

Class notes for CSCI 405 Spring 2024

Kameron Decker Harris

May 12, 2024

Textbook: “Introduction to Algorithms” by Cormen, Leiserson, Rivest, and Stein (4th ed). Note that this is the first quarter adapting to the new textbook. Please let me know if page/formula numbers are incorrect.

Chapter 9: Order stats and medians

Reading: Sections 9.1, 9.2

Intro material and proof for 4th edition

Lecture Notes for Chapter 9:

Medians and Order Statistics

Chapter 9 overview

- ***i th order statistic*** is the i th smallest element of a set of n elements.
- The ***minimum*** is the first order statistic ($i = 1$).
- The ***maximum*** is the n th order statistic ($i = n$).
- A ***median*** is the “halfway point” of the set.
- When n is odd, the median is unique, at $i = (n + 1)/2$.
- When n is even, there are two medians:
 - The ***lower median***, at $i = n/2$, and
 - The ***upper median***, at $i = n/2 + 1$.
 - We mean lower median when we use the phrase “the median.”

The ***selection problem***:

Input: A set A of n distinct numbers and a number i , with $1 \leq i \leq n$.

Output: The element $x \in A$ that is larger than exactly $i - 1$ other elements in A .
In other words, the i th smallest element of A .

Easy to solve the selection problem in $O(n \lg n)$ time:

- Sort the numbers using an $O(n \lg n)$ -time algorithm, such as heapsort or merge sort.
- Then return the i th element in the sorted array.

There are faster algorithms, however.

- First, we’ll look at the problem of selecting the minimum and maximum of a set of elements.
- Then, we’ll look at a simple general selection algorithm with a time bound of $O(n)$ in the average case.
- Finally, we’ll look at a more complicated general selection algorithm with a time bound of $O(n)$ in the worst case.

Minimum and maximum

We can easily obtain an upper bound of $n - 1$ comparisons for finding the minimum of a set of n elements.

- Examine each element in turn and keep track of the smallest one.
- This is the best we can do, because each element, except the minimum, must be compared to a smaller element at least once.

The following pseudocode finds the minimum element in array $A[1 : n]$:

```

MINIMUM( $A, n$ )
   $min = A[1]$ 
  for  $i = 2$  to  $n$ 
    if  $min > A[i]$ 
       $min = A[i]$ 
  return  $min$ 

```

The maximum can be found in exactly the same way by replacing the $>$ with $<$ in the above algorithm.

Simultaneous minimum and maximum

Some applications need both the minimum and maximum of a set of elements.

- For example, a graphics program may need to scale a set of (x, y) data to fit onto a rectangular display. To do so, the program must first find the minimum and maximum of each coordinate.

A simple algorithm to find the minimum and maximum is to find each one independently. There will be $n - 1$ comparisons for the minimum and $n - 1$ comparisons for the maximum, for a total of $2n - 2$ comparisons. This will result in $\Theta(n)$ time. In fact, at most $3 \lfloor n/2 \rfloor$ comparisons suffice to find both the minimum and maximum:

- Maintain the minimum and maximum of elements seen so far.
- Don't compare each element to the minimum and maximum separately.
- Process elements in pairs.
- Compare the elements of a pair to each other.
- Then compare the larger element to the maximum so far, and compare the smaller element to the minimum so far.

This leads to only 3 comparisons for every 2 elements.

Setting up the initial values for the min and max depends on whether n is odd or even.

- If n is even, compare the first two elements and assign the larger to max and the smaller to min. Then process the rest of the elements in pairs.
- If n is odd, set both min and max to the first element. Then process the rest of the elements in pairs.

Analysis of the total number of comparisons

- If n is even, do 1 initial comparison and then $3(n - 2)/2$ more comparisons.

$$\begin{aligned}
 \# \text{ of comparisons} &= \frac{3(n - 2)}{2} + 1 \\
 &= \frac{3n - 6}{2} + 1 \\
 &= \frac{3n}{2} - 3 + 1 \\
 &= \frac{3n}{2} - 2.
 \end{aligned}$$

- If n is odd, do $3(n - 1)/2 = 3 \lfloor n/2 \rfloor$ comparisons.

In either case, the maximum number of comparisons is $\leq 3 \lfloor n/2 \rfloor$.

Selection in expected linear time

Selection of the i th smallest element of the array A can be done in $\Theta(n)$ time.

The function RANDOMIZED-SELECT uses RANDOMIZED-PARTITION from the quicksort algorithm in Chapter 7. RANDOMIZED-SELECT differs from quicksort because it recurses on one side of the partition only.

```

RANDOMIZED-SELECT( $A, p, r, i$ )
    if  $p == r$ 
        return  $A[p]$           //  $1 \leq i \leq r - p + 1$  when  $p == r$  means that  $i = 1$ 
     $q = \text{RANDOMIZED-PARTITION}(A, p, r)$ 
     $k = q - p + 1$ 
    if  $i == k$ 
        return  $A[q]$           // the pivot value is the answer
    elseif  $i < k$ 
        return RANDOMIZED-SELECT( $A, p, q - 1, i$ )
    else return RANDOMIZED-SELECT( $A, q + 1, r, i - k$ )

```

After the call to RANDOMIZED-PARTITION, the array is partitioned into two subarrays $A[p : q - 1]$ and $A[q + 1 : r]$, along with a *pivot* element $A[q]$.

- The elements of subarray $A[p : q - 1]$ are all $\leq A[q]$.
- The elements of subarray $A[q + 1 : r]$ are all $> A[q]$.
- The pivot element is the k th element of the subarray $A[p : r]$, where $k = q - p + 1$.
- If the pivot element is the i th smallest element (i.e., $i = k$), return $A[q]$.
- Otherwise, recurse on the subarray containing the i th smallest element.
 - If $i < k$, this subarray is $A[p : q - 1]$, and we want the i th smallest element.
 - If $i > k$, this subarray is $A[q + 1 : r]$ and, since there are k elements in $A[p : r]$ that precede $A[q + 1 : r]$, we want the $(i - k)$ th smallest element of this subarray.

Analysis

Worst-case running time

$\Theta(n^2)$, because we could be extremely unlucky and always recurse on a subarray that is only one element smaller than the previous subarray.

Expected running time

RANDOMIZED-SELECT works well on average. Because it is randomized, no particular input brings out the worst-case behavior consistently.

Analysis assumes that the recursion goes as deep as possible: until only one element remains.

Intuition: Suppose that each pivot is in the second or third quartiles if the elements were sorted—in the “middle half.” Then at least $1/4$ of the remaining elements are ignored in all future recursive calls \Rightarrow at most $3/4$ of the elements are still *in play*: somewhere within $A[p : r]$. RANDOMIZE-PARTITION takes $\Theta(n)$ time to partition n elements \Rightarrow recurrence would be $T(n) = T(3n/4) + \Theta(n) = \Theta(n)$ by case 3 of the master method.

What if the pivot is not always in the middle half? Probability that it is in the middle half is $1/2$. View selecting a pivot in the middle half as a Bernoulli trial with probability of success $1/2$. Then the number of trials before a success is a geometric distribution with expected value 2. So that half the time, $1/4$ of the elements go out of play, and the other half of the time, as few as one element (the pivot) goes out of play. But that just doubles the running time, so still expect $\Theta(n)$.

Rigorous analysis:

- Define $A^{(j)}$ as the set of elements still in play (within $A[p : r]$) after j recursive calls (i.e., after j th partitioning). $A^{(0)}$ is all the elements in A .
- $|A^{(j)}|$ is a random variable that depends on A and order statistic i , but not on the order of elements in A .
- Each partitioning removes at least one element (the pivot) \Rightarrow sizes of $A^{(j)}$ strictly decrease.
- j th partitioning takes set $A^{(j-1)}$ and produces $A^{(j)}$.
- Assume a 0th “dummy” partitioning that produces $A^{(0)}$.
- j th partitioning is *helpful* if $|A^{(j)}| \leq (3/4)|A^{(j-1)}|$. Not all partitionings are necessarily helpful. Think of a helpful partitioning as a successful Bernoulli trial.

Lemma

A partitioning is helpful with probability $\geq 1/2$.

Proof

- Whether or not a partitioning is helpful depends on the randomly chosen pivot.
- Define “middle half” of an n -element subarray as all but the smallest $\lceil n/4 \rceil - 1$ and greatest $\lceil n/4 \rceil - 1$ elements. That is, all but the first and last $\lceil n/4 \rceil - 1$ if the subarray were sorted.

- Will show that if the pivot is in the middle half, then that pivot leads to a helpful partitioning and that the probability that the pivot is in the middle half is $\geq 1/2$.
- No matter where the pivot lies, either all elements $>$ pivot or all elements $<$ pivot, and the pivot itself, are not in play after partitioning \Rightarrow if the pivot is in the middle half, at least the smallest $\lceil n/4 \rceil - 1$ or greatest $\lceil n/4 \rceil - 1$ elements, plus the pivot, will not be in play after partitioning $\Rightarrow \geq \lceil n/4 \rceil$ elements not in play.
- Then, at most $n - \lceil n/4 \rceil = \lfloor 3n/4 \rfloor < 3n/4$ elements in play \Rightarrow partitioning is helpful. ($n - \lceil n/4 \rceil = \lfloor 3n/4 \rfloor$ is from Exercise 3.3-2.)
- To find a lower bound on the probability that a randomly chosen pivot is in the middle half, find an upper bound on the probability that it is not:

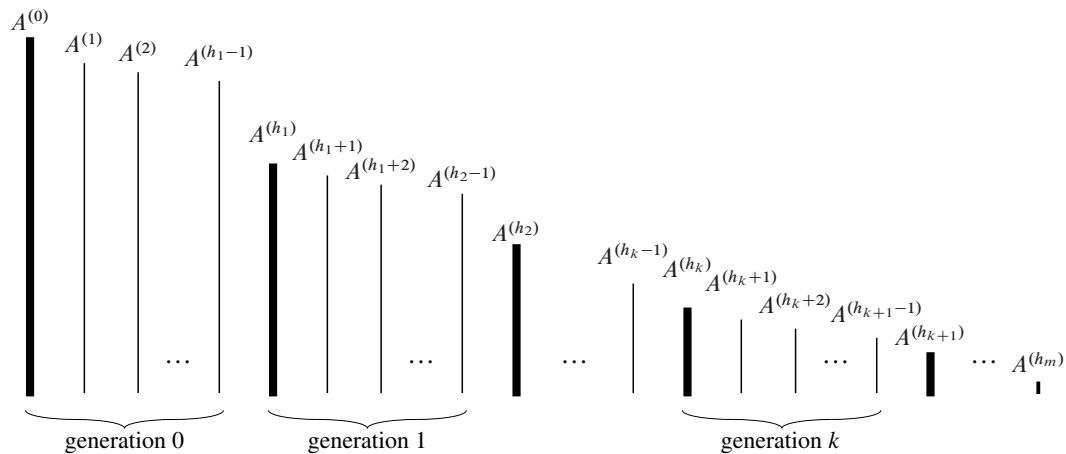
$$\begin{aligned} \frac{2(\lceil n/4 \rceil - 1)}{n} &\leq \frac{2((n/4 + 1) - 1)}{n} \quad (\text{inequality (3.2)}) \\ &= \frac{n/2}{n} \\ &= 1/2. \end{aligned}$$
- Since the pivot has probability $\geq 1/2$ of falling into the middle half, a partitioning is helpful with probability $\geq 1/2$. ■ (lemma)

Theorem

The expected running time of RANDOMIZED-SELECT is $\Theta(n)$.

Proof

- Let the sequence of helpful partitionings be $\langle h_0, h_1, \dots, h_m \rangle$. Consider the 0th partitioning as helpful $\Rightarrow h_0 = 0$. Can bound m , since after at most $\lceil \log_{4/3} n \rceil$ helpful partitionings, only one element remains in play.
- Define $n_k = |A^{(h_k)}|$ and $n_0 = |A^{(0)}|$, the original problem size. $n_k = |A^{(h_k)}| \leq (3/4)|A^{(h_{k-1})}| = (3/4)n_{k-1}$ for $k = 1, 2, \dots, m$.
- Iterating gives $n_k \leq (3/4)^k n_0$.
- Break up sets into m “generations.” The sets in generation k are $A^{(h_k)}, A^{(h_{k+1})}, \dots, A^{(h_{k+1}-1)}$, where $A^{(h_k)}$ is the result of a helpful partitioning and $A^{(h_{k+1}-1)}$ is the last set before the next helpful partitioning.



[Height of each line indicates the size of the set (number of elements in play). Heavy lines are sets $A^{(h_k)}$, resulting from helpful partitionings and are first within their generation. Other lines are not first within their generation. A generation may contain just one set.]

- If $A^{(j)}$ is in the k th generation, then $|A^{(j)}| \leq |A^{(h_k)}| = n_k \leq (3/4)^k n_0$.
- Define random variable $X_k = h_{k+1} - h_k$ as the number of sets in the k th generation $\Rightarrow k$ th generation includes sets $A^{(h_k)}, A^{(h_k+1)}, \dots, A^{(h_k+X_k-1)}$.
- By previous lemma, a partitioning is helpful with probability $\geq 1/2$. The probability is even higher, since a partitioning is helpful even if the pivot doesn't fall into middle half, but the i th smallest element lies in the smaller side. Just use the $1/2$ lower bound $\Rightarrow E[X_k] \leq 2$ for $k = 0, 1, \dots, m-1$ (by equation (C.36), expectation of a geometric distribution).
- The total running time is dominated by the comparisons during partitioning. The j th partitioning takes $A^{(j-1)}$ and compares the pivot with all the other $|A^{(j-1)}| - 1$ elements $\Rightarrow j$ th partitioning makes $< |A^{(j-1)}|$ comparisons.
- The total number of comparisons is less than

$$\begin{aligned} \sum_{k=0}^{m-1} \sum_{j=h_k}^{h_k+X_k-1} |A^{(j)}| &\leq \sum_{k=0}^{m-1} \sum_{j=h_k}^{h_k+X_k-1} |A^{(h_k)}| \\ &= \sum_{k=0}^{m-1} X_k |A^{(h_k)}| \\ &\leq \sum_{k=0}^{m-1} X_k \left(\frac{3}{4}\right)^k n_0. \end{aligned}$$

- Since $E[X_k] \leq 2$, the expected total number of comparisons is less than

$$\begin{aligned} E \left[\sum_{k=0}^{m-1} X_k \left(\frac{3}{4}\right)^k n_0 \right] &= \sum_{k=0}^{m-1} E \left[X_k \left(\frac{3}{4}\right)^k n_0 \right] \quad (\text{linearity of expectation}) \\ &= n_0 \sum_{k=0}^{m-1} \left(\frac{3}{4}\right)^k E[X_k] \\ &\leq 2n_0 \sum_{k=0}^{m-1} \left(\frac{3}{4}\right)^k \\ &< 2n_0 \sum_{k=0}^{\infty} \left(\frac{3}{4}\right)^k \\ &= 8n_0 \quad (\text{infinite geometric series}). \end{aligned}$$

- n_0 is the size of the original array $A \Rightarrow$ an $O(n)$ upper bound on the expected running time. For the lower bound, the first call of RANDOMIZED-PARTITION examines all n elements $\Rightarrow \Theta(n)$. ■ (theorem)

Therefore, we can determine any order statistic in linear time on average, assuming that all elements are distinct.

3rd edition proof (presented in class)

Analysis

Worst-case running time

$\Theta(n^2)$, because we could be extremely unlucky and always recurse on a subarray that is only 1 element smaller than the previous subarray.

Expected running time

RANDOMIZED-SELECT works well on average. Because it is randomized, no particular input brings out the worst-case behavior consistently.

The running time of RANDOMIZED-SELECT is a random variable that we denote by $T(n)$. We obtain an upper bound on $E[T(n)]$ as follows:

- RANDOMIZED-PARTITION is equally likely to return any element of A as the pivot.
- For each k such that $1 \leq k \leq n$, the subarray $A[p..q]$ has k elements (all \leq pivot) with probability $1/n$. [Note that we're now considering a subarray that includes the pivot, along with elements less than the pivot.]
- For $k = 1, 2, \dots, n$, define indicator random variable

$$X_k = I\{\text{subarray } A[p..q] \text{ has exactly } k \text{ elements}\}.$$

- Since $\Pr\{\text{subarray } A[p..q] \text{ has exactly } k \text{ elements}\} = 1/n$, Lemma 5.1 says that $E[X_k] = 1/n$.
- When we call RANDOMIZED-SELECT, we don't know if it will terminate immediately with the correct answer, recurse on $A[p..q-1]$, or recurse on $A[q+1..r]$. It depends on whether the i th smallest element is less than, equal to, or greater than the pivot element $A[q]$.
- To obtain an upper bound, we assume that $T(n)$ is monotonically increasing and that the i th smallest element is always in the larger subarray.
- For a given call of RANDOMIZED-SELECT, $X_k = 1$ for exactly one value of k , and $X_k = 0$ for all other k .
- When $X_k = 1$, the two subarrays have sizes $k-1$ and $n-k$.
- For a subproblem of size n , RANDOMIZED-PARTITION takes $O(n)$ time. [Actually, it takes $\Theta(n)$ time, but $O(n)$ suffices, since we're obtaining only an upper bound on the expected running time.]
- Therefore, we have the recurrence

$$\begin{aligned} T(n) &\leq \sum_{k=1}^n X_k \cdot (T(\max(k-1, n-k)) + O(n)) \\ &= \sum_{k=1}^n X_k \cdot T(\max(k-1, n-k)) + O(n). \end{aligned}$$

- Taking expected values gives

$$\begin{aligned} E[T(n)] &\leq E\left[\sum_{k=1}^n X_k \cdot T(\max(k-1, n-k)) + O(n)\right] \end{aligned}$$

$$\begin{aligned}
&= \sum_{k=1}^n \mathbb{E}[X_k \cdot T(\max(k-1, n-k))] + O(n) \quad (\text{linearity of expectation}) \\
&= \sum_{k=1}^n \mathbb{E}[X_k] \cdot \mathbb{E}[T(\max(k-1, n-k))] + O(n) \quad (\text{equation (C.24)}) \\
&= \sum_{k=1}^n \frac{1}{n} \cdot \mathbb{E}[T(\max(k-1, n-k))] + O(n) .
\end{aligned}$$

- We rely on X_k and $T(\max(k-1, n-k))$ being independent random variables in order to apply equation (C.24).
- Looking at the expression $\max(k-1, n-k)$, we have

$$\max(k-1, n-k) = \begin{cases} k-1 & \text{if } k > \lceil n/2 \rceil , \\ n-k & \text{if } k \leq \lceil n/2 \rceil . \end{cases}$$

- If n is even, each term from $T(\lceil n/2 \rceil)$ up to $T(n-1)$ appears exactly twice in the summation.
- If n is odd, these terms appear twice and $T(\lfloor n/2 \rfloor)$ appears once.
- Either way,

$$\mathbb{E}[T(n)] \leq \frac{2}{n} \sum_{k=\lceil n/2 \rceil}^{n-1} \mathbb{E}[T(k)] + O(n) .$$

- Solve this recurrence by substitution:
 - Guess that $T(n) \leq cn$ for some constant c that satisfies the initial conditions of the recurrence.
 - Assume that $T(n) = O(1)$ for $n < \text{some constant}$. We'll pick this constant later.
 - Also pick a constant a such that the function described by the $O(n)$ term is bounded from above by an for all $n > 0$.
 - Using this guess and constants c and a , we have

$$\begin{aligned}
\mathbb{E}[T(n)] &\leq \frac{2}{n} \sum_{k=\lceil n/2 \rceil}^{n-1} ck + an \\
&= \frac{2c}{n} \left(\sum_{k=1}^{n-1} k - \sum_{k=1}^{\lfloor n/2 \rfloor - 1} k \right) + an \\
&= \frac{2c}{n} \left(\frac{(n-1)n}{2} - \frac{(\lfloor n/2 \rfloor - 1) \lfloor n/2 \rfloor}{2} \right) + an \\
&\leq \frac{2c}{n} \left(\frac{(n-1)n}{2} - \frac{(n/2 - 2)(n/2 - 1)}{2} \right) + an \\
&= \frac{2c}{n} \left(\frac{n^2 - n}{2} - \frac{n^2/4 - 3n/2 + 2}{2} \right) + an \\
&= \frac{c}{n} \left(\frac{3n^2}{4} + \frac{n}{2} - 2 \right) + an
\end{aligned}$$

$$\begin{aligned}
&= c \left(\frac{3n}{4} + \frac{1}{2} - \frac{2}{n} \right) + an \\
&\leq \frac{3cn}{4} + \frac{c}{2} + an \\
&= cn - \left(\frac{cn}{4} - \frac{c}{2} - an \right).
\end{aligned}$$

- To complete this proof, we choose c such that

$$\begin{aligned}
cn/4 - c/2 - an &\geq 0 \\
cn/4 - an &\geq c/2 \\
n(c/4 - a) &\geq c/2 \\
n &\geq \frac{c/2}{c/4 - a} \\
n &\geq \frac{2c}{c - 4a}.
\end{aligned}$$

- Thus, as long as we assume that $T(n) = O(1)$ for $n < 2c/(c - 4a)$, we have $E[T(n)] = O(n)$.

Therefore, we can determine any order statistic in linear time on average.

Selection in worst-case linear time

We can find the i th smallest element in $O(n)$ time *in the worst case*. We'll describe a procedure SELECT that does so.

SELECT recursively partitions the input array.

- **Idea:** Guarantee a good split when the array is partitioned.
- Will use the deterministic procedure PARTITION, but with a small modification. Instead of assuming that the last element of the subarray is the pivot, the modified PARTITION procedure is told which element to use as the pivot.

SELECT works on an array of $n > 1$ elements. It executes the following steps:

1. Divide the n elements into groups of 5. Get $\lceil n/5 \rceil$ groups: $\lfloor n/5 \rfloor$ groups with exactly 5 elements and, if 5 does not divide n , one group with the remaining $n \bmod 5$ elements.
2. Find the median of each of the $\lceil n/5 \rceil$ groups:
 - Run insertion sort on each group. Takes $O(1)$ time per group since each group has ≤ 5 elements.
 - Then just pick the median from each group, in $O(1)$ time.
3. Find the median x of the $\lceil n/5 \rceil$ medians by a recursive call to SELECT. (If $\lceil n/5 \rceil$ is even, then follow our convention and find the lower median.)
4. Using the modified version of PARTITION that takes the pivot element as input, partition the input array around x . Let x be the k th element of the array after partitioning, so that there are $k - 1$ elements on the low side of the partition and $n - k$ elements on the high side.

Chapter 14: Dynamic programming

Reading: Sections 14.1–14.4

Lecture Notes for Chapter 14:

Dynamic Programming

Dynamic Programming

- Not a specific algorithm, but a technique (like divide-and-conquer).
- Developed back in the day when “programming” meant “tabular method” (like linear programming). Doesn’t really refer to computer programming.
- Used for optimization problems:
 - Find *a* solution with *the* optimal value.
 - Minimization or maximization. (We’ll see both.)

