

CIS 520, Machine Learning, Fall 2021  
Homework 10  
Due: Tuesday, November 30th, 11:59pm  
Submit to Gradescope

**Instructions.** Please write up your responses to the following problems clearly and concisely. We require you to write up your responses using L<sup>A</sup>T<sub>E</sub>X; we have provided a L<sup>A</sup>T<sub>E</sub>X template, available on Canvas, to make this easier. **Submit your answers in PDF form to Gradescope. We will not accept paper copies of the homework.**

**Collaboration.** You are allowed and encouraged to work together. You may discuss the **written homework** to understand the problem and reach a solution in groups. However, **it is recommended that each student also write down the solution independently and without referring to written notes from the joint session.** You must understand the solution well enough to reconstruct it by yourself. (This is for your own benefit: you have to take the exams alone.)

## Learning Objectives

After completing this assignment, you will be able to:

- Understanding value iteration in reinforcement learning.
- Understanding Q-Learning and Deep Q-Learning.

## Deliverables

This homework can be completed **individually or in group of 2**. You need to make one submission per group. Make sure to add your team member's name on Gradescope when submitting the homework's written and coding parts.

1. **A PDF compilation of hw10\_template.tex with your solutions**
2. **A hw10.ipynb with the functions implemented**

# 1 Reinforcement Learning: MDPs [50 points]

Recall that in a reinforcement learning problem, an agent tries to learn an optimal behavior policy by interacting with an environment. Such a problem is usually formulated as a Markov decision process (MDP), given by a tuple  $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, p, r, \gamma \rangle$ ; here  $\mathcal{S}$  denotes the set of states in the environment,  $\mathcal{A}$  denotes the set of actions available to the agent,  $p$  is the transition probability function,  $r$  is the reward function, and  $\gamma$  is the discount factor. When the MDP is not fully specified, i.e. when the transition probabilities  $p$  or the reward function  $r$  are unknown, one needs to use reinforcement learning techniques (such as Q-learning). However, if the MDP is fully specified, an optimal policy can be found using dynamic programming.

In this problem, you will consider the **Chain** environment shown in Figure 1. Specifically, in this environment, there are 5 states arranged in a chain, and 2 possible actions available to the agent in each state: **f** (“forward”) and **b** (“backward”). Once an agent takes an action, there is some stochasticity in how the environment responds, and correspondingly, which state the agent ends up in/what reward it receives. Each action in a state is depicted by a small black circle; the arrows from actions to states depict possible transitions, labeled by the probability with which that transition occurs given the action and previous state, and the associated reward. As can be seen, when the agent takes the **f** action, it usually (with probability 0.8) moves one step forward, and receives zero reward (unless it is in state 4, in which case it usually stays in state 4 and receives a reward of 6), but it sometimes (with probability 0.2) falls back all the way to state 0 and receives a reward of 3. When the agent takes the **b** action, the reverse happens: it usually (with probability 0.8) falls back to state 0 and receives a reward of 3, but sometimes (with probability 0.2) moves a step forward/stays in state 4 and receives a different reward.

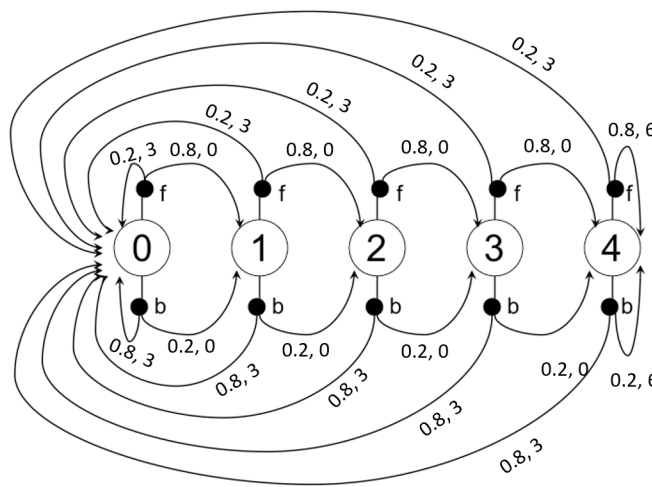


Figure 1: Chain environment

Assume a discount factor of 0.9. The MDP here is then fully specified, and you will use dynamic programming to find an optimal deterministic policy for the agent in this environment.

