

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.

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 \rightarrow 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 (*) from page 171, 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

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 | |

about looking at the edit distance between some *prefix* of the first string, $x[1 \dots i]$, and some *prefix* of the second, $y[1 \dots j]$? 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 $x[1 \dots i]$ and $y[1 \dots j]$? Well, its rightmost column can only be one of three things:

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

The first case incurs a cost of 1 for this particular column, and it remains to align $x[1 \dots i-1]$ with $y[1 \dots j]$. 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 $x[1 \dots i]$ with $y[1 \dots j-1]$. 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 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:

$$\begin{array}{c} O \\ - \end{array} \quad \text{or} \quad \begin{array}{c} - \\ L \end{array} \quad \text{or} \quad \begin{array}{c} O \\ L \end{array}$$

$$\text{Thus, } E(4, 3) = \min\{1 + E(3, 3), 1 + E(4, 2), 1 + E(3, 2)\}.$$

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.

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)

| | | | | | | | |
|-----|--|--|-----|---|--|--|------|
| | | | j-1 | j | | | n |
| | | | | | | | |
| | | | | | | | |
| i-1 | | | | | | | |
| i | | | | | | | |
| | | | | | | | |
| | | | | | | | |
| m | | | | | | | GOAL |

(b)

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

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 , 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)

```

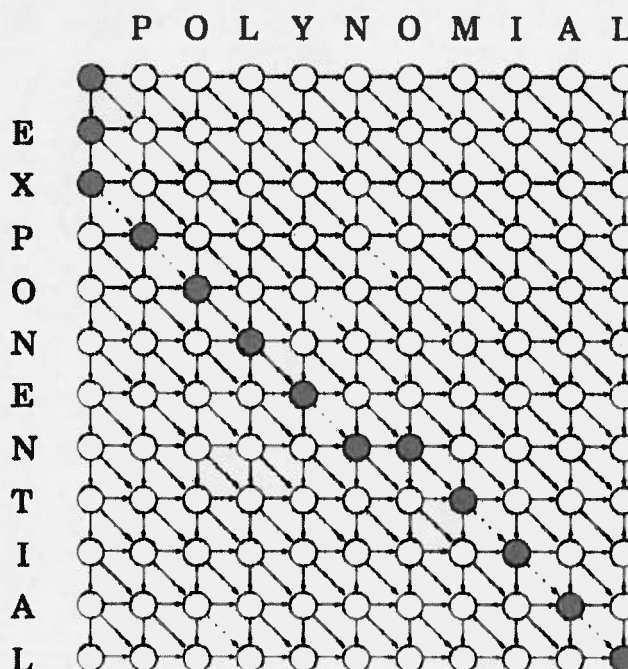
This procedure fills in the table row by row, and left to right within each row. Each entry takes constant time to fill in, so the overall running time is just the size of the table, $O(mn)$.

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

```

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

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, \dots, u_k point to v means “subproblem v can only be solved once the answers to u_1, \dots, 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 $(i-1, j) \rightarrow (i, j)$, $(i, j-1) \rightarrow (i, j)$, and $(i-1, j-1) \rightarrow (i, j)$ (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 $\{(i-1, j-1) \rightarrow (i, j) : x[i] = y[j]\}$ (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.