Four-step method

1. Characterize the structure of an optimal solution.
2. Recursively define the value of an optimal solution.
3. Compute the value of an optimal solution, typically in a bottom-up fashion.
4. Construct an optimal solution from computed information.

Rod cutting

How to cut steel rods into pieces in order to maximize the revenue you can get? Each cut is free. Rod lengths are always an integer number of inches.

Input: A length n and table of prices p_i , for $i = 1, 2, \dots, n$.

Output: The maximum revenue obtainable for rods whose lengths sum to n , computed as the sum of the prices for the individual rods.

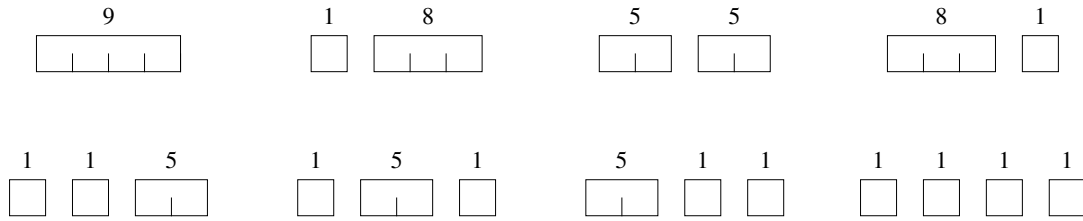
If p_n is large enough, an optimal solution might require no cuts, i.e., just leave the rod as n inches long.

Example: [Using the first 8 values from the example in the textbook.]

length i	1	2	3	4	5	6	7	8
price p_i	1	5	8	9	10	17	17	20

Can cut up a rod in 2^{n-1} different ways, because can choose to cut or not cut after each of the first $n - 1$ inches.

Here are all 8 ways to cut a rod of length 4, with the costs from the example:



The best way is to cut it into two 2-inch pieces, getting a revenue of $p_2 + p_2 = 5 + 5 = 10$.

Let r_i be the maximum revenue for a rod of length i . Can express a solution as a sum of individual rod lengths.

Can determine optimal revenues r_i for the example, by inspection:

i	r_i	optimal solution
1	1	1 (no cuts)
2	5	2 (no cuts)
3	8	3 (no cuts)
4	10	2 + 2
5	13	2 + 3
6	17	6 (no cuts)
7	18	1 + 6 or 2 + 2 + 3
8	22	2 + 6

Can determine optimal revenue r_n by taking the maximum of

- p_n : the revenue from not making a cut,
- $r_1 + r_{n-1}$: the maximum revenue from a rod of 1 inch and a rod of $n - 1$ inches,
- $r_2 + r_{n-2}$: the maximum revenue from a rod of 2 inches and a rod of $n - 2$ inches, ...
- $r_{n-1} + r_1$.

That is,

$$r_n = \max \{p_n, r_1 + r_{n-1}, r_2 + r_{n-2}, \dots, r_{n-1} + r_1\}.$$

Optimal substructure: To solve the original problem of size n , solve subproblems on smaller sizes. After making a cut, two subproblems remain. The optimal solution to the original problem incorporates optimal solutions to the subproblems. May solve the subproblems independently.

Example: For $n = 7$, one of the optimal solutions makes a cut at 3 inches, giving two subproblems, of lengths 3 and 4. Need to solve both of them optimally. The optimal solution for the problem of length 4, cutting into 2 pieces, each of length 2, is used in the optimal solution to the original problem with length 7.

A simpler way to decompose the problem: Every optimal solution has a leftmost cut. In other words, there's some cut that gives a first piece of length i cut off the left end, and a remaining piece of length $n - i$ on the right.

- Need to divide only the remainder, not the first piece.
- Leaves only one subproblem to solve, rather than two subproblems.
- Say that the solution with no cuts has first piece size $i = n$ with revenue p_n , and remainder size 0 with revenue $r_0 = 0$.
- Gives a simpler version of the equation for r_n :

$$r_n = \max \{p_i + r_{n-i} : 1 \leq i \leq n\} .$$

Recursive top-down solution

Direct implementation of the simpler equation for r_n .

The call `CUT-ROD(p, n)` returns the optimal revenue r_n :

CUT-ROD(p, n)

if $n == 0$

```

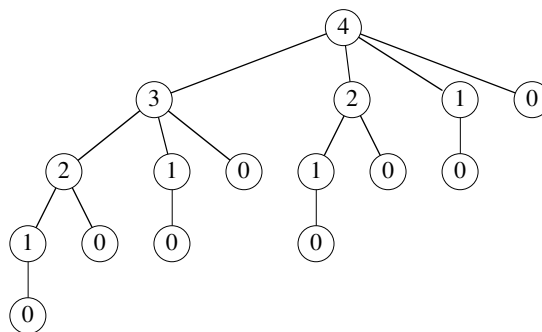
return 0

```

$$q = -\infty$$
for $i = 1$ **to** n
$$q = \max \{q, p[i] + \text{CUT-ROD}(p, n - i)\}$$
return q

This procedure works, but it is terribly *inefficient*. If you code it up and run it, it could take more than an hour for $n = 40$. Running time approximately doubles each time n increases by 1.

Why so inefficient?: CUT-ROD calls itself repeatedly, even on subproblems it has already solved. Here's a tree of recursive calls for $n = 4$. Inside each node is the value of n for the call represented by the node:



Lots of repeated subproblems. Solves the subproblem for size 2 twice, for size 1 four times, and for size 0 eight times.

Exponential growth: Let $T(n)$ equal the number of calls to CUT-ROD with second parameter equal to n . Then

$$T(n) = \begin{cases} 1 & \text{if } n = 0, \\ 1 + \sum_{j=0}^{n-1} T(j) & \text{if } n \geq 1. \end{cases}$$

Summation counts calls where second parameter is $j = n - i$.

Solution to recurrence is $T(n) = 2^n$.

Dynamic-programming solution

Instead of solving the same subproblems repeatedly, arrange to solve each subproblem just once.

Save the solution to a subproblem in a table, and refer back to the table whenever we revisit the subproblem.

“Store, don’t recompute” \Rightarrow time-memory trade-off.

Can turn an exponential-time solution into a polynomial-time solution.

Two basic approaches: top-down with memoization, and bottom-up.

Top-down with memoization

Solve recursively, but store each result in a table.

To find the solution to a subproblem, first look in the table. If the answer is there, use it. Otherwise, compute the solution to the subproblem and then store the solution in the table for future use.

Memoizing is remembering what has been computed previously. [“Memoizing,” not “memorizing.”]

Memoized version of the recursive solution, storing the solution to the subproblem of length i in array entry $r[i]$:

MEMOIZED-CUT-ROD(p, n)

 let $r[0:n]$ be a new array // will remember solution values in r

for $i = 0$ **to** n

$r[i] = -\infty$

return MEMOIZED-CUT-ROD-AUX(p, n, r)

MEMOIZED-CUT-ROD-AUX(p, n, r)

if $r[n] \geq 0$ // already have a solution for length n ?

return $r[n]$

if $n == 0$

$q = 0$

else $q = -\infty$

for $i = 1$ **to** n // i is the position of the first cut

$q = \max \{q, p[i] + \text{MEMOIZED-CUT-ROD-AUX}(p, n - i, r)\}$

$r[n] = q$ // remember the solution value for length n

return q

Bottom-up

Sort the subproblems by size and solve the smaller ones first. That way, when solving a subproblem, have already solved the smaller subproblems needed.

BOTTOM-UP-CUT-ROD(p, n)

```

let  $r[0:n]$  be a new array      // will remember solution values in  $r$ 
 $r[0] = 0$ 
for  $j = 1$  to  $n$                 // for increasing rod length  $j$ 
     $q = -\infty$ 
    for  $i = 1$  to  $j$             //  $i$  is the position of the first cut
         $q = \max \{q, p[i] + r[j - i]\}$ 
     $r[j] = q$                     // remember the solution value for length  $j$ 
return  $r[n]$ 

```

Running time

Both the top-down and bottom-up versions run in $\Theta(n^2)$ time.

- Bottom-up: Doubly nested loops. Number of iterations of inner **for** loop forms an arithmetic series.
- Top-down: MEMOIZED-CUT-ROD solves each subproblem just once, and it solves subproblems for sizes $0, 1, \dots, n$. To solve a subproblem of size n , the **for** loop iterates n times \Rightarrow over all recursive calls, total number of iterations forms an arithmetic series. [Actually using aggregate analysis, which Chapter 16 covers.]

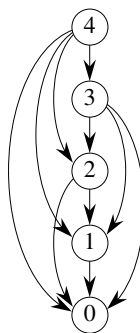
Subproblem graphs

How to understand the subproblems involved and how they depend on each other.

Directed graph:

- One vertex for each distinct subproblem.
- Has a directed edge (x, y) if computing an optimal solution to subproblem x directly requires knowing an optimal solution to subproblem y .

Example: For rod-cutting problem with $n = 4$:



Can think of the subproblem graph as a collapsed version of the tree of recursive calls, where all nodes for the same subproblem are collapsed into a single vertex, and all edges go from parent to child.

Subproblem graph can help determine running time. Because each subproblem is solved just once, running time is sum of times needed to solve each subproblem.

- Time to compute solution to a subproblem is typically linear in the out-degree (number of outgoing edges) of its vertex.
- Number of subproblems equals number of vertices.

When these conditions hold, running time is linear in number of vertices and edges.

Reconstructing a solution

So far, have focused on computing the *value* of an optimal solution, rather than the *choices* that produced an optimal solution.

Extend the bottom-up approach to record not just optimal values, but optimal choices. Save the optimal choices in a separate table. Then use a separate procedure to print the optimal choices.

EXTENDED-BOTTOM-UP-CUT-ROD(p, n)

```

let  $r[0:n]$  and  $s[1:n]$  be new arrays
 $r[0] = 0$ 
for  $j = 1$  to  $n$            // for increasing rod length  $j$ 
     $q = -\infty$ 
    for  $i = 1$  to  $j$          //  $i$  is the position of the first cut
        if  $q < p[i] + r[j - i]$ 
             $q = p[i] + r[j - i]$ 
             $s[j] = i$          // best cut location so far for length  $j$ 
     $r[j] = q$                // remember the solution value for length  $j$ 
return  $r$  and  $s$ 
```

Saves the first cut made in an optimal solution for a problem of size i in $s[i]$.

To print out the cuts made in an optimal solution:

PRINT-CUT-ROD-SOLUTION(p, n)

```

( $r, s$ ) = EXTENDED-BOTTOM-UP-CUT-ROD( $p, n$ )
while  $n > 0$ 
    print  $s[n]$            // cut location for length  $n$ 
     $n = n - s[n]$          // length of the remainder of the rod
```

Example: For the example, EXTENDED-BOTTOM-UP-CUT-ROD returns

i	0	1	2	3	4	5	6	7	8
$r[i]$	0	1	5	8	10	13	17	18	22
$s[i]$			1	2	3	2	2	6	1

A call to PRINT-CUT-ROD-SOLUTION($p, 8$) calls EXTENDED-BOTTOM-UP-CUT-ROD to compute the above r and s tables. Then it prints 2, sets n to 6, prints 6, and finishes (because n becomes 0).

Matrix-chain multiplication

Problem: Given a sequence (chain) $\langle A_1, A_2, \dots, A_n \rangle$ of n matrices, compute the product $A_1 A_2 \cdots A_n$ using standard matrix multiplication (not Strassen's method) while minimizing the number of scalar multiplications.

How to parenthesize the product to minimize the number of scalar multiplications?

Suppose multiplying matrices A and B : $C = A \cdot B$. [The textbook has a procedure to compute $C = C + A \cdot B$, but it's easier in a lecture situation to just use $C = A \cdot B$.] The matrices must be compatible: number of columns of A equals number of rows of B . If A is $p \times q$ and B is $q \times r$, then C is $p \times r$ and takes pqr scalar multiplications.

Example: $A_1 : 10 \times 100$, $A_2 : 100 \times 5$, $A_3 : 5 \times 50$. Compute $A_1 A_2 A_3$, which is 10×50 .

- Try parenthesizing by $((A_1 A_2) A_3)$. First perform $10 \cdot 100 \cdot 5 = 5000$ multiplications, then perform $10 \cdot 5 \cdot 50 = 2500$, for a total of 7500.
- Try parenthesizing by $(A_1 (A_2 A_3))$. First perform $100 \cdot 5 \cdot 50 = 25,000$ multiplications, then perform $10 \cdot 100 \cdot 50 = 50,000$, for a total of 75,000.
- The first way is 10 times faster.

Input to the problem: Let A_i be $p_{i-1} \times p_i$. The input is the sequence of dimensions $\langle p_0, p_1, p_2, \dots, p_n \rangle$.

Note: Not actually multiplying matrices. Just deciding an order with the lowest cost.

Counting the number of parenthesizations

Let $P(n)$ denote the number of ways to parenthesize a product of n matrices. $P(1) = 1$.

When $n \geq 2$, can split anywhere between A_k and A_{k+1} for $k = 1, 2, \dots, n-1$. Then have to split the subproducts. Get

$$P(n) = \begin{cases} 1 & \text{if } n = 1, \\ \sum_{k=1}^{n-1} P(k)P(n-k) & \text{if } n \geq 2. \end{cases}$$

The solution is $P(n) = \Omega(4^n / n^{3/2})$. [The textbook does not prove the solution to this recurrence.] So brute force is a bad strategy.

Step 1: Structure of an optimal solution

Let $A_{i:j}$ be the matrix product $A_i A_{i+1} \dots A_j$.

If $i < j$, then must split between A_k and A_{k+1} for some $i \leq k < j \Rightarrow$ compute $A_{i:k}$ and $A_{k+1:j}$ and then multiply them together. Cost is

- cost of computing $A_{i:k}$
- + cost of computing $A_{k+1:j}$
- + cost of multiplying them together.

Optimal substructure: Suppose that optimal parenthesization of $A_{i:j}$ splits between A_k and A_{k+1} . Then the parenthesization of $A_{i:k}$ must be optimal. Otherwise, if there's a less costly way to parenthesize it, you'd use it and get a parenthesization of $A_{i:j}$ with a lower cost. Same for $A_{k+1:j}$.

Therefore, to build an optimal solution to $A_{i:j}$, split it into how to optimally parenthesize $A_{i:k}$ and $A_{k+1:j}$, find optimal solutions to these subproblems, and then combine the optimal solutions. Need to consider all possible splits.

Step 2: A recursive solution

Define the cost of an optimal solution recursively in terms of optimal subproblem solutions.

Let $m[i, j]$ be the minimum number of scalar multiplications to compute $A_{i:j}$. For the full problem, want $m[1, n]$.

If $i = j$, then just one matrix $\Rightarrow m[i, i] = 0$ for $i = 1, 2, \dots, n$.

If $i < j$, then suppose the optimal split is between A_k and A_{k+1} , where $i \leq k < j$. Then $m[i, j] = m[i, k] + m[k + 1, j] + p_{i-1}p_kp_j$.

But that's assuming you know the value of k . Have to try all possible values and pick the best, so that

$$m[i, j] = \begin{cases} 0 & \text{if } i = j, \\ \min \{m[i, k] + m[k + 1, j] + p_{i-1}p_kp_j : i \leq k < j\} & \text{if } i < j. \end{cases}$$

That formula gives the cost of an optimal solution, but not how to construct it. Define $s[i, j]$ to be a value of k to split $A_{i:j}$ in an optimal parenthesization. Then $s[i, j] = k$ such that $m[i, j] = m[i, k] + m[k + 1, j] + p_{i-1}p_kp_j$.

Step 3: Compute the optimal costs

Could implement a recursive algorithm based on the above equation for $m[i, j]$.

Problem: It would take exponential time.

There are not all that many subproblems: just one for each i, j such that $1 \leq i \leq j \leq n$. There are $\binom{n}{2} + n = \Theta(n^2)$ of them. Thus, a recursive algorithm would solve the same subproblems over and over.

In other words, this problem has overlapping subproblems.

Here is a tabular, bottom-up method to solve the problem. It solves subproblems in order of increasing chain length. The variable $l = j - i + 1$ indicates the chain length.

```

MATRIX-CHAIN-ORDER( $p, n$ )
  let  $m[1:n, 1:n]$  and  $s[1:n-1, 2:n]$  be new tables
  for  $i = 1$  to  $n$  // chain length 1
     $m[i, i] = 0$ 
  for  $l = 2$  to  $n$  //  $l$  is the chain length
    for  $i = 1$  to  $n - l + 1$  // chain begins at  $A_i$ 
       $j = i + l - 1$  // chain ends at  $A_j$ 
       $m[i, j] = \infty$ 
       $m[i, j] = \infty$ 
      for  $k = i$  to  $j - 1$ 
         $q = m[i, k] + m[k + 1, j] + p_{i-1}p_kp_j$ 
        if  $q < m[i, j]$ 
           $m[i, j] = q$  // remember this cost
           $s[i, j] = k$  // remember this index
  return  $m$  and  $s$ 

```

All n chains of length 1 are initialized so that $m[i, i] = 0$ for $i = 1, 2, \dots, n$. Then $n - 1$ chains of length 2 are computed, then $n - 2$ chains of length 3, and so on, up to 1 chain of length n .

[We don't include an example here because the arithmetic is hard for students to process in real time.]

Time: $O(n^3)$, from triply nested loops. Also $\Omega(n^3) \Rightarrow \Theta(n^3)$.

Step 4: Construct an optimal solution

With the s table filled in, recursively print an optimal solution.

```

PRINT-OPTIMAL-PARENS( $s, i, j$ )
  if  $i == j$ 
    print " $A$ " $i$ 
  else print "("
    PRINT-OPTIMAL-PARENS( $s, i, s[i, j]$ )
    PRINT-OPTIMAL-PARENS( $s, s[i, j] + 1, j$ )
  print ")"

```

Initial call is PRINT-OPTIMAL-PARENS($s, 1, n$)

Longest common subsequence

[The textbook has the section on elements of dynamic programming next, but these lecture notes reserve that section for the end of the chapter so that it may refer to two more examples of dynamic programming.]

Problem: Given two sequences, $X = \langle x_1, \dots, x_m \rangle$ and $Y = \langle y_1, \dots, y_n \rangle$. Find a subsequence common to both whose length is longest. A subsequence doesn't have to be consecutive, but it has to be in order.

[To come up with examples of longest common subsequences, search the dictionary for all words that contain the word you are looking for as a subsequence. On a UNIX system, for example, to find all the words with *pine* as a subsequence, use the command `grep '. *p. *i. *n. *e. *'` *dict*, where *dict* is your local dictionary. Then check if that word is actually a longest common subsequence. Working C code for finding a longest common subsequence of two strings appears at <http://www.cs.dartmouth.edu/~thc/code/lcs.c> The comments in the code refer to the second edition of the textbook, but the code is correct.]

Examples

[The examples are of different types of trees.]

s p r i n g t i m e
 / | / | / | /
 p i o n e e r

h o r s e b a c k
 / | / | / | /
 s n o w f l a k e

m a e l s t r o m
 / | / | / | /
 b e c a l m

h e r o i c a l l y
 / | / | / | /
 s c h o l a r l y

Brute-force algorithm:

For every subsequence of X , check whether it's a subsequence of Y .

Time: $\Theta(n2^m)$.

- 2^m subsequences of X to check.
- Each subsequence takes $\Theta(n)$ time to check: scan Y for first letter, from there scan for second, and so on.

Step 1: Characterize an LCS

Notation:

$X_i = \text{prefix } \langle x_1, \dots, x_i \rangle$

$Y_i = \text{prefix } \langle y_1, \dots, y_i \rangle$

Theorem

Let $Z = \langle z_1, \dots, z_k \rangle$ be any LCS of X and Y .

1. If $x_m = y_n$, then $z_k = x_m = y_n$ and Z_{k-1} is an LCS of X_{m-1} and Y_{n-1} .
2. If $x_m \neq y_n$ and $z_k \neq x_m$, then Z is an LCS of X_{m-1} and Y .
3. If $x_m \neq y_n$ and $z_k \neq y_n$, then Z is an LCS of X and Y_{n-1} .

Proof

1. First show that $z_k = x_m = y_n$. Suppose not. Then make a subsequence $Z' = \langle z_1, \dots, z_k, x_m \rangle$. It's a common subsequence of X and Y and has length $k + 1 \Rightarrow Z'$ is a longer common subsequence than $Z \Rightarrow$ contradicts Z being an LCS.

Now show Z_{k-1} is an LCS of X_{m-1} and Y_{n-1} . Clearly, it's a common subsequence. Now suppose there exists a common subsequence W of X_{m-1} and Y_{n-1} that's longer than $Z_{k-1} \Rightarrow$ length of $W \geq k$. Make subsequence W' by appending x_m to W . W' is common subsequence of X and Y , has length $\geq k + 1 \Rightarrow$ contradicts Z being an LCS.

2. If $z_k \neq x_m$, then Z is a common subsequence of X_{m-1} and Y . Suppose there exists a subsequence W of X_{m-1} and Y with length $> k$. Then W is a common subsequence of X and $Y \Rightarrow$ contradicts Z being an LCS.
3. Symmetric to 2. ■ (theorem)

Therefore, an LCS of two sequences contains as a prefix an LCS of prefixes of the sequences.

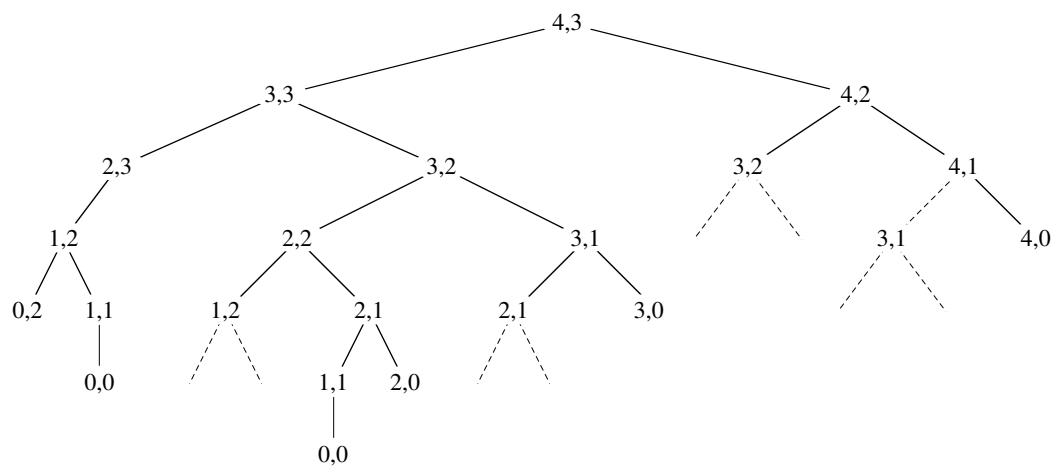
Step 2: Recursively define an optimal solution

Define $c[i, j]$ = length of LCS of X_i and Y_j . Want $c[m, n]$.

