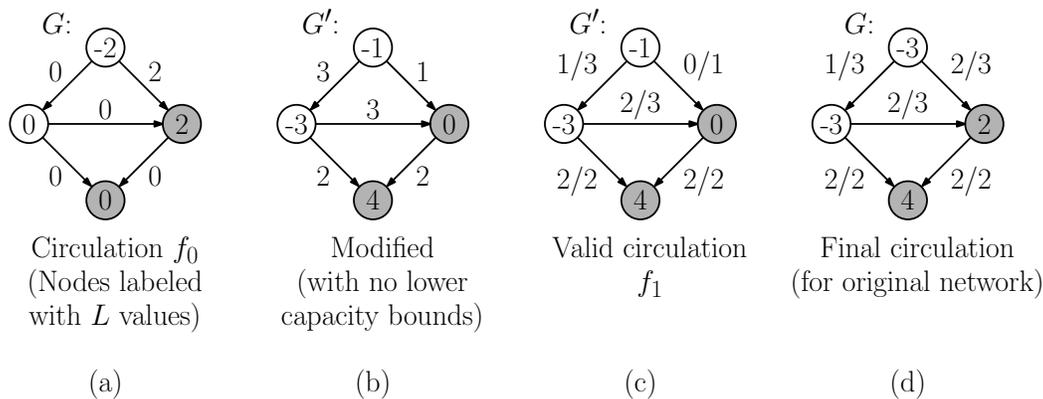


Fig. 63: Reducing the circulation problem with upper and lower flow bounds to a standard circulation problem.

To establish the correctness of our reduction, we prove below that its output  $f_0 + f_1$  is a valid circulation for  $G$  (with lower flow bounds) if and only if  $f_1$  is a valid circulation for  $G'$ .

**Lemma:** The network  $G$  (with both lower and upper flow bounds) has a feasible circulation if and only if  $G'$  (with only upper capacity bounds) has a feasible circulation.

**Proof:** (Sketch. See KL for a formal proof.) Intuitively, if  $G'$  has a feasible circulation  $f'$  then the circulation  $f(u, v) = f'(u, v) + \ell(u, v)$  can be shown to be a valid circulation for  $G$  and it satisfies the lower flow bounds. Conversely, if  $G$  has a feasible circulation (satisfying both the upper and lower flow bounds), then let  $f'(u, v) = f(u, v) - \ell(u, v)$ . As above, it can be shown that  $f'$  is a valid circulation for  $G'$ . (Think of  $f'$  as  $f_1$  and  $f$  as  $f_0 + f_1$ .)

Note that (as in the original circulation problem) we have not presented a new algorithm. Instead, we have shown how to *reduce* the current problem (circulation with lower and upper flow bounds) to a problem we have already solved (circulation with only upper bounds). Again, the running time will be the sum of the time to perform the reduction, which is easily seen to be  $O(n + m)$  plus the time to compute the circulation, which as we have seen reduces to the time to compute a maximum flow, which according to the current best technology is  $O(nm)$  time.

**Application: Survey Design:** To demonstrate the usefulness of circulations with lower flow bounds, let us consider an application problem that arises in the area of data mining. A company sells  $k$  different products, and it maintains a database which stores which customers have bought which products recently. We want to send a survey to a subset of  $n$  customers. We will tailor each survey so it is appropriate for the particular customer it is sent to. Here are some guidelines that we want to satisfy:

- The survey sent to a customer will ask questions only about the products this customer has purchased.
- We want to get as much information as possible, but do not want to annoy the customer by asking too many questions. (Otherwise, they will simply not respond.) Based on our knowledge of how many products customer  $i$  has purchased, and easily they are annoyed, our marketing people have come up with two bounds  $0 \leq c_i \leq c'_i$ . We will ask the  $i$ th customer about at least  $c_i$  products they bought, but (to avoid annoying them) at most  $c'_i$  products.
- Again, our marketing people know that we want more information about some products (e.g., new releases) and less about others. To get a balanced amount of information about each product, for the  $j$ th product we have two bounds  $0 \leq p_j \leq p'_j$ , and we will ask at least  $p_j$  customers about this product and at most  $p'_j$  customers.

We can model this as a bipartite graph  $G$ , in which the customers form one of the parts of the network and products form the other part. There is an edge  $(i, j)$  if customer  $i$  has purchased product  $j$ . The flow through each customer node will reflect the number of products this customer is asked about. The flow through each product node will reflect the number of customers that are asked about this product.

This suggests the following network design. Given the bipartite graph  $G$ , we create a directed network as follows (see Fig. 64).

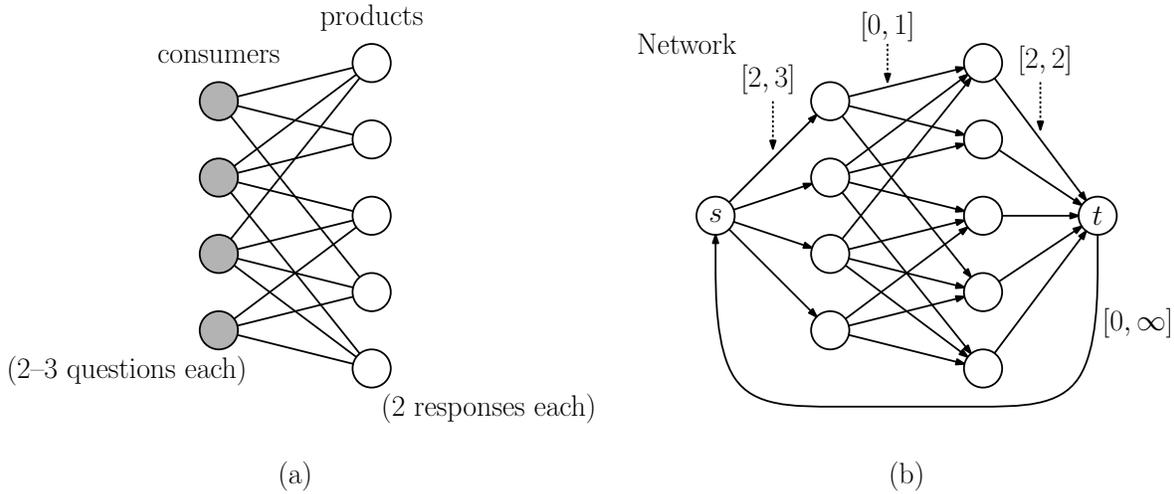


Fig. 64: Reducing the survey design problem to circulation with lower and upper flow bounds. (Assuming  $[c_i, c'_i] = [2, 3]$  for all customers and  $[p_j, p'_j] = [2, 2]$  for all products.)

- For each customer  $i$  who purchased product  $j$  we create a directed edge  $(i, j)$  with an upper flow bounds of 1, respectively. This models the requirement that customer  $i$  will be surveyed at most once about product  $j$ , and customers will be asked only about products they purchased.
- We create a source vertex  $s$  and connect it to all the customers, where the edge from  $s$  to customer  $i$  has lower and upper flow bounds of  $c_i$  and  $c'_i$ , respectively. This models the requirement that customer  $i$  will be asked about at least  $c_i$  products and at most  $c'_i$ .
- We create a sink vertex  $t$ , and create an edge from product  $j$  to  $t$  with lower and upper flow bounds of  $p_j$  and  $p'_j$ . This models the requirement that there are at least  $p_j$  and at most  $p'_j$  customers will be asked about product  $j$ .
- We create an edge  $(s, t)$ . Its lower bound is set to zero and its upper bound can be set to any very large value. This is needed for technical reasons, since we want a circulation.
- All node demands are set to 0.

The correctness of the reduction is established in the following lemma.

**Lemma:** There exists a valid circulation in  $G$  if and only there is a valid survey design.

**Proof:** ( $\Rightarrow$ ) Suppose that  $G$  has a valid (integer-valued) circulation. For each customer-product edge  $(i, j)$  that carries one unit of flow, customer  $i$  is surveyed about product  $j$ . By definition of the edges, customers are surveyed only about products they purchased. From our capacity constraints and the fact that demands are all zero, it follows that the total flow into each customer node is between  $c_i$  and  $c'_i$ , implying that this customer is asked about this many products. Similarly, the

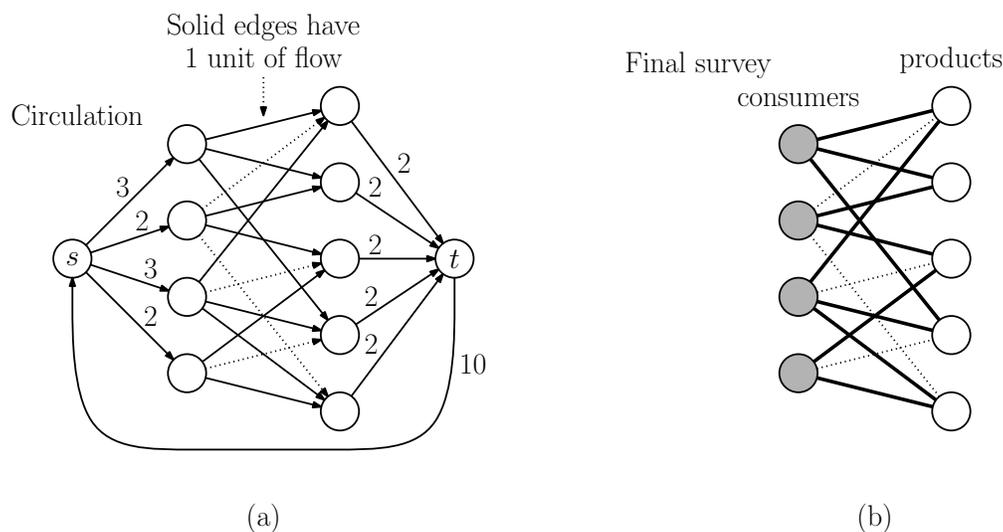


Fig. 65: Correctness of the survey-design reduction to circulation.

total flow out of each product node is between  $p_j$  and  $p'_j$ , implying that this product is involved in this many products surveys (see Fig. 65).

( $\Leftarrow$ ) Suppose that there is a valid survey design. We construct a flow in  $G$  as follows. For each customer-product pair  $(i, j)$  involved in the survey, we create a flow of one unit on edge  $(i, j)$ . We set the flow along the edge  $(s, i)$  to the number of surveys that customer  $i$  answers, we set the flow along the edge  $(j, t)$  to the number of surveys involving product  $j$ , and we set the flow on edge  $(t, s)$  to the total number of surveys. It is straightforward to see that (by the rules of a valid survey design) this is a valid flow in  $G$ , and in particular, it satisfies the lower and upper capacity constraints and the supply and demand constraints.

## Lecture 18: NP-Completeness: General Definitions

**Efficiency and Polynomial Time:** Up to this point of the semester we have been building up your “toolkit” for solving algorithmic problems efficiently. Hopefully when presented with a computational problem, you now have a clearer idea the sorts of techniques that could be used to solve the problem efficiently (such as divide-and-conquer, DFS, greedy, dynamic programming, network flow).

What do we mean when we say “efficient”? If  $n$  is small, a running time of  $2^n$  may be just fine, but when  $n$  is huge, even  $n^2$  may be unacceptably slow. Algorithm designers observed long ago that there are two very general classes of combinatorial problems:

- those involving *brute-force search* of all feasible solutions, whose worst-case running time is an *exponential function* of the input size,
- those that are based on a *systematic solution*, whose worst-case running time is a *polynomial function* of the input size.

An algorithm is said to run in *polynomial time* if its worst-case running time is  $O(n^c)$ , where  $c$  is a nonnegative constant. (Note that running times like  $O(n \log n)$  are polynomial time, since  $n \log n = O(n^2)$ .) By *exponential time* we mean any function that is at least  $\Omega(c^n)$  for a constant  $c > 1$ . Henceforth, we will use the terms “efficient” and “easy” to mean solvable by an algorithm whose

worst-case running time is polynomial in the input size. (irrespective of whether the polynomial is  $n$  or  $n^{1000}$ ).

While the distinction between worst-case polynomial time and worst-case exponential time is quite crude, it has a number of advantages. For example, the composition of any two polynomials is a polynomial. (That is, if  $f(n)$  and  $g(n)$  are both polynomials, then so is  $f(g(n))$ .) This means that, if a program makes a polynomial number of calls to a function that runs in polynomial time, then the overall running time is a polynomial.

**The Emergence of Hard Problems:** Near the end of the 60's, although there was great success in finding efficient solutions to many combinatorial problems, there was also a growing list of problems which were "hard" in the sense that no known efficient algorithmic solutions existed for these problems.

A remarkable discovery was made about this time. Many of these believed hard problems turned out to be equivalent, in the sense that if you could solve *any one* of them in polynomial time, then you could solve *all* of them in polynomial time. Often these hard problems involved slight generalizations to problems that are solvable in polynomial time. A list of some of these problems is shown in Table 2.

Table 2: Computationally hard problems and their (easy) counterparts.

Hard problems (NP-complete)	Easy problems (in P)
3SAT	2SAT
Traveling Salesman Problem (TSP)	Minimum Spanning Tree (MST)
Longest (Simple) Path	Shortest Path
3D Matching	Bipartite Matching
Knapsack	Unary Knapsack
Independent Set in Graphs	Independent Set in Trees
Integer Linear Programming	Linear Programming
Hamiltonian Cycle	Eulerian Cycle
Balanced Cut	Minimum Cut

The mathematical theory, which was developed by Richard Karp and Stephen Cook, gave rise to the notions of P, NP, and NP-completeness. Since then, thousands of problems were identified as being in this equivalence class. It is widely believed that none of them can be solved in polynomial time, but there is no proof of this fact. This has given rise to arguably the biggest open problems in computer science:

$$P = NP?$$

While we will not be able to provide an answer to this question, we will investigate this concept in the next few lectures.

Note that represents a radical departure from what we have been doing so far this semester. The goal is no longer to prove that a problem *can* be solved efficiently by presenting an algorithm for it. Instead we will be trying to show that a problem *cannot* be solved efficiently. The question is how to do this?

**Reasonable Input Encodings:** When trying to show the impossibility of achieving a task efficiently, it is important to define terms precisely. Otherwise, we might be beaten by clever cheats. We will treat the input to our problems as a string over some alphabet that has a constant number, but at least two, characters (e.g., a binary bit string or a Unicode encoding). If you think about it for just a moment, every data structure that we have seen this semester can be *serialized* into such a string, without increasing its size significantly.

How are inputs to be encoded? Observe that if you encode an integer in a very inefficient manner, for example, using *unary notation* (so that 8 is represented as 1111111), rather than an efficient encoding

(say in binary or decimal<sup>13</sup>), the length of the string increases exponentially. Why should we care? Observe that if the input size grows exponentially, then an algorithm that ran in exponential time for the short input size may now run in linear time for the long input size. We consider this a cheat because we haven't devised a faster algorithm, we have just made our measuring yardstick much much longer.

All the representations we have seen this semester (e.g., sets as lists, graphs as adjacency lists or adjacency matrices, etc.) are considered to be reasonable. To determine whether some new representation is reasonable, it should be as concise as possible (in the worst case) and/or it should be possible to convert from an existing reasonable representation to this new form in polynomial time.

**Decision Problem:** Many of the problems that we have discussed involve *optimization* of one form or another: find the shortest path, find the minimum cost spanning tree, find the maximum flow. For rather technical reasons, most NP-complete problems that we will discuss will be phrased as decision problems.

A problem is called a *decision problem* if its output is a simple “yes” or “no” (or you may think of this as True/False, 0/1, accept/reject). For example, the minimum spanning tree decision problem might be: “Given a weighted graph  $G$  and an integer  $z$ , does  $G$  have a spanning tree whose weight is at most  $z$ ?”

This may seem like a less interesting formulation of the problem. It does not ask for the weight of the minimum spanning tree, and it does not even ask for the edges of the spanning tree that achieves this weight. However, our job will be to show that certain problems *cannot* be solved efficiently. If we show that the simple decision problem cannot be solved efficiently, then certainly the more general optimization problem certainly cannot be solved efficiently either. (In fact, if you can solve a decision problem efficiently, it is almost always possible to construct an efficient solution to the optimization problem, but this is a technicality that we won't worry about now.)

**Language Recognition:** Observe that a decision problem can also be thought of as a language recognition problem. Define a *language* to be a set (finite or infinite) of strings. To express a computational problem as a language-recognition problem, we first should be able to express its input as a string. Let us define  $\text{serialize}(X)$  to be a function that maps any combinatorial structure  $X$  into a string. (For example, serializing a graph would involve outputting a string that encodes its vertices, edges, and any additional information such as edge weights.)

Using this, we could define a language MST encoding the Minimum Spanning Tree problem as a language (that is, a collection of strings):

$$\text{MST} = \{\text{serialize}(G, z) \mid G \text{ has a minimum spanning tree of weight at most } z\}.$$

What does it mean to solve the decision problem? When presented with a specific input string  $x = \text{serialize}(G, z)$ , the algorithm would answer “yes” if  $x \in \text{MST}$ , that is, if  $G$  has a spanning tree of weight at most  $z$ , and “no” otherwise. In the first case we say that the algorithm *accepts* the input and otherwise it *rejects* the input. Thus, decision problems are equivalent to language-recognition problems.

Given an input  $x$ , how would we determine whether  $x \in \text{MST}$ ? First, we would decode  $x$  as  $G$  and  $z$ . We would then feed these into any efficient minimum spanning tree algorithm (Kruskal's, say). If the final cost of the spanning tree is at most  $z$ , we accept  $x$  and otherwise we reject it.

**The Class P:** We now present an important definition:

**Definition:** P is the set of all languages (i.e., decision problems) for which membership can be determined in (worst case) polynomial time.

---

<sup>13</sup>The exact choice of the numeric base is not important so long as it is at least 2, since all base representations can be converted to each other with only a constant factor change in the length.

Intuitively,  $P$  corresponds to the set of all decision problems that can be solved efficiently, that is, in polynomial time. Note  $P$  is not a language, rather, it is a set of languages. A set of languages that is defined in terms of how hard it is to determine membership is called a *complexity class*. (Therefore,  $P$  is a complexity class.)

Since Kruskal's algorithm runs in polynomial time, we have  $MST \in P$ . We could define equivalent languages for all of the other optimization problems we have seen this year (e.g., shortest paths, max flow, min cut).

**A Harder Example:** To show that not all languages are (obviously) in  $P$ , consider the following:

$$HC = \{\text{serialize}(G) \mid G \text{ has a simple cycle that visits every vertex of } G\}.$$

Such a cycle is called a *Hamiltonian cycle* and the decision problem is the *Hamiltonian Cycle Problem*.

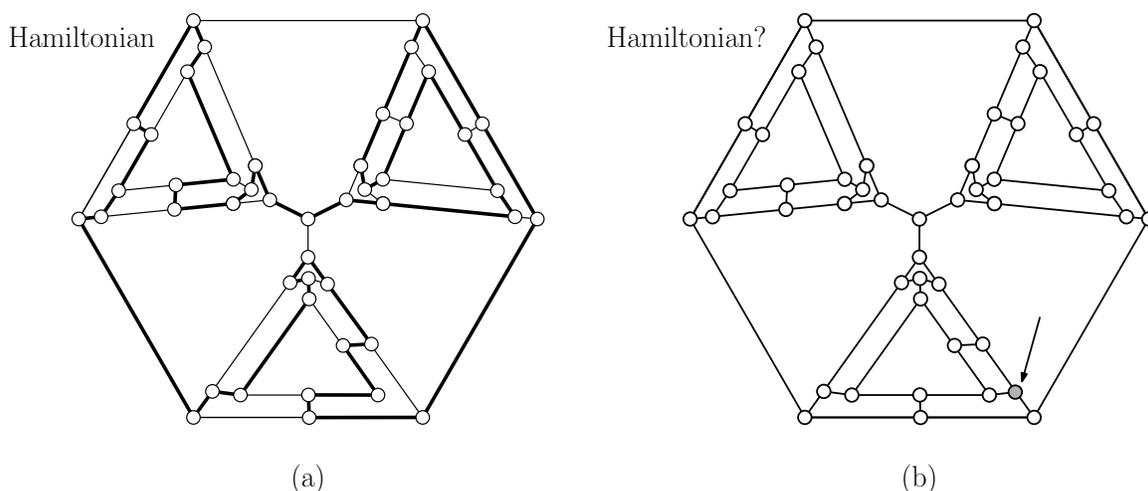


Fig. 66: The Hamiltonian cycle (HC) problem.

In Fig. 66(a) we show an example of a Hamiltonian cycle in a graph. If you think that the problem is easy to solve, try to solve the problem on the graph shown in Fig. 66(b), which has one additional vertex and one additional edge. Either find a Hamiltonian cycle in this graph or show that none exists. To make this even harder, imagine a million-vertex graph with many slight variations of this pattern.

Is  $HC \in P$ ? No one knows the answer for sure, but it is conjectured that it is not. (In fact, we will show that later that  $HC$  is NP-complete.)

In what follows, we will be introducing a number of classes. We will jump back and forth between the terms “language” and “decision problems”, but for our purposes they mean the same things. Before giving all the technical definitions, let us say a bit about what the general classes look like at an intuitive level.

**Polynomial-Time Verification and Certificates:** In order to define NP-completeness, we need to first define NP. Unfortunately, providing a rigorous definition of NP will involve a presentation of the notion of *nondeterministic* models of computation, and will take us away from our main focus. (Formally, NP stands for *nondeterministic polynomial time*.) Instead, we will present a very simple, “hand-wavy” definition, which will suffice for our purposes.

To do so, it is important to first introduce the notion of a verification algorithm. Many language recognition problems that may be *hard to solve*, but they have the property that they are *easy to verify*

that a string is in the language. Recall the Hamiltonian cycle problem defined above. As we saw, there is no obviously efficient way to find a Hamiltonian cycle in a graph. However, suppose that a graph did have a Hamiltonian cycle and someone wanted to convince us of its existence. This person would simply tell us the vertices in the order that they appear along the cycle. It would be a very easy matter for us to inspect the graph and check that this is indeed a legal cycle that it visits all the vertices exactly once. Thus, even though we know of no efficient way to *solve* the Hamiltonian cycle problem, there is a very efficient way to *verify* that a given graph has one. (You might ask, but what if the graph did not have one? Don't worry. A verification process is not required to do anything if the input is not in the language.)

The given cycle in the above example is called a *certificate*. A certificate is a piece of information which allows us to verify that a given string is in a language in polynomial time.

More formally, given a language  $L$ , and given  $x \in L$ , a *verification algorithm* is an algorithm which, given  $x$  and a string  $y$  called the *certificate*, can verify that  $x$  is in the language  $L$  using this certificate as help. If  $x$  is not in  $L$  then there is nothing to verify. If there exists a verification algorithm that runs in polynomial time, we say that  $L$  can be *verified in polynomial time*.

Note that not all languages have the property that they are easy to verify. For example, consider the following languages:

$$\begin{aligned} \text{UHC} &= \{G \mid G \text{ has a unique Hamiltonian cycle}\} \\ \overline{\text{HC}} &= \{G \mid G \text{ has no Hamiltonian cycle}\}. \end{aligned}$$

There is no known polynomial time verification algorithm for either of these. For example, suppose that a graph  $G$  is in the language UHC. What information would someone give us that would allow us to verify that  $G$  is indeed in the language? They could certainly show us one Hamiltonian cycle, but it is unclear that they could provide us with any easily verifiable piece of information that would demonstrate that this is the only one.

**The class NP:** We can now define the complexity class NP.

**Definition:** NP is the set of all languages that can be verified in polynomial time.

Observe that if we can solve a problem efficiently without a certificate, we can certainly solve given the additional help of a certificate. Therefore,  $P \subseteq NP$ . However, it is not known whether  $P = NP$ . It seems unreasonable to think that this should be so. In other words, just being able to verify that you have a correct solution does not help you in finding the actual solution very much. Most experts believe that  $P \neq NP$ , but no one has a proof of this. Next time we will define the notions of NP-hard and NP-complete.

There is one last ingredient that will be needed before defining NP-completeness, namely the notion of a *polynomial time reduction*. We will discuss that next time.

## Lecture 19: NP-Completeness: Reductions

**Recap:** We have introduced a number of concepts on the way to defining NP-completeness:

**Decision Problems/Language recognition:** are problems for which the answer is either yes or no. These can also be thought of as language recognition problems, assuming that the input has been encoded as a string. For example:

$$\begin{aligned} \text{HC} &= \{G \mid G \text{ has a Hamiltonian cycle}\} \\ \text{MST} &= \{(G, c) \mid G \text{ has a MST of cost at most } c\}. \end{aligned}$$

**P:** is the class of all decision problems which can be solved in polynomial time. While  $MST \in P$ , we do not know whether  $HC \in P$  (but we suspect not).

**Certificate:** is a piece of evidence that allows us to *verify* in polynomial time that a string is in a given language. For example, the language  $HC$  above, a certificate could be a sequence of vertices along the cycle. (If the string is not in the language, the certificate can be anything.)

**NP:** is defined to be the class of all languages that can be *verified* in polynomial time. (Formally, it stands for *Nondeterministic Polynomial time*.) Clearly,  $P \subseteq NP$ . It is widely believed that  $P \neq NP$ .

To define NP-completeness, we need to introduce the concept of a reduction.

**Reductions:** The class of NP-complete problems consists of a set of decision problems (languages) (a subset of the class NP) that no one knows how to solve efficiently, but if there were a polynomial time solution for even a single NP-complete problem, then every problem in NP would be solvable in polynomial time.

Before discussing reductions, let us just consider the following question. Suppose that there are two problems,  $H$  and  $U$ . We know (or you strongly believe at least) that  $H$  is *hard*, that is it cannot be solved in polynomial time. On the other hand, the complexity of  $U$  is *unknown*. We want to prove that  $U$  is also hard. How would we do this? We effectively want to show that

$$(H \notin P) \Rightarrow (U \notin P).$$

To do this, we could prove the contrapositive,

$$(U \in P) \Rightarrow (H \in P).$$

To show that  $U$  is not solvable in polynomial time, we will suppose (towards a contradiction) that a polynomial time algorithm for  $U$  did exist, and then we will use this algorithm to solve  $H$  in polynomial time, thus yielding a contradiction.

To make this more concrete, suppose that we have a subroutine<sup>14</sup> that can solve any instance of problem  $U$  in polynomial time. Given an input  $x$  for the problem  $H$ , we could translate it into an *equivalent* input  $x'$  for  $U$ . By “equivalent” we mean that  $x \in H$  if and only if  $x' \in U$  (see Fig. 67). Then we run our subroutine on  $x'$  and output whatever it outputs.

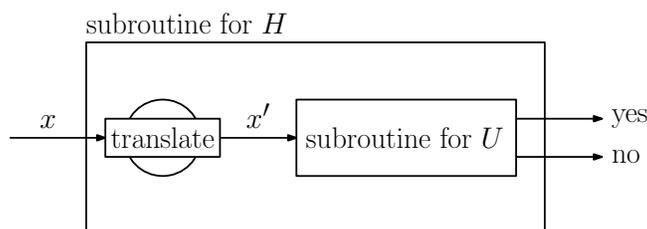


Fig. 67: Reducing  $H$  to  $U$ .

It is easy to see that if  $U$  is solvable in polynomial time, then so is  $H$ . We assume that the translation module runs in polynomial time. If so, we say we have a *polynomial reduction* of problem  $H$  to problem  $U$ , which is denoted  $H \leq_P U$ . This is called a *Karp reduction*.

More generally, we might consider calling the subroutine multiple times. How many times can we call it? Since the composition of two polynomials is a polynomial, we may call it any polynomial

<sup>14</sup>It is important to note here that this supposed subroutine for  $U$  is a *fantasy*. We know (or strongly believe) that  $H$  cannot be solved in polynomial time, thus we are essentially proving that such a subroutine cannot exist, implying that  $U$  cannot be solved in polynomial time.

number of times. A reduction based on making multiple calls to such a subroutine is called a *Cook reduction*. Although Cook reductions are theoretically more powerful than Karp reductions, every NP-completeness proof that I know of is based on the simpler Karp reductions.

**3-Colorability and Clique Cover:** Let us consider an example to make this clearer. The following problem is well-known to be NP-complete, and hence it is strongly believed that the problem cannot be solved in polynomial time.

**3-coloring (3Col):** Given a graph  $G$ , can each of its vertices be labeled with one of three different “colors”, such that no two adjacent vertices have the same label (see Fig. 68(a) and (b)).

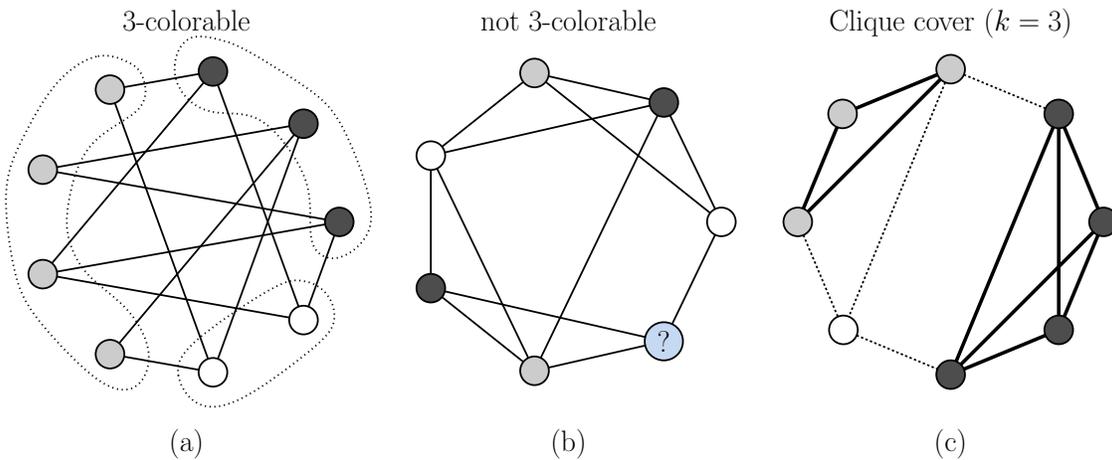


Fig. 68: 3-coloring and Clique Cover.

Coloring arises in various partitioning problems, where there is a constraint that two objects cannot be assigned to the same set of the partition. It is well known that planar graphs can be colored with four colors, and there exists a polynomial time algorithm for doing this. But determining whether three colors are possible (even for planar graphs) seems to be hard, and there is no known polynomial time algorithm.

The 3Col problem will play the role of the known hard problem  $H$ . To play the role of  $U$ , consider the following problem. Given a graph  $G = (V, E)$ , we say that a subset of vertices  $V' \subseteq V$  forms a *clique* if for every pair of distinct vertices  $u, v \in V'$   $(u, v) \in E$ . That is, the subgraph induced by  $V'$  is a complete graph.

**Clique Cover (CCov):** Given a graph  $G = (V, E)$  and an integer  $k$ , can we partition the vertex set into  $k$  subsets of vertices  $V_1, \dots, V_k$  such that each  $V_i$  is a clique of  $G$  (see Fig. 68(c)).

The clique cover problem arises in clustering. We put an edge between two nodes if they are similar enough to be clustered in the same group. We want to know whether it is possible to cluster all the vertices into at most  $k$  groups.

We want to prove that CCov is hard, under the assumption that 3Col is hard, that is,

$$(3\text{Col} \notin \text{P}) \implies (\text{CCov} \notin \text{P}).$$

Again, we'll prove the contrapositive:

$$(\text{CCov} \in \text{P}) \implies (3\text{Col} \in \text{P}).$$

Let us assume that we have access to a polynomial time subroutine  $\text{CCov}(G, k)$ . Given a graph  $G$  and an integer  $k$ , this subroutine returns true (or “yes”) if  $G$  has a clique cover of size  $k$  and false otherwise. How can we use this *alleged* subroutine to solve the well-known hard 3Col problem? We need to find a translation, that maps an instance  $G$  for 3-coloring into an instance  $(G', k)$  for clique cover (see Fig. 69).

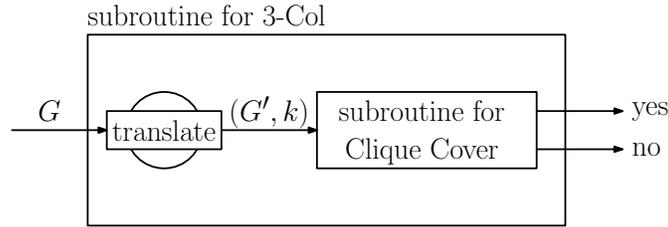


Fig. 69: Reducing 3Col to CCov.

Observe that both problems involve partitioning the vertices up into groups. There are two differences. First, in the 3-coloring problem, the number of groups is fixed at three. In the Clique Cover problem, the number is given as an input. Second, in the 3-coloring problem, in order for two vertices to be in the same group they should *not* have an edge between them. In the Clique Cover problem, for two vertices to be in the same group, they *must* have an edge between them. Our translation therefore, should convert edges into non-edges and vice versa.

This suggests the following idea for reducing the 3-coloring problem to the Clique Cover problem. Given a graph  $G$ , let  $\bar{G}$  denote the *complement graph*, where two distinct nodes are connected by an edge if and only if they are not adjacent in  $G$ . Let  $G$  be the graph for which we want to determine its 3-colorability. The translator outputs the pair  $(\bar{G}, 3)$ . We then feed the pair  $(G', k) = (\bar{G}, 3)$  into a subroutine for clique cover (see Fig. 70).

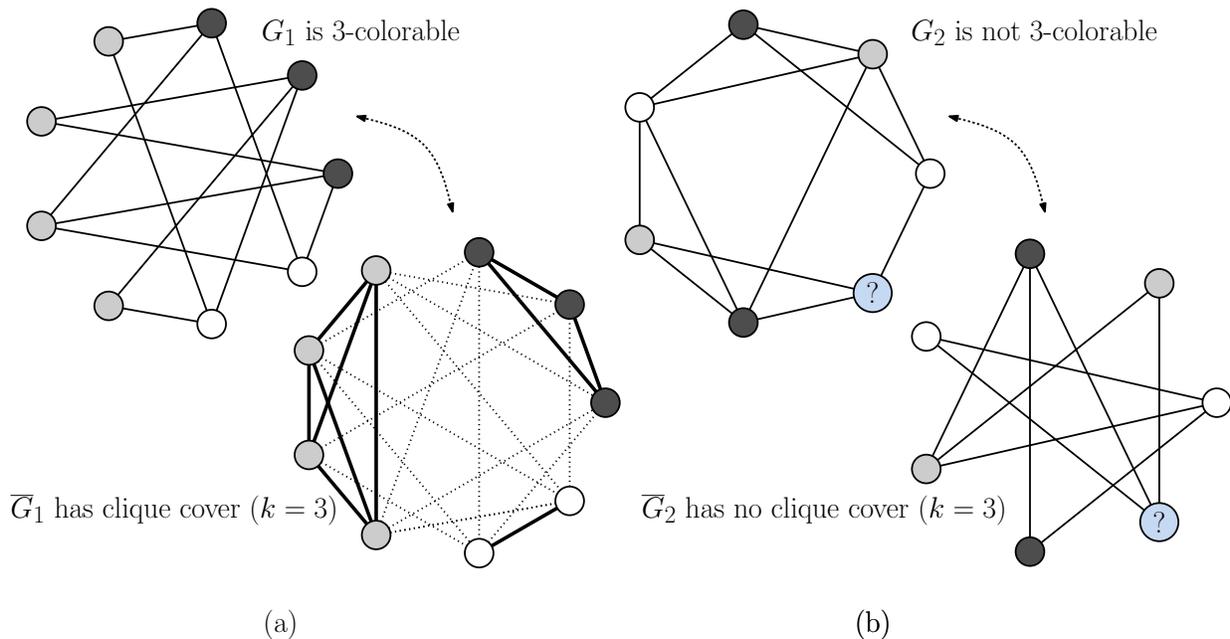


Fig. 70: Clique covers in the complement.

The following formally establishes the correctness of this reduction.

**Claim:** A graph  $G = (V, E)$  is 3-colorable if and only if its complement  $\overline{G} = (V, \overline{E})$  has a clique-cover of size 3. In other words,

$$G \in 3\text{Col} \iff (\overline{G}, 3) \in \text{CCov}.$$

**Proof:** ( $\Rightarrow$ ) If  $G$  3-colorable, then let  $V_1, V_2, V_3$  be the three color classes. We claim that this is a clique cover of size 3 for  $\overline{G}$ , since if  $u$  and  $v$  are distinct vertices in  $V_i$ , then  $\{u, v\} \notin E$  (since adjacent vertices cannot have the same color) which implies that  $\{u, v\} \in \overline{E}$ . Thus every pair of distinct vertices in  $V_i$  are adjacent in  $\overline{G}$ .

( $\Leftarrow$ ) Suppose  $\overline{G}$  has a clique cover of size 3, denoted  $V_1, V_2, V_3$ . For  $i \in \{1, 2, 3\}$  give the vertices of  $V_i$  color  $i$ . We assert that this is a legal coloring for  $G$ , since if distinct vertices  $u$  and  $v$  are both in  $V_i$ , then  $\{u, v\} \in \overline{E}$  (since they are in a common clique), implying that  $\{u, v\} \notin E$ . Hence, two vertices with the same color are not adjacent.

