

Those who cannot remember the past are doomed to repeat it.

— George Santayana, *The Life of Reason, Book I: Introduction and Reason in Common Sense* (1905)

The 1950s were not good years for mathematical research. We had a very interesting gentleman in Washington named Wilson. He was secretary of Defense, and he actually had a pathological fear and hatred of the word 'research'. I'm not using the term lightly; I'm using it precisely. His face would suffuse, he would turn red, and he would get violent if people used the term 'research' in his presence. You can imagine how he felt, then, about the term 'mathematical'. The RAND Corporation was employed by the Air Force, and the Air Force had Wilson as its boss, essentially. Hence, I felt I had to do something to shield Wilson and the Air Force from the fact that I was really doing mathematics inside the RAND Corporation. What title, what name, could I choose?

— Richard Bellman, on the origin of his term 'dynamic programming' (1984)

If we all listened to the professor, we may be all looking for professor jobs.

— Pittsburgh Steelers' head coach Bill Cowher, responding to David Romer's dynamic-programming analysis of football strategy (2003)

5 Dynamic Programming

5.1 Fibonacci Numbers

5.1.1 Recursive Definitions Are Recursive Algorithms

The Fibonacci numbers F_n , named after Leonardo Fibonacci Pisano¹, the mathematician who popularized 'algorism' in Europe in the 13th century, are defined as follows: $F_0 = 0$, $F_1 = 1$, and $F_n = F_{n-1} + F_{n-2}$ for all $n \geq 2$. The recursive definition of Fibonacci numbers immediately gives us a recursive algorithm for computing them:

```

REC FIBO( $n$ ):
  if ( $n < 2$ )
    return  $n$ 
  else
    return REC FIBO( $n - 1$ ) + REC FIBO( $n - 2$ )

```

How long does this algorithm take? Except for the recursive calls, the entire algorithm requires only a constant number of steps: one comparison and possibly one addition. If $T(n)$ represents the number of recursive calls to REC FIBO, we have the recurrence

$$T(0) = 1, \quad T(1) = 1, \quad T(n) = T(n-1) + T(n-2) + 1.$$

This looks an awful lot like the recurrence for Fibonacci numbers! The annihilator method gives us an asymptotic bound of $\Theta(\phi^n)$, where $\phi = (\sqrt{5} + 1)/2 \approx 1.61803398875$, the so-called *golden ratio*, is the largest root of the polynomial $r^2 - r - 1$. But it's fairly easy to prove (hint, hint) the exact solution $T(n) = 2F_{n+1} - 1$. In other words, computing F_n using this algorithm takes more than twice as many steps as just counting to F_n !

Another way to see this is that the REC FIBO is building a big binary tree of additions, with nothing but zeros and ones at the leaves. Since the eventual output is F_n , our algorithm must

¹literally, "Leonardo, son of Bonacci, of Pisa"

call `RECFIBO(1)` (which returns 1) exactly F_n times. A quick inductive argument implies that `RECFIBO(0)` is called exactly F_{n-1} times. Thus, the recursion tree has $F_n + F_{n-1} = F_{n+1}$ leaves, and therefore, because it's a full binary tree, it must have $2F_{n+1} - 1$ nodes.

5.1.2 Memo(r)ization: Remember Everything

The obvious reason for the recursive algorithm's lack of speed is that it computes the same Fibonacci numbers over and over and over. A single call to `RECFIBO(n)` results in one recursive call to `RECFIBO(n-1)`, two recursive calls to `RECFIBO(n-2)`, three recursive calls to `RECFIBO(n-3)`, five recursive calls to `RECFIBO(n-4)`, and in general F_{k-1} recursive calls to `RECFIBO(n-k)` for any integer $0 \leq k < n$. Each call is recomputing some Fibonacci number from scratch.

We can speed up our recursive algorithm considerably just by writing down the results of our recursive calls and looking them up again if we need them later. This process was dubbed *memoization* by Richard Michie in the late 1960s.²

```

MEMFIBO(n):
  if (n < 2)
    return n
  else
    if F[n] is undefined
      F[n] ← MEMFIBO(n-1) + MEMFIBO(n-2)
    return F[n]

```

Memoization clearly decreases the running time of the algorithm, but by how much? If we actually trace through the recursive calls made by `MEMFIBO`, we find that the array $F[\]$ is filled from the bottom up: first $F[2]$, then $F[3]$, and so on, up to $F[n]$. This pattern can be verified by induction: Each entry $F[i]$ is filled only after its predecessor $F[i-1]$. If we ignore the time spent in recursive calls, it requires only constant time to evaluate the recurrence for each Fibonacci number F_i . But by design, the recurrence for F_i is evaluated only once for each index i ! We conclude that `MEMFIBO` performs only $O(n)$ additions, an *exponential* improvement over the naïve recursive algorithm!

5.1.3 Dynamic Programming: Fill Deliberately

But once we see how the array $F[\]$ is filled, we can replace the recursion with a simple loop that intentionally fills the array in order, instead of relying on the complicated recursion to do it for us 'accidentally'.

```

ITERFIBO(n):
  F[0] ← 0
  F[1] ← 1
  for i ← 2 to n
    F[i] ← F[i-1] + F[i-2]
  return F[n]

```

Now the time analysis is immediate: `ITERFIBO` clearly uses $O(n)$ *additions* and stores $O(n)$ *integers*.

This gives us our first explicit *dynamic programming* algorithm. The dynamic programming paradigm was developed by Richard Bellman in the mid-1950s, while working at the RAND

²"My name is Elmer J. Fudd, millionaire. I own a mansion and a yacht."

Corporation. Bellman deliberately chose the name ‘dynamic programming’ to hide the mathematical character of his work from his military bosses, who were actively hostile toward anything resembling mathematical research. Here, the word ‘programming’ does not refer to writing code, but rather to the older sense of *planning* or *scheduling*, typically by filling in a table. For example, sports programs and theater programs are schedules of important events (with ads); television programming involves filling each available time slot with a show (and ads); degree programs are schedules of classes to be taken (with ads). The Air Force funded Bellman and others to develop methods for constructing training and logistics schedules, or as they called them, ‘programs’. The word ‘dynamic’ is meant to suggest that the table is filled in over time, rather than all at once (as in ‘linear programming’, which we will see later in the semester).³

5.1.4 Don’t Remember Everything After All

In many dynamic programming algorithms, it is not necessary to retain *all* intermediate results through the entire computation. For example, we can significantly reduce the space requirements of our algorithm ITERFIBO by maintaining only the two newest elements of the array:

```
ITERFIBO2(n):
  prev ← 1
  curr ← 0
  for i ← 1 to n
    next ← curr + prev
    prev ← curr
    curr ← next
  return curr
```

(This algorithm uses the non-standard but perfectly consistent base case $F_{-1} = 1$ so that ITERFIBO2(0) returns the correct value 0.)

5.1.5 Faster! Faster!

Even this algorithm can be improved further, using the following wonderful fact:

$$\begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} y \\ x + y \end{bmatrix}$$

In other words, multiplying a two-dimensional vector by the matrix $\begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix}$ does exactly the same thing as one iteration of the inner loop of ITERFIBO2. This might lead us to believe that multiplying by the matrix n times is the same as iterating the loop n times:

$$\begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix}^n \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} F_{n-1} \\ F_n \end{bmatrix}.$$

A quick inductive argument proves this fact. So if we want the n th Fibonacci number, we just have to compute the n th power of the matrix $\begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix}$. If we use repeated squaring, computing the n th power of something requires only $O(\log n)$ multiplications. In this case, that means $O(\log n)$ 2×2 matrix multiplications, each of which reduces to a constant number of integer multiplications and additions. Thus, we can compute F_n in only **$O(\log n)$ integer arithmetic operations**.

This is an exponential speedup over the standard iterative algorithm, which was already an exponential speedup over our original recursive algorithm. Right?

³“I thought dynamic programming was a good name. It was something not even a Congressman could object to. So I used it as an umbrella for my activities.”

5.1.6 Whoa! Not so fast!

Well, not exactly. Fibonacci numbers grow exponentially fast. The n th Fibonacci number is approximately $n \log_{10} \phi \approx n/5$ decimal digits long, or $n \log_2 \phi \approx 2n/3$ bits. So we can't possibly compute F_n in logarithmic time — we need $\Omega(n)$ time just to write down the answer!

The way out of this apparent paradox is to observe that **we can't perform arbitrary-precision arithmetic in constant time**. Let $M(n)$ denote the time required to multiply two n -digit numbers. The matrix-based algorithm's actual running time obeys the recurrence $T(n) = T(\lfloor n/2 \rfloor) + M(n)$, which solves to $T(n) = O(M(n))$ using recursion trees. The fastest known multiplication algorithm runs in time $O(n \log n 2^{O(\log^* n)})$, so that is also the running time of the fastest algorithm known to compute Fibonacci numbers.

Is this algorithm slower than our initial “linear-time” iterative algorithm? No! Addition isn't free, either. Adding two n -digit numbers takes $O(n)$ time, so the running time of the iterative algorithm is $O(n^2)$. (Do you see why?) The matrix-squaring algorithm really is faster than the iterative addition algorithm, but not exponentially faster.

In the original recursive algorithm, the extra cost of arbitrary-precision arithmetic is overwhelmed by the huge number of recursive calls. The correct recurrence is $T(n) = T(n-1) + T(n-2) + O(n)$, for which the annihilator method still implies the solution $T(n) = O(\phi^n)$.

5.2 Longest Increasing Subsequence

In a previous lecture, we developed a recursive algorithm to find the length of the longest increasing subsequence of a given sequence of numbers. Given an array $A[1..n]$, the length of the longest increasing subsequence is computed by the function call $\text{LISBIGGER}(-\infty, A[1..n])$, where LISBIGGER is the following recursive algorithm:

```

LISBIGGER(prev,  $A[1..n]$ ):
  if  $n = 0$ 
    return 0
  else
     $max \leftarrow \text{LISBIGGER}(prev, A[2..n])$ 
    if  $A[1] > prev$ 
       $L \leftarrow 1 + \text{LISBIGGER}(A[1], A[2..n])$ 
      if  $L > max$ 
         $max \leftarrow L$ 
    return  $max$ 

```

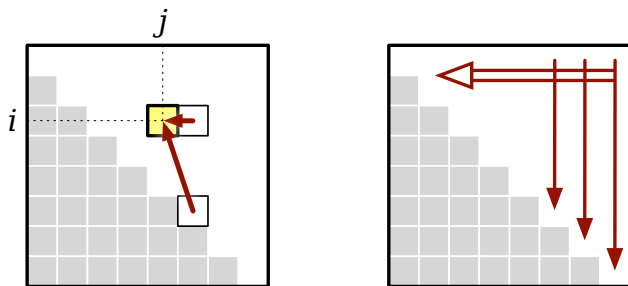
We can simplify our notation slightly with two simple observations. First, the input variable $prev$ is always either $-\infty$ or an element of the input array. Second, the second argument of LISBIGGER is always a *suffix* of the original input array. If we add a new sentinel value $A[0] = -\infty$ to the input array, we can identify any recursive subproblem with two array indices.

Thus, we can rewrite the recursive algorithm as follows. Add the sentinel value $A[0] = -\infty$. Let $LIS(i, j)$ denote the length of the longest increasing subsequence of $A[j..n]$ with all elements larger than $A[i]$. Our goal is to compute $LIS(0, 1)$. For all $i < j$, we have

$$LIS(i, j) = \begin{cases} 0 & \text{if } j > n \\ LIS(i, j+1) & \text{if } A[i] \geq A[j] \\ \max\{LIS(i, j+1), 1 + LIS(j, j+1)\} & \text{otherwise} \end{cases}$$

Because each recursive subproblem can be identified by two indices i and j , we can store the intermediate values in a two-dimensional array $LIS[0..n, 1..n]$.⁴ Since there are $O(n^2)$ entries in the table, our memoized algorithm uses $O(n^2)$ space. Each entry in the table can be computed in $O(1)$ time once we know its predecessors, so our memoized algorithm runs in $O(n^2)$ time.

It's not immediately clear what order the recursive algorithm fills the rest of the table; all we can tell from the recurrence is that each entry $LIS[i, j]$ is filled in *after* the entries $LIS[i, j + 1]$ and $LIS[j, j + 1]$ in the next columns. But just this partial information is enough to give us an explicit evaluation order. If we fill in our table one column at a time, from right to left, then whenever we reach an entry in the table, the entries it depends on are already available.