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

Again, could write a recursive algorithm based on this formulation.

Try with $X = \langle a, t, o, m \rangle$ and $Y = \langle a, n, t \rangle$. Numbers in nodes are values of i, j in each recursive call. Dashed lines indicate subproblems already computed.



- Lots of repeated subproblems.
- Instead of recomputing, store in a table.

Step 3: Compute the length of an LCS

LCS-LENGTH(X, Y, m, n)

let $b[1:m, 1:n]$ and $c[0:m, 0:n]$ be new tables

for $i = 1$ **to** m

$c[i, 0] = 0$

for $j = 0$ **to** n

$c[0, j] = 0$

for $i = 1$ **to** m // compute table entries in row-major order

for $j = 1$ **to** n

if $x_i == y_j$

$c[i, j] = c[i - 1, j - 1] + 1$

$b[i, j] = \nwarrow$

else if $c[i - 1, j] \geq c[i, j - 1]$

$c[i, j] = c[i - 1, j]$

$b[i, j] = \uparrow$

else $c[i, j] = c[i, j - 1]$

$b[i, j] = \leftarrow$

return c and b

PRINT-LCS(b, X, i, j)

if $i == 0$ or $j == 0$

return // the LCS has length 0

if $b[i, j] == \nwarrow$

 PRINT-LCS($b, X, i - 1, j - 1$)

 print x_i // same as y_j

elseif $b[i, j] == \uparrow$

 PRINT-LCS($b, X, i - 1, j$)

else PRINT-LCS($b, X, i, j - 1$)

- Initial call is PRINT-LCS(b, X, m, n).
- $b[i, j]$ points to table entry whose subproblem was used in solving LCS of X_i and Y_j .
- When $b[i, j] = \nwarrow$, LCS extended by one character. So longest common subsequence = entries with \nwarrow in them.

Demonstration

What do spanking and amputation have in common? [Show only $c[i, j]$.]

		<i>j</i>					
<i>w</i>		0	1	2	3	4	5
<i>i</i>	1	0	.25	.45	.5	.7	1.0
	2		0	.2	.25	.45	.75
	3			0	.05	.25	.55
	4				0	.2	.5
	5					0	.3
	6						0

		<i>j</i>				
<i>root</i>		1	2	3	4	5
<i>i</i>	1	1	1	1	2	2
	2		2	2	2	4
	3			3	4	5
	4				4	5
	5					5

Time

$O(n^3)$: for loops nested 3 deep, each loop index takes on $\leq n$ values. Can also show $\Omega(n^3)$. Therefore, $\Theta(n^3)$.

Step 4: Construct an optimal binary search tree

[Exercise 14.5-1 asks to write this pseudocode.]

CONSTRUCT-OPTIMAL-BST(*root*)

$r = \text{root}[1, n]$

print “ k ” _{r} “is the root”

CONSTRUCT-OPT-SUBTREE($1, r - 1, r$, “left”, *root*)

CONSTRUCT-OPT-SUBTREE($r + 1, n, r$, “right”, *root*)

CONSTRUCT-OPT-SUBTREE(*i, j, r, dir, root*)

if $i \leq j$

$t = \text{root}[i, j]$

print “ k ” _{t} “is” *dir* “child of k ” _{r}

CONSTRUCT-OPT-SUBTREE($i, t - 1, t$, “left”, *root*)

CONSTRUCT-OPT-SUBTREE($t + 1, j, t$, “right”, *root*)

Elements of dynamic programming

Mentioned already:

- optimal substructure
- overlapping subproblems

Optimal substructure

- Show that a solution to a problem consists of making a choice, which leaves one or more subproblems to solve.
- Suppose that you are given this last choice that leads to an optimal solution. *[We find that students often have trouble understanding the relationship between optimal substructure and determining which choice is made in an optimal solution. One way that helps them understand optimal substructure is to imagine that the dynamic-programming gods tell you what was the last choice made in an optimal solution.]*
- Given this choice, determine which subproblems arise and how to characterize the resulting space of subproblems.
- Show that the solutions to the subproblems used within the optimal solution must themselves be optimal. Usually use cut-and-paste:
 - Suppose that one of the subproblem solutions is not optimal.
 - *Cut* it out.
 - *Paste* in an optimal solution.
 - Get a better solution to the original problem. Contradicts optimality of problem solution.

That was optimal substructure.

Need to ensure that you consider a wide enough range of choices and subproblems that you get them all. *[The dynamic-programming gods are too busy to tell you what that last choice really was.]* Try all the choices, solve all the subproblems resulting from each choice, and pick the choice whose solution, along with subproblem solutions, is best.

How to characterize the space of subproblems?

- Keep the space as simple as possible.
- Expand it as necessary.

Examples

Rod cutting

- Space of subproblems was rods of length $n - i$, for $1 \leq i \leq n$.
- No need to try a more general space of subproblems.

Matrix-chain multiplication

- Suppose we had tried to constrain the space of subproblems to parenthesizing $A_1 A_2 \cdots A_j$.
- An optimal parenthesization splits at some matrix A_k .
- Get subproblems for $A_1 \cdots A_k$ and $A_{k+1} \cdots A_j$.
- Unless we could guarantee that $k = j - 1$, so that the subproblem for $A_{k+1} \cdots A_j$ has only A_j , then this subproblem is *not* of the form $A_1 A_2 \cdots A_j$.
- Thus, needed to allow the subproblems to vary at both ends—allow both i and j to vary.

Longest common subsequence

- Space of subproblems for $\langle x_1, \dots, x_i \rangle$ and $\langle y_1, \dots, y_j \rangle$ was just $\langle x_1, \dots, x_{i-1} \rangle$ and $\langle y_1, \dots, y_{j-1} \rangle$.
- No need to try a more general space of subproblems.

Optimal binary search trees

- Similar to matrix-chain multiplication.
- Suppose we had tried to constrain space of subproblems to subtrees with keys k_1, k_2, \dots, k_j .
- An optimal BST would have root k_r , for some $1 \leq r \leq j$.
- Get subproblems k_1, \dots, k_{r-1} and k_{r+1}, \dots, k_j .
- Unless we could guarantee that $r = j$, so that subproblem with k_{r+1}, \dots, k_j is empty, then this subproblem is *not* of the form k_1, k_2, \dots, k_j .
- Thus, needed to allow the subproblems to vary at “both ends,” i.e., allow both i and j to vary.

Optimal substructure varies across problem domains:

1. *How many subproblems* are used in an optimal solution.
 2. *How many choices* in determining which subproblem(s) to use.
- Rod cutting:
 - 1 subproblem (of size $n - i$)
 - n choices
 - Matrix-chain multiplication:
 - 2 subproblems ($A_i \cdots A_k$ and $A_{k+1} \cdots A_j$)
 - $j - i$ choices for A_k in $A_i, A_{i+1}, \dots, A_{j-1}$. Having found optimal solutions to subproblems, choose from among the $j - i$ candidates for A_k .
 - Longest common subsequence:
 - 1 subproblem
 - Either
 - 1 choice (if $x_i = y_j$, LCS of X_{i-1} and Y_{j-1}), or
 - 2 choices (if $x_i \neq y_j$, LCS of X_{i-1} and Y , and LCS of X and Y_{j-1})
 - Optimal binary search tree:
 - 2 subproblems (k_i, \dots, k_{r-1} and k_{r+1}, \dots, k_j)
 - $j - i + 1$ choices for k_r in k_i, \dots, k_j . Having found optimal solutions to subproblems, choose from among the $j - i + 1$ candidates for k_r .

Informally, running time depends on (# of subproblems overall) \times (# of choices).

- Rod cutting: $\Theta(n)$ subproblems, $\leq n$ choices for each
 $\Rightarrow O(n^2)$ running time.
- Matrix-chain multiplication: $\Theta(n^2)$ subproblems, $O(n)$ choices for each
 $\Rightarrow O(n^3)$ running time.

- Longest common subsequence: $\Theta(mn)$ subproblems, ≤ 2 choices for each $\Rightarrow \Theta(mn)$ running time.
- Optimal binary search tree: $\Theta(n^2)$ subproblems, $O(n)$ choices for each $\Rightarrow O(n^3)$ running time.

Can use the subproblem graph to get the same analysis: count the number of edges.

- Each vertex corresponds to a subproblem.
- Choices for a subproblem are vertices that the subproblem has edges going to.
- For rod cutting, subproblem graph has n vertices and $\leq n$ edges per vertex $\Rightarrow O(n^2)$ running time.
In fact, can get an exact count of the edges: for $i = 0, 1, \dots, n$, vertex for subproblem size i has out-degree $i \Rightarrow \# \text{ of edges} = \sum_{i=0}^n i = n(n+1)/2$.
- Subproblem graph for matrix-chain multiplication has $\Theta(n^2)$ vertices, each with degree $\leq n-1 \Rightarrow O(n^3)$ running time.

Dynamic programming uses optimal substructure *bottom up*.

- *First* find optimal solutions to subproblems.
- *Then* choose which to use in optimal solution to the problem.

When we look at greedy algorithms, we'll see that they work *top down*: *first* make a choice that looks best, *then* solve the resulting subproblem.

Don't be fooled into thinking optimal substructure applies to all optimization problems. It doesn't.

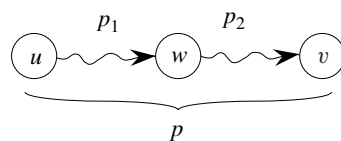
Here are two problems that look similar. In both, we're given an *unweighted, directed graph* $G = (V, E)$.

- V is a set of *vertices*.
- E is a set of *edges*.

And we ask about finding a **path** (sequence of connected edges) from vertex u to vertex v .

- **Shortest path**: find a path $u \rightsquigarrow v$ with fewest edges. Must be **simple** (no *cycles*), since removing a cycle from a path gives a path with fewer edges.
- **Longest simple path**: find a *simple* path $u \rightsquigarrow v$ with most edges. If didn't require simple, could repeatedly traverse a cycle to make an arbitrarily long path.

Shortest path has optimal substructure.



- Suppose p is shortest path $u \rightsquigarrow v$.
- Let w be any vertex on p .
- Let p_1 be the portion of p going $u \rightsquigarrow w$.
- Then p_1 is a shortest path $u \rightsquigarrow w$.

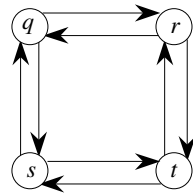
Proof Suppose there exists a shorter path p'_1 going $u \rightsquigarrow w$. Cut out p_1 , replace it with p'_1 , get path $u \xrightarrow{p'_1} w \xrightarrow{p_2} v$ with fewer edges than p . ■

Therefore, can find shortest path $u \rightsquigarrow v$ by considering all intermediate vertices w , then finding shortest paths $u \rightsquigarrow w$ and $w \rightsquigarrow v$.

Same argument applies to p_2 .

Does longest path have optimal substructure?

- It seems like it should.
- It does *not*.



Consider $q \rightarrow r \rightarrow t =$ longest path $q \rightsquigarrow t$. Are its subpaths longest paths?
No!

- Subpath $q \rightsquigarrow r$ is $q \rightarrow r$.
- Longest simple path $q \rightsquigarrow r$ is $q \rightarrow s \rightarrow t \rightarrow r$.
- Subpath $r \rightsquigarrow t$ is $r \rightarrow t$.
- Longest simple path $r \rightsquigarrow t$ is $r \rightarrow q \rightarrow s \rightarrow t$.

Not only isn't there optimal substructure, but can't even assemble a legal solution from solutions to subproblems.

Combine longest simple paths:

$q \rightarrow s \rightarrow t \rightarrow r \rightarrow q \rightarrow s \rightarrow t$

Not simple!

In fact, this problem is NP-complete (so it probably has no optimal substructure to find.)

What's the big difference between shortest path and longest path?

- Shortest path has **independent** subproblems.
- Solution to one subproblem does not affect solution to another subproblem of the same problem.
- Longest simple path: subproblems are *not* independent.
- Consider subproblems of longest simple paths $q \rightsquigarrow r$ and $r \rightsquigarrow t$.
- Longest simple path $q \rightsquigarrow r$ uses s and t .
- Cannot use s and t to solve longest simple path $r \rightsquigarrow t$, since if you do, the path isn't simple.
- But you *have* to use t to find longest simple path $r \rightsquigarrow t$!

- Using resources (vertices) to solve one subproblem renders them unavailable to solve the other subproblem.

[For shortest paths, for a shortest path $u \xrightarrow{p_1} w \xrightarrow{p_2} v$, no vertex other than w can appear in p_1 and p_2 . Otherwise, get a cycle.]

Independent subproblems in our examples:

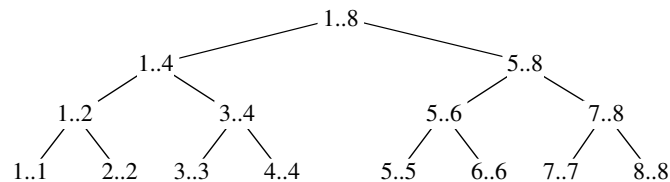
- Rod cutting and longest common subsequence
 - 1 subproblem \Rightarrow automatically independent.
- Matrix-chain multiplication
 - $A_i \cdots A_k$ and $A_{k+1} \cdots A_j \Rightarrow$ independent.
- Optimal binary search tree
 - k_i, \dots, k_{r-1} and $k_{r+1}, \dots, k_j \Rightarrow$ independent.

Overlapping subproblems

These occur when a recursive algorithm revisits the same problem over and over.

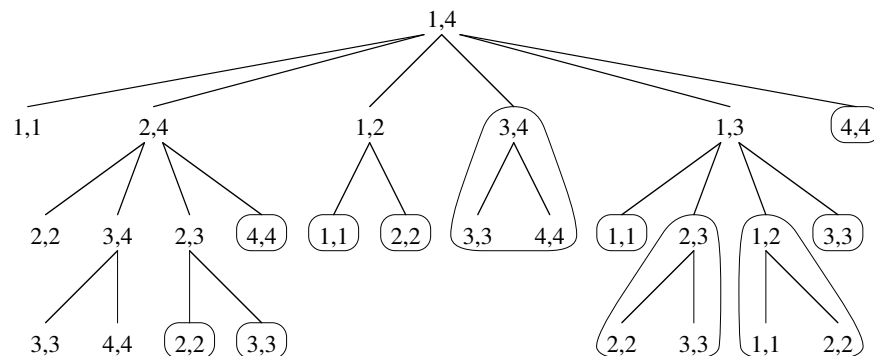
Good divide-and-conquer algorithms usually generate a brand new problem at each stage of recursion.

Example: merge sort



Alternative approach to dynamic programming: *memoization*

- “Store, don’t recompute.”
- Make a table indexed by subproblem.
- When solving a subproblem:
 - Lookup in table.
 - If answer is there, use it.
 - Else, compute answer, then store it.
- For matrix-chain multiplication:



Each node has the parameters i and j . Computations performed in highlighted subtrees are replaced by a single table lookup if computing recursively with memoization.

- In bottom-up dynamic programming, we go one step further. Determine in what order to access the table, and fill it in that way.

Chapter 15: Greedy algorithms

Reading: 15.1–15.3

Lecture Notes for Chapter 15:

Greedy Algorithms

[The fourth edition removed the starred sections on matroids and task scheduling (an application of matroids). These sections were replaced by a new, unstarred section covering offline caching, which had been the subject of Problem 16-5 in the third edition.]

Chapter 15 overview

Similar to dynamic programming.

Used for optimization problems.

Idea

When you have a choice to make, make the one that looks best *right now*. Make a *locally optimal choice* in hope of getting a *globally optimal solution*.

Greedy algorithms don't always yield an optimal solution. But sometimes they do. We'll see a problem for which they do. Then we'll look at some general characteristics of when greedy algorithms give optimal solutions. We then study two other applications of the greedy method: Huffman coding and offline caching. *[Later chapters use the greedy method as well: minimum spanning tree, Dijkstra's algorithm for single-source shortest paths, and a greedy set-covering heuristic.]*

Activity selection

n **activities** require *exclusive* use of a common resource. For example, scheduling the use of a classroom.

Set of activities $S = \{a_1, \dots, a_n\}$.

a_i needs resource during period $[s_i, f_i)$, which is a half-open interval, where s_i = start time and f_i = finish time.

Goal

Select the largest possible set of nonoverlapping (***mutually compatible***) activities.

Could have many other objectives:

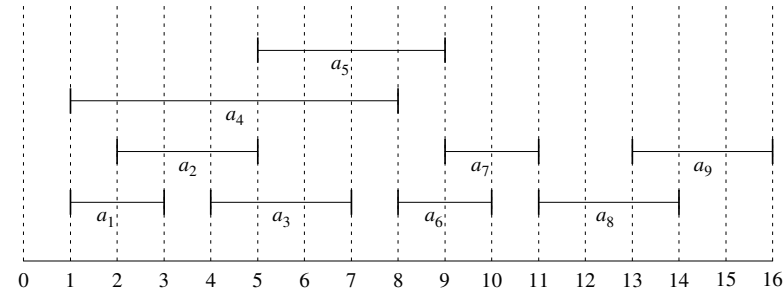
- Schedule room for longest time.
- Maximize income rental fees.

Assume that activities are sorted by finish time: $f_1 \leq f_2 \leq f_3 \leq \dots \leq f_{n-1} \leq f_n$.

Example

S sorted by finish time: [Leave on board]

i	1	2	3	4	5	6	7	8	9
s_i	1	2	4	1	5	8	9	11	13
f_i	3	5	7	8	9	10	11	14	16



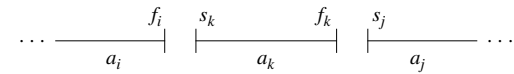
Maximum-size mutually compatible set: $\{a_1, a_3, a_6, a_8\}$.

Not unique: also $\{a_1, a_3, a_6, a_9\}$, $\{a_1, a_3, a_7, a_8\}$, $\{a_1, a_3, a_7, a_9\}$, $\{a_1, a_5, a_7, a_8\}$, $\{a_1, a_5, a_7, a_9\}$, $\{a_2, a_5, a_7, a_8\}$, $\{a_2, a_5, a_7, a_9\}$.

Optimal substructure of activity selection

$$S_{ij} = \{a_k \in S : f_i \leq s_k < f_k \leq s_j\} \quad [\text{Leave on board}]$$

= activities that start after a_i finishes and finish before a_j starts.



Activities in S_{ij} are compatible with

- all activities that finish by f_i , and
- all activities that start no earlier than s_j .

Let A_{ij} be a maximum-size set of mutually compatible activities in S_{ij} .

Let $a_k \in A_{ij}$ be some activity in A_{ij} . Then we have two subproblems:

- Find mutually compatible activities in S_{ik} (activities that start after a_i finishes and that finish before a_k starts).
- Find mutually compatible activities in S_{kj} (activities that start after a_k finishes and that finish before a_j starts).

Let

$A_{ik} = A_{ij} \cap S_{ik}$ = activities in A_{ij} that finish before a_k starts,

$A_{kj} = A_{ij} \cap S_{kj}$ = activities in A_{ij} that start after a_k finishes.

Then $A_{ij} = A_{ik} \cup \{a_k\} \cup A_{kj}$

$$\Rightarrow |A_{ij}| = |A_{ik}| + |A_{kj}| + 1.$$

Claim

Optimal solution A_{ij} must include optimal solutions for the two subproblems for S_{ik} and S_{kj} .

Proof of claim Use the usual cut-and-paste argument. Will show the claim for S_{kj} ; proof for S_{ik} is symmetric.

Suppose we could find a set A'_{kj} of mutually compatible activities in S_{kj} , where $|A'_{kj}| > |A_{kj}|$. Then use A'_{kj} instead of A_{kj} when solving the subproblem for S_{ij} . Size of resulting set of mutually compatible activities would be $|A_{ik}| + |A'_{kj}| + 1 > |A_{ik}| + |A_{kj}| + 1 = |A|$. Contradicts assumption that A_{ij} is optimal. ■ (claim)

One recursive solution

Since optimal solution A_{ij} must include optimal solutions to the subproblems for S_{ik} and S_{kj} , could solve by dynamic programming.

Let $c[i, j]$ = size of optimal solution for S_{ij} . Then

$$c[i, j] = c[i, k] + c[k, j] + 1 .$$

But we don't know which activity a_k to choose, so we have to try them all:

$$c[i, j] = \begin{cases} 0 & \text{if } S_{ij} = \emptyset , \\ \max \{c[i, k] + c[k, j] + 1 : a_k \in S_{ij}\} & \text{if } S_{ij} \neq \emptyset . \end{cases}$$

Could then develop a recursive algorithm and memoize it. Or could develop a bottom-up algorithm and fill in table entries.

Instead, we will look at a greedy approach.

Making the greedy choice

Choose an activity to add to optimal solution *before* solving subproblems. For activity-selection problem, we can get away with considering only the greedy choice: the activity that leaves the resource available for as many other activities as possible.

Question: Which activity leaves the resource available for the most other activities?

Answer: The first activity to finish. (If more than one activity has earliest finish time, can choose any such activity.)

Since activities are sorted by finish time, just choose activity a_1 .

That leaves only one subproblem to solve: finding a maximum size set of mutually compatible activities that start after a_1 finishes. (Don't have to worry about activities that finish before a_1 starts, because $s_1 < f_1$ and no activity a_i has finish time $f_i < f_1 \Rightarrow$ no activity a_i has $f_i \leq s_1$.)

Since have only subproblem to solve, simplify notation:

$$S_k = \{a_i \in S : s_i \geq f_k\} = \text{activities that start after } a_k \text{ finishes} .$$

Making greedy choice of $a_1 \Rightarrow S_1$ remains as only subproblem to solve. [Slight abuse of notation: referring to S_k not only as a set of activities but as a subproblem consisting of these activities.]

By optimal substructure, if a_1 is in an optimal solution, then an optimal solution to the original problem consists of a_1 plus all activities in an optimal solution to S_1 .

But need to prove that a_1 is always part of some optimal solution.

Theorem

If S_k is nonempty and a_m has the earliest finish time in S_k , then a_m is included in some optimal solution.

Proof Let A_k be an optimal solution to S_k , and let a_j have the earliest finish time of any activity in A_k . If $a_j = a_m$, done. Otherwise, let $A'_k = A_k - \{a_j\} \cup \{a_m\}$ be A_k but with a_m substituted for a_j .

Claim

Activities in A'_k are disjoint.

Proof of claim Activities in A_k are disjoint, a_j is first activity in A_k to finish, and $f_m \leq f_j$. ■ (claim)

Since $|A'_k| = |A_k|$, conclude that A'_k is an optimal solution to S_k , and it includes a_m . ■ (theorem)

So, don't need full power of dynamic programming. Don't need to work bottom-up.

Instead, can just repeatedly choose the activity that finishes first, keep only the activities that are compatible with that one, and repeat until no activities remain.

Can work top-down: make a choice, then solve a subproblem. Don't have to solve subproblems before making a choice.

Recursive greedy algorithm

Start and finish times are represented by arrays s and f , where f is assumed to be already sorted in monotonically increasing order.

