

EN7: Function Approximation

Friday, November 18, 2022 8:56 PM

1. point. (RL2e 10.1) On-policy Monte-Carlo control with approximation.
Written: We have not explicitly considered or given pseudocode for any Monte Carlo methods in this chapter.

- What would they be like?

- Why is it reasonable not to give pseudocode for them?

- How would they perform on the Mountain Car task?

Ans 1
In MonteCarlo calculate return at each timestep & then has to back-track & update the state of the episode. Making Q to include reward but not the non-state.

The book hasn't provided the pseudocode for the MonteCarlo methods explicitly in the function approximation because the computational cost of the MonteCarlo method is high in the complexity which makes it unreasonable to use.

Markov method is a variation of TD method so writing the pseudocode will be similar to n-step SARSA not considering terminal state.

In the Mountain Car task, the MonteCarlo method will not learn until the end of the episode, the amount of step increases the computational complexity.

2. 1 point. (RL2e 10.2) Semi-gradient expected SARSA and Q-learning.

Written:

- (a) Give pseudocode for semi-gradient one-step Expected Sarsa for control.

- (b) What changes to your pseudocode are necessary to derive semi-gradient Q-learning?

Ans 2(a) Expected SARSA + Control

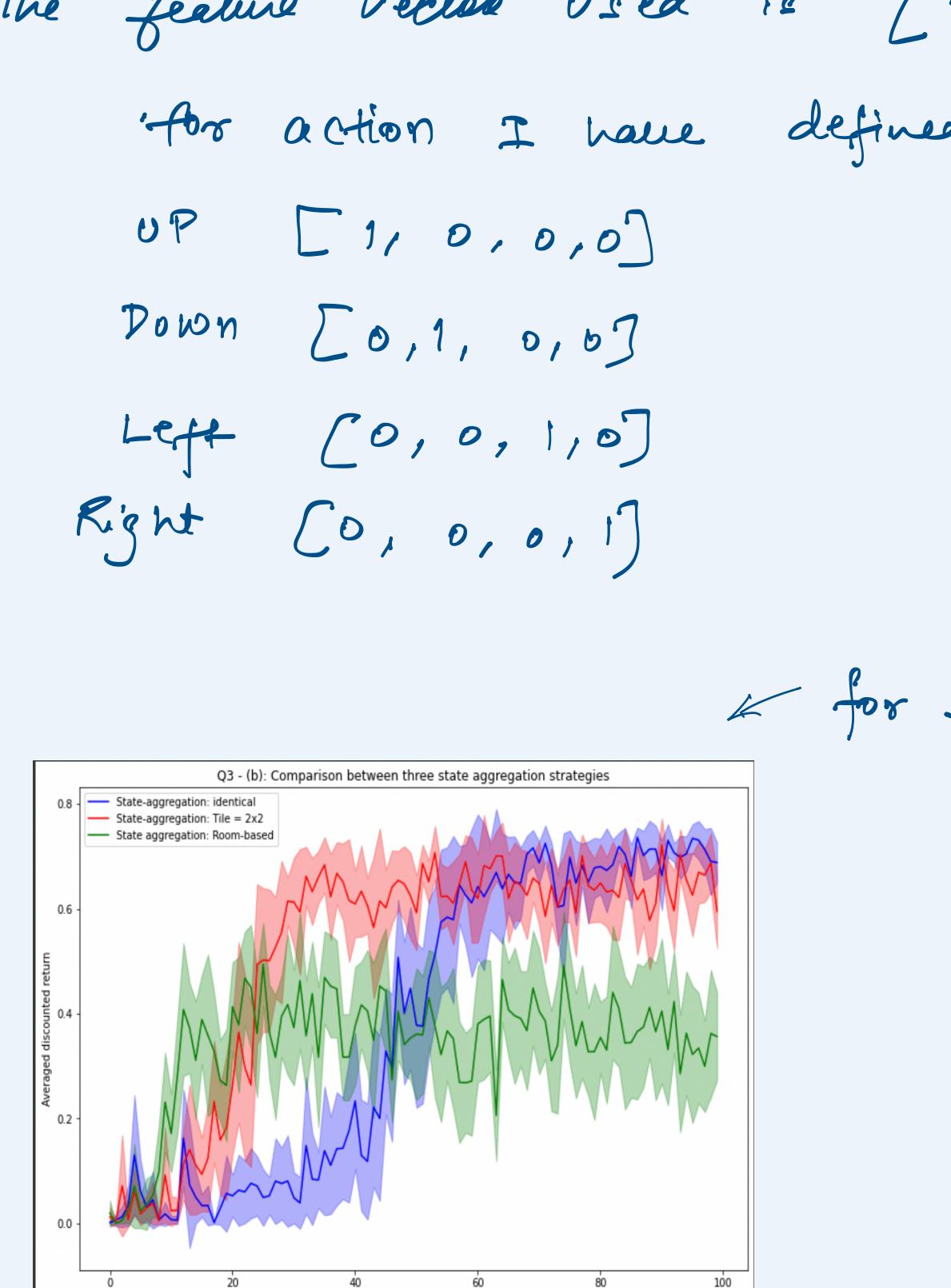
Input: a differentiable action-value function parameterization $\hat{q}: \mathcal{S} \times \mathcal{A} \times \mathbb{R}^d \rightarrow \mathbb{R}$
Algorithm parameters: step size $\alpha > 0$, small $\epsilon > 0$
Initialize value-function weights $\mathbf{w} \in \mathbb{R}^d$ arbitrarily (e.g., $\mathbf{w} = 0$)
Loop for each episode:
 $S, A \leftarrow$ initial state and action of episode (e.g., ϵ -greedy)
Loop for each step of episode:
Take action A , observe R, S'
If S' is terminal:
 $\mathbf{w} \leftarrow \mathbf{w} + \alpha[R - \hat{q}(S, A, \mathbf{w})] \nabla \hat{q}(S, A, \mathbf{w})$
Go to next episode
else
 $\mathbf{w} \leftarrow \mathbf{w} + \alpha[R + \gamma \max_a \hat{q}(S', a, \mathbf{w})] \nabla \hat{q}(S, A, \mathbf{w})$
 $S \leftarrow S'$
 $\downarrow \hat{q}(S', \max_a \hat{q}(S', a, \mathbf{w}))$