Dependencies in the memoization table for longest increasing subsequence, and a legal evaluation order

Finally, putting everything together, we obtain the following dynamic programming algorithm:

```

LIS(A[1..n]):
  A[0] ← -∞           ⟨⟨Add a sentinel⟩⟩
  for i ← 0 to n      ⟨⟨Base cases⟩⟩
    LIS[i, n + 1] ← 0
  for j ← n downto 1
    for i ← 0 to j - 1
      if A[i] ≥ A[j]
        LIS[i, j] ← LIS[i, j + 1]
      else
        LIS[i, j] ← max{LIS[i, j + 1], 1 + LIS[j, j + 1]}
  return LIS[0, 1]
    
```

As expected, the algorithm clearly uses $O(n^2)$ time and space. However, we can reduce the space to $O(n)$ by only maintaining the two most recent columns of the table, $LIS[·, j]$ and $LIS[·, j + 1]$.⁵

This is not the only recursive strategy we could use for computing longest increasing subsequences efficiently. Here is another recurrence that gives us the $O(n)$ space bound for free. Let $LIS'(i)$ denote the length of the longest increasing subsequence of $A[i..n]$ that starts with $A[i]$. Our goal is to compute $LIS'(0) - 1$; we subtract 1 to ignore the sentinel value $-\infty$. To define $LIS'(i)$ recursively, we only need to specify the *second* element in subsequence; the Recursion Fairy will do the rest.

$$LIS'(i) = 1 + \max \{ LIS'(j) \mid j > i \text{ and } A[j] > A[i] \}$$

Here, I'm assuming that $\max \emptyset = 0$, so that the base case is $L'(n) = 1$ falls out of the recurrence automatically. Memoizing this recurrence requires only $O(n)$ space, and the resulting algorithm

⁴In fact, we only need half of this array, because we always have $i < j$. But even if we cared about constant factors in this class (we don't), this would be the wrong time to worry about them. The first order of business is to find an algorithm that actually works; once we have that, then we can think about optimizing it.

⁵See, I told you not to worry about constant factors yet!

runs in $O(n^2)$ time. To transform this memoized recurrence into a dynamic programming algorithm, we only need to guarantee that $LIS'(j)$ is computed before $LIS'(i)$ whenever $i < j$.

```

LIS2(A[1..n]):
  A[0] = -∞           ⟨⟨Add a sentinel⟩⟩
  for i ← n downto 0
    LIS'[i] ← 1
    for j ← i + 1 to n
      if A[j] > A[i] and 1 + LIS'[j] > LIS'[i]
        LIS'[i] ← 1 + LIS'[j]
  return LIS'[0] - 1   ⟨⟨Don't count the sentinel⟩⟩

```

5.3 The Pattern: Smart Recursion

In a nutshell, dynamic programming is *recursion without repetition*. Dynamic programming algorithms store the solutions of intermediate subproblems, often *but not always* in some kind of array or table. Many algorithms students make the mistake of focusing on the table (because tables are easy and familiar) instead of the *much* more important (and difficult) task of finding a correct recurrence. As long as we memoize the correct recurrence, an explicit table isn't really necessary, but if the recursion is incorrect, nothing works.

**Dynamic programming is *not* about filling in tables.
It's about smart recursion!**

Dynamic programming algorithms are almost always developed in two distinct stages.

1. **Formulate the problem recursively.** Write down a recursive formula or algorithm for the whole problem in terms of the answers to smaller subproblems. This is the hard part. It generally helps to think in terms of a recursive definition of the object you're trying to construct. A complete recursive formulation has two parts:
 - (a) Describe the precise function you want to evaluate, in coherent English. Without this specification, it is impossible, even in principle, to determine whether your solution is correct.
 - (b) Give a formal recursive definition of that function.
2. **Build solutions to your recurrence from the bottom up.** Write an algorithm that starts with the base cases of your recurrence and works its way up to the final solution, by considering intermediate subproblems in the correct order. This stage can be broken down into several smaller, relatively mechanical steps:
 - (a) **Identify the subproblems.** What are all the different ways can your recursive algorithm call itself, starting with some initial input? For example, the argument to RECFIBO is always an integer between 0 and n .
 - (b) **Analyze space and running time.** The number of possible distinct subproblems determines the space complexity of your memoized algorithm. To compute the time complexity, add up the running times of all possible subproblems, *ignoring the recursive calls*. For example, if we already know F_{i-1} and F_{i-2} , we can compute F_i in $O(1)$ time, so computing the first n Fibonacci numbers takes $O(n)$ time.

- (c) **Choose a data structure to memoize intermediate results.** For most problems, each recursive subproblem can be identified by a few integers, so you can use a multidimensional array. For some problems, however, a more complicated data structure is required.
- (d) **Identify dependencies between subproblems.** Except for the base cases, every recursive subproblem depends on other subproblems—which ones? Draw a picture of your data structure, pick a generic element, and draw arrows from each of the other elements it depends on. Then formalize your picture.
- (e) **Find a good evaluation order.** Order the subproblems so that each subproblem comes *after* the subproblems it depends on. Typically, this means you should consider the base cases first, then the subproblems that depends only on base cases, and so on. More formally, the dependencies you identified in the previous step define a partial order over the subproblems; in this step, you need to find a linear extension of that partial order. ***Be careful!***
- (f) **Write down the algorithm.** You know what order to consider the subproblems, and you know how to solve each subproblem. So do that! If your data structure is an array, this usually means writing a few nested for-loops around your original recurrence. ***You don't need to do this on homework or exams.***

Of course, you have to prove that each of these steps is correct. If your recurrence is wrong, or if you try to build up answers in the wrong order, your algorithm won't work!

5.4 Warning: Greed is Stupid

If we're very very very very lucky, we can bypass all the recurrences and tables and so forth, and solve the problem using a *greedy* algorithm. The general greedy strategy is look for the best first step, take it, and then continue. While this approach seems very natural, it almost never works; optimization problems that can be solved correctly by a greedy algorithm are *very* rare. Nevertheless, for many problems that should be solved by dynamic programming, many students' first intuition is to apply a greedy strategy.

For example, a greedy algorithm for the longest increasing subsequence problem might look for the smallest element of the input array, accept that element as the start of the target subsequence, and then recursively look for the longest increasing subsequence to the right of that element. If this sounds like a stupid hack to you, pat yourself on the back. It isn't even *close* to the correct solution.

Everyone should tattoo the following sentence on the back of their hands, right under all the rules about logarithms and big-Oh notation:

**Greedy algorithms never work!
Use dynamic programming instead!**

What, never?

No, never!

What, *never*?

Well. . . hardly ever!⁶

A different lecture note describes the effort required to prove that greedy algorithms are correct, in the rare instances when they are. **You will not receive any credit for any greedy algorithm for any problem in this class without a formal proof of correctness.** We'll push through the formal proofs for several greedy algorithms later in the semester.

5.5 Edit Distance

The *edit distance* between two words is the minimum number of letter insertions, letter deletions, and letter substitutions required to transform one word into another. For example, the edit distance between **FOOD** and **MONEY** is at most four:

FOOD → MO**O**D → MON**A**D → MON**E**D → M**O**NEY

This distance function was independently proposed by Vladimir Levenshtein in 1964 (working on coding theory), T. K. Vintsyuk in 1968 (working on speech recognition), and Stanislaw Ulam in 1972 (working with biological sequences). For this reason, edit distance is sometimes called *Levenstein distance* or *Ulam distance*.

A good way to display this editing process is to place the words one above the other, with a gap in the first word for every insertion and a gap in the second word for every deletion. Columns with two *different* characters correspond to substitutions. In this representation, the number of editing steps is just the number of columns that do not contain the same character twice.

```

F O O D
M O N E Y

```

It's fairly obvious that you can't get from **FOOD** to **MONEY** in three steps, so their edit distance is exactly four. Unfortunately, this is not so easy in general. Here's a longer example, showing that the distance between **ALGORITHM** and **ALTRUISTIC** is at most six. Is this optimal?

```

A L G O R I T H M
A L T R U I S T I C

```

To develop a dynamic programming algorithm to compute the edit distance between two strings, we first need to develop a recursive definition. Our gap representation for edit sequences has a crucial “optimal substructure” property. Suppose we have the gap representation for the shortest edit sequence for two strings. **If we remove the last column, the remaining columns must represent the shortest edit sequence for the remaining substrings.** We can easily prove this by contradiction. If the substrings had a shorter edit sequence, we could just glue the last column back on and get a shorter edit sequence for the original strings. Once we figure out what should go in the last column, the Recursion Fairy will magically give us the rest of the optimal gap representation.

So let's recursively define the edit distance between two strings $A[1..m]$ and $B[1..n]$, which we denote by $Edit(A[1..m], B[1..n])$. If neither string is empty, there are three possibilities for the last column in the shortest edit sequence:

- **Insertion:** The last entry in the bottom row is empty. In this case, the edit distance is equal to $Edit(A[1..m-1], B[1..n]) + 1$. The +1 is the cost of the final insertion, and the recursive expression gives the minimum cost for the other columns.

⁶Greedy methods hardly ever work! So give three cheers, and one cheer more, for the careful Captain of the *Pinafore*! Then give three cheers, and one cheer more, for the Captain of the *Pinafore*!

- **Deletion:** The last entry in the top row is empty. In this case, the edit distance is equal to $Edit(A[1..m], B[1..n-1]) + 1$. The +1 is the cost of the final deletion, and the recursive expression gives the minimum cost for the other columns.
- **Substitution:** Both rows have characters in the last column. If the characters are the same, the substitution is free, so the edit distance is equal to $Edit(A[1..m-1], B[1..n-1])$. If the characters are different, then the edit distance is equal to $Edit(A[1..m-1], B[1..n-1]) + 1$.

The edit distance between A and B is the smallest of these three possibilities:⁷

$$Edit(A[1..m], B[1..n]) = \min \left\{ \begin{array}{l} Edit(A[1..m-1], B[1..n]) + 1 \\ Edit(A[1..m], B[1..n-1]) + 1 \\ Edit(A[1..m-1], B[1..n-1]) + [A[m] \neq B[n]] \end{array} \right\}$$

This recurrence has two easy base cases. The only way to convert the empty string into a string of n characters is by performing n insertions. Similarly, the only way to convert a string of m characters into the empty string is with m deletions. Thus, if ε denotes the empty string, we have

$$Edit(A[1..m], \varepsilon) = m \quad \text{and} \quad Edit(\varepsilon, B[1..n]) = n.$$

Both of these expressions imply the trivial base case $Edit(\varepsilon, \varepsilon) = 0$.

Now notice that the arguments to our recursive subproblems are always **prefixes** of the original strings A and B . We can simplify our notation by using the lengths of the prefixes, instead of the prefixes themselves, as the arguments to our recursive function.

Let $Edit(i, j)$ denote the edit distance between the prefixes $A[1..i]$ and $B[1..j]$.

This function satisfies the following recurrence:

$$Edit(i, j) = \begin{cases} i & \text{if } j = 0 \\ j & \text{if } i = 0 \\ \min \left\{ \begin{array}{l} Edit(i-1, j) + 1, \\ Edit(i, j-1) + 1, \\ Edit(i-1, j-1) + [A[i] \neq B[j]] \end{array} \right\} & \text{otherwise} \end{cases}$$

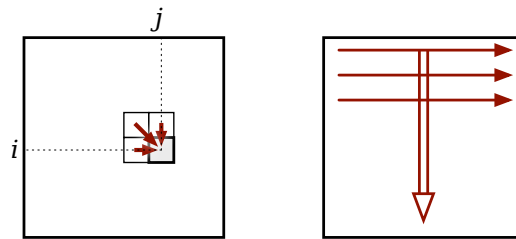
The edit distance between the original strings A and B is just $Edit(m, n)$. This recurrence translates directly into a recursive algorithm; the precise running time is not obvious, but it's clearly exponential in m and n . **Fortunately, we don't care about the precise running time of the recursive algorithm.** The recursive running time wouldn't tell us anything about our eventual dynamic programming algorithm, so we're just not going to bother computing it.⁸

Because each recursive subproblem can be identified by two indices i and j , we can memoize intermediate values in a two-dimensional array $Edit[0..m, 0..n]$. Note that the index ranges

⁷Once again, I'm using Iverson's bracket notation $[P]$ to denote the *indicator variable* for the logical proposition P , which has value 1 if P is true and 0 if P is false.

⁸The running time of the pure recursive algorithm obeys the recurrence $T(m, n) = O(1) + T(m, n-1) + T(m-1, n) + T(n-1, m-1)$. I don't know how to solve this recurrence exactly for all possible values of m and n . However, the function $T'(N) := \max_{n+m=N} T(n, m)$ satisfies the simpler recurrence $T'(N) = O(1) + 2T'(N-1) + T'(N-2)$, which implies the upper bound $T(m, n) \leq T'(n+m) = O((1 + \sqrt{2})^{n+m})$ by the annihilator method.

start at zero to accommodate the base cases. Since there are $\Theta(mn)$ entries in the table, our memoized algorithm uses $\Theta(mn)$ space. Since each entry in the table can be computed in $\Theta(1)$ time once we know its predecessors, our memoized algorithm runs in $\Theta(mn)$ time.



