

9 *Online Planning*

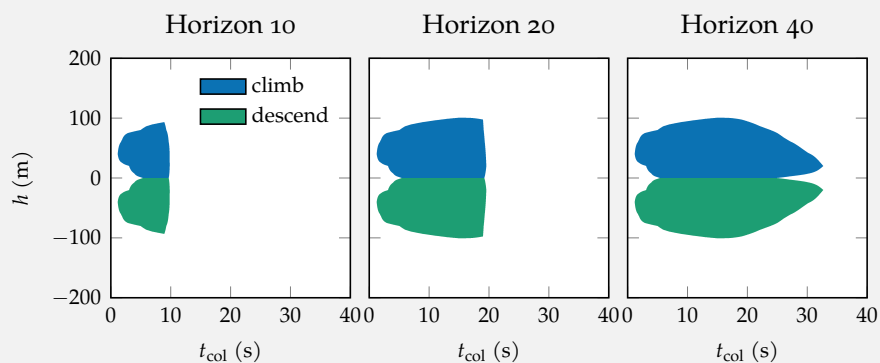
The solution methods we have discussed so far compute policies offline, before any actions are executed in the real problem. Even offline approximation methods can be intractable in many high-dimensional problems. This chapter discusses *online planning* methods that find actions based on reasoning about states that are reachable from the current state. The *reachable state space* is often orders of magnitude smaller than the full state space, which can significantly reduce storage and computational requirements compared to offline methods. We will discuss a variety of algorithms that aim to make online planning efficient, including pruning the state space, sampling, and planning more deeply along trajectories that appear more promising.

9.1 *Receding Horizon Planning*

In *receding horizon planning*, we plan from the current state to a maximum fixed horizon or depth d . We then execute the action from our current state, transition to the next state, and replan. The online planning methods discussed in this chapter follow this receding horizon planning scheme. They differ in how they explore different courses of action.

A challenge in applying receding horizon planning is determining the appropriate depth. Deeper planning generally requires more computation. For some problems, a shallow depth can be quite effective; the fact that we replan at each step can compensate for our lack of longer-term modeling. In other problems, greater planning depths may be necessary so that our planner can be driven toward goals or away from unsafe states, as illustrated in example 9.1.

Suppose we want to apply receding horizon planning to aircraft collision avoidance. The objective is to provide descend or climb advisories when necessary to avoid collision. A collision occurs when our altitude relative to the intruder h is within ± 50 m and the time to potential collision t_{col} is zero. We want to plan deeply enough so that we can provide an advisory sufficiently early to avoid collision with a high degree of confidence. The plots here show the actions that would be taken by a receding horizon planner with different depths.



If the depth is $d = 10$, we provide advisories only within 10 s of collision. Due to the limitations of the vehicle dynamics and the uncertainty of the behavior of the other aircraft, providing advisories this late compromises safety. With $d = 20$, we can do better, but there are cases where we would want to alert a little earlier to further reduce collision risk. There is no motivation to plan deeper than $d = 40$ because we do not need to advise any maneuvers that far ahead of potential collision.

Example 9.1. Receding horizon planning for collision avoidance to different planning depths. In this problem, there are four state variables. These plots show slices of the state space under the assumption that the aircraft is currently level and there has not yet been an advisory. The horizontal axis is the time to collision t_{col} , and the vertical axis is our altitude h relative to the intruder. Appendix F.6 provides additional details about this problem.

9.2 Lookahead with Rollouts

Chapter 8 involved extracting policies that are greedy with respect to an approximate value function U through the use of one-step lookahead.¹ A simple online strategy involves acting greedily with respect to values estimated through simulation to depth d . To run a simulation, we need a policy to simulate. Of course, we do not know the optimal policy, but we can use what is called a *rollout policy* instead. Rollout policies are typically stochastic, with actions drawn from a distribution $a \sim \pi(s)$. To produce these rollout simulations, we use a *generative model* $s' \sim T(s, a)$ to generate successor states s' from the distribution $T(s' | s, a)$. This generative model can be implemented through draws from a random number generator, which can be easier to implement in practice compared to explicitly representing the distribution $T(s' | s, a)$.

Algorithm 9.1 combines one-step lookahead with values estimated through rollout. This approach often results in better behavior than that of the original rollout policy, but optimality is not guaranteed. It can be viewed as an approximate form of policy improvement used in the policy iteration algorithm (section 7.4). A simple variation of this algorithm is to use multiple rollouts to arrive at a better estimate of the expected discounted return. If we run m simulations for each action and resulting state, the time complexity is $O(m \times |\mathcal{A}| \times |\mathcal{S}| \times d)$.

9.3 Forward Search

Forward search determines the best action to take from an initial state s by expanding all possible transitions up to depth d . These expansions form a *search tree*.² Such search trees have a worst-case branching factor of $|\mathcal{S}| \times |\mathcal{A}|$, yielding a computational complexity of $O((|\mathcal{S}| \times |\mathcal{A}|)^d)$. Figure 9.1 shows a search tree applied to a problem with three states and two actions. Figure 9.2 visualizes the states visited during forward search on the hex world problem.

Algorithm 9.2 calls itself recursively to the specified depth. Once reaching the specified depth, it uses an estimate of the utility provided by the function U . If we simply want to plan to the specified horizon, we set $U(s) = 0$. If our problem requires planning beyond the depth that we can afford to compute online, we can use an estimate of the value function obtained offline using, for example, one of the value function approximations described in the previous chapter. Combining

¹ The lookahead strategy was originally introduced in algorithm 7.2 as part of our discussion of exact solution methods.

² The exploration of the tree occurs as a *depth-first search*. Appendix E reviews depth-first search and other standard search algorithms in the deterministic context.

```

struct RolloutLookahead
   $\mathcal{P}$  # problem
   $\pi$  # rollout policy
  d # depth
end

randstep( $\mathcal{P}$ ::MDP, s, a) =  $\mathcal{P}$ .TR(s, a)

function rollout( $\mathcal{P}$ , s,  $\pi$ , d)
  ret = 0.0
  for t in 1:d
    a =  $\pi$ (s)
    s, r = randstep( $\mathcal{P}$ , s, a)
    ret +=  $\mathcal{P}$ . $\gamma$ ^(t-1) * r
  end
  return ret
end

function ( $\pi$ ::RolloutLookahead)(s)
  U(s) = rollout( $\pi$ . $\mathcal{P}$ , s,  $\pi$ . $\pi$ ,  $\pi$ .d)
  return greedy( $\pi$ . $\mathcal{P}$ , U, s).a
end

```

Algorithm 9.1. A function that runs a rollout of policy π in problem \mathcal{P} from state s to depth d . It returns the total discounted reward. This function can be used with the `greedy` function (introduced in algorithm 7.5) to generate an action that is likely to be an improvement over the original rollout policy. We will use this algorithm later for problems other than MDPs, requiring us to only have to modify `randstep` appropriately.

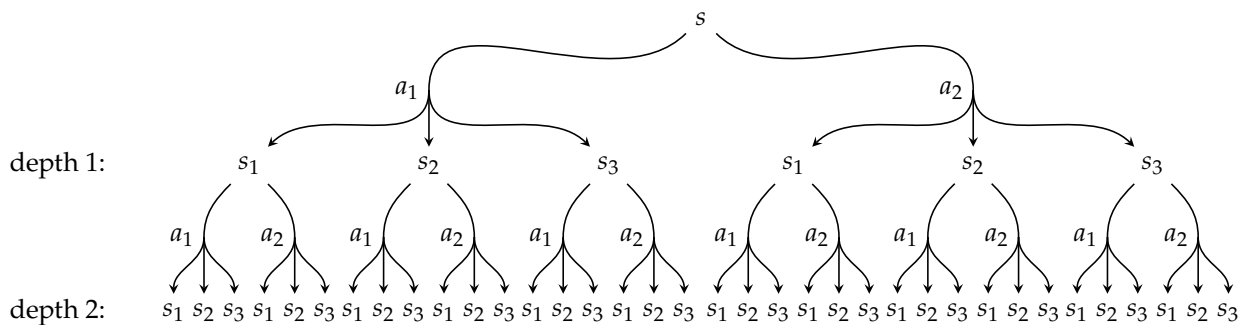


Figure 9.1. A forward search tree for a problem with three states and two actions.

online and offline approaches in this way is sometimes referred to as *hybrid planning*.

```

struct ForwardSearch
   $\mathcal{P}$  # problem
  d # depth
  U # value function at depth d
end

function forward_search( $\mathcal{P}$ , s, d, U)
  if d  $\leq$  0
    return (a=nothing, u=U(s))
  end
  best = (a=nothing, u=-Inf)
  U'(s) = forward_search( $\mathcal{P}$ , s, d-1, U).u
  for a in  $\mathcal{P}.\mathcal{A}$ 
    u = lookahead( $\mathcal{P}$ , U', s, a)
    if u > best.u
      best = (a=a, u=u)
    end
  end
  return best
end

( $\pi$ ::ForwardSearch)(s) = forward_search( $\pi.\mathcal{P}$ , s,  $\pi.d$ ,  $\pi.U$ ).a

```

Algorithm 9.2. The forward search algorithm for finding an approximately optimal action online for a problem \mathcal{P} from a current state s . The search is performed to depth d , at which point the terminal value is estimated with an approximate value function U . The returned named tuple consists of the best action a and its finite-horizon expected value u . The problem type is not constrained to be an MDP; section 22.2 uses this same algorithm in the context of partially observable problems with a different implementation for `lookahead`.