**Polynomial-time reduction:** We now take this intuition of reducing one problem to another through the use of a subroutine call, and place it on more formal footing. Notice that in the example above, we converted an instance of the 3-coloring problem ( $G$ ) into an equivalent instance of the Clique Cover problem  $(\overline{G}, 3)$ .

**Definition:** We say that a language (i.e. decision problem)  $L_1$  is *polynomial-time reducible* to language  $L_2$  (written  $L_1 \leq_P L_2$ ) if there is a polynomial time computable function  $f$ , such that for all  $x$ ,  $x \in L_1$  if and only if  $f(x) \in L_2$ .

In the previous example we showed that  $3\text{Col} \leq_P \text{CCov}$ , and in particular,  $f(G) = (\overline{G}, 3)$ . Note that it is easy to complement a graph in  $O(n^2)$  (i.e. polynomial) time (e.g. flip 0's and 1's in the adjacency matrix). Thus  $f$  is computable in polynomial time.

Intuitively, saying that  $L_1 \leq_P L_2$  means that “if  $L_2$  is solvable in polynomial time, then so is  $L_1$ .” This is because a polynomial time subroutine for  $L_2$  could be applied to  $f(x)$  to determine whether  $f(x) \in L_2$ , or equivalently whether  $x \in L_1$ . Thus, in sense of polynomial time computability,  $L_1$  is “no harder” than  $L_2$ .

The way in which this is used in NP-completeness is exactly the converse. We usually have strong evidence that  $L_1$  is not solvable in polynomial time, and hence the reduction is effectively equivalent to saying “since  $L_1$  is not likely to be solvable in polynomial time, then  $L_2$  is also not likely to be solvable in polynomial time.” Thus, this is how polynomial time reductions can be used to show that problems are as hard to solve as known difficult problems.

**Lemma:** If  $L_1 \leq_P L_2$  and  $L_2 \in P$  then  $L_1 \in P$ .

**Lemma:** If  $L_1 \leq_P L_2$  and  $L_1 \notin P$  then  $L_2 \notin P$ .

Because the composition of two polynomials is a polynomial, we can chain reductions together.

**Lemma:** If  $L_1 \leq_P L_2$  and  $L_2 \leq_P L_3$  then  $L_1 \leq_P L_3$ .

**NP-completeness:** The set of NP-complete problems are all problems in the complexity class NP, for which it is known that if any one is solvable in polynomial time, then they all are, and conversely, if any one is not solvable in polynomial time, then none are. This is made mathematically formal using the notion of polynomial time reductions.

**Definition:** A language  $L$  is *NP-hard* if  $L' \leq_P L$ , for all  $L' \in \text{NP}$ . (Note that  $L$  does not need to be in NP.)

**Definition:** A language  $L$  is *NP-complete* if:

- (1)  $L \in \text{NP}$  (that is, it can be verified in polynomial time), and
- (2)  $L$  is NP-hard (that is, every problem in NP is polynomially reducible to it).

An alternative (and usually easier way) to show that a problem is NP-complete is to use transitivity.

**Lemma:**  $L$  is NP-complete if

- (1)  $L \in \text{NP}$  and
- (2)  $L' \leq_P L$  for some *known* NP-complete language  $L'$ .

The reason is that all  $L'' \in \text{NP}$  are reducible to  $L'$  (since  $L'$  is NP-complete and hence NP-hard) and hence by transitivity  $L''$  is reducible to  $L$ , implying that  $L$  is NP-hard.

This gives us a way to prove that problems are NP-complete, once we know that *one* problem is NP-complete. Unfortunately, it appears to be almost impossible to prove that one problem is NP-complete, because the definition says that we have to be able to reduce *every* problem in NP to this problem. There are infinitely many such problems, so how can we ever hope to do this?

We will talk about this next time with Cook's theorem. Cook showed that there is one problem called SAT (short for *boolean satisfiability*) that is NP-complete. To prove a second problem is NP-complete, all we need to do is to show that our problem is in NP (and hence it is reducible to SAT), and then to show that we can reduce SAT (or generally some known NPC problem) to our problem. It follows that our problem is equivalent to SAT (with respect to solvability in polynomial time). This is illustrated in Fig. 71 below.

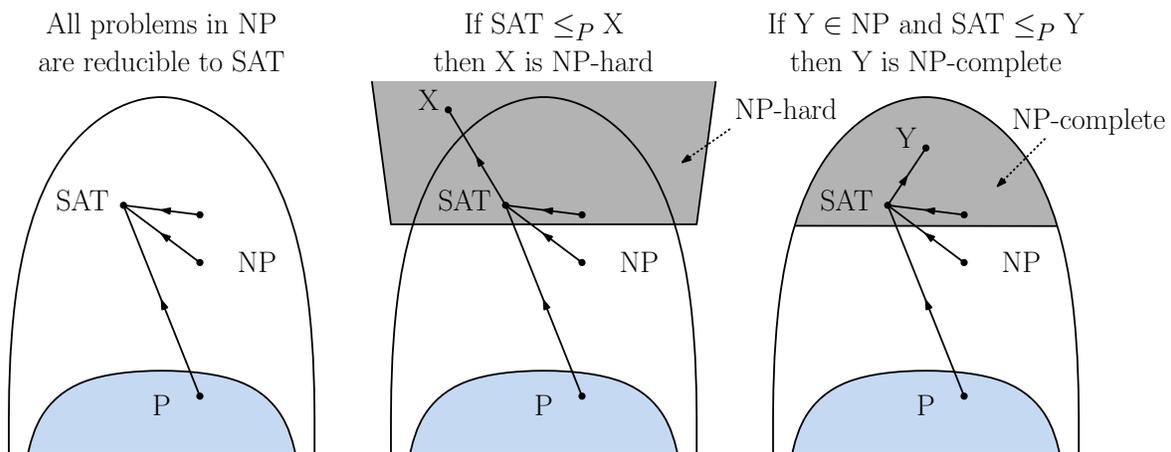


Fig. 71: Structure of NPC and reductions.

## Lecture 20: NP-Completeness: 3SAT and Independent Set

**Recap:** Recall the following definitions, which were given in earlier lectures.

**P:** is the set of decisions problems solvable in polynomial time, or equivalently, the set of languages for which membership can be determined in polynomial time.

**NP:** is the set of languages that can be *verified* in polynomial time, or equivalently, that can be solved in polynomial time by a “guessing computer”, whose guesses are guaranteed to produce an output of “yes” if at all possible.

**Polynomial reduction:**  $L_1 \leq_P L_2$  means that there is a polynomial time computable function  $f$  such that  $x \in L_1$  if and only if  $f(x) \in L_2$ . A more intuitive way to think about this is that if we had a subroutine to solve  $L_2$  in polynomial time, then we could use it to solve  $L_1$  in polynomial time. Polynomial reductions are *transitive*, that is,  $L_1 \leq_P L_2$  and  $L_2 \leq_P L_3$  implies  $L_1 \leq_P L_3$ .

**NP-Hard:**  $L$  is NP-hard if for all  $L' \in \text{NP}$ ,  $L' \leq_P L$ . By transitivity of  $\leq_P$ , we can say that  $L$  is NP-hard if  $L' \leq_P L$  for some known NP-hard problem  $L'$ .

**NP-Complete:**  $L$  is NP-complete if (1)  $L \in \text{NP}$  and (2)  $L$  is NP-hard.

It follows from these definitions that:

- If *any* NP-hard problems is solvable in polynomial time, then *every* NP-complete problem (in fact, every problem in NP) is also solvable in polynomial time.
- If *any* NP-complete problem cannot be solved in polynomial time, then *every* NP-complete problem (in fact, every NP-hard problem) cannot be solved in polynomial time.

Thus all NP-complete problems are equivalent to one another (in that they are either all solvable in polynomial time, or none are).

**Cook's Theorem:** To get the ball rolling, we need to prove that there is *at least one* NP-complete problem. Stephen Cook achieved this task. This first NP-complete problem involves boolean formulas. A boolean formula consists of variables (say  $x$ ,  $y$ , and  $z$ ) and the logical operations *not* (denoted  $\bar{x}$ ), *and* (denoted  $x \wedge y$ ), and *or* (denoted  $x \vee y$ ).

Given a boolean formula, we say that it is *satisfiable* if there is a way to assign truth values (0 or 1) to the variables such that it evaluates to 1. (As opposed to the case where every variable assignment results in 0.) For example, consider the following formula:

$$F_1(x, y, z) = (x \wedge (y \vee \bar{z})) \wedge ((\bar{y} \wedge \bar{z}) \vee \bar{x}).$$

$F_1$  is satisfiable, by the assignment  $x = 1$  and  $y = z = 0$ . On the other hand, the formula

$$F_2(x, y) = (\bar{z} \vee x) \wedge (z \vee y) \wedge (\bar{x} \wedge \bar{y})$$

is not satisfiable since every possible assignment of 0-1 values to  $x$ ,  $y$ , and  $z$  evaluates to 0.

The *boolean satisfiability problem* (SAT) is as follows: given a boolean formula  $F$ , is it possible to assign truth values (0/1, true/false) to  $F$ 's variables, so that it evaluates to true?

**Cook's Theorem:** SAT is NP-complete.

A complete proof would take about a full lecture (not counting the week or so of background on nondeterminism and Turing machines). Here is an intuitive justification.

**SAT is in NP:** We nondeterministically guess truth values to the variables. (In the context of verification, the certificate consists of the assignment of values to the variables.) We then plug the values into the formula and evaluate it. Clearly, this can be done in polynomial time.

**SAT is NP-Hard:** To show that the 3SAT is NP-hard, Cook reasoned as follows. First, every NP-problem can be encoded as a program that runs in polynomial time on a given input, subject to a number of nondeterministic guesses. Since the program runs in polynomial time, we can express its execution on a specific input as straight-line program (that is, one containing no loops or function calls) that contains a polynomial number of lines of code in your favorite programming language. We then compile each line of code into machine code, and convert each machine code instruction into an equivalent boolean circuit. Finally, we can express each of these circuits equivalently as a boolean formula.

The nondeterministic choices can be implemented as boolean variables in this formula, whose values take on the possible values of 0 and 1. By definition of nondeterminism, the program answers “yes” if there is some choice of decisions that leads to an output of “yes”. In our context, this means that there is some way of assigning 0-1 values to the variables so that our circuit produces an output of 1, that is, if the associated boolean formula is satisfied.

Therefore, if you *could* determine the satisfiability of this formula in polynomial time, you could determine whether the original nondeterministic program output “yes” in polynomial time.

Cook proved that satisfiability is NP-hard even for boolean formulas of a special form. To define this form, we start by defining a *literal* to be either a variable or its negation, that is,  $x$  or  $\bar{x}$ . A formula is said to be in *3-conjunctive normal form* (3-CNF) if it is the boolean-and of clauses where each clause is the boolean-or of exactly three literals. For example

$$(x_1 \vee x_2 \vee \bar{x}_3) \wedge (\bar{x}_1 \vee x_3 \vee x_4) \wedge (x_2 \vee \bar{x}_3 \vee \bar{x}_4)$$

is in 3-CNF form. The *3-CNF satisfiability problem* (3SAT) is the problem of determining whether a 3-CNF<sup>15</sup> boolean formula is satisfiable.

**NP-completeness proofs:** Now that we know that 3SAT is NP-complete, we can use this fact to prove that other problems are NP-complete. We will start with the independent set problem.

**Independent Set (IS):** Given an undirected graph  $G = (V, E)$  and an integer  $k$  does  $G$  contain a subset  $V'$  of  $k$  vertices such that no two vertices in  $V'$  are adjacent to one another.

For example, the graph  $G$  shown in Fig. 72 has an independent set (shown with shaded nodes) of size 4, but there is no independent set of size 5. Therefore  $(G, 4) \in \text{IS}$  but  $(G, 5) \notin \text{IS}$ . The independent set problem arises when there is some sort of selection problem, but there are mutual restrictions pairs that cannot both be selected. (For example, you want to invite as many of your friends to your party, but many pairs do not get along, represented by edges between them, and you do not want to invite two enemies.)

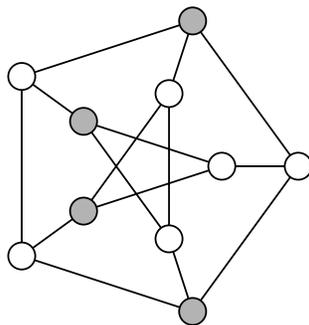


Fig. 72: A graph with an independent set of size  $k = 4$ .

**Claim:** IS is NP-complete.

**Proof:** As with all NP-completeness proofs, there are two parts.

**IS is in NP:** Recall that this means that it is possible to present a certificate that proves when a given instance has an independent set. (If the instance does not have an independent set, then we don't care what the certificate contains.) In this case, the certificate consists of the

<sup>15</sup>Is there something special about the number 3? 1SAT is trivial to solve. 2SAT is trickier, but it can be solved in polynomial time (by reduction to DFS on an appropriate directed graph).  $k$ SAT is NP-complete for any  $k \geq 3$ .

$k$  vertices of  $V'$ . In polynomial time we can verify that, for each pair of vertices  $u, v \in V'$ , there is no edge between them. (In particular, if  $G$  is given as an adjacency matrix, this can be done in  $O(n^2)$  time.)

**IS is NP hard:** It suffices to show that some known NP-complete problem (3SAT) is polynomially reducible to IS, that is,  $3SAT \leq_P IS$ . Let  $F$  be a boolean formula in 3-CNF form. We wish to find a polynomial time computable function  $f$  that maps  $F$  into a input for the IS problem, a graph  $G$  and integer  $k$ . (This is shown schematically in Fig. 73.) That is,  $f(F) = (G, k)$ , such that  $F$  is satisfiable if and only if  $G$  has an independent set of size  $k$ .

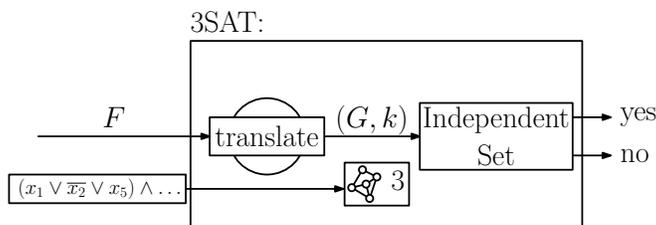


Fig. 73: Reduction of 3-SAT to IS.

This will imply that if we could solve the independent set problem for  $G$  and  $k$  in polynomial time, then we would be able to solve 3SAT in polynomial time. The rest of this section presents this reduction in detail.

Since this is the first nontrivial reduction we will do, let's take a moment to think about the process by which we develop a reduction. An important aspect to reductions is that we *do not* know whether the formula is satisfiable, we *don't know* which variables should be true or false, and we *don't have time* to determine this. (Remember: It is NP-complete!) The translation function  $f$  must operate without knowledge of the answer.

#### What is to be selected?

**3SAT:** Which variables are assigned to be true. Equivalently, which literals are assigned true.

**IS:** Which vertices are to be placed in  $V'$ .

**Idea:** Let's create a vertex in  $G$  for each literal in each clause. A natural approach would be that if a literal is true, then it will correspond to putting the vertex in the independent set. Unfortunately, this will not quite work. Instead, we observe that *at least one* vertex of each clause must be true. We will take *exactly one* such literal from each clause to put into our independent set.

#### Requirements:

**3SAT:** Each clause must contain at least one literal whose value it true.

**IS:**  $V'$  must contain at least  $k$  vertices.

**Idea:** Let's group vertices into groups of three, one group per clause. As mentioned above, exactly one vertex of each group must be in any independent set. We'll set  $k$  equal to the number of clauses to enforce this condition.

#### Restrictions:

**3SAT:** If  $x_i$  is assigned true, then  $\bar{x}_i$  must be false, and vice versa.

**IS:** If  $u$  and  $v$  are adjacent, then both  $u$  and  $v$  cannot be in the independent set.

**Conclusion:** We'll put an edge between two vertices if they correspond to complimentary literals.

In summary, our strategy will be to create clusters of three vertices, one for each literal in each clause. We call these *clause clusters* (see Fig. 74). Since each clause must have at least one true literal, we will model this by forcing the IS algorithm to select one (and only one) vertex per clause cluster. Let's set  $k$  to the number of clauses. But, this does not force us to select one true literal from each clause, since we might take two from some clause cluster and zero from another. To prevent this, we will connect all the vertices within each clause cluster to each other. At most one can be taken to be in any independent set. Since we need to select  $k$  vertices, this will force us to pick exactly one from each cluster.

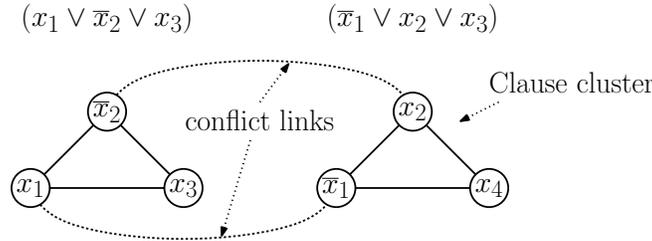


Fig. 74: Clause clusters for the clauses  $(x_1 \vee \bar{x}_2 \vee x_3)$  and  $(\bar{x}_1 \vee x_2 \vee x_5)$ .

To enforce the restriction that only one of  $x$  and  $\bar{x}$  can be set to 1, we create edges between all vertices associated with  $x$  to all vertices associated with  $\bar{x}$ . We call these *conflict links*. A formal description of the reduction is given below. The input is a boolean formula  $F$  in 3-CNF, and the output is a graph  $G$  and integer  $k$ .

3SAT to IS reduction

---

```

 $k \leftarrow$  number of clauses in  $F$ 
for each (clause  $(x_1 \vee x_2 \vee x_3)$  in  $F$ )
    create a clause cluster consisting of three vertices labeled  $x_1, x_2,$  and  $x_3$ 
    create edges  $(x_1, x_2), (x_2, x_3), (x_3, x_1)$  between all pairs of vertices in the cluster
for each (vertex  $x_i$ )
    create edges between  $x_i$  and all its complement vertices  $\bar{x}_i$  (conflict links)
return  $(G, k)$ 

```

---

Given any reasonable encoding of  $F$ , it is an easy programming exercise to create  $G$  in polynomial time. As an example, suppose that we are given the 3-CNF formula:

$$F = (x_1 \vee \bar{x}_2 \vee \bar{x}_3) \wedge (\bar{x}_1 \vee x_2 \vee x_3) \wedge (\bar{x}_1 \vee x_2 \vee \bar{x}_3) \wedge (x_1 \vee \bar{x}_2 \vee x_3).$$

The reduction produces the graph shown in Fig. 75. The clauses clusters appear in clockwise order starting from the top.

In our example, the formula is satisfied by the assignment  $x_1 = 1, x_2 = 1,$  and  $x_3 = 0$ . Note that the literal  $x_1$  satisfies the first and last clauses,  $x_2$  satisfies the second, and  $\bar{x}_3$  satisfies the third. Observe that by selecting the corresponding vertices from the clusters, we obtain an independent set of size  $k = 4$ .

**Correctness:** We'll show that  $F$  is satisfiable if and only if  $G$  has an independent set of size  $k$ .

( $\Rightarrow$ ): If  $F$  is satisfiable, then each of the  $k$  clauses of  $F$  must have at least one true literal. Select such a literal from each clause. Let  $V'$  denote the corresponding vertices from each of the clause clusters (one from each cluster). We claim that  $V'$  is an independent set of size  $k$ . Since there are  $k$  clauses, clearly  $|V'| = k$ . We only take one vertex from each clause cluster, and we cannot take two conflicting literals to be in  $V'$ . For each edge of  $G$ , both of its endpoints cannot be in  $V'$ . Therefore  $V'$  is an independent set of size  $k$ .

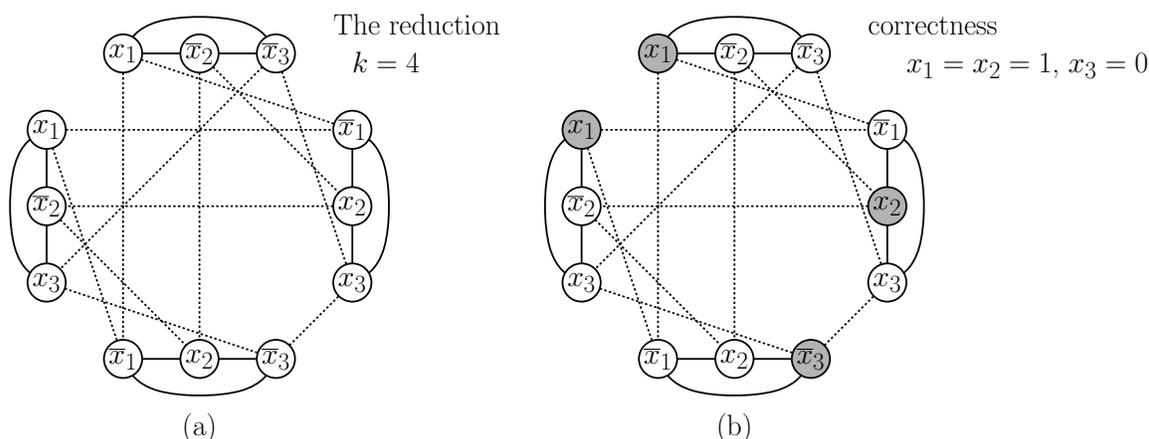


Fig. 75: 3SAT to IS reduction.

( $\Leftarrow$ ): Suppose that  $G$  has an independent set  $V'$  of size  $k$ . We cannot select two vertices from a clause cluster, and since there are  $k$  clusters,  $V'$  has exactly one vertex from each clause cluster. Note that if a vertex labeled  $x$  is in  $V'$  then the adjacent vertex  $\bar{x}$  cannot also be in  $V'$ . Therefore, there exists an assignment in which every literal corresponding to a vertex appearing in  $V'$  is set to 1. Such an assignment satisfies one literal in each clause, and therefore the entire formula is satisfied.

Let us emphasize a few things about this reduction:

- Every NP-complete problem has three similar elements: (a) something is being selected, (b) something is forcing us to select a sufficient number of such things (requirements), and (c) something is limiting our ability to select these things (restrictions). A reduction's job is to determine how to map these similar elements to each other.
- Our reduction did not attempt to solve the 3SAT problem. (As a sign of this, observe that whatever we did for one literal, we did for all.) Remember this rule! If your reduction treats some entities different other, based on what you think the final answer may be, you are very likely making a mistake. Remember, these problems are NP-complete!

## Lecture 21: NP-Completeness: Clique, Vertex Cover, and Dominating Set

**Recap:** Last time we gave a reduction from 3SAT (satisfiability of boolean formulas in 3-CNF form) to IS (independent set in graphs). Today we give a few more examples of reductions. Recall that to show that a decision problem (language)  $L$  is NP-complete we need to show:

- $L \in \text{NP}$ . (That is, given an input and an appropriate certificate, we can guess the solution and verify whether the input is in the language), and
- $L$  is NP-hard, which we can show by giving a reduction from some known NP-complete problem  $L'$  to  $L$ , that is,  $L' \leq_P L$ . (That is, there is a polynomial time function that transforms an instance  $L'$  into an equivalent instance of  $L$  for the other problem).

**Some Easy Reductions:** Next, let us consider some closely related NP-complete problems:

**Clique (CLIQUE):** The *clique problem* is: given an undirected graph  $G = (V, E)$  and an integer  $k$ , does  $G$  have a subset  $V'$  of  $k$  vertices such that for each distinct  $u, v \in V'$ ,  $\{u, v\} \in E$ . In other words, does  $G$  have a  $k$  vertex subset whose induced subgraph is complete.

**Vertex Cover (VC):** A *vertex cover* in an undirected graph  $G = (V, E)$  is a subset of vertices  $V' \subseteq V$  such that every edge in  $G$  has at least one endpoint in  $V'$ . The *vertex cover problem* (VC) is: given an undirected graph  $G$  and an integer  $k$ , does  $G$  have a vertex cover of size  $k$ ?

**Dominating Set (DS):** A *dominating set* in a graph  $G = (V, E)$  is a subset of vertices  $V'$  such that every vertex in the graph is either in  $V'$  or is adjacent to some vertex in  $V'$ . The *dominating set problem* (DS) is: given a graph  $G = (V, E)$  and an integer  $k$ , does  $G$  have a dominating set of size  $k$ ?

Don't confuse the clique (CLIQUE) problem with the clique-cover (CC) problem that we discussed in an earlier lecture. The clique problem seeks to find a single clique of size  $k$ , and the clique-cover problem seeks to partition the vertices into  $k$  groups, each of which is a clique.

We have discussed the facts that cliques are of interest in applications dealing with clustering. The vertex cover problem arises in various servicing applications. For example, you have a compute network and a program that checks the integrity of the communication links. To save the space of installing the program on every computer in the network, it suffices to install it on all the computers forming a vertex cover. From these nodes all the links can be tested. Dominating set is useful in facility location problems. For example, suppose we want to select where to place a set of fire stations such that every house in the city is within two minutes of the nearest fire station. We create a graph in which two locations are adjacent if they are within two minutes of each other. A minimum sized dominating set will be a minimum set of locations such that every other location is reachable within two minutes from one of these sites.

The CLIQUE problem is obviously closely related to the independent set problem (IS): Given a graph  $G$  does it have a  $k$  vertex subset that is completely *disconnected*. It is not quite as clear that the vertex cover problem is related. However, the following lemma makes this connection clear as well (see Fig. 76). Given a graph  $G$ , recall that  $\bar{G}$  is the *complement graph* where edges and non-edges are reverse. Also, recall that  $A \setminus B$  denotes set resulting by removing the elements of  $B$  from  $A$ .

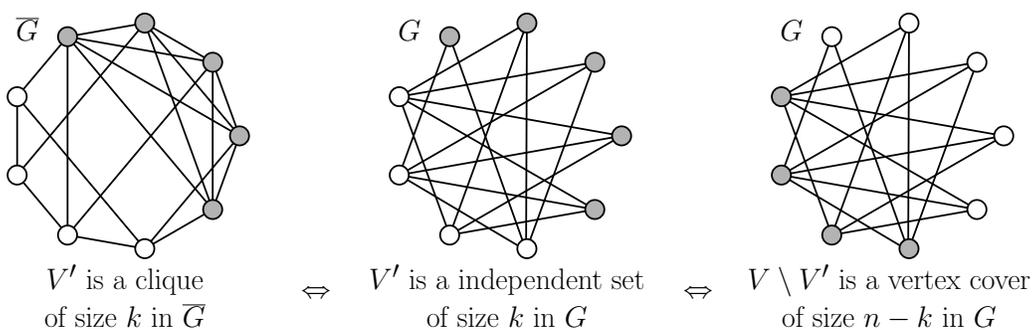


Fig. 76: Clique, Independent set, and Vertex Cover.

**Lemma:** Given an undirected graph  $G = (V, E)$  with  $n$  vertices and a subset  $V' \subseteq V$  of size  $k$ . The following are equivalent:

- (i)  $V'$  is a clique of size  $k$  for the complement,  $\bar{G}$
- (ii)  $V'$  is an independent set of size  $k$  for  $G$

(iii)  $V \setminus V'$  is a vertex cover of size  $n - k$  for  $G$ , (where  $n = |V|$ )

**Proof:**

(i)  $\Rightarrow$  (ii): If  $V'$  is a clique for  $\overline{G}$ , then for each  $u, v \in V'$ ,  $\{u, v\}$  is an edge of  $\overline{G}$  implying that  $\{u, v\}$  is not an edge of  $G$ , implying that  $V'$  is an independent set for  $G$ .

(ii)  $\Rightarrow$  (iii): If  $V'$  is an independent set for  $G$ , then for each  $u, v \in V'$ ,  $\{u, v\}$  is not an edge of  $G$ , implying that every edge in  $G$  is incident to a vertex in  $V \setminus V'$ , implying that  $V \setminus V'$  is a vertex cover for  $G$ .

(iii)  $\Rightarrow$  (i): If  $V \setminus V'$  is a vertex cover for  $G$ , then for any  $u, v \in V'$  there is no edge  $\{u, v\}$  in  $G$ , implying that there is an edge  $\{u, v\}$  in  $\overline{G}$ , implying that  $V'$  is a clique in  $\overline{G}$ .

Thus, if we had an algorithm for solving any one of these problems, we could easily translate it into an algorithm for the others. In particular, we have the following.

**Theorem:** CLIQUE is NP-complete.

**CLIQUE  $\in$  NP:** Given an instance  $(G, k)$  for CLIQUE, we guess the  $k$  vertices that will form the clique. (These vertices form the *certificate*.) We can easily verify in polynomial time that all pairs of vertices in the set are adjacent (e.g., by inspection of  $O(k^2)$  entries of the adjacency matrix). If so, we output “yes” and otherwise “no”. (If the graph has a CLIQUE of size  $k$ , one of these guesses will work, and so we will correctly classify  $G$  as having a clique of size  $k$ . If not, all guesses will fail, and we will correctly classify  $G$  as not having such a clique.)

**IS  $\leq_P$  CLIQUE:** We want to show that given an instance of the IS problem  $(G, k)$ , we can produce an equivalent instance of the CLIQUE problem in polynomial time. The reduction function  $f$  inputs  $G$  and  $k$ , and outputs the pair  $(\overline{G}, k)$ . Clearly this can be done in polynomial time. By the above lemma, this instance is equivalent.

**Theorem:** VC is NP-complete.

**VC  $\in$  NP:** Given an instance  $(G, k)$  for VC, we guess the  $k$  vertices that will form the vertex cover. (Again, these vertices form the *certificate*.) We then verify that these vertices form a vertex cover, by checking that every edge of  $G$  is incident to one of these vertices. If so, we output “yes” and otherwise “no”. (Again, if  $G$  has a vertex cover of size  $k$ , one of these guesses will work, and we correctly classify  $G$  as having a vertex cover of size  $k$ . Otherwise, all fail and we classify  $G$  as not having such a vertex cover.)

**IS  $\leq_P$  VC:** We want to show that given an instance of the IS problem  $(G, k)$ , we can produce an equivalent instance of the VC problem in polynomial time. The reduction function  $f$  inputs  $G$  and  $k$ , computes the number of vertices,  $n$ , and then outputs  $(G, n - k)$ . Clearly this can be done in polynomial time. By the lemma above, these instances are equivalent.

**Note:** Note that in each of the above reductions, the reduction function did not know whether  $G$  has an independent set or not. It must run in polynomial time, and IS is an NP-complete problem. So it does not have time to determine whether  $G$  has an independent set or which vertices are in the set.

**Dominating Set:** As with vertex cover, dominating set is an example of a graph covering problem. Here the condition is a little different, each *vertex* is *adjacent* to at least one member of the dominating set, as opposed to each *edge* being *incident* to at least one member of the vertex cover. Obviously, if  $G$  is connected and has a vertex cover of size  $k$ , then it has a dominating set of size  $k$  (the same set of vertices), but the converse is not necessarily true. However, the similarity suggests that if VC is NP-complete, then DS is likely to be NP-complete as well. We will show this fact next.

As usual the proof has two parts. First we show that DS  $\in$  NP. The certificate just consists of the subset  $V'$  in the dominating set. In polynomial time we can determine whether every vertex is in  $V'$  or is adjacent to a vertex in  $V'$ .

**Vertex Cover to Dominating Set:** Next, we show that a known NP-complete problem is reducible to dominating set. We choose vertex cover and show that  $VC \leq_P DS$ . We want a polynomial time function, which given an instance of the vertex cover problem  $(G, k)$ , produces an instance  $(G', k')$  of the dominating set problem, such that  $G$  has a vertex cover of size  $k$  if and only if  $G'$  has a dominating set of size  $k'$ .

How to we translate between these problems? The key difference is the condition. In VC: “every edge is incident to a vertex in  $V'$ ”. In DS: “every vertex is either in  $V'$  or is adjacent to a vertex in  $V'$ ”. Thus the translation must somehow map the notion of “incident” to “adjacent”. Because incidence is a property of edges, and adjacency is a property of vertices, this suggests that the reduction function maps edges of  $G$  into vertices in  $G'$ , such that an incident edge in  $G$  is mapped to an adjacent vertex in  $G'$ .

This suggests the following idea (which does not quite work). We will insert a vertex into the middle of each edge of the graph. In other words, for each edge  $\{u, v\}$ , we will create a new *special vertex*, called  $w_{uv}$ , and replace the edge  $\{u, v\}$  with the two edges  $\{u, w_{uv}\}$  and  $\{v, w_{uv}\}$ . The fact that  $u$  was incident to edge  $\{u, v\}$  has now been replaced with the fact that  $u$  is adjacent to the corresponding vertex  $w_{uv}$ . We still need to dominate the neighbor  $v$ . To do this, we will leave the edge  $\{u, v\}$  in the graph as well. Let  $G'$  be the resulting graph.

This is still not quite correct though. Define an *isolated vertex* to be one that is incident to no edges. If  $u$  is isolated it can only be dominated if it is included in the dominating set. Since it is not incident to any edges, it does not need to be in the vertex cover. Let  $V_I$  denote the isolated vertices in  $G$ , and let  $n_I$  denote the number of isolated vertices. The number of vertices to request for the dominating set will be  $k' = k + n_I$ . Okay, we are now ready to state the result and prove it.

**Theorem:** DS is NP-complete.

**DS  $\in$  NP:** Given an instance  $(G, k)$  for DS, we guess the certificate, which consists of the  $k$  vertices that will form the dominating set. We then verify that these vertices form a dominating set, by checking that every vertex of  $G$  is either in this set or is adjacent to a vertex in this set. If so, we output “yes” and otherwise “no”. (Again, if  $G$  has a dominating set of size  $k$ , one of these guesses will work, and we correctly classify  $G$  as having a dominating set of size  $k$ . Otherwise, all fail and we classify  $G$  as not having such a dominating set.)

**VC  $\leq_P$  DS:** We want to show that given an instance of the VC problem  $(G, k)$ , we can produce an equivalent instance of the DS problem in polynomial time. We create a graph  $G'$  as follows. Initially  $G' = G$ . For each edge  $\{u, v\}$  in  $G$  we create a new vertex  $w_{uv}$  in  $G'$  and add edges  $\{u, w_{uv}\}$  and  $\{v, w_{uv}\}$  in  $G'$ . Let  $I$  denote the number of isolated vertices and set  $k' = k + n_I$ . Output  $(G', k')$ . This reduction illustrated in Fig. 77. Note that every step can be performed in polynomial time.

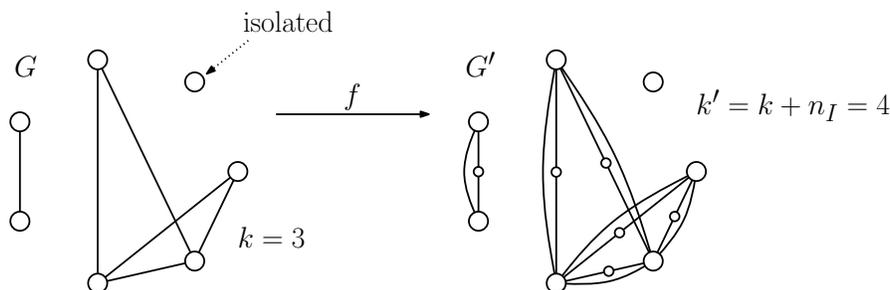


Fig. 77: Dominating set reduction with  $k = 3$  and one isolated vertex.

To establish the correctness of the reduction, we need to show that  $G$  has a vertex cover of size  $k$  if and only if  $G'$  has a dominating set of size  $k'$ .

( $\Rightarrow$ ) First we argue that if  $V'$  is a vertex cover for  $G$ , then  $V'' = V' \cup V_I$  is a dominating set for  $G'$ . Observe that

$$|V''| = |V' \cup V_I| \leq k + n_I = k'.$$

Note that  $|V' \cup V_I|$  might be of size less than  $k + n_I$ , if there are any isolated vertices in  $V'$ . If so, we can add any vertices we like to make the size equal to  $k'$ .

To see that  $V''$  is a dominating set, first observe that all the isolated vertices are in  $V''$  and so they are dominated. Second, each of the special vertices  $w_{uv}$  in  $G'$  corresponds to an edge  $\{u, v\}$  in  $G$  implying that either  $u$  or  $v$  is in the vertex cover  $V'$ . Thus  $w_{uv}$  is dominated by the same vertex in  $V''$ . Finally, each of the nonisolated original vertices  $v$  is incident to at least one edge in  $G$ , and hence either it is in  $V'$  or else all of its neighbors are in  $V'$ . In either case,  $v$  is either in  $V''$  or adjacent to a vertex in  $V''$ . This is shown in the top part of the following Fig. 78.