Dependencies in the memoization table for edit distance, and a legal evaluation order

Each entry $Edit[i, j]$ depends only on its three neighboring entries $Edit[i - 1, j]$, $Edit[i, j - 1]$, and $Edit[i - 1, j - 1]$. If we fill in our table in the standard row-major order—row by row from top down, each row from left to right—then whenever we reach an entry in the table, the entries it depends on are already available. Putting everything together, we obtain the following dynamic programming algorithm:

```

EDITDISTANCE( $A[1..m], B[1..n]$ ):
  for  $j \leftarrow 1$  to  $n$ 
     $Edit[0, j] \leftarrow j$ 
  for  $i \leftarrow 1$  to  $m$ 
     $Edit[i, 0] \leftarrow i$ 
    for  $j \leftarrow 1$  to  $n$ 
      if  $A[i] = B[j]$ 
         $Edit[i, j] \leftarrow \min \{Edit[i - 1, j] + 1, Edit[i, j - 1] + 1, Edit[i - 1, j - 1]\}$ 
      else
         $Edit[i, j] \leftarrow \min \{Edit[i - 1, j] + 1, Edit[i, j - 1] + 1, Edit[i - 1, j - 1] + 1\}$ 
  return  $Edit[m, n]$ 

```

This dynamic programming algorithm is most commonly attributed to Robert Wagner and Michael Fischer, who described the algorithm in 1974. However, in full compliance with Stigler's Law of Eponymy, either identical or more general algorithms were independently discovered by Taras Vintsyuk in 1968, V. M. Velichko and N. G. Zagoruyko in 1970, David Sankoff in 1972, Peter Sellers⁹ in 1974, and almost certainly several others.¹⁰ Interestingly, none of these papers cite either Levenshtein or Ulam.

The resulting table for **ALGORITHM** \rightarrow **ALTRUISTIC** is shown on the next page. Bold numbers indicate places where characters in the two strings are equal. The arrows represent the predecessor(s) that actually define each entry. Each direction of arrow corresponds to a different edit operation: horizontal=deletion, vertical=insertion, and diagonal=substitution. Bold diagonal arrows indicate “free” substitutions of a letter for itself. Any path of arrows from the top

⁹“Gentlemen! You can't fight in here! This is the War Room!” Okay, not *that* Peter Sellers.

¹⁰The Vintsyuk–Velichko–Zagoruyko–Sankoff–Sellers–Wagner–Fischer edit-distance algorithm is occasionally also attributed to Saul Needleman and Christian Wunsch in 1970, but this attribution is incorrect. “The Needleman–Wunsch algorithm” more commonly refers to the standard dynamic programming algorithm for computing the longest common subsequence of two strings (or equivalently, the edit distance where only insertions and deletions are permitted) in $O(mn)$ time, but this attribution is *also* incorrect!! In fact, Needleman and Wunsch's algorithm computes (weighted) longest common subsequences (possibly with gap costs) in $O(m^2n^2)$ time, using a different recurrence. Sankoff explicitly describes his $O(mn)$ -time algorithm as an improvement of Needleman and Wunsch's.

left corner to the bottom right corner of this table represents an optimal edit sequence between the two strings. (There can be many such paths.) Moreover, since we can compute these arrows in a post-processing phase from the values stored in the table, we can reconstruct the actual optimal editing sequence in $O(n + m)$ additional time.

		A	L	G	O	R	I	T	H	M
	0	→1	→2	→3	→4	→5	→6	→7	→8	→9
A	1	↓	↘	↘	↘	↘	↘	↘	↘	↘
L	2	1	↓	↘	↘	↘	↘	↘	↘	↘
T	3	2	1	↘	↘	↘	↘	↘	↘	↘
R	4	3	2	2	2	↓	↘	↘	↘	↘
U	5	4	3	3	3	3	3	↘	↘	↘
I	6	5	4	4	4	4	3	↘	↘	↘
S	7	6	5	5	5	5	4	4	5	6
T	8	7	6	6	6	6	5	4	5	6
I	9	8	7	7	7	7	6	5	5	6
C	10	9	8	8	8	8	7	6	6	6

The memoization table for $Edit(ALGORITHM, ALTRUISTIC)$

The edit distance between **ALGORITHM** and **ALTRUISTIC** is indeed six. There are three paths through this table from the top left to the bottom right, so there are three optimal edit sequences:

A L G O R I T H M
A L T R U I S T I C

A L G O R I T H M
A L T R U I S T I C

A L G O R I T H M
A L T R U I S T I C

5.6 More Examples

In the previous note on backtracking algorithms, we saw two other examples of recursive algorithms that we can significantly speed up via dynamic programming.

5.6.1 Subset Sum

Recall that the *Subset Sum* problem asks, given a set X of positive integers (represented as an array $X[1..n]$) and an integer T , whether any subset of X sums to T . In that lecture, we developed a recursive algorithm which can be reformulated as follows. Fix the original input array $X[1..n]$ and the original target sum T , and define the boolean function

$$SS(i, t) = \text{some subset of } X[i..n] \text{ sums to } t.$$

Our goal is to compute $S(1, T)$, using the recurrence

$$SS(i, t) = \begin{cases} \text{TRUE} & \text{if } t = 0, \\ \text{FALSE} & \text{if } t < 0 \text{ or } i > n, \\ SS(i + 1, t) \vee SS(i + 1, t - X[i]) & \text{otherwise.} \end{cases}$$

There are only nT possible values for the input parameters that lead to the interesting case of this recurrence, and we can memoize all such values in an $n \times T$ array. If $S(i + 1, t)$ and $S(i + 1, t - X[i])$ are already known, we can compute $S(i, t)$ in constant time, so memoizing this recurrence gives us an algorithm that runs in **$O(nT)$ time**.¹¹ To turn this into an explicit dynamic programming algorithm, we only need to consider the subproblems $S(i, t)$ in the proper order:

```

SUBSETSUM( $X[1..n], T$ ):
   $S[n + 1, 0] \leftarrow \text{TRUE}$ 
  for  $t \leftarrow 1$  to  $T$ 
     $S[n + 1, t] \leftarrow \text{FALSE}$ 

  for  $i \leftarrow n$  downto 1
     $S[i, 0] = \text{TRUE}$ 
    for  $t \leftarrow 1$  to  $X[i] - 1$ 
       $S[i, t] \leftarrow S[i + 1, t]$     «Avoid the case  $t < 0$ »
    for  $t \leftarrow X[i]$  to  $T$ 
       $S[i, t] \leftarrow S[i + 1, t] \vee S[i + 1, t - X[i]]$ 
  return  $S[1, T]$ 

```

This iterative algorithm clearly always uses **$O(nT)$ time and space**. In particular, if T is significantly larger than 2^n , this algorithm is actually slower than our naïve recursive algorithm. Dynamic programming isn't *always* an improvement!

5.6.2 NFA acceptance

The other problem we considered in the previous lecture note was determining whether a given NFA $M = (\Sigma, Q, s, A, \delta)$ accepts a given string $w \in \Sigma^*$. To make the problem concrete, we can assume without loss of generality that the alphabet is $\Sigma = \{1, 2, \dots, |\Sigma|\}$, the state set is $Q = \{1, 2, \dots, |Q|\}$, the start state is state 1, and our input consists of three arrays:

- A boolean array $A[1..|Q|]$, where $A[q] = \text{TRUE}$ if and only if $q \in A$.
- A boolean array $\delta[1..|Q|, 1..|\Sigma|, 1..|Q|]$, where $\delta[p, a, q] = \text{TRUE}$ if and only if $p \in \delta(q, a)$.
- An array $w[1..n]$ of symbols, representing the input string.

Now consider the boolean function

$\text{Accepts?}(q, i) = \text{TRUE}$ if and only if M accepts the suffix $w[i..n]$ starting in state q ,

or equivalently,

$\text{Accepts?}(q, i) = \text{TRUE}$ if and only if $\delta^*(q, w[i..n])$ contains at least one state in A .

¹¹Even though *SubsetSum* is NP-complete, this bound does *not* imply that $P=NP$, because T is not necessarily bounded by a polynomial function of the input size.

We need to compute $Accepts(1, 1)$. The recursive definition of the string transition function δ^* implies the following recurrence for $Accepts?$:

$$Accepts?(q, i) := \begin{cases} \text{TRUE} & \text{if } i > n \text{ and } q \in A \\ \text{FALSE} & \text{if } i > n \text{ and } q \notin A \\ \bigvee_{r \in \delta(q, a)} Accepts?(r, x) & \text{if } w = ax \end{cases}$$

Rewriting this recurrence in terms of our input representation gives us the following:

$$Accepts?(q, i) := \begin{cases} \text{TRUE} & \text{if } i > n \text{ and } A[q] = \text{TRUE} \\ \text{FALSE} & \text{if } i > n \text{ and } A[q] = \text{FALSE} \\ \bigvee_{r=1}^{|Q|} (\delta[q, w[i], r] \wedge Accepts?(r, i + 1)) & \text{otherwise} \end{cases}$$

We can memoize this function into a two-dimensional array $Accepts?[1..|Q|, 1..n+1]$. Each entry $Accepts?[q, i]$ depends on some subset of entries of the form $Accepts?[r, i+1]$. So we can fill the memoization table by considering the possible indices i in decreasing order in the outer loop, and consider states q in arbitrary order in the inner loop. Evaluating each entry $Accepts?[q, i]$ requires $O(|Q|)$ time, using an even deeper loop over all states r , and there are $O(n|Q|)$ such entries. Thus, the entire dynamic programming algorithm requires $O(n|Q|^2)$ time.

```

NFAACCEPTS?(A[1..|Q|],  $\delta[1..|Q|, 1..|\Sigma|, 1..|Q|]$ , w[1..n]):
  for q  $\leftarrow$  1 to |Q|
    Accepts?[q, n+1]  $\leftarrow$  A[q]
  for i  $\leftarrow$  n down to 1
    for q  $\leftarrow$  1 to |Q|
      Accepts?[q, i]  $\leftarrow$  FALSE
      for r  $\leftarrow$  1 to |Q|
        if  $\delta[q, w[i], r]$  and Accepts?[r, i+1]
          Accepts?[q, i]  $\leftarrow$  TRUE
  return Accepts?[1, 1]
```

5.7 Optimal Binary Search Trees

In an earlier lecture, we developed a recursive algorithm for the optimal binary search tree problem. We are given a sorted array $A[1..n]$ of search keys and an array $f[1..n]$ of frequency counts, where $f[i]$ is the number of searches to $A[i]$. Our task is to construct a binary search tree for that set such that the total cost of all the searches is as small as possible. We developed the following recurrence for this problem:

$$OptCost(f[1..n]) = \min_{1 \leq r \leq n} \left\{ OptCost(f[1..r-1]) + \sum_{i=1}^n f[i] + OptCost(f[r+1..n]) \right\}$$

To put this recurrence in more standard form, fix the frequency array f , and let $OptCost(i, j)$ denote the total search time in the optimal search tree for the subarray $A[i..j]$. To simplify notation a bit, let $F(i, j)$ denote the total frequency count for all the keys in the interval $A[i..j]$:

$$F(i, j) := \sum_{k=i}^j f[k]$$

We can now write

$$OptCost(i, j) = \begin{cases} 0 & \text{if } j < i \\ F(i, j) + \min_{i \leq r \leq j} (OptCost(i, r-1) + OptCost(r+1, j)) & \text{otherwise} \end{cases}$$

The base case might look a little weird, but all it means is that the total cost for searching an empty set of keys is zero.

The algorithm will be somewhat simpler and more efficient if we precompute all possible values of $F(i, j)$ and store them in an array. Computing each value $F(i, j)$ using a separate for-loop would $O(n^3)$ time. A better approach is to turn the recurrence

$$F(i, j) = \begin{cases} f[i] & \text{if } i = j \\ F(i, j-1) + f[j] & \text{otherwise} \end{cases}$$

into the following $O(n^2)$ -time dynamic programming algorithm:

```
INITF( $f[1..n]$ ):
  for  $i \leftarrow 1$  to  $n$ 
     $F[i, i-1] \leftarrow 0$ 
    for  $j \leftarrow i$  to  $n$ 
       $F[i, j] \leftarrow F[i, j-1] + f[j]$ 
```

This will be used as an initialization subroutine in our final algorithm.

So now let's compute the optimal search tree cost $OptCost(1, n)$ from the bottom up. We can store all intermediate results in a table $OptCost[1..n, 0..n]$. Only the entries $OptCost[i, j]$ with $j \geq i-1$ will actually be used. The base case of the recurrence tells us that any entry of the form $OptCost[i, i-1]$ can immediately be set to 0. For any other entry $OptCost[i, j]$, we can use the following algorithm fragment, which comes directly from the recurrence:

```
COMPUTEOPTCOST( $i, j$ ):
   $OptCost[i, j] \leftarrow \infty$ 
  for  $r \leftarrow i$  to  $j$ 
     $tmp \leftarrow OptCost[i, r-1] + OptCost[r+1, j]$ 
    if  $OptCost[i, j] > tmp$ 
       $OptCost[i, j] \leftarrow tmp$ 
   $OptCost[i, j] \leftarrow OptCost[i, j] + F[i, j]$ 
```

The only question left is what order to fill in the table.