9.4 Branch and Bound

Branch and bound (algorithm 9.3) attempts to avoid the exponential computational complexity of forward search. It prunes branches by reasoning about bounds on the value function. The algorithm requires knowing a lower bound on the value function $\underline{U}(s)$ and an upper bound on the action value function $\overline{Q}(s, a)$. The lower bound is used to evaluate the states at the maximum depth. This lower bound is propagated upward through the tree through Bellman updates. If we find that the upper bound of an action at a state is lower than the lower bound of a previously explored action from that state, then we need not explore that action, allowing us to *prune* the associated subtree from consideration.

Branch and bound will give the same result as forward search, but it can be more efficient depending on how many branches are pruned. The worst-case complexity of branch and bound is still the same as forward search. To facilitate

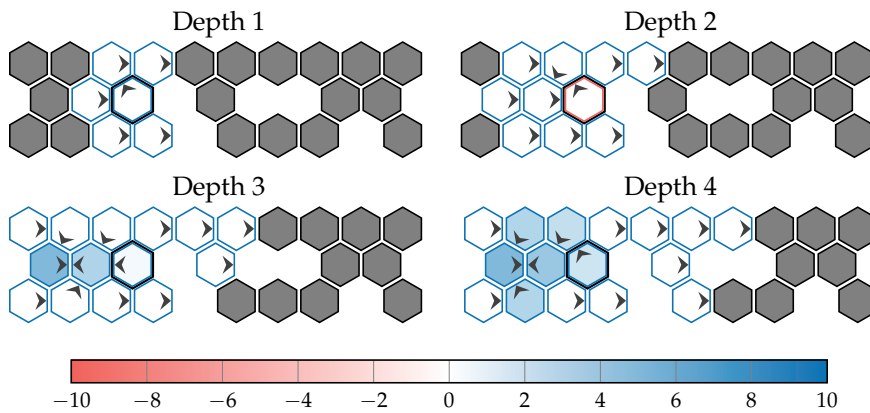


Figure 9.2. Forward search applied to the hex world problem with four maximum depths. The search can visit a node multiple times. The actions and colors for visited states were chosen according to the shallowest, highest-value node in the search tree for that state. The initial state has an additional black border.

```

struct BranchAndBound
   $\mathcal{P}$  # problem
  d # depth
  Ulo # lower bound on value function at depth d
  Qhi # upper bound on action value function
end

function branch_and_bound( $\mathcal{P}$ , s, d, Ulo, Qhi)
  if d ≤ 0
    return (a=nothing, u=Ulo(s))
  end
  U'(s) = branch_and_bound( $\mathcal{P}$ , s, d-1, Ulo, Qhi).u
  best = (a=nothing, u=-Inf)
  for a in sort( $\mathcal{P}.\mathcal{A}$ , by=a→Qhi(s,a), rev=true)
    if Qhi(s, a) < best.u
      return best # safe to prune
    end
    u = lookahead( $\mathcal{P}$ , U', s, a)
    if u > best.u
      best = (a=a, u=u)
    end
  end
  return best
end

( $\pi$ ::BranchAndBound)(s) = branch_and_bound( $\pi.\mathcal{P}$ , s,  $\pi.d$ ,  $\pi.Ulo$ ,  $\pi.Qhi$ ).a

```

Algorithm 9.3. The branch and bound algorithm for finding an approximately optimal action online for a discrete MDP \mathcal{P} from a current state s . The search is performed to depth d with value function lower bound Ulo and action value function upper bound Qhi . The returned named tuple consists of the best action a and its finite-horizon expected value u . This algorithm is also used for POMDPs.

pruning, actions are traversed in descending order by upper bound. Tighter bounds will generally result in more pruning, as shown in example 9.2.

Consider applying branch and bound to the mountain car problem. We can use the value function of a heuristic policy for the lower bound $\underline{U}(s)$, such as a heuristic policy that always accelerates in the direction of motion. For our upper bound $\overline{Q}([x, v], a)$, we can use the return expected when accelerating toward the goal with no hill. Branch and bound visits about a third as many states as forward search.

Example 9.2. Branch and bound applied to the mountain car problem (appendix F.4). Branch and bound can achieve a significant speedup over forward search.

9.5 Sparse Sampling

A method known as *sparse sampling*³ (algorithm 9.4) attempts to reduce the branching factor of forward search and branch and bound. Instead of branching on all possible next states, we consider only a limited number of samples of the next state. Although the sampling of the next state results in an approximation, this method can work well in practice and can significantly reduce computation. If we draw m samples of the next state for each action node in the search tree, the computational complexity is $O((m \times |\mathcal{A}|)^d)$, which is still exponential in the depth but no longer depends on the size of the state space. Figure 9.3 shows an example.

³M.J. Kearns, Y. Mansour, and A. Y. Ng, "A Sparse Sampling Algorithm for Near-Optimal Planning in Large Markov Decision Processes," *Machine Learning*, vol. 49, no. 2–3, pp. 193–208, 2002.

9.6 Monte Carlo Tree Search

Monte Carlo tree search (algorithm 9.5) avoids the exponential complexity in the horizon by running m simulations from the current state.⁴ During these simulations, the algorithm updates estimates of the action value function $Q(s, a)$ and a record of the number of times a particular state-action pair has been selected, $N(s, a)$. After running these m simulations from our current state s , we simply choose the action that maximizes our estimate of $Q(s, a)$.

A simulation (algorithm 9.6) begins by traversing the explored state space, consisting of the states for which we have estimates of Q and N . We follow an exploration strategy to choose actions from the various states. A common approach is to select the action that maximizes the *UCB1 exploration heuristic*:⁵

⁴For a survey, see C. B. Browne, E. Powley, D. Whitehouse, S. M. Lucas, P. I. Cowling, P. Rohlfshagen, S. Tavener, D. Perez, S. Samothrakis, and S. Colton, "A Survey of Monte Carlo Tree Search Methods," *IEEE Transactions on Computational Intelligence and AI in Games*, vol. 4, no. 1, pp. 1–43, 2012.

⁵UCB stands for upper confidence bound. This is one of many strategies discussed by P. Auer, N. Cesa-Bianchi, and P. Fischer, "Finite-Time Analysis of the Multiarmed Bandit Problem," *Machine Learning*, vol. 47, no. 2–3, pp. 235–256, 2002. The equation is derived from the Chernoff-Hoeffding bound.

```

struct SparseSampling
   $\mathcal{P}$  # problem
  d # depth
  m # number of samples
  U # value function at depth d
end

function sparse_sampling( $\mathcal{P}$ , s, d, m, U)
  if d  $\leq$  0
    return (a=nothing, u=U(s))
  end
  best = (a=nothing, u=-Inf)
  for a in  $\mathcal{P}.\mathcal{A}$ 
    u = 0.0
    for i in 1:m
      s', r = randstep( $\mathcal{P}$ , s, a)
      a', u' = sparse_sampling( $\mathcal{P}$ , s', d-1, m, U)
      u += (r +  $\mathcal{P}.\gamma*u'$ ) / m
    end
    if u > best.u
      best = (a=a, u=u)
    end
  end
  return best
end

( $\pi$ ::SparseSampling)(s) = sparse_sampling( $\pi.\mathcal{P}$ , s,  $\pi.d$ ,  $\pi.m$ ,  $\pi.U$ ).a

```

Algorithm 9.4. The sparse sampling algorithm for finding an approximately optimal action online for a discrete problem \mathcal{P} from a current state s to depth d with m samples per action. The returned named tuple consists of the best action a and its finite-horizon expected value u .

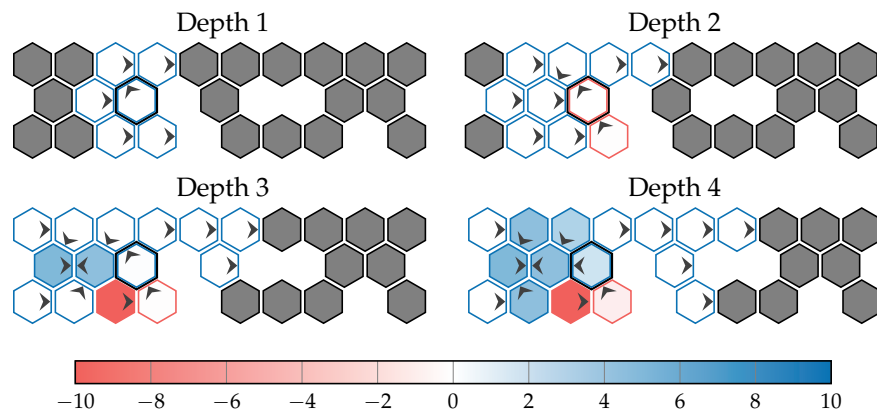


Figure 9.3. Sparse sampling with $m = 10$ applied to the hex world problem. Visited tiles are colored according to their estimated value. The bordered tile is the initial state. Compare to forward search in figure 9.2.


```

struct MonteCarloTreeSearch
  P # problem
  N # visit counts
  Q # action value estimates
  d # depth
  m # number of simulations
  c # exploration constant
  U # value function estimate
end

function (π::MonteCarloTreeSearch)(s)
  for k in 1:π.m
    simulate!(π, s)
  end
  return argmax(a→π.Q[(s,a)], π.P.A)
end

```

Algorithm 9.5. The Monte Carlo tree search policy for finding an approximately optimal action from a current state s .

$$Q(s, a) + c \sqrt{\frac{\log N(s)}{N(s, a)}} \quad (9.1)$$

where $N(s) = \sum_a N(s, a)$ is the total visit count to s and c is an exploration parameter that scales the value of unexplored actions. The second term corresponds to an *exploration bonus*. If $N(s, a) = 0$, the bonus is defined to be infinity. With $N(s, a)$ in the denominator, the exploration bonus is higher for actions that have not been tried as frequently. Algorithm 9.7 implements this exploration strategy. We will discuss many other exploration strategies later in chapter 15.