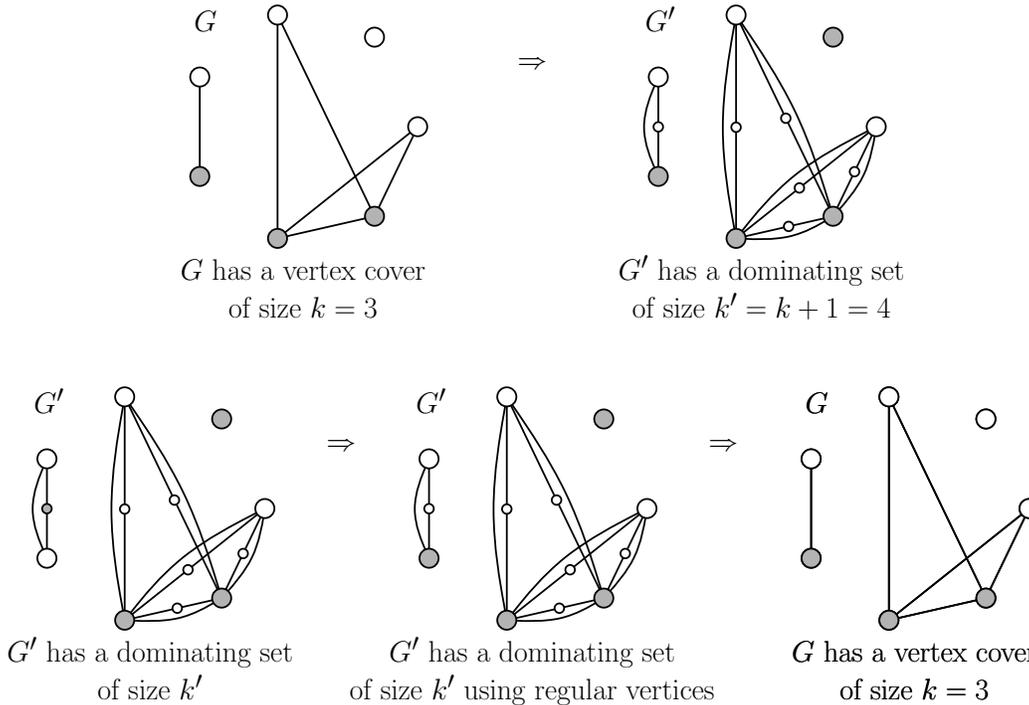


Fig. 78: Correctness of the VC to DS reduction (where  $k = 3$  and  $I = 1$ ).

( $\Leftarrow$ ) Conversely, we claim that if  $G'$  has a dominating set  $V''$  of size  $k' = k + n_I$  then  $G$  has a vertex cover  $V'$  of size  $k$ . Note that all  $n_I$  isolated vertices of  $G'$  must be in the dominating set. First, let  $V''' = V'' \setminus V_I$  be the remaining  $k$  vertices. We might try to claim something like:  $V'''$  is a vertex cover for  $G$ . But this will not necessarily work, because  $V'''$  may have vertices that are not part of the original graph  $G$ .

However, we claim that we never need to use any of the newly created special vertices in  $V'''$ . In particular, if some vertex  $w_{uv} \in V'''$ , then modify  $V'''$  by replacing  $w_{uv}$  with  $u$ . (We could have just as easily replaced it with  $v$ .) Observe that the vertex  $w_{uv}$  is adjacent to only  $u$  and  $v$ , so it dominates itself and these other two vertices. By using  $u$  instead, we still dominate  $u$ ,  $v$ , and  $w_{uv}$  (because  $u$  has edges going to  $v$  and  $w_{uv}$ ). Thus by replacing  $w_{u,v}$  with  $u$  we dominate the same vertices (and potentially more). Let  $V'$  denote the resulting set after this modification. (This is shown in the lower middle part of Fig 78.)

We claim that  $V'$  is a vertex cover for  $G$ . If, to the contrary there were an edge  $\{u, v\}$  of  $G$  that

was not covered (neither  $u$  nor  $v$  was in  $V'$ ) then the special vertex  $w_{uv}$  would not be adjacent to any vertex of  $V''$  in  $G'$ , contradicting the hypothesis that  $V''$  was a dominating set for  $G'$ .

Whew! So, this completes the proof of the correctness of the reduction, so we conclude that DS is NP-complete.

## Lecture 22: Approximation Algorithms: Vertex Cover and TSP

**Coping with NP-completeness:** With NP-completeness we have seen that there are many important optimization problems that are likely to be quite hard to solve exactly. Since these are important problems, we cannot simply give up at this point, since people do need solutions to these problems. How do we cope with NP-completeness:

**Brute-force search:** This is usually only a viable option for small input sizes (e.g.,  $n \leq 20$ ).

**Heuristics:** This is a strategy for producing a valid solution, but may be there no guarantee on how close it is to optimal.

**General Search Algorithms:** There are a number of very powerful techniques for solving general combinatorial optimization problems. These go under various names such as *branch-and-bound*, *Metropolis-Hastings*, *simulated annealing*, and *genetic algorithms*. The performance of these approaches varies considerably from one problem to problem and instance to instance. But in some cases they can perform quite well.

**Approximation Algorithms:** This is an algorithm that runs in polynomial time (ideally), and produces a solution that is guaranteed to be within some factor of the optimum solution.

**Approximation Bounds:** Most NP-complete problems have been stated as decision problems for theoretical reasons. However underlying most of these problems is a natural optimization problem. For example, find the vertex cover of *minimum* size, the independent set of *maximum* size, color a graph with the *minimum* number of colors, or find the traveling salesman tour with the *shortest* cost. An approximation algorithm is one that returns a valid (or feasible) answer, but not necessarily one of the optimal size/weight/cost.

How do we measure how good an approximation algorithm is? We define the *approximation ratio* of an approximation algorithm as the worst-case ratio between the answer produced by the approximation algorithm and the optimum solution. For minimization problems (where the approximate solution will usually be larger than the optimum) this is expressed as approx/opt, and for maximization problems (where the approximate solution will be smaller than the optimum) this is expressed as opt/approx. Thus, in either case, the ratio will be at least 1, and the smaller the better.

Although NP-complete problems are equivalent with respect to whether they can be solved exactly in polynomial time in the worst case, their approximability varies considerably. Here are some possibilities:

- Some NP-complete problems are *inapproximable* in the sense no polynomial time algorithm achieves a ratio bound smaller than  $\infty$  unless  $P = NP$ . (For example, there is no bound on the approximability of either independent set or graph coloring, unless  $P = NP$ .)
- Some NP-complete problems can be approximated, but the ratio bound is a *function of  $n$* . (For example, the set cover problem, can be approximated to within a factor of  $O(\log n)$ , and this is believed to be the best possible.)
- Some NP-complete problems can be approximated and the ratio bound is a *constant*. (For example, we will see below that Vertex Cover can be approximated to within a factor of 2.)

- Some NP-complete problems can be approximated *arbitrarily well*. In particular, the user provides a parameter  $\varepsilon > 0$  and the algorithm achieves a ratio bound of  $(1 + \varepsilon)$ . Of course, as  $\varepsilon$  approaches 0 the algorithm's running time gets worse. If such an algorithm runs in polynomial time for any fixed  $\varepsilon$ , it is called a *polynomial time approximation scheme*, or PTAS. (For example, there is a PTAS for the Euclidean traveling salesman problem.)

**Vertex Cover:** We begin by showing that there is an approximation algorithm for vertex cover with a ratio bound of 2, that is, this algorithm will be guaranteed to find a vertex cover whose size is at most twice that of the optimum. Recall that a vertex cover is a subset of vertices such that every edge in the graph is incident to at least one of these vertices. The *vertex cover optimization problem* is to find a vertex cover of minimum size (See Fig. 79).

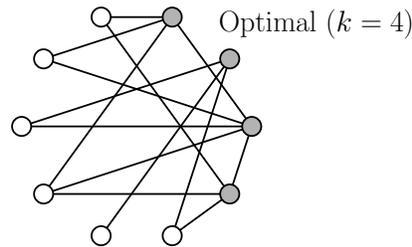


Fig. 79: Vertex cover (optimal solution).

How does one go about finding an approximation algorithm. The first approach is to try something that seems like a “reasonably” good strategy, a *heuristic*. It turns out that many simple heuristics, when not optimal, can often be proved to be close to optimal.

Here is a very simple algorithm, that guarantees an approximation within a factor of 2 for the vertex cover problem. First, we compute a maximal matching for the graph. Recall that a *matching* for a graph  $G = (V, E)$  is a subset of edges  $M \subseteq E$  such that each vertex is incident to at most one edge of  $M$ . A matching is *maximal* (not necessarily *maximum*) if it is not possible to add any more edges to the matching.

To compute a maximal matching, repeatedly find an unmarked edge  $(u, v)$ , add it to the matching, and then mark all the edges incident to either  $u$  or  $v$ . Clearly, the result is a matching, and it cannot be extended to a larger matching, so it is maximal. To obtain a vertex cover, observe that for every edge  $(u, v)$  an *any* matching, either  $u$  or  $v$  must be included in *any* vertex cover. Thus, a simple (and rather dumb) solution is to add *both* of them to the vertex cover. (The algorithm shown in the following code fragment, and it is illustrated in Fig. 80.)

VC approximation via Maximal Matching

```

VC-approx(V,E) {
  C = emptyset
  mark all edges as uncovered
  while (there exists an unmarked edge (u,v)) do { // (u,v) is in a maximal matching
    add both u and v to C // add u and v to the cover
    mark all the edges incident to u and to v // mark all incident edges
  }
  return C as the approximate vertex cover
}

```

**Claim:** Above algorithm achieves an approximation ratio of 2.

**Proof:** Let  $M$  be the set of edges  $(u, v)$  identified by the algorithm, and let  $C$  be the set of endpoints of  $M$ . Let  $k = |M|$ . Since  $M$  is a matching,  $|C| = 2k$ . Because  $M$  is maximal, every edge in  $G$  is incident to some endpoint in  $C$ , implying that  $C$  is a vertex cover of size  $2k$ . By definition, any vertex cover must contain at least one endpoint from each of  $M$ 's edges. Since no two edges of  $M$  share the same endpoint, any vertex cover has size at least  $|M| = k$ . Thus,  $C$  is within a factor of 2 of the optimum size.

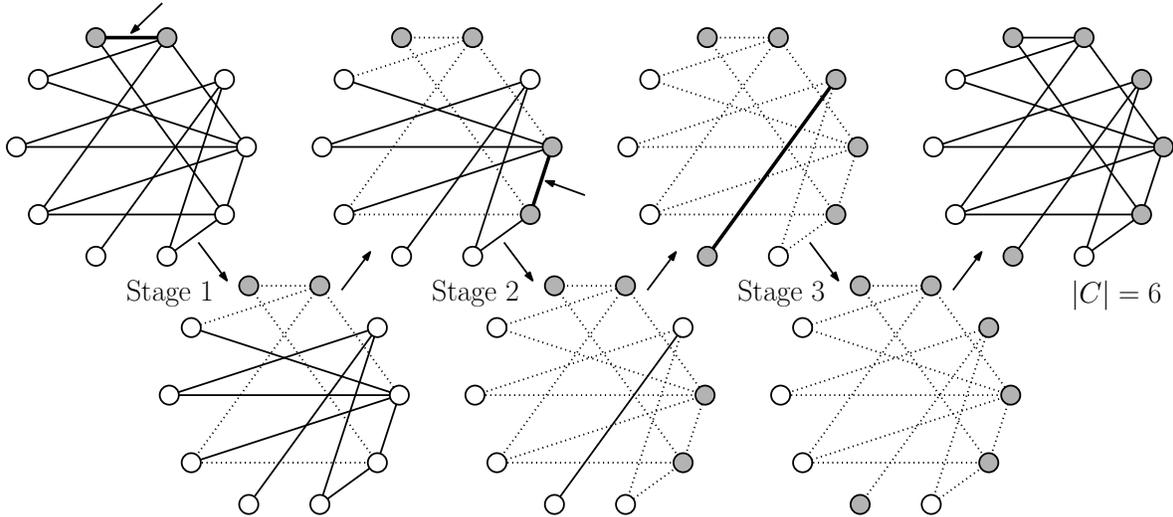


Fig. 80: The matching-based heuristic for vertex cover.

This proof illustrates one of the main features of the analysis of any approximation algorithm. Namely, that we need some way of finding a bound on the optimal solution. (For minimization problems we want a lower bound, for maximization problems an upper bound.) The bound should be related to something that we can compute in polynomial time. In this case, the bound is related to the set of edges  $A$ , which form a maximal independent set of edges.

**Traveling Salesman with Triangle Inequality:** In the Traveling Salesperson Problem (TSP) we are given a complete undirected graph with nonnegative edge weights, and we want to find a cycle that visits all vertices and is of minimum cost. Let  $w(u, v)$  denote the weight on edge  $(u, v)$ . Given a set of edges  $A$  forming a tour we define  $W(A)$  to be the sum of edge weights in  $A$ .

For many of the applications of TSP, the edge weights satisfy a natural property called the *triangle inequality*. Intuitively, this says that the direct path from  $u$  to  $x$ , is never longer than an indirect path. More formally, for all  $u, v, x \in V$

$$w(u, v) \leq w(u, x) + w(x, v).$$

There are many examples of graphs that satisfy the triangle inequality. For example, given any weighted graph, if we define  $w(u, v)$  to be the shortest path length between  $u$  and  $v$ , then it will satisfy the triangle inequality. Another example is if we are given a set of points in the plane, and define a complete graph on these points, where  $w(u, v)$  is defined to be the Euclidean distance between these points, then the triangle inequality is also satisfied.

When the underlying cost function satisfies the triangle inequality there is an approximation algorithm for TSP with a ratio-bound of 2. (In fact, there is a slightly more complex version of this algorithm that has a ratio bound of 1.5, but we will not discuss it.) Thus, although this algorithm does not

produce an optimal tour, the tour that it produces cannot be worse than twice the cost of the optimal tour.

The key insight is to observe that a TSP with one edge removed is just a spanning tree. However it is not necessarily a minimum spanning tree. Therefore, the cost of the minimum TSP tour is at least as large as the cost of the MST. We can compute MST's efficiently, using, for example, Kruskal's algorithm. If we can find some way to convert the MST into a TSP tour while increasing its cost by at most a constant factor, then we will have an approximation for TSP. We shall see that if the edge weights satisfy the triangle inequality, then this is possible.

Here is how the algorithm works. Given any free tree there is a tour of the tree called a *twice around tour* that traverses the edges of the tree twice, once in each direction (see Fig. 81).

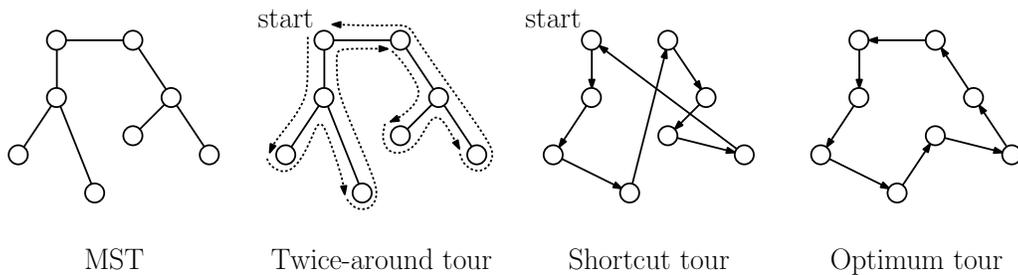


Fig. 81: TSP Approximation.

This path is not simple because it revisits vertices, but we can make it simple by *short-cutting*, that is, we skip over previously visited vertices. Notice that the final order in which vertices are visited using the short-cuts is exactly the same as a preorder traversal of the MST. (In fact, any subsequence of the twice-around tour which visits each vertex exactly once will suffice.) The triangle inequality assures us that the path length will not increase when we take short-cuts.

TSP Approximation

---

```

approx-TSP(G=(V,E)) {
  T = minimum spanning tree for G
  r = any vertex
  H = list of vertices visited by a preorder walk of T starting at r
  return H
}

```

---

**Claim:** Approx-TSP achieves a approximation ratio of 2.

**Proof:** Let  $H$  denote the tour produced by this algorithm and let  $H^*$  be the optimum tour. Let  $T$  be the minimum spanning tree. As we said before, since we can remove any edge of  $H^*$  resulting in a spanning tree, and since  $T$  is the minimum cost spanning tree we have

$$W(T) \leq W(H^*).$$

Now observe that the twice-around tour of  $T$  has cost  $2 \cdot W(T)$ , since every edge in  $T$  is hit twice. By the triangle inequality, whenever we short-cut an edge of  $T$  to form  $H$  we do not increase the cost of the tour, and so we have

$$W(H) \leq 2 \cdot W(T).$$

Combining these we have

$$W(H) \leq 2 \cdot W(T) \leq 2 \cdot W(H^*) \Rightarrow \frac{W(H)}{W(H^*)} \leq 2,$$

as desired.

## Lecture 23: Max Dominance

**Max-Dominance:** Let us consider a natural problem, that arises in a number of financial applications. Suppose you are considering buying a new car. You would like a sporty car with fast acceleration, but you are frugal and also want a car that has good gas mileage. You study various models of cars and for each you record the acceleration and mileage as points on an  $(x, y)$  plot. It is not surprising that some cars have excellent mileage but poor acceleration and vice versa. One thing you can determine is that if model  $A$  has both lower mileage and lower acceleration than model  $B$ , you are not interested in model  $A$ . Formally, we say that a point  $A = (A.x, A.y)$  is *dominated* by  $B = (B.x, B.y)$  if  $A.x < B.x$  and  $A.y < B.y$ . The points that are not dominated by any other point are said to be the *dominant points*. (In economics, this is related to the concept of *Pareto optima*.) The *max dominance problem* is, given a set  $P$  of  $n$  points (see Fig. (a)), compute the set of dominant points from the set (see Fig. (b)).

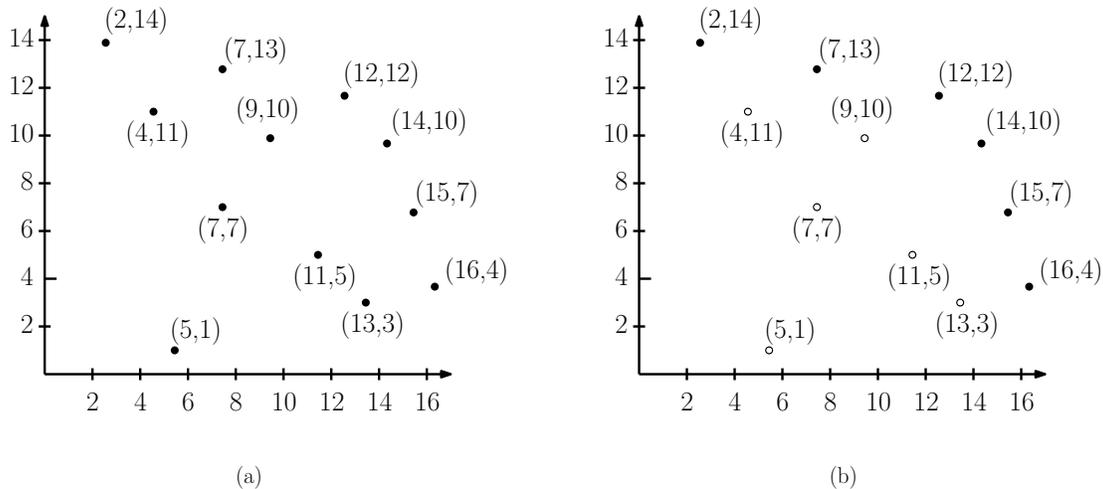


Fig. 82: Max-Dominance.

There is an obvious brute-force algorithm that runs in  $O(n^2)$  time, which operates by comparing all pairs of points. The question we consider here is whether there is an approach that is significantly better.

**A Major Improvement:** The problem with the previous algorithm is that, even though we have cut the number of comparisons roughly in half, each point is still making lots of comparisons. Can we save time by making only one comparison for each point? The inner while loop is testing to see whether *any* point that follows  $P[i]$  in the sorted list has a larger  $y$ -coordinate. This suggests, that if we knew which point among  $P[i+1, \dots, n]$  had the maximum  $y$ -coordinate, we could just test against that point.

How can we do this? Here is a simple observation. For any set of points, the point with the maximum  $y$ -coordinate is the maximal point with the smallest  $x$ -coordinate. This suggests that we can sweep the points backwards, from right to left. We keep track of the index  $j$  of the most recently seen maximal point. (Initially the rightmost point is maximal.) When we encounter the point  $P[i]$ , it is maximal if and only if  $P[i].y \geq P[j].y$ . This suggests the following algorithm.

The running time of the for-loop is obviously  $O(n)$ , because there is just a single loop that is executed  $n - 1$  times, and the code inside takes constant time. The total running time is dominated by the  $O(n \log n)$  sorting time, for a total of  $O(n \log n)$  time.

---

```

MaxDom3(P, n) {
  Sort P in ascending order by x-coordinate;
  output P[n];                // last point is always maximal
  j = n;
  for i = n-1 downto 1 {
    if (P[i].y >= P[j].y) {   // is P[i] maximal?
      output P[i];           // yes..output it
      j = i;                 // P[i] has the largest y so far
    }
  }
}

```

---

How much of an improvement is this? Probably the most accurate way to find out would be to code the two up, and compare their running times. But just to get a feeling, let's look at the ratio of the running times, ignoring constant factors:

$$\frac{n^2}{n \lg n} = \frac{n}{\lg n}.$$

(I use the notation  $\lg n$  to denote the logarithm base 2,  $\ln n$  to denote the natural logarithm (base  $e$ ) and  $\log n$  when I do not care about the base. Note that a change in base only affects the value of a logarithm function by a constant amount, so inside of  $O$ -notation, we will usually just write  $\log n$ .)

For relatively small values of  $n$  (e.g., less than 100), both algorithms are probably running fast enough that the difference will be practically negligible. (Rule 1 of algorithm optimization: Don't optimize code that is already fast enough.) On larger inputs, say,  $n = 1,000$ , the ratio of  $n$  to  $\log n$  is about  $1000/10 = 100$ , so there is a 100-to-1 ratio in running times. Of course, we would need to factor in constant factors, but since we are not using any really complex data structures, it is hard to imagine that the constant factors will differ by more than, say, 10. For even larger inputs, say,  $n = 1,000,000$ , we are looking at a ratio of roughly  $1,000,000/20 = 50,000$ . This is quite a significant difference, irrespective of the constant factors.

**Divide and Conquer Approach:** One problem with the previous algorithm is that it relies on sorting. This is nice and clean (since it is usually easy to get good code for sorting without troubling yourself to write your own). However, if you really wanted to squeeze the most efficiency out of your code, you might consider whether you can solve this problem without invoking a sorting algorithm.

One of the basic maxims of algorithm design is to first approach any problem using one of the standard algorithm design paradigms, e.g. divide and conquer, dynamic programming, greedy algorithms, depth-first search. We will talk more about these methods as the semester continues. For this problem, divide-and-conquer is a natural method to choose. What is this paradigm?

**Divide:** Divide the problem into two subproblems (ideally of approximately equal sizes),

**Conquer:** Solve each subproblem recursively, and

**Combine:** Combine the solutions to the two subproblems into a global solution.

How shall we divide the problem? I can think of a couple of ways. One is similar to how *MergeSort* operates. Just take the array of points  $P[1..n]$ , and split into two subarrays of equal size  $P[1..n/2]$  and  $P[n/2+1..n]$ . Because we do not sort the points, there is no particular relationship between the points in one side of the list from the other.

Another approach, which is more reminiscent of *QuickSort* is to select a random element from the list, called a *pivot*,  $x = P[r]$ , where  $r$  is a random integer in the range from 1 to  $n$ , and then partition the

list into two sublists, those elements whose  $x$ -coordinates are less than or equal to  $x$  and those that greater than  $x$ . This will not be guaranteed to split the list into two equal parts, but on average it can be shown that it does a pretty good job.

Let's consider the first method. (The quicksort method will also work, but leads to a tougher analysis.) Here is more concrete outline. We will describe the algorithm at a very high level. The input will be a point array, and a point array will be returned. The key ingredient is a function that takes the maxima of two sets, and merges them into an overall set of maxima.

```

MaxDom4(P, n) {
    if (n == 1) return {P[1]};           // one point is trivially maximal
    m = n/2;                             // midpoint of list
    M1 = MaxDom4(P[1..m], m);           // solve for first half
    M2 = MaxDom4(P[m+1..n], n-m);       // solve for second half
    return MaxMerge(M1, M2);           // merge the results
}

```

---

The general process is illustrated in Fig. .

The main question is how the procedure `MaxMerge()` is implemented, because it does all the work. Let us assume that it returns a list of points in *sorted order* according to  $x$ -coordinates of the maximal points. Observe that if a point is to be maximal overall, then it must be maximal in one of the two sublists. However, just because a point is maximal in some list, does not imply that it is globally maximal. (Consider point (7, 10) in the example.) However, if it dominates all the points of the other sublist, then we can assert that it is maximal.

I will describe the procedure at a very high level. It operates by walking through each of the two sorted lists of maximal points. It maintains two pointers, one pointing to the next unprocessed item in each list. Think of these as *fingers*. Take the finger pointing to the point with the smaller  $x$ -coordinate. If its  $y$ -coordinate is larger than the  $y$ -coordinate of the point under the other finger, then this point is maximal, and is copied to the next position of the result list. Otherwise it is not copied. In either case, we move to the next point in the same list, and repeat the process. The result list is returned.

The details will be left as an exercise. Observe that because we spend a constant amount of time processing each point (either copying it to the result list or skipping over it) the total execution time of this procedure is  $O(n)$ .

**Recurrences:** How do we analyze recursive procedures like this one? If there is a simple pattern to the sizes of the recursive calls, then the best way is usually by setting up a *recurrence*, that is, a function which is defined recursively in terms of itself.

We break the problem into two subproblems of size roughly  $n/2$  (we will say exactly  $n/2$  for simplicity), and the additional overhead of merging the solutions is  $O(n)$ . We will ignore constant factors, writing  $O(n)$  just as  $n$ , giving:

$$\begin{aligned}
 T(n) &= 1 && \text{if } n = 1, \\
 T(n) &= 2T(n/2) + n && \text{if } n > 1.
 \end{aligned}$$

**Solving Recurrences by The Master Theorem:** There are a number of methods for solving the sort of recurrences that show up in divide-and-conquer algorithms. The easiest method is to apply the *Master Theorem* that is given in CLRS. Here is a slightly more restrictive version, but adequate for a lot of instances. See CLRS for the more complete version of the Master Theorem and its proof.

**Theorem:** (Simplified Master Theorem) Let  $a \geq 1$ ,  $b > 1$  be constants and let  $T(n)$  be the recurrence

$$T(n) = aT(n/b) + cn^k,$$

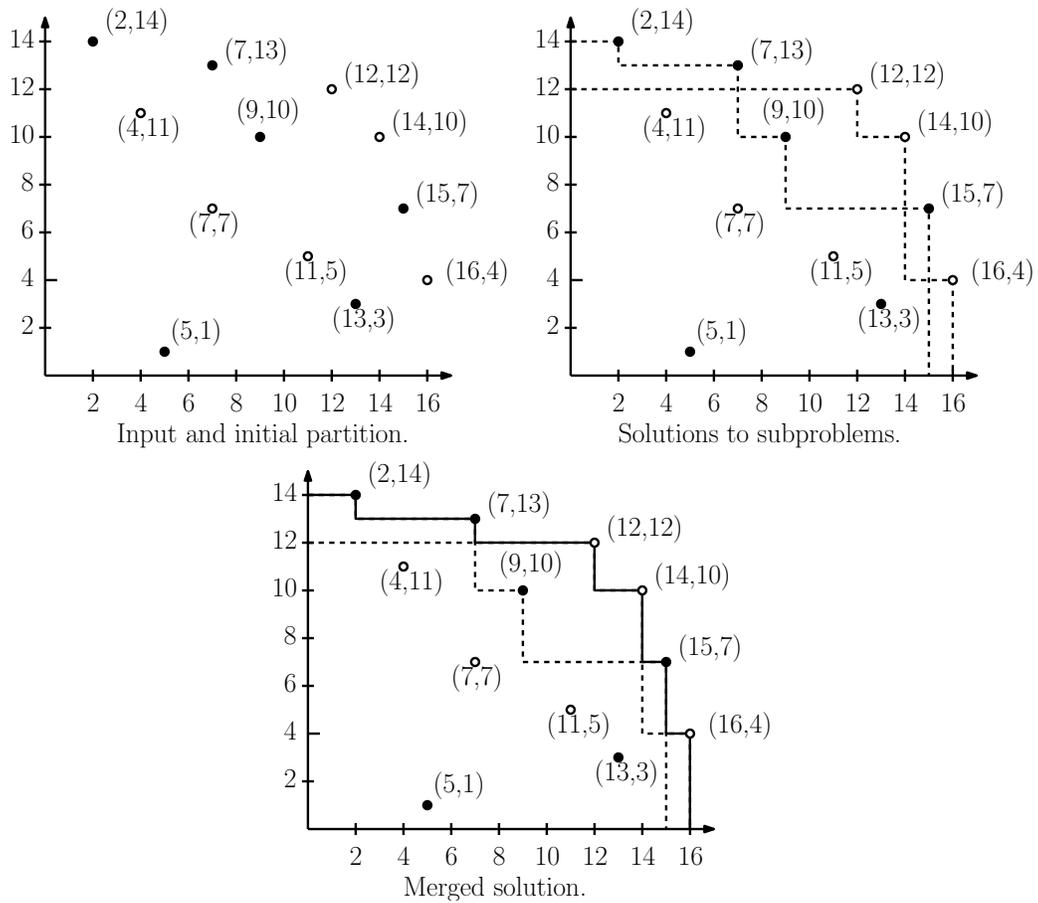


Fig. 83: Divide and conquer approach.

defined for  $n \geq 0$ .

**Case (1):**  $a > b^k$  then  $T(n)$  is  $\Theta(n^{\log_b a})$ .

**Case (2):**  $a = b^k$  then  $T(n)$  is  $\Theta(n^k \log n)$ .

**Case (3):**  $a < b^k$  then  $T(n)$  is  $\Theta(n^k)$ .

Using this version of the Master Theorem we can see that in our recurrence  $a = 2$ ,  $b = 2$ , and  $k = 1$ , so  $a = b^k$  and case (2) applies. Thus  $T(n)$  is  $\Theta(n \log n)$ .

There many recurrences that cannot be put into this form. For example, the following recurrence is quite common:  $T(n) = 2T(n/2) + n \log n$ . This solves to  $T(n) = \Theta(n \log^2 n)$ , but the Master Theorem (either this form or the one in CLRS will not tell you this.) For such recurrences, other methods are needed.

**Expansion:** A more basic method for solving recurrences is that of *expansion* (which CLRS calls *iteration*). This is a rather painstaking process of repeatedly applying the definition of the recurrence until (hopefully) a simple pattern emerges. This pattern usually results in a summation that is easy to solve. If you look at the proof in CLRS for the Master Theorem, it is actually based on expansion.

Let us consider applying this to the following recurrence. We assume that  $n$  is a power of 3.

$$\begin{aligned} T(1) &= 1 \\ T(n) &= 2T\left(\frac{n}{3}\right) + n \quad \text{if } n > 1 \end{aligned}$$

First we expand the recurrence into a summation, until seeing the general pattern emerge.

$$\begin{aligned} T(n) &= 2T\left(\frac{n}{3}\right) + n \\ &= 2\left(2T\left(\frac{n}{9}\right) + \frac{n}{3}\right) + n = 4T\left(\frac{n}{9}\right) + \left(n + \frac{2n}{3}\right) \\ &= 4\left(2T\left(\frac{n}{27}\right) + \frac{n}{9}\right) + \left(n + \frac{2n}{3}\right) = 8T\left(\frac{n}{27}\right) + \left(n + \frac{2n}{3} + \frac{4n}{9}\right) \\ &\vdots \\ &= 2^k T\left(\frac{n}{3^k}\right) + \sum_{i=0}^{k-1} \frac{2^i n}{3^i} = 2^k T\left(\frac{n}{3^k}\right) + n \sum_{i=0}^{k-1} (2/3)^i. \end{aligned}$$

The parameter  $k$  is the number of expansions (not to be confused with the value of  $k$  we introduced earlier on the overhead). We want to know how many expansions are needed to arrive at the basis case. To do this we set  $n/(3^k) = 1$ , meaning that  $k = \log_3 n$ . Substituting this in and using the identity  $a^{\log_b c} = b^{\log_a c}$  we have:

$$T(n) = 2^{\log_3 n} T(1) + n \sum_{i=0}^{\log_3 n - 1} (2/3)^i = n^{\log_3 2} + n \sum_{i=0}^{\log_3 n - 1} (2/3)^i.$$

Next, we can apply the formula for the geometric series and simplify to get:

$$\begin{aligned} T(n) &= n^{\log_3 2} + n \frac{1 - (2/3)^{\log_3 n}}{1 - (2/3)} \\ &= n^{\log_3 2} + 3n(1 - (2/3)^{\log_3 n}) = n^{\log_3 2} + 3n(1 - n^{\log_3(2/3)}) \\ &= n^{\log_3 2} + 3n(1 - n^{(\log_3 2) - 1}) = n^{\log_3 2} + 3n - 3n^{\log_3 2} \\ &= 3n - 2n^{\log_3 2}. \end{aligned}$$

Since  $\log_3 2 \approx 0.631 < 1$ ,  $T(n)$  is dominated by the  $3n$  term asymptotically, and so it is  $\Theta(n)$ .

**Induction and Constructive Induction:** Another technique for solving recurrences (and this works for summations as well) is to guess the solution, or the general form of the solution, and then attempt to verify its correctness through induction. Sometimes there are parameters whose values you do not know. This is fine. In the course of the induction proof, you will usually find out what these values must be. We will consider a famous example, that of the *Fibonacci numbers*.

$$\begin{aligned} F_0 &= 0 \\ F_1 &= 1 \\ F_n &= F_{n-1} + F_{n-2} \quad \text{for } n \geq 2. \end{aligned}$$

The Fibonacci numbers arise in data structure design. If you study AVL (height balanced) trees in data structures, you will learn that the minimum-sized AVL trees are produced by the recursive construction given below. Let  $L(i)$  denote the number of leaves in the minimum-sized AVL tree of height  $i$ . To construct a minimum-sized AVL tree of height  $i$ , you create a root node whose children consist of a minimum-sized AVL tree of heights  $i-1$  and  $i-2$ . Thus the number of leaves obeys  $L(0) = L(1) = 1$ ,  $L(i) = L(i-1) + L(i-2)$ . It is easy to see that  $L(i) = F_{i+1}$ .

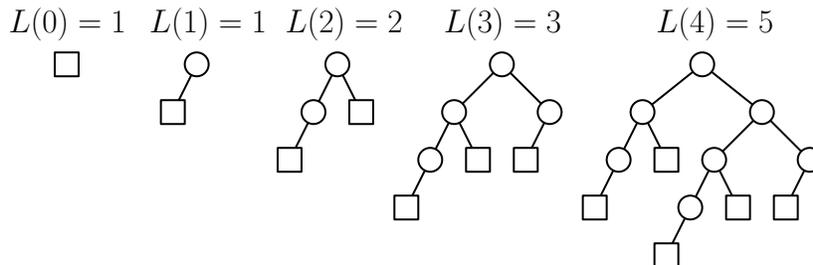


Fig. 84: Minimum-sized AVL trees.

If you expand the Fibonacci series for a number of terms, you will observe that  $F_n$  appears to grow exponentially, but not as fast as  $2^n$ . It is tempting to conjecture that  $F_n \leq \phi^{n-1}$ , for some real parameter  $\phi$ , where  $1 < \phi < 2$ . We can use induction to prove this and derive a bound on  $\phi$ .

**Lemma:** For all integers  $n \geq 1$ ,  $F_n \leq \phi^{n-1}$  for some constant  $\phi$ ,  $1 < \phi < 2$ .

**Proof:** We will try to derive the tightest bound we can on the value of  $\phi$ .

**Basis:** For the basis cases we consider  $n = 1$ . Observe that  $F_1 = 1 \leq \phi^0$ , as desired.

**Induction step:** For the induction step, let us assume that  $F_m \leq \phi^{m-1}$  whenever  $1 \leq m < n$ .

Using this *induction hypothesis* we will show that the lemma holds for  $n$  itself, whenever  $n \geq 2$ .

Since  $n \geq 2$ , we have  $F_n = F_{n-1} + F_{n-2}$ . Now, since  $n-1$  and  $n-2$  are both strictly less than  $n$ , we can apply the induction hypothesis, from which we have

$$F_n \leq \phi^{n-2} + \phi^{n-3} = \phi^{n-3}(1 + \phi).$$

We want to show that this is at most  $\phi^{n-1}$  (for a suitable choice of  $\phi$ ). Clearly this will be true if and only if  $(1 + \phi) \leq \phi^2$ . This is not true for all values of  $\phi$  (for example it is not true when  $\phi = 1$  but it is true when  $\phi = 2$ .)

