

RBE595 - Week 13 Assignment

Keith Chester

Due date: April 24, 2023

Problem 1

What are the two sources of error in Deep RL with function approximation?

The two common sources are approximation error and the exploration/exploitation trade-off (which we see in a number of other RL algorithms).

When using neural nets to approximate value functions or policy functions in complex environments, the resulting neural net is a function approximation of the environment's true function. The difference is that there can be errors, biases, or edge cases not represented by the approximation. This can be combated by additional training, simulation, and more data depending on your approach; but in a sufficiently complex environment it may be possible that your net can never fully represent every possible situation.

Exploration-exploitation trade off is similar to the issue we faced in earlier algorithms we've explored. If an agent's value function overestimates the value of a set of actions, the agent will become too greedy and exploit this misattributed high-value actions repeatedly, missing out on needed exploration. Conversely, under-estimating means the agent can spend too much time exploring and miss out on opportunities to focus on high returning actions, ultimately sinking the agent's performance.

Problem 2

In TD learning with a neural network what are we trying to minimize? What are we trying to maximize?

In temporal difference (TD) learning with neural networks we are aiming to minimize the error between the predicted and actual values of a particular state/action pair.

We are trying to maximize the ultimate reward we get from the agent by minimizing this error of prediction.

Problem 3

What is the main benefit of deep neural networks in function approximation compared to linear method? What is the linear function approximation useful for?

A linear method to solving a system involves creating a equation that directly describes the behavior of a system. From it we can quickly and efficiently generate agents to fit a system and solve it. Neural networks can approximate any function given a deep enough network, but require training in order to approximate a given system.

When a system is simple enough, a linear method of developing a linear equation is easily done and significantly more efficient to learn than any neural net. The problem is that systems that become more complex become harder to represent, until you reach a level of complexity where a neural net's ability to approximate any function becomes the only method efficient enough to represent and solve the system as an agent.

Problem 4

In DQN, what is the purpose of the target network and value network?

In Deep-Q Networks (DQN), the target network and value network are each a Q-function approximator hoping to correctly evaluate the value of a given state/action pair. The target network is used to generate the target values for the value network to train on. The target network is updated less frequently than the value network, and is used to generate the target values for the value network to train on. By having the two networks update separately and not in lockstep with eachother we have a more stable training process for our agent.

Problem 5

In the Deep Q-learning method, which are true:

- **a.** *epsilon-greedy is required to ensure exploration.*
- **b.** *exploration is taken care of by the randomization provided by experience replay?*

Epsilon-greedy is a method for determining and balancing choosing an action at random rather than choosing the greedy option of the best possible choice of an action presented to the agent at a given state. Experience replay is a method of training networks like DQNs that record momentary states and the resulting values/actions for the agent to repeatedly learn from.

Experience replay isn't strictly for exploration but more for stability of training through allowing the agent to avoid "forgetting" past experiences and situations. The random presentation of states and experiences from the agent's past can be viewed as a form of exploration, but does not allow the agent to choose new unexplored paths quite like epsilon-greedy would.

Thus I would argue that A is true between the two options presented.

Problem 6

True or False: *Value function-based RL methods are oriented towards finding deterministic policies whereas policy search methods are geared towards finding stochastic policies.*

This is true; when we are trying to find a value function we're hoping to create a deterministic and ultimately accurate function that gives us a numerical score for the value of a given state. Policies, however, are stochastic as they allow appropriate exploitation/exploration during training. Likewise, if we look at the popular method of dealing with continuous (versus discrete) action spaces, a popular method is to have networks generate the prerequisites for a distribution (such as mean and variance fed into a normal distribution) to create a stochastic policy for a given state.

Problem 7

What makes the optimization space smooth in policy gradient methods?