As we take actions specified by algorithm 9.7, we step into new states sampled from the generative model $T(s, a)$, similar to the sparse sampling method. We increment the visit count $N(s, a)$ and update $Q(s, a)$ to maintain the mean value.

At some point, we will either reach the maximum depth or a state that we have not yet explored. If we reach an unexplored state s , we initialize $N(s, a)$ and $Q(s, a)$ to zero for each action a . We may modify algorithm 9.6 to initialize these counts and value estimates to some other values based on prior expert knowledge of the problem. After initializing N and Q , we then return a value estimate at the state s . It is common to estimate this value through a rollout of some policy using the process outlined in section 9.2.

Examples 9.3 to 9.7 work through an illustration of Monte Carlo tree search applied to the 2048 problem. Figure 9.4 shows a search tree generated by running Monte Carlo tree search on 2048. Example 9.8 discusses the impact of using different strategies for estimating values.

```

function simulate!( $\pi$ ::MonteCarloTreeSearch, s, d= $\pi$ .d)
    if d  $\leq$  0
        return  $\pi$ .U(s)
    end
     $\mathcal{P}$ , N, Q, c =  $\pi$ . $\mathcal{P}$ ,  $\pi$ .N,  $\pi$ .Q,  $\pi$ .c
     $\mathcal{A}$ , TR,  $\gamma$  =  $\mathcal{P}$ . $\mathcal{A}$ ,  $\mathcal{P}$ .TR,  $\mathcal{P}$ . $\gamma$ 
    if !haskey(N, (s, first( $\mathcal{A}$ )))
        for a in  $\mathcal{A}$ 
            N[(s,a)] = 0
            Q[(s,a)] = 0.0
        end
        return  $\pi$ .U(s)
    end
    a = explore( $\pi$ , s)
    s', r = TR(s,a)
    q = r +  $\gamma$ *simulate!( $\pi$ , s', d-1)
    N[(s,a)] += 1
    Q[(s,a)] += (q-Q[(s,a)])/N[(s,a)]
    return q
end

```

Algorithm 9.6. A method for running a Monte Carlo tree search simulation starting from state s to depth d .

```

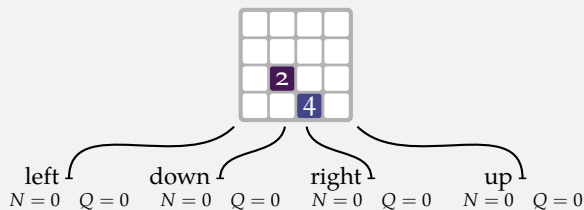
bonus(Nsa, Ns) = Nsa == 0 ? Inf : sqrt(log(Ns)/Nsa)

function explore( $\pi$ ::MonteCarloTreeSearch, s)
     $\mathcal{A}$ , N, Q, c =  $\pi$ . $\mathcal{P}$ . $\mathcal{A}$ ,  $\pi$ .N,  $\pi$ .Q,  $\pi$ .c
    Ns = sum(N[(s,a)] for a in  $\mathcal{A}$ )
    return argmax(a $\rightarrow$ Q[(s,a)] + c*bonus(N[(s,a)], Ns),  $\mathcal{A}$ )
end

```

Algorithm 9.7. An exploration policy used in Monte Carlo tree search when determining which nodes to traverse through the search tree. The policy is determined by a dictionary of state-action visitation counts N and values Q , as well as an exploration parameter c . When $N[(s,a)] = 0$, the policy returns infinity.

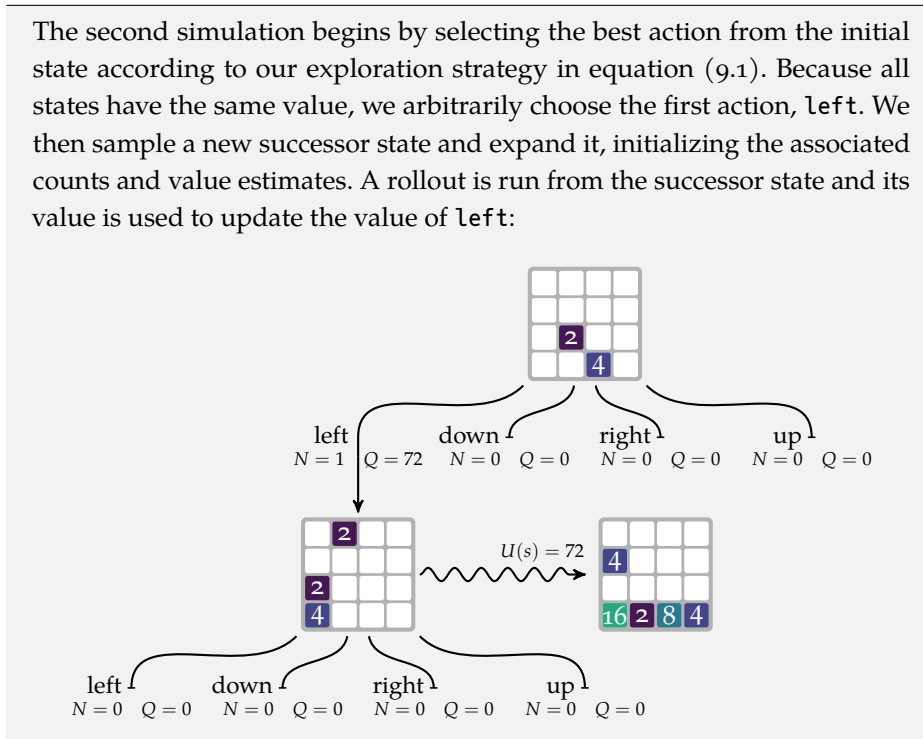
Consider using Monte Carlo tree search to play 2048 (appendix F.2) with a maximum depth $d = 10$, an exploration parameter $c = 100$, and a 10-step random rollout to estimate $U(s)$. Our first simulation expands the starting state. The count and value are initialized for each action from the initial state:



Example 9.3. An example of solving 2048 with Monte Carlo tree search.

The second simulation begins by selecting the best action from the initial state according to our exploration strategy in equation (9.1). Because all states have the same value, we arbitrarily choose the first action, **left**. We then sample a new successor state and expand it, initializing the associated counts and value estimates. A rollout is run from the successor state and its value is used to update the value of **left**:

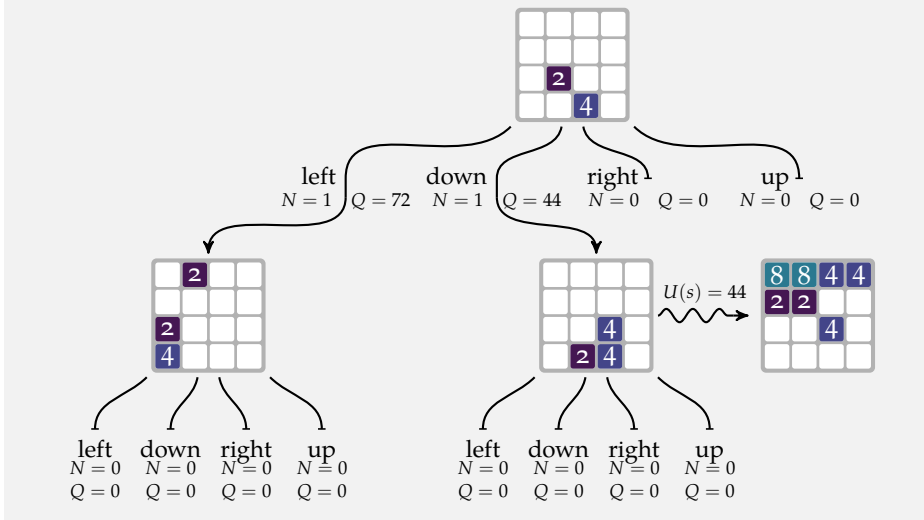
Example 9.4. A (continued) example of solving 2048 with Monte Carlo tree search.



