
CS6700 : Reinforcement Learning

Written Assignment #2

Deadline: 3 May 2024, 11:59 pm

Name:

Roll Number:

- 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^AT_EX template file.
 - **Please start early.**
-

1. (3 marks) Ego-centric representations are based on an agent's current position in the world. In a sense the agent says, I don't care where I am, but I am only worried about the position of the objects in the world relative to me. You could think of the agent as being at the origin always. Comment on the suitability (advantages and disadvantages) of using an ego-centric representation in RL.

Solution:

2. (4 marks) One of the goals of using options is to be able to cache away policies that caused interesting behaviors. These could be to rare state transitions, or access to a new part of the state space, etc. While people have looked at generating options from frequently occurring states in a goal-directed trajectory, such an approach would not work in this case, without a lot of experience. Suggest a method to learn about interesting behaviors in the world while exploring. [*Hint: Think about pseudo rewards.*]

Solution:

3. (2 marks) Consider a navigation task from a fixed starting point to a fixed goal in an arbitrarily complex (with respect to number of rooms, walls, layout, etc) grid-world in which apart from the standard 4 actions of moving North, South, East and West you also have 2 additional options which take you from fixed points A_1 to A_2 and B_1 to B_2 . Now you learn that in any optimal trajectory from the starting point to the goal, you never visit a state that belongs to the set of states that either of the options can visit. Now would you expect your learning to be faster or slower in this domain by including these options in your learning? If it is case dependant, give a potential reason of why in some cases it could slow down learning and also a potential reason for why in some cases it could speed up learning.

Solution:

4. (3 marks) This question requires you to do some additional reading. Dietterich specifies certain conditions for safe-state abstraction for the MaxQ framework. I had mentioned in class that even if we do not use the MaxQ value function decomposition, the hierarchy provided is still useful. So, which of the safe-state abstraction conditions are still necessary when we do not use value function decomposition.

Solution:

5. (1 mark) In any model-based methods, the two main steps are **Planning** and **Model Update**. Now suppose you plan at a rate of F_P (F_P times per time-step) and update the model at a rate of F_M , compare the performance of the algorithm in the following scenarios:

1. $F_P \gg F_M$
2. $F_P \ll F_M$

Solution:

6. (3 marks) In the class, we discussed 2 main learning methods, policy gradient methods and value function methods. Suppose that a policy gradient method uses a class of policies that do not contain the optimal policy; and a value function based method uses a function approximator that can represent the values of the policies of this class, but not that of the optimal policy. Which method would you prefer and why?

Solution:

7. (3 marks) The generalized advantage estimation equation (\hat{A}_t^{GAE}) is defined as below:

$$\hat{A}_t^{GAE} = (1 - \lambda) \left(\hat{A}_t^1 + \lambda \hat{A}_t^2 + \lambda^2 \hat{A}_t^3 + \dots \right)$$

where, \hat{A}_t^n is the n-step estimate of the advantage function and $\lambda \in [0, 1]$.

Show that

$$\hat{A}_t^{GAE} = \sum_{l=0}^{\infty} (\gamma \lambda)^l \delta_{t+l}$$

where δ_t is the TD error at time t , i.e.

$$\delta_t = r_t + \gamma V(s_{t+1}) - V(s_t)$$

Solution:

8. (3 marks) In complex environments, the Monte Carlo Tree Search (MCTS) algorithm may not be effective since it relies on having a known model of the environment. How can we address this issue by combining MCTS with Model-Based Reinforcement Learning (MBRL) technique? Please provide a pseudocode for the algorithm, describing the loss function and update equations.

Solution: