an expression which involves only smaller subproblems. How long does this step take? It requires the predecessors of $j$ to be known; for this the adjacency list of the reverse graph $G^R$, constructible in linear time (recall Exercise 3.5), is handy. The computation of $L(j)$ then takes time proportional to the indegree of $j$, giving an overall running time linear in $|E|$. This is at most $O(n^2)$, the maximum being when the input array is sorted in increasing order. Thus the dynamic programming solution is both simple and efficient.

There is one last issue to be cleared up: the $L$-values only tell us the *length* of the optimal subsequence, so how do we recover the subsequence itself? This is easily managed with the same bookkeeping device we used for shortest paths in Chapter 4. While computing $L(j)$, we should also note down prev$(j)$, the next-to-last node on the longest path to $j$. The optimal subsequence can then be reconstructed by following these backpointers.

## 6.3 Edit distance

When a spell checker encounters a possible misspelling, it looks in its dictionary for other words that are close by. What is the appropriate notion of closeness in this case?

A natural measure of the distance between two strings is the extent to which they can be *aligned*, or matched up. Technically, an alignment is simply a way of writing the strings one above the other. For instance, here are two possible alignments of SNOWY and SUNNY:

```
        S  —  N  O  W  Y              —  S  N  O  W  —  Y
        S  U  N  N  —  Y              S  U  N  —  —  N  Y
             Cost: 3                        Cost: 5
```

The "−" indicates a "gap"; any number of these can be placed in either string. The *cost* of an alignment is the number of columns in which the letters differ. And the *edit distance* between two strings is the cost of their best possible alignment. Do you see that there is no better alignment of SNOWY and SUNNY than the one shown here with a cost of 3?

Edit distance is so named because it can also be thought of as the minimum number of *edits*—insertions, deletions, and substitutions of characters—needed to transform the first string into the second. For instance, the alignment shown on the left corresponds to three edits: insert U, substitute O → N, and delete W.

In general, there are so many possible alignments between two strings that it would be terribly inefficient to search through all of them for the best one. Instead we turn to dynamic programming.

### A dynamic programming solution

When solving a problem by dynamic programming, the most crucial question is, *What are the subproblems?* As long as they are chosen so as to have the property
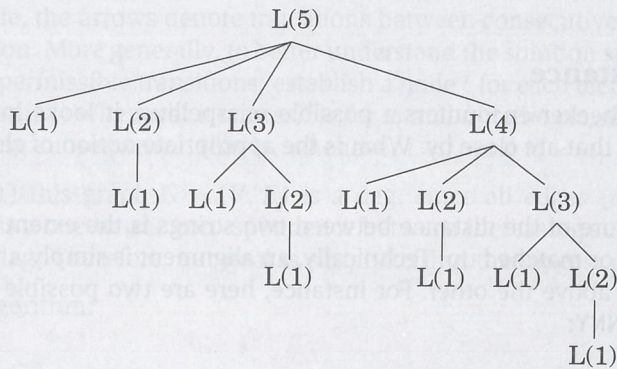
## Recursion? No, thanks.

Returning to our discussion of longest increasing subsequences: the formula for $L(j)$ also suggests an alternative, recursive algorithm. Wouldn't that be even simpler?

Actually, recursion is a very bad idea: the resulting procedure would require exponential time! To see why, suppose that the dag contains edges $(i, j)$ for *all* $i < j$—that is, the given sequence of numbers $a_1, a_2, \ldots, a_n$ is sorted. In that case, the formula for subproblem $L(j)$ becomes

$$L(j) = 1 + \max\{L(1), L(2), \ldots, L(j-1)\}.$$

The following figure unravels the recursion for $L(5)$. Notice that the same subproblems get solved over and over again!



For $L(n)$ this tree has exponentially many nodes (can you bound it?), and so a recursive solution is disastrous.

Then why did recursion work so well with divide-and-conquer? The key point is that in divide-and-conquer, a problem is expressed in terms of subproblems that are *substantially smaller*, say half the size. For instance, mergesort sorts an array of size $n$ by recursively sorting two subarrays of size $n/2$. Because of this sharp drop in problem size, the full recursion tree has only logarithmic depth and a polynomial number of nodes.

In contrast, in a typical dynamic programming formulation, a problem is reduced to subproblems that are only slightly smaller—for instance, $L(j)$ relies on $L(j-1)$. Thus the full recursion tree generally has polynomial depth and an exponential number of nodes. However, it turns out that most of these nodes are repeats, that there are not too many *distinct* subproblems among them. Efficiency is therefore obtained by explicitly enumerating the distinct subproblems and solving them in the right order.

> **Programming?**
> _____
>
> The origin of the term *dynamic programming* has very little to do with writing code. It was first coined by Richard Bellman in the 1950s, a time when computer programming was an esoteric activity practiced by so few people as to not even merit a name. Back then programming meant "planning," and "dynamic programming" was conceived to optimally plan multistage processes. The dag of Figure 6.2 can be thought of as describing the possible ways in which such a process can evolve: each node denotes a state, the leftmost node is the starting point, and the edges leaving a state represent possible actions, leading to different states in the next unit of time.
>
> The etymology of *linear programming*, the subject of Chapter 7, is similar.

