

WRITTEN ASSIGNMENT 8

Due: Friday 04/25/2025 @ 11:59pm EST

Disclaimer

I encourage you to work together, I am a firm believer that we are at our best (and learn better) when we communicate with our peers. Perspective is incredibly important when it comes to solving problems, and sometimes it takes talking to other humans (or rubber ducks in the case of programmers) to gain a perspective we normally would not be able to achieve on our own. The only thing I ask is that you report who you work with: this is **not** to punish anyone, but instead will help me figure out what topics I need to spend extra time on/who to help. When you turn in your solution (please use some form of typesetting: do **NOT** turn in handwritten solutions), please note who you worked with.

Question 1: Reward Function Flavors (25 points)

In lecture, we talked about MDPs that are formulated with a reward function $R(s)$ (i.e. the reward only depends on the current state). However, sometimes MDPs are formulated with a reward function $R(s, a)$ (i.e. a reward function that depends on the action taken), or even $R(s, a, s')$ (i.e. a reward function that depends on the action taken and the way the action is resolved). In this problem, you will show that even though someone may choose one flavor of reward function over another, they are all actually the same:

1. Write the bellman equation that uses $R(s, a)$ and write the bellman equation that uses $R(s, a, s')$
2. Show how an MDP with reward function $R(s, a, s')$ can be converted into a different MDP with reward $R(s, a)$ such that optimal policies in the new MDP correspond exactly to optimal policies in the original MDP.
3. Show how an MDP with reward function $R(s, a)$ can be converted into a different MDP with reward $R(s)$ such that optimal policies in the new MDP correspond exactly to optimal policies in the original MDP.

Question 2: Sum of Discounted Rewards vs. Max Reward (25 points)

In lecture we defined the utility of a trajectory to be some additive combination of the rewards along that trajectory. So far this has taken two forms: additive rewards and discounted rewards. However, what happens if we define the utility of a trajectory as the maximum reward observed in that trajectory? Show that this utility function does not result in stationary preferences between trajectories (i.e. that such an agent may change its preference for the optimal trajectory as a function of time). Is it still possible to define a utility function on trajectories such that a policy which maximizes the expected trajectory utility results in optimal behavior?

Extra Credit: Proof that the Bellman Equation is a Contraction Function (30 points)

In lecture we claimed that the bellman equation is a contraction function. Specifically, we said that, for any two vectors of utilities \vec{u} and \vec{u}' :

$$\|B(\vec{u}) - B(\vec{u}')\|_\infty \leq \gamma \|\vec{u} - \vec{u}'\|_\infty$$

1. Show that, for any functions f and g :

$$|\max_a f(a) - \max_a g(a)| \leq \max_a |f(a) - g(a)|$$

2. Derive an expression for $\left| \left(B(\vec{u}) - B(\vec{u}') \right)(s) \right|$ and then apply the result from part 1 to complete the proof that the bellman equation is a contraction function.