Ans 2(b) SEMI Gradient Q Learning

Input: a differentiable action-value function parameterization $\hat{q}: \mathcal{S} \times \mathcal{A} \times \mathbb{R}^d \rightarrow \mathbb{R}$
Algorithm parameters: step size $\alpha > 0$, small $\epsilon > 0$
Initialize value-function weights $\mathbf{w} \in \mathbb{R}^d$ arbitrarily (e.g., $\mathbf{w} = 0$)
Loop for each episode:
 $S, A \leftarrow$ initial state and action of episode (e.g., ϵ -greedy)
Loop for each step of episode:
Take action A , observe R, S'
If S' is terminal:
 $\mathbf{w} \leftarrow \mathbf{w} + \alpha[R - \hat{q}(S, A, \mathbf{w})] \nabla \hat{q}(S, A, \mathbf{w})$
Go to next episode
else
 $\mathbf{w} \leftarrow \mathbf{w} + \alpha[R + \gamma \max_a \hat{q}(S', a, \mathbf{w})] \nabla \hat{q}(S, A, \mathbf{w})$
 $S \leftarrow S'$

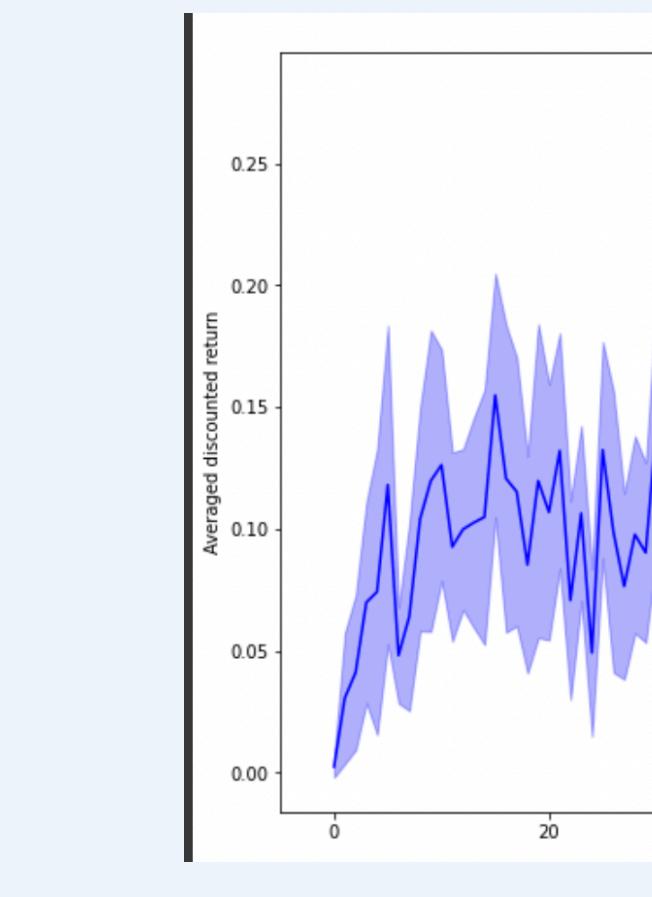
Q3.

a) PLOT + Q3(a)

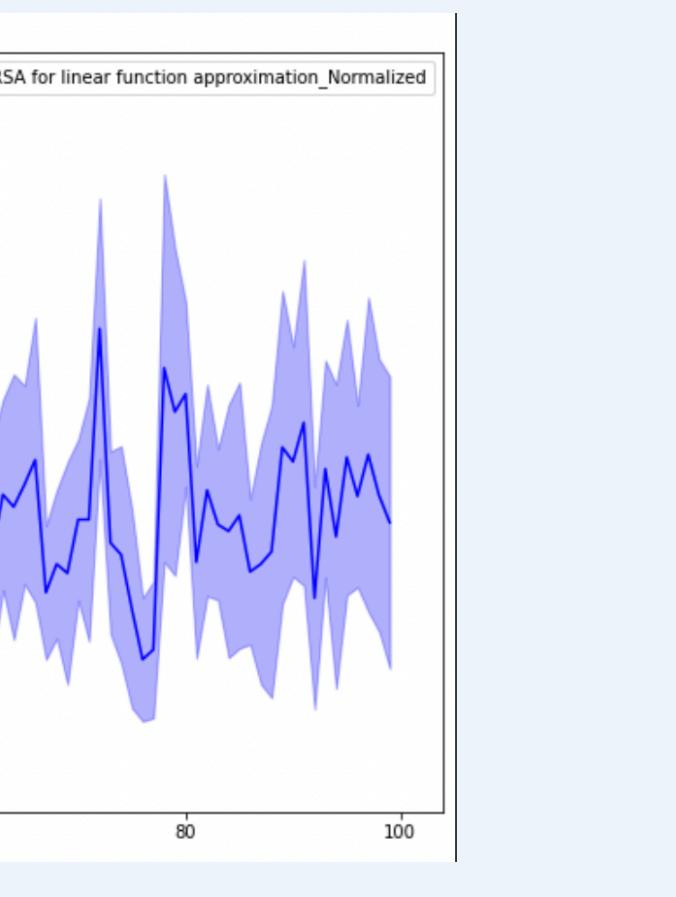


Q3(c)