Each entry $OptCost[i, j]$ depends on all entries $OptCost[i, r-1]$ and $OptCost[r+1, j]$ with $i \leq k \leq j$. In other words, every entry in the table depends on all the entries directly to the left or directly below. In order to fill the table efficiently, we must choose an order that computes all those entries before $OptCost[i, j]$. There are at least three different orders that satisfy this constraint. The one that occurs to most people first is to scan through the table one diagonal at a time, starting with the trivial base cases $OptCost[i, i-1]$. The complete algorithm looks like this:

```
OPTIMALSEARCHTREE( $f[1..n]$ ):
  INITF( $f[1..n]$ )
  for  $i \leftarrow 1$  to  $n$ 
     $OptCost[i, i-1] \leftarrow 0$ 
  for  $d \leftarrow 0$  to  $n-1$ 
    for  $i \leftarrow 1$  to  $n-d$ 
      COMPUTEOPTCOST( $i, i+d$ )
  return  $OptCost[1, n]$ 
```

We could also traverse the array row by row from the bottom up, traversing each row from left to right, or column by column from left to right, traversing each columns from the bottom up.

```

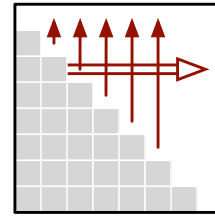
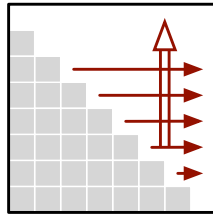
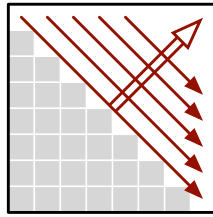
OPTIMALSEARCHTREE2( $f[1..n]$ ):
  INITF( $f[1..n]$ )
  for  $i \leftarrow n$  downto 1
     $OptCost[i, i-1] \leftarrow 0$ 
    for  $j \leftarrow i$  to  $n$ 
      COMPUTEOPTCOST( $i, j$ )
  return  $OptCost[1, n]$ 

```

```

OPTIMALSEARCHTREE3( $f[1..n]$ ):
  INITF( $f[1..n]$ )
  for  $j \leftarrow 0$  to  $n$ 
     $OptCost[j+1, j] \leftarrow 0$ 
    for  $i \leftarrow j$  downto 1
      COMPUTEOPTCOST( $i, j$ )
  return  $OptCost[1, n]$ 

```



Three different evaluation orders for the table $OptCost[i, j]$.

No matter which of these orders we actually use, the resulting algorithm runs in $\Theta(n^3)$ time and uses $\Theta(n^2)$ space. We could have predicted these space and time bounds directly from the original recurrence.

$$OptCost(i, j) = \begin{cases} 0 & \text{if } j = i - 1 \\ F(i, j) + \min_{i \leq r \leq j} (OptCost(i, r-1) + OptCost(r+1, j)) & \text{otherwise} \end{cases}$$

First, the function has two arguments, each of which can take on any value between 1 and n , so we probably need a table of size $O(n^2)$. Next, there are *three* variables in the recurrence (i , j , and r), each of which can take any value between 1 and n , so it should take us $O(n^3)$ time to fill the table.

5.8 The CYK Parsing Algorithm

In the same earlier lecture, we developed a recursive backtracking algorithm for parsing context-free languages. The input consists of a string w and a context-free grammar G in Chomsky normal form—meaning every production has the form $A \rightarrow a$, for some symbol a , or $A \rightarrow BC$, for some non-terminals B and C . Our task is to determine whether w is in the language generated by G .

Our backtracking algorithm recursively evaluates the boolean function $Generates?(A, x)$, which equals TRUE if and only if string x can be derived from non-terminal A , using the following recurrence:

$$Generates?(A, x) = \begin{cases} \text{TRUE} & \text{if } |x| = 1 \text{ and } A \rightarrow x \\ \text{FALSE} & \text{if } |x| = 1 \text{ and } A \not\rightarrow x \\ \bigvee_{A \rightarrow BC} \bigvee_{y \cdot z = x} Generates?(B, y) \wedge Generates?(C, z) & \text{otherwise} \end{cases}$$

This recurrence was transformed into a dynamic programming algorithm by Tadao Kasami in 1965, and again independently by Daniel Younger in 1967, and again independently by John

Cocke in 1970, so naturally the resulting algorithm is known as “Cocke-Younger-Kasami”, or more commonly *the CYK algorithm*.

We can derive the CYK algorithm from the previous recurrence as follows. As usual for recurrences involving strings, we need to modify the function slightly to ease memoization. Fix the input string w , and then let $Gen?(A, i, j) = \text{TRUE}$ if and only if the substring $w[i..j]$ can be derived from non-terminal A . Now our earlier recurrence can be rewritten as follows:

$$Gen?(A, i, j) = \begin{cases} \text{TRUE} & \text{if } i = j \text{ and } A \rightarrow w[i] \\ \text{FALSE} & \text{if } i = j \text{ and } A \not\rightarrow w[i] \\ \bigvee_{A \rightarrow BC} \bigvee_{k=i}^{j-1} Gen?(B, i, k) \wedge Gen?(C, k+1, j) & \text{otherwise} \end{cases}$$

This recurrence can be memoized into a three-dimensional boolean array $Gen[1..|\Gamma|, 1..n, 1..n]$, where the first dimension is indexed by the non-terminals Γ in the input grammar. Each entry $Gen[A, i, j]$ in this array depends on entries of the form $Gen[\cdot, i, k]$ for some $k < j$, or $Gen[\cdot, k+1, j]$ for some $k \geq i$. Thus, we can fill the array by increasing j in the outer loop, decreasing i in the middle loop, and considering non-terminals A in arbitrary order in the inner loop. The resulting dynamic programming algorithm runs in $O(n^3 \cdot |\Gamma|)$ time.

```

CYK(w, G):
  for i ← 1 to n
    for all non-terminals A
      if G contains the production A → w[i]
        Gen[A, i, i] ← TRUE
      else
        Gen[A, i, i] ← FALSE
  for j ← 1 to n
    for i ← n down to j + 1
      for all non-terminals A
        Gen[A, i, j] ← FALSE
      for all production rules A → BC
        for k ← i to j - 1
          if Gen[B, i, k] and Gen[C, k + 1, j]
            Gen[A, i, j] ← TRUE
  return Gen[S, 1, n]

```

5.9 Dynamic Programming on Trees

So far, all of our dynamic programming examples use a multidimensional array to store the results of recursive subproblems. However, as the next example shows, this is not always the most appropriate data structure to use.

A **independent set** in a graph is a subset of the vertices that have no edges between them. Finding the largest independent set in an arbitrary graph is extremely hard; in fact, this is one of the canonical NP-hard problems described in another lecture note. But from some special cases of graphs, we can find the largest independent set efficiently. In particular, when the input graph is a tree (a connected and acyclic graph) with n vertices, we can compute the largest independent set in $O(n)$ time.

In the recursion notes, we saw a recursive algorithm for computing the size of the largest independent set in an arbitrary graph:


```

MAXIMUMINDSETSIZE(G):
  if  $G = \emptyset$ 
    return 0

   $v \leftarrow$  any node in  $G$ 
   $withv \leftarrow 1 + \text{MAXIMUMINDSETSIZE}(G \setminus N(v))$ 
   $withoutv \leftarrow \text{MAXIMUMINDSETSIZE}(G \setminus \{v\})$ 
  return  $\max\{withv, withoutv\}$ .

```

Here, $N(v)$ denotes the *neighborhood* of v : the set containing v and all of its neighbors. As we observed in the other lecture notes, this algorithm has a worst-case running time of $O(2^n \text{poly}(n))$, where n is the number of vertices in the input graph.

Now suppose we require that the input graph is a tree; we will call this tree T instead of G from now on. We need to make a slight change to the algorithm to make it truly recursive. The subgraphs $T \setminus \{v\}$ and $T \setminus N(v)$ are forests, which may have more than one component. But the largest independent set in a disconnected graph is just the union of the largest independent sets in its components, so we can separately consider each tree in these forests. Fortunately, this has the added benefit of making the recursive algorithm more efficient, especially if we can choose the node v such that the trees are all significantly smaller than T . Here is the modified algorithm:

```

MAXIMUMINDSETSIZE(T):
  if  $T = \emptyset$ 
    return 0

   $v \leftarrow$  any node in  $T$ 
   $withv \leftarrow 1$ 
  for each tree  $T'$  in  $T \setminus N(v)$ 
     $withv \leftarrow withv + \text{MAXIMUMINDSETSIZE}(T')$ 
   $withoutv \leftarrow 0$ 
  for each tree  $T'$  in  $T \setminus \{v\}$ 
     $withoutv \leftarrow withoutv + \text{MAXIMUMINDSETSIZE}(T')$ 
  return  $\max\{withv, withoutv\}$ .

```

Now let's try to memoize this algorithm. Each recursive subproblem considers a subtree (that is, a connected subgraph) of the original tree T . Unfortunately, a single tree T can have exponentially many subtrees, so we seem to be doomed from the start!

Fortunately, there's a degree of freedom that we have not yet exploited: *We get to choose the vertex v .* We need a recipe—an algorithm!—for choosing v in each subproblem that limits the number of different subproblems the algorithm considers. To make this work, we impose some additional structure on the original input tree. Specifically, we declare one of the vertices of T to be the *root*, and we orient all the edges of T away from that root. Then we let v be the root of the input tree; this choice guarantees that each recursive subproblem considers a *rooted* subtree of T . Each vertex in T is the root of exactly one subtree, so now the number of distinct subproblems is exactly n . We can further simplify the algorithm by only passing a single node instead of the entire subtree:

```

MAXIMUMINDSETSIZE(v):
   $withv \leftarrow 1$ 
  for each grandchild  $x$  of  $v$ 
     $withv \leftarrow withv + \text{MAXIMUMINDSETSIZE}(x)$ 
   $withoutv \leftarrow 0$ 
  for each child  $w$  of  $v$ 
     $withoutv \leftarrow withoutv + \text{MAXIMUMINDSETSIZE}(w)$ 
  return  $\max\{withv, withoutv\}$ .

```

What data structure should we use to store intermediate results? The most natural choice is the tree itself! Specifically, for each node v , we store the result of $\text{MAXIMUMINDSETSIZE}(v)$ in a new field $v.MIS$. (We *could* use an array, but then we'd have to add a new field to each node anyway, pointing to the corresponding array entry. Why bother?)

What's the running time of the algorithm? The non-recursive time associated with each node v is proportional to the number of children and grandchildren of v ; this number can be very different from one vertex to the next. But we can turn the analysis around: Each vertex contributes a constant amount of time to its parent and its grandparent! Since each vertex has at most one parent and at most one grandparent, the total running time is $O(n)$.

What's a good order to consider the subproblems? The subproblem associated with any node v depends on the subproblems associated with the children and grandchildren of v . So we can visit the nodes in any order, provided that all children are visited before their parent. In particular, we can use a straightforward post-order traversal.

Here is the resulting dynamic programming algorithm. Yes, it's still recursive. I've swapped the evaluation of the with- v and without- v cases; we need to visit the kids first anyway, so why not consider the subproblem that depends directly on the kids first?

```

MAXIMUMINDSETSIZE( $v$ ):
  without $v$   $\leftarrow$  0
  for each child  $w$  of  $v$ 
    without $v$   $\leftarrow$  without $v$  + MAXIMUMINDSETSIZE( $w$ )
  with $v$   $\leftarrow$  1
  for each grandchild  $x$  of  $v$ 
    with $v$   $\leftarrow$  with $v$  +  $x.MIS$ 
   $v.MIS$   $\leftarrow$  max{with $v$ , without $v$ }
  return  $v.MIS$ 

```

Another option is to store *two* values for each rooted subtree: the size of the largest independent set *that includes the root*, and the size of the largest independent set *that excludes the root*. This gives us an even simpler algorithm, with the same $O(n)$ running time.

```

MAXIMUMINDSETSIZE( $v$ ):
   $v.MISno$   $\leftarrow$  0
   $v.MISyes$   $\leftarrow$  1
  for each child  $w$  of  $v$ 
     $v.MISno$   $\leftarrow$   $v.MISno$  + MAXIMUMINDSETSIZE( $w$ )
     $v.MISyes$   $\leftarrow$   $v.MISyes$  +  $w.MISno$ 
  return max{ $v.MISyes$ ,  $v.MISno$ }

```

Exercises

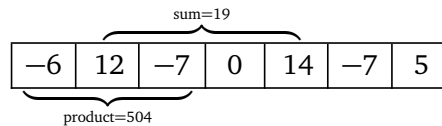
Sequences/Arrays

1. In a previous life, you worked as a cashier in the lost Antarctic colony of Nadira, spending the better part of your day giving change to your customers. Because paper is a very rare and valuable resource in Antarctica, cashiers were required by law to use the fewest bills possible whenever they gave change. Thanks to the numerological predilections of one of its founders, [the currency of Nadira, called Dream Dollars](#), was available in the following denominations: \$1, \$4, \$7, \$13, \$28, \$52, \$91, \$365.¹²

¹²For more details on the history and culture of Nadira, including images of the various denominations of Dream Dollars, see <http://moneyart.biz/dd/>.

