
CS 188 Introduction to Written HW 11 Sol.
Spring 2022 Artificial Intelligence

Solutions for HW 11 (Written)

Q1. [26 pts] MDPs and RL

The agent is in a 2×4 gridworld as shown in the figure. We start from square 1 and finish in square 8. When square 8 is reached, we receive a reward of +10 at the game end. For anything else, we receive a constant reward of -1 (you can think of this as a time penalty).

1	2	3	4
5	6	7	8

The actions in this MDP include: up, down, left and right. The agent cannot take actions that take them off the board. In the table below, we provide initial non-zero estimates of Q values (Q values for invalid actions are left as blanks):

Table 1

	action=up	action=down	action=left	action=right
state=1		Q(1, down)=4		Q(1, right)=3
state=2		Q(2, down)=6	Q(2, left)=4	Q(2, right)=5
state=3		Q(3, down)=8	Q(3, left)=5	Q(3, right)=7
state=4		Q(4, down)=9	Q(4, left)=6	
state=5	Q(5, up)=5			Q(5, right)=6
state=6	Q(6, up)=4		Q(6, left)=5	Q(6, right)=7
state=7	Q(7, up)=6		Q(7, left)=6	Q(7, right)=8

- (a) Your friend Adam guesses that the actions in this MDP are fully deterministic (e.g. taking down from 2 will land you in 6 with probability 1 and everywhere else with probability 0). Since we have full knowledge of T and R , we can thus use the Bellman equation to improve (i.e., further update) the initial Q estimates.

Adam tells you to use the following update rule for Q values, where he assumes that your policy is greedy and thus does $\max_a Q(s, a)$. The update rule he prescribes is as follows:

$$Q_{k+1}(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma \max_{a'} Q_k(s', a')]$$

- (i) [1 pt] Perform one update of $Q(3, \text{left})$ using the equation above, where $\gamma = 0.8$. You may break ties in any way.

$-1 + 0.8 \times (1 \times 6)$ because we are in a deterministic grid world with a greedy policy.

- (ii) [1 pt] Perform one update of $Q(3, \text{down})$ using the equation above, where $\gamma = 0.8$.

$-1 + 0.8 \times (1 \times 8)$.

- (iii) [3 pts] For the Q update rule prescribed above, how is it different from the Q learning update that we saw in lecture, which is $Q_{k+1}(s, a) = (1 - \alpha)Q_k(s, a) + \alpha * \text{sample}$?

Adam's Q update rule is the update rule for Q-value iteration (Bellman update). The difference between Q-value iteration and Q learning is that Q learning does not require knowing the transition function T .

- (b) After observing the agent for a while, Adam realized that his assumption of T being deterministic is wrong in one specific way: **when the agent tries to legally move down, it occasionally ends up moving left instead** (except from grid 1 where moving left results in out-of-bound). All other movements are still deterministic.

Suppose we have run the Q updates outlined in the equation above until convergence, to get $Q_{\text{wrong}}^*(s, a)$ under the original assumption of the wrong (deterministic) T . Suppose $Q_{\text{correct}}^*(s, a)$ denotes the Q values under the new correct T . Note that you don't explicitly know the exact probabilities associated with this new T , but you know that it qualitatively differs in the way described above. As prompted below, list the set of (s, a) pairs where $Q_{\text{wrong}}^*(s, a)$ is either an over-estimate or under-estimate of $Q_{\text{correct}}^*(s, a)$.

- (i) [3 pts] List of (s, a) where $Q_{wrong}^*(s, a)$ is an over-estimate. Explain why.

There are two types of (s, a) pairs that are over-estimated.

First, all the (s, a) pairs that have non-zero probability of landing in $s' = 2, 3$, or 4 . This is because $V(2)$, $V(3)$, $V(4)$ will all end up being overestimated. So all (s, a) pairs that use $2, 3, 4$ as s' will be overestimation.

Second, all the (s, a) pairs that start from state $s = 2, 3, 4$ with $a = \text{down}$. This is because they have probability of moving left to farther away from 8 .

To sum up, the (s, a) pairs are: $(1, \text{right})$, $(2, \text{right})$, $(2, \text{down})$, $(3, \text{left})$, $(3, \text{right})$, $(3, \text{down})$, $(4, \text{left})$, $(4, \text{down})$, $(6, \text{up})$, $(7, \text{up})$.