1. [10 points] Formulate an MDP for the above **Chain** environment by specifying each component of the tuple  $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, p, r, \gamma \rangle$ . For the components  $p$  and  $r$ , write down  $p(s'|s, a)$  and  $r(s, a, s')$  for each setting of  $s, a, s'$  for which the probability/reward is non-zero.
2. [25 points] Find the optimal state-value function  $V^*$ , i.e. find the optimal value  $V^*(s)$  for each state  $s$ .

(Hint: You will need to solve the Bellman optimality equation for  $V^*$ :

$$V^*(s) = \max_a \sum_{s'} p(s'|s, a) (r(s, a, s') + \gamma V^*(s')).$$

You may write a small piece of code to solve this using value iteration. However, you do not need to turn in the code, and we will not grade your code. You can start with any initial value function  $V^0$ , iteratively compute  $V^{t+1}$  by evaluating the RHS above on  $V^t$ , and repeat until convergence (until  $V^{t+1}$  becomes indistinguishable from  $V^t$ .)

3. [15 points] Using your solution to the second part above, find an optimal deterministic policy  $\pi^*$ , i.e. find an optimal action  $\pi^*(s)$  for each state  $s$ .

Hint: Recall that an optimal policy is given by  $\pi^*(s) \in \arg \max_a \underbrace{\sum_{s'} p(s'|s, a) (r(s, a, s') + \gamma V^*(s'))}_{Q^*(s, a)}$ . You

may write a small piece of code to find such an optimal policy. However, you do not need to turn in the code, and we will not grade your code.

## 2 Reinforcement Learning: Q-Learning [50 points]

In this question, we will implement Q-learning for a classic reinforcement learning problem, MountainCar. See the provided notebook for more details.

### 2.1 “Traditional” Q-Learning

Implement the `choose_action`, `update_epsilon`, `update_Q`, and `Qlearning` functions in the notebook. We will implement Q-learning as discussed in lecture, with an  $\epsilon$ -greedy strategy for choosing actions. We will use a simple annealing strategy (`update_epsilon`) by decaying  $\epsilon$  by a constant rate in each episode.

[20 points] After running Q-learning using the parameters provided in the notebook, use the code provided to plot the **Average Reward, Success Rate, and Car Final Position vs Episodes**. Also give an image of your car reaching the flag in the final episode. Include the plots in your typeset submission.

### 2.2 Deep Q-Learning

Now we will take a deep learning approach by using a neural network to represent the Q values. The network will be as follows:

Layers	Size
Fully Connected Layer 1	Input size: 2 (size of state space), Output Size: 100
Tanh 1	—
Fully Connected (Hidden) Layer	Input size: 100, Output size: 100
Tanh 2	—
Fully Connected Layer 2	Input size: 100, Output size: 3 (size of action space)

Table 1: Architecture for the QNetwork

For this architecture, do not use bias terms for the fully connected layers.

Implement the `QNetwork` class, and the `choose_action`, `reward_shaping`, `DeepQlearning` functions.

With respect to the `reward_shaping` function, implement a new reward function such that: The reward is always  $0.5 + \text{position}$  of the car in the next state, and if the position of the car in the next state is greater than 0.5 (a successful episode), add +1 on top of the reward you just computed. Return this reward.

[15 points] After running deep Q-learning using the parameters provided in the notebook, plot the **Success Rate**, and **Car Final Position vs Episodes**. Include the plots in your typeset submission.

## 2.3 Questions

1. [7 points] As part of the Deep QLearning implementation, you implemented a `reward_shaping` function to aid in the learning process. Compare this to the original reward structure in part 2.1 – why do you think this modification of the reward is helpful?
2. [8 points] Compare the **Success Rate** and **Car Final Position** plots between your two implementations. Which algorithm is learning a successful policy more quickly? Briefly comment on potential reasons for any differences in performance.