At the critical value of  $\phi$  this inequality will be an equality, implying that we want to find the roots of the equation

$$\phi^2 - \phi - 1 = 0.$$

By the quadratic formula we have

$$\phi = \frac{1 \pm \sqrt{1+4}}{2} = \frac{1 \pm \sqrt{5}}{2}.$$

Since  $\sqrt{5} \approx 2.24$ , observe that one of the roots is negative, and hence would not be a possible candidate for  $\phi$ . The positive root is

$$\phi = \frac{1 + \sqrt{5}}{2} \approx 1.618.$$

There is a very subtle bug in the preceding proof. Can you spot it? The error occurs in the case  $n = 2$ . Here we claim that  $F_2 = F_1 + F_0$  and then we apply the induction hypothesis to both  $F_1$  and  $F_0$ . But the induction hypothesis only applies for  $m \geq 1$ , and hence cannot be applied to  $F_0$ ! To fix it we could include  $F_2$  as part of the basis case as well.

Notice not only did we prove the lemma by induction, but we actually determined the value of  $\phi$  which makes the lemma true. This is why this method is called *constructive induction*.

By the way, the value  $\phi = \frac{1}{2}(1 + \sqrt{5})$  is a famous constant in mathematics, architecture and art. It is the *golden ratio*. Two numbers  $A$  and  $B$  satisfy the golden ratio if

$$\frac{A}{B} = \frac{A+B}{A}.$$

It is easy to verify that  $A = \phi$  and  $B = 1$  satisfies this condition. This proportion occurs throughout the world of art and architecture.

## Lecture 24: Recurrences and Generating Functions

**Generating Functions:** The method of constructive induction provided a way to get a bound on  $F_n$ , but we did not get an exact answer, and we had to generate a good guess before we were even able to start.

Let us consider an approach to determine an exact representation of  $F_n$ , which requires no guesswork. This method is based on a very elegant concept, called a *generating function*. Consider any infinite sequence:

$$a_0, a_1, a_2, a_3, \dots$$

If we would like to “encode” this sequence succinctly, we could define a polynomial function such that these are the coefficients of the function:

$$G(z) = a_0 + a_1z + a_2z^2 + a_3z^3 + \dots$$

This is called the *generating function* of the sequence. What is  $z$ ? It is just a symbolic variable. We will (almost) never assign it a specific value. Thus, every infinite sequence of numbers has a corresponding generating function, and vice versa. What is the advantage of this representation? It turns out that we can perform arithmetic transformations on these functions (e.g., adding them, multiplying them, differentiating them) and this has a corresponding effect on the underlying transformations. It turns out that some nicely-structured sequences (like the Fibonacci numbers, and many sequences arising from linear recurrences) have generating functions that are easy to write down and manipulate.

Let’s consider the generating function for the Fibonacci numbers:

$$\begin{aligned} G(z) &= F_0 + F_1z + F_2z^2 + F_3z^3 + \dots \\ &= z + z^2 + 2z^3 + 3z^4 + 5z^5 + \dots \end{aligned}$$

The trick in dealing with generating functions is to figure out how various manipulations of the generating function to generate algebraically equivalent forms. For example, notice that if we multiply the generating function by a factor of  $z$ , this has the effect of shifting the sequence to the right:

$$\begin{aligned} G(z) &= F_0 + F_1z + F_2z^2 + F_3z^3 + F_4z^4 + \dots \\ zG(z) &= F_0z + F_1z^2 + F_2z^3 + F_3z^4 + \dots \\ z^2G(z) &= F_0z^2 + F_1z^3 + F_2z^4 + \dots \end{aligned}$$

Now, let's try the following manipulation. Compute  $G(z) - zG(z) - z^2G(z)$ , and see what we get

$$\begin{aligned} (1 - z - z^2)G(z) &= F_0 + (F_1 - F_0)z + (F_2 - F_1 - F_0)z^2 + (F_3 - F_2 - F_1)z^3 \\ &\quad + \dots + (F_i - F_{i-1} - F_{i-2})z^i + \dots \\ &= z. \end{aligned}$$

Observe that every term except the second is equal to zero by the definition of  $F_i$ . (The particular manipulation we picked was chosen to cause this cancellation to occur.) From this we may conclude that

$$G(z) = \frac{z}{1 - z - z^2}.$$

So, now we have an alternative representation for the Fibonacci numbers, as the coefficients of this function if expanded as a power series. So what good is this? The main goal is to get at the coefficients of its power series expansion. There are certain common tricks that people use to manipulate generating functions.

The first is to observe that there are some functions for which it is very easy to get a power series expansion. For example, the following is a simple consequence of the formula for the geometric series. If  $0 < c < 1$  then

$$\sum_{i=0}^{\infty} c^i = \frac{1}{1 - c}.$$

Setting  $z = c$ , we have

$$\frac{1}{1 - z} = 1 + z + z^2 + z^3 + \dots$$

(In other words,  $1/(1 - z)$  is the generating function for the sequence  $(1, 1, 1, \dots)$ ). In general, given a constant  $a$  we have

$$\frac{1}{1 - az} = 1 + az + a^2z^2 + a^3z^3 + \dots$$

is the generating function for  $(1, a, a^2, a^3, \dots)$ . It would be great if we could modify our generating function to be in the form of  $1/(1 - az)$  for some constant  $a$ , since then we could then extract the coefficients of the power series easily.

In order to do this, we would like to rewrite the generating function in the following form:

$$G(z) = \frac{z}{1 - z - z^2} = \frac{A}{1 - az} + \frac{B}{1 - bz},$$

for some  $A, B, a, b$ . We will skip the steps in doing this, but it is not hard to verify the roots of  $(1 - az)(1 - bz)$  (which are  $1/a$  and  $1/b$ ) must be equal to the roots of  $1 - z - z^2$ . We can then solve for  $a$  and  $b$  by taking the reciprocals of the roots of this quadratic. Then by some simple algebra we can plug these values in and solve for  $A$  and  $B$  yielding:

$$G(z) = \frac{z}{1 - z - z^2} = \left( \frac{1/\sqrt{5}}{1 - \phi z} + \frac{-1/\sqrt{5}}{1 - \hat{\phi} z} \right) = \frac{1}{\sqrt{5}} \left( \frac{1}{1 - \phi z} - \frac{1}{1 - \hat{\phi} z} \right),$$

where  $\phi = (1 + \sqrt{5})/2$  and  $\hat{\phi} = (1 - \sqrt{5})/2$ . (In particular, to determine  $A$ , multiply the equation by  $1 - \phi z$ , and then consider what happens when  $z = 1/\phi$ . A similar trick can be applied to get  $B$ . In general, this is called the method of *partial fractions*.)

Now we are in good shape, because we can extract the coefficients for these two fractions from the above function. From this we have the following:

$$G(z) = \frac{1}{\sqrt{5}} \left( \begin{array}{cccc} 1 & + & \phi z & + & \phi^2 z^2 & + & \dots \\ -1 & + & -\hat{\phi} z & + & -\hat{\phi}^2 z^2 & + & \dots \end{array} \right)$$

Combining terms we have

$$G(z) = \frac{1}{\sqrt{5}} \sum_{i=0}^{\infty} (\phi^i - \hat{\phi}^i) z^i.$$

We can now read off the coefficients easily. In particular it follows that

$$F_n = \frac{1}{\sqrt{5}} (\phi^n - \hat{\phi}^n).$$

This is an exact result, and no guesswork was needed. The only parts that involved some cleverness (beyond the invention of generating functions) was (1) coming up with the simple closed form formula for  $G(z)$  by taking appropriate differences and applying the rule for the recurrence, and (2) applying the method of partial fractions to get the generating function into one for which we could easily read off the final coefficients.

This is a rather remarkable, because it says that we can express the integer  $F_n$  as the sum of two powers of irrational numbers  $\phi$  and  $\hat{\phi}$ . You might try this for a few specific values of  $n$  to see why this is true. By the way, when you observe that  $\hat{\phi} < 1$ , it is clear that the first term is the dominant one. Thus we have, for large enough  $n$ ,  $F_n = \phi^n / \sqrt{5}$ , rounded to the nearest integer.

## Lecture 25: Algorithm Design: The Stable Marriage Problem

**Stable Marriage:** As an introduction to algorithm design, we will consider a well known discrete computational problem, called the *stable marriage problem*. In spite of the name, the problem's original formulation had nothing to do with the institution of marriage, but it was motivated by a number of practical applications where it was desired to set up pairings between entities, e.g., assigning medical school graduates to hospitals for residence training, assigning interns to companies, or assigning students to fraternities or sororities.

In all these applications we may have two groups of entities (e.g., students and university admission slots) where we wish to make an assignment from one to the other and where each side has some notion of preference. For example, each student has a ranking of the universities he/she wishes to attend and each university has a ranking of students it wants to admit. The goal is to produce a pairing that is in some sense "stable" in the sense that matched pairs should not have an obvious incentive to split up in order to form a different partnership.

Following tradition, we will couch this problem abstract in terms of a group of  $n$  men and  $n$  women that wish to be paired, that is, to *marry*.<sup>16</sup> We will place the algorithm in the role of a metaphorical matchmaker. First, we will use the traditional notion of marriage, the outcome of our process will be

---

<sup>16</sup>It is worth noting that in the above applications there is an asymmetrical relationship between the groups. We shall see that the algorithm that we will develop will *not* be gender-neutral. In particular, one gender will play a more active role and the other a more passive role. While this analogy may have made reasonable sense in early 1960's American culture, when the algorithm was first developed, many aspects of the algorithm's description seem out of place with modern culture. If this bothers you, please feel free to swap all references to "men" and "women".

a full pairing, one man to one woman and vice versa. Second, we assume that there is some notion of preference involved. This will be modeled by assuming that each man provides a rank ordering of the women according to decreasing preference level and vice versa.

Consider the following example. There are three women: Anny (A), Betty (B), and Carry (C), and there are three men: Eddy (E), Freddy (F), and Gerry (G). Here are their preferences (highest to lowest).

Men			Women		
Eddy (E)	Freddy (F)	Gerry (G)	Anny (A)	Betty (B)	Carry (C)
B	B	C	G	G	E
A	C	B	F	E	F
C	A	A	E	F	G

**Stability:** There are many ways in which we might define the notion of stable pairing of men to women. Clearly, we cannot guarantee that everyone will get their first preference. (Both Eddy and Freddy list Betty first.) There is a very weak condition that we would like to place on our matching. Intuitively, it should not be the case that there is a single unmarried pair would find it in their simultaneous best interest to ignore the pairing set up by the matchmaker and elope together. That is, there should be no man who can say to another woman, “We each prefer each other to our assigned partners—let’s elope!” If no such *instability* exists, the pairing is said to be *stable*.

**Definition 1:** Given a pair of sets  $X$  and  $Y$ , a *matching*, is a collection of pairs  $(x, y)$ , where  $x \in X$  and  $y \in Y$ , and each element of  $X$  appears in at most one pair, and each element of  $Y$  appears in at most one pair. A matching is *perfect* if every element of  $X$  and  $Y$  occurs in some pair. (**Beware:** Perfectness in a matching has nothing to do with optimality or stability. It simply means that everyone has a mate.)

**Definition 2:** Given sets  $X$  and  $Y$  of equal size and a preference ordering for each element of each set, a perfect matching is *stable* if there is no pair  $(x, y)$  that is *not* in the matching and  $x$  prefers  $y$  to its current match and  $y$  prefers  $x$  to its current match.

For example, among the following, can you spot which are stable and which are unstable? To make it easier to spot instabilities, after each person I have listed in brackets the people that they would have preferred over their assigned choice.

Assignment I	Assignment II	Assignment III
E [B] ↔ A [G, F]	E [B, A] ↔ C [ ]	E [B, A] ↔ C [ ]
F [B] ↔ C [E]	F [ ] ↔ B [G, E]	F [B, C] ↔ A [G]
G [C] ↔ B [ ]	G [C, B] ↔ A [ ]	G [C] ↔ B [ ]

The answer appears in the footnote below<sup>17</sup> You might wonder whether among all stable matchings, are some better than others? What would “better” mean? (More stable?) We will not consider this issue here, but it is an interesting one.

**The Gale-Shapley Algorithm:** The algorithm that we will describe is essentially due to Gale and Shapley, who considered this problem back in 1962. The algorithm is based on two basic primitive actions:

**Proposal:** An unengaged man makes a proposal to a woman

<sup>17</sup>The only unstable one is II. Observe that Eddy would have preferred Betty over his assigned mate Carry, and Betty would have preferred Eddy to her assigned mate Freddy. Thus, the unmarried pair  $(E, B)$  is an example of an instability. It is easy to verify that assignments I and III are stable.

**Decision:** A woman who receives a proposal can either accept or reject it. If she is already engaged and accepts a proposal, her existing engagement is broken off, and her old mate becomes unengaged.

There is an obvious sexual bias here, since men do the proposing and women do the deciding. It is interesting to consider a more balanced system where either side can offer proposals. (Not surprisingly, it does make a difference whether men or women do the proposing, from the perspective of who tends to get assigned mates of higher preference. We'll leave this question as an exercise.)

The original Gale-Shapley algorithm was presented as occurring over a sequence of *rounds*, during which all the unengaged men make proposals all at once, followed by the women either accepting or rejecting these proposals. However, in our book this is simplified by observing that the loop structure is simpler (and the results no different) if we process one man at a time, repeating the process until every man is either engaged or has exhausted everyone on his preference list.

We present the code for the Gale-Shapley algorithm in the following code block. Our presentation is not based on the above rounds-structure, but rather in the form that Kleinberg and Tardos present it, where in each iteration a single proposal is made and decided upon. (The two algorithms are essentially no different, and the order of events and final results are the same for both.) An example of this algorithm on the preferences given above is shown in Fig. 85.

The Gale-Shapley Algorithm

```
// Input: 2n preference lists, each consisting of n names.
// Output: A matching that pairs each man with each woman.
Initially all men and all women are unengaged
while (there is an unengaged man who hasn't yet proposed to every woman) {
  Let m be any such man
  Let w be the highest woman on his list to whom he has not yet proposed
  if (w is unengaged) then she accepts ((m, w) are now engaged)
  else {
    Let m' be the man w is engaged to currently
    if (w prefers m to m') {
      Break off the engagement (m', w)
      Create the new engagement (m, w) (upgrade)
      Man m' is now unengaged
    }
  }
}
```

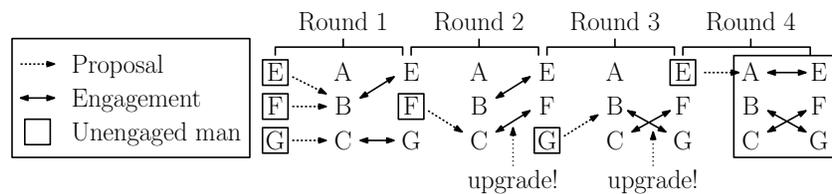


Fig. 85: Example of the round form version of the GS Algorithm on the preference lists given earlier. The final matching is (Eddy ↔ Anny), (Freddy ↔ Betty), (Gerry ↔ Carry).

**Correctness of the Gale-Shapley Algorithm:** Here are some easy observations regarding the Gale-Shapley (GS) Algorithm.

**Lemma 1:** Once a woman becomes engaged, she remains engaged for the remainder of the algorithm (although her mate may change), and her mate can only get better over time in terms of her preference list.

**Lemma 2:** The mates assigned to each man decrease over time in terms of his preference list.

Lemma 1 follows from the fact that a woman only breaks off an engagement to form a new one to a man of higher preference. Lemma 2 follows from the fact that each man makes offers in decreasing preference order.

Next we show that the algorithm terminates.

**Lemma 3:** The GS Algorithm terminates after at most  $n^2$  iterations of the while loop.

**Proof:** Consider the pairs  $(m, w)$  in which man  $m$  has not yet proposed to woman  $w$ . Initially there are  $n^2$  such pairs, but with each iteration of the while loop, at least one man proposes to one woman. Once a man proposes to a woman, he will never propose to her again (by Lemma 2). Thus, after  $n^2$  iterations, no one is left to propose.

The above lemma does not imply that the algorithm succeeds in finding a pairing between all the pairs (stable or not), and so we prove this next. Recall that a 1-to-1 pairing is called a *perfect matching*.

**Lemma 4:** On termination of the GS algorithm, the set of engagements form a perfect matching.

**Proof:** Every time we create a new engagement we break an old one. Thus, at any time, each woman is engaged to exactly one man, and vice versa. The only thing that could go wrong is that, at the end of the algorithm, some man  $m$  is unengaged after exhausting his list. Since there is a 1-to-1 correspondence between engaged men and engaged women, this would imply that some woman  $w$  is also unengaged. From Lemma 1 we know that once a woman is asked, she will become engaged and will remain engaged henceforth (although possibly to different mates). This implies that  $w$  has never been asked. But she appears on  $m$ 's list, and therefore she must have been asked, a contradiction.

Finally, we show that the resulting perfect matching is indeed stable. This establishes the correctness of the GS algorithm formally.

**Lemma 5:** The matching output by the GS algorithm is a stable matching.

**Proof:** Suppose to the contrary that there is some instability in the final output. This means that there is an unmarried pair  $(m, w)$  with the following properties. Let  $w'$  denote the assigned mate of  $m$  and let  $m'$  denote the assigned mate to  $w$ .

- $m$  prefers  $w$  to his assigned mate  $w'$ , and
- $w$  prefers  $m$  to her assigned mate  $m'$  (see Fig. 86),

(and hence  $m$  and  $w$  have the incentive to elope).

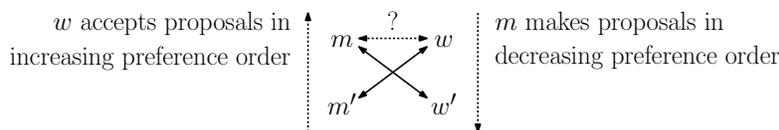


Fig. 86: The proof of Lemma 5.

Let's see why this cannot happen. Observe that since  $m$  prefers  $w$  he proposed to his preferred mate  $w$  before  $w'$ . What went wrong with his plans? Either  $w$  was already engaged to someone she preferred over  $m$  and rejected the offer outright, or she took his offer initially but later opted for someone whom she preferred and broke off the engagement with  $m$ . (Recall from Lemma 1 that once engaged, a woman's assigned mate only improves over time with respect to her preferences.) In either case,  $w$  ends up with someone she prefers over  $m$ . This means that she ends up with someone that she prefers over  $m'$ , whom she ranked even lower than  $m$ . Thus, the pair  $(m', w)$  could never have been generated by the algorithm, which yields the desired contradiction.

In summary, we have shown that the algorithm terminates, and it generates a correct result.

**Algorithm Efficiency:** Is this an efficient algorithm? Observe that this is much more efficient than a brute-force algorithm, which simply enumerates all the possible matchings, testing whether each is stable. This algorithm would take at least  $\Omega(n!)$  running time. Given how fast the factorial function grows, such an approach would only be useable for very small input sizes.

As observed earlier in Lemma 3, the GS algorithm runs in  $O(n^2)$  time. While normally, we would be inclined to call an algorithm running in  $O(n^2)$  time a *quadratic time* algorithm, notice that this is deceptively inaccurate. When we express running time, we do so in terms of the input size. In this case, the input for  $n$  men and  $n$  women consists of  $2n$  preference lists, each consisting of  $n$  elements. Thus the input size is  $N = 2n^2$ . Since the algorithm runs in  $O(n^2) = O(N)$  time, this is really a linear-time algorithm!

Note that in the practical applications where the GS algorithm is used, the input size is actually only  $O(n)$ . The reason is that, when very large input sizes are involved, it may not be practical to ask every many to rank order every woman, and vice versa. Typically, an individual is asked to rank just the top three or top five items in their preference list, and we hope that we can come up with a reasonably stable matching. Of course, if the preference lists are incomplete in this manner, then the algorithm may fail to produce a stable matching.

## Lecture 26: Dijkstra's Algorithm for Shortest Paths

**Shortest Paths:** Today we consider the problem of computing shortest paths in a directed graph. We have already seen that breadth-first search is an  $O(V + E)$  algorithm for finding shortest paths from a single source vertex to all other vertices, assuming that the graph has no edge weights. (Thus, distance is the number of edges on a path.) Suppose that each edge  $(u, v) \in E$  is associated with an edge weight  $w(u, v)$ . We define the *length* of a path to be the sum of weights along the edges of the path. We define the *distance* between any two vertices  $u$  and  $v$  to be the minimum length of any path between the vertices. We will denote this by  $\delta(u, v)$ . Because a vertex is joined to itself by an empty path, we have  $\delta(u, u) = 0$ , for all  $u \in V$ .

There are many ways in which to formulate the shortest path problem. For example, we may want be interested in the shortest path between a single source vertex and a single sink vertex, or we might be given a collection of source-sink pairs. Alternately, in the *single source* shortest-path problem, we are given a source vertex  $s \in V$ , and we wish to compute shortest paths to all other vertices. (The *single-sink* is a simple variant, which can be obtained by reversing all the edge directions.) Finally, the *all-pairs* shortest path problem involves computing the distances between all pairs of vertices. Of course, in addition to computing the distance between vertices, we will want to provide some intermediate structure that makes it possible to reconstruct the shortest path. Today, we will consider an algorithm for the single-source problem.

**Single Source Shortest Paths:** The *single source shortest path* problem is as follows. We are given a digraph  $G = (V, E)$  with numeric edge weights and a distinguished *source vertex*,  $s \in V$ . The objective is to determine the distance  $\delta(s, v)$  from  $s$  to every vertex  $v$  in the graph.

An important issue in the design of a shortest path algorithm is whether negative-valued edge weights are allowed. (Negative edges weights do not usually arise in transportation networks, but they can arise in financial transaction networks, where a transaction (edge) may result in either a lost or a profit.) In general, the shortest path problem is well defined, even if the graph has negative edge weights, provided that there are no negative cost cycles. (Otherwise you can make the path arbitrarily "short" by iterating forever around such a cycle.) Today, we will present a simple greedy algorithm for the single-source problem, which assumes that the edge weights are nonnegative. The algorithm, called

*Dijkstra's algorithm*, was invented by the famous Dutch computer scientist, Edsger Dijkstra in 1959. It is among the most famous algorithms in Computer Science.

In our presentation of the algorithm, we will stress the task of computing just the distance from the source to each vertex (not the path itself). As we did in the breadth-first search algorithm, it will be possible to make a minor modification to compute the paths themselves. As in BFS, we will use *predecessor link*, that point the route back to the source. By reversing the resulting path, we can obtain the shortest path. Since we store one predecessor link per vertex, the total space needed is only  $O(n)$ .

**Shortest Paths and Relaxation:** The basic structure of Dijkstra's algorithm is to maintain an *estimate* of the shortest path for each vertex, call this  $d[v]$ . Intuitively  $d[v]$  stores the length of the shortest path from  $s$  to  $v$  that the algorithm currently knows of. Indeed, there will always exist a path of length  $d[v]$ , but it might not be the ultimate shortest path. Initially, we know of no paths, so  $d[v] = \infty$ , and  $d[s] = 0$ . As the algorithm proceeds and sees more and more vertices, it updates  $d[v]$  for each vertex in the graph, until all the  $d[v]$  values “converge” to the true shortest distances.

The process by which an estimate is updated is sometimes called *relaxation*. Here is how relaxation works. Intuitively, if you can see that your solution is not yet reached an optimum value, then push it a little closer to the optimum. In particular, if you discover a path from  $s$  to  $v$  shorter than  $d[v]$ , then you need to update  $d[v]$ . This notion is common to many optimization algorithms.

Consider an edge from a vertex  $u$  to  $v$  whose weight is  $w(u, v)$ . Suppose that we have already computed current estimates on  $d[u]$  and  $d[v]$ . We know that there is a path from  $s$  to  $u$  of weight  $d[u]$ . By taking this path and following it with the edge  $(u, v)$  we get a path to  $v$  of length  $d[u] + w(u, v)$ . If this path is better than the existing path of length  $d[v]$  to  $v$ , we should update  $d[v]$  to the value  $d[u] + w(u, v)$  (see Fig. 89.) We should also remember that the shortest path to  $v$  passes through  $u$ , which we do by setting  $\text{pred}[v]$  to  $u$  (see the code block below).

```

relax(u, v) {
  if (d[u] + w(u, v) < d[v]) { // is the path through u shorter?
    d[v] = d[u] + w(u, v) // yes, then take it
    pred[v] = u // record that we go through u
  }
}

```

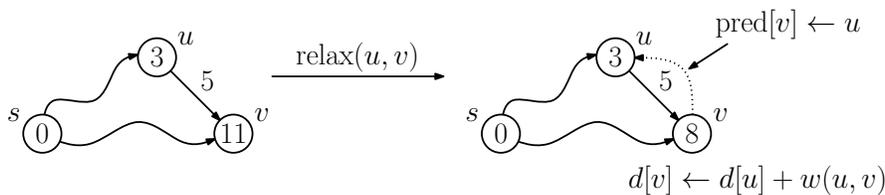


Fig. 87: Relaxation.

Observe that whenever we set  $d[v]$  to a finite value, there is always evidence of a path of that length. Therefore  $d[v] \geq \delta(s, v)$ . If  $d[v] = \delta(s, v)$ , then further relaxations cannot change its value.

It is not hard to see that if we perform  $\text{relax}(u, v)$  repeatedly over all edges of the graph, the  $d[v]$  values will eventually converge to the final true distance value from  $s$ . The cleverness of any shortest path algorithm is to perform the updates in a judicious manner, so the convergence is as fast as possible. In particular, the best possible would be to order relaxation operations in such a way that each edge is relaxed exactly once. Assuming that the edge weights are nonnegative, Dijkstra's algorithm achieves this objective.<sup>18</sup>

<sup>18</sup>Note, by the way that while this objective is optimal in the worst case, there are instances where you might hope for much

**Dijkstra’s Algorithm:** Dijkstra’s algorithm operates by maintaining a subset of vertices,  $S \subseteq V$ , for which we claim we “know” the true distance, that is  $d[v] = \delta(s, v)$ . Initially  $S = \emptyset$ , the empty set, and we set  $d[s] = 0$  and all others to  $+\infty$ . One by one, we select vertices from  $V \setminus S$  to add to  $S$ . (If you haven’t seen it before, the notation “ $A \setminus B$ ” means the set  $A$  excluding the elements of set  $B$ . Thus  $V \setminus S$  consists of the vertices that are not in  $S$ .)

The set  $S$  can be implemented using an array of vertex marks. Initially all vertices are marked as “undiscovered,” and we set  $\text{mark}[v] = \text{finished}$  to indicate that  $v \in S$ .

How do we select which vertex among the vertices of  $V \setminus S$  to add next to  $S$ ? Here is where greedy selection comes in. Dijkstra recognized that the best way in which to perform relaxations is by increasing order of distance from the source. This way, whenever a relaxation is being performed, it is possible to infer that result of the relaxation yields the final distance value. To implement this, we take the vertex of  $V \setminus S$  for which  $d[u]$  is minimum. That is, we take the unprocessed vertex that is closest (by our estimate) to  $s$ . Later we will justify why this is the proper choice.

In order to perform this selection efficiently, we store the vertices of  $V \setminus S$  in a *priority queue* (e.g. a heap), where the key value of each vertex  $u$  is  $d[u]$ . We will need to make use of three basic operations that are provided by the priority queue:

**Build:** Create a priority queue from a list of  $n$  elements, each with an associated key value.

**Extract min:** Remove (and return a reference to) the element with the smallest key value.

**Decrease key:** Given a reference to an element in the priority queue, decrease its key value to a specified value, and reorganize if needed.

For example, using a standard binary heap (as in heapsort) the first operation can be done in  $O(n)$  time, and the other two can be done in  $O(\log n)$  time each. Dijkstra’s algorithm is given in the code block below, and see Fig. 87 for an example.

Notice that the marking is not really used by the algorithm, but it has been included to make the connection with the correctness proof a little clearer.

To analyze Dijkstra’s algorithm, recall that  $n = |V|$  and  $m = |E|$ . We account for the time spent on each vertex after it is extracted from the priority queue. It takes  $O(\log n)$  to extract this vertex from the queue. For each incident edge, we spend potentially  $O(\log n)$  time if we need to decrease the key of the neighboring vertex. Thus the time is  $O(\log n + \deg(u) \cdot \log n)$  time. The other steps of the update run in constant time. Recalling that the sum of degrees of the vertices in a graph is  $O(m)$ , the overall running time is given by  $T(n, m)$ , where

$$\begin{aligned} T(n, m) &= \sum_{u \in V} (\log n + \deg(u) \cdot \log n) = \sum_{u \in V} (1 + \deg(u)) \log n \\ &= \log n \sum_{u \in V} (1 + \deg(u)) = (\log n)(n + 2m) = \Theta((n + m) \log n). \end{aligned}$$

Since  $G$  is connected,  $n$  is asymptotically no greater than  $m$ , so this is  $O(m \log n)$ .

**Correctness:** Recall that  $d[v]$  is the distance value assigned to vertex  $v$  by Dijkstra’s algorithm, and let  $\delta(s, v)$  denote the length of the true shortest path from  $s$  to  $v$ . To see that Dijkstra’s algorithm correctly gives the final true distances, we need to show that  $d[v] = \delta(s, v)$  when the algorithm terminates. This is a consequence of the following lemma, which states that once a vertex  $u$  has been added to  $S$  (i.e., has been marked “finished”),  $d[u]$  is the true shortest distance from  $s$  to  $u$ . Since at the end of the algorithm, all vertices are in  $S$ , then all distance estimates are correct.

---

better performance. For example, given a cartographic road map of the entire United States, computing the shortest path between two locations near Washington DC should not require relaxing every edge of this road map. A better approach to this problem is provided by another greedy algorithm, called  $A^*$ -search.

```

dijkstra(G,w,s) {
  for each (u in V) {                                // initialization
    d[u] = +infinity
    mark[u] = undiscovered
    pred[u] = null
  }
  d[s] = 0                                           // distance to source is 0
  Q = a priority queue of all vertices u sorted by d[u]
  while (Q is nonEmpty) {                            // until all vertices processed
    u = extract vertex with minimum d[u] from Q
    for each (v in Adj[u]) {
      if (d[u] + w(u,v) < d[v]) { // relax(u,v)
        d[v] = d[u] + w(u,v)
        decrease v's key in Q to d[v]
        pred[v] = u
      }
    }
    mark[u] = finished
  }
  [The pred pointers define an "inverted" shortest path tree]
}

```

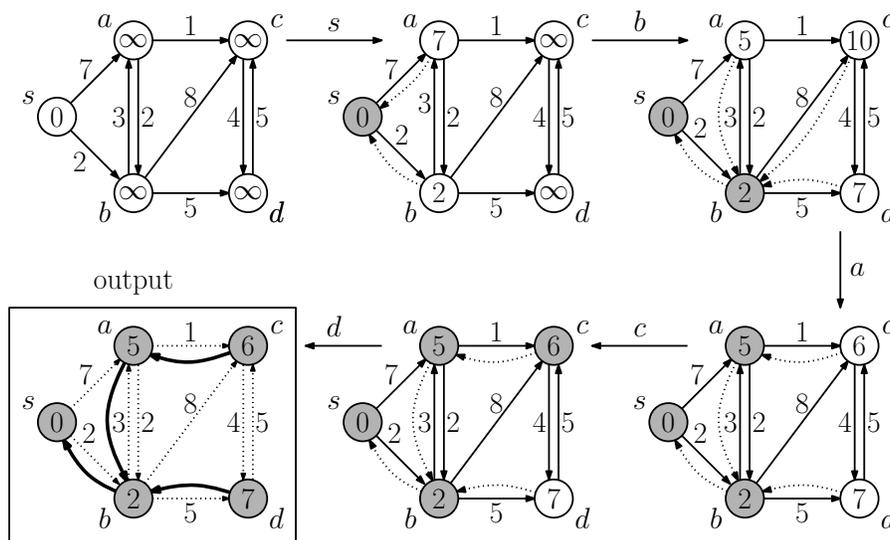


Fig. 88: Dijkstra's Algorithm example.

**Lemma:** When a vertex  $u$  is added to  $S$ ,  $d[u] = \delta(s, u)$ .

**Proof:** Suppose to the contrary that at some point Dijkstra's algorithm *first* attempts to add a vertex  $u$  to  $S$  for which  $d[u] \neq \delta(s, u)$ . By our observations about relaxation,  $d[u]$  is never less than  $\delta(s, u)$ , thus we have  $d[u] > \delta(s, u)$ . Consider the situation just prior to the insertion of  $u$ , and consider the true shortest path from  $s$  to  $u$ . Because  $s \in S$  and  $u \in V \setminus S$ , at some point this path must first jump out of  $S$ . Let  $(x, y)$  be the first edge taken by the shortest path, where  $x \in S$  and  $y \in V \setminus S$  (see Fig. 89). (Note that it may be that  $x = s$  and/or  $y = u$ ).

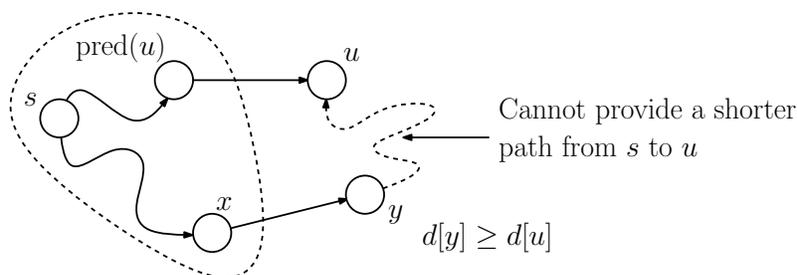


Fig. 89: Correctness of Dijkstra's Algorithm.

Because  $u$  is the first vertex where we made a mistake and since  $x$  was already processed, we have  $d[x] = \delta(s, x)$ . Since we applied relaxation to  $x$  when it was processed, we must have

$$d[y] = d[x] + w(x, y) = \delta(s, x) + w(x, y) = \delta(s, y).$$

Since  $y$  appears before  $u$  along the shortest path and edge weights are nonnegative, we have  $\delta(s, y) \leq \delta(s, u)$ . Also, because  $u$  (not  $y$ ) was chosen next for processing, we know that  $d[u] \leq d[y]$ . Putting this together, we have

$$\delta(s, u) < d[u] \leq d[y] = \delta(s, y) \leq \delta(s, u).$$

Clearly we cannot have  $\delta(s, u) < \delta(s, u)$ , which establishes the desired contradiction.

**Variants:** Dijkstra's algorithm is robust to a number of variations. Here are some variants of the single-source shortest path problem that can be solved either by making a small modification to Dijkstra's algorithm and/or by modifying the underlying graph. (I'll leave their solutions as an exercise.)

**Vertex weights:** There is a cost associated with each vertex. The overall cost is the sum of vertex and/or edge weights on the path.

**Single-Sink Shortest Path:** Find the shortest path from each vertex to a sink vertex  $t$ .

**Multi-Source/Multi-Sink:** You are given a collection of source vertices  $\{s_1, \dots, s_k\}$ . For each vertex find the shortest path from its nearest source. (Analogous for multi-sink.)

**Multiplicative Cost:** Define the cost of a path to be the product of the edge weights (rather than the sum.) If all the edge weights are at least 1, find the single-source shortest path.

## Lecture 27: Greedy Algorithms for Minimum Spanning Trees

**Minimum Spanning Trees:** A common problem in communications networks and circuit design is that of connecting together a set of nodes (communication sites or circuit components) by a network of minimal total length (where length is the sum of the lengths of connecting wires). We assume that the network is undirected. To minimize the length of the connecting network, it never pays to have any

cycles (since we could break any cycle without destroying connectivity and decrease the total length). Since the resulting connection graph is connected, undirected, and acyclic, it is a *free tree*.

The computational problem is called the *minimum spanning tree* problem (MST for short). More formally, given a connected, undirected graph  $G = (V, E)$ , a *spanning tree* is an acyclic subset of edges  $T \subseteq E$  that connects all the vertices together. Assuming that each edge  $(u, v)$  of  $G$  has a numeric weight or cost,  $w(u, v)$ , (may be zero or negative) we define the cost of a spanning tree  $T$  to be the sum of edges in the spanning tree

$$w(T) = \sum_{(u,v) \in T} w(u,v).$$

A *minimum spanning tree* (MST) is a spanning tree of minimum weight. Note that the minimum spanning tree may not be unique, but it is true that if all the edge weights are distinct, then the MST will be distinct (this is a rather subtle fact, which we will not prove). Fig. 90 shows three spanning trees for the same graph, where the shaded rectangles indicate the edges in the spanning tree. The spanning tree shown in Fig. 90(a) is not a minimum spanning tree (in fact, it is a maximum weight spanning tree), while the other two are MSTs.

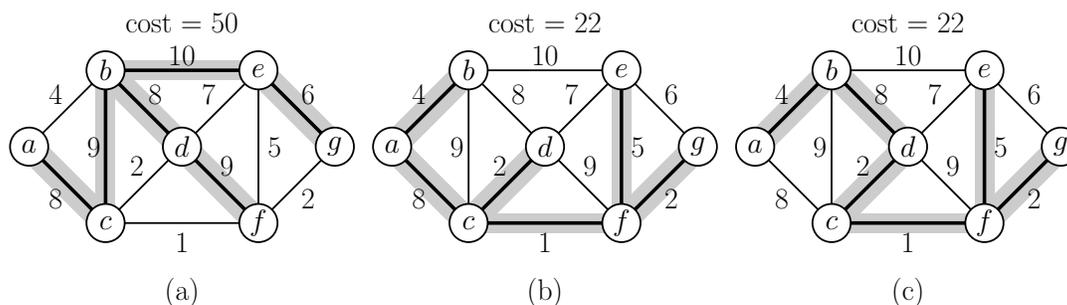


Fig. 90: Spanning trees (the middle and right are minimum spanning trees).

**Generic approach:** We will present three *greedy* algorithms (Kruskal’s, Prim’s, and Boruvka’s) for computing a minimum spanning tree. Recall that a *greedy algorithm* is one that builds a solution by repeated selecting the cheapest (or generally locally optimal choice) among all options at each stage. An important characteristic of greedy algorithms is that once they make a choice, they never “unmake” this choice. Before presenting these algorithms, let us review some basic facts about free trees. They are all quite easy to prove.

- Lemma:**
- (i) A free tree with  $n$  vertices has exactly  $n - 1$  edges.
  - (ii) There exists a unique path between any two vertices of a free tree.
  - (iii) Adding any edge to a free tree creates a unique cycle. Breaking *any* edge on this cycle restores a free tree.

Let  $G = (V, E)$  be the input graph. The intuition behind the greedy MST algorithms is simple, we maintain a subset of edges  $A$ , which will initially be empty, and we will add edges one at a time, until  $A$  is a spanning tree. We say that a subset  $A \subseteq E$  is *viable* if  $A$  is a subset of edges in some MST. We say that an edge  $(u, v) \in E \setminus A$  is *safe* if  $A \cup \{(u, v)\}$  is viable. (Recall that  $E \setminus A$  means the edges of  $E$  that are *not* in  $A$ .) In other words, the choice  $(u, v)$  is a safe choice to add so that  $A$  can still be extended to form an MST. Note that if  $A$  is viable it cannot contain a cycle. A generic greedy algorithm operates by repeatedly adding any *safe* edge to the current spanning tree. (Note that viability is a property of subsets of edges and safety is a property of a single edge.)

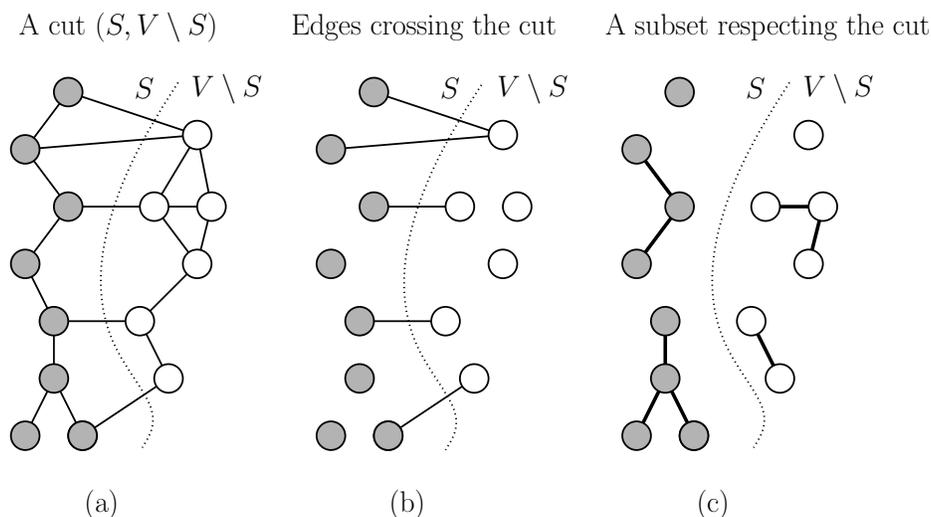


Fig. 91: MST-related terminology.

**When is an edge safe?** Let  $S$  be a subset of the vertices  $S \subseteq V$ . Here are a few useful definitions (see Fig. 91):

- A *cut*  $(S, V \setminus S)$  is a partition of the vertices into two disjoint subsets (see Fig. 91(a)).
- An edge  $(u, v)$  *crosses* the cut if  $u \in S$  and  $v \notin S$  (see Fig. 91(b)).
- Given a subset of edges  $A$ , we say that a cut *respects*  $A$  if no edge in  $A$  crosses the cut (see Fig. 91(c)).

It is not hard to see why respecting cuts are important to this problem. If we have computed a partial MST, and we wish to know which edges can be added that do *not* induce a cycle in the current MST, any edge that crosses a respecting cut is a possible candidate.

An edge of  $E$  is a *light edge* crossing a cut, if among all edges crossing the cut, it has the minimum weight (the light edge may not be unique if there are duplicate edge weights). Intuition says that since all the edges that cross a respecting cut do not induce a cycle, then the lightest edge crossing a cut is a natural choice. The main theorem which drives both algorithms is the following. It essentially says that we can always augment  $A$  by adding the minimum weight edge that crosses a cut which respects  $A$ . (It is stated in complete generality, so that it can be applied to both algorithms.)

**MST Lemma:** Let  $G = (V, E)$  be a connected, undirected graph with real-valued weights on the edges. Let  $A$  be a viable subset of  $E$  (i.e. a subset of some MST), let  $(S, V \setminus S)$  be any cut that respects  $A$ , and let  $(u, v)$  be a light edge crossing this cut. Then the edge  $(u, v)$  is *safe* for  $A$ .

**Proof:** It will simplify the proof to assume that all the edge weights are distinct. Let  $T$  be any MST for  $G$  (see Fig. 92). If  $T$  contains  $(u, v)$  then we are done. Suppose that no MST contains  $(u, v)$ . We will derive a contradiction.

Add the edge  $(u, v)$  to  $T$ , thus creating a cycle. Since  $u$  and  $v$  are on opposite sides of the cut and since any cycle must cross the cut an even number of times, there must be at least one other edge  $(x, y)$  in  $T$  that crosses the cut.

The edge  $(x, y)$  is not in  $A$  (because the cut respects  $A$ ). By removing  $(x, y)$  we restore a spanning tree, call it  $T'$ . We have

$$w(T') = w(T) - w(x, y) + w(u, v).$$

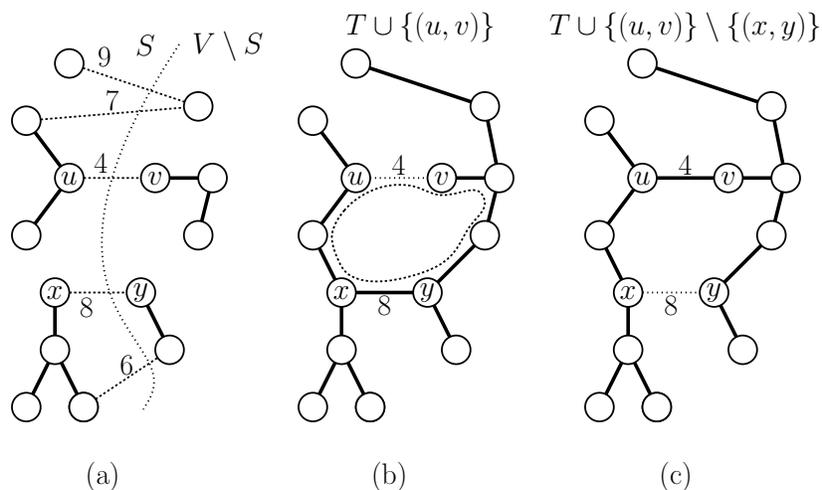


Fig. 92: Proof of the MST Lemma. Edge  $(u, v)$  is the light edge crossing cut  $(S, V \setminus S)$ .

Since  $(u, v)$  is lightest edge crossing the cut, we have  $w(u, v) < w(x, y)$ . Thus  $w(T') < w(T)$ . This contradicts the assumption that  $T$  was an MST.

**Kruskal's Algorithm:** Kruskal's algorithm works by attempting to add edges to the  $A$  in increasing order of weight (lightest edges first). If the next edge does not induce a cycle among the current set of edges, then it is added to  $A$ . If it does, then this edge is passed over, and we consider the next edge in order. Note that as this algorithm runs, the edges of  $A$  will induce a forest on the vertices. As the algorithm continues, the trees of this forest are merged together, until we have a single tree containing all the vertices.

Observe that this strategy leads to a correct algorithm. Why? Consider the edge  $(u, v)$  that Kruskal's algorithm seeks to add next, and suppose that this edge does not induce a cycle in  $A$ . Let  $A'$  denote the tree of the forest  $A$  that contains vertex  $u$ . Consider the cut  $(A', V \setminus A')$ . Every edge crossing the cut is not in  $A$ , and so this cut respects  $A$ , and  $(u, v)$  is the light edge across the cut (because any lighter edge would have been considered earlier by the algorithm). Thus, by the MST Lemma,  $(u, v)$  is safe.

The only tricky part of the algorithm is how to detect efficiently whether the addition of an edge will create a cycle in  $A$ . We could perform a DFS on subgraph induced by the edges of  $A$ , but this will take too much time. We want a fast test that tells us whether  $u$  and  $v$  are in the same tree of  $A$ .

This can be done by a data structure (which we have not studied) called the *disjoint set union-find data structure*. This data structure supports three operations:

**create( $u$ ):** Create a set containing a single item  $v$ .

**find( $u$ ):** Find the set that contains a given item  $u$ .

**union( $u, v$ ):** Merge the set containing  $u$  and the set containing  $v$  into a common set.

**Theorem:** Given a collection of  $n$  elements, each initially in its own set, the union-find data structure can perform any sequence of up to  $n$  union and find operations in total  $O(n \cdot \alpha(n))$  time, where  $\alpha(n)$  is the (*extremely slow growing*) inverse Ackerman's function.

You are not responsible for knowing how this data structure works, since we will use it as a "black-box". In Kruskal's algorithm, the vertices of the graph will be the elements to be stored in the sets, and the

```

KruskalMST(G=(V,E), w) {
  A = {} // initially A is empty
  Place each vertex u in a set by itself
  Sort E in increasing order by weight w
  for each ((u, v) in this order) {
    if (find(u) != find(v)) { // u and v in different trees
      add (u, v) to A // join subtrees together
      union(u, v) // merge these two components
    }
  }
  return A
}

```

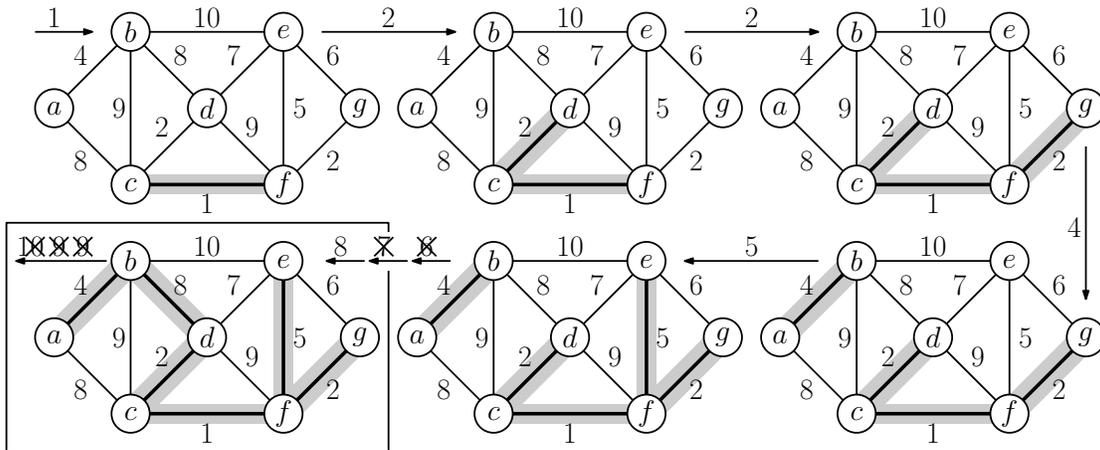


Fig. 93: Kruskal's Algorithm. Each vertex is labeled according to the set that contains it.

sets will be vertices in each tree of  $A$ . The set  $A$  can be stored as a simple list of edges. The algorithm is shown in the code fragment below, and an example is shown in Fig. 93.

**Analysis:** How long does Kruskal's algorithm take? As usual, let  $n$  be the number of vertices and  $m$  be the number of edges. Since the graph is connected, we may assume that  $m \geq n - 1$ . Observe that it takes  $\Theta(m \log m)$  time to sort the edges. The for-loop is iterated  $m$  times, and each iteration involves a constant number of accesses to the Union-Find data structure on a collection of  $n$  items. Thus each access is  $\Theta(n)$  time, for a total of  $\Theta(m \log n)$ . Thus the total running time is the sum of these, which is  $\Theta((n + m) \log n)$ . Since  $n$  is asymptotically no larger than  $m$ , we could write this more simply as  $\Theta(m \log n)$ .

**Prim's Algorithm:** Prim's algorithm is another greedy algorithm for computing minimum spanning trees. It differs from Kruskal's algorithm only in how it selects the next *safe edge* to add at each step. Its running time is essentially the same as Kruskal's algorithm,  $O((n + m) \log n)$ . There are two reasons for studying Prim's algorithm. The first is to show that there is more than one way to solve a problem (an important lesson to learn in algorithm design), and the second is that Prim's algorithm looks very much like another greedy algorithm, called Dijkstra's algorithm, that we will study for a completely different problem, shortest paths. Thus, not only is Prim's a different way to solve the same MST problem, it is also the same way to solve a different problem. (Whatever that means!)

**Different ways to grow a tree:** Kruskal's algorithm worked by ordering the edges, and inserting them one by one into the spanning tree, taking care never to introduce a cycle. Intuitively Kruskal's works by merging or splicing two trees together, until all the vertices are in the same tree.

In contrast, Prim's algorithm builds the tree up by adding leaves one at a time to the current tree. We start with a root vertex  $s$  (it can be *any* vertex). At any time, the subset of edges  $A$  forms a single tree (in Kruskal's it formed a forest). We look to add a single vertex as a leaf to the tree. The process is illustrated in the following figure.

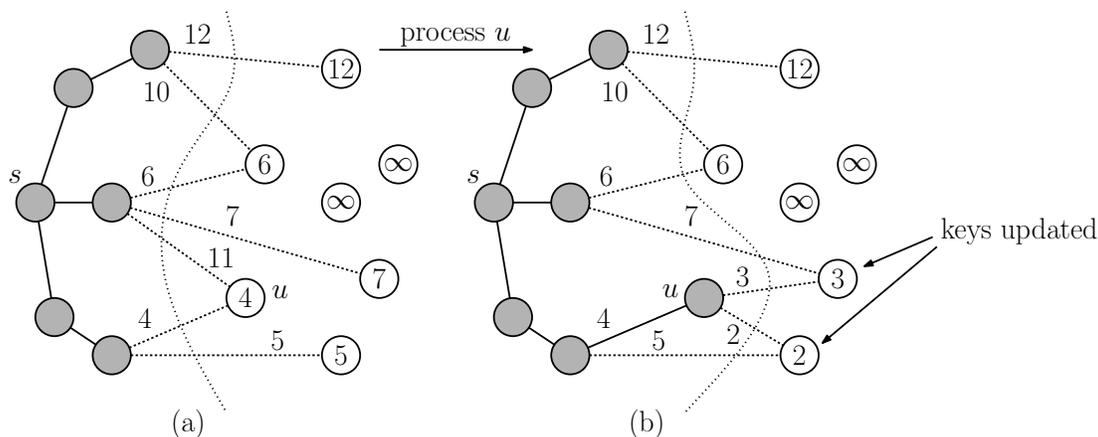


Fig. 94: Prim's Algorithm.

Observe that if we consider the set of vertices  $S$  currently part of the tree, and its complement  $(V \setminus S)$ , we have a cut of the graph and the current set of tree edges  $A$  respects this cut. Which edge should we add next? The MST Lemma from the previous lecture tells us that it is safe to add the *light edge*. In the figure, this is the edge of weight 4 going to vertex  $u$ . Then  $u$  is added to the vertices of  $S$ , and the cut changes. Note that some edges that crossed the cut before are no longer crossing it, and others that were not crossing the cut are.

It is easy to see, that the key questions in the efficient implementation of Prim's algorithm is how to update the cut efficiently, and how to determine the light edge quickly. To do this, we will make use

of a *priority queue* data structure. Recall that this is the data structure used in HeapSort. This is a data structure that stores a set of items, where each item is associated with a *key* value. The priority queue supports three operations.

**insert**( $u, k$ ): Insert  $u$  with the key value  $k$  in  $Q$ .

**extract-min**(): Extract the item with the minimum key value in  $Q$ .

**decrease-key**( $u, k'$ ): Decrease the value of  $u$ 's key value to  $k'$ .

A priority queue can be implemented using the same heap data structure used in heapsort. All of the above operations can be performed in  $O(\log n)$  time, where  $n$  is the number of items in the heap.

What do we store in the priority queue? At first you might think that we should store the edges that cross the cut, since this is what we are removing with each step of the algorithm. The problem is that when a vertex is moved from one side of the cut to the other, this results in a complicated sequence of updates.

There is a much more elegant solution, and this is what makes Prim's algorithm so nice. For each vertex in  $u \in V \setminus S$  (not part of the current spanning tree) we associate  $u$  with a key value  $key[u]$ , which is the weight of the lightest edge going from  $u$  to any vertex in  $S$ . We also store in  $pred[u]$  the end vertex of this edge in  $S$ . If there is not edge from  $u$  to a vertex in  $V \setminus S$ , then we set its key value to  $+\infty$ . We will also need to know which vertices are in  $S$  and which are not. We do this by coloring the vertices in  $S$  black.

Here is Prim's algorithm. The root vertex  $s$  can be any vertex in  $V$ .

---

Prim's Algorithm

```

PrimMST(G=(V,E), w, s) {
  for each (u in V) {                                // initialization
    key[u] = +infinity
    color[u] = undiscovered
  }
  key[s] = 0                                         // start at root
  pred[s] = null
  add all vertices to priority queue Q
  while (Q is nonEmpty) {                            // until all vertices in MST
    u = extract-min from Q                          // vertex with lightest edge
    for each (v in Adj[u]) {
      if ((color[v] == undiscovered) && (w(u,v) < key[v])) {
        key[v] = w(u,v)                            // new lighter edge out of v
        decrease key value of v to key[v]
        pred[v] = u
      }
    }
    color[u] = finished
  }
  [The pred pointers define the MST as an inverted tree rooted at s]
}

```

---

The following figure illustrates Prim's algorithm. The arrows on edges indicate the predecessor pointers, and the numeric label in each vertex is the key value.

To analyze Prim's algorithm, we account for the time spent on each vertex as it is extracted from the priority queue. It takes  $O(\log n)$  to extract this vertex from the queue. For each incident edge, we spend potentially  $O(\log n)$  time decreasing the key of the neighboring vertex. Thus the time is

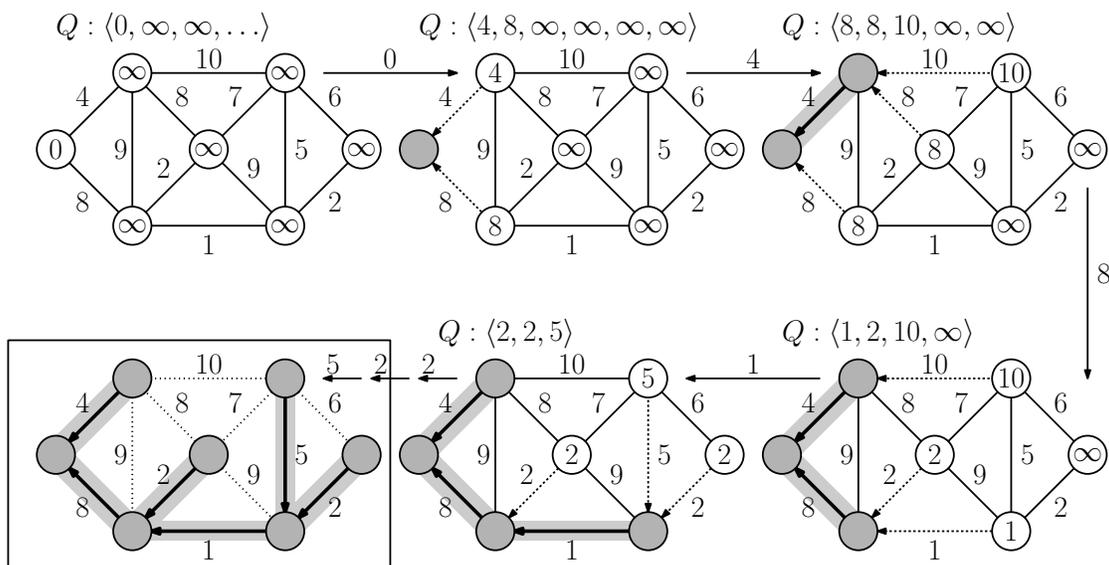


Fig. 95: Prim's algorithm example.

$O(\log n + \text{degree}(u) \log n)$  time. The other steps of the update are constant time. So the overall running time is

$$\begin{aligned}
 T(n, m) &= \sum_{u \in V} (\log n + \text{degree}(u) \log n) = \sum_{u \in V} (1 + \text{degree}(u)) \log n \\
 &= \log n \sum_{u \in V} (1 + \text{degree}(u)) = (\log n)(n + 2E) = \Theta((n + m) \log n).
 \end{aligned}$$

Since  $G$  is connected,  $n$  is asymptotically no greater than  $m$ , so this is  $\Theta(m \log n)$ . This is exactly the same as Kruskal's algorithm.

**Boruvka's Algorithm:** Given that we have seen two algorithms (Kruskal's and Prim's) for solving the MST problem, it may seem like complete overkill to consider yet another algorithm. This one is called Boruvka's algorithm. It is actually the oldest of the three algorithms (invented in 1926 by the Czech mathematician Otakar Břuvka, well before the first digital computers!). The reason for studying this algorithm is that of the three algorithms, it is the easiest to implement on a parallel computer. Unlike Kruskal's and Prim's algorithms, which add edges one at a time, Boruvka's algorithm adds a whole set of edges all at once to the MST.

Boruvka's algorithm is similar to Kruskal's algorithm, in the sense that it works by maintaining a collection of disconnected trees. Let us call each subtree a *component*. Initially, each vertex is by itself in a one-vertex component. Recall that with each stage of Kruskal's algorithm, we add the (globally) lightest (minimum weight) edge that connects two different components together. To prove Kruskal's algorithm correct, we argued (from the MST Lemma) that the lightest such edge will be *safe* to add to the MST.

In fact, a closer inspection of the proof reveals that the lightest edge exiting *any* component is always safe. This suggests a more parallel way to grow the MST. Each component determines the lightest edge that goes from inside the component to outside the component (we don't care where). We say that such an edge *leaves* the component.

Note that two components might select the same edge by this process. By the above observation, all of these edges are safe, so we may add them all at once to the set  $A$  of edges in the MST. If there are edges of equal weight, the algorithm might attempt to add two such edges between the same two components, which would result in the generation of a cycle. (Each edge individually is safe, but adding both simultaneously generates a cycle.) We make the assumption that edge weights are distinct (or at least, there is some uniform rule for breaking ties so the above problem cannot arise).

Note that in a single step of Boruvka's algorithm many components can be merged together into a single component. We then apply DFS to the edges of  $A$ , to identify the new components. This process is repeated until only one component remains. A fairly high-level description of Boruvka's algorithm is given below.

Boruvka's Algorithm

---

```

BoruvkaMST(G=(V,E), w) {
    initialize each vertex to be its own component
    A = {} // A holds edges of the MST
    while (there are two or more components) {
        for (each component C) {
            find the lightest edge (u,v) with u in C and v not in C
            add {u,v} to A (unless it is already there)
        }
        apply DFS to graph (V, A), to compute the new components
    }
    return A // return final MST edges

```

---

There are a number of unspecified details in Boruvka's algorithm, which we will not spell out in detail, except to note that they can be solved in  $\Theta(n + m)$  time through DFS. First, we may apply DFS, but only traversing the edges of  $A$  to compute the components. Each DFS tree will correspond to a separate component. We label each vertex with its component number as part of this process. With these labels it is easy to determine which edges go between components (since their endpoints have different labels). Then we can traverse each component again to determine the lightest edge that leaves the component. (In fact, with a little more cleverness, we can do all this without having to perform two separate DFS's.) The algorithm is illustrated in the figure below.

**Analysis:** How long does Boruvka's algorithm take? Observe that because each iteration involves doing a DFS, each iteration (of the outer do-while loop) can be performed in  $\Theta(n + m)$  time. The question is how many iterations are required in general? We claim that there are never more than  $O(\log n)$  iterations needed. To see why, let  $m$  denote the number of components at some stage. Each of the  $m$  components, will merge with at least one other component. Afterwards the number of remaining components could be as low as 1 (if they all merge together), but never higher than  $m/2$  (if they merge in pairs). Thus, the number of components decreases by at least half with each iteration. Since we start with  $n$  components, this can happen at most  $\lg n$  time, until only one component remains. Thus, the total running time is  $\Theta((n + m) \log n)$  time. Again, since  $G$  is connected,  $n$  is asymptotically no larger than  $m$ , so we can write this more succinctly as  $\Theta(m \log n)$ . Thus all three algorithms have the same asymptotic running time.

## Lecture 28: Medians and Selection

**Selection:** We have discussed recurrences and the divide-and-conquer method of solving problems. Today we will give a rather surprising (and very tricky) algorithm which shows the power of these techniques.

The problem that we will consider is very easy to state, but surprisingly difficult to solve optimally. Suppose that you are given a set of  $n$  numbers. Define the *rank* of an element to be one plus the

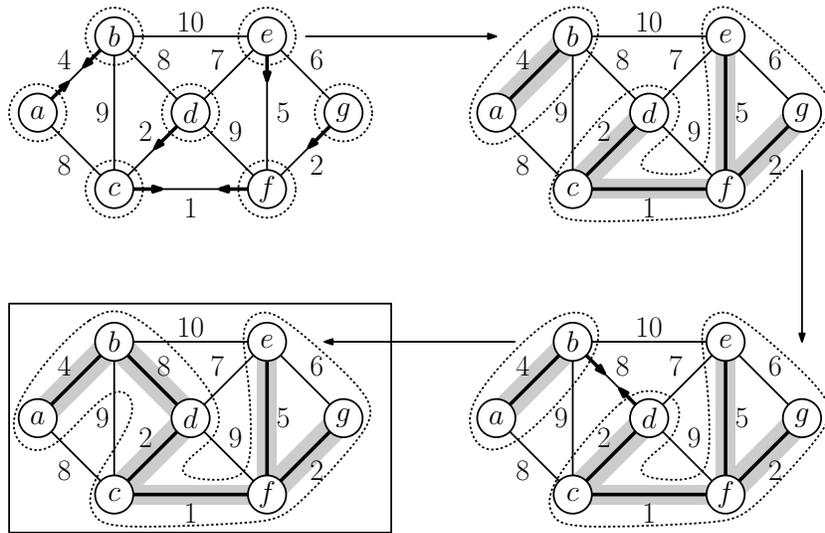


Fig. 96: Boruvka's Algorithm.

number of elements that are smaller than this element. Since duplicate elements make our life more complex (by creating multiple elements of the same rank), we will make the simplifying assumption that all the elements are distinct for now. It will be easy to get around this assumption later. Thus, the rank of an element is its final position if the set is sorted. The minimum is of rank 1 and the maximum is of rank  $n$ .

Of particular interest in statistics is the *median*. If  $n$  is odd then the median is defined to be the element of rank  $(n + 1)/2$ . When  $n$  is even there are two natural choices, namely the elements of ranks  $n/2$  and  $(n/2) + 1$ . In statistics it is common to return the average of these two elements. We will define the median to be either of these elements.

Medians are useful as measures of the *central tendency* of a set, especially when the distribution of values is highly skewed. For example, the median income in a community is likely to be more meaningful measure of the central tendency than the average is, since if Bill Gates lives in your community then his gigantic income may significantly bias the average, whereas it cannot have a significant influence on the median. They are also useful, since in divide-and-conquer applications, it is often desirable to partition a set about its median value, into two sets of roughly equal size. Today we will focus on the following generalization, called the *selection problem*.

**Selection:** Given a set  $A$  of  $n$  distinct numbers and an integer  $k$ ,  $1 \leq k \leq n$ , output the element of  $A$  of rank  $k$ .

The selection problem can easily be solved in  $\Theta(n \log n)$  time, simply by sorting the numbers of  $A$ , and then returning  $A[k]$ . The question is whether it is possible to do better. In particular, is it possible to solve this problem in  $\Theta(n)$  time? We will see that the answer is yes, and the solution is far from obvious.

**The Sieve Technique:** The reason for introducing this algorithm is that it illustrates a very important special case of divide-and-conquer, which I call the *sieve technique*. We think of divide-and-conquer as breaking the problem into a small number of smaller subproblems, which are then solved recursively. The sieve technique is a special case, where the number of subproblems is just 1.

The sieve technique works in phases as follows. It applies to problems where we are interested in finding a single item from a larger set of  $n$  items. We do not know which item is of interest, however after

doing some amount of analysis of the data, taking say  $\Theta(n^k)$  time, for some constant  $k$ , we find that we do not know what the desired item is, but we can identify a large enough number of elements that *cannot* be the desired value, and can be eliminated from further consideration. In particular “large enough” means that the number of items is at least some fixed constant fraction of  $n$  (e.g.  $n/2$ ,  $n/3$ ,  $0.0001n$ ). Then we solve the problem recursively on whatever items remain. Each of the resulting recursive solutions then do the same thing, eliminating a constant fraction of the remaining set.

**Applying the Sieve to Selection:** To see more concretely how the sieve technique works, let us apply it to the selection problem. Recall that we are given an array  $A[1..n]$  and an integer  $k$ , and want to find the  $k$ -th smallest element of  $A$ . Since the algorithm will be applied inductively, we will assume that we are given a subarray  $A[p..r]$  as we did in MergeSort, and we want to find the  $k$ th smallest item (where  $k \leq r - p + 1$ ). The initial call will be to the entire array  $A[1..n]$ .

There are two principal algorithms for solving the selection problem, but they differ only in one step, which involves judiciously choosing an item from the array, called the *pivot element*, which we will denote by  $x$ . Later we will see how to choose  $x$ , but for now just think of it as a random element of  $A$ . We then partition  $A$  into three parts.  $A[q]$  contains the element  $x$ , subarray  $A[p..q-1]$  will contain all the elements that are less than  $x$ , and  $A[q+1..r]$ , will contain all the element that are greater than  $x$ . (Recall that we assumed that all the elements are distinct.) Within each subarray, the items may appear in any order (see Fig. 97).

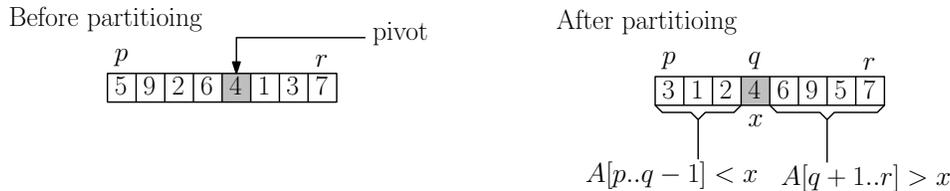


Fig. 97: Partitioning

It is easy to see that the rank of the pivot  $x$  is  $q - p + 1$  in  $A[p..r]$ . Let  $x\text{Rank} = q - p + 1$ . If  $k = x\text{Rank}$ , then the pivot is the  $k$ th smallest, and we may just return it. If  $k < x\text{Rank}$ , then we know that we need to recursively search in  $A[p..q-1]$  and if  $k > x\text{Rank}$  then we need to recursively search  $A[q+1..r]$ . In this latter case we have eliminated  $q$  smaller elements, so we want to find the element of rank  $k - q$ . Here is the complete pseudocode (see Fig. 98).

---

```
Selection by the Sieve Technique
```

```

Select(array A, int p, int r, int k) { // return kth smallest of A[p..r]
  if (p == r) return A[p]           // only 1 item left, return it
  else {
    x = ChoosePivot(A, p, r)        // choose the pivot element
    q = Partition(A, p, r, x)       // <A[p..q-1], x, A[q+1..r]>
    xRank = q - p + 1              // rank of the pivot
    if (k == xRank) return x       // the pivot is the kth smallest
    else if (k < xRank)
      return Select(A, p, q-1, k) // select from left
    else
      return Select(A, q+1, r, k-xRank) // select from right
  }
}

```

---

Notice that this algorithm satisfies the basic form of a sieve algorithm. It analyzes the data (by choosing the pivot element and partitioning) and it eliminates some part of the data set, and recurses on the

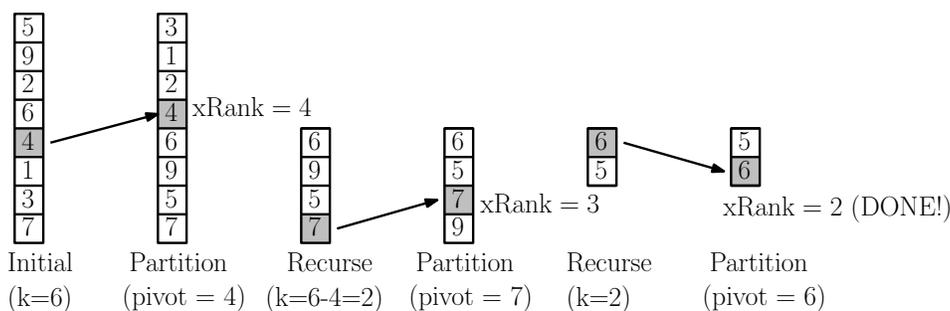


Fig. 98: Selection Algorithm.

rest. When  $k = \text{xRank}$  then we get lucky and eliminate everything. Otherwise we either eliminate the pivot and the right subarray or the pivot and the left subarray.

We will discuss the details of choosing the pivot and partitioning later, but assume for now that they both take  $\Theta(n)$  time. The question that remains is how many elements did we succeed in eliminating? If  $x$  is the largest or smallest element in the array, then we may only succeed in eliminating one element with each phase. In fact, if  $x$  is one of the smallest elements of  $A$  or one of the largest, then we get into trouble, because we may only eliminate it and the few smaller or larger elements of  $A$ . Ideally  $x$  should have a rank that is neither too large nor too small.

Let us suppose for now (optimistically) that we are able to design the procedure `Choose.Pivot` in such a way that is eliminates exactly half the array with each phase, meaning that we recurse on the remaining  $n/2$  elements. This would lead to the following recurrence.

$$T(n) = \begin{cases} 1 & \text{if } n = 1, \\ T(n/2) + n & \text{otherwise.} \end{cases}$$

We can solve this either by expansion (iteration) or the Master Theorem. If we expand this recurrence level by level we see that we get the summation

$$T(n) = n + \frac{n}{2} + \frac{n}{4} + \dots \leq \sum_{i=0}^{\infty} \frac{n}{2^i} = n \sum_{i=0}^{\infty} \frac{1}{2^i}.$$

Recall the formula for the infinite geometric series. For any  $c$  such that  $|c| < 1$ ,  $\sum_{i=0}^{\infty} c^i = 1/(1 - c)$ . Using this we have

$$T(n) \leq 2n \in O(n).$$

(This only proves the upper bound on the running time, but it is easy to see that it takes at least  $\Omega(n)$  time, so the total running time is  $\Theta(n)$ .)

This is a bit counterintuitive. Normally you would think that in order to design a  $\Theta(n)$  time algorithm you could only make a single, or perhaps a constant number of passes over the data set. In this algorithm we make many passes (it could be as many as  $\lg n$ ). However, because we eliminate a constant fraction of elements with each phase, we get this convergent geometric series in the analysis, which shows that the total running time is indeed linear in  $n$ . This lesson is well worth remembering. It is often possible to achieve running times in ways that you would not expect.

Note that the assumption of eliminating half was not critical. If we eliminated even one per cent, then the recurrence would have been  $T(n) = T(99n/100) + n$ , and we would have gotten a geometric series involving  $99/100$ , which is still less than 1, implying a convergent series. Eliminating *any* constant fraction would have been good enough.

**Choosing the Pivot:** There are two issues that we have left unresolved. The first is how to choose the pivot element, and the second is how to partition the array. Both need to be solved in  $\Theta(n)$  time. The second problem is a rather easy programming exercise. Later, when we discuss QuickSort, we will discuss partitioning in detail.

For the rest of the lecture, let's concentrate on how to choose the pivot. Recall that before we said that we might think of the pivot as a random element of  $A$ . Actually this is not such a bad idea. Let's see why.

The key is that we want the procedure to eliminate at least some constant fraction of the array after each partitioning step. Let's consider the top of the recurrence, when we are given  $A[1..n]$ . Suppose that the pivot  $x$  turns out to be of rank  $q$  in the array. The partitioning algorithm will split the array into  $A[1..q-1] < x$ ,  $A[q] = x$  and  $A[q+1..n] > x$ . If  $k = q$ , then we are done. Otherwise, we need to search one of the two subarrays. They are of sizes  $q-1$  and  $n-q$ , respectively. The subarray that contains the  $k$ th smallest element will generally depend on what  $k$  is, so in the worst case,  $k$  will be chosen so that we have to recurse on the larger of the two subarrays. Thus if  $q > n/2$ , then we may have to recurse on the left subarray of size  $q-1$ , and if  $q < n/2$ , then we may have to recurse on the right subarray of size  $n-q$ . In either case, we are in trouble if  $q$  is very small, or if  $q$  is very large.

If we could select  $q$  so that it is roughly of middle rank, then we will be in good shape. For example, if  $n/4 \leq q \leq 3n/4$ , then the larger subarray will never be larger than  $3n/4$ . Earlier we said that we might think of the pivot as a random element of the array  $A$ . Actually this works pretty well in practice. The reason is that roughly half of the elements lie between ranks  $n/4$  and  $3n/4$ , so picking a random element as the pivot will succeed about half the time to eliminate at least  $n/4$ . Of course, we might be continuously unlucky, but a careful analysis will show that the expected running time is still  $\Theta(n)$ . We will return to this later.

Instead, we will describe a rather complicated method for computing a pivot element that achieves the desired properties. Recall that we are given an array  $A[1..n]$ , and we want to compute an element  $x$  whose rank is (roughly) between  $n/4$  and  $3n/4$ . We will have to describe this algorithm at a very high level, since the details are rather involved. Here is the description for `Select_Pivot`:

**Groups of 5:** Partition  $A$  into groups of 5 elements, e.g.  $A[1..5]$ ,  $A[6..10]$ ,  $A[11..15]$ , etc. There will be exactly  $m = \lceil n/5 \rceil$  such groups (the last one might have fewer than 5 elements). This can easily be done in  $\Theta(n)$  time.

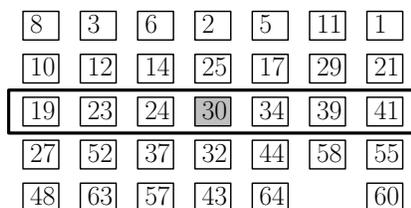
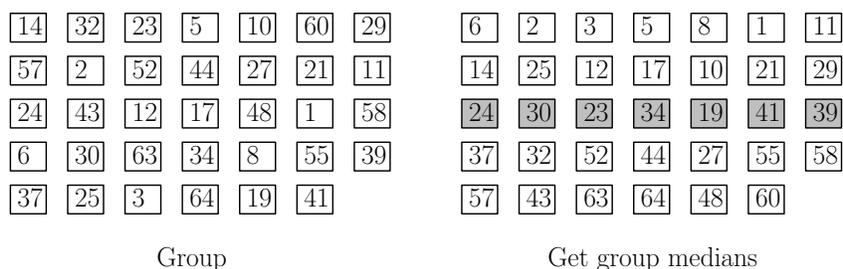
**Group medians:** Compute the median of each group of 5. There will be  $m$  group medians. We do not need an intelligent algorithm to do this, since each group has only a constant number of elements. For example, we could just BubbleSort each group and take the middle element. Each will take  $\Theta(1)$  time, and repeating this  $\lceil n/5 \rceil$  times will give a total running time of  $\Theta(n)$ . Copy the group medians to a new array  $B$ .

**Median of medians:** Compute the median of the group medians. For this, we will have to call the selection algorithm recursively on  $B$ , e.g. `Select(B, 1, m, k)`, where  $m = \lceil n/5 \rceil$ , and  $k = \lfloor (m+1)/2 \rfloor$ . Let  $x$  be this median of medians. Return  $x$  as the desired pivot.

The algorithm is illustrated in the figure below. To establish the correctness of this procedure, we need to argue that  $x$  satisfies the desired rank properties.

**Lemma:** The element  $x$  is of rank at least  $n/4$  and at most  $3n/4$  in  $A$ .

**Proof:** We will show that  $x$  is of rank at least  $n/4$ . The other part of the proof is essentially symmetrical. To do this, we need to show that there are at least  $n/4$  elements that are less than or equal to  $x$ . This is a bit complicated, due to the floor and ceiling arithmetic, so to simplify things we will assume that  $n$  is evenly divisible by 5. Consider the groups shown in the tabular form above. Observe that at least half of the group medians are less than or equal to  $x$ . (Because  $x$  is



Get median of medians  
(Sorting of group medians is not really performed)

Fig. 99: Choosing the Pivot. (30 is the final pivot.)

their median.) And for each group median, there are three elements that are less than or equal to this median within its group (because it is the median of its group). Therefore, there are at least  $3((n/5)/2) = 3n/10 \geq n/4$  elements that are less than or equal to  $x$  in the entire array.

**Analysis:** The last order of business is to analyze the running time of the overall algorithm. We achieved the main goal, namely that of eliminating a constant fraction (at least  $1/4$ ) of the remaining list at each stage of the algorithm. The recursive call in `Select()` will be made to list no larger than  $3n/4$ . However, in order to achieve this, within `Select_Pivot()` we needed to make a recursive call to `Select()` on an array  $B$  consisting of  $\lceil n/5 \rceil$  elements. Everything else took only  $\Theta(n)$  time. As usual, we will ignore floors and ceilings, and write the  $\Theta(n)$  as  $n$  for concreteness. The running time is

$$T(n) \leq \begin{cases} 1 & \text{if } n = 1, \\ T(n/5) + T(3n/4) + n & \text{otherwise.} \end{cases}$$

This is a very strange recurrence because it involves a mixture of different fractions ( $n/5$  and  $3n/4$ ). This mixture will make it impossible to use the Master Theorem, and difficult to apply iteration. However, this is a good place to apply constructive induction. We know we want an algorithm that runs in  $\Theta(n)$  time.

**Theorem:** There is a constant  $c$ , such that  $T(n) \leq cn$ .

**Proof:** (by strong induction on  $n$ )

**Basis:** ( $n = 1$ ) In this case we have  $T(n) = 1$ , and so  $T(n) \leq cn$  as long as  $c \geq 1$ .

**Step:** We assume that  $T(n') \leq cn'$  for all  $n' < n$ . We will then show that  $T(n) \leq cn$ . By definition we have

$$T(n) = T(n/5) + T(3n/4) + n.$$

Since  $n/5$  and  $3n/4$  are both less than  $n$ , we can apply the induction hypothesis, giving

$$\begin{aligned} T(n) &\leq c\frac{n}{5} + c\frac{3n}{4} + n = cn\left(\frac{1}{5} + \frac{3}{4}\right) + n \\ &= cn\frac{19}{20} + n = n\left(\frac{19c}{20} + 1\right). \end{aligned}$$

This last expression will be  $\leq cn$ , provided that we select  $c$  such that  $c \geq (19c/20) + 1$ . Solving for  $c$  we see that this is true provided that  $c \geq 20$ .

Combining the constraints that  $c \geq 1$ , and  $c \geq 20$ , we see that by letting  $c = 20$ , we are done.

A natural question is why did we pick groups of 5? If you look at the proof above, you will see that it works for any value that is strictly greater than 4. (You might try it replacing the 5 with 3, 4, or 6 and see what happens.)

## Lecture 29: Long Integer Multiplication

**Long Integer Multiplication:** The following little algorithm shows a bit more about the surprising applications of divide-and-conquer. The problem that we want to consider is how to perform arithmetic on long integers, and multiplication in particular. The reason for doing arithmetic on long numbers stems from cryptography. Most techniques for encryption are based on number-theoretic techniques. For example, the character string to be encrypted is converted into a sequence of numbers, and encryption keys are stored as long integers. Efficient encryption and decryption depends on being able to perform arithmetic on long numbers, typically containing hundreds of digits.

Addition and subtraction on large numbers is relatively easy. If  $n$  is the number of digits, then these algorithms run in  $\Theta(n)$  time. (Go back and analyze your solution to the problem on Homework 1). But the standard algorithm for multiplication runs in  $\Theta(n^2)$  time, which can be quite costly when lots of long multiplications are needed.

This raises the question of whether there is a more efficient way to multiply two very large numbers. It would seem surprising if there were, since for centuries people have used the same algorithm that we all learn in grade school. In fact, we will see that it is possible.

**Divide-and-Conquer Algorithm:** We know the basic grade-school algorithm for multiplication. We normally think of this algorithm as applying on a digit-by-digit basis, but if we partition an  $n$  digit number into two “super digits” with roughly  $n/2$  each into longer sequences, the same multiplication rule still applies.

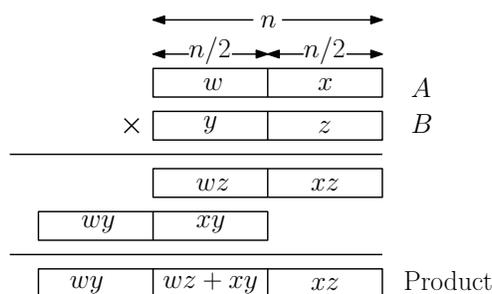


Fig. 100: Long integer multiplication.

To avoid complicating things with floors and ceilings, let's just assume that the number of digits  $n$  is a power of 2. Let  $A$  and  $B$  be the two numbers to multiply. Let  $A[0]$  denote the least significant digit and let  $A[n - 1]$  denote the most significant digit of  $A$ . Because of the way we write numbers, it is more natural to think of the elements of  $A$  as being indexed in decreasing order from left to right as  $A[n - 1..0]$  rather than the usual  $A[0..n - 1]$ .

Let  $m = n/2$ . Let

$$\begin{aligned} w &= A[n - 1..m] & x &= A[m - 1..0] \quad \text{and} \\ y &= B[n - 1..m] & z &= B[m - 1..0]. \end{aligned}$$

If we think of  $w$ ,  $x$ ,  $y$  and  $z$  as  $n/2$  digit numbers, we can express  $A$  and  $B$  as

$$\begin{aligned} A &= w \cdot 10^m + x \\ B &= y \cdot 10^m + z, \end{aligned}$$

and their product is

$$\text{mult}(A, B) = \text{mult}(w, y)10^{2m} + (\text{mult}(w, z) + \text{mult}(x, y))10^m + \text{mult}(x, z).$$

The operation of multiplying by  $10^m$  should be thought of as simply shifting the number over by  $m$  positions to the right, and so is not really a multiplication. Observe that all the additions involve numbers involving roughly  $n/2$  digits, and so they take  $\Theta(n)$  time each. Thus, we can express the multiplication of two long integers as the result of four products on integers of roughly half the length of the original, and a constant number of additions and shifts, each taking  $\Theta(n)$  time. This suggests that if we were to implement this algorithm, its running time would be given by the following recurrence

$$T(n) = \begin{cases} 1 & \text{if } n = 1, \\ 4T(n/2) + n & \text{otherwise.} \end{cases}$$

If we apply the Master Theorem, we see that  $a = 4$ ,  $b = 2$ ,  $k = 1$ , and  $a > b^k$ , implying that Case 1 holds and the running time is  $\Theta(n^{\lg 4}) = \Theta(n^2)$ . Unfortunately, this is no better than the standard algorithm.

**Faster Divide-and-Conquer Algorithm:** Even though this exercise appears to have gotten us nowhere, it actually has given us an important insight. It shows that the critical element is the number of multiplications on numbers of size  $n/2$ . The number of additions (as long as it is a constant) does not affect the running time. So, if we could find a way to arrive at the same result algebraically, but by trading off multiplications in favor of additions, then we would have a more efficient algorithm. (Of course, we cannot simulate multiplication through repeated additions, since the number of additions must be a constant, independent of  $n$ .)

The key turns out to be an algebraic “trick”. The quantities that we need to compute are  $C = wy$ ,  $D = xz$ , and  $E = (wz + xy)$ . Above, it took us four multiplications to compute these. However, observe that if instead we compute the following quantities, we can get everything we want, using only three multiplications (but with more additions and subtractions).

$$\begin{aligned} C &= \text{mult}(w, y) \\ D &= \text{mult}(x, z) \\ E &= \text{mult}((w + x), (y + z)) - C - D = (wy + wz + xy + xz) - wy - xz = (wz + xy). \end{aligned}$$

Finally we have

$$\text{mult}(A, B) = C \cdot 10^{2m} + E \cdot 10^m + D.$$

Altogether we perform 3 multiplications, 4 additions, and 2 subtractions all of numbers with  $n/2$  digits. We still need to shift the terms into their proper final positions. The additions, subtractions, and shifts take  $\Theta(n)$  time in total. So the total running time is given by the recurrence:

$$T(n) = \begin{cases} 1 & \text{if } n = 1, \\ 3T(n/2) + n & \text{otherwise.} \end{cases}$$

Now when we apply the Master Theorem, we have  $a = 3$ ,  $b = 2$  and  $k = 1$ , yielding  $T(n) \in \Theta(n^{\lg 3}) \approx \Theta(n^{1.585})$ .

Is this really an improvement? This algorithm carries a larger constant factor because of the overhead of recursion and the additional arithmetic operations. But asymptotics says that if  $n$  is large enough, then this algorithm will be superior. For example, if we assume that the clever algorithm has overheads that are 5 times greater than the simple algorithm (e.g.  $5n^{1.585}$  versus  $n^2$ ) then this algorithm beats the simple algorithm for  $n \geq 50$ . If the overhead was 10 times larger, then the crossover would occur for  $n \geq 260$ . Although this may seem like a very large number, recall that in cryptography applications, encryption keys of this length and longer are quite reasonable.

## Lecture 30: Divide and Conquer: Mergesort and Inversion Counting

**Divide and Conquer:** So far, we have been studying a basic algorithm design technique called greedy algorithms. Today, we begin study of a different technique, called *divide and conquer*. The ancient Roman rulers understood this principle well (although they were probably not thinking about algorithms at the time). You divide your enemies (by getting them to distrust each other) and then conquer them one by one. In algorithm design, the idea is to take a problem on a large input, break the input into smaller pieces, solve the pieces individually (usually recursively), and then combine the piecewise solutions into a global solution.

Summarizing, the main elements to a divide-and-conquer solution are

- *Divide* (the problem into a small number of pieces),
- *Conquer* (solve each piece, by applying divide-and-conquer recursively to it), and
- *Combine* (the pieces together into a global solution).

There are a huge number computational problems that can be solved efficiently using divide-and-conquer. Divide-and-conquer algorithms typically involve recursion, since this is usually the most natural way to deal with the “conquest” part of the algorithm. Analyzing the running times of recursive programs is usually done by solving a *recurrence*.

**MergeSort:** Perhaps the simplest example of a divide-and-conquer algorithm is MergeSort. I am sure you are familiar with this algorithm, but for the sake of completeness, let’s recall how it works. We are given an sequence of  $n$  numbers, which we denote by  $A$ . The objective is to permute the array elements into non-decreasing order.  $A$  may be stored as an array or a linked list. Let’s not worry about these implementaiton details for now. We will need to assume that, whatever representation we use, we can determine the lists size in constant time, and we can enumerate the elements from left to right.

Here is the basic structure of MergeSort. Let  $\text{size}(A)$  denote the number of elements of  $A$ .

**Basis case:** If  $\text{size}(A) = 1$ , then the array is trivially sorted and we are done.

**General case:** Otherwise:

**Divide:** Split  $A$  into two subsequences, each of size roughly  $n/2$ . (More precisely, one will be of size  $\lfloor n/2 \rfloor$  and the other of size  $\lceil n/2 \rceil$ .)

**Conquer:** Sort each subsequence (by calling MergeSort recursively on each).

**Combine:** Merge the two sorted subsequences into a single sorted list.

**MergeSort:** The key to the algorithm is the merging process. Let us assume inductively that the sequence has been split into two, which are presented as two subarrays,  $A[p..m]$  and  $A[m+1..r]$ , each of which has been sorted. The merging process copies the elements of these two subarrays into temporary array  $B$ . We maintain two indices  $i$  and  $j$ , indicating the current elements of the left and right subarrays, respectively. At each step, we copy whichever element is smaller  $A[i]$  or  $A[j]$  to the next position of  $B$ . (Ties are broken in favor of  $A$ .) See Fig. 101.)

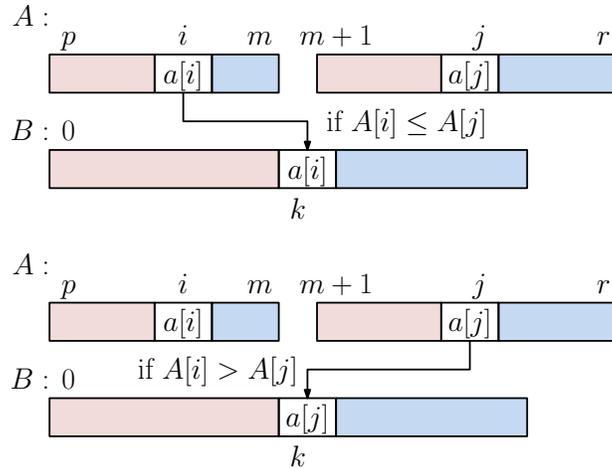


Fig. 101: Merging two sorted lists.

The two code blocks below present the MergeSort algorithm and the merging utility, which merges two sorted lists. Assuming that the input is in the array  $A[1..n]$ , the initial call is  $\text{MergeSort}(A, 1, n)$ .

---

```

MergeSort(A, p, r) {                                     // sort A[p..r]
    if (p < r) {                                        // we have at least 2 items
        m = (p + r)/2                                  // midpoint
        MergeSort(A, p, m)                             // sort the left half
        MergeSort(A, m+1, r)                           // sort the right half
        merge(A, p, m, r)                              // merge the two halves
    }
}

merge(A, p, m, r) {                                    // merge A[p..m] and A[m+1..r]
    new array B[0..r-p]
    i = p; j = m+1; k = 0;                             // initialize indices
    while (i <= m and j <= r) {                       // while both are nonempty
        if (A[i] <= A[j]) B[k++] = A[i++]           // next item from left
        else B[k++] = A[j++]                         // next item from right
    }
    while (i <= m) B[k++] = A[i++]                   // copy any extras to B
    while (j <= r) B[k++] = A[j++]
    for (k = 0 to r-p) A[p+k] = B[k]                 // copy B back to A
}

```

---

This completes the description of the algorithm. Observe that of the last two while-loops in the merge procedure, only one will be executed. (Do you see why?) Another question worth considering is the following. Suppose that in the merge function, the statement “ $A[i] \leq A[j]$ ” had instead been written “ $A[i] < A[j]$ ”? Would the algorithm still be correct? Can you see any reason for preferring one version over the other? (Hint: Consider what happens when  $A$  contains duplicate copies of the same element.)

Fig. 102 shows an example of the execution of MergeSort. The dividing part of the algorithm is shown on the left and the merging part is shown on the right.

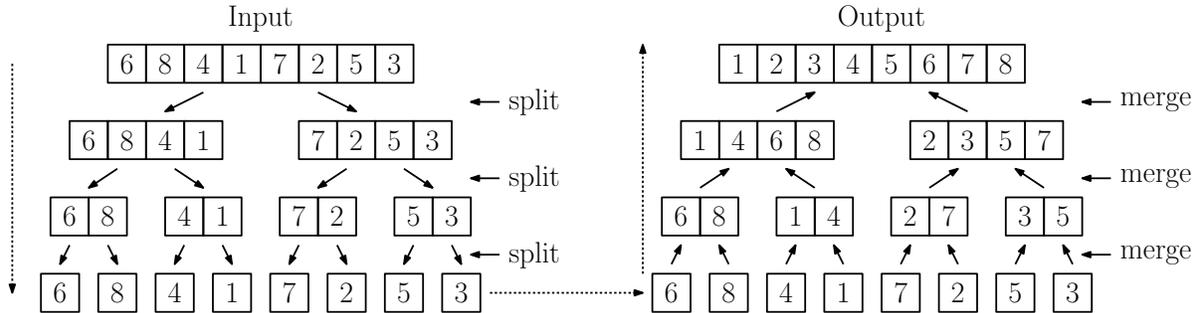


Fig. 102: MergeSort example.

**Analysis:** Next, let us analyze the running time of MergeSort. First observe that the running time of the procedure  $\text{merge}(A, p, m, r)$  is easily seen to be  $O(r - p + 1)$ , that is, it is proportional to the total size of the two lists being merged. The reason is that, each time through the loop, we succeed in copying one element from  $A[p..r]$  to the final output.

Now, how do we describe the running time of the entire MergeSort algorithm? We will do this through the use of a *recurrence*, that is, a function that is defined recursively in terms of itself. Let's see how to apply this to MergeSort. Let  $T(n)$  denote the worst case running time of MergeSort on an input of length  $n \geq 1$ . First observe that if we call MergeSort with a list containing a single element, then the running time is a constant. Since we are ignoring constant factors, we can just write  $T(n) = 1$ . When we call MergeSort with a list of length  $n \geq 2$ , e.g.  $\text{merge}(A, p, r)$ , where  $r - p + 1 = n$ , the algorithm first computes  $m = \lfloor (p + r)/2 \rfloor$ . The subarray  $A[p..r]$ , which contains  $r - p + 1$  elements. We'll ignore the floors and ceilings, and simply declare that each subarray is of size  $n/2$ . Thus, we have

$$T(n) = \begin{cases} 1 & \text{if } n = 1, \\ 2T\left(\frac{n}{2}\right) + n & \text{otherwise.} \end{cases}$$

**Solving the Recurrence:** In order to complete the analysis, we need to solve the above recurrence. There are a few ways to solve recurrences. My favorite method is to apply repeated expansion until a pattern emerges. Then, express the result in terms of the number of iterations performed.

$$\begin{aligned} T(n) &= 2T(n/2) + n \\ &= 2(2T(n/4) + (n/2)) + n = 4T(n/4) + 2n \\ &= 4(2T(n/8) + n/4) + 2n = 8T(n/8) + 3n \\ &= \dots \\ &= 2^k T(n/2^k) + kn. \end{aligned}$$

The above expression as a function of  $k$  is messy, but it is useful. We know that  $T(1) = 1$ . To use that

fact, we need to determine what value to set  $k$  so that  $n/2^k = 1$ . Therefore, we have  $k = \lg n$ .<sup>19</sup> By substituting this value for  $k$ , we have  $T(n/2^k) = T(1) = 1$  and plugging this into the above formula, we obtain

$$T(n) = 2^{\lg n} \cdot T(1) + n \lg n = n \cdot 1 + n \lg n = O(n \log n),$$

Therefore, the running time of MergeSort is  $O(n \log n)$ .

Many of the recurrences that arise in divide-and-conquer algorithms have a similar structure. The following theorem is useful for compute asymptotic bounds for these recurrences.

**Theorem:** (Master Theorem) Let  $a \geq 1$ ,  $b > 1$  be constants and let  $T(n)$  be the recurrence

$$T(n) = aT\left(\frac{n}{b}\right) + n^k,$$

defined for  $n \geq 0$ . (Let us assume that  $n$  is a power of  $b$ . This doesn't affect the asymptotics. The basis case,  $T(1)$  can be any constant value.) Then:

**Case 1:** if  $a > b^k$ , then  $T(n) \in \Theta(n^{\log_b a})$

**Case 2:** if  $a = b^k$ , then  $T(n) \in \Theta(n^k \log n)$

**Case 3:** if  $a < b^k$ , then  $T(n) \in \Theta(n^k)$ .

**Inversion Counting:** Let's consider a variant on this. Although the problem description does not appear to have anything to do with sorting or Mergesort, we will see that the solutions to these problems are closely related. Suppose that you are given two rank ordered lists of preferences. For example, suppose that you and bunch of your friends are given a list of 50 popular movies, and you are rank order them from most favorite to least favorite. After this exercise, you want to know which people tended to rank movies in roughly the same way that you did. Here is an example:

Movie Title	Alice	Bob	Carol
Gone with the Wind	1	4	6
Citizen Kane	2	1	8
The Seven Samurai	3	3	4
The Godfather	4	2	1
Titanic	5	5	7
My Cousin Vinny	6	7	2
Star Wars	7	8	5
Plan 9 from Outer Space	8	6	3

Given two such lists, how would you determine their degree of similarity? Here is one possible approach. Given two lists of preferences,  $L_1$  and  $L_2$ , define an *inversion* to be a pair of movies  $x$  and  $y$ , such that  $L_1$  has  $x$  before  $y$  and  $L_2$  has  $y$  before  $x$ . Since there are  $\binom{n}{2} = n(n-1)/2$  unordered pairs, the maximum number of inversions is  $\binom{n}{2}$ , which is  $O(n^2)$ . If the two rankings are the same, then there are no inversions. Thus, the number of inversions can be seen as one possible measure of similarity between two lists of  $n$  numbers. (An example is shown in in Fig. 103.)

We can reduce this problem from one involving two lists to one involving just one. In particular, assume that the first list consists of the sequence  $\langle 1, \dots, n \rangle$ . Let the other list be denoted by  $\langle a_1, \dots, a_n \rangle$ . (More generally, you can relabel the elements so that the index of the element is its position in the first list.) An *inversion* is a pair of indices  $(i, j)$  such that  $i < j$ , but  $a_i > a_j$ . Given a list of  $n$  (distinct) numbers, our objective is to count the number of inversions.

Naively, we can solve this problem in  $O(n^2)$  time. For each  $a_i$ , we search all  $i + 1 \leq j \leq n$ , and increment a counter for every  $j$  such that  $a_i > a_j$ . We will investigate a more efficient method based on divide-and-conquer.

<sup>19</sup>Recall that "lg" means logarithm base 2. This worked because we ignored the floors and ceilings, and hence, treated  $n$  as if it were a power of 2. More accurately, we have  $k = \lceil \lg n \rceil$ .

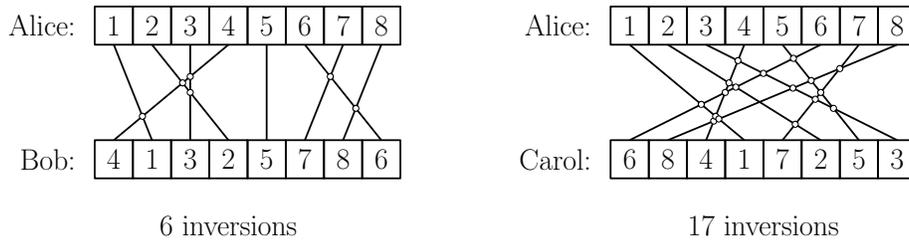


Fig. 103: Movie preferences and inversions.

**Divide-and-conquer solution:** How would we approach solving this problem using divide-and-conquer? Here is one way of doing it:

**Basis case:** If  $\text{size}(A) = 1$ , then there are no inversions.

**General case:** Otherwise:

**Divide:** Split  $A$  into two subsequences, each of size roughly  $n/2$ .

**Conquer:** Compute the number of inversions *within* each of the subsequences.

**Combine:** Count the number of inversions occurring *between* the two sequences.

The computation of the inversions within each subsequence is solved by recursion. The key to an efficient implementation of the algorithm is the step where we count the inversions between the two lists. It will be much easier to count inversions if we first sort the list. In fact, our approach will be to both sort and count inversions at the same time.

Let us assume that the input is given as an array  $A[p..r]$ . Let us assume inductive that it has been split into two subarrays,  $A[p..m]$  and  $A[m+1..r]$ , each of which has already been sorted. During the merging process, we maintain two indices  $i$  and  $j$ , indicating the current elements of the left and right subarrays, respectively (see Fig. 104).

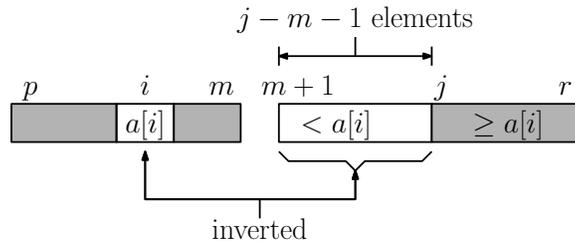


Fig. 104: Counting inversions when  $A[i] \leq A[j]$ .

Whenever  $A[i] > A[j]$  the algorithm advances  $j$ . It follows, therefore, that if  $A[i] \leq A[j]$ , then every element of the subarray  $A[m+1..j-1]$  is strictly smaller than  $A[i]$ . Since the elements of the left subarray appear in the original array before all the elements of the right subarray, it follows that  $A[i]$  generates an inversion with *all* the elements of the subarray  $A[m+1..j-1]$ . The number of elements in this subarray is  $(j-1) - (m+1) + 1 = j - m - 1$ . Therefore, before when we process  $A[i]$ , we increment an inversion counter by  $j - m - 1$ .

The other part of the code that is affected is when we copy elements from the end of the left subarray to the final array. In this case, each element that is copied generates an inversion with respect to all the elements of the right subarray, that is,  $A[m+1..r]$ . There are  $r - m$  such elements. We add this value to the inversion counter.

The algorithm is modeled on the same pseudo-code as that used for MergeSort and is presented in the following code block. Assuming that the input is stored in the array  $A[1..n]$ , the initial call is  $\text{InvCount}(A, 1, n)$ .

Inversion Counting

---

```

InvCount(A, p, r) {
    if (p >= r) return 0
    m = (p + r)/2
    x1 = InvCount(A, p, m)
    x2 = InvCount(A, m+1, r)
    x3 = invMerge(A, p, m, r)
    return x1 + x2 + x3
}

invMerge(A, p, m, r) {
    new array B[0..r-p]
    i = p; j = m+1; k = 0;
    ct = 0
    while (i <= m and j <= r) {
        if (A[i] <= A[j]) {
            B[k++] = A[i++]
            ct += j - m - 1
        }
        else B[k++] = A[j++]
    }
    while (i <= m) {
        B[k++] = A[i++]
        ct += r - m
    }
    while (j <= r) B[k++] = A[j++]
    for (k = 0 to r-p) A[p+k] = B[k]
}

```

---

This approach is illustrated in Fig. 105. Observe that inversions are counted in the merging process (shown as small white circles in the figure).

## Lecture 31: Divide-and-Conquer: Closest Pair

**Closest Pair:** Today, we consider another application of divide-and-conquer, which comes from the field of computational geometry. We are given a set  $P$  of  $n$  points in the plane, and we wish to find the closest pair of points  $p, q \in P$  (see Fig. 106(a)). This problem arises in a number of applications. For example, in air-traffic control, you may want to monitor planes that come too close together, since this may indicate a possible collision. Recall that, given two points  $p = (p_x, p_y)$  and  $q = (q_x, q_y)$ , their (Euclidean) distance is

$$\|pq\| = \sqrt{(p_x - q_x)^2 + (p_y - q_y)^2}.$$

Clearly, this problem can be solved by brute force in  $O(n^2)$  time, by computing the distance between each pair, and returning the smallest. Today, we will present an  $O(n \log n)$  time algorithm, which is based a clever use of divide-and-conquer.

Before getting into the solution, it is worth pointing out a simple strategy that fails to work. If two points are very close together, then clearly both their  $x$ -coordinates and their  $y$ -coordinates are close

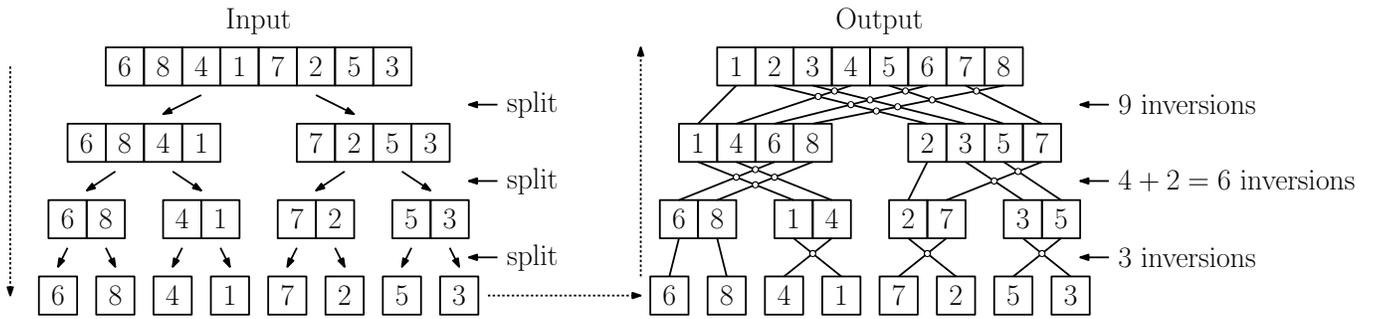


Fig. 105: Inversion counting by divide and conquer.

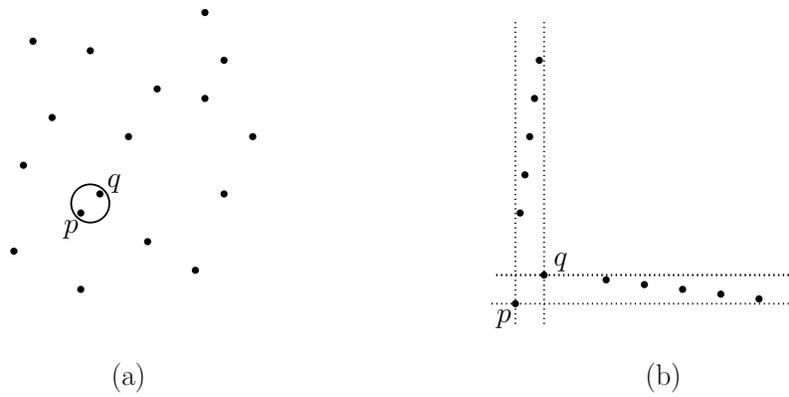


Fig. 106: (a) The closest pair problem and (b) why sorting on  $x$ - or  $y$ -alone doesn't work.

together. So, how about if we *sort* the points based on their  $x$ -coordinates and, for each point of the set, we'll consider just *nearby* points in the list. It would seem that (subject to figuring out exactly what “nearby” means) such a strategy might be made to work. The problem is that it could fail miserably. In particular, consider the point set of Fig. 106(b). The points  $p$  and  $q$  are the closest points, but we can place an arbitrarily large number of points between them in terms of their  $x$ -coordinates. We need to separate these points sufficiently far in terms of their  $y$ -coordinates that  $p$  and  $q$  remain the closest pair. As a result, the positions of  $p$  and  $q$  can be arbitrarily far apart in the sorted order. Of course, we can do the same with respect to the  $y$ -coordinate. Clearly, we cannot focus on one coordinate alone.<sup>20</sup>

**Divide-and-Conquer Algorithm:** Let us investigate how to design an  $O(n \log n)$  time divide-and-conquer approach to the problem. The input consists of a set of points  $P$ , represented, say, as an array of  $n$  elements, where each element stores the  $(x, y)$  coordinates of the point. (For simplicity, let's assume there are no duplicate  $x$ -coordinates.) The output will consist of a single number, being the closest distance. It is easy to modify the algorithm to also produce the pair of points that achieves this distance.

For reasons that will become clear later, in order to implement the algorithm efficiently, it will be helpful to begin by *presorting* the points, both with respect to their  $x$ - and  $y$ -coordinates. Let  $P_x$  be an array of points sorted by  $x$ , and let  $P_y$  be an array of points sorted by  $y$ . We can compute these sorted arrays in  $O(n \log n)$  time. Note that this initial sorting is done only *once*. In particular, the recursive calls do not repeat the sorting process.

Like any divide-and-conquer algorithm, after the initial basis case, our approach involves three basic elements: divide, conquer, and combine.

**Basis:** If  $|P| \leq 3$ , then just solve the problem by brute force in  $O(1)$  time.

**Divide:** Otherwise, partition the points into two subarrays  $P_L$  and  $P_R$  based on their  $x$ -coordinates.

In particular, imagine a vertical line  $\ell$  that splits the points roughly in half (see Fig. 107). Let  $P_L$  be the points to the left of  $\ell$  and  $P_R$  be the points to the right of  $\ell$ .

In the same way that we represented  $P$  using two sorted arrays, we do the same for  $P_L$  and  $P_R$ . Since we have presorted  $P_x$  by  $x$ -coordinates, we can determine the median element for  $\ell$  in constant time. After this, we can partition each of arrays  $P_x$  and  $P_y$  in  $O(n)$  time each.

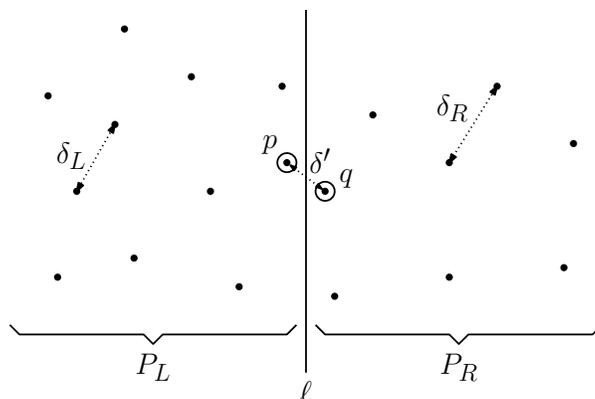


Fig. 107: Divide-and-conquer closest pair algorithm.