The third simulation begins by selecting the second action, **down**, because it has infinite value due to the exploration bonus given for unexplored actions. The first action has finite value:

$$Q(s_0, \text{left}) + c \sqrt{\frac{\log N(s_0)}{N(s_0, \text{left})}} = 72 + 100 \sqrt{\frac{\log 1}{1}} = 72$$

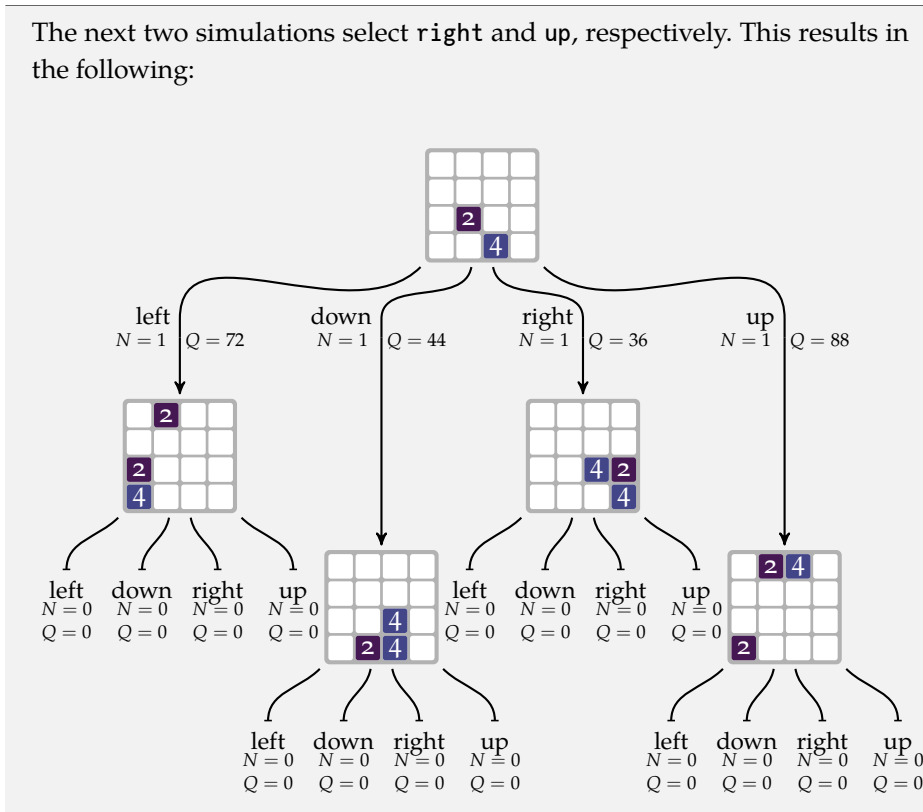
We take the **down** action and sample a new successor state, which is expanded. A rollout is run from the successor state and its value is used to update the value of **down**:



Example 9.5. A (continued) example of solving 2048 with Monte Carlo tree search.

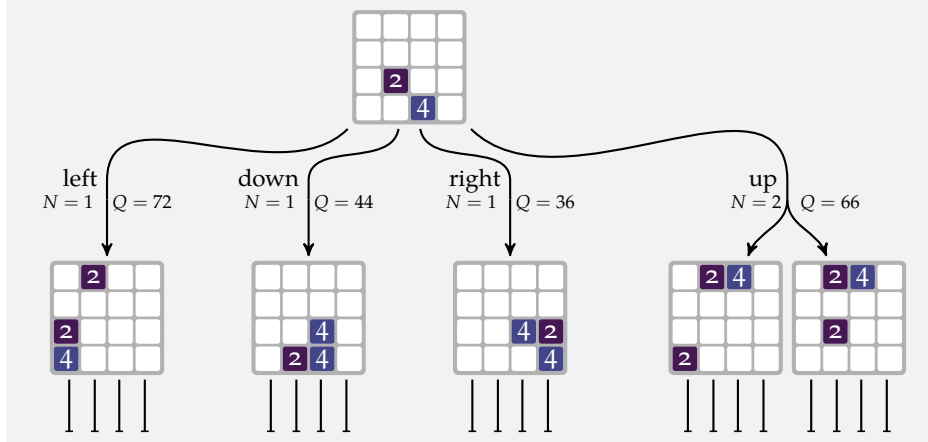
The next two simulations select **right** and **up**, respectively. This results in the following:

Example 9.6. A (continued) example of solving 2048 with Monte Carlo tree search.



In the fifth simulation, up has the highest value. The successor state after taking up in the source state will not necessarily be the same as the first time up was selected. We evaluate $U(s) = 44$ and update our visitation count to 2 and our estimated value to $Q \leftarrow 88 + (44 - 88)/2 = 66$. A new successor node is created:

Example 9.7. A (continued) example of solving 2048 with Monte Carlo tree search.



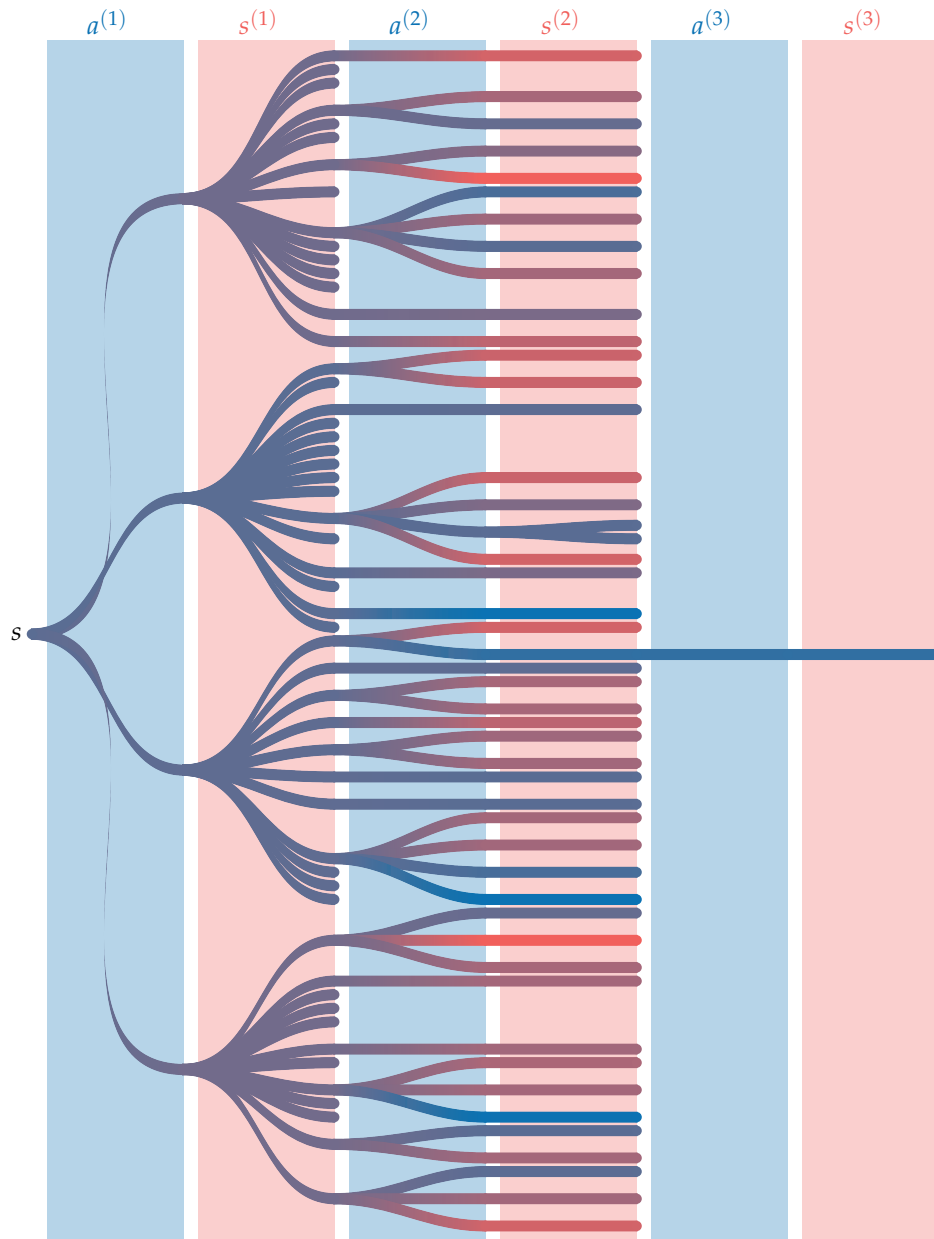


Figure 9.4. A Monte Carlo tree search tree on 2048 after 100 simulations. In general, Monte Carlo tree search for MDPs produces a search graph because there can be multiple ways to reach the same state. The colors in the tree indicate the estimated values at the nodes, with high values in blue and low values in red. The tree is shallow, with a fairly high branching factor, because 2048 has many reachable states for each action.

Rollouts are not the only means by which we can estimate utilities in Monte Carlo tree search. Custom evaluation functions can often be crafted for specific problems to help guide the algorithm. For example, we can encourage Monte Carlo tree search to order its tiles in 2048 using evaluation functions that return the weighted sum across tile values:

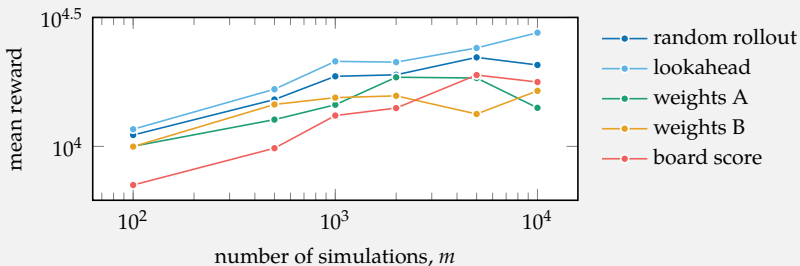
1	2	3	4
2	3	4	5
3	4	5	6
4	5	6	7

heuristic A weights

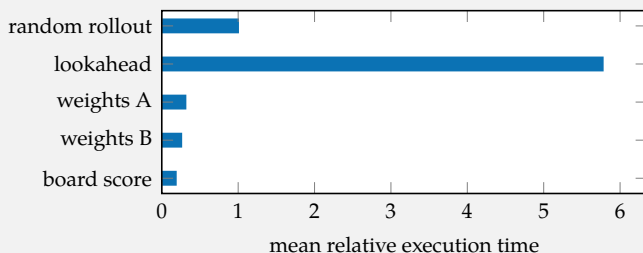
0	1	2	3
7	6	5	4
8	9	10	11
15	14	13	12

heuristic B weights

The plot here compares Monte Carlo tree search on 2048 using rollouts with a uniform random policy, rollouts with a one-step lookahead policy, the two evaluation functions, and using the current board score:



Rollouts perform well but require more execution time. Here we plot the average execution time relative to random rollouts for $m = 100$ from a starting state:



Example 9.8. The performance of Monte Carlo tree search varies with the number of simulations and as the board evaluation method is changed. Heuristic board evaluations tend to be efficient and can more effectively guide the search when run counts are low. Lookahead rollout evaluations take about 18 times longer than heuristic evaluations.

There are variations of this basic Monte Carlo tree search algorithm that can better handle large action and state spaces. Instead of expanding all the actions, we can use *progressive widening*. The number of actions considered from state s is limited to $\theta_1 N(s)^{\theta_2}$, where θ_1 and θ_2 are hyperparameters. Similarly, we can limit the number of states that result from taking action a from state s in the same way, using what is called *double progressive widening*. If the number of states that have been simulated from state s after action a is below $\theta_3 N(s, a)^{\theta_4}$, then we sample a new state; otherwise, we sample one of the previously sampled states with probability proportional to the number of times it has been visited. This strategy can be used to handle large as well as continuous action and state spaces.⁶

9.7 Heuristic Search

Heuristic search (algorithm 9.8) uses m simulations of a greedy policy with respect to a value function U from the current state s .⁷ The value function U is initialized to an upper bound of the value function \bar{U} , which is referred to as a *heuristic*. As we run these simulations, we update our estimate of U through lookahead. After running these simulations, we simply select the greedy action from s with respect to U . Figure 9.5 shows how U and the greedy policy changes with the number of simulations.

Heuristic search is guaranteed to converge to the optimal utility function so long as the heuristic \bar{U} is indeed an upper bound on the value function.⁸ The efficiency of the search depends on the tightness of the upper bound. Unfortunately, tight bounds can be difficult to obtain in practice. While a heuristic that is not a true upper bound may not converge to the optimal policy, it may still converge to a policy that performs well. The time complexity is $O(m \times d \times |S| \times |A|)$.

9.8 Labeled Heuristic Search

Labeled heuristic search (algorithm 9.9) is a variation of heuristic search that runs simulations with value updates while labeling states based on whether their value is solved.⁹ We say that a state s is solved if its utility residual falls below a threshold $\delta > 0$:

$$|U_{k+1}(s) - U_k(s)| < \delta \quad (9.2)$$

⁶ A. Couëtoux, J.-B. Hoock, N. Sokolovska, O. Teytaud, and N. Bonnard, “Continuous Upper Confidence Trees,” in *Learning and Intelligent Optimization (LION)*, 2011.

⁷ A. G. Barto, S. J. Bradtke, and S. P. Singh, “Learning to Act Using Real-Time Dynamic Programming,” *Artificial Intelligence*, vol. 72, no. 1–2, pp. 81–138, 1995. Other forms of heuristic search are discussed by Mausam and A. Kolobov, *Planning with Markov Decision Processes: An AI Perspective*. Morgan & Claypool, 2012.

⁸ Such a heuristic is referred to as an *admissible heuristic*.

⁹ B. Bonet and H. Geffner, “Labeled RTDP: Improving the Convergence of Real-Time Dynamic Programming,” in *International Conference on Automated Planning and Scheduling (ICAPS)*, 2003.

```

struct HeuristicSearch
   $\mathcal{P}$  # problem
  Uhi # upper bound on value function
  d # depth
  m # number of simulations
end

function simulate!( $\pi$ ::HeuristicSearch, U, s)
   $\mathcal{P} = \pi.\mathcal{P}$ 
  for d in 1: $\pi.d$ 
    a, u = greedy( $\mathcal{P}$ , U, s)
    U[s] = u
    s = rand( $\mathcal{P}.T(s, a)$ )
  end
end

function ( $\pi$ ::HeuristicSearch)(s)
  U = [ $\pi.Uhi(s)$  for s in  $\pi.\mathcal{P}.S$ ]
  for i in 1: $\pi.m$ 
    simulate!( $\pi$ , U, s)
  end
  return greedy( $\pi.\mathcal{P}$ , U, s).a
end

```

Algorithm 9.8. Heuristic search runs m simulations starting from an initial state s to a depth d . The search is guided by a heuristic initial value function Uhi , which leads to optimality in the limit of simulations if it is an upper bound on the optimal value function.

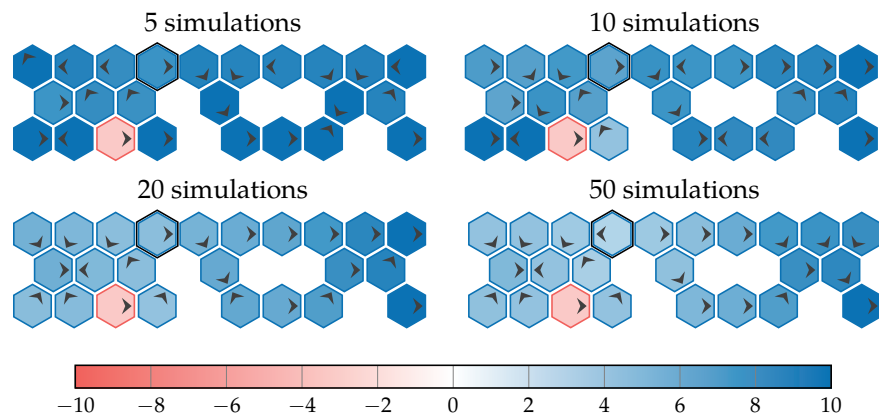


Figure 9.5. Heuristic search runs simulations with Bellman updates to improve a value function on the hex world problem to obtain a policy from an initial state, shown here with an additional black border. These simulations are run to depth 8 with heuristic $\bar{U}(s) = 10$. Each hex is colored according to the utility function value in that iteration. We see that the algorithm eventually finds an optimal policy.

We run simulations with value updates until the current state is solved. In contrast with the heuristic search in the previous section, which runs a fixed number of iterations, this labeling process focuses computational effort on the most important areas of the state space.

```

struct LabeledHeuristicSearch
     $\mathcal{P}$  # problem
    Uhi # upper bound on value function
    d # depth
     $\delta$  # gap threshold
end

function ( $\pi$ ::LabeledHeuristicSearch)(s)
    U, solved = [ $\pi$ .Uhi(s) for s in  $\mathcal{P}$ .S], Set()
    while s  $\notin$  solved
        simulate!( $\pi$ , U, solved, s)
    end
    return greedy( $\pi$ . $\mathcal{P}$ , U, s).a
end

```

Algorithm 9.9. Labeled heuristic search, which runs simulations starting from the current state to depth d until the current state is solved. The search is guided by a heuristic upper bound on the value function U_{hi} and maintains a growing set of solved states. States are considered solved when their utility residuals fall below δ . A value function policy is returned.

Simulations in labeled heuristic search (algorithm 9.10) begin by running to a maximum depth of d by following a policy that is greedy with respect to our estimated value function U , similar to the heuristic search in the previous section. We may stop a simulation before a depth of d if we reach a state that has been labeled as solved in a prior simulation.

```

function simulate!( $\pi$ ::LabeledHeuristicSearch, U, solved, s)
    visited = []
    for d in 1: $\pi$ .d
        if s  $\in$  solved
            break
        end
        push!(visited, s)
        a, u = greedy( $\pi$ . $\mathcal{P}$ , U, s)
        U[s] = u
        s = rand( $\pi$ . $\mathcal{P}$ .T(s, a))
    end
    while !isempty(visited)
        if label!( $\pi$ , U, solved, pop!(visited))
            break
        end
    end
end

```

Algorithm 9.10. Simulations are run from the current state to a maximum depth d . We stop a simulation at depth d or if we encounter a state that is in the set `solved`. After a simulation, we call `label!` on the states we visited in reverse order.

After each simulation, we iterate over the states we visited during that simulation in reverse order, performing a labeling routine on each state and stopping if a state is found that is not solved. The labeling routine (algorithm 9.11) searches the states in the *greedy envelope* of s , which is defined to be the states reachable from s under a greedy policy with respect to U . The state s is considered not solved if there is a state in the greedy envelope of s whose utility residual is greater than threshold δ . If no such state is found, then s is marked as solved—as well as all states in the greedy envelope of s because they must have converged as well. If a state with a sufficiently large utility residual is found, then the utilities of all states traversed during the search of the greedy envelope are updated.

Figure 9.6 shows several different greedy envelopes. Figure 9.7 shows the states traversed in a single iteration of labeled heuristic search. Figure 9.8 shows the progression of heuristic search on the hex world problem.

9.9 Open-Loop Planning

The online methods discussed in this chapter, as well as the offline methods discussed in the previous chapters, are examples of *closed-loop planning*, which involves accounting for future state information in the planning process.¹⁰ Often, *open-loop planning* can provide a satisfactory approximation of an optimal closed-loop plan while greatly enhancing computational efficiency by avoiding having to reason about the acquisition of future information. Sometimes this open-loop planning approach is referred to as *model predictive control*.¹¹ As with receding horizon control, model predictive control solves the open-loop problem, executes the action from our current state, transitions to the next state, and then replans.

Open-loop plans can be represented as a sequence of actions up to a depth d . The planning process reduces to an optimization problem:

$$\underset{a_{1:d}}{\text{maximize}} \quad U(a_{1:d}) \quad (9.3)$$

where $U(a_{1:d})$ is the expected return when executing the sequence of actions $a_{1:d}$. Depending on the application, this optimization problem may be convex or lend itself to a convex approximation, meaning that it can be solved quickly using a variety of algorithms.¹² Later in this section, we will discuss a few different formulations that can be used to transform equation (9.3) into a convex problem.

¹⁰ The loop in this context refers to the observe-act loop introduced in section 1.1.

¹¹ F. Borrelli, A. Bemporad, and M. Morari, *Predictive Control for Linear and Hybrid Systems*. Cambridge University Press, 2019.

¹² Appendix A.6 reviews convexity. An introduction to convex optimization is provided by S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge University Press, 2004.

```

function expand( $\pi$ ::LabeledHeuristicSearch, U, solved, s)
     $\mathcal{P}, \delta = \pi.\mathcal{P}, \pi.\delta$ 
     $\mathcal{S}, \mathcal{A}, \mathcal{T} = \mathcal{P}.\mathcal{S}, \mathcal{P}.\mathcal{A}, \mathcal{P}.\mathcal{T}$ 
    found, toexpand, envelope = false, Set(s), []
    while !isempty(toexpand)
        s = pop!(toexpand)
        push!(envelope, s)
        a, u = greedy( $\mathcal{P}$ , U, s)
        if abs(U[s] - u) >  $\delta$ 
            found = true
        else
            for s' in  $\mathcal{S}$ 
                if  $\mathcal{T}(s, a, s') > 0$  && s'  $\notin$  (solved  $\cup$  envelope)
                    push!(toexpand, s')
                end
            end
        end
    end
    return (found, envelope)
end

function label!( $\pi$ ::LabeledHeuristicSearch, U, solved, s)
    if s  $\in$  solved
        return false
    end
    found, envelope = expand( $\pi$ , U, solved, s)
    if found
        for s  $\in$  reverse(envelope)
            U[s] = greedy( $\pi.\mathcal{P}$ , U, s).u
        end
    else
        union!(solved, envelope)
    end
    return found
end

```

Algorithm 9.11. The `label!` function will attempt to find a state in the greedy envelope of `s` whose utility residual exceeds a threshold δ . The function `expand` computes the greedy envelope of `s` and determines whether any of those states have utility residuals above the threshold. If a state has a residual that exceeds the threshold, then we update the utilities of the states in the envelope. Otherwise, we add that envelope to the set of solved states.

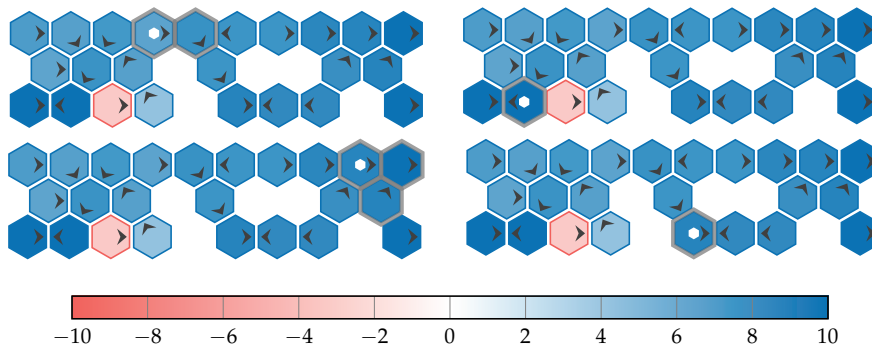


Figure 9.6. The greedy envelope for $\delta = 1$ for several states visualized for a value function on the hex world problem. The value function was obtained by running basic heuristic search for 10 iterations from an initial state, shown with a white hex center, to a maximum depth of 8. We find that the size of the greedy envelope, outlined in gray, can vary widely depending on the state.



Figure 9.7. A single iteration of labeled heuristic search conducts an exploratory run (arrows), followed by labeling (hexagonal border). Only two states are labeled in this iteration: the hidden terminal state and the state with a hexagonal border. Both the exploratory run and the labeling step update the value function.

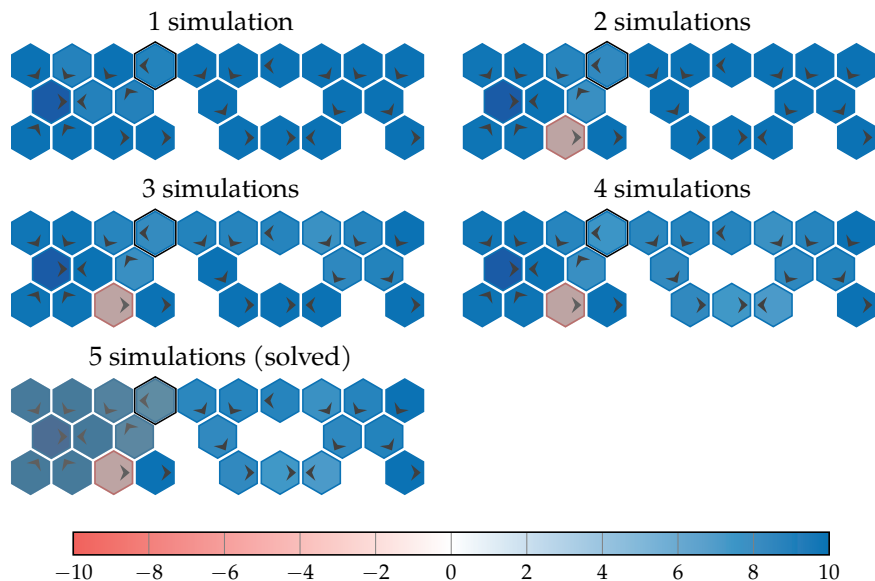


Figure 9.8. A progression of heuristic search on the hex world problem using $\delta = 1$ and a heuristic $\bar{U}(s) = 10$. The solved states in each iteration are covered in a gray wash. The set of solved states grows from the terminal reward state back toward the initial state with the dark border.

Open-loop planning can often allow us to devise effective decision strategies in high-dimensional spaces where closed-loop planning is computationally infeasible. This type of planning gains this efficiency by not accounting for future information. Example 9.9 provides a simple instance of where open-loop planning can result in poor decisions, even when we account for stochasticity.

Consider a problem with nine states, as shown in the margin, with two decision steps starting from the initial state s_1 . In our decisions, we must decide between going up (blue arrows) and going down (green arrows). The effects of these actions are deterministic, except that if we go up from s_1 , then we end up in state s_2 half the time and in state s_3 half the time. We receive a reward of 30 in states s_5 and s_7 and a reward of 20 in states s_8 and s_9 , as indicated in the illustration.

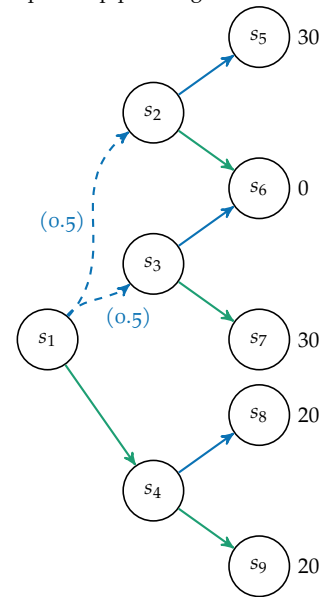
There are exactly four open-loop plans: (up, up), (up, down), (down, up), and (down, down). In this simple example, it is easy to compute their expected utilities:

- $U(\text{up, up}) = 0.5 \times 30 + 0.5 \times 0 = 15$
- $U(\text{up, down}) = 0.5 \times 0 + 0.5 \times 30 = 15$
- $U(\text{down, up}) = 20$
- $U(\text{down, down}) = 20$

According to the set of open-loop plans, it is best to choose down from s_1 because our expected reward is 20 instead of 15.

Closed-loop planning, in contrast, takes into account the fact that we can base our next decision on the observed outcome of our first action. If we choose to go up from s_1 , then we can choose to go down or up depending on whether we end up in s_2 or s_3 , thereby guaranteeing a reward of 30.

Example 9.9. Suboptimality of open-loop planning.



9.9.1 Deterministic Model Predictive Control

A common approximation to make $U(a_{1:d})$ amenable to optimization is to assume deterministic dynamics:

$$\begin{aligned} & \underset{a_{1:d}, s_{2:d}}{\text{maximize}} && \sum_{t=1}^d \gamma^t R(s_t, a_t) \\ & \text{subject to} && s_{t+1} = T(s_t, a_t), \quad t \in 1:d-1 \end{aligned} \quad (9.4)$$

where s_1 is the current state and $T(s, a)$ is a deterministic transition function that returns the state that results from taking action a from state s . A common strategy for producing a suitable deterministic transition function from a stochastic transition function is to use the most likely transition. If the dynamics in equation (9.4) are linear and the reward function is convex, then the problem is convex.

Example 9.10 provides an instance involving navigating to a goal state while avoiding an obstacle and minimizing acceleration effort. Both the state space and action space are continuous, and we can find a solution in well under a second. Replanning after every step can help compensate for stochasticity or unexpected events. For example, if the obstacle moves, we can readjust our plan, as illustrated in figure 9.9.

9.9.2 Robust Model Predictive Control

We can change the problem formulation to provide robustness to outcome uncertainty. There are many *robust model predictive control* formulations,¹³ but one involves choosing the best open-loop plan given the worst-case state transitions. This formulation defines $T(s, a)$ to be an *uncertainty set* consisting of all possible states that can result from taking action a in state s . In other words, the uncertainty set is the support of the distribution $T(\cdot | s, a)$. Optimizing with respect to worst-case state transitions requires transforming the optimization problem in equation (9.4) into a *minimax* problem:

$$\begin{aligned} & \underset{a_{1:d}}{\text{maximize}} & \underset{s_{2:d}}{\text{minimize}} & \sum_{t=1}^d \gamma^t R(s_t, a_t) \\ & \text{subject to} && s_{t+1} \in T(s_t, a_t), \quad t \in 1:d-1 \end{aligned} \quad (9.5)$$

¹³ A. Bemporad and M. Morari, "Robust Model Predictive Control: A Survey," in *Robustness in Identification and Control*, A. Garulli, A. Tesi, and A. Vicino, eds., Springer, 1999, pp. 207–226.

In this problem, our state \mathbf{s} represents our agent's two-dimensional position concatenated with its two-dimensional velocity vector, with \mathbf{s} initially set to $[0, 0, 0, 0]$. Our action \mathbf{a} is an acceleration vector, where each component must be between ± 1 . At each step, we use our action to update our velocity, and we use our velocity to update our position. Our objective is to reach a goal state of $\mathbf{s}_{\text{goal}} = [10, 10, 0, 0]$. We plan up to $d = 10$ steps with no discounting. With each step, we accumulate a cost of $\|\mathbf{a}_t\|_2^2$ to minimize acceleration effort. At the last step, we want to be as close to the goal state as possible, with a penalty of $100\|\mathbf{s}_d - \mathbf{s}_{\text{goal}}\|_2^2$. We also have to ensure that we avoid a circular obstacle with radius 2 centered at $[3, 4]$. We can formulate this problem as follows and extract the first action from the plan:

```

model = Model(Ipopt.Optimizer)
d = 10
current_state = zeros(4)
goal = [10,10,0,0]
obstacle = [3,4]
@variables model begin
    s[1:4, 1:d]
    -1 ≤ a[1:2,1:d] ≤ 1
end
# velocity update
@constraint(model, [i=2:d,j=1:2], s[2+j,i] == s[2+j,i-1] + a[j,i-1])
# position update
@constraint(model, [i=2:d,j=1:2], s[j,i] == s[j,i-1] + s[2+j,i-1])
# initial condition
@constraint(model, s[:,1] .== current_state)
# obstacle
@constraint(model, [i=1:d], sum((s[1:2,i] - obstacle).^2) ≥ 4)
@objective(model, Min, 100*sum((s[:,d] - goal).^2) + sum(a.^2))
optimize!(model)
action = value.(a[:,1])

```

Example 9.10. Open-loop planning in a deterministic environment. We attempt to find a path around a circular obstacle. This implementation uses the JuMP.jl interface to the Ipopt solver. A. Wächter and L. T. Biegler, "On the Implementation of an Interior-Point Filter Line-Search Algorithm for Large-Scale Nonlinear Programming," *Mathematical Programming*, vol. 106, no. 1, pp. 25–57, 2005.

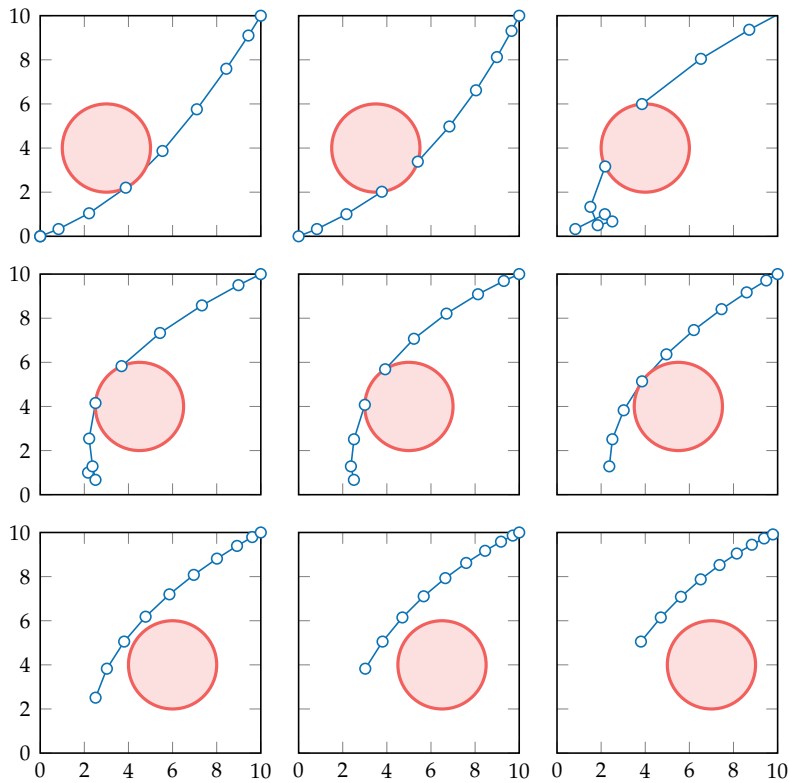


Figure 9.9. Model predictive control applied to the problem in example 9.10, with the addition of a moving obstacle. The sequence progresses left-to-right, and top-to-bottom. Initially, we have a plan that passes to the right of the obstacle, but in the third cell, we see that we must change our mind and pass to the left. We have to maneuver around a little to adjust our velocity vector appropriately with minimal effort. Of course, we could have created a better path (in terms of our utility function) if our planning process had known that the obstacle was moving in a particular direction.

Unfortunately, this formulation can result in extremely conservative behavior. If we adapt example 9.10 to model the uncertainty in the motion of the obstacle, the accumulation of uncertainty may become quite large, even when planning with a relatively short horizon. One way to reduce the accumulation of uncertainty is to restrict the uncertainty set output by $T(s, a)$ to contain only, say, 95% of the probability mass. Another issue with this approach is that the minimax optimization problem is often not convex and difficult to solve.

9.9.3 Multiforecast Model Predictive Control

One way to address the computational challenge within the minimax problem in equation (9.5) is to use m forecast scenarios, each of which follows its own deterministic transition function.¹⁴ There are various formulations of this kind of *multiforecast model predictive control*, which is a type of *hindsight optimization*.¹⁵ One common approach is to have the deterministic transition functions depend on the step k , $T_i(s, a, k)$, which is the same as augmenting the state space to include depth. Example 9.11 demonstrates how this might be done for a linear Gaussian model.

¹⁴S. Garatti and M.C. Campi, “Modulating Robustness in Control Design: Principles and Algorithms,” *IEEE Control Systems Magazine*, vol. 33, no. 2, pp. 36–51, 2013.

¹⁵It is called hindsight optimization because it represents a solution that is optimizing using knowledge about action outcomes that can only be known in hindsight. E.K.P. Chong, R.L. Givan, and H.S. Chang, “A Framework for Simulation-Based Network Control via Hindsight Optimization,” in *IEEE Conference on Decision and Control (CDC)*, 2000.

Suppose we have a problem with linear Gaussian dynamics:

$$T(\mathbf{s}' \mid \mathbf{s}, \mathbf{a}) = \mathcal{N}(\mathbf{T}_s \mathbf{s} + \mathbf{T}_a \mathbf{a}, \Sigma)$$

The problem in figure 9.9 is linear, with no uncertainty, but if we allow the obstacle to move according to a Gaussian distribution at each step, then the dynamics become linear Gaussian. We can approximate the dynamics using a set of m forecast scenarios, each consisting of d steps. We can pull $m \times d$ samples $\epsilon_{ik} \sim \mathcal{N}(\mathbf{0}, \Sigma)$ and define the deterministic transition functions:

$$T_i(\mathbf{s}, \mathbf{a}, k) = \mathbf{T}_s \mathbf{s} + \mathbf{T}_a \mathbf{a} + \epsilon_{ik}$$

Example 9.11. Modeling linear Gaussian transition dynamics in multiforecast model predictive control.

We try to find the best sequence of actions for the worst sampled scenario:

$$\begin{aligned} & \underset{a_{1:d}}{\text{maximize}} && \underset{i, s_{2:d}}{\text{minimize}} && \sum_{k=1}^d \gamma^k R(s_k, a_k) && (9.6) \\ & \text{subject to} && && s_{k+1} = T_i(s_k, a_k, k), && k \in 1:d-1 \end{aligned}$$

This problem can be much easier to solve than the original robust problem.

We can also use a multiforecast approach to optimize the average case.¹⁶ The formulation is similar to equation (9.6), except that we replace the minimization with an expectation and allow different action sequences to be taken for different scenarios, with the constraint that the first action must agree:

$$\begin{aligned}
 & \underset{a_{1:d}^{(1:m)}, s_{2:d}^{(i)}}{\text{maximize}} && \frac{1}{m} \sum_{i=1}^m \sum_{k=1}^d \gamma^k R(s_k^{(i)}, a_k^{(i)}) \\
 & \text{subject to} && s_{k+1}^{(i)} = T_i(s_k^{(i)}, a_k^{(i)}, k), \quad k \in 1:d-1, i \in 1:m \\
 & && a_1^{(i)} = a_1^{(j)}, \quad i \in 1:m, j \in 1:m
 \end{aligned} \tag{9.7}$$

This formulation can result in robust behavior without being overly conservative, while still maintaining computational tractability. Both formulations in equations (9.6) and (9.7) can be made more robust by increasing the number of forecast scenarios m at the expense of additional computation.

9.10 Summary

- Online methods plan from the current state, focusing computation on states that are reachable.
- Receding horizon planning involves planning to a certain horizon and then replanning with each step.
- Lookahead with rollouts involves acting greedily with respect to values estimated using simulations of a rollout policy; it is computationally efficient compared to other algorithms, but there are no guarantees on performance.
- Forward search considers all state-action transitions up to a certain depth, resulting in computational complexity that grows exponentially in both the number of states and the number of actions.
- Branch and bound uses upper and lower bound functions to prune portions of the search tree that will not lead to a better outcome in expectation.
- Sparse sampling avoids the exponential complexity in the number of states by limiting the number of sampled transitions from every search node.

¹⁶ This approach was applied to optimizing power flow policies by N. Moehle, E. Busseti, S. Boyd, and M. Wytoczek, “Dynamic Energy Management,” in *Large Scale Optimization in Supply Chains and Smart Manufacturing*, Springer, 2019, pp. 69–126.

- Monte Carlo tree search guides search to promising areas of the search space by taking actions that balance exploration with exploitation.
- Heuristic search runs simulations of a policy that is greedy with respect to a value function that is updated along the way using lookahead.
- Labeled heuristic search reduces computation by not reevaluating states whose values have converged.
- Open-loop planning aims to find the best possible sequence of actions and can be computationally efficient if the optimization problem is convex.

9.11 Exercises

Exercise 9.1. Why does branch and bound have the same worst-case computational complexity as forward search?

Solution: In the worst case, branch and bound will never prune, resulting in a traversal of the same search tree as forward search with the same complexity.

Exercise 9.2. Given two admissible heuristics h_1 and h_2 , how can we use both of them in heuristic search?

Solution: Create a new heuristic $h(s) = \min\{h_1(s), h_2(s)\}$ and use it instead. This new heuristic is guaranteed to be admissible and cannot be a worse bound than either h_1 or h_2 . Both $h_1(s) \geq U^*(s)$ and $h_2(s) \geq U^*(s)$ imply that $h(s) \geq U^*(s)$.

Exercise 9.3. Given two inadmissible heuristics h_1 and h_2 , describe a way we can use both of them in heuristic search.

Solution: We could define a new heuristic $h_3(s) = \max(h_1(s), h_2(s))$ to get a potentially admissible, or “less-inadmissible,” heuristic. It may be slower to converge, but it may be more likely to not miss out on a better solution.

Exercise 9.4. Suppose we have a discrete MDP with state space \mathcal{S} and action space \mathcal{A} and we want to perform forward search to depth d . Due to computational constraints and the requirement that we must simulate to depth d , we decide to generate new, smaller state and action spaces by re-discretizing the original state and action spaces on a coarser scale with $|\mathcal{S}'| < |\mathcal{S}|$ and $|\mathcal{A}'| < |\mathcal{A}|$. In terms of the original state and action spaces, what would the size of the new state and action spaces need to be in order to make the computational complexity of forward search approximately depth-invariant with respect to the size of our original state and action spaces, that is, $O(|\mathcal{S}'||\mathcal{A}'|)$?

Solution: We need

$$|\mathcal{S}'| = |\mathcal{S}|^{\frac{1}{d}} \quad \text{and} \quad |\mathcal{A}'| = |\mathcal{A}|^{\frac{1}{d}}$$

This results in the following complexity:

$$O(|\mathcal{S}'|^d |\mathcal{A}'|^d) = O\left(\left(|\mathcal{S}|^{\frac{1}{d}}\right)^d \left(|\mathcal{A}|^{\frac{1}{d}}\right)^d\right) = O(|\mathcal{S}||\mathcal{A}|)$$

Exercise 9.5. Building on the previous exercise, suppose now that we want to keep all the original actions in our action space and only re-discretize the state space. What would the size of the new state space need to be to make the computational complexity of forward search approximately depth-invariant with respect to the size of our original state and action spaces?

Solution: The computational complexity of forward search is given by $O(|\mathcal{S}||\mathcal{A}|^d)$, which can also be written as $O(|\mathcal{S}|^d |\mathcal{A}|^d)$. Thus, in order for our coarser state space to lead to forward search that is approximately depth-invariant with respect to the size of our original state and action spaces, we need

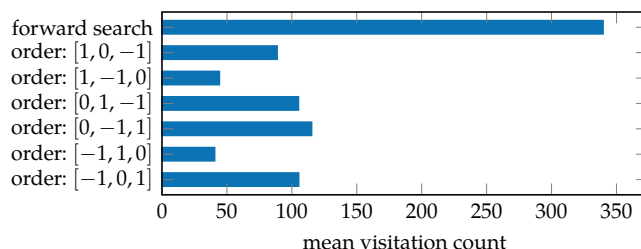
$$|\mathcal{S}'| = \left(\frac{|\mathcal{S}|}{|\mathcal{A}|^{d-1}}\right)^{\frac{1}{d}}$$

This gives us:

$$O(|\mathcal{S}'|^d |\mathcal{A}'|^d) = O\left(\left[\left(\frac{|\mathcal{S}|}{|\mathcal{A}|^{d-1}}\right)^{\frac{1}{d}}\right]^d |\mathcal{A}|^d\right) = O\left(|\mathcal{S}| \frac{|\mathcal{A}|^d}{|\mathcal{A}|^{d-1}}\right) = O(|\mathcal{S}||\mathcal{A}|)$$

Exercise 9.6. Will changing the ordering of the action space cause forward search to take different actions? Will changing the ordering of the action space cause branch and bound to take different actions? Can the ordering of the action space affect how many states are visited by branch and bound?

Solution: Forward search enumerates over all possible future actions. It may return different actions if there are ties in their expected utilities. Branch and bound maintains the same optimality guarantee over the horizon as forward search by sorting by upper bound. The ordering of the action space can affect branch and bound's visitation rate when the upper bound produces the same expected value for two or more actions. Below we show this effect on the modified mountain car problem from example 9.2. The plot compares the number of states visited in forward search to that of branch and bound for different action orderings to depth 6. Branch and bound consistently visits far fewer states than forward search, but action ordering can still affect state visitation.



Exercise 9.7. Is sparse sampling with $m = |\mathcal{S}|$ equivalent to forward search?

Solution: No. While the computational complexities are identical at $O(|\mathcal{S}|^d |\mathcal{A}|^d)$, forward search will branch on all states in the state space, while sparse sampling will branch on $|\mathcal{S}|$ randomly sampled states.

Exercise 9.8. Given an MDP with $|\mathcal{S}| = 10$, $|\mathcal{A}| = 3$, and a uniform transition distribution $T(s' | s, a) = 1/|\mathcal{S}|$ for all s and a , what is the probability that sparse sampling with $m = |\mathcal{S}|$ samples and depth $d = 1$ yields the exact same search tree produced by forward search with depth $d = 1$?

Solution: For both forward search and sparse sampling, we branch on all actions from the current state node. For forward search, at each of these action nodes, we branch on all states, while for sparse sampling, we will branch on $m = |\mathcal{S}|$ sampled states. If these sampled states are exactly equal to the state space, that action branch is equivalent to the branch produced in forward search. Thus, for a single action branch we have:

the probability the first state is unique	$\frac{10}{10}$
the probability the second state is unique (not equal to the first state)	$\frac{9}{10}$
the probability the third state is unique (not equal to the first or second state)	$\frac{8}{10}$
\vdots	\vdots

Since each of these sampled states is independent, this leads to the probability of all unique states in the state space being selected with probability

$$\frac{10 \times 9 \times 8 \times \dots}{10 \times 10 \times 10 \times \dots} = \frac{10!}{10^{10}} \approx 0.000363$$

Since each of the sampled states across different action branches is independent, the probability that all three action branches sample the unique states in the state space is

$$\left(\frac{10!}{10^{10}}\right)^3 \approx (0.000363)^3 \approx 4.78 \times 10^{-11}$$

Exercise 9.9. Given the following tables of $Q(s, a)$ and $N(s, a)$, use the upper confidence bound in equation (9.1) to compute the MCTS traversal action for each state with an exploration parameter of $c_1 = 10$ and again for $c_2 = 20$.

	$Q(s, a_1)$	$Q(s, a_2)$		$N(s, a_1)$	$N(s, a_2)$
s_1	10	-5	s_1	27	4
s_2	12	10	s_2	32	18

Solution: For the first exploration parameter $c_1 = 10$, we tabulate the upper confidence bound of each state-action pair and select the action maximizing the bound for each state:

	$UCB(s, a_1)$	$UCB(s, a_2)$	$\arg \max_a UCB(s, a)$
s_1	$10 + 10\sqrt{\frac{\log 31}{27}} \approx 13.566$	$-5 + 10\sqrt{\frac{\log 31}{4}} \approx 4.266$	a_1
s_2	$12 + 10\sqrt{\frac{\log 50}{32}} \approx 15.496$	$10 + 10\sqrt{\frac{\log 50}{18}} \approx 14.662$	a_1

And for $c_2 = 20$, we have:

	$UCB(s, a_1)$	$UCB(s, a_2)$	$\arg \max_a UCB(s, a)$
s_1	$10 + 20\sqrt{\frac{\log 31}{27}} \approx 17.133$	$-5 + 20\sqrt{\frac{\log 31}{4}} \approx 13.531$	a_1
s_2	$12 + 20\sqrt{\frac{\log 50}{32}} \approx 18.993$	$10 + 20\sqrt{\frac{\log 50}{18}} \approx 19.324$	a_2