- (a) The greedy change algorithm repeatedly takes the largest bill that does not exceed the target amount. For example, to make \$122 using the greedy algorithm, we first take a \$91 bill, then a \$28 bill, and finally three \$1 bills. Give an example where this greedy algorithm uses more Dream Dollar bills than the minimum possible.
 - (b) Describe and analyze a recursive algorithm that computes, given an integer k , the minimum number of bills needed to make k Dream Dollars. (Don't worry about making your algorithm fast; just make sure it's correct.)
 - (c) Describe a dynamic programming algorithm that computes, given an integer k , the minimum number of bills needed to make k Dream Dollars. (This one needs to be fast.)
2. Suppose you are given an array $A[1..n]$ of numbers, which may be positive, negative, or zero, and which are *not* necessarily integers.
 - (a) Describe and analyze an algorithm that finds the largest sum of elements in a contiguous subarray $A[i..j]$.
 - (b) Describe and analyze an algorithm that finds the largest *product* of elements in a contiguous subarray $A[i..j]$.

For example, given the array $[-6, 12, -7, 0, 14, -7, 5]$ as input, your first algorithm should return the integer 19, and your second algorithm should return the integer 504.



For the sake of analysis, assume that comparing, adding, or multiplying any pair of numbers takes $O(1)$ time.

[Hint: Problem (a) has been a standard computer science interview question since at least the mid-1980s. You can find many correct solutions on the web; the problem even has its own [Wikipedia page](#)! But at least in 2013, the few solutions I found on the web for problem (b) were all either slower than necessary or actually incorrect.]

3. This series of exercises asks you to develop efficient algorithms to find optimal *subsequences* of various kinds. A subsequence is anything obtained from a sequence by extracting a subset of elements, but keeping them in the same order; the elements of the subsequence need not be contiguous in the original sequence. For example, the strings **C**, **DAMN**, **YAIIOAI**, and **DYNAMICPROGRAMMING** are all subsequences of the string **DYNAMICPROGRAMMING**.
 - (a) Let $A[1..m]$ and $B[1..n]$ be two arbitrary arrays. A **common subsequence** of A and B is another sequence that is a subsequence of both A and B . Describe an efficient algorithm to compute the length of the *longest* common subsequence of A and B .
 - (b) Let $A[1..m]$ and $B[1..n]$ be two arbitrary arrays. A **common supersequence** of A and B is another sequence that contains both A and B as subsequences. Describe an efficient algorithm to compute the length of the *shortest* common supersequence of A and B .

- (c) Call a sequence $X[1..n]$ of numbers **bitonic** if there is an index i with $1 < i < n$, such that the prefix $X[1..i]$ is increasing and the suffix $X[i..n]$ is decreasing. Describe an efficient algorithm to compute the length of the longest bitonic subsequence of an arbitrary array A of integers.
- (d) Call a sequence $X[1..n]$ of numbers **oscillating** if $X[i] < X[i+1]$ for all even i , and $X[i] > X[i+1]$ for all odd i . Describe an efficient algorithm to compute the length of the longest oscillating subsequence of an arbitrary array A of integers.
- (e) Describe an efficient algorithm to compute the length of the shortest oscillating supersequence of an arbitrary array A of integers.
- (f) Call a sequence $X[1..n]$ of numbers **convex** if $2 \cdot X[i] < X[i-1] + X[i+1]$ for all i . Describe an efficient algorithm to compute the length of the longest convex subsequence of an arbitrary array A of integers.
- (g) Call a sequence $X[1..n]$ of numbers **weakly increasing** if each element is larger than the average of the two previous elements; that is, $2 \cdot X[i] > X[i-1] + X[i-2]$ for all $i > 2$. Describe an efficient algorithm to compute the length of the longest weakly increasing subsequence of an arbitrary array A of integers.
- (h) Call a sequence $X[1..n]$ of numbers **double-increasing** if $X[i] > X[i-2]$ for all $i > 2$. (In other words, a semi-increasing sequence is a perfect shuffle of two increasing sequences.) Describe an efficient algorithm to compute the length of the longest double-increasing subsequence of an arbitrary array A of integers.
- * (i) Recall that a sequence $X[1..n]$ of numbers is *increasing* if $X[i] < X[i+1]$ for all i . Describe an efficient algorithm to compute the length of the *longest common increasing subsequence* of two given arrays of integers. For example, $\langle 1, 4, 5, 6, 7, 9 \rangle$ is the longest common increasing subsequence of the sequences $\langle 3, 1, 4, 1, 5, 9, 2, 6, 5, 3, 5, 8, 9, 7, 9, 3 \rangle$ and $\langle 1, 4, 1, 4, 2, 1, 3, 5, 6, 2, 3, 7, 3, 0, 9, 5 \rangle$.
4. Describe an algorithm to compute the number of times that one given array $X[1..k]$ appears as a subsequence of another given array $Y[1..n]$. For example, if all characters in X and Y are equal, your algorithm should return $\binom{n}{k}$. For purposes of analysis, assume that adding two ℓ -bit integers requires $\Theta(\ell)$ time.
5. You and your eight-year-old nephew Elmo decide to play a simple card game. At the beginning of the game, the cards are dealt face up in a long row. Each card is worth a different number of points. After all the cards are dealt, you and Elmo take turns removing either the leftmost or rightmost card from the row, until all the cards are gone. At each turn, you can decide which of the two cards to take. The winner of the game is the player that has collected the most points when the game ends.

Having never taken an algorithms class, Elmo follows the obvious greedy strategy—when it's his turn, Elmo *always* takes the card with the higher point value. Your task is to find a strategy that will beat Elmo whenever possible. (It might seem mean to beat up on a little kid like this, but Elmo absolutely *hates* it when grown-ups let him win.)

- (a) Prove that you should not also use the greedy strategy. That is, show that there is a game that you can win, but only if you do *not* follow the same greedy strategy as Elmo.

- (b) Describe and analyze an algorithm to determine, given the initial sequence of cards, the maximum number of points that you can collect playing against Elmo.
- (c) When Elmo was four, he used an even simple strategy—on his turn, he always chose his next card uniformly at random. That is, if there was more than one card left on his turn, he would take the leftmost card with probability $1/2$, and the rightmost card with probability $1/2$. Describe an algorithm to determine, given the initial sequence of cards, the maximum *expected* number of points you can collect playing against four-year-old-Elmo.
- (d) Five years later, Elmo has become a *much* stronger player. Describe and analyze an algorithm to determine, given the initial sequence of cards, the maximum number of points that you can collect playing against a *perfect* opponent.
6. It's almost time to show off your flippin' sweet dancing skills! Tomorrow is the big dance contest you've been training for your entire life, except for that summer you spent with your uncle in Alaska hunting wolverines. You've obtained an advance copy of the the list of n songs that the judges will play during the contest, in chronological order.

You know all the songs, all the judges, and your own dancing ability extremely well. For each integer k , you know that if you dance to the k th song on the schedule, you will be awarded exactly $Score[k]$ points, but then you will be physically unable to dance for the next $Wait[k]$ songs (that is, you cannot dance to songs $k + 1$ through $k + Wait[k]$). The dancer with the highest total score at the end of the night wins the contest, so you want your total score to be as high as possible.

Describe and analyze an efficient algorithm to compute the maximum total score you can achieve. The input to your sweet algorithm is the pair of arrays $Score[1..n]$ and $Wait[1..n]$.

7. You are driving a bus along a highway, full of rowdy, hyper, thirsty students and a soda fountain machine. Each minute that a student is on your bus, that student drinks one ounce of soda. Your goal is to drop the students off quickly, so that the total amount of soda consumed by all students is as small as possible.

You know how many students will get off of the bus at each exit. Your bus begins somewhere along the highway (probably not at either end) and moves at a constant speed of 37.4 miles per hour. You must drive the bus along the highway; however, you may drive forward to one exit then backward to an exit in the opposite direction, switching as often as you like. (You can stop the bus, drop off students, and turn around instantaneously.)

Describe an efficient algorithm to drop the students off so that they drink as little soda as possible. Your input consists of the bus route (a list of the exits, together with the travel time between successive exits), the number of students you will drop off at each exit, and the current location of your bus (which you may assume is an exit).

8. A palindrome is any string that is exactly the same as its reversal, like **I**, or **DEED**, or **RACECAR**, or **AMANAPLANACATACANALPANAMA**.
- (a) Describe and analyze an algorithm to find the length of the *longest subsequence* of a given string that is also a palindrome. For example, the longest palindrome

subsequence of MAHDYNAMICPROGRAMZLETMESHOWYOUTHEM is MHYMRORMYHM, so given that string as input, your algorithm should output the number 11.

- (b) Describe and analyze an algorithm to find the length of the *shortest supersequence* of a given string that is also a palindrome. For example, the shortest palindrome supersequence of TWENTYONE is TWENTYOYOTNEWT, so given the string TWENTYONE as input, your algorithm should output the number 13.
- (c) Any string can be decomposed into a sequence of palindromes. For example, the string BUBBASEESABANANA (“Bubba sees a banana.”) can be broken into palindromes in the following ways (and many others):

BUB • BASEESAB • ANANA
B • U • BB • A • SEES • ABA • NAN • A
B • U • BB • A • SEES • A • B • ANANA
B • U • B • B • A • S • E • E • S • A • B • A • N • ANA

Describe and analyze an efficient algorithm to find the smallest number of palindromes that make up a given input string. For example, given the input string BUBBASEESABANANA, your algorithm would return the integer 3.

9. Suppose you have a black-box subroutine `QUALITY` that can compute the ‘quality’ of any given string $A[1..k]$ in $O(k)$ time. For example, the quality of a string might be 1 if the string is a Québécois curse word, and 0 otherwise.

Given an arbitrary input string $T[1..n]$, we would like to break it into contiguous substrings, such that the total quality of all the substrings is as large as possible. For example, the string SAINTCIBOIREDESACRAMENTDECRISE can be decomposed into the substrings SAINT • CIBOIRE • DE • SACRAMENT • DE • CRISSE, of which three (or possibly four) are *sacres*.

Describe an algorithm that breaks a string into substrings of maximum total quality, using the `QUALITY` subroutine.

10. (a) Suppose we are given a set L of n line segments in the plane, where each segment has one endpoint on the line $y = 0$ and one endpoint on the line $y = 1$, and all $2n$ endpoints are distinct. Describe and analyze an algorithm to compute the largest subset of L in which no pair of segments intersects.
- (b) Suppose we are given a set L of n line segments in the plane, where each segment has one endpoint on the line $y = 0$ and one endpoint on the line $y = 1$, and all $2n$ endpoints are distinct. Describe and analyze an algorithm to compute the largest subset of L in which *every* pair of segments intersects.
- (c) Suppose we are given a set L of n line segments in the plane, where the endpoints of each segment lie on the unit circle $x^2 + y^2 = 1$, and all $2n$ endpoints are distinct. Describe and analyze an algorithm to compute the largest subset of L in which no pair of segments intersects.
- (d) Suppose we are given a set L of n line segments in the plane, where the endpoints of each segment lie on the unit circle $x^2 + y^2 = 1$, and all $2n$ endpoints are distinct. Describe and analyze an algorithm to compute the largest subset of L in which *every* pair of segments intersects.

11. Let P be a set of n points evenly distributed on the unit circle, and let S be a set of m line segments with endpoints in P . The endpoints of the m segments are *not* necessarily distinct; n could be significantly smaller than $2m$.
- Describe an algorithm to find the size of the largest subset of segments in S such that every pair is disjoint. Two segments are disjoint if they do not intersect even at their endpoints.
 - Describe an algorithm to find the size of the largest subset of segments in S such that every pair is interior-disjoint. Two segments are interior-disjoint if their intersection is either empty or an endpoint of both segments.
 - Describe an algorithm to find the size of the largest subset of segments in S such that every pair intersects.
 - Describe an algorithm to find the size of the largest subset of segments in S such that every pair crosses. Two segments cross if they intersect but not at their endpoints.

For full credit, all four algorithms should run in $O(mn)$ time.

12. A *shuffle* of two strings X and Y is formed by interspersing the characters into a new string, keeping the characters of X and Y in the same order. For example, the string **BANANAANANAS** is a shuffle of the strings **BANANA** and **ANANAS** in several different ways.

BANANAANANAS
 BANANAANANAS
 BANANANANAS

Similarly, the strings **PRODGYRNAMAMMIINCG** and **DYPRONGARMAMMICING** are both shuffles of **DYNAMIC** and **PROGRAMMING**:

PRODGYRNAMAMMIINCG
 DYPRONGARMAMMICING

Given three strings $A[1..m]$, $B[1..n]$, and $C[1..m+n]$, describe and analyze an algorithm to determine whether C is a shuffle of A and B .

13. Let's define a *summary* of two strings A and B to be a concatenation of substrings of the following form:
- ▲**SNA** indicates a substring **SNA** of only the first string A .
 - ◆**F00** indicates a common substring **F00** of both strings.
 - ▼**BAR** indicates a substring **BAR** of only the second string B .

A summary is *valid* if we can recover the original strings A and B by concatenating the appropriate substrings of the summary in order and discarding the delimiters ▲, ◆, and ▼. Each regular character has length 1, and each delimiter ▲, ◆, or ▼ has some fixed non-negative length Δ . The *length* of a summary is the sum of the lengths of its symbols.

