



**Figure 4.21** An environment in which a random walk will take exponentially many steps to find the goal.

estimate  $h(s)$  and is updated as the agent gains experience in the state space. Figure 4.22 shows a simple example in a one-dimensional state space. In (a), the agent seems to be stuck in a flat local minimum at the shaded state. Rather than staying where it is, the agent should follow what seems to be the best path to the goal based on the current cost estimates for its neighbors. The estimated cost to reach the goal through a neighbor  $s'$  is the cost to get to  $s'$  plus the estimated cost to get to a goal from there—that is,  $c(s, a, s') + H(s')$ . In the example, there are two actions with estimated costs  $1 + 9$  and  $1 + 2$ , so it seems best to move right. Now, it is clear that the cost estimate of 2 for the shaded state was overly optimistic. Since the best move cost 1 and led to a state that is at least 2 steps from a goal, the shaded state must be at least 3 steps from a goal, so its  $H$  should be updated accordingly, as shown in Figure 4.22(b). Continuing this process, the agent will move back and forth twice more, updating  $H$  each time and “flattening out” the local minimum until it escapes to the right.

An agent implementing this scheme, which is called learning real-time A\* (**LRTA\***), is shown in Figure 4.23. Like **ONLINE-DFS-AGENT**, it builds a map of the environment using the *result* table. It updates the cost estimate for the state it has just left and then chooses the “apparently best” move according to its current cost estimates. One important detail is that actions that have not yet been tried in a state  $s$  are always assumed to lead immediately to the goal with the least possible cost, namely  $h(s)$ . This **optimism under uncertainty** encourages the agent to explore new, possibly promising paths.

An **LRTA\*** agent is guaranteed to find a goal in any finite, safely explorable environment. Unlike **A\***, however, it is not complete for infinite state spaces—there are cases where it can be led infinitely astray. It can explore an environment of  $n$  states in  $O(n^2)$  steps in the worst case, but often does much better. The **LRTA\*** agent is just one of a large family of online agents that can be defined by specifying the action selection rule and the update rule in different ways. We will discuss this family, which was developed originally for stochastic environments, in Chapter 21.

### Learning in online search

The initial ignorance of online search agents provides several challenges. As the agents learn a “map” of the environment—more precisely, they learn the cost of each state—simply by recording each of their experiences. (In deterministic environments means that one experience is enough to learn the cost of a state.) As the local search agents acquire more accurate estimates of the value of each state, the local updating rules, as in **LRTA\***. In Chapter 21 we will see that