

---

## CS6700 : Reinforcement Learning

### Written Assignment #1

**Topics:** Intro, Bandits, MDP, Q-learning, SARSA

**Deadline:** 11 March 2022, 11:55 pm

**Name:** <your name here>

**Roll number:** <your roll no. here>

---

- This is an individual assignment. Collaborations and discussions are strictly prohibited.
  - Be precise with your explanations. Unnecessary verbosity will be penalized.
  - Check the Moodle discussion forums regularly for updates regarding the assignment.
  - Type your solutions in the provided L<sup>A</sup>T<sub>E</sub>X template file.
  - **Please start early.**
- 

1. (5 marks) [MDP as Bandit] Consider the following grid world task. The environment is a  $10 \times 10$  grid. The aim is to learn a policy to go from the start state to the goal state in the fewest possible steps. The 4 deterministic actions available are to move one step up, down, left or right. Standard grid world dynamics apply. The agent receives a reward of 0 at each time step and 1 when it reaches the goal. There is a discount factor  $0 < \gamma < 1$ . Formulate this problem as a family of bandit tasks. These tasks are obviously related to one another. Describe the structure of the set up and the rewards associated with each action for each of the tasks, to make it perform similarly to a  $Q$ -learning agent.

**Solution:**

2. (4 marks) [Delayed reward] Consider the task of controlling a system when the control actions are delayed. The control agent takes an action on observing the state at time  $t$ . The action is applied to the system at time  $t + \tau$ . The agent receives a reward at each time step.

- (a) (2 marks) What is an appropriate notion of return for this task?

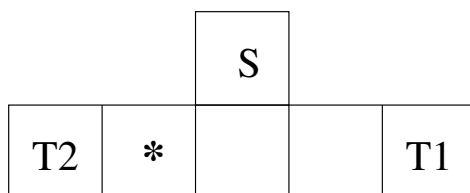
**Solution:**

- (b) (2 marks) Give the TD(0) backup equation for estimating the value function of a given policy.

**Solution:**

3. (5 marks) [Blackwell Optimality] For some policy  $\pi$  in an MDP, if there exists a constant  $k$ , such that for all  $\gamma \in [k, 1)$ ,  $\pi$  is optimal for the discounted reward formulation, then  $\pi$  is said to be *Blackwell optimal*. Consider the gridworld problem shown below, where S denotes a starting state and T1 and T2 are terminal states. The reward for terminating in T1 is +5 and for terminating in T2 is +10. Any transition into the state marked \* has a reward of  $a \in (-\infty, 0)$ . All other transitions have a reward of 0.

For this problem, give a characterization of Blackwell optimal policies, in particular the value  $k$ , parameterized by  $a$ . In other words, for different ranges of  $a$ , give the Blackwell optimal policy, along with the value of  $k$ .



**Solution:**

4. (10 marks) [Jack's Car Rental] Jack manages two locations for a nationwide car rental company. Each day, some number of customers arrive at each location to rent cars. If Jack has a car available, he rents it out and is credited \$ 10 by the national company. If he is out of cars at that location, then the business is lost. Cars become available for renting the day after they are returned. To help ensure that cars are available where they are needed, Jack can move them between the two locations overnight, at a cost of \$ 2 per car moved. We assume that the number of cars requested and returned at each location are Poisson random variables, meaning that the probability that the number  $n$  is  $\frac{\lambda^n}{n!}e^{-\lambda}$ , where  $\lambda$  is the expected number. Suppose  $\lambda$  is 3 and 4 for rental requests at the first and second locations and 3 and 2 for returns. To simplify the problem slightly, we assume that there can be no more than 20 cars at each location (any additional cars are returned to the nationwide company, and thus disappear from the problem) and a maximum of five cars can be moved from one location to the other in one night.
- (a) (4 marks) Formulate this as an MDP. What are the state and action sets? What is the reward function? Describe the transition probabilities (you can use a formula rather than a tabulation of them, but be as explicit as you can about the probabilities.) Give a definition of return and describe why it makes sense.

**Solution:**

- (b) (3 marks) One of Jack's employees at the first location rides a bus home each night and lives near the second location. She is happy to shuttle one car to the second location for free. Each additional car still costs \$ 2, as do all cars moved in the other direction. In addition, Jack has limited parking space at each location. If more than 10 cars are kept overnight at a location (after any moving of cars), then an additional cost of \$ 4 must be incurred to use a second parking lot (independent of how many cars are kept there). These sorts of nonlinearities and arbitrary dynamics often occur in real problems and cannot easily be handled by optimization methods other than dynamic programming. Can you think of a way to incrementally change your MDP formulation above to account for these changes?

**Solution:**

- (c) (3 marks) Describe how the task of Jack's Car Rental could be reformulated in terms of *afterstates*. Why, in terms of this specific task, would such a reformulation be likely to speed convergence? (*Hint:- Refer page 136-137 in RL book 2nd edition. You can also refer to the video at <https://www.youtube.com/watch?v=w3wGvwi336I>*)

**Solution:**

5. (7 marks) [Organ Playing] You receive the following letter:

Dear Friend, Some time ago, I bought this old house, but found it to be haunted by ghostly sardonic laughter. As a result it is hardly habitable. There is hope, however, for by actual testing I have found that this haunting is subject to certain laws, obscure but infallible, and that the laughter can be affected by my playing the organ or burning incense. In each minute, the laughter occurs or not, it shows no degree. What it will do during the ensuing minute depends, in the following exact way, on what has been happening during the preceding minute: Whenever there is laughter, it will continue in the succeeding minute unless I play the organ, in which case it will stop. But continuing to play the organ does not keep the house quiet. I notice, however, that whenever I burn incense when the house is quiet and do not play the organ it remains quiet for the next minute. At this minute of writing, the laughter is going on. Please tell me what manipulations of incense and organ I should make to get that house quiet, and to keep it so.

Sincerely,

At Wits End

- (a) (4 marks) Formulate this problem as an MDP (for the sake of uniformity, formulate it as a continuing discounted problem, with  $\gamma = 0.9$ . Let the reward be +1 on any

transition into the silent state, and -1 on any transition into the laughing state.) Explicitly give the state set, action sets, state transition, and reward function.

**Solution:**

- (b) (2 marks) Starting with simple policy of **always** burning incense, and not playing organ, perform a couple of policy iterations.

**Solution:**

- (c) (2 marks) Finally, what is your advice to “At Wits End”?

**Solution:**

6. (4 marks) [Stochastic Gridworld] An  $\epsilon$ -greedy version of a policy means that with probability  $1-\epsilon$  we follow the policy action and for the rest we uniformly pick an action. Design a stochastic gridworld where a deterministic policy will produce the same trajectories as a  $\epsilon$ -greedy policy in a deterministic gridworld. In other words, for every trajectory under the same policy, the probability of seeing it in each of the worlds is the same. By the same policy I mean that in the stochastic gridworld, you have a deterministic policy and in the deterministic gridworld, you use the same policy, except for  $\epsilon$  fraction of the actions, which you choose uniformly randomly.

- (a) (2 marks) Give the complete specification of the world.

**Solution:**

- (b) (2 marks) Will SARSA on the two worlds converge to the same policy? Justify.

**Solution:**

7. (5 marks) [Contextual Bandits] Consider the standard multi class classification task (Here, the goal is to construct a function which, given a new data point, will correctly predict the class to which the new point belongs). Can we formulate this as contextual bandit problem (Multi armed Bandits with side information) instead of standard supervised learning setting? What are the pros/cons over the supervised learning method. Justify your answer. Also describe the complete Contextual Bandit formulation.

**Solution:**