(*) from page 158. it is an easy matter to write down the algorithm: iteratively solve one subproblem after the other, in order of increasing size.

Our goal is to find the edit distance between two strings $x[1 \cdots m]$ and $y[1 \cdots n]$. What is a good subproblem? Well, it should go part of the way toward solving the whole problem; so how about looking at the edit distance between some *prefix* of the first string, $x[1 \cdots i]$, and some *prefix* of the second, $y[1 \cdots j]$? Call this subproblem $E(i, j)$ (see Figure 6.3). Our final objective, then, is to compute $E(m, n)$.

For this to work, we need to somehow express $E(i, j)$ in terms of smaller subproblems. Let's see—what do we know about the best alignment between $x[1 \cdots i]$ and $y[1 \cdots j]$? Well, its rightmost column can only be one of three things:

$$\begin{matrix} x[i] \\ - \end{matrix} \quad \text{or} \quad \begin{matrix} - \\ y[j] \end{matrix} \quad \text{or} \quad \begin{matrix} x[i] \\ y[j] \end{matrix}$$

The first case incurs a cost of 1 for this particular column, and it remains to align $x[1 \cdots i-1]$ with $y[1 \cdots j]$. But this is exactly the subproblem $E(i-1, j)$! We seem to be getting somewhere. In the second case, also with cost 1, we still need to align $x[1 \cdots i]$ with $y[1 \cdots j-1]$. This is again another subproblem, $E(i, j-1)$. And in the final case, which either costs 1 (if $x[i] \neq y[j]$) or 0 (if $x[i] = y[j]$), what's left is the subproblem $E(i-1, j-1)$. In short, we have expressed $E(i, j)$ in terms of

**Figure 6.3**   The subproblem $E(7, 5)$.

| E | X | P | O | N | E | N | T | I | A | L |

| P | O | L | Y | N | O | M | I | A | L |

three *smaller* subproblems $E(i-1, j)$, $E(i, j-1)$, $E(i-1, j-1)$. We have no idea which of them is the right one, so we need to try them all and pick the best:

$$E(i, j) = \min\{1 + E(i-1, j),\ 1 + E(i, j-1),\ \text{diff}(i, j) + E(i-1, j-1)\}$$

where for convenience $\text{diff}(i, j)$ is defined to be 0 if $x[i] = y[j]$ and 1 otherwise.

For instance, in computing the edit distance between EXPONENTIAL and POLYNOMIAL, subproblem $E(4, 3)$ corresponds to the prefixes EXPO and POL. The rightmost column of their best alignment must be one of the following:

$$\begin{matrix} \text{0} \\ \text{\textemdash} \end{matrix} \quad \text{or} \quad \begin{matrix} \text{\textemdash} \\ \text{L} \end{matrix} \quad \text{or} \quad \begin{matrix} \text{0} \\ \text{L} \end{matrix}$$
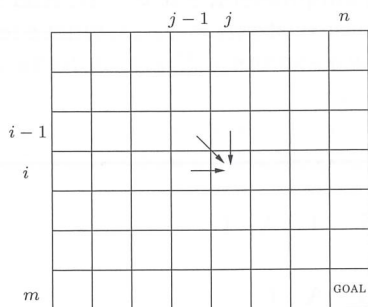
Thus, $E(4, 3) = \min\{1 + E(3, 3),\ 1 + E(4, 2),\ 1 + E(3, 2)\}$.

The answers to all the subproblems $E(i, j)$ form a two-dimensional table, as in Figure 6.4. In what order should these subproblems be solved? Any order is fine, as long as $E(i-1, j)$, $E(i, j-1)$, and $E(i-1, j-1)$ are handled *before* $E(i, j)$. For instance, we could fill in the table one row at a time, from top row to bottom row, and moving left to right across each row. Or alternatively, we could fill it in column by column. Both methods would ensure that by the time we get around to computing a particular table entry, all the other entries we need are already filled in.

With both the subproblems and the ordering specified, we are almost done. There just remain the "base cases" of the dynamic programming, the very smallest subproblems. In the present situation, these are $E(0, \cdot)$ and $E(\cdot, 0)$, both of which are easily solved. $E(0, j)$ is the edit distance between the 0-length prefix of $x$,

**Figure 6.4**   (a) The table of subproblems. Entries $E(i-1, j-1)$, $E(i-1, j)$, and $E(i, j-1)$ are needed to fill in $E(i, j)$. (b) The final table of values found by dynamic programming.

