

# Markov Decision Processes & Reinforcement Learning

## 1. Value Iteration

Say an MDP has state space  $S$  and reward  $R$  where all rewards are positive. If you run value iteration, what is the largest  $k$  for which  $V_k(s)$  is zero?

$k = 0$ .  $V_0(s) = 0$  for all states  $s$ . Since  $V_1(s) = \max\{r_1, r_2, \dots\}$  and all rewards are positive, all  $V_k(s) > 0$  when  $k > 0$ .

## 2. Infinite Horizon

You are on a linear space and can move only right or left. Each position has reward  $r_i$  and  $\gamma \approx 1$ . Describe your optimal policy given any state  $i$  (don't forget about ties).

You will always move towards state  $j$  where  $j \in \arg \max_i \{r_i\}$  and stay there forever. When there are ties, you move to the closest state  $j$  and then stay there. "Closest" is calculated as the state  $j$  which will get you  $\max \sum_{k=i}^j r_k$  (i.e., you get the maximum reward by going to that state).

### 3. Value Iteration vs Expectimax Search

What is the running time of Expectimax Search and Value Iteration as a function of the horizon,  $T$ , the number of actions,  $M$ , the the number of transitions for any state and action,  $L$ , and the number of steps,  $N$ ? Why don't we always choose the algorithm with better asymptotic running time?

Expectimax Search takes time proportional to the height of the game tree (see Figure 3 in the course notes on MDP's), so we see that the running time is  $O((ML)^T)$ . Value Iteration uses dynamic programming to ensure that it doesn't revisit the same state multiple times, and we get running time linear in the inputs  $O(NMLT)$ . Although value iteration is a linear-time algorithm, it visits states that might not actually be reachable, while Expectimax only runs computations on states that can be reached. For some problem domains, this results in higher efficiency.

#### 4. Setting Rewards to Compute Shortest Path

Suppose we have a grid-world where each state is represented as a point in  $\{(x, y) | 0 \leq x \leq 4, 0 \leq y \leq 4\}$ . Suppose you have a robot that starts in the lower left corner and is given a goal point that it needs to reach. At each state, the robot can move one point to the left, right, up, or down, and each action has a 90% chance of success (otherwise you stay at the current point).

- What is the size of the state space in the resulting MDP?
- Suppose we want to minimize the length of path to the goal, but we have no preference for which path the robot should take (and all points are reachable). What rewards could we assign to each state so that we recover a shortest path to the goal as our optimal policy?

- Our grid is composed of each point in a  $5 \times 5$  grid, so we have 25 possible states in the state-space.
- We can set the reward of each state, new state pair (that does not end with the goal node) to be the same negative value. Then, when we try to maximize our reward we will choose the shortest path, and since each transition has the same reward, we do not preferentially choose any paths with the same length.

## 5. Q learning

Suppose you are standing on a linear board: you can take action  $L$  or  $R$  (walk left or right). If you walk, you have probability  $p_a$  that you actually walk to the next square, where  $a \in L, R$ . Otherwise, your cat distracted you and you are still on the same square. Staying a square gives you reward  $r_i$ . Your learning rate is  $\alpha = .5$  and  $\gamma = .5$ .

- (a) You are on square 1, you choose  $a = L$  and receive  $r = 4$ . What is your updated value of  $Q(1, L)$ ?
- (b) In the next step, you are on square 0, you choose  $a = R$ , receive  $r = 3$  and end up in  $s = 1$ . What is your updated value of  $Q(0, R)$ ?

All Q values are initially 0.

(a)  $Q(1, L) = 0 + .5(4 + .5 \times 0 - 0) = 2$

(b)  $Q(0, R) = 0 + .5(3 + .5 \times 2 - 0) = 2$

Notice that we do not need to make use of  $p_a$ , the transition probability!

## 6. $\epsilon$ -Greedy

Why do we generally use an  $\epsilon$ -Greedy algorithm when choosing the current action during the Q-Learning algorithm? Describe how you would change the value of  $\epsilon$  over time to encourage early exploration/learning and good decision making after the state space is sufficiently explored.

First we recall that in the  $\epsilon$ -greedy algorithm we choose the best action, via

$$\max_{a' \in A} Q(s, a'),$$

with probability  $1 - \epsilon$  and a random action otherwise. We use this because it encourages exploration of the state space by introducing randomness into the choice of action at each step, which we recall is necessary for RL to work with. In the beginning, we would choose a high value of  $\epsilon$ , which would encourage a lot of random exploration of the state space, and as time went on, we could slowly decrease the value of  $\epsilon$  to encourage taking the action that maximizes the Q value in the current state, which we expect to be a good action given that we've explored the state space sufficiently to closely approximate the true Q function.

## 7. Convergence of Q-Learning

Suppose we have a deterministic world where each state-action pair  $(s, a)$  is visited infinitely often. Consider an interval during which every such pair is visited. Suppose that the largest error in the approximation of  $\hat{Q}_n$  after  $n$  iterations is  $e_n$ , i.e.

$$e_n = \max_{s,a} |\hat{Q}_n(s, a) - Q(s, a)|$$

Show that  $e_{n+1}$  is bounded above by  $\gamma e_n$ , where  $\gamma$  is the usual discount factor.

$$\begin{aligned} e_{n+1} &= \max_{s,a} |\hat{Q}_{n+1}(s, a) - Q(s, a)| \\ &= \max_{s,a} \left| (R(s, a) + \gamma \max_{a'} \hat{Q}_n(s', a')) - (R(s, a) + \gamma \max_{a'} Q(s', a')) \right| \\ &= \left| \gamma \max_{a'} \hat{Q}_n(s', a') - \gamma \max_{a'} Q(s', a') \right| \\ &= \gamma \left| \max_{a'} \hat{Q}_n(s', a') - \max_{a'} Q(s', a') \right| \\ &\leq \gamma \max_{a'} |\hat{Q}_n(s', a') - Q(s', a')| \\ &\leq \gamma \max_{s'', a'} |\hat{Q}_n(s'', a) - Q(s'', a')| \\ &= \gamma e_n \end{aligned}$$