Side note: the state values for $1, 5, 6, 7, 8$ are not affected because the optimal value of those state can be obtained through a sequence of nodes without having to take a "down" action that is affected by this noisy failure

- (ii) [3 pts] List of (s, a) where $Q_{wrong}^*(s, a)$ is an under-estimate (and why):

None

- (c) [2 pts] Suppose that we have a mysterious oracle that can give us either all the correct Q-values $Q(s, a)$ or all the correct state values $V(s)$. Which one do you prefer to be given if you want to use it to find the optimal policy, and why?

Q. When you are using Q, you only need the value itself to determine a good action. However, with value function, you also need the transition function to determine this.

- (d) [2 pts] Suppose that you perform actions in this grid and observe the following episode: $3, \text{right}, 4, \text{down}, 8$ (terminal).

With learning rate $\alpha = 0.2$, discount $\gamma = 0.8$, perform an update of $Q(3, \text{right})$ and $Q(4, \text{down})$. Note that here, we update Q values based on the sampled actions as in TD learning, rather than the greedy actions.

$$Q(3, \text{right}): 7 * (1 - 0.2) + 0.2 * (-1 + 0.8 * 9)$$

$$Q(4, \text{down}): 9 * (1 - 0.2) + 0.2 * (10)$$

Please note that the -1 and 10 comes from the first paragraph of the problem statement (which describes the $R(s, a, s')$)

- (e) [2 pts] One way to encourage an agent to perform more exploration in the world is known as the " ϵ -greedy" algorithm. For any given policy $\pi(s)$, this algorithm says to take the original action $a = \pi(s)$ with probability $(1 - \epsilon)$, and to take a random action (drawn from a uniform distribution over all legal actions) with probability ϵ . If ϵ can be tuned, would you assign it to be a high or low value at the beginning of training? What about at the end of the training? Please answer both questions and justify your choices.

Higher at beginning because want more exploration. We should use that exploration to converge to more optimal strategies at the end, but near the end of training, lower epsilon so we can exploit instead of explore.

- (f) Instead of using the " ϵ -greedy" algorithm, we will now do some interesting exploration with softmax. We first introduce a new type of policy: A stochastic policy $\pi(a|s)$ represents the probability of action a being prescribed, conditioned on the current state. In other words, the policy is now a distribution over possible actions, rather than a function that outputs a deterministic action.

Let's define a new policy as follows:

$$\pi(a|s) = \frac{e^{Q(s,a)}}{\sum_{a'} e^{Q(s,a')}}$$

- (i) [2 pts] Suppose we are at square 3 in the grid and we want to use the originally provided Q values from the table. What is the probability that this policy will tell us to go right? What is the probability that this policy will tell us to go left? Note that the sum over actions prescribed above refers to a sum over legal actions. You may leave your answer in terms of e .

$$\pi(3, \text{right}) = \frac{e^7}{e^8 + e^5 + e^7} = \frac{e^2}{e^3 + e^2 + 1} \approx 0.259$$

$$\pi(3, \text{left}) = \frac{e^5}{e^8 + e^5 + e^7} = \frac{1}{e^3 + e^2 + 1} \approx 0.035$$

(ii) [2 pts] How is this exploration strategy qualitatively different from “ ϵ -greedy”?

This exploration is guided by Q value rather than purely random, so you can explore while still taking some amount of goodness (value) into account.

(g) Your friend Cody argues that we could still explicitly calculate Q updates (like Adam’s approach in part (a)) even if we don’t know the true underlying transition function $T(s, a, s')$, because he believes that our T can be roughly approximated from samples.

(i) [3 pts] Suppose you collect 1,000 transitions from $s = 3, a = \text{Down}$, in the form of (s_{start}, a, s_{end}) . Describe how you can use these samples to compute $T_{approx}(s = 3, a = \text{Down}, s')$, which is an approximation of the true underlying (unknown) $T(s, a, s')$.

(s =3, a = Down, s'= 6)	(s = 3, a= Down, s'=7)
99	901

We can approximate transition function using samples. $p(s = 3, a = \text{Down}, s' = 6) = 0.099, p(s = 3, a = \text{Down}, s' = 7) = 0.901$

(ii) [2 pts] Now perform one step of q-value iteration based on your transition model computed above.

$Q(s = 3, a = \text{Down}) = -1 + 0.099 * Q(s = 6, \text{right}) + 0.901 * Q(S = 7, \text{right}) = -1 + 0.099 * 7 + 0.901 * 8$