<sup>20</sup>While the above example shows that sorting along any one coordinate axis may fail, there is a variant of this strategy that can be used for computing nearest neighbors approximately. This approach is based on the observation that if two points are close together, their projections onto a *randomly oriented vector* will be close, and if they are far apart, their projections onto a randomly oriented vector will be far apart in expectation. This observation underlies a popular nearest neighbor algorithm called *locality sensitive hashing*.

**Conquer:** Compute the closest pair *within* each of the subsets  $P_L$  and  $P_R$  each, by invoking the algorithm recursively. Let  $\delta_L$  and  $\delta_R$  be the closest pair distances in each case (see Fig. 107). Let  $\delta = \min(\delta_L, \delta_R)$ .

**Combine:** Note that  $\delta$  is not necessarily the final answer, because there may be two points that are very close to one another but are on opposite sides of  $\ell$ . To complete the algorithm, we want to determine the closest pair of points *between* the sets, that is, the closest points  $p \in P_L$  and  $q \in P_R$  (see Fig. 107). Since we already have an upper bound  $\delta$  on the closest pair, it suffices to solve the following *restricted problem*: if the closest pair  $(p, q)$  are within distance  $\delta$ , then we will return such a pair, otherwise, we may return any pair. (This restriction is very important to the algorithm's efficiency.) In the next section, we'll show how to solve this restricted problem in  $O(n)$  time. Given the closest such pair  $(p, q)$ , let  $\delta' = \|pq\|$ . We return  $\min(\delta, \delta')$  as the final result.

Assuming that we can solve the “Combine” step in  $O(n)$  time, it will follow that the algorithm's running time is given by the recurrence  $T(n) = 2T(n/2) + n$ , and (as in Mergesort) the overall running time is  $O(n \log n)$ , as desired.

**Closest Pair Between the Sets:** To finish up the algorithm, we need to compute the closest pair  $p$  and  $q$ , where  $p \in P_L$  and  $q \in P_R$ . As mentioned above, because we already know of the existence of two points within distance  $\delta$  of each other, this algorithm is allowed to fail, if there is no such pair that is closer than  $\delta$ . The input to our algorithm consists of the point set  $P$ , the  $x$ -coordinate of the vertical splitting line  $\ell$ , and the value of  $\delta = \min(\delta_L, \delta_R)$ . Recall that our goal is to do this in  $O(n)$  time.

This is where the real creativity of the algorithm enters. Observe that if such a pair of points exists, we may assume that both points lie within distance  $\delta$  of  $\ell$ , for otherwise the resulting distance would exceed  $\delta$ . Let  $S$  denote this subset of  $P$  that lies within a vertical strip of width  $2\delta$  centered about  $\ell$  (see Fig. 108(a)).<sup>21</sup>

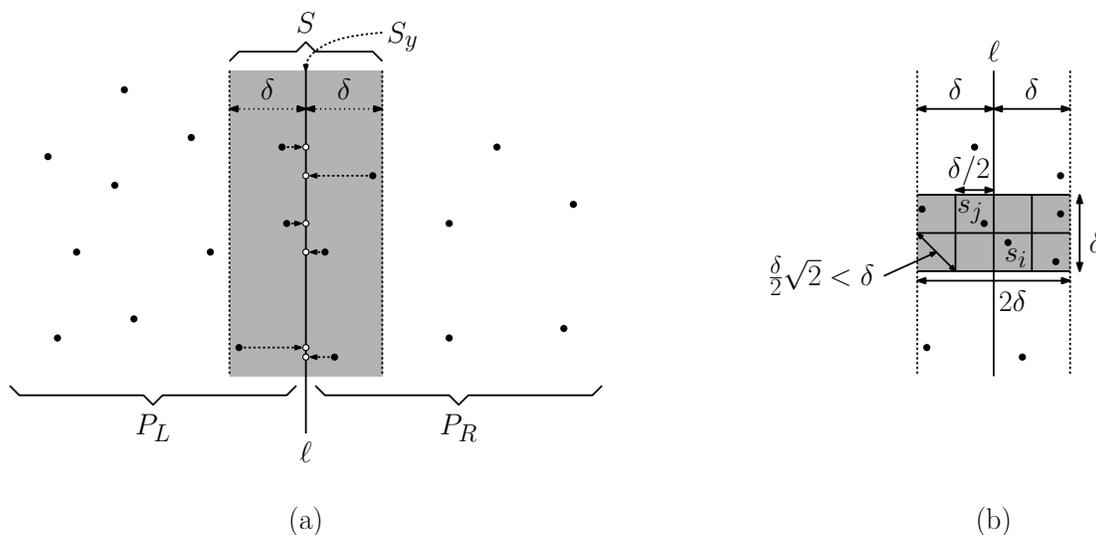


Fig. 108: Closest pair in the strip.

How do we find the closest pair within  $S$ ? Sorting comes to our rescue. Let  $S_y = \langle s_1, \dots, s_m \rangle$  denote the points of  $S$  sorted by their  $y$ -coordinates (see Fig. 108(a)). At the start of the lecture, we asserted

<sup>21</sup>You might be tempted to think that we have pruned away many of the points of  $P$ , and this is the source of efficiency, but this is not generally true. It might very well be that *every* point of  $P$  lies within the strip, and so we cannot afford to apply a brute-force solution to our problem.

that considering the points that are close according to their  $x$ - or  $y$ -coordinate alone is not sufficient. It is rather surprising, therefore, that this *does* work for the set  $S_y$ .

The key observation is that if  $S_y$  contains two points that are within distance  $\delta$  of each other, these two points *must* be within a constant number of positions of each other in the sorted array  $S_y$ . The following lemma formalizes this observation.

**Lemma:** Given any two points  $s_i, s_j \in S_y$ , if  $\|s_i s_j\| \leq \delta$ , then  $|j - i| \leq 7$ .

**Proof:** Suppose that  $\|s_i s_j\| \leq \delta$ . Since they are in  $S$  they are each within distance  $\delta$  of  $\ell$ . Clearly, the  $y$ -coordinates of these two points can differ by at most  $\delta$ . So they must both reside in a rectangle of width  $2\delta$  and height  $\delta$  centered about  $\ell$  (see Fig. 108(b)). Split this rectangle into eight identical squares each of side length  $\delta/2$ . A square of side length  $x$  has a diagonal of length  $x\sqrt{2}$ , and no two points within such a square can be farther away than this. Therefore, the distance between any two points lying within one of these eight squares is at most

$$\frac{\delta\sqrt{2}}{2} = \frac{\delta}{\sqrt{2}} < \delta.$$

Since each square lies entirely on one side of  $\ell$ , no square can contain two or more points of  $P$ , since otherwise, these two points would contradict the fact that  $\delta$  is the closest pair seen so far. Thus, there can be at most eight points of  $S$  in this rectangle, one for each square. Therefore,  $|j - i| \leq 7$ .

**Avoiding Repeated Sorting:** One issue that we have not yet addressed is how to compute  $S_y$ . Recall that we cannot afford to sort these points explicitly, because we may have  $n$  points in  $S$ , and this part of the algorithm needs to run in  $O(n)$  time<sup>22</sup> This is where presorting comes in. Recall that the points of  $P_y$  are already sorted by  $y$ -coordinates. To compute  $S_y$ , we enumerate the points of  $P_y$ , and each time we find a point that lies within the strip, we copy it to the next position of array  $S_y$ . This runs in  $O(n)$  time, and preserves the  $y$ -ordering of the points.

By the way, it is natural to wonder whether the value “8” in the statement of the lemma is optimal. Getting the best possible value is likely to be a tricky geometric exercise. Our textbook proves a weaker bound of “16”. Of course, from the perspective of asymptotic complexity, the exact constant does not matter.

The final algorithm is presented in the code fragment below.

## Lecture 32: Dynamic Programming: 0-1 Knapsack Problem

**0-1 Knapsack Problem:** Imagine that a burglar breaks into a museum and finds  $n$  items. Let  $v_i$  denote the value of the  $i$ -th item, and let  $w_i$  denote the weight of the  $i$ -th item. The burglar carries a knapsack capable of holding total weight  $W$ . The burglar wishes to carry away the most valuable subset items subject to the weight constraint.

For example, a burglar would rather steal diamonds before gold because the value per pound is better. But he would rather steal gold before lead for the same reason. We assume that the burglar cannot take a fraction of an object, so he/she must make a decision to take the object entirely or leave it behind. (There is a version of the problem where the burglar can take a fraction of an object for a fraction of the value and weight. This is much easier to solve.)

<sup>22</sup>If we were to pay the full sorting cost with each recursive call, the running time would be given by the recurrence  $T(n) = 2T(n/2) + n \log n$ . Solving this recurrence leads to the solution  $T(n) = O(n \log^2 n)$ , thus we would miss our target running time by an  $O(\log n)$  factor.

---

```
closestPair(P = (Px, Py)) {
  n = |P|
  if (n <= 3) solve by brute force          // basis case
  else {
    Find the vertical line L through P's median // divide
    Split P into PL and PR (split Px and Py as well)
    dL = closestPair(PL)                    // conquer
    dR = closestPair(PR)
    d = min(dL, dR)
    for (i = 1 to n) {                      // create Sy
      if (Py[i] is within distance d of L) {
        append Py[i] to Sy
      }
    }
    d' = stripClosest(Sy)                   // closest in strip
    return min(d, d')                       // overall closest
  }
}

stripClosest(Sy) {                          // closest in strip
  m = |Sy|
  d' = infinity
  for (i = 1 to m) {
    for (j = i+1 to min(m, i+7)) {         // search neighbors
      if (dist(Sy[i], Sy[j]) <= d') {
        d' = dist(Sy[i], Sy[j])           // new closest found
      }
    }
  }
  return d'
}
```

---

More formally, given  $\langle v_1, v_2, \dots, v_n \rangle$  and  $\langle w_1, w_2, \dots, w_n \rangle$ , and  $W > 0$ , we wish to determine the subset  $T \subseteq \{1, 2, \dots, n\}$  (of objects to “take”) that maximizes

$$\sum_{i \in T} v_i,$$

subject to

$$\sum_{i \in T} w_i \leq W.$$

Let us assume that the  $v_i$ 's,  $w_i$ 's and  $W$  are all positive integers. It turns out that this problem is NP-complete, and so we cannot really hope to find an efficient solution. However if we make the same sort of assumption that we made in counting sort, we can come up with an efficient solution.

We assume that the  $w_i$ 's are small integers, and that  $W$  itself is a small integer. We show that this problem can be solved in  $O(nW)$  time. (Note that this is not very good if  $W$  is a large integer. But if we truncate our numbers to lower precision, this gives a reasonable approximation algorithm.)

Here is how we solve the problem. We construct an array  $V[0..n, 0..W]$ . For  $1 \leq i \leq n$ , and  $0 \leq j \leq W$ , the entry  $V[i, j]$  we will store the maximum value of any subset of objects  $\{1, 2, \dots, i\}$  that can fit into a knapsack of weight  $j$ . If we can compute all the entries of this array, then the array entry  $V[n, W]$  will contain the maximum value of all  $n$  objects that can fit into the entire knapsack of weight  $W$ .

To compute the entries of the array  $V$  we will imply an inductive approach. As a basis, observe that  $V[0, j] = 0$  for  $0 \leq j \leq W$  since if we have no items then we have no value. We consider two cases:

**Leave object  $i$ :** If we choose to not take object  $i$ , then the optimal value will come about by considering how to fill a knapsack of size  $j$  with the remaining objects  $\{1, 2, \dots, i - 1\}$ . This is just  $V[i - 1, j]$ .

**Take object  $i$ :** If we take object  $i$ , then we gain a value of  $v_i$  but have used up  $w_i$  of our capacity. With the remaining  $j - w_i$  capacity in the knapsack, we can fill it in the best possible way with objects  $\{1, 2, \dots, i - 1\}$ . This is  $v_i + V[i - 1, j - w_i]$ . This is only possible if  $w_i \leq j$ .

Since these are the only two possibilities, we can see that we have the following rule for constructing the array  $V$ . The ranges on  $i$  and  $j$  are  $i \in [0..n]$  and  $j \in [0..W]$ .

$$\begin{aligned} V[0, j] &= 0 \\ V[i, j] &= \begin{cases} V[i - 1, j] & \text{if } w_i > j \\ \max(V[i - 1, j], v_i + V[i - 1, j - w_i]) & \text{if } w_i \leq j \end{cases} \end{aligned}$$

The first line states that if there are no objects, then there is no value, irrespective of  $j$ . The second line implements the rule above.

It is very easy to take these rules and produce an algorithm that computes the maximum value for the knapsack in time proportional to the size of the array, which is  $O((n + 1)(W + 1)) = O(nW)$ . The algorithm is given below.

An example is shown in the figure below. The final output is  $V[n, W] = V[4, 10] = 90$ . This reflects the selection of items 2 and 4, of values \$40 and \$50, respectively and weights  $4 + 3 \leq 10$ .

The only missing detail is what items should we select to achieve the maximum. We will leave this as an exercise. The key is to record for each entry  $V[i, j]$  in the matrix whether we got this entry by taking the  $i$ th item or leaving it. With this information, it is possible to reconstruct the optimum knapsack contents.

---

```

KnapSack(v[1..n], w[1..n], n, W) {
  allocate V[0..n][0..W];
  for j = 0 to W do V[0, j] = 0;           // initialization
  for i = 1 to n do {
    for j = 0 to W do {
      leave_val = V[i-1, j];             // value if we leave i
      if (j >= w[i])                     // enough capacity for i
        take_val = v[i] + V[i-1, j - w[i]]; // value if we take i
      else
        take_val = -INFINITY;           // cannot take i
      V[i,j] = max(leave_val, take_val); // final value is max
    }
  }
  return V[n, W];
}

```

---

Values of the objects are  $\langle 10, 40, 30, 50 \rangle$ .  
Weights of the objects are  $\langle 5, 4, 6, 3 \rangle$ .

Capacity $\rightarrow$			$j = 0$	1	2	3	4	5	6	7	8	9	10
Item	Value	Weight	0	0	0	0	0	0	0	0	0	0	0
1	10	5	0	0	0	0	0	10	10	10	10	10	10
2	40	4	0	0	0	0	40	40	40	40	40	50	50
3	30	6	0	0	0	0	40	40	40	40	40	50	70
4	50	3	0	0	0	50	50	50	50	90	90	90	90

Final result is  $V[4, 10] = 90$  (for taking items 2 and 4).

Fig. 109: 0-1 Knapsack Example.

## Lecture 33: Dynamic Programming: Minimum Weight Triangulation

**Polygons and Triangulations:** Let's consider a geometric problem that outwardly appears to be quite different from chain-matrix multiplication, but actually has remarkable similarities. We begin with a number of definitions. Define a *polygon* to be a piecewise linear closed curve in the plane. In other words, we form a cycle by joining line segments end to end. The line segments are called the *sides* of the polygon and the endpoints are called the *vertices*. A polygon is *simple* if it does not cross itself, that is, if the sides do not intersect one another except for two consecutive sides sharing a common vertex. A simple polygon subdivides the plane into its *interior*, its *boundary* and its *exterior*. A simple polygon is said to be *convex* if every interior angle is at most 180 degrees. Vertices with interior angle equal to 180 degrees are normally allowed, but for this problem we will assume that no such vertices exist.

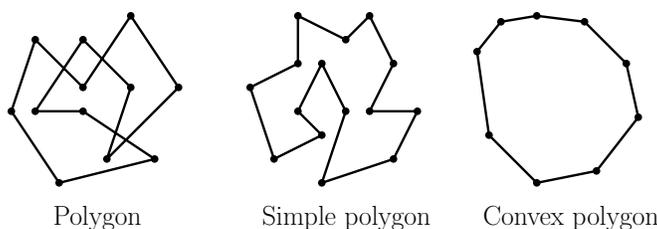


Fig. 110: Polygons.

Given a convex polygon, we assume that its vertices are labeled in counterclockwise order  $P = \langle v_1, \dots, v_n \rangle$ . We will assume that indexing of vertices is done modulo  $n$ , so  $v_0 = v_n$ . This polygon has  $n$  sides,  $\overline{v_{i-1}v_i}$ .

Given two nonadjacent sides  $v_i$  and  $v_j$ , where  $i < j - 1$ , the line segment  $\overline{v_i v_j}$  is a *chord*. (If the polygon is simple but not convex, we include the additional requirement that the interior of the segment must lie entirely in the interior of  $P$ .) Any chord subdivides the polygon into two polygons:  $\langle v_i, v_{i+1}, \dots, v_j \rangle$ , and  $\langle v_j, v_{j+1}, \dots, v_i \rangle$ . A *triangulation* of a convex polygon  $P$  is a subdivision of the interior of  $P$  into a collection of triangles with disjoint interiors, whose vertices are drawn from the vertices of  $P$ . Equivalently, we can define a triangulation as a maximal set  $T$  of nonintersecting chords. (In other words, every chord that is not in  $T$  intersects the interior of some chord in  $T$ .) It is easy to see that such a set of chords subdivides the interior of the polygon into a collection of triangles with pairwise disjoint interiors (and hence the name *triangulation*). It is not hard to prove (by induction) that every triangulation of an  $n$ -sided polygon consists of  $n - 3$  chords and  $n - 2$  triangles. Triangulations are of interest for a number of reasons. Many geometric algorithms operate by first decomposing a complex polygonal shape into triangles.

In general, given a convex polygon, there are many possible triangulations. In fact, the number is exponential in  $n$ , the number of sides. Which triangulation is the “best”? There are many criteria that are used depending on the application. One criterion is to imagine that you must “pay” for the ink you use in drawing the triangulation, and you want to minimize the amount of ink you use. (This may sound fanciful, but minimizing wire length is an important condition in chip design. Further, this is one of many properties which we could choose to optimize.) This suggests the following optimization problem:

**Minimum-weight convex polygon triangulation:** Given a convex polygon determine the triangulation that minimizes the sum of the perimeters of its triangles. (See Fig. 111.)

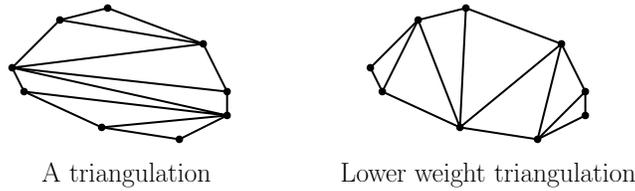


Fig. 111: Triangulations of convex polygons, and the minimum weight triangulation.

Given three distinct vertices  $v_i, v_j, v_k$ , we define the *weight* of the associated triangle by the weight function

$$w(v_i, v_j, v_k) = |v_i v_j| + |v_j v_k| + |v_k v_i|,$$

where  $|v_i v_j|$  denotes the length of the line segment  $\overline{v_i v_j}$ .

**Dynamic Programming Solution:** Let us consider an  $(n + 1)$ -sided polygon  $P = \langle v_0, v_1, \dots, v_n \rangle$ . Let us assume that these vertices have been numbered in counterclockwise order. To derive a DP formulation we need to define a set of subproblems from which we can derive the optimum solution. For  $0 \leq i < j \leq n$ , define  $t[i, j]$  to be the weight of the minimum weight triangulation for the subpolygon that lies to the right of directed chord  $\overline{v_i v_j}$ , that is, the polygon with the counterclockwise vertex sequence  $\langle v_i, v_{i+1}, \dots, v_j \rangle$ . Observe that if we can compute this quantity for all such  $i$  and  $j$ , then the weight of the minimum weight triangulation of the entire polygon can be extracted as  $t[0, n]$ . (As usual, we only compute the minimum weight. But, it is easy to modify the procedure to extract the actual triangulation.)

As a basis case, we define the weight of the trivial “2-sided polygon” to be zero, implying that  $t[i, i + 1] = 0$ . In general, to compute  $t[i, j]$ , consider the subpolygon  $\langle v_i, v_{i+1}, \dots, v_j \rangle$ , where  $j > i + 1$ . One of the chords of this polygon is the side  $\overline{v_i v_j}$ . We may split this subpolygon by introducing a triangle whose base is this chord, and whose third vertex is any vertex  $v_k$ , where  $i < k < j$ . This subdivides the polygon into the subpolygons  $\langle v_i, v_{i+1}, \dots, v_k \rangle$  and  $\langle v_k, v_{k+1}, \dots, v_j \rangle$  whose minimum weights are already known to us as  $t[i, k]$  and  $t[k, j]$ . In addition we should consider the weight of the newly added triangle  $\triangle v_i v_k v_j$ . Thus, we have the following recursive rule:

$$t[i, j] = \begin{cases} 0 & \text{if } j = i + 1 \\ \min_{i < k < j} (t[i, k] + t[k, j] + w(v_i v_k v_j)) & \text{if } j > i + 1. \end{cases}$$

The final output is the overall minimum weight, which is,  $t[0, n]$ . This is illustrated in Fig. 112

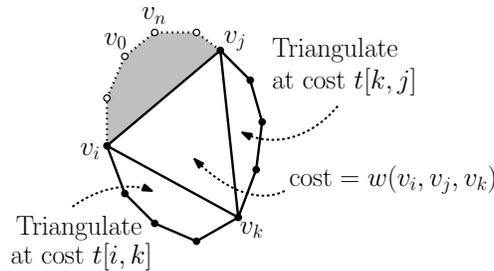


Fig. 112: Triangulations and tree structure.

Note that this has almost exactly the same structure as the recursive definition used in the chain matrix multiplication algorithm (except that some indices are different by 1.) The same  $\Theta(n^3)$  algorithm can be applied with only minor changes.

**Relationship to Binary Trees:** One explanation behind the similarity of triangulations and the chain matrix multiplication algorithm is to observe that both are fundamentally related to binary trees. In the case of the chain matrix multiplication, the associated binary tree is the evaluation tree for the multiplication, where the leaves of the tree correspond to the matrices, and each node of the tree is associated with a product of a sequence of two or more matrices. To see that there is a similar correspondence here, consider an  $(n + 1)$ -sided convex polygon  $P = \langle v_0, v_1, \dots, v_n \rangle$ , and fix one side of the polygon (say  $\overline{v_0 v_n}$ ). Now consider a rooted binary tree whose root node is the triangle containing side  $\overline{v_0 v_n}$ , whose internal nodes are the nodes of the dual tree, and whose leaves correspond to the remaining sides of the tree. Observe that partitioning the polygon into triangles is equivalent to a binary tree with  $n$  leaves, and vice versa. This is illustrated in Fig. 113. Note that every triangle is associated with an internal node of the tree and every edge of the original polygon, except for the distinguished starting side  $\overline{v_0 v_n}$ , is associated with a leaf node of the tree.

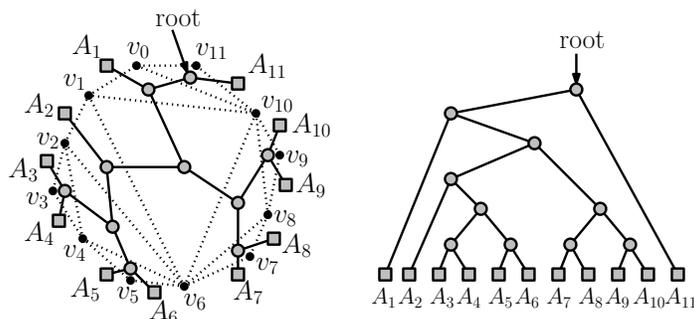


Fig. 113: Triangulations and tree structure.

Once you see this connection. Then the following two observations follow easily. Observe that the associated binary tree has  $n$  leaves, and hence (by standard results on binary trees)  $n - 1$  internal nodes. Since each internal node other than the root has one edge entering it, there are  $n - 2$  edges between the internal nodes. Each internal node corresponds to one triangle, and each edge between internal nodes corresponds to one chord of the triangulation.

## Lecture 34: Hamiltonian Path

**Hamiltonian Cycle:** Today we consider a collection of problems related to finding paths in graphs and digraphs. Recall that given a graph (or digraph) a *Hamiltonian cycle* is a simple cycle that visits every vertex in the graph (exactly once). A *Hamiltonian path* is a simple path that visits every vertex in the graph (exactly once). The Hamiltonian cycle (HC) and Hamiltonian path (HP) problems ask whether a given graph (or digraph) has such a cycle or path, respectively. There are four variations of these problems depending on whether the graph is directed or undirected, and depending on whether you want a path or a cycle, but all of these problems are NP-complete.

An important related problem is the traveling salesman problem (TSP). Given a complete graph (or digraph) with integer edge weights, determine the cycle of minimum weight that visits all the vertices. Since the graph is complete, such a cycle will always exist. The decision problem formulation is, given a complete weighted graph  $G$ , and integer  $X$ , does there exist a Hamiltonian cycle of total weight at most  $X$ ? Today we will prove that Hamiltonian Cycle is NP-complete. We will leave TSP as an easy exercise. (It is done in Section 36.5.5 in CLRS.)

**Component Design:** Up to now, most of the reductions that we have seen (for Clique, VC, and DS in particular) are of a relatively simple variety. They are sometimes called *local replacement* reductions, because they operate by making some local change throughout the graph.

We will present a much more complex style of reduction for the Hamiltonian path problem on directed graphs. This type of reduction is called a *component design* reduction, because it involves designing special subgraphs, sometimes called *components* or *gadgets* (also called *widgits*) whose job it is to enforce a particular constraint. Very complex reductions may involve the creation of many gadgets. This one involves the construction of only one. (See CLRS's or KT's presentation of HP for other examples of gadgets.)

The gadget that we will use in the directed Hamiltonian path reduction, called a *DHP-gadget*, is shown in the figure below. It consists of three incoming edges labeled  $i_1, i_2, i_3$  and three outgoing edges, labeled  $o_1, o_2, o_3$ . It was designed so it satisfied the following property, which you can verify. Intuitively it says that if you enter the gadget on any subset of 1, 2 or 3 input edges, then there is a way to get through the gadget and hit every vertex exactly once, and in doing so each path must end on the corresponding output edge.

**Claim:** Given the DHP-gadget:

- For any subset of input edges, there exists a set of paths which join each input edge  $i_1, i_2$ , or  $i_3$  to its respective output edge  $o_1, o_2$ , or  $o_3$  such that together these paths visit every vertex in the gadget exactly once.
- Any subset of paths that start on the input edges and end on the output edges, and visit all the vertices of the gadget exactly once, must join corresponding inputs to corresponding outputs. (In other words, a path that starts on input  $i_1$  must exit on output  $o_1$ .)

The proof is not hard, but involves a careful inspection of the gadget. It is probably easiest to see this on your own, by starting with one, two, or three input paths, and attempting to get through the gadget without skipping vertex and without visiting any vertex twice. To see whether you really understand the gadget, answer the question of why there are 6 groups of triples. Would some other number work?

**DHP is NP-complete:** This gadget is an essential part of our proof that the directed Hamiltonian path problem is NP-complete.

**Theorem:** The directed Hamiltonian Path problem is NP-complete.

**Proof: DHP  $\in$  NP:** The certificate consists of the sequence of vertices (or edges) in the path. It is an easy matter to check that the path visits every vertex exactly once.

**3SAT  $\leq_P$  DHP:** This will be the subject of the rest of this section.

Let us consider the similar elements between the two problems. In 3SAT we are selecting a truth assignment for the variables of the formula. In DHP, we are deciding which edges will be a part of the path. In 3SAT there must be at least one true literal for each clause. In DHP, each vertex must be visited exactly once.

We are given a boolean formula  $F$  in 3-CNF form (three literals per clause). We will convert this formula into a digraph. Let  $x_1, x_2, \dots, x_m$  denote the variables appearing in  $F$ . We will construct one DHP-gadget for each clause in the formula. The inputs and outputs of each gadget correspond to the literals appearing in this clause. Thus, the clause  $(\bar{x}_2 \vee x_5 \vee \bar{x}_8)$  would generate a clause gadget with inputs labeled  $\bar{x}_2, x_5$ , and  $\bar{x}_8$ , and the same outputs.

The general structure of the digraph will consist of a series vertices, one for each variable. Each of these vertices will have two outgoing paths, one taken if  $x_i$  is set to true and one if  $x_i$  is set to false. Each of these paths will then pass through some number of DHP-gadgets. The true path for  $x_i$  will pass through all the clause gadgets for clauses in which  $x_i$  appears, and the false path will pass through all the gadgets for clauses in which  $\bar{x}_i$  appears. (The order in which the path passes through the gadgets is unimportant.) When the paths for  $x_i$  have passed through their last gadgets, then they are joined to the next variable vertex,  $x_{i+1}$ . This is illustrated in the following figure. (The figure only shows a

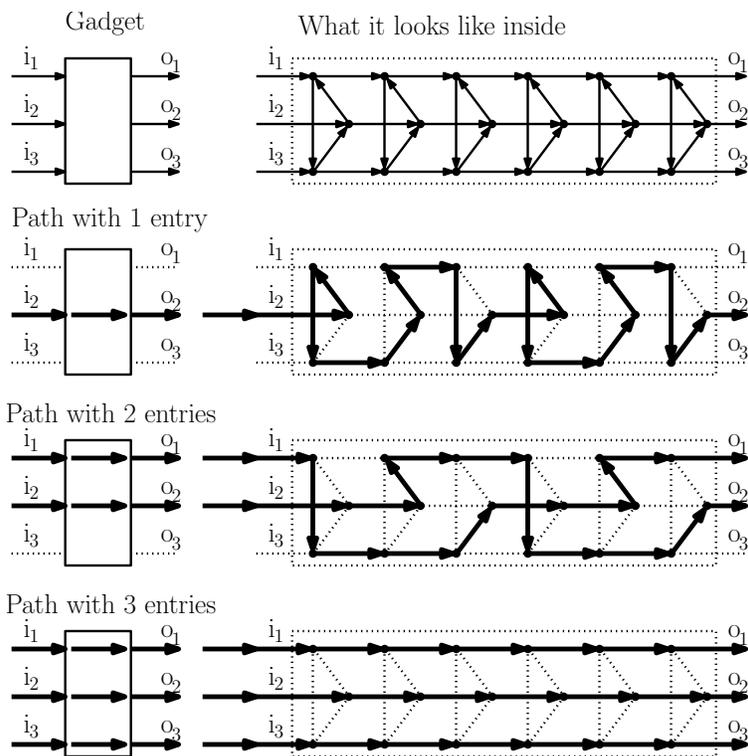


Fig. 114: DHP-Gadget and examples of path traversals.

portion of the construction. There will be paths coming into these same gadgets from other variables as well.) We add one final vertex  $x_e$ , and the last variable's paths are connected to  $x_e$ . (If we wanted to reduce to Hamiltonian cycle, rather than Hamiltonian path, we could join  $x_e$  back to  $x_1$ .)

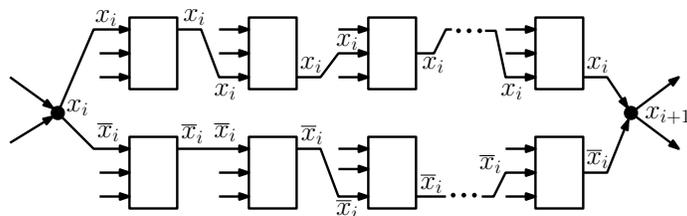


Fig. 115: General structure of reduction from 3SAT to DHP.

Note that for each variable, the Hamiltonian path must either use the true path or the false path, but it cannot use both. If we choose the true path for  $x_i$  to be in the Hamiltonian path, then we will have at least one path passing through each of the gadgets whose corresponding clause contains  $x_i$ , and if we chose the false path, then we will have at least one path passing through each gadget for  $\bar{x}_i$ .

For example, consider the following boolean formula in 3-CNF. The construction yields the digraph shown in the following figure.

$$(\bar{x}_1 \vee x_2 \vee x_3) \wedge (x_1 \vee \bar{x}_2 \vee \bar{x}_3) \wedge (x_2 \vee \bar{x}_1 \vee \bar{x}_3) \wedge (x_1 \vee x_3 \vee \bar{x}_2).$$

**The Reduction:** Let us give a more formal description of the reduction. Recall that we are given a boolean formula  $F$  in 3-CNF. We create a digraph  $G$  as follows. For each variable  $x_i$  appearing in  $F$ , we create

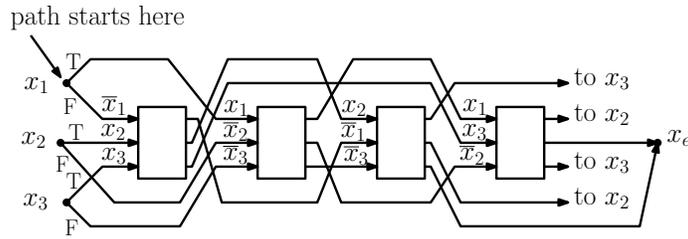


Fig. 116: Example of the 3SAT to DHP reduction.

a *variable vertex*, named  $x_i$ . We also create a vertex named  $x_e$  (the ending vertex). For each clause  $c$ , we create a DHP-gadget whose inputs and outputs are labeled with the three literals of  $c$ . (The order is unimportant, as long as each input and its corresponding output are labeled the same.)

We join these vertices with the gadgets as follows. For each variable  $x_i$ , consider all the clauses  $c_1, c_2, \dots, c_k$  in which  $x_i$  appears as a literal (uncomplemented). Join  $x_i$  by an edge to the input labeled with  $x_i$  in the gadget for  $c_1$ , and in general join the output of gadget  $c_j$  labeled  $x_i$  with the input of gadget  $c_{j+1}$  with this same label. Finally, join the output of the last gadget  $c_k$  to the next vertex variable  $x_{i+1}$ . (If this is the last variable, then join it to  $x_e$  instead.) The resulting chain of edges is called the *true path* for variable  $x_i$ . Form a second chain in exactly the same way, but this time joining the gadgets for the clauses in which  $\bar{x}_i$  appears. This is called the *false path* for  $x_i$ . The resulting digraph is the output of the reduction. Observe that the entire construction can be performed in polynomial time, by simply inspecting the formula, creating the appropriate vertices, and adding the appropriate edges to the digraph. The following lemma establishes the correctness of this reduction.

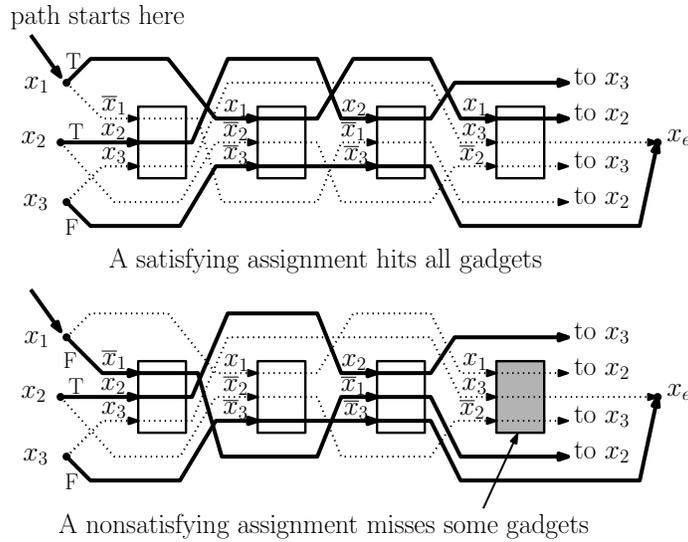


Fig. 117: Correctness of the 3SAT to DHP reduction. The upper figure shows the Hamiltonian path resulting from the satisfying assignment,  $x_1 = 1, x_2 = 1, x_3 = 0$ , and the lower figure shows the non-Hamiltonian path resulting from the non-satisfying assignment  $x_1 = 0, x_2 = 1, x_3 = 0$ .

**Lemma:** The boolean formula  $F$  is satisfiable if and only if the digraph  $G$  produced by the above reduction has a Hamiltonian path.

**Proof:** We need to prove both the “only if” and the “if”.

$\Rightarrow$ : Suppose that  $F$  has a satisfying assignment. We claim that  $G$  has a Hamiltonian path. This path will start at the variable vertex  $x_1$ , then will travel along either the true path or false path for  $x_1$ , depending on whether it is 1 or 0, respectively, in the assignment, and then it will continue with  $x_2$ , then  $x_3$ , and so on, until reaching  $x_e$ . Such a path will visit each variable vertex exactly once. Because this is a satisfying assignment, we know that for each clause, either 1, 2, or 3 of its literals will be true. This means that for each clause, either 1, 2, or 3, paths will attempt to travel through the corresponding gadget. However, we have argued in the above claim that in this case it is possible to visit every vertex in the gadget exactly once. Thus every vertex in the graph is visited exactly once, implying that  $G$  has a Hamiltonian path.

$\Leftarrow$ : Suppose that  $G$  has a Hamiltonian path. We assert that the form of the path must be essentially the same as the one described in the previous part of this proof. In particular, the path must visit the variable vertices in increasing order from  $x_1$  until  $x_e$ , because of the way in which these vertices are joined together.

Also observe that for each variable vertex, the path will proceed along either the true path or the false path. If it proceeds along the true path, set the corresponding variable to 1 and otherwise set it to 0. We will show that the resulting assignment is a satisfying assignment for  $F$ .