To start, add fictitious activity a_0 with $f_0 = 0$, so that $S_0 = S$, the entire set of activities.

Procedure RECURSIVE-ACTIVITY-SELECTOR takes as parameters the arrays s and f , index k of current subproblem, and number n of activities in the original problem.

RECURSIVE-ACTIVITY-SELECTOR(s, f, k, n)

$m = k + 1$

while $m \leq n$ and $s[m] < f[k]$ // find the first activity in S_k to finish

$m = m + 1$

if $m \leq n$

return $\{a_m\} \cup \text{RECURSIVE-ACTIVITY-SELECTOR}(s, f, m, n)$

else return \emptyset

Initial call

RECURSIVE-ACTIVITY-SELECTOR($s, f, 0, n$).

Idea

The **while** loop checks $a_{k+1}, a_{k+2}, \dots, a_n$ until it finds an activity a_m that is compatible with a_k (need $s_m \geq f_k$).

- If the loop terminates because a_m is found ($m \leq n$), then recursively solve S_m , and return this solution, along with a_m .
- If the loop never finds a compatible a_m ($m > n$), then just return empty set.

Go through example given earlier. Should get $\{a_1, a_3, a_6, a_8\}$.

Time

$\Theta(n)$ —each activity examined exactly once, assuming that activities are already sorted by finish times.

Iterative greedy algorithm

Can convert the recursive algorithm to an iterative one. It's already almost tail recursive.

GREEDY-ACTIVITY-SELECTOR(s, f, n)

$A = \{a_1\}$

$k = 1$

for $m = 2$ **to** n

if $s[m] \geq f[k]$ // is a_m in S_k ?

$A = A \cup \{a_m\}$ // yes, so choose it

$k = m$ // and continue from there

return A

Go through example given earlier. Should again get $\{a_1, a_3, a_6, a_8\}$.

Time

$\Theta(n)$, if activities are already sorted by finish times.

For both the recursive and iterative algorithms, add $O(n \lg n)$ time if activities need to be sorted.

Elements of the greedy strategy

The choice that seems best at the moment is the one we go with.

What did we do for activity selection?

1. Determine the optimal substructure.

2. Develop a recursive solution.
3. Show that if you make the greedy choice, only one subproblem remains.
4. Prove that it's always safe to make the greedy choice.
5. Develop a recursive greedy algorithm.
6. Convert it to an iterative algorithm.

At first, it looked like dynamic programming. In the activity-selection problem, we started out by defining subproblems S_{ij} , where both i and j varied. But then found that making the greedy choice allowed us to restrict the subproblems to be of the form S_k .

Could instead have gone straight for the greedy approach: in our first crack at defining subproblems, use the S_k form. Could then have proven that the greedy choice a_m (the first activity to finish), combined with optimal solution to the remaining compatible activities S_m , gives an optimal solution to S_k .

Typically, we streamline these steps:

1. Cast the optimization problem as one in which we make a choice and are left with one subproblem to solve.
2. Prove that there's always an optimal solution that makes the greedy choice, so that the greedy choice is always safe.
3. Demonstrate optimal substructure by showing that, having made the greedy choice, combining an optimal solution to the remaining subproblem with the greedy choice gives an optimal solution to the original problem.

No general way to tell whether a greedy algorithm is optimal, but two key ingredients are

1. greedy-choice property and
2. optimal substructure.

Greedy-choice property

Can assemble a globally optimal solution by making locally optimal (greedy) choices.

Dynamic programming

- Make a choice at each step.
- Choice depends on knowing optimal solutions to subproblems. Solve subproblems *first*.
- Solve *bottom-up* (unless memoizing).

Greedy

- Make a choice at each step.
- Make the choice *before* solving the subproblems.
- Solve *top-down*.

Typically show the greedy-choice property by what we did for activity selection:

- Look at an optimal solution.
- If it includes the greedy choice, done.
- Otherwise, modify the optimal solution to include the greedy choice, yielding another solution that's just as good.

Can get efficiency gains from greedy-choice property.

- Preprocess input to put it into greedy order.
- Or, if dynamic data, use a priority queue.

Optimal substructure

Just show that optimal solution to subproblem and greedy choice \Rightarrow optimal solution to problem.

Greedy vs. dynamic programming

The knapsack problem is a good example of the difference.

0-1 knapsack problem

- n items.
- Item i is worth v_i , weighs w_i pounds.
- Find a most valuable subset of items with total weight $\leq W$.
- Have to either take an item or not take it—can't take part of it.

Fractional knapsack problem

Like the 0-1 knapsack problem, but can take fraction of an item.

Both have optimal substructure.

But the fractional knapsack problem has the greedy-choice property, and the 0-1 knapsack problem does not.

To solve the fractional problem, rank items by value/weight: v_i/w_i . Let $v_i/w_i \geq v_{i+1}/w_{i+1}$ for all i . Take items in decreasing order of value/weight. Will take all of the items with the greatest value/weight, and possibly a fraction of the next item.

FRACTIONAL-KNAPSACK(v, w, W)

$load = 0$

$i = 1$

while $load < W$ and $i \leq n$

if $w_i \leq W - load$

 take all of item i

else take $(W - load)/w_i$ of item i

 add what was taken to $load$

$i = i + 1$

Time: $O(n \lg n)$ to sort, $O(n)$ thereafter.

Greedy doesn't work for the 0-1 knapsack problem. Might get empty space, which lowers the average value per pound of the items taken.

i	1	2	3
v_i	60	100	120
w_i	10	20	30
v_i/w_i	6	5	4

$W = 50$.

Greedy solution:

- Take items 1 and 2.
- value = 160, weight = 30.

Have 20 pounds of capacity left over.

Optimal solution:

- Take items 2 and 3.
- value = 220, weight = 50.

No leftover capacity.

Huffman codes

Goal: Compress a data file made up of characters. You know how often each character appears in the file—its *frequency*. Each character is represented by some bit sequence: a *codeword*. Use as few bits as possible to represent the file.

Fixed-length code: All codewords have the same number of bits. For $n \geq 2$ characters, need $\lceil \lg n \rceil$ bits.

Variable-length code: Represent different characters with differing numbers of bits. In particular, give frequently occurring characters shorter codewords and infrequently occurring characters longer codewords.

Example: For a data file of 100,000 characters:

	a	b	c	d	e	f
Frequency (in thousands)	45	13	12	16	9	5
Fixed-length codeword	000	001	010	011	100	101
Variable-length codeword	0	101	100	111	1101	1100

For a fixed-length code, need 3 bits per character. For 100,000 characters, need 300,000 bits. For this variable-length code, need

$$\begin{array}{rcl}
 45,000 \cdot 1 & = & 45,000 \\
 + 13,000 \cdot 3 & = & 39,000 \\
 + 12,000 \cdot 3 & = & 36,000 \\
 + 16,000 \cdot 3 & = & 48,000 \\
 + 9,000 \cdot 4 & = & 36,000 \\
 + 5,000 \cdot 4 & = & 20,000 \\
 \hline
 & = & 224,000 \text{ bits}
 \end{array}$$

Prefix-free codes

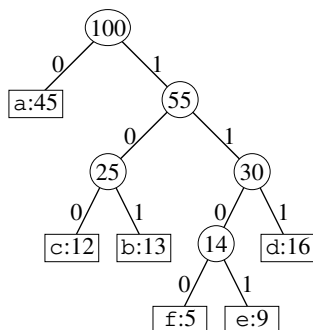
No codeword is also a prefix of any other codeword. [Called “prefix codes” in earlier editions of the book. Changed to “prefix-free codes” in the fourth edition because each codeword is free of prefixes of other codes.] A prefix-free code can always achieve the optimal compression.

Encoding: Just concatenate codewords for each character in the file. **Example:** To encode *face*: $1100 \cdot 0 \cdot 100 \cdot 1101 = 110001001101$, where \cdot is concatenation.

Decoding: Since no codeword is a prefix of any other codeword, just process bits until you get a match. Then discard the bits and go from the rest of the compressed file. **Example:** If encoding is 100011001101 , get a match on $100 = c$. That leaves 011001101 . Get a match on $0 = a$. That leaves 11001101 . Get a match on $1100 = f$. That leaves 1101 . Get a match on $1101 = e$. So the encoded file represents *cafe*.

Binary tree representation

Use a binary tree whose leaves are the characters. The codeword for a character is given by the simple path from the root down to that character’s leaf, where going left is 0 and going right is 1.



Here, each leaf has its character and frequency (in thousands). Each internal node holds the sum of the frequencies of the leaves in its subtree.

An optimal code is always given by a full binary tree: each internal node has 2 children \Rightarrow if C is the alphabet for the characters, then the tree has $|C|$ leaves and $|C| - 1$ internal nodes.

How to compute the number of bits to encode a file for alphabet C given tree T :

For each character $c \in C$, denote its frequency by $c.freq$. Denote the depth of c in T by $d_T(c)$, which equals the length of c ’s codeword. Then the number of bits to encode the file, the *cost* of T , is

$$B(T) = \sum_{c \in C} c.freq \cdot d_T(c) .$$

Constructing a Huffman code

[Named after David Huffman.] The algorithm builds tree T bottom-up. It repeatedly selects two nodes with the lowest frequency and makes them children of

a new node whose frequency is the sum of the two nodes' frequencies. It uses a min-priority queue Q keyed on the *freq* attribute, which all nodes have.

HUFFMAN(C)

$n = |C|$

$Q = C$

for $i = 1$ **to** $n - 1$

 allocate a new node z

$x = \text{EXTRACT-MIN}(Q)$

$y = \text{EXTRACT-MIN}(Q)$

$z.\text{left} = x$

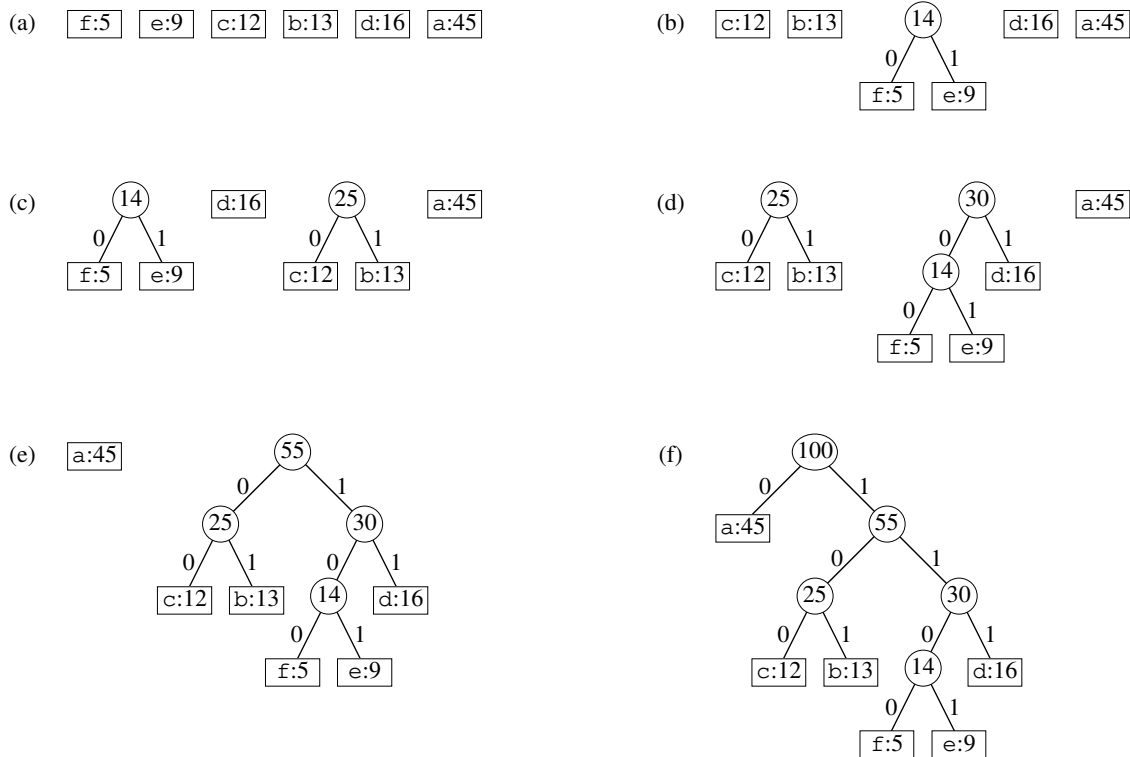
$z.\text{right} = y$

$z.\text{freq} = x.\text{freq} + y.\text{freq}$

$\text{INSERT}(Q, z)$

return $\text{EXTRACT-MIN}(Q)$ // the root of the tree is the only node left

Example: Using the frequencies from before:



Running time: Let $n = |C|$. The running time depends on how the min-priority queue Q is implemented. If with a binary min-heap, can initialize Q in $O(n)$ time. The **for** loop runs $n - 1$ times, and each INSERT and EXTRACT-MIN call takes $O(\lg n)$ time $\Rightarrow O(n \lg n)$ time in all.

Correctness

Show the greedy-choice and optimal-substructure properties.

Lemma (Greedy-choice property)

For alphabet C , let x and y be the two characters with the lowest frequencies. Then there exists an optimal prefix-free code for C where the codewords for x and y have the same length and differ only in the last bit.

Proof Given a tree T for some optimal prefix-free code, modify it so that x and y are sibling leaves of maximum depth. Then the codewords for x and y will have the same length and differ in the last bit.

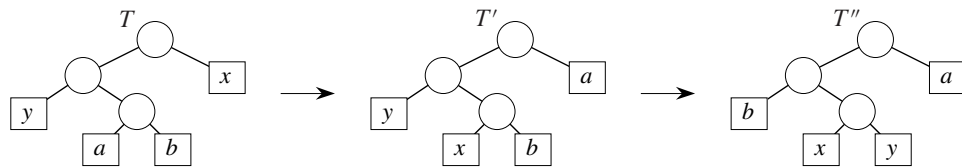
Let a, b be two characters that are sibling leaves of maximum depth in T . Assume wlog that $a.\text{freq} \leq b.\text{freq}$ and $x.\text{freq} \leq y.\text{freq}$. Must have $x.\text{freq} \leq a.\text{freq}$ and $y.\text{freq} \leq b.\text{freq}$.

Could have $x.\text{freq} = a.\text{freq}$ or $y.\text{freq} = b.\text{freq}$. If $x.\text{freq} = b.\text{freq}$, then $a.\text{freq} = b.\text{freq} = x.\text{freq} = y.\text{freq}$ (Exercise 15.3-1), and the lemma is trivially true. So assume that $x.\text{freq} \neq b.\text{freq} \Rightarrow x \neq b$.

In T : exchange a and x , producing T' .

In T' : exchange b and y , producing T'' .

In T'' , x and y are sibling leaves of maximum depth.

**Claim**

$B(T') \leq B(T)$. (Exchanging a and x does not increase the cost.)

Proof of claim

$$\begin{aligned}
 B(T) - B(T') &= \sum_{c \in C} c.\text{freq} \cdot d_T(c) - \sum_{c \in C} c.\text{freq} \cdot d_{T'}(c) \\
 &= x.\text{freq} \cdot d_T(x) + a.\text{freq} \cdot d_T(a) - x.\text{freq} \cdot d_{T'}(x) - a.\text{freq} \cdot d_{T'}(a) \\
 &= x.\text{freq} \cdot d_T(x) + a.\text{freq} \cdot d_T(a) - x.\text{freq} \cdot d_T(a) - a.\text{freq} \cdot d_T(x) \\
 &= (a.\text{freq} - x.\text{freq})(d_T(a) - d_T(x)) \\
 &\geq 0.
 \end{aligned}$$

The last line follows because $x.\text{freq} \leq a.\text{freq}$ and a is a maximum-depth leaf $\Rightarrow d_T(a) \geq d_T(x)$. ■ (claim)

Similarly, $B(T'') \leq B(T')$ because exchanging y and b doesn't increase the cost. Therefore, $B(T'') \leq B(T') \leq B(T)$. T is optimal $\Rightarrow B(T) \leq B(T'') \Rightarrow B(T'') = B(T) \Rightarrow T''$ is optimal, and x and y are sibling leaves of maximum depth. ■

The lemma shows that to build up an optimal tree, can begin with the greedy choice of merging the two characters with lowest frequency. Greedy because the cost of a merger is the sum of the frequencies of its children and the cost of a tree equals the sum of the costs of its mergers (Exercise 15.3-4).

Lemma (Optimal-substructure property)

For alphabet C , let x, y be the two characters with minimum frequency. Let $C' = (C - \{x, y\}) \cup z$ for a new character z with $z.freq = x.freq + y.freq$. Let T' be a tree representing an optimal prefix-free code for C' , and T be T' with the leaf for z replaced by an internal node with children x and y . Then T represents an optimal prefix-free code for C .

Proof $c \in C - \{x, y\} \Rightarrow d_T(c) = d_{T'}(c) \Rightarrow c.freq \cdot d_T(c) = c.freq \cdot d_{T'}(c)$.
 $d_T(x) = d_T(y) = d_{T'}(z) + 1 \Rightarrow$

$$\begin{aligned} x.freq \cdot d_T(x) + y.freq \cdot d_T(y) &= (x.freq + y.freq)(d_{T'}(z) + 1) \\ &= z.freq \cdot d_{T'}(z) + (x.freq + y.freq), \end{aligned}$$

so that $B(T) = B(T') + x.freq + y.freq$, which is equivalent to $B(T') = B(T) - x.freq - y.freq$.

Now suppose T doesn't represent an optimal prefix-free code for C . Then $B(T'') < B(T)$ for some optimal tree T'' . By the previous lemma, without loss of generality, T'' has x and y as siblings. Replace the common parent of x and y by a leaf z with $z.freq = x.freq + y.freq$ and call the resulting tree T''' . Then,

$$\begin{aligned} B(T''') &= B(T'') - x.freq - y.freq \\ &< B(T) - x.freq - y.freq \\ &= B(T'), \end{aligned}$$

so that T' was not optimal, a contradiction. ■

Theorem

HUFFMAN produces an optimal prefix-free code.

Proof The greedy-choice and optimal-substructure properties both apply. ■

Offline caching

In a computer, a **cache** is memory that is smaller but faster than main memory. It holds a small subset of what's in main memory. Caches store data in **blocks**, also known as **cache lines**, usually 32, 64, or 128 bytes. [We use the term blocks in this discussion, rather than cache lines.]

A program makes a sequence of memory requests to blocks. Each block usually has several requests to some data that it holds.

The cache size is limited to k blocks, starting out empty before the first request. Each request causes either 0 or 1 block to enter the cache, and either 0 or 1 block to be evicted. A request for block b may have one of three outcomes:

1. b is already in the cache due to some previous request \Rightarrow **cache hit**. The cache remains unchanged.
2. b is not already in the cache, but the cache is not yet full (contains $< k$ blocks). b goes into the cache, so that the cache now contains one more block than before the request.

Chapter 20: Basic graph algorithms

Reading: Chapter 20

Lecture Notes for Chapter 20: Elementary Graph Algorithms

Graph representation

Given graph $G = (V, E)$. In pseudocode, represent vertex set by $G.V$ and edge set by $G.E$.

- G may be either directed or undirected.
- Two common ways to represent graphs for algorithms:
 1. Adjacency lists.
 2. Adjacency matrix.

When expressing the running time of an algorithm, it's often in terms of both $|V|$ and $|E|$. In asymptotic notation—and *only* in asymptotic notation—we'll drop the cardinality. Example: $O(V + E)$ really means $O(|V| + |E|)$.

[The introduction to Part VI talks more about this.]

Adjacency lists

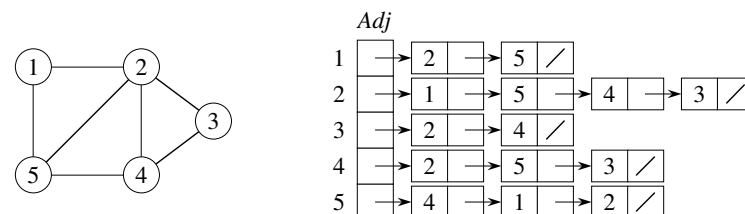
Array Adj of $|V|$ lists, one per vertex.

Vertex u 's list has all vertices v such that $(u, v) \in E$. (Works for both directed and undirected graphs.)

In pseudocode, denote the array as attribute $G.Adj$, so will see notation such as $G.Adj[u]$.

Example

For an undirected graph:



If edges have *weights*, can put the weights in the lists.

Weight: $w : E \rightarrow \mathbb{R}$

We'll use weights later on for spanning trees and shortest paths.

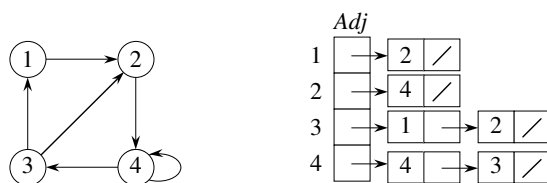
Space: $\Theta(V + E)$.

Time: to list all vertices adjacent to u : $\Theta(\text{degree}(u))$.

Time: to determine whether $(u, v) \in E$: $O(\text{degree}(u))$.

Example

For a directed graph:



Same asymptotic space and time.

Adjacency matrix

$|V| \times |V|$ matrix $A = (a_{ij})$

$$a_{ij} = \begin{cases} 1 & \text{if } (i, j) \in E, \\ 0 & \text{otherwise.} \end{cases}$$

	1	2	3	4	5
1	0	1	0	0	1
2	1	0	1	1	1
3	0	1	0	1	0
4	0	1	1	0	1
5	1	1	0	1	0

	1	2	3	4
1	0	1	0	0
2	0	0	0	1
3	1	1	0	0
4	0	0	1	1

Space: $\Theta(V^2)$.

Time: to list all vertices adjacent to u : $\Theta(V)$.

Time: to determine whether $(u, v) \in E$: $\Theta(1)$.

Can store weights instead of bits for weighted graph.

We'll use both representations in these lecture notes.

Representing graph attributes

Graph algorithms usually need to maintain attributes for vertices and/or edges. Use the usual dot-notation: denote attribute d of vertex v by $v.d$.

Use the dot-notation for edges, too: denote attribute f of edge (u, v) by $(u, v).f$.

Implementing graph attributes

No one best way to implement. Depends on the programming language, the algorithm, and how the rest of the program interacts with the graph.

If representing the graph with adjacency lists, can represent vertex attributes in additional arrays that parallel the *Adj* array, e.g., $d[1 : |V|]$, so that if vertices adjacent to u are in $Adj[u]$, store $u.d$ in array entry $d[u]$.

But can represent attributes in other ways. Example: represent vertex attributes as instance variables within a subclass of a `Vertex` class.

Breadth-first search

Input: Graph $G = (V, E)$, either directed or undirected, and **source vertex** $s \in V$.

Output:

- $v.d$ = distance (smallest # of edges) from s to v , for all $v \in V$.
- $v.\pi$ is v 's **predecessor** on a shortest path (smallest # of edges) from s .
 (u, v) is last edge on shortest path $s \rightsquigarrow v$.
Predecessor subgraph contains edges (u, v) such that $v.\pi = u$.
 The predecessor subgraph forms a tree, called the **breadth-first tree**.

Later, we'll see a generalization of breadth-first search, with edge weights. For now, we'll keep it simple.

[Omitting colors of vertices. Used in book to reason about the algorithm.]

Intuition

Breadth-first search expands the frontier between discovered and undiscovered vertices uniformly across the breadth of the frontier.

Discovers vertices in waves, starting from s .

- First visits all vertices 1 edge from s .
- From there, visits all vertices 2 edges from s .
- Etc.

Use FIFO queue Q to maintain wavefront.

- $v \in Q$ if and only if wave has visited v but has not come out of v yet.
- Q contains vertices at a distance k , and possibly some vertices at a distance $k + 1$. Therefore, at any time Q contains portions of two consecutive waves.

```

BFS( $V, E, s$ )
  for each vertex  $u \in V - \{s\}$ 
     $u.d = \infty$ 
     $u.\pi = \text{NIL}$ 
   $s.d = 0$ 
   $Q = \emptyset$ 
  ENQUEUE( $Q, s$ )
  while  $Q \neq \emptyset$ 
     $u = \text{DEQUEUE}(Q)$ 
    for each vertex  $v$  in  $G.\text{Adj}[u]$  // search the neighbors of  $u$ 
      if  $v.d == \infty$  // is  $v$  being discovered now?
         $v.d = u.d + 1$ 
         $v.\pi = u$ 
        ENQUEUE( $Q, v$ ) //  $v$  is now on the frontier
    //  $u$  is now behind the frontier.

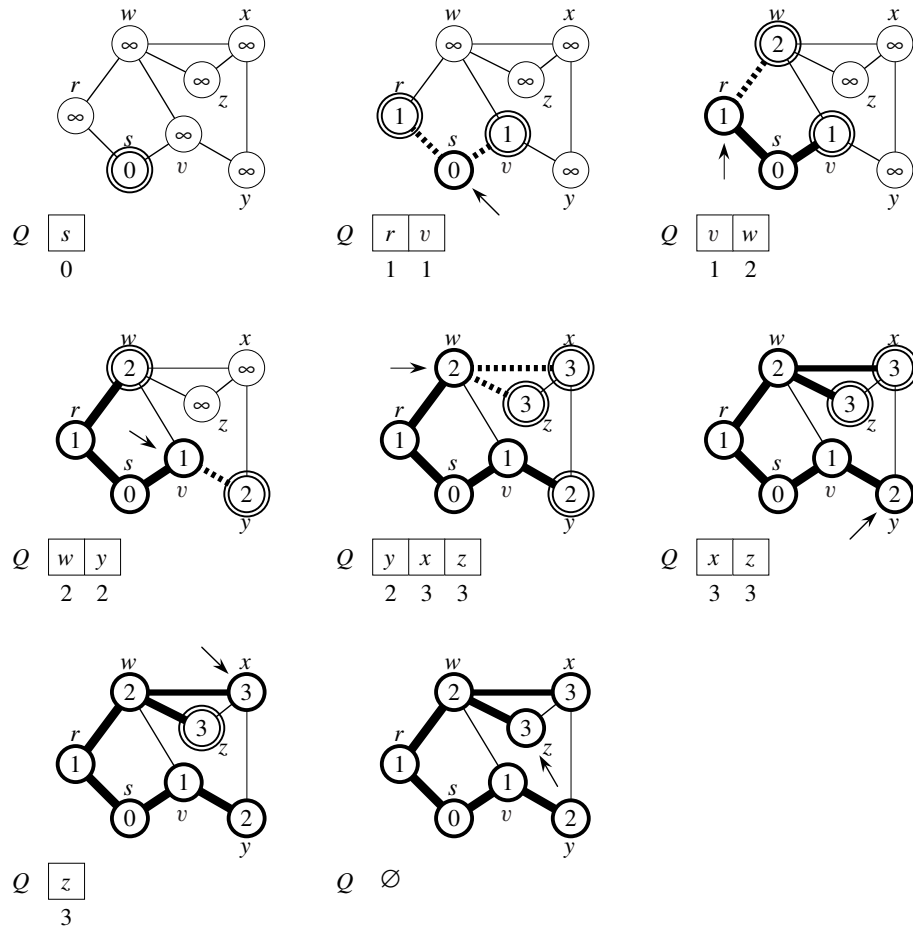
```

[In the book, the test for whether v is being newly discovered uses the colors. Checking whether $v.d$ is finite or infinite works just as well, since once v is discovered it gets a finite d value. Can also check for whether $v.\pi$ equals NIL.]

Example

BFS on an undirected graph: *[There is a more detailed, colorized example in book. Go through this example, showing how vertices are discovered and Q is updated].*

- Arrows point to the vertex being visited.
- Edges drawn with heavy lines are in the predecessor subgraph.
- Dashed lines go to newly discovered vertices. They are drawn with heavy lines because they are also now in the predecessor subgraph.
- Double-outline vertices have been discovered and are in Q , waiting to be visited.
- Heavy-outline vertices have been discovered, dequeued from Q , and visited.



Can show that Q consists of vertices with d values.

$k \quad k \quad k \quad \dots \quad k \quad k+1 \quad k+1 \quad \dots \quad k+1$

- Only 1 or 2 values.
- If 2, differ by 1 and all smallest are first.

Since each vertex gets a finite d value at most once, values assigned to vertices are monotonically increasing over time.

[Actual proof of correctness is a bit trickier. See book.]

BFS may not reach all vertices.

Time = $O(V + E)$.

- $O(V)$ because every vertex enqueued at most once.
- $O(E)$ because every vertex dequeued at most once and edge (u, v) is examined only when u is dequeued. Therefore, every edge examined at most once if directed, at most twice if undirected.

To print the vertices on a shortest path from s to v :

```

PRINT-PATH( $G, s, v$ )
  if  $v == s$ 
    print  $s$ 
  elseif  $v.\pi == \text{NIL}$ 
    print “no path from”  $s$  “to”  $v$  “exists”
  else PRINT-PATH( $G, s, v.\pi$ )
    print  $v$ 

```

Depth-first search

Input: $G = (V, E)$, directed or undirected. No source vertex given.

Output:

- 2 *timestamps* on each vertex:

- $v.d = \text{discovery time}$
- $v.f = \text{finish time}$

These will be useful for other algorithms later on.

- $v.\pi$ is v 's predecessor in the *depth-first forest* of ≥ 1 *depth-first trees*.
If $u = v.\pi$, then (u, v) is a *tree edge*.

Methodically explores *every* edge.

- Start over from different vertices as necessary.

As soon as a vertex is discovered, explore from it.

- Unlike BFS, which puts a vertex on a queue so that it's explored from later.

As DFS progresses, every vertex has a *color*:

- WHITE = undiscovered
- GRAY = discovered, but not finished (not done exploring from it)
- BLACK = finished (have found everything reachable from it)

Discovery and finish times:

- Unique integers from 1 to $2|V|$.
- For all v , $v.d < v.f$.

In other words, $1 \leq v.d < v.f \leq 2|V|$.

Pseudocode

Uses a global timestamp *time*.

DFS(G)

```

for each vertex  $u \in G.V$ 
   $u.color = \text{WHITE}$ 
   $u.\pi = \text{NIL}$ 
 $time = 0$ 
for each vertex  $u \in G.V$ 
  if  $u.color == \text{WHITE}$ 
    DFS-VISIT( $G, u$ )

```

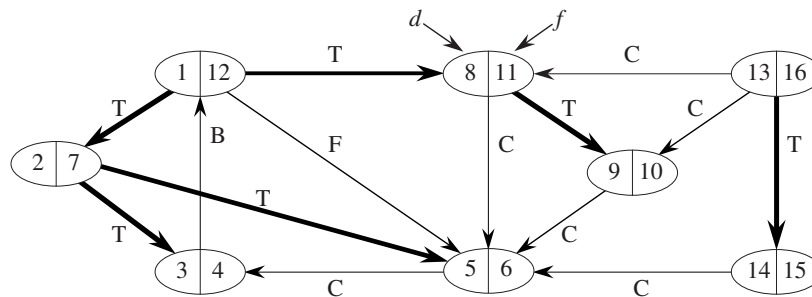
```

DFS-VISIT( $G, u$ )
     $time = time + 1$                 // white vertex  $u$  has just been discovered
     $u.d = time$ 
     $u.color = GRAY$ 
    for each vertex  $v$  in  $G.Adj[u]$  // explore each edge  $(u, v)$ 
        if  $v.color == WHITE$ 
             $v.\pi = u$ 
            DFS-VISIT( $G, v$ )
     $time = time + 1$ 
     $u.f = time$ 
     $u.color = BLACK$                 // blacken  $u$ ; it is finished

```

Example

[Go through this example of DFS on a directed graph, adding in the d and f values as they're computed. Show colors as they change. Don't put in the edge types yet, except that the tree edges are drawn with heavy lines.]



Time = $\Theta(V + E)$.

- Similar to BFS analysis.
- Θ , not just O , since guaranteed to examine every vertex and edge.

Each depth-first tree is made of edges (u, v) such that u is gray and v is white when (u, v) is explored.

Theorem (Parenthesis theorem)

[Proof omitted.]

For all u, v , exactly one of the following holds:

1. $u.d < u.f < v.d < v.f$ or $v.d < v.f < u.d < u.f$ (i.e., the intervals $[u.d, u.f]$ and $[v.d, v.f]$ are disjoint) and neither of u and v is a descendant of the other.
2. $u.d < v.d < v.f < u.f$ and v is a descendant of u . (v is discovered after and finished before u .)
3. $v.d < u.d < u.f < v.f$ and u is a descendant of v . (u is discovered after and finished before v .)

So $u.d < v.d < u.f < v.f$ (v is both discovered and finished after u) *cannot* happen.

Like parentheses:

- OK: $() []$ $([])$ $[()]$
- Not OK: $([])$ $[()]$

Corollary

v is a proper descendant of u if and only if $u.d < v.d < v.f < u.f$.

Theorem (White-path theorem)

[Proof omitted.]

v is a descendant of u if and only if at time $u.d$, there is a path $u \rightsquigarrow v$ consisting of only white vertices. (Except for u , which was *just* colored gray.)

Classification of edges

- **Tree edge:** in the depth-first forest. Found by exploring (u, v) .
- **Back edge:** (u, v) , where u is a descendant of v .
- **Forward edge:** (u, v) , where v is a descendant of u , but not a tree edge.
- **Cross edge:** any other edge. Can go between vertices in same depth-first tree or in different depth-first trees.

[Now label the example from above with edge types.]

In an undirected graph, there may be some ambiguity since (u, v) and (v, u) are the same edge. Classify by the first type above that matches.

Theorem

[Proof omitted.]

A DFS of an *undirected* graph yields only tree and back edges. No forward or cross edges.

Topological sort**Directed acyclic graph (dag)**

A directed graph with no cycles.

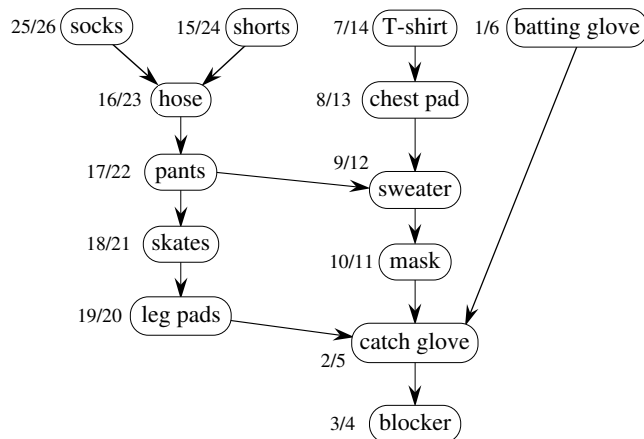
Good for modeling processes and structures that have a **partial order**:

- $a > b$ and $b > c \Rightarrow a > c$.
- But may have a and b such that neither $a > b$ nor $b > c$.

Can always make a **total order** (either $a > b$ or $b > a$ for all $a \neq b$) from a partial order. In fact, that's what a topological sort will do.

Example

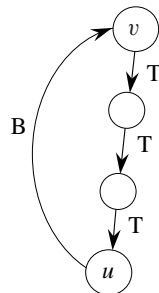
Dag of dependencies for putting on goalie equipment for ice hockey: [Leave on board, but show without discovery and finish times. Will put them in later.]

**Lemma**

A directed graph G is acyclic if and only if a DFS of G yields no back edges.

Proof \Rightarrow : Show that back edge \Rightarrow cycle.

Suppose there is a back edge (u, v) . Then v is ancestor of u in depth-first forest.



Therefore, there is a path $v \rightsquigarrow u$, so $v \rightsquigarrow u \rightarrow v$ is a cycle.

\Leftarrow : Show that cycle \Rightarrow back edge.

Suppose G contains cycle c . Let v be the first vertex discovered in c , and let (u, v) be the preceding edge in c . At time $v.d$, vertices of c form a white path $v \rightsquigarrow u$ (since v is the first vertex discovered in c). By white-path theorem, u is descendant of v in depth-first forest. Therefore, (u, v) is a back edge. ■ (lemma)

Topological sort of a dag: a linear ordering of vertices such that if $(u, v) \in E$, then u appears somewhere before v . (Not like sorting numbers.)

TOPOLOGICAL-SORT(G)

call DFS(G) to compute finish times $v.f$ for all $v \in G.V$
 output vertices in order of *decreasing* finish times

Don't need to sort by finish times.

- Can just output vertices as they're finished and understand that we want the *reverse* of this list.
- Or put them onto the *front* of a linked list as they're finished. When done, the list contains vertices in topologically sorted order.

Time

$\Theta(V + E)$.

Do example. [Now write discovery and finish times in goalie equipment example.]

Order:

```

26  socks
24  shorts
23  hose
22  pants
21  skates
20  leg pads
14  t-shirt
13  chest pad
12  sweater
11  mask
6   batting glove
5   catch glove
4   blocker

```

Correctness

Just need to show if $(u, v) \in E$, then $v.f < u.f$.

When edge (u, v) is explored, what are the colors of u and v ?

- u is gray.
- Is v gray, too?
 - No, because then v would be ancestor of u .
 $\Rightarrow (u, v)$ is a back edge.
 \Rightarrow contradiction of previous lemma (dag has no back edges).
- Is v white?
 - Then becomes descendant of u .
 By parenthesis theorem, $u.d < v.d < \underline{v.f} < u.f$.
- Is v black?
 - Then v is already finished.
 Since exploring (u, v) , u is not yet finished.
 Therefore, $v.f < u.f$. ■

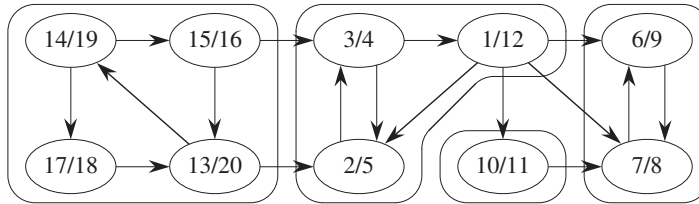
Strongly connected components

Given directed graph $G = (V, E)$.

A **strongly connected component (SCC)** of G is a maximal set of vertices $C \subseteq V$ such that for all $u, v \in C$, both $u \rightsquigarrow v$ and $v \rightsquigarrow u$.

Example

[Don't show discovery/finish times yet.]



Algorithm uses $G^T = \text{transpose of } G$.

- $G^T = (V, E^T)$, $E^T = \{(u, v) : (v, u) \in E\}$.
- G^T is G with all edges reversed.

Can create G^T in $\Theta(V + E)$ time if using adjacency lists.

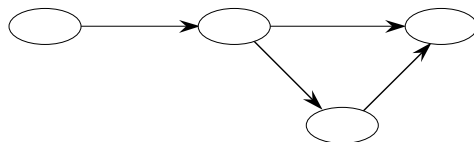
Observation

G and G^T have the *same* SCC's. (u and v are reachable from each other in G if and only if reachable from each other in G^T .)

Component graph

- $G^{\text{SCC}} = (V^{\text{SCC}}, E^{\text{SCC}})$.
- V^{SCC} has one vertex for each SCC in G .
- E^{SCC} has an edge if there's an edge between the corresponding SCC's in G .

For our example:



Lemma

G^{SCC} is a dag. More formally, let C and C' be distinct SCC's in G , let $u, v \in C$, $u', v' \in C'$, and suppose there is a path $u \rightsquigarrow u'$ in G . Then there cannot also be a path $v' \rightsquigarrow v$ in G .

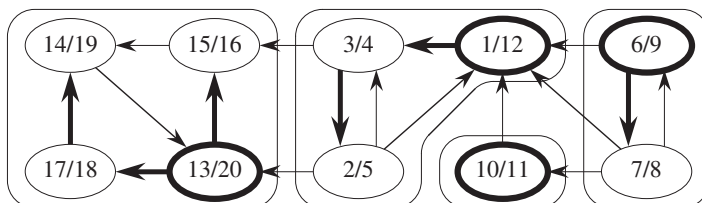
Proof Suppose there is a path $v' \rightsquigarrow v$ in G . Then there are paths $u \rightsquigarrow u' \rightsquigarrow v'$ and $v' \rightsquigarrow v \rightsquigarrow u$ in G . Therefore, u and v' are reachable from each other, so they are not in separate SCC's. ■ (lemma)

SCC(G)

call DFS(G) to compute finish times $u.f$ for each vertex u
 create G^T
 call DFS(G^T), but in the main loop, consider vertices in order of decreasing $u.f$
 (as computed in first DFS)
 output the vertices in each tree of the depth-first forest formed in second DFS
 as a separate SCC

Example:

1. Do DFS in G . [Now show discovery and finish times in G .]
2. G^T .
3. DFS in G^T . [Discovery and finish times are from first DFS in G . Roots in second DFS in G^T are drawn with heavy outlines, tree edges in second DFS are drawn with heavy lines.]



Time: $\Theta(V + E)$.

How can this possibly work?

Idea

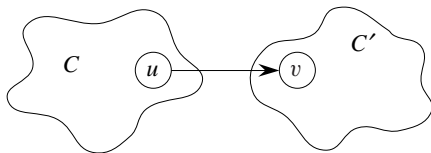
By considering vertices in second DFS in decreasing order of finish times from first DFS, visiting vertices of the component graph in topological sort order.

To prove that it works, first deal with 2 notational issues:

- Will be discussing $u.d$ and $u.f$. These always refer to *first* DFS.
- Extend notation for d and f to sets of vertices $U \subseteq V$:
 - $d(U) = \min \{u.d : u \in U\}$ (earliest discovery time in U)
 - $f(U) = \max \{u.f : u \in U\}$ (latest finish time in U)

Lemma

Let C and C' be distinct SCC's in $G = (V, E)$. Suppose that there is an edge $(u, v) \in E$ such that $u \in C$ and $v \in C'$.



Then $f(C) > f(C')$.

Proof Two cases, depending on which SCC had the first discovered vertex during the first DFS.

- If $d(C) < d(C')$, let x be the first vertex discovered in C . At time $x.d$, all vertices in C and C' are white. Thus, there exist paths of white vertices from x to all vertices in C and C' .

By the white-path theorem, all vertices in C and C' are descendants of x in depth-first tree.

By the parenthesis theorem, $x.f = f(C) > f(C')$.

- If $d(C) > d(C')$, let y be the first vertex discovered in C' . At time $y.d$, all vertices in C' are white and there is a white path from y to each vertex in $C' \Rightarrow$ all vertices in C' become descendants of y . Again, $y.f = f(C')$.

At time $y.d$, all vertices in C are white.

By earlier lemma, since there is an edge (u, v) , we cannot have a path from C' to C .

So no vertex in C is reachable from y .

Therefore, at time $y.f$, all vertices in C are still white.

Therefore, for all $w \in C$, $w.f > y.f$, which implies that $f(C) > f(C')$.

■ (lemma)

Corollary

Let C and C' be distinct SCC's in $G = (V, E)$. Suppose there is an edge $(u, v) \in E^T$, where $u \in C$ and $v \in C'$. Then $f(C) < f(C')$.

Proof $(u, v) \in E^T \Rightarrow (v, u) \in E$. Since SCC's of G and G^T are the same, $f(C') > f(C)$. ■ (corollary)

Corollary

Let C and C' be distinct SCC's in $G = (V, E)$, and suppose that $f(C) > f(C')$. Then there cannot be an edge from C to C' in G^T .

Proof It's the contrapositive of the previous corollary. ■

Now we have the intuition to understand why the SCC procedure works.

The second DFS, on G^T , starts with an SCC C such that $f(C)$ is maximum. The second DFS starts from some $x \in C$, and it visits all vertices in C . The corollary says that since $f(C) > f(C')$ for all $C' \neq C$, there are no edges from C to C' in G^T .

Therefore, the second DFS visits *only* vertices in C .

Which means that the depth-first tree rooted at x contains *exactly* the vertices of C .

The next root chosen in the second DFS is in SCC C' such that $f(C')$ is maximum over all SCC's other than C . DFS visits all vertices in C' , but the only edges out of C' go to C , *which we've already visited*.

Therefore, the only tree edges will be to vertices in C' .

The process continues.

Each root chosen for the second DFS can reach only

- vertices in its SCC—get tree edges to these,
- vertices in SCC's *already visited* in second DFS—get *no* tree edges to these.

Visiting vertices of $(G^T)^{\text{SCC}}$ in reverse of topologically sorted order.

[The book has a formal proof.]