For example, each of the following strings is a valid summary of the strings **KITTEN** and **KNITTING**:

- ◆**K**▼**N**◆**ITT**▲**E**▼**I**◆**N**▼**G** has length $9 + 7\Delta$.
- ◆**K**▼**N**◆**ITT**▲**EN**▼**ING** has length $10 + 5\Delta$.

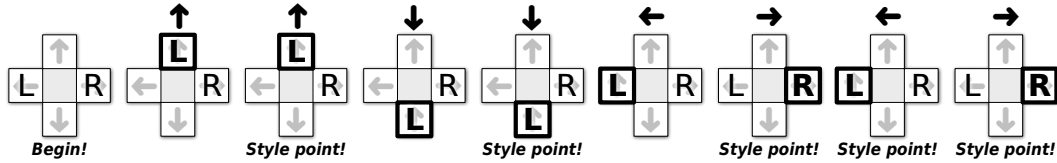
- \blacklozenge K \blacktriangle ITTEN \blacktriangledown NITTING has length $13 + 3\Delta$.
- \blacktriangle KITTEN \blacktriangledown KNITTING has length $14 + 2\Delta$.

Describe and analyze an algorithm that computes the length of the shortest summary of two given strings $A[1..m]$ and $B[1..n]$. The delimiter length Δ is also part of the input to your algorithm. For example:

- Given strings KITTEN and KNITTING and $\Delta = 0$, your algorithm should return 9.
 - Given strings KITTEN and KNITTING and $\Delta = 1$, your algorithm should return 15.
 - Given strings KITTEN and KNITTING and $\Delta = 2$, your algorithm should return 18.
14. Describe and analyze an efficient algorithm to find the length of the longest contiguous substring that appears both forward and backward in an input string $T[1..n]$. The forward and backward substrings must not overlap. Here are several examples:
- Given the input string ALGORITHM, your algorithm should return 0.
 - Given the input string RECURSION, your algorithm should return 1, for the substring R.
 - Given the input string REDIVIDE, your algorithm should return 3, for the substring EDI. (The forward and backward substrings must not overlap!)
 - Given the input string DYNAMICPROGRAMMINGMANYTIMES, your algorithm should return 4, for the substring YNAM. (In particular, it should *not* return 6, for the subsequence YNAMIR).
15. **Dance Dance Revolution** is a dance video game, first introduced in Japan by Konami in 1998. Players stand on a platform marked with four arrows, pointing forward, back, left, and right, arranged in a cross pattern. During play, the game plays a song and scrolls a sequence of n arrows (\leftarrow , \uparrow , \downarrow , or \rightarrow) from the bottom to the top of the screen. At the precise moment each arrow reaches the top of the screen, the player must step on the corresponding arrow on the dance platform. (The arrows are timed so that you'll step with the beat of the song.)

You are playing a variant of this game called “Vogue Vogue Revolution”, where the goal is to play perfectly but move as little as possible. When an arrow reaches the top of the screen, if one of your feet is already on the correct arrow, you are awarded one style point for maintaining your current pose. If neither foot is on the right arrow, you must move one (and *only* one) of your feet from its current location to the correct arrow on the platform. If you ever step on the wrong arrow, or fail to step on the correct arrow, or move more than one foot at a time, or move either foot when you are already standing on the correct arrow, all your style points are taken away and you lose the game.

How should you move your feet to maximize your total number of style points? For purposes of this problem, assume you always start with you left foot on \leftarrow and you right foot on \rightarrow , and that you've memorized the entire sequence of arrows. For example, if the sequence is $\uparrow\uparrow\downarrow\downarrow\leftarrow\rightarrow\leftarrow\rightarrow$, you can earn 5 style points by moving you feet as shown below:



- (a) **Prove** that for any sequence of n arrows, it is possible to earn at least $n/4 - 1$ style points.
- (b) Describe an efficient algorithm to find the maximum number of style points you can earn during a given VVR routine. The input to your algorithm is an array $Arrow[1..n]$ containing the sequence of arrows.

16. Consider the following solitaire form of Scrabble. We begin with a fixed, finite sequence of tiles; each tile contains a letter and a numerical value. At the start of the game, we draw the seven tiles from the sequence and put them into our hand. In each turn, we form an English word from some or all of the tiles in our hand, place those tiles on the table, and receive the total value of those tiles as points. If no English word can be formed from the tiles in our hand, the game immediately ends. Then we repeatedly draw the next tile from the start of the sequence until either (a) we have seven tiles in our hand, or (b) the sequence is empty. (Sorry, no double/triple word/letter scores, bingos, blanks, or passing.) Our goal is to obtain as many points as possible.

For example, suppose we are given the tile sequence



Then we can earn 68 points as follows:

- We initially draw $I_2, N_2, X_8, A_1, N_2, A_1, D_3$.
 - Play the word N_2, A_1, I_2, A_1, D_3 for 9 points, leaving N_2, X_8 in our hand.
 - Draw the next five tiles U_5, D_3, I_2, D_3, K_8 .
 - Play the word U_5, N_2, D_3, I_2, D_3 for 15 points, leaving K_8, X_8 in our hand.
 - Draw the next five tiles U_5, B_4, L_2, A_1, K_8 .
 - Play the word B_4, U_5, L_2, K_8 for 19 points, leaving K_8, X_8, A_1 in our hand.
 - Draw the next three tiles H_5, A_1, N_2 , emptying the list.
 - Play the word A_1, N_2, K_8, H_5 for 16 points, leaving X_8, A_1 in our hand.
 - Play the word A_1, X_8 for 9 points, emptying our hand and ending the game.
- (a) Suppose you are given as input two arrays $Letter[1..n]$, containing a sequence of letters between **A** and **Z**, and $Value[A..Z]$, where $Value[\ell]$ is the value of letter ℓ . Design and analyze an efficient algorithm to compute the maximum number of points that can be earned from the given sequence of tiles.

- (b) Now suppose two tiles with the same letter can have different values; you are given two arrays $Letter[1..n]$ and $Value[1..n]$. Design and analyze an efficient algorithm to compute the maximum number of points that can be earned from the given sequence of tiles.

In both problems, the output is a single number: the maximum possible score. Assume that you can find all English words that can be made from any set of at most seven tiles, along with the point values of those words, in $O(1)$ time.

17. *Vankin's Mile* is an American solitaire game played on an $n \times n$ square grid. The player starts by placing a token on any square of the grid. Then on each turn, the player moves the token either one square to the right or one square down. The game ends when player moves the token off the edge of the board. Each square of the grid has a numerical value, which could be positive, negative, or zero. The player starts with a score of zero; whenever the token lands on a square, the player adds its value to his score. The object of the game is to score as many points as possible.

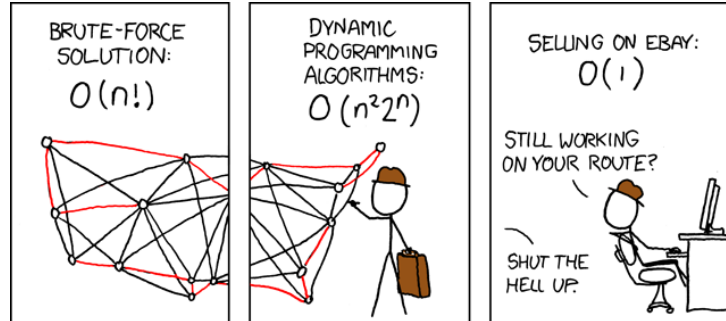
For example, given the grid below, the player can score $8 - 6 + 7 - 3 + 4 = 10$ points by placing the initial token on the 8 in the second row, and then moving down, down, right, down, down. (This is *not* the best possible score for these values.)

-1	7	-8	10	-5
-4	-9	8	-6	0
5	-2	-6	-6	7
-7	4	7	-3	-3
7	1	-6	4	-9

- (a) Describe and analyze an efficient algorithm to compute the maximum possible score for a game of Vankin's Mile, given the $n \times n$ array of values as input.
- (b) In the European version of this game, appropriately called *Vankin's Kilometer*, the player can move the token either one square down, one square right, or one square left in each turn. However, to prevent infinite scores, the token cannot land on the same square more than once. Describe and analyze an efficient algorithm to compute the maximum possible score for a game of Vankin's Kilometer, given the $n \times n$ array of values as input.¹³
18. Suppose you are given an $m \times n$ bitmap, represented by an array $M[1..n, 1..n]$ of 0s and 1s. A *solid block* in M is a subarray of the form $M[i..i', j..j']$ containing only 1-bits. A solid block is square if it has the same number of rows and columns.
- (a) Describe an algorithm to find the maximum area of a solid *square* block in M in $O(n^2)$ time.
- (b) Describe an algorithm to find the maximum area of a solid block in M in $O(n^3)$ time.
- * (c) Describe an algorithm to find the maximum area of a solid block in M in $O(n^2)$ time.

¹³If we also allowed upward movement, the resulting game (Vankin's Fathom?) would be Ebay-hard.

- *19. Describe and analyze an algorithm to solve the traveling salesman problem in $O(2^n \text{poly}(n))$ time. Given an undirected n -vertex graph G with weighted edges, your algorithm should return the weight of the lightest cycle in G that visits every vertex exactly once, or ∞ if G has no such cycles. [Hint: The obvious recursive algorithm takes $O(n!)$ time.]



— Randall Munroe, *xkcd* (<http://xkcd.com/399/>)
Reproduced under a Creative Commons Attribution-NonCommercial 2.5 License

- *20. Let $\mathcal{A} = \{A_1, A_2, \dots, A_n\}$ be a finite set of strings over some fixed alphabet Σ . An *edit center* for \mathcal{A} is a string $C \in \Sigma^*$ such that the maximum edit distance from C to any string in \mathcal{A} is as small as possible. The *edit radius* of \mathcal{A} is the maximum edit distance from an edit center to a string in \mathcal{A} . A set of strings may have several edit centers, but its edit radius is unique.

$$\text{EditRadius}(\mathcal{A}) = \min_{C \in \Sigma^*} \max_{A \in \mathcal{A}} \text{Edit}(A, C) \quad \text{EditCenter}(\mathcal{A}) = \arg \min_{C \in \Sigma^*} \max_{A \in \mathcal{A}} \text{Edit}(A, C)$$

- (a) Describe and analyze an efficient algorithm to compute the edit radius of three given strings.
- (b) Describe and analyze an efficient algorithm to approximate the edit radius of an arbitrary set of strings within a factor of 2. (Computing the *exact* edit radius is NP-hard unless the number of strings is fixed.)
- ★21. Let $D[1..n]$ be an array of digits, each an integer between 0 and 9. A **digital subsequence** of D is a sequence of positive integers composed in the usual way from disjoint substrings of D . For example, 3, 4, 5, 6, 8, 9, 32, 38, 46, 64, 83, 279 is a digital subsequence of the first several digits of π :

3, 1, 4, 1, 5, 9, 2, 6, 5, 3, 5, 8, 9, 7, 9, 3, 2, 3, 8, 4, 6, 2, 6, 4, 3, 3, 8, 3, 2, 7, 9

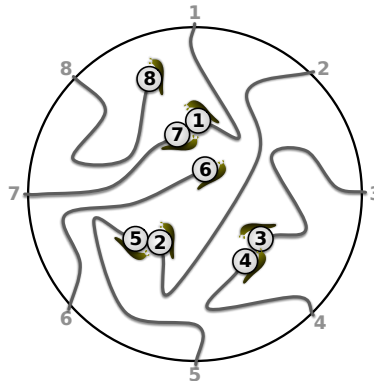
The *length* of a digital subsequence is the number of integers it contains, *not* the number of digits; the preceding example has length 12. As usual, a digital subsequence is **increasing** if each number is larger than its predecessor.

Describe and analyze an efficient algorithm to compute the longest increasing digital subsequence of D . [Hint: Be careful about your computational assumptions. How long does it take to compare two k -digit numbers?]

For full credit, your algorithm should run in $O(n^4)$ time; faster algorithms are worth extra credit. The fastest algorithm I know for this problem runs in $O(n^2 \log n)$ time; achieving this bound requires several tricks, both in the algorithm and in its analysis.

Splitting Sequences/Arrays

22. Every year, as part of its annual meeting, the Antarctic Snail Lovers of Upper Glacierville hold a Round Table Mating Race. Several high-quality breeding snails are placed at the edge of a round table. The snails are numbered in order around the table from 1 to n . During the race, each snail wanders around the table, leaving a trail of slime behind it. The snails have been specially trained never to fall off the edge of the table or to cross a slime trail, even their own. If two snails meet, they are declared a breeding pair, removed from the table, and whisked away to a romantic hole in the ground to make little baby snails. Note that some snails may never find a mate, even if the race goes on forever.



The end of a typical Antarctic SLUG race. Snails 6 and 8 never find mates.
The organizers must pay $M[3, 4] + M[2, 5] + M[1, 7]$.

