# CS6700 : Reinforcement Learning
## Written Assignment #2

**Topics**: Adv. Value-based methods, FA, PG, AC, POMDP, HRL      **Deadline**: 26 April 2022, 11:59 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 LaTeXtemplate file.

- **Please start early.**

---

1. (4 marks) Recall the four advanced value-based methods we studied in class: PER, Double DQN, Dueling DQN, Expected SARSA. While solving some RL tasks, you encounter the problems given below. Which advanced value-based method would you use to overcome it and why? Give one or two lines of explanation for 'why'.

    (a) (1 mark) Problem 1: In most states of the environment, choice of action doesn't matter.

    > **Solution:**

    (b) (1 mark) Problem 2: Sparse rewards.

    > **Solution:**

    (c) (1 mark) Problem 3: Agent seems to be consistently picking sub-optimal actions during exploitation.

    > **Solution:**

    (d) (1 mark) Problem 4: Environment is stochastic with high negative reward and low positive reward, like in cliff-walking.

    > **Solution:**

2. (4 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:**

3. (12 marks) Santa decides that he no longer has the memory to store every good and bad deed for every child in the world. Instead, he implements a feature-based linear function approximator to determine if a child gets toys or coal. Assume for simplicity that he uses only the following few features:

- Is the child a girl? (0 for no, 1 for yes)
- Age? (real number from $0 - 12$)
- Was the child good last year? (0 for no, 1 for yes)
- Number of good deeds this year
- Number of bad deeds this year

Santa uses his function approximator to output a real number. If that number is greater than his good threshold, the child gets toys. Otherwise, the child gets coal.

(a) (4 marks) Write the full equation to calculate the value for a given child (i.e., $f(s, \vec{\theta}) = \ldots$), where $s$ is a child's name and $\vec{\theta}$ is a weight vector $\vec{\theta} = (\theta(1), \theta(2), \ldots, \theta(5))^{\mathrm{T}}$. Assume child $s$ is described by the features given above, and that the feature values are respectively written as $\phi_s^{\mathrm{girl}}, \phi_s^{\mathrm{age}}, \phi_s^{\mathrm{last}}, \phi_s^{\mathrm{good}}$, and $\phi_s^{\mathrm{bad}}$.

> **Solution:**

(b) (4 marks) What is the gradient $\left( \nabla_{\vec{\theta}} f(s, \vec{\theta}) \right)$ ? I.e. give the vector of partial derivatives

$$\left( \frac{\partial f(s, \vec{\theta})}{\partial \theta(1)}, \frac{\partial f(s, \vec{\theta})}{\partial \theta(2)}, \ldots, \frac{\partial f(s, \vec{\theta})}{\partial \theta(n)} \right)^{\mathrm{T}}$$

based on your answer to the previous question.

> **Solution:**

(c) (4 marks) Using the feature names given above, describe in words something about a function that would make it impossible to represent it adequately using the above linear function approximator. Can you define a new feature in terms of the original ones that would make it linearly representable?

> **Solution:**

4. (6 marks) Recent advances in computational learning theory, have led to the development of very powerful classification engines. One way to take advantage of these classifiers is to turn the reinforcement learning problem into a classification problem. Here the policy is treated as a labeling on the states and a suitable classifier is trained to learn the labels from a few samples. Once the policy is adequately represented, it can be then used in a policy evaluation stage. Can this method be considered a policy gradient method? Justify your answer. Describe a complete method that generates appropriate targets for the classifier.

> **Solution:**

5. (5 marks) Suppose that the system that you are trying to learn about (estimation or control) is not perfectly Markov. Comment on the suitability of using different solution approaches for such a task, namely, Temporal Difference learning, Monte Carlo methods, and Policy Gradient algorithms. Explicitly state any assumptions that you are making. Suppose that the system that you are trying to learn about (estimation or control) is not perfectly Markov. Comment on the suitability of using different solution approaches for such a task, namely, Temporal Difference learning, Monte Carlo methods, and Policy Gradient algorithms. Explicitly state any assumptions that you are making.

**Solution:**

6. (2 marks) We discussed two different motivations for actor-critic algorithms: the original motivation was as an extension of reinforcement comparison, and the modern motivation is as a variance reduction mechanism for policy gradient algorithms. Why is the original version of actor-critic not a policy gradient method?

**Solution:**

7. (5 marks) We typically assume tabula rasa learning in RL and that beyond the states and actions, you have no knowledge about the dynamics of the system. What if you had a partially specified approximate model of the world - one that tells you about the effects of the actions from certain states, i.e., the possible next states, but not the exact probabilities. Nor is the model specified for all states. How will you modify Q learning or SARSA to make effective use of the model? Specifically describe how you can reduce the number of *real* samples drawn from the *world*.

**Solution:**

8. (4 marks) We discussed Q-MDPs in the class as a technique for solving the problem of behaving in POMDPs. It was mentioned that the behavior produced by this approximation would not be optimal. In what sense is it not optimal? Are there circumstances under which it can be optimal?

**Solution:**

9. (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:**

10. (4 marks) One of the goals of using options is to be able to cache away policies that caused interesting behaviors. These could be 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:**