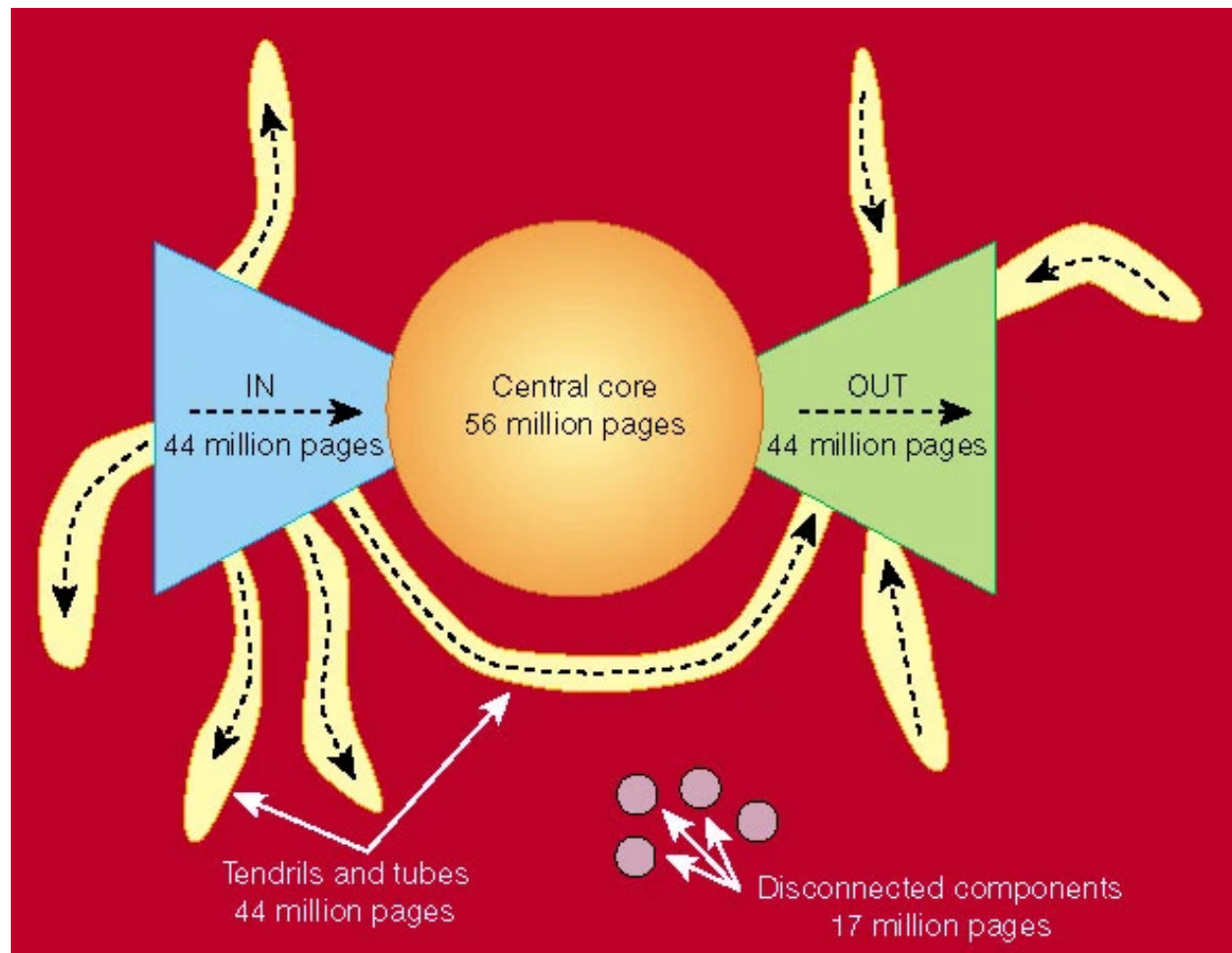


Figure 1: The Web circa 2000

Chapter 21: MSTs

Reading: Chapter 21

Lecture Notes for Chapter 21:

Minimum Spanning Trees

Chapter 21 overview

Problem

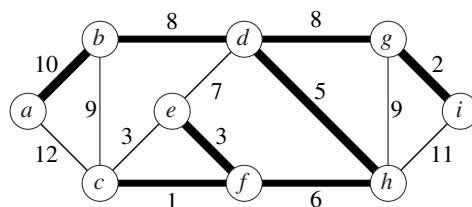
- A town has a set of houses and a set of roads.
- A road connects 2 and only 2 houses.
- A road connecting houses u and v has a repair cost $w(u, v)$.
- **Goal:** Repair enough (and no more) roads such that
 1. everyone stays connected: can reach every house from all other houses, and
 2. total repair cost is minimum.

Model as a graph:

- Undirected graph $G = (V, E)$.
- **Weight** $w(u, v)$ on each edge $(u, v) \in E$.
- Find $T \subseteq E$ such that
 1. T connects all vertices (T is a *spanning tree*), and
 2. $w(T) = \sum_{(u,v) \in T} w(u, v)$ is minimized.

A spanning tree whose weight is minimum over all spanning trees is called a *minimum spanning tree*, or *MST*.

Example of such a graph [Differs from Figure 21.1 in the textbook. Edges in the MST are drawn with heavy lines.] :



In this example, there is more than one MST. Replace edge (e, f) in the MST by (c, e) . Get a different spanning tree with the same weight.

Growing a minimum spanning tree

Some properties of an MST:

- It has $|V| - 1$ edges.
- It has no cycles.
- It might not be unique.

Building up the solution

- Build a set A of edges.
- Initially, A has no edges.
- As edges are added to A , maintain a loop invariant:
Loop invariant: A is a subset of some MST.
- Add only edges that maintain the invariant. If A is a subset of some MST, an edge (u, v) is *safe* for A if and only if $A \cup \{(u, v)\}$ is also a subset of some MST. So add only safe edges.

Generic MST algorithm

GENERIC-MST(G, w)

```

 $A = \emptyset$ 
while  $A$  does not form a spanning tree
    find an edge  $(u, v)$  that is safe for  $A$ 
     $A = A \cup \{(u, v)\}$ 
return  $A$ 

```

Use the loop invariant to show that this generic algorithm works.

Initialization: The empty set trivially satisfies the loop invariant.

Maintenance: Since only safe edges are added, A remains a subset of some MST.

Termination: The loop must terminate by the time it considers all edges. All edges added to A are in an MST, so upon termination, A is a spanning tree that is also an MST.

Finding a safe edge

How to find safe edges?

Let's look at the example. Edge (c, f) has the lowest weight of any edge in the graph. Is it safe for $A = \emptyset$?

Intuitively: Let $S \subset V$ be any proper subset of vertices that includes c but not f (so that f is in $V - S$). In any MST, there has to be one edge (at least) that connects S with $V - S$. Why not choose the edge with minimum weight? (Which would be (c, f) in this case.)

Some definitions: Let $S \subset V$ and $A \subseteq E$.

- A **cut** $(S, V - S)$ is a partition of vertices into disjoint sets S and $V - S$.
- Edge $(u, v) \in E$ **crosses** cut $(S, V - S)$ if one endpoint is in S and the other is in $V - S$.
- A cut **respects** A if and only if no edge in A crosses the cut.
- An edge is a **light edge** crossing a cut if and only if its weight is minimum over all edges crossing the cut. For a given cut, there can be > 1 light edge crossing it.

Theorem

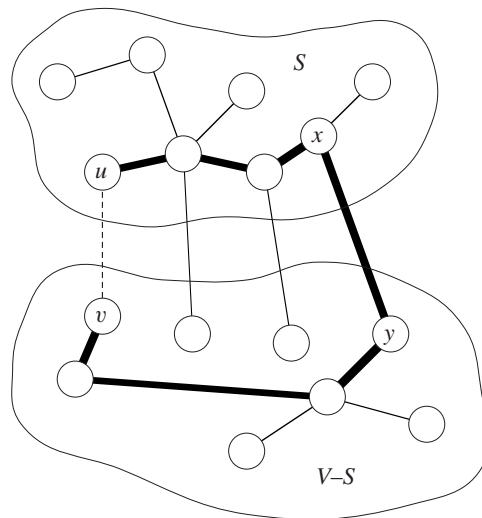
Let A be a subset of some MST, $(S, V - S)$ be a cut that respects A , and (u, v) be a light edge crossing $(S, V - S)$. Then (u, v) is safe for A .

Proof Let T be an MST that includes A .

If T contains (u, v) , done.

So now assume that T does not contain (u, v) . Construct a different MST T' that includes $A \cup \{(u, v)\}$.

Recall: a tree has unique path between each pair of vertices. Since T is an MST, it contains a unique path p between u and v . Path p must cross the cut $(S, V - S)$ at least once. Let (x, y) be an edge of p that crosses the cut. From how we chose (u, v) , must have $w(u, v) \leq w(x, y)$.



[Except for the dashed edge (u, v) , all edges shown are in T . A is some subset of the edges of T , but A cannot contain any edges that cross the cut $(S, V - S)$, since this cut respects A . Edges with heavy lines are the path p .]

Since the cut respects A , edge (x, y) is not in A .

To form T' from T :

- Remove (x, y) . Breaks T into two components.
- Add (u, v) . Reconnects.

So $T' = T - \{(x, y)\} \cup \{(u, v)\}$.

T' is a spanning tree.

$$\begin{aligned} w(T') &= w(T) - w(x, y) + w(u, v) \\ &\leq w(T), \end{aligned}$$

since $w(u, v) \leq w(x, y)$. Since T' is a spanning tree, $w(T') \leq w(T)$, and T is an MST, then T' must be an MST.

Need to show that (u, v) is safe for A :

- $A \subseteq T$ and $(x, y) \notin A \Rightarrow A \subseteq T'$.
- $A \cup \{(u, v)\} \subseteq T'$.
- Since T' is an MST, (u, v) is safe for A . ■ (theorem)

So, in GENERIC-MST:

- A is a forest containing connected components. Initially, each component is a single vertex.
- Any safe edge merges two of these components into one. Each component is a tree.
- Since an MST has exactly $|V| - 1$ edges, the **for** loop iterates $|V| - 1$ times. Equivalently, after adding $|V| - 1$ safe edges, we're down to just one component.

Corollary

If $C = (V_C, E_C)$ is a connected component in the forest $G_A = (V, A)$ and (u, v) is a light edge connecting C to some other component in G_A (i.e., (u, v) is a light edge crossing the cut $(V_C, V - V_C)$), then (u, v) is safe for A .

Proof Set $S = V_C$ in the theorem. ■ (corollary)

This idea naturally leads to the algorithm known as Kruskal's algorithm to solve the minimum-spanning-tree problem.

Kruskal's algorithm

$G = (V, E)$ is a connected, undirected, weighted graph. $w : E \rightarrow \mathbb{R}$.

- Starts with each vertex being its own component.
- Repeatedly merges two components into one by choosing the light edge that connects them (i.e., the light edge crossing the cut between them).
- Scans the set of edges in monotonically increasing order by weight.
- Uses a disjoint-set data structure to determine whether an edge connects vertices in different components.

MST-KRUSKAL(G, w)

$A = \emptyset$

for each vertex $v \in G.V$

 MAKE-SET(v)

 create a single list of the edges in $G.E$

 sort the list of edges into nondecreasing order by weight w

for each edge (u, v) taken from the sorted list in order

if FIND-SET(u) \neq FIND-SET(v)

$A = A \cup \{(u, v)\}$

 UNION(u, v)

return A

Run through the above example to see how Kruskal's algorithm works on it:

(c, f) : safe

(g, i) : safe

(e, f) : safe

(c, e) : reject

(d, h) : safe

(f, h) : safe

(e, d) : reject

(b, d) : safe

(d, g) : safe

(b, c) : reject

(g, h) : reject

(a, b) : safe

At this point, there is only one component, so that all other edges will be rejected.
[Could add a test to the main loop of KRUSKAL to stop once $|V| - 1$ edges have been added to A .]

Get the heavy edges shown in the figure.

Suppose (c, e) had been examined *before* (e, f) . Then would have found (c, e) safe and would have rejected (e, f) .

Analysis

Initialize A : $O(1)$

First **for** loop: $|V|$ MAKE-SETs

Sort E : $O(E \lg E)$

Second **for** loop: $O(E)$ FIND-SETs and UNIONS

- Assuming the implementation of disjoint-set data structure, already seen in Chapter 19, that uses union by rank and path compression:

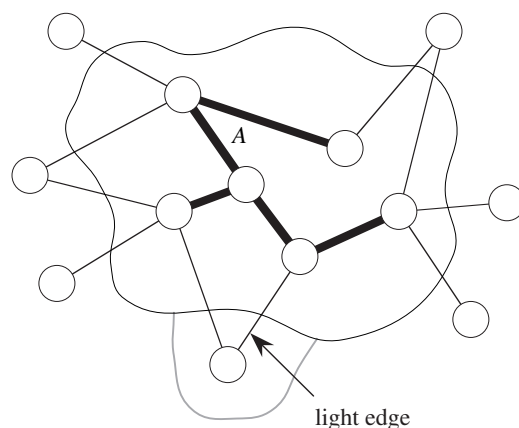
$$O((V + E) \alpha(V)) + O(E \lg E).$$

- Since G is connected, $|E| \geq |V| - 1 \Rightarrow O(E \alpha(V)) + O(E \lg E)$.
- $\alpha(|V|) = O(\lg V) = O(\lg E)$.
- Therefore, total time is $O(E \lg E)$.

- $|E| \leq |V|^2 \Rightarrow \lg |E| = O(2 \lg V) = O(\lg V)$.
- Therefore, $O(E \lg V)$ time. (If edges are already sorted, $O(E \alpha(V))$, which is almost linear.)

Prim's algorithm

- Builds one tree, so A is always a tree.
- Starts from an arbitrary “root” r .
- At each step, find a light edge connecting A to an isolated vertex. Such an edge must be safe for A . Add this edge to A .



[Edges of A are drawn with heavy lines.]

How to find the light edge quickly?

Use a priority queue Q :

- Each object is a vertex *not* in A .
- $v.key$ is the minimum weight of any edge connecting v to a vertex in A . $v.key = \infty$ if no such edge.
- $v.\pi$ is v 's parent in A .
- Maintain A implicitly as $A = \{(v, v.\pi) : v \in V - \{r\} - Q\}$.
- At completion, Q is empty and the minimum spanning tree is $A = \{(v, v.\pi) : v \in V - \{r\}\}$.

Chapter 22: Single-Source Shortest Paths

Reading: Chapter 22 intro, 22.1, 22.3, 22.5 (skip DAGs and difference constraints)

Lecture Notes for Chapter 22: Single-Source Shortest Paths

Shortest paths

How to find the shortest route between two points on a map.

Input:

- Directed graph $G = (V, E)$
- Weight function $w : E \rightarrow \mathbb{R}$

Weight of path $p = \langle v_0, v_1, \dots, v_k \rangle$

$$= \sum_{i=1}^k w(v_{i-1}, v_i)$$

= sum of edge weights on path p .

Shortest-path weight u to v :

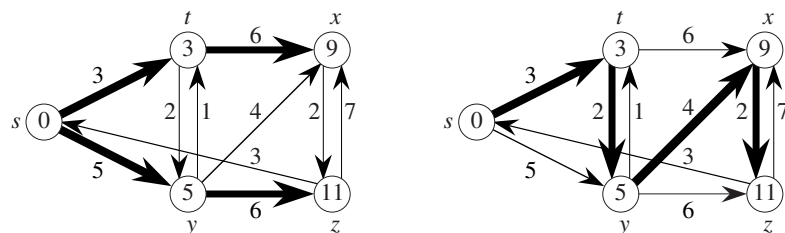
$$\delta(u, v) = \begin{cases} \min\{w(p) : u \xrightarrow{p} v\} & \text{if there exists a path } u \rightsquigarrow v, \\ \infty & \text{otherwise.} \end{cases}$$

Shortest path u to v is any path p such that $w(p) = \delta(u, v)$.

Example

shortest paths from s

[δ values appear inside vertices. Heavy edges show shortest paths.]



This example shows that a shortest path might not be unique.

It also shows that when we look at shortest paths from one vertex to all other vertices, the shortest paths are organized as a tree.

Can think of weights as representing any measure that

- accumulates linearly along a path, and
- we want to minimize.

Examples: time, cost, penalties, loss.

Generalization of breadth-first search to weighted graphs.

Variants

- **Single-source:** Find shortest paths from a given **source** vertex $s \in V$ to every vertex $v \in V$.
- **Single-destination:** Find shortest paths to a given destination vertex.
- **Single-pair:** Find shortest path from u to v . No way known that's better in worst case than solving single-source.
- **All-pairs:** Find shortest path from u to v for all $u, v \in V$. We'll see algorithms for all-pairs in the next chapter.

Negative-weight edges

OK, as long as no negative-weight cycles are reachable from the source.

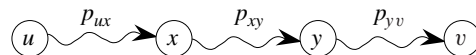
- If we have a negative-weight cycle, we can just keep going around it, and get $w(s, v) = -\infty$ for all v on the cycle.
- But OK if the negative-weight cycle is not reachable from the source.
- Some algorithms work only if there are no negative-weight edges in the graph. We'll be clear when they're allowed and not allowed.

Optimal substructure

Lemma

Any subpath of a shortest path is a shortest path.

Proof Cut-and-paste.



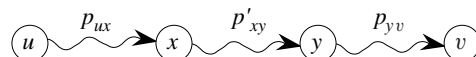
Suppose this path p is a shortest path from u to v .

Then $\delta(u, v) = w(p) = w(p_{ux}) + w(p_{xy}) + w(p_{yv})$.

Now suppose there exists a shorter path $x \xrightarrow{p'_{xy}} y$.

Then $w(p'_{xy}) < w(p_{xy})$.

Construct p' :



Then

$$\begin{aligned} w(p') &= w(p_{ux}) + w(p'_{xy}) + w(p_{yv}) \\ &< w(p_{ux}) + w(p_{xy}) + w(p_{yv}) \\ &= w(p). \end{aligned}$$

Contradicts the assumption that p is a shortest path.

■ (lemma)

Cycles

Shortest paths can't contain cycles:

- Already ruled out negative-weight cycles.
- Positive-weight \Rightarrow we can get a shorter path by omitting the cycle.
- 0-weight: no reason to use them \Rightarrow assume that our solutions won't use them.

Output of single-source shortest-path algorithm

For each vertex $v \in V$:

- $v.d = \delta(s, v)$.
 - Initially, $v.d = \infty$.
 - Reduces as algorithms progress. But always maintain $v.d \geq \delta(s, v)$.
 - Call $v.d$ a **shortest-path estimate**.
- $v.\pi$ = predecessor of v on a shortest path from s .
 - If no predecessor, $v.\pi = \text{NIL}$.
 - π induces a tree—**shortest-path tree**.
 - We won't prove properties of π in lecture—see text.

Initialization

All the shortest-paths algorithms start with INITIALIZE-SINGLE-SOURCE.

INITIALIZE-SINGLE-SOURCE(G, s)

for each vertex $v \in G.V$

$v.d = \infty$

$v.\pi = \text{NIL}$

$s.d = 0$

Relaxing an edge (u, v)

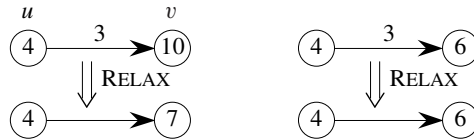
Can the shortest-path estimate for v be improved by going through u and taking (u, v) ?

RELAX(u, v, w)

if $v.d > u.d + w(u, v)$

$v.d = u.d + w(u, v)$

$v.\pi = u$



For all the single-source shortest-paths algorithms we'll look at,

- start by calling INITIALIZE-SINGLE-SOURCE,
- then relax edges.

The algorithms differ in the order and how many times they relax each edge.

Shortest-paths properties

[The textbook states these properties in the chapter introduction and proves them in a later section. You might elect to just state these properties at first and prove them later.]

Based on calling INITIALIZE-SINGLE-SOURCE once and then calling RELAX zero or more times.

Triangle inequality: For all $(u, v) \in E$, we have $\delta(s, v) \leq \delta(s, u) + w(u, v)$.

Proof Weight of shortest path $s \rightsquigarrow v$ is \leq weight of any path $s \rightsquigarrow v$. Path $s \rightsquigarrow u \rightarrow v$ is a path $s \rightsquigarrow v$, and if we use a shortest path $s \rightsquigarrow u$, its weight is $\delta(s, u) + w(u, v)$. ■

Upper-bound property: Always have $v.d \geq \delta(s, v)$ for all v . Once $v.d$ gets down to $\delta(s, v)$, it never changes.

Proof Initially true.

Suppose there exists a vertex such that $v.d < \delta(s, v)$.

Without loss of generality, v is first vertex for which this happens.

Let u be the vertex that causes $v.d$ to change.

Then $v.d = u.d + w(u, v)$.

So,

$$\begin{aligned}
 v.d &< \delta(s, v) \\
 &\leq \delta(s, u) + w(u, v) \quad (\text{triangle inequality}) \\
 &\leq u.d + w(u, v) \quad (v \text{ is first violation}) \\
 \Rightarrow v.d &< u.d + w(u, v) .
 \end{aligned}$$

Contradicts $v.d = u.d + w(u, v)$.

Once $v.d$ reaches $\delta(s, v)$, it never goes lower. It never goes up, since relaxations only lower shortest-path estimates. ■

No-path property: If $\delta(s, v) = \infty$, then $v.d = \infty$ always.

Proof $v.d \geq \delta(s, v) = \infty \Rightarrow v.d = \infty$. ■

Convergence property: If $s \rightsquigarrow u \rightarrow v$ is a shortest path, $u.d = \delta(s, u)$, and edge (u, v) is relaxed, then $v.d = \delta(s, v)$ afterward.

Proof After relaxation:

$$\begin{aligned} v.d &\leq u.d + w(u, v) && \text{(RELAX code)} \\ &= \delta(s, u) + w(u, v) \\ &= \delta(s, v) && \text{(lemma—optimal substructure)} \end{aligned}$$

Since $v.d \geq \delta(s, v)$, must have $v.d = \delta(s, v)$. ■

Path-relaxation property: Let $p = \langle v_0, v_1, \dots, v_k \rangle$ be a shortest path from $s = v_0$ to v_k . If the edges of p are relaxed, *in the order*, $(v_0, v_1), (v_1, v_2), \dots, (v_{k-1}, v_k)$, even intermixed with other relaxations, then $v_k.d = \delta(s, v_k)$.

Proof Induction to show that $v_i.d = \delta(s, v_i)$ after (v_{i-1}, v_i) is relaxed.

Basis: $i = 0$. Initially, $v_0.d = 0 = \delta(s, v_0) = \delta(s, s)$.

Inductive step: Assume $v_{i-1}.d = \delta(s, v_{i-1})$. Relax (v_{i-1}, v_i) . By convergence property, $v_i.d = \delta(s, v_i)$ afterward and $v_i.d$ never changes. ■

The Bellman-Ford algorithm

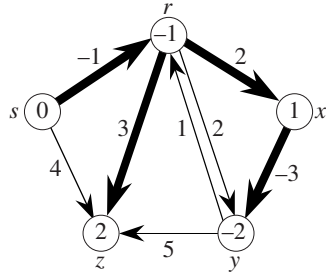
- Allows negative-weight edges.
- Computes $v.d$ and $v.\pi$ for all $v \in V$.
- Returns TRUE if no negative-weight cycles reachable from s , FALSE otherwise.

```

BELLMAN-FORD( $G, w, s$ )
  INITIALIZE-SINGLE-SOURCE( $G, s$ )
  for  $i = 1$  to  $|G.V| - 1$ 
    for each edge  $(u, v) \in G.E$ 
      RELAX( $u, v, w$ )
  for each edge  $(u, v) \in G.E$ 
    if  $v.d > u.d + w(u, v)$ 
      return FALSE
  return TRUE

```

Time: $O(V^2 + VE)$. The first **for** loop makes $|V| - 1$ passes over the edges, and each pass takes $\Theta(V + E)$ time. We use O rather than Θ because sometimes $< |V| - 1$ passes are enough (Exercise 22.1-3).

Example

Values you get on each pass and how quickly it converges depends on order of relaxation.

But guaranteed to converge after $|V| - 1$ passes, assuming no negative-weight cycles.

Proof Use path-relaxation property.

Let v be reachable from s , and let $p = \langle v_0, v_1, \dots, v_k \rangle$ be a shortest path from s to v , where $v_0 = s$ and $v_k = v$. Since p is acyclic, it has $\leq |V| - 1$ edges, so that $k \leq |V| - 1$.

Each iteration of the **for** loop relaxes all edges:

- First iteration relaxes (v_0, v_1) .
- Second iteration relaxes (v_1, v_2) .
- k th iteration relaxes (v_{k-1}, v_k) .

By the path-relaxation property, $v.d = v_k.d = \delta(s, v_k) = \delta(s, v)$. ■

How about the TRUE/FALSE return value?

- Suppose there is no negative-weight cycle reachable from s .

At termination, for all $(u, v) \in E$,

$$\begin{aligned}
 v.d &= \delta(s, v) \\
 &\leq \delta(s, u) + w(u, v) \quad (\text{triangle inequality}) \\
 &= u.d + w(u, v) .
 \end{aligned}$$

So BELLMAN-FORD returns TRUE.

- Now suppose there exists negative-weight cycle $c = \langle v_0, v_1, \dots, v_k \rangle$, where $v_0 = v_k$, reachable from s .

$$\text{Then } \sum_{i=1}^k w(v_{i-1}, v_i) < 0 .$$

Suppose (for contradiction) that BELLMAN-FORD returns TRUE.

Then $v_i.d \leq v_{i-1}.d + w(v_{i-1}, v_i)$ for $i = 1, 2, \dots, k$.

Sum around c :

$$\begin{aligned}
 \sum_{i=1}^k v_i.d &\leq \sum_{i=1}^k (v_{i-1}.d + w(v_{i-1}, v_i)) \\
 &= \sum_{i=1}^k v_{i-1}.d + \sum_{i=1}^k w(v_{i-1}, v_i)
 \end{aligned}$$

Each vertex appears once in each summation $\sum_{i=1}^k v_i \cdot d$ and $\sum_{i=1}^k v_{i-1} \cdot d \Rightarrow$

$$0 \leq \sum_{i=1}^k w(v_{i-1}, v_i) \ .$$

Contradicts c being a negative-weight cycle. ■

Single-source shortest paths in a directed acyclic graph

Since a dag, we're guaranteed no negative-weight cycles.

DAG-SHORTEST-PATHS(G, w, s)

topologically sort the vertices of G

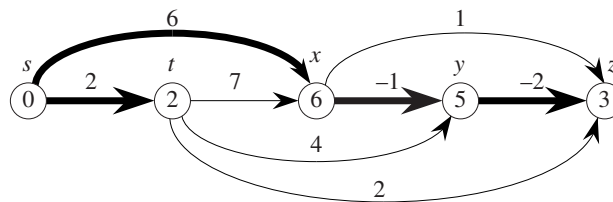
INITIALIZE-SINGLE-SOURCE(G, s)

for each vertex $u \in G.V$, taken in topologically sorted order

for each vertex v in $G.Adj[u]$

 RELAX(u, v, w)

Example



Time

$\Theta(V + E)$.

Correctness

Because vertices are processed in topologically sorted order, edges of *any* path must be relaxed in order of appearance in the path.

\Rightarrow Edges on any shortest path are relaxed in order.

\Rightarrow By path-relaxation property, correct. ■

Dijkstra's algorithm

No negative-weight *edges*.

Essentially a weighted version of breadth-first search.

- Instead of a FIFO queue, uses a priority queue.
- Keys are shortest-path weights ($v.d$).
- Can think of waves, like BFS.

- A wave emanates from the source.
- The first time that a wave arrives at a vertex, a new wave emanates from that vertex.
- The time it takes for the wave to arrive at a neighboring vertex equals the weight of the edge. (In BFS, each wave takes unit time to arrive at each neighbor.)

Have two sets of vertices:

- S = vertices whose final shortest-path weights are determined,
- Q = priority queue = $V - S$.

DIJKSTRA(G, w, s)

INITIALIZE-SINGLE-SOURCE(G, s)

$S = \emptyset$

$Q = \emptyset$

for each vertex $u \in G.V$

 INSERT(Q, u)

while $Q \neq \emptyset$

$u = \text{EXTRACT-MIN}(Q)$

$S = S \cup \{u\}$

for each vertex v in $G.Adj[u]$

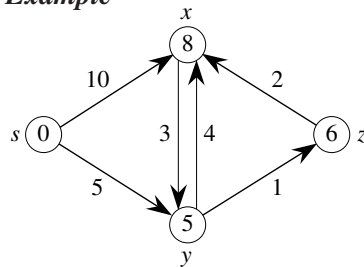
 RELAX(u, v, w)

if the call of RELAX decreased $v.d$

 DECREASE-KEY($Q, v, v.d$)

- Looks a lot like Prim's algorithm, but computing $v.d$, and using shortest-path weights as keys.
- Dijkstra's algorithm can be viewed as greedy, since it always chooses the "lightest" ("closest"?) vertex in $V - S$ to add to S .

Example



Order of adding to S : s, y, z, x .

Correctness

We will show that at the start of each iteration of the **while** loop, $v.d = \delta(s, v)$ for all $v \in S$. The algorithm terminates when $S = V$, so that $v.d = \delta(s, v)$ for all $v \in V$.

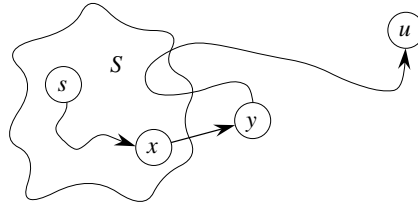
The proof is by induction on the number of iterations of the **while** loop, i.e., on $|S|$. The bases are for $|S| = 0$, so that $S = \emptyset$ and the claim is trivially true, and for $|S| = 1$, so that $S = \{s\}$ and $s.d = \delta(s, s) = 0$.

Inductive hypothesis: $v.d = \delta(s, v)$ for all $v \in S$.

Inductive step: The algorithm extracts vertex u from $V - S$. Because the algorithm adds u into S , we need to show that $u.d = \delta(s, u)$ at that time. If there is no path from s to u , then we are done, by the no-path property.

If there is a path from s to u :

- Let y be the first vertex on a shortest path from s to u that is *not* in S .
- Let $x \in S$ be the predecessor of y on that shortest path.
- Could have $y = u$ or $x = s$.



- y appears no later than u on the shortest path and all edge weights are nonnegative $\Rightarrow \delta(s, y) \leq \delta(s, u)$.
- How we chose $u \Rightarrow u.d \leq y.d$ at the time u is extracted from $V - S$.
- Upper-bound property $\Rightarrow \delta(s, u) \leq u.d$.
- $x \in S \Rightarrow x.d = \delta(s, x)$. Edge (x, y) was relaxed when x was added into S . Convergence property \Rightarrow set $y.d = \delta(s, y)$ at that time.
- Thus, we have $\delta(s, y) \leq \delta(s, u) \leq u.d \leq y.d$ and $y.d = \delta(s, y) \Rightarrow \delta(s, y) = \delta(s, u) = u.d = y.d$.
- Hence, $u.d = \delta(s, u)$. Upper-bound property $\Rightarrow u.d$ doesn't change afterward. ■

Analysis

$|V|$ INSERT and EXTRACT-MIN operations.

$\leq |E|$ DECREASE-KEY operations.

Like Prim's algorithm, depends on implementation of priority queue.

- If binary heap, each operation takes $O(\lg V)$ time $\Rightarrow O(E \lg V)$.
- If a Fibonacci heap:
 - Each EXTRACT-MIN takes $O(1)$ amortized time.
 - There are $\Theta(V)$ INSERT and EXTRACT-MIN operations, taking $O(\lg V)$ amortized time each.
 - Therefore, time is $O(V \lg V + E)$.

Difference constraints

Special case of linear programming.

Given a set of inequalities of the form $x_j - x_i \leq b_k$.

Chapter 23: All-Pairs Shortest Paths

Reading: Chapter 23 intro, 23.1, 23.2 (focus on Floyd-Warshall algorithm)

Lecture Notes for Chapter 23:

All-Pairs Shortest Paths

Chapter 23 overview

Given a directed graph $G = (V, E)$, weight function $w : E \rightarrow \mathbb{R}$, $|V| = n$. Assume that the vertices are numbered $1, 2, \dots, n$.

Goal: create an $n \times n$ matrix $D = (d_{ij})$ of shortest-path distances, so that $d_{ij} = \delta(i, j)$ for all vertices i and j .

Could run BELLMAN-FORD once from each vertex:

- $O(V^2 E)$ —which is $O(V^4)$ if the graph is *dense* ($E = \Theta(V^2)$).

If no negative-weight edges, could run Dijkstra's algorithm once from each vertex:

- $O(VE \lg V)$ with binary heap— $O(V^3 \lg V)$ if dense,
- $O(V^2 \lg V + VE)$ with Fibonacci heap— $O(V^3)$ if dense.

We'll see how to do in $O(V^3)$ in all cases, with no fancy data structure.

Shortest paths and matrix multiplication

Assume that G is given as adjacency matrix of weights: $W = (w_{ij})$, with vertices numbered 1 to n .

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

Won't worry about predecessors—see book.

Will use dynamic programming at first.

Optimal substructure

Recall: subpaths of shortest paths are shortest paths.

Recursive solution

Let $l_{ij}^{(r)}$ = weight of shortest path $i \rightsquigarrow j$ that contains $\leq r$ edges.

- $r = 0$
 \Rightarrow there is a shortest path $i \rightsquigarrow j$ with $\leq r$ edges if and only if $i = j$
 $\Rightarrow l_{ij}^{(0)} = \begin{cases} 0 & \text{if } i = j, \\ \infty & \text{if } i \neq j. \end{cases}$
- $r \geq 1$
 $\Rightarrow l_{ij}^{(r)} = \min \left\{ l_{ij}^{(r-1)}, \min \{ l_{ik}^{(r-1)} + w_{kj} : 1 \leq k \leq n \} \right\}$
(k ranges over all possible predecessors of j)
 $= \min \{ l_{ik}^{(r-1)} + w_{kj} : 1 \leq k \leq n \}$ (since $w_{jj} = 0$ for all j).
- Observe that when $r = 1$, must have $l_{ij}^{(1)} = w_{ij}$.
 Conceptually, when the path is restricted to at most 1 edge, the weight of the shortest path $i \rightsquigarrow j$ must be w_{ij} .
 And the math works out, too:

$$\begin{aligned} l_{ij}^{(1)} &= \min \{ l_{ik}^{(0)} + w_{kj} : 1 \leq k \leq n \} \\ &= l_{ii}^{(0)} + w_{ij} \quad (l_{ii}^{(0)} \text{ is the only non-}\infty \text{ among } l_{ik}^{(0)}) \\ &= w_{ij}. \end{aligned}$$

All simple shortest paths contain $\leq n - 1$ edges

$$\Rightarrow \delta(i, j) = l_{ij}^{(n-1)} = l_{ij}^{(n)} = l_{ij}^{(n+1)} = \dots$$

Compute a solution bottom-up

Compute $L^{(1)}, L^{(2)}, \dots, L^{(n-1)}$.

Start with $L^{(1)} = W$, since $l_{ij}^{(1)} = w_{ij}$.

Go from $L^{(r-1)}$ to $L^{(r)}$:

EXTEND-SHORTEST-PATHS($L^{(r-1)}, W, L^{(r)}, n$)

// Assume that the elements of $L^{(r)}$ are initialized to ∞ .

for $i = 1$ **to** n

for $j = 1$ **to** n

for $k = 1$ **to** n

$$l_{ij}^{(r)} = \min \{ l_{ij}^{(r)}, l_{ik}^{(r-1)} + w_{kj} \}$$

Compute each $L^{(r)}$:

SLOW-APSP($W, L^{(0)}, n$)

let $L = (l_{ij})$ and $M = (m_{ij})$ be new $n \times n$ matrices

$L = L^{(0)}$

for $r = 1$ **to** $n - 1$

$M = \infty$ // initialize M

 EXTEND-SHORTEST-PATHS(L, W, M, n)

$L = M$

return L

Time

- EXTEND-SHORTEST-PATHS: $\Theta(n^3)$.
- SLOW-ALL-APSP: $\Theta(n^4)$.

Observation

EXTEND-SHORTEST-PATHS is like matrix multiplication. Make the following substitutions in the equation $l_{ij}^{(r)} = \min \{l_{ik}^{(r-1)} + w_{kj} : 1 \leq k \leq n\}$:

$$l^{(r-1)} \rightarrow a$$

$$w \rightarrow b$$

$$l^{(r)} \rightarrow c$$

$$\min \rightarrow +$$

$$+ \rightarrow \cdot$$

You get $c_{ij} = \sum_{k=1}^n a_{ik} \cdot b_{kj}$, which is the equation for computing c_{ij} in matrix multiplication.

Making these changes to EXTEND-SHORTEST-PATHS and replacing ∞ (identity for min) with 0 (identity for +) gives a procedure for matrix multiplication:

```

for  $i = 1$  to  $n$ 
  for  $j = 1$  to  $n$ 
    for  $k = 1$  to  $n$ 
       $c_{ij} = c_{ij} + a_{ik} \cdot b_{kj}$ 

```

So, we can view EXTEND-SHORTEST-PATHS as just like matrix multiplication!

Why do we care?

Because our goal is to compute $L^{(n-1)}$ as fast as we can. Don't need to compute *all* the intermediate $L^{(1)}, L^{(2)}, L^{(3)}, \dots, L^{(n-2)}$.

Suppose we had a matrix A and we wanted to compute A^{n-1} (like calling EXTEND-SHORTEST-PATHS $n - 1$ times).

Could compute A, A^2, A^4, A^8, \dots

If we knew $A^r = A^{n-1}$ for all $r \geq n - 1$, could just finish with A^r , where r is the smallest power of 2 that is $\geq n - 1$ ($r = 2^{\lceil \lg(n-1) \rceil}$).

FASTER-ALL-PAIRS-SHORTEST-PATHS(W, n)

let L and M be new $n \times n$ matrices

$L = W$

$r = 1$

while $r < n - 1$

$M = \infty$ // initialize M

 EXTEND-SHORTEST-PATHS(L, L, M, n) // compute $M = L^2$

$r = 2r$

$L = M$ // ready for the next iteration

return L

OK to overshoot, since products don't change after $L^{(n-1)}$.

Time

$$\Theta(n^3 \lg n).$$

Floyd-Warshall algorithm

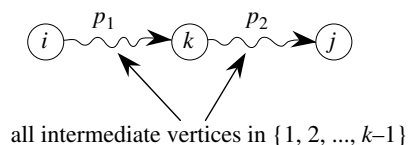
A different dynamic-programming approach.

For path $p = \langle v_1, v_2, \dots, v_l \rangle$, an *intermediate vertex* is any vertex of p other than v_1 or v_l .

Let $d_{ij}^{(k)}$ = shortest-path weight of any path $i \rightsquigarrow j$ with all intermediate vertices in $\{1, 2, \dots, k\}$.

Consider a shortest path $i \overset{p}{\rightsquigarrow} j$ with all intermediate vertices in $\{1, 2, \dots, k\}$:

- If k is not an intermediate vertex, then all intermediate vertices of p are in $\{1, 2, \dots, k-1\}$.
- If k is an intermediate vertex:

**Recursive formulation**

$$d_{ij}^{(k)} = \begin{cases} w_{ij} & \text{if } k = 0, \\ \min \{d_{ij}^{(k-1)}, d_{ik}^{(k-1)} + d_{kj}^{(k-1)}\} & \text{if } k \geq 1. \end{cases}$$

Have $d_{ij}^{(0)} = w_{ij}$ because can't have intermediate vertices $\Rightarrow \leq 1$ edge.

Want $D^{(n)} = (d_{ij}^{(n)})$, since all vertices numbered $\leq n$.

Compute bottom-up

Compute in increasing order of k :

FLOYD-WARSHALL(W, n)

$D^{(0)} = W$

for $k = 1$ **to** n

 let $D^{(k)} = (d_{ij}^{(k)})$ be a new $n \times n$ matrix

for $i = 1$ **to** n

for $j = 1$ **to** n

$d_{ij}^{(k)} = \min \{d_{ij}^{(k-1)}, d_{ik}^{(k-1)} + d_{kj}^{(k-1)}\}$

return $D^{(n)}$

Can drop superscripts. (See Exercise 23.2-4 in text.)

Time $\Theta(n^3)$.**Computing predecessors**

Can compute predecessor matrix Π while computing the D matrices. Let $\Pi^{(k)} = (\pi_{ij}^{(k)})$ for $k = 0, 1, \dots, n$.

Define $\pi_{ij}^{(k)}$ recursively. For $k = 0$, a shortest path from i to j has no intermediate vertices:

$$\pi_{ij}^{(0)} = \begin{cases} \text{NIL} & \text{if } i = j \text{ or } w_{ij} = \infty, \\ i & \text{if } i \neq j \text{ and } w_{ij} < \infty. \end{cases}$$

For $k \geq 1$:

- If shortest path from i to j has k as an intermediate vertex, then it's $i \rightsquigarrow k \rightsquigarrow j$ where $k \neq j$. Choose j 's predecessor to be the predecessor of j on a shortest path from k to j with all intermediate vertices $< k$: $\pi_{ij}^{(k)} = \pi_{kj}^{(k-1)}$.
- Otherwise, shortest path from i to j does not have k as an intermediate vertex. Keep the same predecessor as shortest path from i to j with all intermediate vertices $< k$: $\pi_{ij}^{(k)} = \pi_{ij}^{(k-1)}$.

Transitive closure

Given $G = (V, E)$, directed.

Compute $G^* = (V, E^*)$.

- $E^* = \{(i, j) : \text{there is a path } i \rightsquigarrow j \text{ in } G\}$.

Could assign weight of 1 to each edge, then run FLOYD-WARSHALL.

- If $d_{ij} < n$, then there is a path $i \rightsquigarrow j$.
- Otherwise, $d_{ij} = \infty$ and there is no path.

Simpler way

Substitute other values and operators in FLOYD-WARSHALL.

- Use unweighted adjacency matrix
- $\min \rightarrow \vee$ (OR)
- $+$ $\rightarrow \wedge$ (AND)
- $t_{ij}^{(k)} = \begin{cases} 1 & \text{if there is path } i \rightsquigarrow j \text{ with all intermediate vertices in } \{1, 2, \dots, k\}, \\ 0 & \text{otherwise.} \end{cases}$
- $t_{ij}^{(0)} = \begin{cases} 0 & \text{if } i \neq j \text{ and } (i, j) \notin E, \\ 1 & \text{if } i = j \text{ or } (i, j) \in E. \end{cases}$
- $t_{ij}^{(k)} = t_{ij}^{(k-1)} \vee (t_{ik}^{(k-1)} \wedge t_{kj}^{(k-1)})$.

TRANSITIVE-CLOSURE(G, n)

```

let  $T^{(0)} = (t_{ij}^{(0)})$  be a new  $n \times n$  matrix
for  $i = 1$  to  $n$ 
  for  $j = 1$  to  $n$ 
    if  $i = j$  or  $(i, j) \in G.E$ 
       $t_{ij}^{(0)} = 1$ 
    else  $t_{ij}^{(0)} = 0$ 
for  $k = 1$  to  $n$ 
  let  $T^{(k)} = (t_{ij}^{(k)})$  be a new  $n \times n$  matrix
  for  $i = 1$  to  $n$ 
    for  $j = 1$  to  $n$ 
       $t_{ij}^{(k)} = t_{ij}^{(k-1)} \vee (t_{ik}^{(k-1)} \wedge t_{kj}^{(k-1)})$ 
return  $T^{(n)}$ 

```

Time

$\Theta(n^3)$, but simpler operations than FLOYD-WARSHALL.

Johnson's algorithm

Idea

If the graph is sparse, it pays to run Dijkstra's algorithm once from each vertex.

If we use a Fibonacci heap for the priority queue, the running time is down to $O(V^2 \lg V + VE)$, which is better than FLOYD-WARSHALL's $\Theta(V^3)$ time if $E = o(V^2)$.

But Dijkstra's algorithm requires that all edge weights be nonnegative.

Donald Johnson figured out how to make an equivalent graph that *does* have all edge weights ≥ 0 .

Reweighting

Compute a new weight function \hat{w} such that

1. For all $u, v \in V$, p is a shortest path $u \rightsquigarrow v$ using w if and only if p is a shortest path $u \rightsquigarrow v$ using \hat{w} .
2. For all $(u, v) \in E$, $\hat{w}(u, v) \geq 0$.

Property (1) says that it suffices to find shortest paths with \hat{w} . Property (2) says we can do so by running Dijkstra's algorithm from each vertex.

How to come up with \hat{w} ?

Lemma shows it's easy to get property (1):

Chapter 24: Network Flows

Reading: Chapter 24

Lecture Notes for Chapter 24:

Maximum Flow

Chapter 24 overview

Network flow

[The third and fourth editions treat flow networks differently from the first two editions. The concept of net flow is gone, except that we do discuss net flow across a cut. Skew symmetry is also gone, as is implicit summation notation. The third and fourth editions count flows on edges directly. We find that although the mathematics is not quite as slick as in the first two editions, the approach in the newer editions matches intuition more closely, and therefore students tend to pick it up more quickly.]

Use a graph to model material that flows through conduits.

Each edge represents one conduit, and has a **capacity**, which is an upper bound on the **flow rate** = units/time.

Can think of edges as pipes of different sizes. But flows don't have to be of liquids. The textbook has an example where a flow is how many trucks per day can ship hockey pucks between cities.

Want to compute the maximum rate that material can be shipped from a designated **source** to a designated **sink**.

Flow networks

$G = (V, E)$ directed.

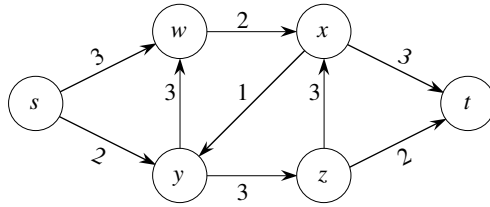
Each edge (u, v) has a **capacity** $c(u, v) \geq 0$.

If $(u, v) \notin E$, then $c(u, v) = 0$.

If $(u, v) \in E$, then reverse edge $(v, u) \notin E$. (Can work around this restriction.)

Source vertex s , **sink** vertex t , assume $s \rightsquigarrow v \rightsquigarrow t$ for all $v \in V$, so that each vertex lies on a path from source to sink.

Example: *[Edges are labeled with capacities.]*

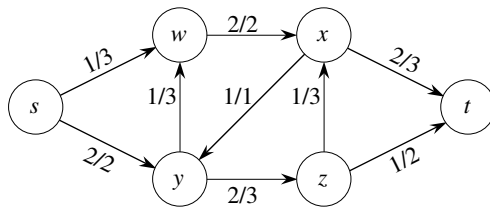
**Flow**

A function $f : V \times V \rightarrow \mathbb{R}$ satisfying

- **Capacity constraint:** For all $u, v \in V$, $0 \leq f(u, v) \leq c(u, v)$,
- **Flow conservation:** For all $u \in V - \{s, t\}$, $\underbrace{\sum_{v \in V} f(v, u)}_{\text{flow into } u} = \underbrace{\sum_{v \in V} f(u, v)}_{\text{flow out of } u}$.

$$\text{Equivalently, } \sum_{v \in V} f(u, v) - \sum_{v \in V} f(v, u) = 0.$$

[Add flows to previous example. Edges here are labeled as flow/capacity. Leave on board.]