For every pair of snails, the Antarctic SLUG race organizers have posted a monetary reward, to be paid to the owners if that pair of snails meets during the Mating Race. Specifically, there is a two-dimensional array $M[1..n, 1..n]$ posted on the wall behind the Round Table, where $M[i, j] = M[j, i]$ is the reward to be paid if snails i and j meet.

Describe and analyze an algorithm to compute the maximum total reward that the organizers could be forced to pay, given the array M as input.

23. Suppose you are given a sequence of integers separated by $+$ and \times signs; for example:

$$1 + 3 \times 2 \times 0 + 1 \times 6 + 7$$

You can change the value of this expression by adding parentheses in different places. For example:

$$(1 + (3 \times 2)) \times 0 + (1 \times 6) + 7 = 13$$

$$((1 + (3 \times 2 \times 0) + 1) \times 6) + 7 = 19$$

$$(1 + 3) \times 2 \times (0 + 1) \times (6 + 7) = 208$$

- (a) Describe and analyze an algorithm to compute the maximum possible value the given expression can take by adding parentheses, assuming all integers in the input are positive. [Hint: This is easy.]
- (b) Describe and analyze an algorithm to compute the maximum possible value the given expression can take by adding parentheses, assuming all integers in the input are non-negative.

- (c) Describe and analyze an algorithm to compute the maximum possible value the given expression can take by adding parentheses, with no further restrictions on the input.

Assume any arithmetic operation takes $O(1)$ time.

24. Suppose you are given a sequence of integers separated by + and - signs; for example:

$$1 + 3 - 2 - 5 + 1 - 6 + 7$$

You can change the value of this expression by adding parentheses in different places. For example:

$$\begin{aligned} 1 + 3 - 2 - 5 + 1 - 6 + 7 &= -1 \\ (1 + 3 - (2 - 5)) + (1 - 6) + 7 &= 9 \\ (1 + (3 - 2)) - (5 + 1) - (6 + 7) &= -17 \end{aligned}$$

Describe and analyze an algorithm to compute, given a list of integers separated by + and - signs, the maximum possible value the expression can take by adding parentheses.

You may only use parentheses to group additions and subtractions; in particular, you are not allowed to create implicit multiplication as in $1 + 3(-2)(-5) + 1 - 6 + 7 = 33$.

25. A **basic arithmetic expression** is composed of characters from the set $\{1, +, \times\}$ and parentheses. Almost every integer can be represented by more than one basic arithmetic expression. For example, all of the following basic arithmetic expression represent the integer 14:

$$\begin{aligned} 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 \\ ((1 + 1) \times (1 + 1 + 1 + 1 + 1)) + ((1 + 1) \times (1 + 1)) \\ (1 + 1) \times (1 + 1 + 1 + 1 + 1 + 1 + 1) \\ (1 + 1) \times (((1 + 1 + 1) \times (1 + 1)) + 1) \end{aligned}$$

Describe and analyze an algorithm to compute, given an integer n as input, the minimum number of 1's in a basic arithmetic expression whose value is n . The number of parentheses doesn't matter, just the number of 1's. For example, when $n = 14$, your algorithm should return 8, for the final expression above. For full credit, the running time of your algorithm should be bounded by a small polynomial function of n .

26. After graduating from Illinois, you decide to interview for a position at the Wall Street bank **Long Live Boole**. The managing director of the bank, Eloob Egroeg, is a genius mathematician who worships George Boole (the inventor of Boolean Logic) every morning before leaving for the office. The first day of every hired employee is a 'solve-or-die' day where s/he has to solve one of the problems posed by Eloob within 24 hours. Those who fail to solve the problem are fired immediately!

Entering the bank for the first time, you notice that the employee offices are organized in a straight row, with a large T or F printed on the door of each office. Furthermore, between each adjacent pair of offices, there is a board marked by one of the symbols \wedge , \vee , or \oplus .

When you ask about these arcane symbols, Eloob confirms that T and F represent the boolean values TRUE and FALSE, and the symbols on the boards represent the standard boolean operators AND, OR, and XOR. He also explains that these letters and symbols describe whether certain combinations of employees can work together successfully. At the start of any new project, Eloob hierarchically clusters his employees by adding parentheses to the sequence of symbols, to obtain an unambiguous boolean expression. The project is successful if this parenthesized boolean expression evaluates to T .

For example, if the bank has three employees, and the sequence of symbols on and between their doors is $T \wedge F \oplus T$, there is exactly one successful parenthesization scheme: $(T \wedge (F \oplus T))$. However, if the list of door symbols is $F \wedge T \oplus F$, there is no way to add parentheses to make the project successful.

Eloob finally poses your solve-or-die interview question: Describe an algorithm to decide whether a given sequence of symbols can be parenthesized so that the resulting boolean expression evaluates to T . The input to your algorithm is an array $S[0..2n]$, where $S[i] \in \{T, F\}$ when i is even, and $S[i] \in \{\vee, \wedge, \oplus\}$ when i is odd.

27. Suppose we want to display a paragraph of text on a computer screen. The text consists of a sequence of n words, where the i th word has length $\ell[i]$. We want to break the paragraph into several lines of total length exactly L . For example, according to T_EX, the program used to typeset these notes, *the paragraph you are reading right now* is approximately 14.27585 cm \approx 5.62231 inches wide.

Depending on how the paragraph is broken into lines of text, we must insert different amounts of white space between the words. The paragraph should be fully justified, meaning that the first character on each line starts at the left margin, and *except for the last line*, the last character on each line ends at the right margin. There must be at least one unit of white space between any two words on the same line.

Define the *slop* of a paragraph layout as the sum over all lines, *except the last*, of the cube of the amount of extra white-space in each line, not counting the one unit of required space between each adjacent pair of words. Specifically, if a line contains words i through j , then the slop of that line is defined to be $(L - j + i - \sum_{k=i}^j \ell[k])^3$. Describe a dynamic programming algorithm to print the paragraph with minimum slop.

28. You have mined a large slab of marble from your quarry. For simplicity, suppose the marble slab is a rectangle measuring n inches in height and m inches in width. You want to cut the slab into smaller rectangles of various sizes—some for kitchen countertops, some for large sculpture projects, others for memorial headstones. You have a marble saw that can make either horizontal or vertical cuts across any rectangular slab. At any time, you can query the spot price $P[x, y]$ of an x -inch by y -inch marble rectangle, for any positive integers x and y . These prices depend on customer demand, and people who buy marble countertops are weird, so do not make any assumptions about them; in particular, larger rectangles may have significantly smaller spot prices. Given the array of spot prices and the integers m and n as input, describe an algorithm to compute how to subdivide an $n \times m$ marble slab to maximize your profit.

29. A string w of parentheses **(** and **)** and brackets **[** and **]** is *balanced* if it satisfies one of the following conditions:
- w is the empty string.
 - $w = \mathbf{(x)}$ for some balanced string x
 - $w = \mathbf{[x]}$ for some balanced string x
 - $w = xy$ for some balanced strings x and y

For example, the string

$$w = \mathbf{([()])} \mathbf{[()]} \mathbf{([()])} \mathbf{([()])} \mathbf{([()])}$$

is balanced, because $w = xy$, where

$$x = \mathbf{([()])} \mathbf{[()]} \mathbf{([()])} \quad \text{and} \quad y = \mathbf{([()])} \mathbf{([()])} \mathbf{([()])}$$

- Describe and analyze an algorithm to determine whether a given string of parentheses and brackets is balanced.
- Describe and analyze an algorithm to compute the length of a longest balanced subsequence of a given string of parentheses and brackets.
- Describe and analyze an algorithm to compute the length of a shortest balanced supersequence of a given string of parentheses and brackets.
- Describe and analyze an algorithm to compute the minimum edit distance from a given string of parentheses and brackets to a balanced string of parentheses and brackets.

For each problem, your input is an array $w[1..n]$, where $w[i] \in \{\mathbf{(,), [,]}\}$ for every index i . (You may prefer to use different symbols instead of parentheses and brackets—for example, L, R, l, r or $\triangleleft, \triangleright, \blacktriangleleft, \blacktriangleright$ —but please tell us what symbols you're using!)

30. Congratulations! Your research team has just been awarded a \$50M multi-year project, jointly funded by DARPA, Google, and McDonald's, to produce DWIM: The first compiler to read programmers' minds! Your proposal and your numerous press releases all promise that DWIM will automatically correct errors in any given piece of code, while modifying that code as little as possible. Unfortunately, now it's time to start actually making the damn thing work.

As a warmup exercise, you decide to tackle the following necessary subproblem. Recall that the *edit distance* between two strings is the minimum number of single-character insertions, deletions, and replacements required to transform one string into the other. An *arithmetic expression* is a string w such that

- w is a string of one or more decimal digits,
- $w = \mathbf{(x)}$ for some arithmetic expression x , or
- $w = x \diamond y$ for some arithmetic expressions x and y and some binary operator \diamond .

Suppose you are given a string of tokens from the alphabet $\{\#, \diamond, (,)\}$, where $\#$ represents a decimal digit and \diamond represents a binary operator. Describe an algorithm to compute the minimum edit distance from the given string to an arithmetic expression.

31. Let P be a set of points in the plane in *convex position*. Intuitively, if a rubber band were wrapped around the points, then every point would touch the rubber band. More formally, for any point p in P , there is a line that separates p from the other points in P . Moreover, suppose the points are indexed $P[1], P[2], \dots, P[n]$ in counterclockwise order around the 'rubber band', starting with the leftmost point $P[1]$.

This problem asks you to solve a special case of the traveling salesman problem, where the salesman must visit every point in P , and the cost of moving from one point $p \in P$ to another point $q \in P$ is the Euclidean distance $|pq|$.

- Describe a simple algorithm to compute the shortest *cyclic* tour of P .
 - A *simple* tour is one that never crosses itself. Prove that the shortest tour of P must be simple.
 - Describe and analyze an efficient algorithm to compute the shortest tour of P that starts at the leftmost point $P[1]$ and ends at the rightmost point $P[r]$.
 - Describe and analyze an efficient algorithm to compute the shortest tour of P , with no restrictions on the endpoints.
32. (a) Describe and analyze an efficient algorithm to determine, given a string w and a regular expression R , whether $w \in L(R)$.
- (b) *Generalized* regular expressions allow the binary operator \cap (intersection) and the unary operator \neg (complement), in addition to the usual \cdot (concatenation), $+$ (or), and $*$ (Kleene closure) operators. NFA constructions and Kleene's theorem imply that any generalized regular expression E represents a regular language $L(E)$.
- Describe and analyze an efficient algorithm to determine, given a string w and a generalized regular expression E , whether $w \in L(E)$.

In both problems, assume that you are actually given a parse tree for the (generalized) regular expression, not just a string.

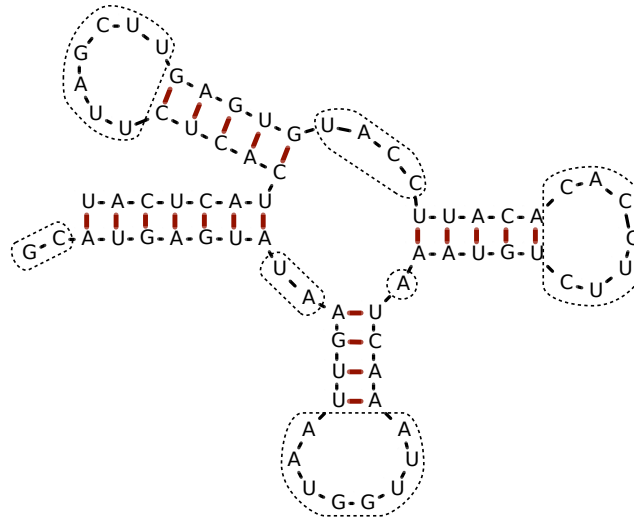
33. Ribonucleic acid (RNA) molecules are long chains of millions of nucleotides or *bases* of four different types: adenine (A), cytosine (C), guanine (G), and uracil (U). The *sequence* of an RNA molecule is a string $b[1..n]$, where each character $b[i] \in \{A, C, G, U\}$ corresponds to a base. In addition to the chemical bonds between adjacent bases in the sequence, hydrogen bonds can form between certain pairs of bases. The set of bonded base pairs is called the *secondary structure* of the RNA molecule.

We say that two base pairs (i, j) and (i', j') with $i < j$ and $i' < j'$ **overlap** if $i < i' < j < j'$ or $i' < i < j' < j$. In practice, most base pairs are non-overlapping. Overlapping base pairs create so-called *pseudoknots* in the secondary structure, which are essential for some RNA functions, but are more difficult to predict.

Suppose we want to predict the best possible secondary structure for a given RNA sequence. We will adopt a drastically simplified model of secondary structure:

- Each base can be paired with at most one other base.
- Only A-U pairs and C-G pairs can bond.
- Pairs of the form $(i, i + 1)$ and $(i, i + 2)$ cannot bond.
- Overlapping base pairs cannot bond.

The last restriction allows us to visualize RNA secondary structure as a sort of fat tree.



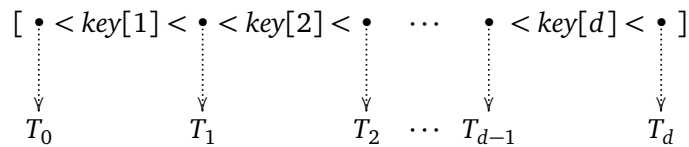
Example RNA secondary structure with 21 base pairs, indicated by heavy red lines. Gaps are indicated by dotted curves. This structure has score $2^2 + 2^2 + 8^2 + 1^2 + 7^2 + 4^2 + 7^2 = 187$

- (a) Describe and analyze an algorithm that computes the maximum possible *number* of bonded base pairs in a secondary structure for a given RNA sequence.
- (b) A *gap* in a secondary structure is a maximal substring of unpaired bases. Large gaps lead to chemical instabilities, so secondary structures with smaller gaps are more likely. To account for this preference, let's define the *score* of a secondary structure to be the sum of the *squares* of the gap lengths. (This score function is utterly fictional; real RNA structure prediction requires *much* more complicated scoring functions.) Describe and analyze an algorithm that computes the minimum possible score of a secondary structure for a given RNA sequence.
34. A standard method to improve the cache performance of search trees is to pack more search keys and subtrees into each node. A **B-tree** is a rooted tree in which each internal node stores up to B keys and pointers to up to $B + 1$ children, each the root of a smaller B -tree. Specifically, each node v stores three fields:
- a positive integer $v.d \leq B$,
 - a *sorted* array $v.key[1..v.d]$, and
 - an array $v.child[0..v.d]$ of child pointers.

In particular, the number of child pointers is always exactly one more than the number of keys.¹⁴

Each pointer $v.child[i]$ is either NULL or a pointer to the root of a B -tree whose keys are all larger than $v.key[i]$ and smaller than $v.key[i + 1]$. In particular, all keys in the leftmost subtree $v.child[0]$ are smaller than $v.key[1]$, and all keys in the rightmost subtree $v.child[v.d]$ are larger than $v.key[v.d]$.

Intuitively, you should have the following picture in mind:



Here T_i is the subtree pointed to by $child[i]$.

The **cost** of searching for a key x in a B -tree is the number of nodes in the path from the root to the node containing x as one of its keys. A 1-tree is just a standard binary search tree.

Fix an arbitrary positive integer $B > 0$. (I suggest $B = 8$.) Suppose you are given a sorted array $A[1, \dots, n]$ of search keys and a corresponding array $F[1, \dots, n]$ of frequency counts, where $F[i]$ is the number of times that we will search for $A[i]$. Your task is to describe and analyze an efficient algorithm to find a B -tree that minimizes the total cost of searching for the given keys with the given frequencies.

- Describe a polynomial-time algorithm for the special case $B = 2$.
- Describe an algorithm for arbitrary B that runs in $O(n^{B+c})$ time for some fixed integer c .
- Describe an algorithm for arbitrary B that runs in $O(n^c)$ time for some fixed integer c that does *not* depend on B .

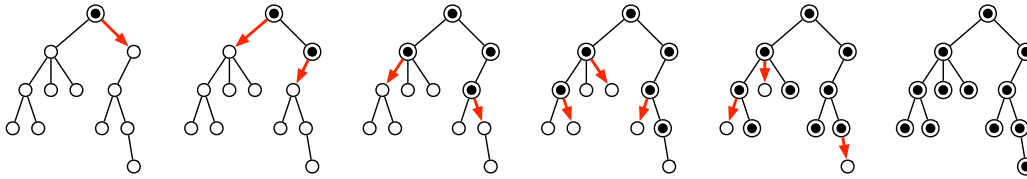
¹⁴**A few comments about B -trees.** Normally, B -trees are required to satisfy two additional constraints, which guarantee a worst-case search cost of $O(\log_B n)$: Every leaf must have exactly the same depth, and every node except possibly the root must contain at least $B/2$ keys. However, in this problem, we are not interested in optimizing the *worst-case* search cost, but rather the *total* cost of a sequence of searches, so we will not impose these additional constraints.

In most large database systems, the parameter B is chosen so that each node exactly fits in a cache line. Since the entire cache line is loaded into cache anyway, and the cost of loading a cache line exceeds the cost of searching within the cache, the running time is dominated by the number of cache faults. This effect is even more noticeable if the data is too big to fit in RAM; then the cost is dominated by the number of *page faults*, and B should be roughly the size of a page. In extreme cases, the data is too large even to fit on disk (or flash-memory “disk”) and is instead distributed on a bank of magnetic tape cartridges, in which case the cost is dominated by the number of *tape faults*. (I invite anyone who thinks tape is dead to visit a supercomputing center like Blue Waters.) In principle, your data might be so large that the cost of searching is actually dominated by the number of *FedEx faults*. (See <https://what-if.xkcd.com/31/>.)

Don’t worry about the cache/disk/tape/FedEx performance in your solutions; just analyze the CPU time as usual. Designing algorithms with few cache misses or page faults is a interesting pastime; simultaneously optimizing CPU time *and* cache misses *and* page faults *and* FedEx faults is a topic of active research. Sadly, this kind of design and analysis requires tools we won’t see in this class.

Trees and Subtrees

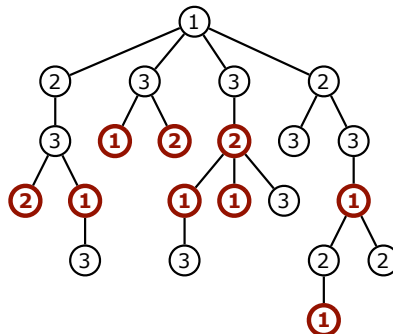
35. Suppose we need to distribute a message to all the nodes in a rooted tree. Initially, only the root node knows the message. In a single round, any node that knows the message can forward it to at most one of its children. Design an algorithm to compute the minimum number of rounds required for the message to be delivered to all nodes in a given tree.



A message being distributed through a tree in five rounds.

36. Oh, no! You have been appointed as the organizer of Giggle, Inc.’s annual mandatory holiday party! The employees at Giggle are organized into a strict hierarchy, that is, a tree with the company president at the root. The all-knowing oracles in Human Resources have assigned a real number to each employee measuring how “fun” the employee is. In order to keep things social, there is one restriction on the guest list: an employee cannot attend the party if their immediate supervisor is also present. On the other hand, the president of the company *must* attend the party, even though she has a negative fun rating; it’s her company, after all. Give an algorithm that makes a guest list for the party that maximizes the sum of the “fun” ratings of the guests.
37. Since so few people came to last year’s holiday party, the president of Giggle, Inc. decides to give each employee a present instead this year. Specifically, each employee must receive one of the three gifts: (1) an all-expenses-paid six-week vacation anywhere in the world, (2) an all-the-pancakes-you-can-eat breakfast for two at Jumping Jack Flash’s Flapjack Stack Shack, or (3) a burning paper bag full of dog poop. Corporate regulations prohibit any employee from receiving exactly the same gift as his/her direct supervisor. Any employee who receives a better gift than his/her direct supervisor will almost certainly be fired in a fit of jealousy.

As Giggle, Inc.’s official party czar, it’s *your* job to decide which gift each employee receives. Describe an algorithm to distribute gifts so that the minimum number of people are fired. Yes, you may send the president a flaming bag of dog poop.



A tree labeling with cost 9. The nine bold nodes have smaller labels than their parents. This is *not* the optimal labeling for this tree.

More formally, you are given a rooted tree T , representing the company hierarchy, and you want to label each node in T with an integer 1, 2, or 3, so that every node has a different label from its parent. The *cost* of an labeling is the number of nodes that have smaller labels than their parents. Describe and analyze an algorithm to compute the minimum cost of any labeling of the given tree T .

38. After the Flaming Dog Poop Holiday Debacle, you were strongly encouraged to seek other employment, and so you left Giggle for its competitor Yeehaw! Unfortunately, the new president of Yeehaw! just decided to imitate Giggle by throwing her own holiday party, and in light of your past experience, appointed you as the official party organizer. The president demands that you invite exactly k employees, including the president herself, and everyone who is invited is required to attend. Yeah, that'll be fun.

Just like at Giggle, employees at Yeehaw! are organized into a strict hierarchy: a tree with the company president at the root. The all-knowing oracles in Human Resources have assigned a real number to each employee indicating the *awkwardness* of inviting both that employee and their immediate supervisor; a negative value indicates that the employee and their supervisor actually like each other. Your goal is to choose a subset of exactly k employees to invite, so that the total awkwardness of the resulting party is as small as possible. For example, if the guest list does not include both an employee and their immediate supervisor, the total awkwardness is zero. The input to your algorithm is the tree T , the integer k , and the awkwardness of each node in T .

- (a) Describe an algorithm that computes the total awkwardness of the least awkward subset of k employees, assuming the company hierarchy is described by a *binary* tree. That is, assume that each employee directly supervises at most two others.
- * (b) Describe an algorithm that computes the total awkwardness of the least awkward subset of k employees, with no restrictions on the company hierarchy.

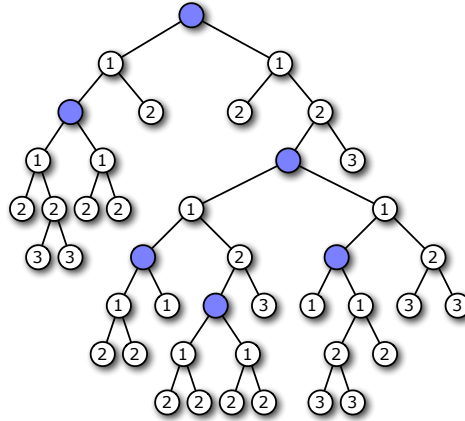
39. Let T be a rooted binary tree with n vertices, and let $k \leq n$ be a positive integer. We would like to mark k vertices in T so that every vertex has a nearby marked ancestor. More formally, we define the *clustering cost* of any subset K of vertices as

$$\text{cost}(K) = \max_v \text{cost}(v, K),$$

where the maximum is taken over all vertices v in the tree, and

$$\text{cost}(v, K) = \begin{cases} 0 & \text{if } v \in K \\ \infty & \text{if } v \text{ is the root of } T \text{ and } v \notin K \\ 1 + \text{cost}(\text{parent}(v)) & \text{otherwise} \end{cases}$$

Describe and analyze a dynamic-programming algorithm to compute the minimum clustering cost of any subset of k vertices in T . For full credit, your algorithm should run in $O(n^2k^2)$ time.



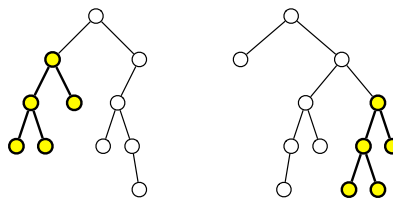
A subset of 5 vertices with clustering cost 3

To make the following problem statements precise, we must distinguish between several different types of trees and subtrees:

- By default, a *tree* is just a connected, acyclic, undirected graph.
- A *rooted tree* has a distinguished vertex called the *root*. A tree without a distinguished root vertex is called an *unrooted tree* or a *free tree*.
- In an *ordered tree*, the neighbors of every vertex have a well-defined cyclic order; a tree without these orders is called an *unordered tree*. Equivalently, each node in an ordered rooted tree has a *sequence* of children, which are the roots of ordered subtrees. In contrast, each node in an unordered rooted tree has a *set* of children, which are the roots of unordered subtrees.
- A *binary tree* is a rooted tree in which every node has a (possibly empty) *left* subtree and a (possibly empty) *right* subtree. Two binary trees are isomorphic if they are both empty, or if their left subtrees are isomorphic and their right subtrees are isomorphic.
- A *free subtree* of a tree is any connected subgraph; a *rooted subtree* consists of a node and all its descendants. By default, a subtree of an *unrooted tree* means a free subtree, and a subtree of a *rooted tree* means a rooted subtree. By default, subtrees of ordered rooted trees are themselves ordered trees.

40. This question asks you to find efficient algorithms to compute the **largest common rooted subtree** of two given rooted trees. The precise definition of “common” depends on which rooted trees we consider to be isomorphic.

(a) Describe an algorithm to find the largest common *binary* subtree of two given *binary* trees.



Two binary trees, with their largest common (rooted) subtree emphasized

your algorithm should return the largest rooted minor M such that every node in M has a smaller label than its children in M .

- (c) Suppose we are given a *binary tree* T whose nodes are labeled with numbers. Describe an algorithm to find the largest *binary-search-ordered rooted minor* of T . That is, your algorithm should return a rooted minor M such that every node in M has at most two children, and an inorder traversal of M is an increasing subsequence of an inorder traversal of T .
- (d) Recall that a rooted tree is *ordered* if the children of each node have a well-defined left-to-right order. Describe an algorithm to find the largest binary-search-ordered minor of an *arbitrary ordered tree* T whose nodes are labeled with numbers. Again, the left-to-right order of nodes in M should be consistent with their order in T .
- * (e) Describe an algorithm to find the largest common *ordered* rooted minor of two *ordered* labeled rooted trees.
- ★ (f) Describe an algorithm to find the largest common *unordered* rooted minor of two *unordered* labeled rooted trees. [Hint: This problem will be much easier after you've seen flows.]