Any Hamiltonian path must visit all the vertices in every gadget. By the above claim about DHP-gadgets, if a path visits all the vertices and enters along input edge then it must exit along the corresponding output edge. Therefore, once the Hamiltonian path starts along the true or false path for some variable, it must remain on edges with the same label. That is, if the path starts along the true path for  $x_i$ , it must travel through all the gadgets with the label  $x_i$  until arriving at the variable vertex for  $x_{i+1}$ . If it starts along the false path, then it must travel through all gadgets with the label  $\bar{x}_i$ .

Since all the gadgets are visited and the paths must remain true to their initial assignments, it follows that for each corresponding clause, at least one (and possibly 2 or three) of the literals must be true. Therefore, this is a satisfying assignment.

## Lecture 35: Subset Sum

**Subset Sum:** The Subset Sum problem (SS) is the following. Given a finite set  $S$  of positive integers  $S = \{w_1, w_2, \dots, w_n\}$  and a *target value*,  $t$ , we want to know whether there exists a subset  $S' \subseteq S$  that sums exactly to  $t$ .

This problem is a simplified version of the 0-1 Knapsack problem, presented as a decision problem. Recall that in the 0-1 Knapsack problem, we are given a collection of objects, each with an associated weight  $w_i$  and associated value  $v_i$ . We are given a knapsack of capacity  $W$ . The objective is to take as many objects as can fit in the knapsack's capacity so as to maximize the value. (In the fractional knapsack we could take a portion of an object. In the 0-1 Knapsack we either take an object entirely or leave it.) In the simplest version, suppose that the value is the same as the weight,  $v_i = w_i$ . (This would occur for example if all the objects were made of the same material, say, gold.) Then, the best we could hope to achieve would be to fill the knapsack entirely. By setting  $t = W$ , we see that the subset sum problem is equivalent to this simplified version of the 0-1 Knapsack problem. It follows that if we can show that this simpler version is NP-complete, then certainly the more general 0-1 Knapsack problem (stated as a decision problem) is also NP-complete.

Consider the following example.

$$S = \{3, 6, 9, 12, 15, 23, 32\} \quad \text{and} \quad t = 33.$$

The subset  $S' = \{6, 12, 15\}$  sums to  $t = 33$ , so the answer in this case is yes. If  $t = 34$  the answer would be no.

**Dynamic Programming Solution:** There is a dynamic programming algorithm which solves the Subset Sum problem in  $O(n \cdot t)$  time.<sup>23</sup>

The quantity  $n \cdot t$  is a polynomial function of  $n$ . This would seem to imply that the Subset Sum problem is in P. But there is an important catch. Recall that in all NP-complete problems we assume (1) running time is measured as a function of input size (number of bits) and (2) inputs must be encoded in a reasonable succinct manner. Let us assume that the numbers  $w_i$  and  $t$  are all  $b$ -bit numbers represented in base 2, using the fewest number of bits possible. Then the input size is  $O(nb)$ . The value of  $t$  may be as large as  $2^b$ . So the resulting algorithm has a running time of  $O(n2^b)$ . This is polynomial in  $n$ , but exponential in  $b$ . Thus, this running time is not polynomial as a function of the input size.

Note that an important consequence of this observation is that the SS problem is not hard when the numbers involved are small. If the numbers involved are of a fixed number of bits (a constant independent of  $n$ ), then the problem is solvable in polynomial time. However, we will show that in the general case, this problem is NP-complete.

**SS is NP-complete:** The proof that Subset Sum (SS) is NP-complete involves the usual two elements.

- (i)  $SS \in NP$ .
- (ii) Some known NP-complete problem is reducible to SS. In particular, we will show that Vertex Cover (VC) is reducible to SS, that is,  $VC \leq_P SS$ .

To show that SS is in NP, we need to give a verification procedure. Given  $S$  and  $t$ , the certificate is just the indices of the numbers that form the subset  $S'$ . We can add two  $b$ -bit numbers together in  $O(b)$  time. So, in polynomial time we can compute the sum of elements in  $S'$ , and verify that this sum equals  $t$ .

For the remainder of the proof we show how to reduce vertex cover to subset sum. We want a polynomial time computable function  $f$  that maps an instance of the vertex cover (a graph  $G$  and integer  $k$ ) to an instance of the subset sum problem (a set of integers  $S$  and target integer  $t$ ) such that  $G$  has a vertex cover of size  $k$  if and only if  $S$  has a subset summing to  $t$ . Thus, if subset sum were solvable in polynomial time, so would vertex cover.

How can we encode the notion of selecting a subset of vertices that cover all the edges to that of selecting a subset of numbers that sums to  $t$ ? In the vertex cover problem we are selecting vertices, and in the subset sum problem we are selecting numbers, so it seems logical that the reduction should map vertices into numbers. The constraint that these vertices should cover all the edges must be mapped to the constraint that the sum of the numbers should equal the target value.

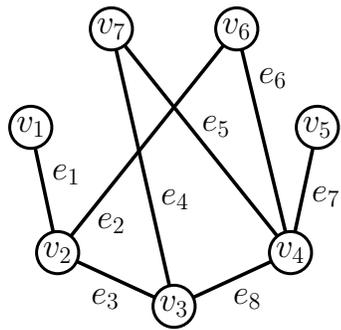
**An Initial Approach:** Here is an idea, which does not work, but gives a sense of how to proceed. Let  $E$  denote the number of edges in the graph. First number the edges of the graph from 1 through  $E$ . Then represent each vertex  $v_i$  as an  $E$ -element bit vector, where the  $j$ -th bit from the left is set to 1 if and only if the edge  $e_j$  is incident to vertex  $v_i$ . (Another way to think of this is that these bit vectors form the rows of an *incidence matrix* for the graph.) An example is shown below, in which  $k = 3$ .

Now, suppose we take any subset of vertices and form the logical-or of the corresponding bit vectors. If the subset is a vertex cover, then every edge will be covered by at least one of these vertices, and so the logical-or will be a bit vector of all 1's,  $1111 \dots 1$ . Conversely, if the logical-or is a bit vector of 1's, then each edge has been covered by some vertex, implying that the vertices form a vertex cover. (Later we will consider how to encode the fact that there only allowed  $k$  vertices in the cover.)

Since bit vectors can be thought of as just a way of representing numbers in binary, this is starting to feel more like the subset sum problem. The target would be the number whose bit vector is all 1's.

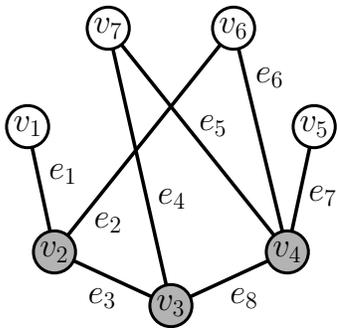
---

<sup>23</sup>We will leave this as an exercise, but the formulation is, for  $0 \leq i \leq n$  and  $0 \leq t' \leq t$ ,  $S[i, t'] = 1$  if there is a subset of  $\{w_1, w_2, \dots, w_i\}$  that sums to  $t'$ , and 0 otherwise. The  $i$ th row of this table can be computed in  $O(t)$  time, given the contents of the  $(i - 1)$ -st row.



	$e_1$	$e_2$	$e_3$	$e_4$	$e_5$	$e_6$	$e_7$	$e_8$
$v_1$	1	0	0	0	0	0	0	0
$v_2$	1	1	1	0	0	0	0	0
$v_3$	0	0	1	1	0	0	0	1
$v_4$	0	0	0	0	1	1	1	1
$v_5$	0	0	0	0	0	0	1	0
$v_6$	0	1	0	0	0	1	0	0
$v_7$	0	0	0	1	1	0	0	0

Fig. 118: Encoding a graph as a collection of bit vectors.



	$e_1$	$e_2$	$e_3$	$e_4$	$e_5$	$e_6$	$e_7$	$e_8$
$v_1$	1	0	0	0	0	0	0	0
$v_2$	1	1	1	0	0	0	0	0
$v_3$	0	0	1	1	0	0	0	1
$v_4$	0	0	0	0	1	1	1	1
$v_5$	0	0	0	0	0	0	1	0
$v_6$	0	1	0	0	0	1	0	0
$v_7$	0	0	0	1	1	0	0	0

$$t = 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 = v_2 \vee v_3 \vee v_4$$

Fig. 119: The logical-or of a vertex cover equals 1111...1.

There are a number of problems, however. First, logical-or is not the same as addition. For example, if both of the endpoints of some edge are in the vertex cover, then its value in the corresponding column would be 2, not 1. Second, we have no way of controlling how many vertices go into the vertex cover. (We could just take the logical-or of all the vertices, and then the logical-or would certainly be a bit vectors of 1's.)

There are two ways in which addition differs significantly from logical-or. The first is the issue of carries. For example, the  $1101 \vee 0011 = 1111$ , but in binary  $1101 + 0011 = 1000$ . To fix this, we recognize that we do not have to use a binary (base-2) representation. In fact, we can assume any base system we want. Observe that each column of the incidence matrix has at most two 1's in any column, because each edge is incident to at most two vertices. Thus, if use any base that is at least as large as base 3, we will never generate a carry to the next position. In fact we will use base 4 (for reasons to be seen below). Note that the base of the number system is just for own convenience of notation. Once the numbers have been formed, they will be converted into whatever form our machine assumes for its input representation, e.g. decimal or binary.

The second difference between logical-or and addition is that an edge may generally be covered either once or twice in the vertex cover. So, the final sum of these numbers will be a number consisting of 1 and 2 digits, e.g.  $1211 \dots 112$ . This does not provide us with a unique target value  $t$ . We know that no digit of our sum can be a zero. To fix this problem, we will create a set of  $E$  additional *slack values*. For  $1 \leq i \leq E$ , the  $i$ th slack value will consist of all 0's, except for a single 1-digit in the  $i$ th position, e.g.,  $00000100000$ . Our target will be the number  $2222 \dots 222$  (all 2's). To see why this works, observe that from the numbers of our vertex cover, we will get a sum consisting of 1's and 2's. For each position where there is a 1, we can supplement this value by adding in the corresponding slack value. Thus we can boost any value consisting of 1's and 2's to all 2's. On the other hand, note that if there are any 0 values in the final sum, we will not have enough slack values to convert this into a 2.

There is one last issue. We are only allowed to place only  $k$  vertices in the vertex cover. We will handle this by adding an additional column. For each number arising from a vertex, we will put a 1 in this additional column. For each slack variable we will put a 0. In the target, we will require that this column sum to the value  $k$ , the size of the vertex cover. Thus, to form the desired sum, we must select exactly  $k$  of the vertex values. Note that since we only have a base-4 representation, there might be carries out of this last column (if  $k \geq 4$ ). But since this is the last column, it will not affect any of the other aspects of the construction.

**The Final Reduction:** Here is the final reduction, given the graph  $G = (V, E)$  and integer  $k$  for the vertex cover problem.

- (1) Create a set of  $n$  vertex values,  $x_1, x_2, \dots, x_n$  using base-4 notation. The value  $x_i$  is equal a 1 followed by a sequence of  $E$  base-4 digits. The  $j$ -th digit is a 1 if edge  $e_j$  is incident to vertex  $v_i$  and 0 otherwise.
- (2) Create  $E$  slack values  $y_1, y_2, \dots, y_E$ , where  $y_i$  is a 0 followed by  $E$  base-4 digits. The  $i$ -th digit of  $y_i$  is 1 and all others are 0.
- (3) Let  $t$  be the base-4 number whose first digit is  $k$  (this may actually span multiple base-4 digits), and whose remaining  $E$  digits are all 2.
- (4) Convert the  $x_i$ 's, the  $y_j$ 's, and  $t$  into whatever base notation is used for the subset sum problem (e.g. base 10). Output the set  $S = \{x_1, \dots, x_n, y_1, \dots, y_E\}$  and  $t$ .

Observe that this can be done in polynomial time, in  $O(E^2)$ , in fact. The construction is illustrated in Fig. 120.

**Correctness:** We claim that  $G$  has a vertex cover of size  $k$  if and only if  $S$  has a subset that sums to  $t$ . If  $G$  has a vertex cover  $V'$  of size  $k$ , then we take the vertex values  $x_i$  corresponding to the vertices of  $V'$ ,

		e <sub>1</sub>	e <sub>2</sub>	e <sub>3</sub>	e <sub>4</sub>	e <sub>5</sub>	e <sub>6</sub>	e <sub>7</sub>	e <sub>8</sub>	
x <sub>1</sub>	1	1	0	0	0	0	0	0	0	}
x <sub>2</sub>	1	1	1	1	0	0	0	0	0	
x <sub>3</sub>	1	0	0	1	1	0	0	0	1	
x <sub>4</sub>	1	0	0	0	0	1	1	1	1	
x <sub>5</sub>	1	0	0	0	0	0	0	1	0	
x <sub>6</sub>	1	0	1	0	0	0	1	0	0	
x <sub>7</sub>	1	0	0	0	1	1	0	0	0	
y <sub>1</sub>	0	1	0	0	0	0	0	0	0	}
y <sub>2</sub>	0	0	1	0	0	0	0	0	0	
y <sub>3</sub>	0	0	0	1	0	0	0	0	0	
y <sub>4</sub>	0	0	0	0	1	0	0	0	0	
y <sub>5</sub>	0	0	0	0	0	1	0	0	0	
y <sub>6</sub>	0	0	0	0	0	0	1	0	0	
y <sub>7</sub>	0	0	0	0	0	0	0	1	0	
y <sub>8</sub>	0	0	0	0	0	0	0	0	1	
t	3	2	2	2	2	2	2	2	2	

↑ vertex cover size (k=3)

Fig. 120: Vertex cover to subset sum reduction.

and for each edge that is covered only once in  $V'$ , we take the corresponding slack variable. It follows from the comments made earlier that the lower-order  $E$  digits of the resulting sum will be of the form  $222 \dots 2$  and because there are  $k$  elements in  $V'$ , the leftmost digit of the sum will be  $k$ . Thus, the resulting subset sums to  $t$ .

Conversely, if  $S$  has a subset  $S'$  that sums to  $t$  then we assert that it must select exactly  $k$  values from among the vertex values, since the first digit must sum to  $k$ . We claim that these vertices  $V'$  form a vertex cover. In particular, no edge can be left uncovered by  $V'$ , since (because there are no carries) the corresponding column would be 0 in the sum of vertex values. Thus, no matter what slack values we add, the resulting digit position could not be equal to 2, and so this cannot be a solution to the subset sum problem.

It is worth noting again that in this reduction, we needed to have large numbers. For example, the target value  $t$  is at least as large as  $4^E \geq 4^n$  (where  $n$  is the number of vertices in  $G$ ). In our dynamic programming solution  $W = t$ , so the DP algorithm would run in  $\Omega(n4^n)$  time, which is not polynomial time.

## Lecture 36: Subset-Sum Approximation

**Polynomial Approximation Schemes:** Last time we saw that for some NP-complete problems, it is possible to approximate the problem to within a fixed constant ratio bound. For example, the approximation algorithm produces an answer that is within a factor of 2 of the optimal solution. However, in practice, people would like to control the precision of the approximation. This is done by specifying a parameter  $\epsilon > 0$  as part of the input to the approximation algorithm, and requiring that the algorithm produce an answer that is within a *relative error* of  $\epsilon$  of the optimal solution. It is understood that as  $\epsilon$  tends to 0, the running time of the algorithm will increase. Such an algorithm is called a *polynomial approximation scheme*.

	e <sub>1</sub>	e <sub>2</sub>	e <sub>3</sub>	e <sub>4</sub>	e <sub>5</sub>	e <sub>6</sub>	e <sub>7</sub>	e <sub>8</sub>	
X <sub>1</sub>	1	1	0	0	0	0	0	0	} Vertex values (take those in vertex cover)
X <sub>2</sub>	1	1	1	1	0	0	0	0	
X <sub>3</sub>	1	0	0	1	1	0	0	1	
X <sub>4</sub>	1	0	0	0	0	1	1	1	
X <sub>5</sub>	1	0	0	0	0	0	0	1	
X <sub>6</sub>	1	0	1	0	0	0	1	0	
X <sub>7</sub>	1	0	0	0	1	1	0	0	
X <sub>8</sub>	1	0	0	0	0	0	0	0	
Y <sub>1</sub>	0	1	0	0	0	0	0	0	} Slack values (take one for each edge that has only one endpoint in the cover)
Y <sub>2</sub>	0	0	1	0	0	0	0	0	
Y <sub>3</sub>	0	0	0	1	0	0	0	0	
Y <sub>4</sub>	0	0	0	0	1	0	0	0	
Y <sub>5</sub>	0	0	0	0	0	1	0	0	
Y <sub>6</sub>	0	0	0	0	0	0	1	0	
Y <sub>7</sub>	0	0	0	0	0	0	0	1	
Y <sub>8</sub>	0	0	0	0	0	0	0	1	
t	3	2	2	2	2	2	2	2	} vertex cover size

Fig. 121: Correctness of the reduction.

For example, the running time of the algorithm might be  $O(2^{(1/\epsilon)n^2})$ . It is easy to see that in such cases the user pays a big penalty in running time as a function of  $\epsilon$ . (For example, to produce a 1% error, the “constant” factor would be  $2^{100}$  which would be around 4 quadrillion centuries on your 100 Mhz Pentium.) A *fully polynomial approximation scheme* is one in which the running time is polynomial in both  $n$  and  $1/\epsilon$ . For example, a running time of  $O((n/\epsilon)^2)$  would satisfy this condition. In such cases, reasonably accurate approximations are computationally feasible.

Unfortunately, there are very few NP-complete problems with fully polynomial approximation schemes. In fact, recently there has been strong evidence that many NP-complete problems do not have polynomial approximation schemes (fully or otherwise). Today we will study one that does.

**Subset Sum:** Recall that in the subset sum problem we are given a set  $S$  of positive integers  $\{x_1, x_2, \dots, x_n\}$  and a target value  $t$ , and we are asked whether there exists a subset  $S' \subseteq S$  that sums exactly to  $t$ . The optimization problem is to determine the subset whose sum is as large as possible but not larger than  $t$ .

This problem is basic to many packing problems, and is indirectly related to processor scheduling problems that arise in operating systems as well. Suppose we are also given  $0 < \epsilon < 1$ . Let  $z^* \leq t$  denote the optimum sum. The approximation problem is to return a value  $z \leq t$  such that

$$z \geq z^*(1 - \epsilon).$$

If we think of this as a knapsack problem, we want our knapsack to be within a factor of  $(1 - \epsilon)$  of being as full as possible. So, if  $\epsilon = 0.1$ , then the knapsack should be at least 90% as full as the best possible.

What do we mean by polynomial time here? Recall that the running time should be polynomial in the size of the input length. Obviously  $n$  is part of the input length. But  $t$  and the numbers  $x_i$  could also be huge binary numbers. Normally we just assume that a binary number can fit into a word of our computer, and do not count their length. In this case we will be on the safe side. Clearly  $t$  requires  $O(\log t)$  digits to be store in the input. We will take the input size to be  $n + \log t$ .

Intuitively it is not hard to believe that it should be possible to determine whether we can fill the knapsack to within 90% of optimal. After all, we are used to solving similar sorts of packing problems all the time in real life. But the mental heuristics that we apply to these problems are not necessarily easy to convert into efficient algorithms. Our intuition tells us that we can afford to be a little “sloppy” in keeping track of exactly full the knapsack is at any point. The value of  $\epsilon$  tells us just how sloppy we can be. Our approximation will do something similar. First we consider an exponential time algorithm, and then convert it into an approximation algorithm.

**Exponential Time Algorithm:** This algorithm is a variation of the dynamic programming solution we gave for the knapsack problem. Recall that there we used a 2-dimensional array to keep track of whether we could fill a knapsack of a given capacity with the first  $i$  objects. We will do something similar here. As before, we will concentrate on the question of which sums are possible, but determining the subsets that give these sums will not be hard.

Let  $L_i$  denote a list of integers that contains the sums of all  $2^i$  subsets of  $\{x_1, x_2, \dots, x_i\}$  (including the empty set whose sum is 0). For example, for the set  $\{1, 4, 6\}$  the corresponding list of sums contains  $\langle 0, 1, 4, 5(= 1 + 4), 6, 7(= 1 + 6), 10(= 4 + 6), 11(= 1 + 4 + 6) \rangle$ . Note that  $L_i$  can have as many as  $2^i$  elements, but may have fewer, since some subsets may have the same sum.

There are two things we will want to do for efficiency. (1) Remove any duplicates from  $L_i$ , and (2) only keep sums that are less than or equal to  $t$ . Let us suppose that we a procedure `MergeLists(L1, L2)` which merges two sorted lists, and returns a sorted lists with all duplicates removed. This is essentially the procedure used in MergeSort but with the added duplicate element test. As a bit of notation, let  $L + x$  denote the list resulting by adding the number  $x$  to every element of list  $L$ . Thus  $\langle 1, 4, 6 \rangle + 3 = \langle 4, 7, 9 \rangle$ . This gives the following procedure for the subset sum problem.

---

Exact Subset Sum

```
Exact_SS(x[1..n], t) {
  L = <0>;
  for i = 1 to n do {
    L = MergeLists(L, L+x[i]);
    remove for L all elements greater than t;
  }
  return largest element in L;
}
```

---

For example, if  $S = \{1, 4, 6\}$  and  $t = 8$  then the successive lists would be

$$\begin{aligned} L_0 &= \langle 0 \rangle \\ L_1 &= \langle 0 \rangle \cup \langle 0 + 1 \rangle = \langle 0, 1 \rangle \\ L_2 &= \langle 0, 1 \rangle \cup \langle 0 + 4, 1 + 4 \rangle = \langle 0, 1, 4, 5 \rangle \\ L_3 &= \langle 0, 1, 4, 5 \rangle \cup \langle 0 + 6, 1 + 6, 4 + 6, 5 + 6 \rangle = \langle 0, 1, 4, 5, 6, 7, 10, 11 \rangle. \end{aligned}$$

The last list would have the elements 10 and 11 removed, and the final answer would be 7. The algorithm runs in  $\Omega(2^n)$  time in the worst case, because this is the number of sums that are generated if there are no duplicates, and no items are removed.

**Approximation Algorithm:** To convert this into an approximation algorithm, we will introduce a “trim” the lists to decrease their sizes. The idea is that if the list  $L$  contains two numbers that are very close to one another, e.g. 91,048 and 91,050, then we should not need to keep both of these numbers in the list. One of them is good enough for future approximations. This will reduce the size of the lists that the algorithm needs to maintain. But, how much trimming can we allow and still keep our approximation bound? Furthermore, will we be able to reduce the list sizes from exponential to polynomial?

The answer to both these questions is yes, provided you apply a proper way of trimming the lists. We will trim elements whose values are sufficiently close to each other. But we should define close in manner that is relative to the sizes of the numbers involved. The trimming must also depend on  $\epsilon$ . We select  $\delta = \epsilon/n$ . (Why? We will see later that this is the value that makes everything work out in the end.) Note that  $0 < \delta < 1$ . Assume that the elements of  $L$  are sorted. We walk through the list. Let  $z$  denote the last untrimmed element in  $L$ , and let  $y \geq z$  be the next element to be considered. If

$$\frac{y - z}{y} \leq \delta$$

then we trim  $y$  from the list. Equivalently, this means that the final trimmed list cannot contain two value  $y$  and  $z$  such that

$$(1 - \delta)y \leq z \leq y.$$

We can think of  $z$  as *representing*  $y$  in the list.

For example, given  $\delta = 0.1$  and given the list

$$L = \langle 10, 11, 12, 15, 20, 21, 22, 23, 24, 29 \rangle,$$

the trimmed list  $L'$  will consist of

$$L' = \langle 10, 12, 15, 20, 23, 29 \rangle.$$

Another way to visualize trimming is to break the interval from  $[1, t]$  into a set of *buckets* of exponentially increasing size. Let  $d = 1/(1 - \delta)$ . Note that  $d > 1$ . Consider the intervals

$$[1, d], [d, d^2], [d^2, d^3], \dots, [d^{k-1}, d^k],$$

where  $d^k \geq t$ . If  $z \leq y$  are in the same interval  $[d^{i-1}, d^i]$  then

$$\frac{y - z}{y} \leq \frac{d^i - d^{i-1}}{d^i} = 1 - \frac{1}{d} = \delta.$$

Thus, we cannot have more than one item within each bucket. We can think of trimming as a way of enforcing the condition that items in our lists are not relatively too close to one another, by enforcing the condition that no bucket has more than one item.

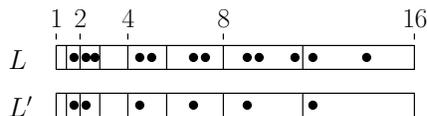


Fig. 122: Trimming Lists for Approximate Subset Sum.

**Claim:** The number of distinct items in a trimmed list is  $O((n \log t)/\epsilon)$ , which is polynomial in input size and  $1/\epsilon$ .

**Proof:** We know that each pair of consecutive elements in a trimmed list differ by a ratio of at least  $d = 1/(1 - \delta) > 1$ . Let  $k$  denote the number of elements in the trimmed list, ignoring the element of value 0. Thus, the smallest nonzero value and maximum value in the trimmed list differ by a ratio of at least  $d^{k-1}$ . Since the smallest (nonzero) element is at least as large as 1, and the largest is no larger than  $t$ , then it follows that  $d^{k-1} \leq t/1 = t$ . Taking the natural log of both sides we

have  $(k - 1) \ln d \leq \ln t$ . Using the facts that  $\delta = \epsilon/n$  and the log identity that  $\ln(1 + x) \leq x$ , we have

$$\begin{aligned} k - 1 &\leq \frac{\ln t}{\ln d} = \frac{\ln t}{-\ln(1 - \delta)} \\ &\leq \frac{\ln t}{\delta} = \frac{n \ln t}{\epsilon} \\ k &= O\left(\frac{n \log t}{\epsilon}\right). \end{aligned}$$

Observe that the input size is at least as large as  $n$  (since there are  $n$  numbers) and at least as large as  $\log t$  (since it takes  $\log t$  digits to write down  $t$  on the input). Thus, this function is polynomial in the input size and  $1/\epsilon$ .

The approximation algorithm operates as before, but in addition we call the procedure `Trim` given below.

---

Approximate Subset Sum

```
Trim(L, delta) {
  let the elements of L be denoted y[1..m];
  L' = <y[1]>; // start with first item
  last = y[1]; // last item to be added
  for i = 2 to m do {
    if (last < (1-delta) y[i]) { // different enough?
      append y[i] to end of L';
      last = y[i];
    }
  }
}

Approx_SS(x[1..n], t, eps) {
  delta = eps/n; // approx factor
  L = <0>; // empty sum = 0
  for i = 1 to n do {
    L = MergeLists(L, L+x[i]); // add in next item
    L = Trim(L, delta); // trim away "near" duplicates
    remove for L all elements greater than t;
  }
  return largest element in L;
}
```

---

For example, consider the set  $S = \{104, 102, 201, 101\}$  and  $t = 308$  and  $\epsilon = 0.20$ . We have  $\delta = \epsilon/4 = 0.05$ . An example of the algorithm's execution is shown in Fig. 123.

The final output is 302. The optimum is  $307 = 104 + 102 + 101$ . So our actual relative error in this case is within 2%.

The running time of the procedure is  $O(n|L|)$  which is  $O(n^2 \ln t/\epsilon)$  by the earlier claim.

**Approximation Analysis:** The final question is why the algorithm achieves an relative error of at most  $\epsilon$  over the optimum solution. Let  $Y^*$  denote the optimum (largest) subset sum and let  $Y$  denote the value returned by the algorithm. We want to show that  $Y$  is not too much smaller than  $Y^*$ , that is,

$$Y \geq Y^*(1 - \epsilon).$$

Our proof will make use of an important inequality from real analysis.

init:  $L_0 = \langle 0 \rangle$   
  
 merge:  $L_1 = \langle 0, 104 \rangle$   
 trim:  $L_1 = \langle 0, 104 \rangle$   
 remove:  $L_1 = \langle 0, 104 \rangle$   
  
 merge:  $L_2 = \langle 0, 102, 104, 206 \rangle$   
 trim:  $L_2 = \langle 0, 102, 206 \rangle$   
 remove:  $L_2 = \langle 0, 102, 206 \rangle$   
  
 merge:  $L_3 = \langle 0, 102, 201, 206, 303, 407 \rangle$   
 trim:  $L_3 = \langle 0, 102, 201, 303, 407 \rangle$   
 remove:  $L_3 = \langle 0, 102, 201, 303 \rangle$   
  
 merge:  $L_4 = \langle 0, 101, 102, 201, 203, 302, 303, 404 \rangle$   
 trim:  $L_4 = \langle 0, 101, 201, 302, 404 \rangle$   
 remove:  $L_4 = \langle 0, 101, 201, 302 \rangle$

Fig. 123: Subset-sum approximation example.

**Lemma:** For  $n > 0$  and  $a$  real numbers,

$$(1 + a) \leq \left(1 + \frac{a}{n}\right)^n \leq e^a.$$

Recall that our intuition was that we would allow a relative error of  $\epsilon/n$  at each stage of the algorithm. Since the algorithm has  $n$  stages, then the total relative error should be (obviously?)  $n(\epsilon/n) = \epsilon$ . The catch is that these are relative, not absolute errors. These errors do not accumulate additively, but rather by multiplication. So we need to be more careful.

Let  $L_i^*$  denote the  $i$ -th list in the exponential time (optimal) solution and let  $L_i$  denote the  $i$ -th list in the approximate algorithm. We claim that for each  $y \in L_i^*$  there exists a representative item  $z \in L_i$  whose relative error from  $y$  that satisfies

$$(1 - \epsilon/n)^i y \leq z \leq y.$$

The proof of the claim is by induction on  $i$ . Initially  $L_0 = L_0^* = \langle 0 \rangle$ , and so there is no error. Suppose by induction that the above equation holds for each item in  $L_{i-1}^*$ . Consider an element  $y \in L_{i-1}^*$ . We know that  $y$  will generate two elements in  $L_i^*$ :  $y$  and  $y + x_i$ . We want to argue that there will be a representative that is “close” to each of these items.

By our induction hypothesis, there is a representative element  $z$  in  $L_{i-1}$  such that

$$(1 - \epsilon/n)^{i-1} y \leq z \leq y.$$

When we apply our algorithm, we will form two new items to add (initially) to  $L_i$ :  $z$  and  $z + x_i$ . Observe that by adding  $x_i$  to the inequality above and a little simplification we get

$$(1 - \epsilon/n)^{i-1} (y + x_i) \leq z + x_i \leq y + x_i.$$

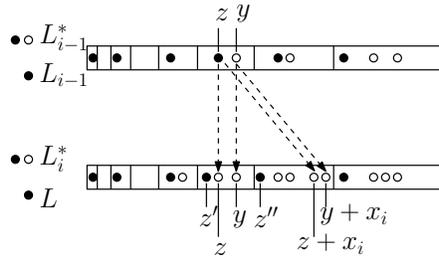


Fig. 124: Subset sum approximation analysis.

The items  $z$  and  $z + x_i$  might not appear in  $L_i$  because they may be trimmed. Let  $z'$  and  $z''$  be their respective representatives. Thus,  $z'$  and  $z''$  are elements of  $L_i$ . We have

$$\begin{aligned} (1 - \epsilon/n)z &\leq z' \leq z \\ (1 - \epsilon/n)(z + x_i) &\leq z'' \leq z + x_i. \end{aligned}$$

Combining these with the inequalities above we have

$$\begin{aligned} (1 - \epsilon/n)^{i-1} (1 - \epsilon/n)y &\leq (1 - \epsilon/n)^i y \leq z' \leq y \\ (1 - \epsilon/n)^{i-1} (1 - \epsilon/n)(y + x_i) &\leq (1 - \epsilon/n)^i (y + x_i) \leq z'' \leq z + y_i. \end{aligned}$$

Since  $z$  and  $z''$  are in  $L_i$  this is the desired result. This ends the proof of the claim.

Using our claim, and the fact that  $Y^*$  (the optimum answer) is the largest element of  $L_n^*$  and  $Y$  (the approximate answer) is the largest element of  $L_n$  we have

$$(1 - \epsilon/n)^n Y^* \leq Y \leq Y^*.$$

This is not quite what we wanted. We wanted to show that  $(1 - \epsilon)Y^* \leq Y$ . To complete the proof, we observe from the lemma above (setting  $a = -\epsilon$ ) that

$$(1 - \epsilon) \leq \left(1 - \frac{\epsilon}{n}\right)^n.$$

This completes the approximate analysis.

## Lecture 37: Approximations: Bin Packing

**Bin Packing:** Bin packing is another well-known NP-complete problem. This is a partitioning problem where we are given a set of objects that are to be partitioned among a collection of containers, called *bins*. Each bin has the same capacity, and the objective is to use the smallest number of bins to hold all the objects.

More formally, we are given a set of  $n$  objects, where  $s_i$  denotes the *size* of the  $i$ th object. It will simplify the presentation to assume that the sizes have been normalized so that  $0 < s_i < 1$ . We want to put these objects into a set of bins. Each bin can hold a subset of objects whose total size is at most 1. The problem is to partition the objects among the bins so as to use the fewest possible bins. (Note that if your bin size is not 1, then you can reduce the problem into this form by simply dividing all sizes by the size of the bin.)

Bin packing arises in many applications. Many of these applications involve not only the size of the object but their geometric shape as well. For example, these include packing boxes into a truck, or cutting the maximum number of pieces of certain shapes out of a piece of sheet metal. However, even if we ignore the geometry, and just consider the sizes of the objects, the decision problem is still NP-complete. (The reduction is from the knapsack problem.)

Here is a simple heuristic algorithm for the bin packing problem, called the *first-fit heuristic*. We start with an unlimited number of empty bins. We take each object in turn, and find the first bin that has space to hold this object. We put this object in this bin. The algorithm is illustrated in Fig. 125. We claim that first-fit uses at most twice as many bins as the optimum. That is, if the optimal solution uses  $b_{\text{opt}}$  bins, and first-fit uses  $b_{\text{ff}}$  bins, then we show below that

$$\frac{b_{\text{ff}}}{b_{\text{opt}}} \leq 2.$$

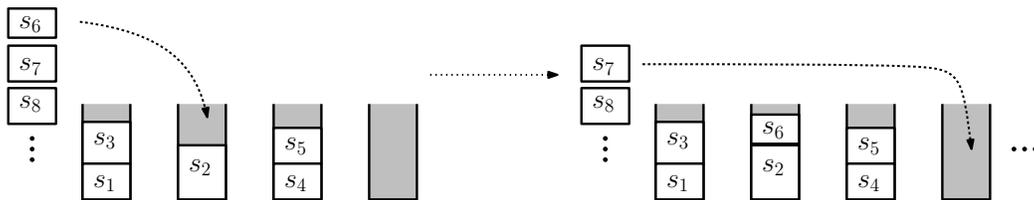


Fig. 125: First-fit Heuristic.

**Theorem:** The first-fit heuristic achieves a ratio bound of 2.

**Proof:** Consider an instance  $\{s_1, \dots, s_n\}$  of the bin packing problem. Let  $S = \sum_i s_i$  denote the sum of all the object sizes. Let  $b_{\text{opt}}$  denote the optimal number of bins, and  $b_{\text{ff}}$  denote the number of bins used by first-fit.

First, observe that since no bin can hold more than one unit's worth of items, and we have a total of  $S$  units to be stored, it follows that we need a minimum of  $S$  bins to store everything. (And this would be achieved only if every bin were filled exactly to the top.) Thus,  $b_{\text{opt}} \geq S$ .

Next, we claim that  $b_{\text{ff}} \leq 2S$ . To see this, let  $t_i$  denote the total size of the objects that first-fit puts into bin  $i$ . There cannot be two bins  $i < j$  such that  $t_i + t_j < 1$ . The reason is that any item we decided to put into bin  $j$  must be small enough to fit into bin  $i$ . Thus, the first-fit algorithm would never put such an item into bin  $j$ . In particular, this implies that for all  $i$ ,  $t_i + t_{i+1} \geq 1$  (where indices are taken circularly modulo the number of bins). Thus we have

$$b_{\text{ff}} = \sum_{i=1}^{b_{\text{ff}}} 1 \leq \sum_{i=1}^{b_{\text{ff}}} (t_i + t_{i+1}) = \sum_{i=1}^{b_{\text{ff}}} t_i + \sum_{i=1}^{b_{\text{ff}}} t_{i+1} = S + S = 2S \leq 2b_{\text{opt}},$$

which completes the proof.

There are in fact a number of other heuristics for bin packing. Another example is *best-fit*, which attempts to put the object into the bin in which it fits most closely with the available space (assuming that there is sufficient available space). This is not necessarily a good idea, since it might tend to create very small spaces that will be hard to fill. There is also a variant of first-fit, called *first-fit-decreasing*, in which the objects are first sorted in decreasing order of size. (This makes intuitive sense, because it is best to first load the big items, and then try to squeeze the smaller objects into the remaining space.)

A more careful (an complicated) proof establishes that first-fit has a approximation ratio that is a bit smaller than 2, and in fact  $17/10 = 1.7$  is possible. Best-fit has a very similar bound. It can be shown that first-fit-decreasing has a significantly better bound than either of these. In particular, it achieves a ratio bound of  $11/9 \approx 1.222$ .