- Note that all flows are \leq capacities.
- Verify flow conservation by adding up flows at a couple of vertices.
- Note that all flows = 0 is legitimate.

$$\begin{aligned} \text{Value of flow } f &= |f| \\ &= \sum_{v \in V} f(s, v) - \sum_{v \in V} f(v, s) \\ &= \text{flow out of source} - \text{flow into source} . \end{aligned}$$

In the example above, value of flow $f = |f| = 3$.

Antiparallel edges

Definition of flow network does not allow both (u, v) and (v, u) to be edges. These edges would be **antiparallel**.

What if really need antiparallel edges?

- Choose one of them, say (u, v) .
- Create a new vertex v' .
- Replace (u, v) by two new edges (u, v') and (v', v) , with $c(u, v') = c(v', v) = c(u, v)$.
- Get an equivalent flow network with no antiparallel edges.

Multiple sources and sinks

If you need multiple sources s_1, \dots, s_m , create a single **supersource** s with edges (s, s_i) for $i = 1, \dots, m$ and infinite capacity on each edge.

Same idea for multiple sinks: create a single **supersink** with an infinite-capacity edge from each sink to the supersink.

Now there are just one source and one sink.

Cuts

A **cut** (S, T) of flow network $G = (V, E)$ is a partition of V into S and $T = V - S$ such that $s \in S$ and $t \in T$.

- Similar to cut used in minimum spanning trees, except that here the graph is directed, and require $s \in S$ and $t \in T$.

For flow f , the **net flow** across cut (S, T) is

$$f(S, T) = \sum_{u \in S} \sum_{v \in T} f(u, v) - \sum_{u \in S} \sum_{v \in T} f(v, u) .$$

Capacity of cut (S, T) is

$$c(S, T) = \sum_{u \in S} \sum_{v \in T} c(u, v) .$$

A **minimum cut** of G is a cut whose capacity is minimum over all cuts of G .

Asymmetry between net flow across a cut and capacity of a cut: For capacity, count only capacities of edges going from S to T . Ignore edges going in the reverse direction. For net flow, count flow on all edges across the cut: flow on edges going from S to T minus flow on edges going from T to S .

In previous example, consider the cut $S = \{s, w, y\}$, $T = \{x, z, t\}$.

$$\begin{aligned} f(S, T) &= \underbrace{f(w, x) + f(y, z)}_{\text{from } S \text{ to } T} - \underbrace{f(x, y)}_{\text{from } T \text{ to } S} \\ &= 2 + 2 - 1 \\ &= 3 . \end{aligned}$$

$$\begin{aligned} c(S, T) &= \underbrace{c(w, x) + c(y, z)}_{\text{from } S \text{ to } T} \\ &= 2 + 3 \\ &= 5 . \end{aligned}$$

Now consider the cut $S = \{s, w, x, y\}$, $T = \{z, t\}$.

$$\begin{aligned} f(S, T) &= \underbrace{f(x, t) + f(y, z)}_{\text{from } S \text{ to } T} - \underbrace{f(z, x)}_{\text{from } T \text{ to } S} \\ &= 2 + 2 - 1 \\ &= 3 . \end{aligned}$$

$$\begin{aligned} c(S, T) &= \underbrace{c(x, t) + c(y, z)}_{\text{from } S \text{ to } T} \\ &= 3 + 3 \\ &= 6 . \end{aligned}$$

Same flow as previous cut, higher capacity.

Lemma

For any cut (S, T) , $f(S, T) = |f|$.

(Net flow across the cut equals value of the flow.)

[Leave on board.]

[This proof is much more involved than the proof in the first two editions. You might want to omit it, or just give the intuition that no matter where you cut the pipes in a network, you'll see the same flow volume coming out of the openings.]

Proof Rewrite flow conservation: for any $u \in V - \{s, t\}$,

$$\sum_{v \in V} f(u, v) - \sum_{v \in V} f(v, u) = 0.$$

Take definition of $|f|$ and add in left-hand side of above equation, summed over all vertices in $S - \{s\}$. Above equation applies to each vertex in $S - \{s\}$ (since $t \notin S$ and obviously $s \notin S - \{s\}$), so just adding in lots of 0s:

$$|f| = \sum_{v \in V} f(s, v) - \sum_{v \in V} f(v, s) + \sum_{u \in S - \{s\}} \left(\sum_{v \in V} f(u, v) - \sum_{v \in V} f(v, u) \right).$$

Expand right-hand summation and regroup terms:

$$\begin{aligned} |f| &= \sum_{v \in V} f(s, v) - \sum_{v \in V} f(v, s) + \sum_{u \in S - \{s\}} \sum_{v \in V} f(u, v) - \sum_{u \in S - \{s\}} \sum_{v \in V} f(v, u) \\ &= \sum_{v \in V} \left(f(s, v) + \sum_{u \in S - \{s\}} f(u, v) \right) - \sum_{v \in V} \left(f(v, s) + \sum_{u \in S - \{s\}} f(v, u) \right) \\ &= \sum_{v \in V} \sum_{u \in S} f(u, v) - \sum_{v \in V} \sum_{u \in S} f(v, u). \end{aligned}$$

Partition V into $S \cup T$ and split each summation over V into summations over S and T :

$$\begin{aligned} |f| &= \sum_{v \in S} \sum_{u \in S} f(u, v) + \sum_{v \in T} \sum_{u \in S} f(u, v) - \sum_{v \in S} \sum_{u \in S} f(v, u) - \sum_{v \in T} \sum_{u \in S} f(v, u) \\ &= \sum_{v \in T} \sum_{u \in S} f(u, v) - \sum_{v \in T} \sum_{u \in S} f(v, u) \\ &\quad + \left(\sum_{v \in S} \sum_{u \in S} f(u, v) - \sum_{v \in S} \sum_{u \in S} f(v, u) \right). \end{aligned}$$

Summations within parentheses are the same, since $f(x, y)$ appears once in each summation, for any $x, y \in V$. These summations cancel:

$$\begin{aligned} |f| &= \sum_{u \in S} \sum_{v \in T} f(u, v) - \sum_{u \in S} \sum_{v \in T} f(v, u) \\ &= f(S, T). \end{aligned}$$

■ (lemma)

Corollary

The value of any flow \leq capacity of any cut.

[Leave on board.]

Proof Let (S, T) be any cut, f be any flow.

$$\begin{aligned}
 |f| &= f(S, T) && \text{(lemma)} \\
 &= \sum_{u \in S} \sum_{v \in T} f(u, v) - \sum_{u \in S} \sum_{v \in T} f(v, u) && \text{(definition of } f(S, T)) \\
 &\leq \sum_{u \in S} \sum_{v \in T} f(u, v) && (f(v, u) \geq 0) \\
 &\leq \sum_{u \in S} \sum_{v \in T} c(u, v) && \text{(capacity constraint)} \\
 &= c(S, T) . && \text{(definition of } c(S, T)) \quad \blacksquare \text{ (corollary)}
 \end{aligned}$$

Therefore, maximum flow \leq capacity of minimum cut.

Will see a little later that this is in fact an equality.

The Ford-Fulkerson method

Residual network

Given a flow f in network $G = (V, E)$.

Consider a pair of vertices $u, v \in V$.

How much additional flow can be pushed directly from u to v ?

That's the **residual capacity**,

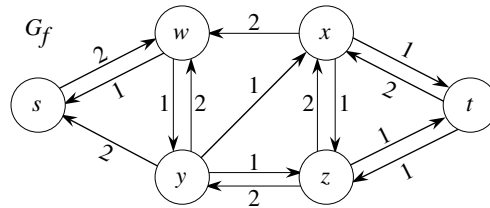
$$c_f(u, v) = \begin{cases} c(u, v) - f(u, v) & \text{if } (u, v) \in E , \\ f(v, u) & \text{if } (v, u) \in E , \\ 0 & \text{otherwise (i.e., } (u, v), (v, u) \notin E \text{)} . \end{cases}$$

The **residual network** is $G_f = (V, E_f)$, where

$$E_f = \{(u, v) \in V \times V : c_f(u, v) > 0\} .$$

Each edge of the residual network can admit a positive flow.

For our example:



Every edge $(u, v) \in E_f$ corresponds to an edge $(u, v) \in E$ or $(v, u) \in E$.

Therefore, $|E_f| \leq 2|E|$.

Residual network is similar to a flow network, except that it may contain antiparallel edges $((u, v)$ and $(v, u))$. Can define a flow in a residual network that satisfies the definition of a flow, but with respect to capacities c_f in G_f .

Given flows f in G and f' in G_f , define $(f \uparrow f')$, the **augmentation** of f by f' , as a function $V \times V \rightarrow \mathbb{R}$:

$$(f \uparrow f')(u, v) = \begin{cases} f(u, v) + f'(u, v) - f'(v, u) & \text{if } (u, v) \in E, \\ 0 & \text{otherwise} \end{cases}$$

for all $u, v \in V$.

Intuition: Increase the flow on (u, v) by $f'(u, v)$ but decrease it by $f'(v, u)$ because pushing flow on the reverse edge in the residual network decreases the flow in the original network. Also known as **cancellation**.

Lemma

Given a flow network G , a flow f in G , and the residual network G_f , let f' be a flow in G_f . Then $f \uparrow f'$ is a flow in G with value $|f \uparrow f'| = |f| + |f'|$.

[See textbook for proof. It has a lot of summations in it. Probably not worth writing on the board.]

Augmenting path

A simple path $s \rightsquigarrow t$ in G_f .

- Admits more flow along each edge.
- Like a sequence of pipes through which can squirt more flow from s to t .

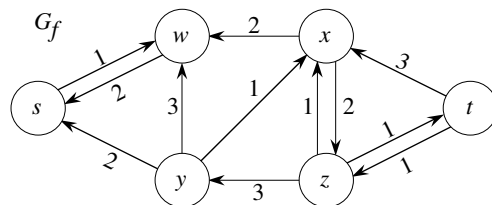
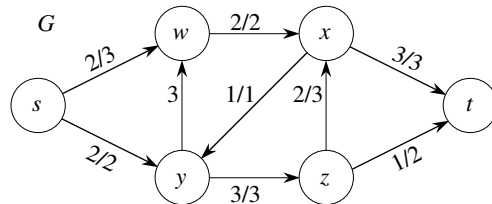
How much more flow can be pushed from s to t along augmenting path p ? That is the **residual capacity** of p :

$$c_f(p) = \min \{c_f(u, v) : (u, v) \text{ is on } p\}.$$

For our example, consider the augmenting path $p = \langle s, w, y, z, x, t \rangle$.

Minimum residual capacity is 1.

After pushing 1 additional unit along p : [Continue from G left on board from before. Edge (y, w) has $f(y, w) = 0$, which we omit, showing only $c(y, w) = 3$.]



Observe that G_f now has no augmenting path. Why? No edges cross the cut $(\{s, w\}, \{x, y, z, t\})$ in the forward direction in G_f . So no path can get from s to t .

Claim that the flow shown in G is a maximum flow.

Lemma

Given flow network G , flow f in G , residual network G_f . Let p be an augmenting path in G_f . Define $f_p : V \times V \rightarrow \mathbb{R}$:

$$f_p(u, v) = \begin{cases} c_f(p) & \text{if } (u, v) \text{ is on } p, \\ 0 & \text{otherwise.} \end{cases}$$

Then f_p is a flow in G_f with value $|f_p| = c_f(p) > 0$.

Corollary

Given flow network G , flow f in G , and an augmenting path p in G_f , define f_p as in lemma. Then $f \uparrow f_p$ is a flow in G with value $|f \uparrow f_p| = |f| + |f_p| > |f|$.

Theorem (Max-flow min-cut theorem)

The following are equivalent:

1. f is a maximum flow.
2. G_f has no augmenting path.
3. $|f| = c(S, T)$ for some cut (S, T) .

Proof

(1) \Rightarrow (2): Show the contrapositive: if G_f has an augmenting path, then f is not a maximum flow. If G_f has augmenting path p , then by the above corollary, $f \uparrow f_p$ is a flow in G with value $|f| + |f_p| > |f|$, so that f was not a maximum flow.

(2) \Rightarrow (3): Suppose G_f has no augmenting path. Define

$$S = \{v \in V : \text{there exists a path } s \rightsquigarrow v \text{ in } G_f\},$$

$$T = V - S.$$

Must have $t \in T$; otherwise there is an augmenting path.

Therefore, (S, T) is a cut.

Consider $u \in S$ and $v \in T$:

- If $(u, v) \in E$, must have $f(u, v) = c(u, v)$; otherwise, $(u, v) \in E_f \Rightarrow v \in S$.
- If $(v, u) \in E$, must have $f(v, u) = 0$; otherwise, $c_f(u, v) = f(v, u) > 0 \Rightarrow (u, v) \in E_f \Rightarrow v \in S$.
- If $(u, v), (v, u) \notin E$, must have $f(u, v) = f(v, u) = 0$.

Then,

$$\begin{aligned} f(S, T) &= \sum_{u \in S} \sum_{v \in T} f(u, v) - \sum_{v \in T} \sum_{u \in S} f(v, u) \\ &= \sum_{u \in S} \sum_{v \in T} c(u, v) - \sum_{v \in T} \sum_{u \in S} 0 \\ &= c(S, T). \end{aligned}$$

By lemma, $|f| = f(S, T) = c(S, T)$.

(3) \Rightarrow (1): An earlier corollary says that the value of any flow is \leq the capacity of any cut, so that $|f| \leq c(S, T)$.

Therefore, $|f| = c(S, T) \Rightarrow f$ is a max flow.

■ (theorem)

Ford-Fulkerson algorithm

Keep augmenting flow along an augmenting path until there is no augmenting path. Represent the flow attribute using the usual dot-notation, but on an edge: $(u, v).f$.

```

FORD-FULKERSON( $G, s, t$ )
  for each edge  $(u, v) \in G.E$ 
     $(u, v).f = 0$ 
  while there is an augmenting path  $p$  in  $G_f$ 
     $c_f(p) = \min \{c_f(u, v) : (u, v) \text{ is in } p\}$ 
    for each edge  $(u, v)$  in  $p$       // augment  $f$  by  $c_f(p)$ 
      if  $(u, v) \in G.E$ 
         $(u, v).f = (u, v).f + c_f(p)$ 
      else  $(v, u).f = (v, u).f - c_f(p)$ 
  return  $f$ 

```

Analysis

If capacities are all integer, then each augmenting path raises $|f|$ by ≥ 1 . If max flow is f^* , then need $\leq |f^*|$ iterations \Rightarrow time is $O(E |f^*|)$.

[Handwaving—see textbook for better explanation.]

Note that this running time is *not* polynomial in input size. It depends on $|f^*|$, which is not a function of $|V|$ and $|E|$.

If capacities are rational, can scale them to integers.

If irrational, FORD-FULKERSON might never terminate!

Edmonds-Karp algorithm

Do FORD-FULKERSON, but compute augmenting paths by BFS of G_f . Augmenting paths are shortest paths $s \rightsquigarrow t$ in G_f , with all edge weights = 1.

Edmonds-Karp runs in $O(VE^2)$ time.

To prove, need to look at distances to vertices in G_f .

Let $\delta_f(u, v)$ = shortest path distance u to v in G_f , with unit edge weights.

Lemma

For all $v \in V - \{s, t\}$, $\delta_f(s, v)$ increases monotonically with each flow augmentation.

Proof Suppose there exists $v \in V - \{s, t\}$ such that some flow augmentation causes $\delta_f(s, v)$ to decrease. Will derive a contradiction.

Let f be the flow before the first augmentation that causes a shortest-path distance to decrease, f' be the flow afterward.

Let v be a vertex with minimum $\delta_{f'}(s, v)$ whose distance was decreased by the augmentation, so that $\delta_{f'}(s, v) < \delta_f(s, v)$.

Let a shortest path s to v in $G_{f'}$ be $s \rightsquigarrow u \rightarrow v$, so that $(u, v) \in E_{f'}$ and $\delta_{f'}(s, v) = \delta_{f'}(s, u) + 1$. (Or $\delta_{f'}(s, u) = \delta_{f'}(s, v) - 1$.)

Since $\delta_{f'}(s, u) < \delta_{f'}(s, v)$ and how we chose v , we have $\delta_{f'}(s, u) \geq \delta_f(s, u)$.

Claim

$(u, v) \notin E_f$.

Proof of claim If $(u, v) \in E_f$, then

$$\begin{aligned} \delta_f(s, v) &\leq \delta_f(s, u) + 1 \quad (\text{triangle inequality}) \\ &\leq \delta_{f'}(s, u) + 1 \\ &= \delta_{f'}(s, v), \end{aligned}$$

contradicting $\delta_{f'}(s, v) < \delta_f(s, v)$.

■ (claim)

How can $(u, v) \notin E_f$ and $(u, v) \in E_{f'}$?

The augmentation must increase flow v to u .

Since Edmonds-Karp augments along shortest paths, the shortest path s to u in G_f has (v, u) as its last edge.

Therefore,

$$\begin{aligned} \delta_f(s, v) &= \delta_f(s, u) - 1 \\ &\leq \delta_{f'}(s, u) - 1 \\ &= \delta_{f'}(s, v) - 2, \end{aligned}$$

contradicting $\delta_{f'}(s, v) < \delta_f(s, v)$.

Therefore, v cannot exist.

■ (lemma)

Theorem

Edmonds-Karp performs $O(VE)$ augmentations.

Proof Suppose p is an augmenting path and $c_f(u, v) = c_f(p)$. Then call (u, v) a **critical** edge in G_f , and it disappears from the residual network after augmenting along p .

≥ 1 edge on any augmenting path is critical.

Will show that each of the $|E|$ edges can become critical $\leq |V|/2$ times.

Consider $u, v \in V$ such that either $(u, v) \in E$ or $(v, u) \in E$. Since augmenting paths are shortest paths, when (u, v) becomes critical the first time, $\delta_f(s, v) = \delta_f(s, u) + 1$.

Augment flow, so that (u, v) disappears from the residual network. This edge cannot reappear in the residual network until flow from u to v decreases, which happens only if (v, u) is on an augmenting path in $G_{f'}$: $\delta_{f'}(s, u) = \delta_{f'}(s, v) + 1$. (f' is flow when this occurs.)

By lemma, $\delta_f(s, v) \leq \delta_{f'}(s, v) \Rightarrow$

$$\begin{aligned} \delta_{f'}(s, u) &= \delta_{f'}(s, v) + 1 \\ &\geq \delta_f(s, v) + 1 \\ &= \delta_f(s, u) + 2. \end{aligned}$$

Therefore, from the time (u, v) becomes critical to the next time, distance of u from s increases by ≥ 2 . Initially, distance to u is ≥ 0 , and augmenting path can't have s, u , and t as intermediate vertices.

Therefore, until u becomes unreachable from source, its distance is $\leq |V| - 2$
 \Rightarrow after (u, v) becomes critical the first time, it can become critical
 $\leq (|V| - 2)/2 = |V|/2 - 1$ times more
 $\Rightarrow (u, v)$ can become critical $\leq |V|/2$ times.

Since $O(E)$ pairs of vertices can have an edge between them in residual network, total # of critical edges during execution of Edmonds-Karp is $O(VE)$. Since each augmenting path has ≥ 1 critical edge, have $O(VE)$ augmentations. ■ (theorem)

Use BFS to find each augmenting path in $O(E)$ time $\Rightarrow O(VE^2)$ time.

Can get better bounds. [Push-relabel algorithms in the first three editions of the textbook give $O(V^3)$. The two sections on push-relabel algorithm were dropped from the fourth edition but are available from the MIT Press website for the book.]

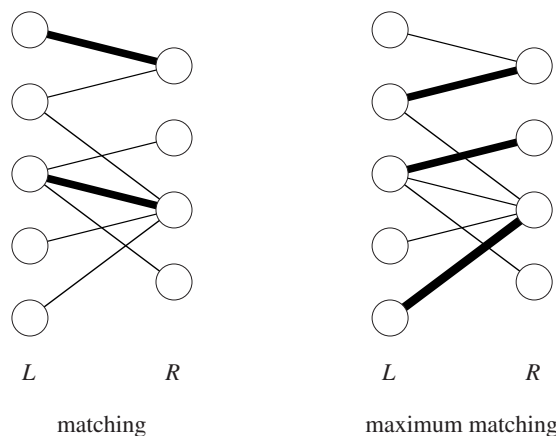
Maximum bipartite matching

Example of a problem that can be solved by turning it into a flow problem.

$G = (V, E)$ (undirected) is **bipartite** if there is a partition of the vertices $V = L \cup R$ such that all edges in E go between L and R .

A **matching** is a subset of edges $M \subseteq E$ such that for all $v \in V$, ≤ 1 edge of M is incident on v . (Vertex v is **matched** if an edge of M is incident on it; otherwise **unmatched**).

Maximum matching: a matching of maximum cardinality. (M is a maximum matching if $|M| \geq |M'|$ for all matchings M' .)



[Edges in matchings are drawn with heavy lines.]

Problem

Given a bipartite graph (with the partition), find a maximum matching.

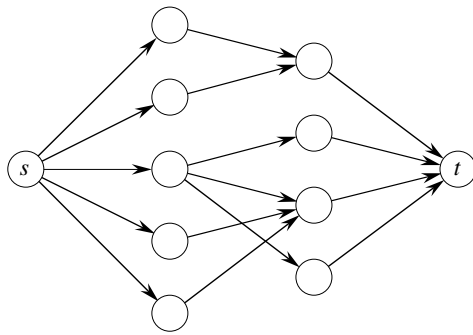
Application

Matching planes to routes.

- L = set of planes.
- R = set of routes.
- $(u, v) \in E$ if plane u can fly route v .
- Want maximum # of routes to be served by planes.

Given G , define flow network $G' = (V', E')$.

- $V' = V \cup \{s, t\}$.
- $E' = \{(s, u) : u \in L\} \cup \{(u, v) : (u, v) \in E\} \cup \{(v, t) : v \in R\}$.
- $c(u, v) = 1$ for all $(u, v) \in E'$.



Each vertex in V has ≥ 1 incident edge $\Rightarrow |E| \geq |V|/2$.

Therefore, $|E| \leq |E'| = |E| + |V| \leq 3|E|$.

Therefore, $|E'| = \Theta(E)$.

Find a max flow in G' . Textbook shows that it will have integer values for all (u, v) .

Use edges (u, v) such that $u \in L$ and $v \in R$ that carry flow of 1 in matching.

Textbook proves that this method produces a maximum matching.

[The next chapter, Chapter 25, has a better algorithm (Hopcroft-Karp) to find a maximum matching, as well as other algorithms based on bipartite matchings.]