Semi Gradient SARSA



SEMI GRADIENT Q Learning



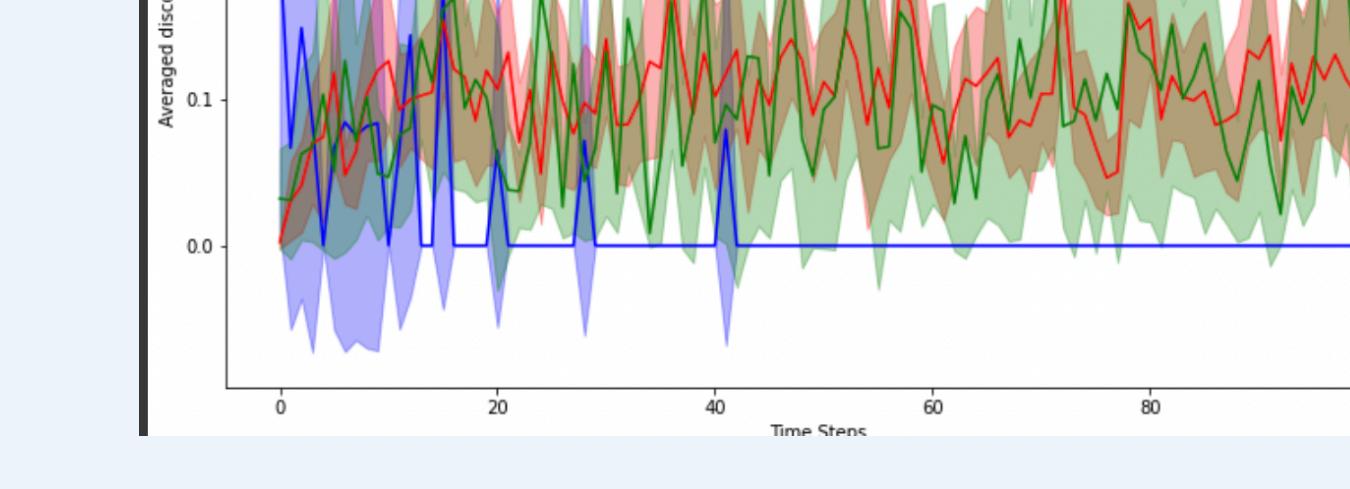
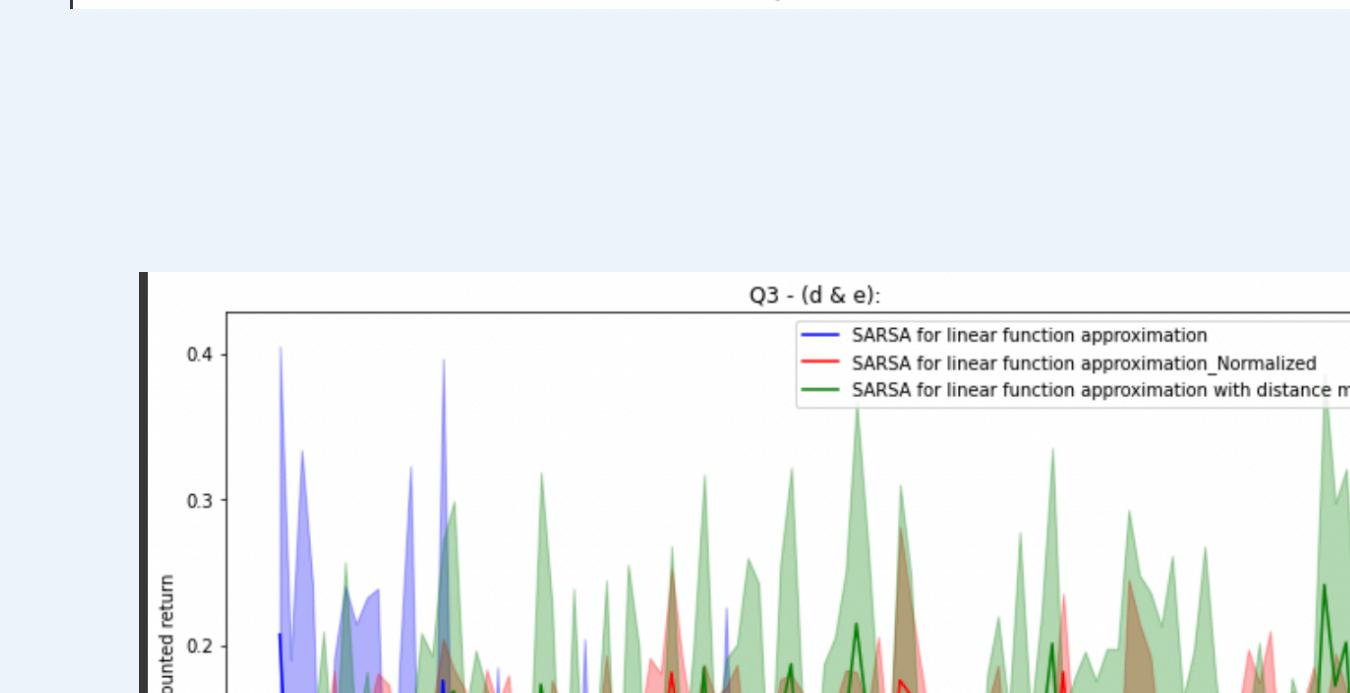
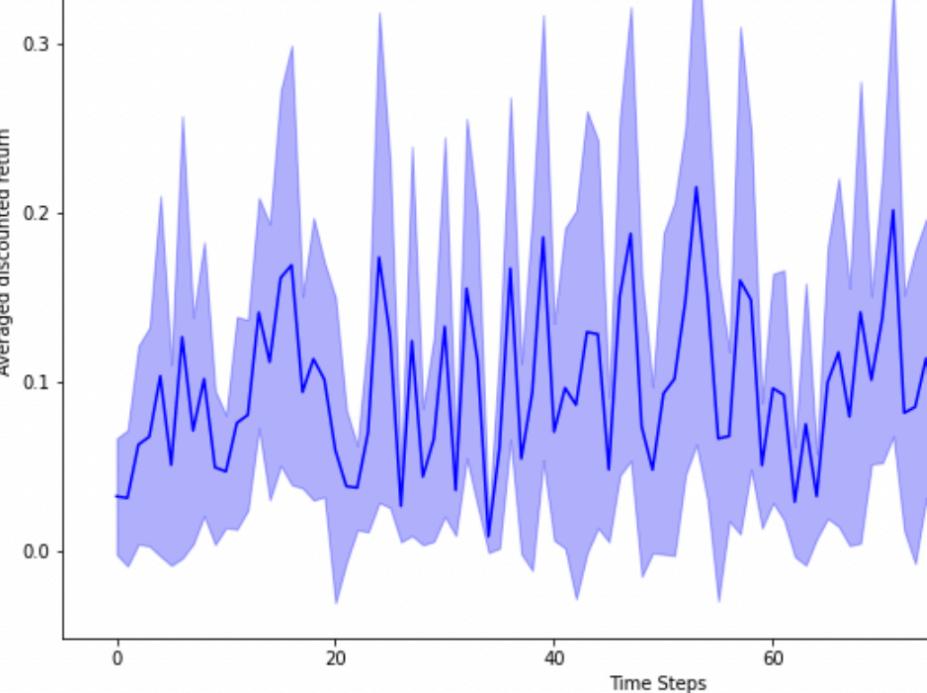
As we can observe with the non aggregated semi gradient SARSA & semi gradient Q learning.

On the semi gradient Q learning the graph has more smooth learning curve than the semi-gradient SARSA, as we can expect as the Q learning chooses from the value of max. but for the state aggregation much improvement is observed as the states get reduced & choosing the max among them also limits.

This shows us we need a good feature design.

Ans 3(e)

Normalized state coordinates: feature $[x/10, y/10, 1]$



On this value seems to vary more compared to the above graph but the returns are better.

When I changed the hyperparameters $\gamma = 0.999$ much better learning curve is observed.