(a)



(b)

|   |    | P | O | L | Y | N | O | M | I | A | L |
|---|----|---|---|---|---|---|---|---|---|---|---|
|   | 0  | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| E | 1  | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| X | 2  | 2 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| P | 3  | 2 | 3 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| O | 4  | 3 | 2 | 3 | 4 | 5 | 5 | 6 | 7 | 8 | 9 |
| N | 5  | 4 | 3 | 3 | 4 | 4 | 5 | 6 | 7 | 8 | 9 |
| E | 6  | 5 | 4 | 4 | 4 | 5 | 5 | 6 | 7 | 8 | 9 |
| N | 7  | 6 | 5 | 5 | 5 | 4 | 5 | 6 | 7 | 8 | 9 |
| T | 8  | 7 | 6 | 6 | 6 | 5 | 5 | 6 | 7 | 8 | 9 |
| I | 9  | 8 | 7 | 7 | 7 | 6 | 6 | 6 | 6 | 7 | 8 |
| A | 10 | 9 | 8 | 8 | 8 | 7 | 7 | 7 | 7 | 6 | 7 |
| L | 11 | 10 | 9 | 8 | 9 | 8 | 8 | 8 | 8 | 7 | 6 |

namely the empty string, and the first $j$ letters of $y$: clearly, $j$. And similarly, $E(i, 0) = i$.

At this point, the algorithm for edit distance basically writes itself.

```
for i = 0, 1, 2, ..., m:
    E(i, 0) = i
for j = 1, 2, ..., n:
    E(0, j) = j
for i = 1, 2, ..., m:
    for j = 1, 2, ..., n:
        E(i, j) = min{E(i − 1, j) + 1, E(i, j − 1) + 1, E(i − 1, j − 1) + diff(i, j)}
return E(m, n)
```

This procedure fills in the table row by row, and left to right within each row. Each entry takes constant time to fill in, so the overall running time is just the size of the table, $O(mn)$.

And in our example, the edit distance turns out to be 6:

```
E   X   P   O   N   E   N   −   T   I   A   L
−   −   P   O   L   Y   N   O   M   I   A   L
```
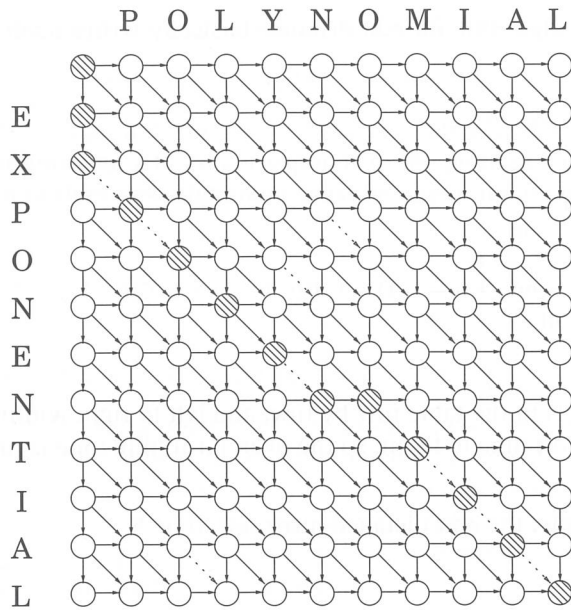
**The underlying dag**

Every dynamic program has an underlying dag structure: think of each node as representing a subproblem, and each edge as a precedence constraint on the order in which the subproblems can be tackled. Having nodes $u_1, \ldots, u_k$ point to $v$ means "subproblem $v$ can only be solved once the answers to $u_1, \ldots, u_k$ are known."

In our present edit distance application, the nodes of the underlying dag correspond to subproblems, or equivalently, to positions $(i, j)$ in the table. Its edges are the precedence constraints, of the form $(i − 1, j) \to (i, j)$, $(i, j − 1) \to (i, j)$, and $(i − 1, j − 1) \to (i, j)$ (Figure 6.5). In fact, we can take things a little further and put weights on the edges so that the edit distances are given by shortest paths in the dag! To see this, set all edge lengths to 1, except for $\{(i − 1, j − 1) \to (i, j) : x[i] = y[j]\}$ (shown dotted in the figure), whose length is 0. The final answer is then simply the distance between nodes $s = (0, 0)$ and $t = (m, n)$. One possible shortest path is shown, the one that yields the alignment we found earlier. On this path, each move down is a deletion, each move right is an insertion, and each diagonal move is either a match or a substitution.

By altering the weights on this dag, we can allow generalized forms of edit distance, in which insertions, deletions, and substitutions have different associated costs.

**Figure 6.5** The underlying dag, and a path of length 6.



## 6.4 Knapsack

During a robbery, a burglar finds much more loot than he had expected and has to decide what to take. His bag (or "knapsack") will hold a total weight of at most $W$ pounds. There are $n$ items to pick from, of weight $w_1, \ldots, w_n$ and dollar value $v_1, \ldots, v_n$. What's the most valuable combination of items he can fit into his bag?[1]

For instance, take $W = 10$ and

| Item | Weight | Value |
|------|--------|-------|
| 1 | 6 | $30 |
| 2 | 3 | $14 |
| 3 | 4 | $16 |
| 4 | 2 | $9 |

---

[1]If this application seems frivolous, replace "weight" with "CPU time" and "only $W$ pounds can be taken" with "only $W$ units of CPU time are available." Or use "bandwidth" in place of "CPU time," etc. The knapsack problem generalizes a wide variety of resource-constrained selection tasks.