

Ex-5 Temporal Difference Learning

Thursday, November 3, 2022 6:31 PM

Q.1 Temporal Vs Monte-Carlo

MonteCarlo - Requiring end of the episode to determine state value function

Temporal Difference - Only waits until the next step

meaning at time $t+1 \rightarrow$ TD method uses observed reward R_{t+1}

$$TD \text{ target} = R_{t+1} + \gamma V(s_{t+1})$$

The Scenario where MC would be better

where we have the episodic scenarios

like **Black Jack** where reward only

affects after the termination of Episode,

& updating intermediate state value function will not be of any benefit, it will only take the computational resource

Q-Learning vs SARSA

SARSA

$$\delta(s,a) = \hat{V}(s,a) + \alpha [\tau + \gamma \hat{V}(s',a') - \hat{V}(s,a)]$$

Q-Learning

$$\delta(s,a) = \hat{V}(s,a) + \alpha [\tau + \gamma \hat{V}(s',a') - \hat{V}(s,a)]$$

Q-Learning \leftarrow selects best next action

& SARSA next action depends on the policy it follows

both converges to true value but at different rates.

Q.2 Q-Learning is considered off policy because

it follows the greedy approach i.e goes for the best action instead of sticking to the current policy

therefore Q-Learning is the off policy control method

Q.3 If the action is chosen greedily, still the SARSA & Q-Learning are not same, as

SARSA will choose Greedy action based on the Q value of the previous iterations.

Q-Learning on the other hand take the best action based on the current Q

Hence Q-Learning & SARSA are not same even after making the action greedy

Q.3

Which algorithm is better would be affected if a wider range of α values are used

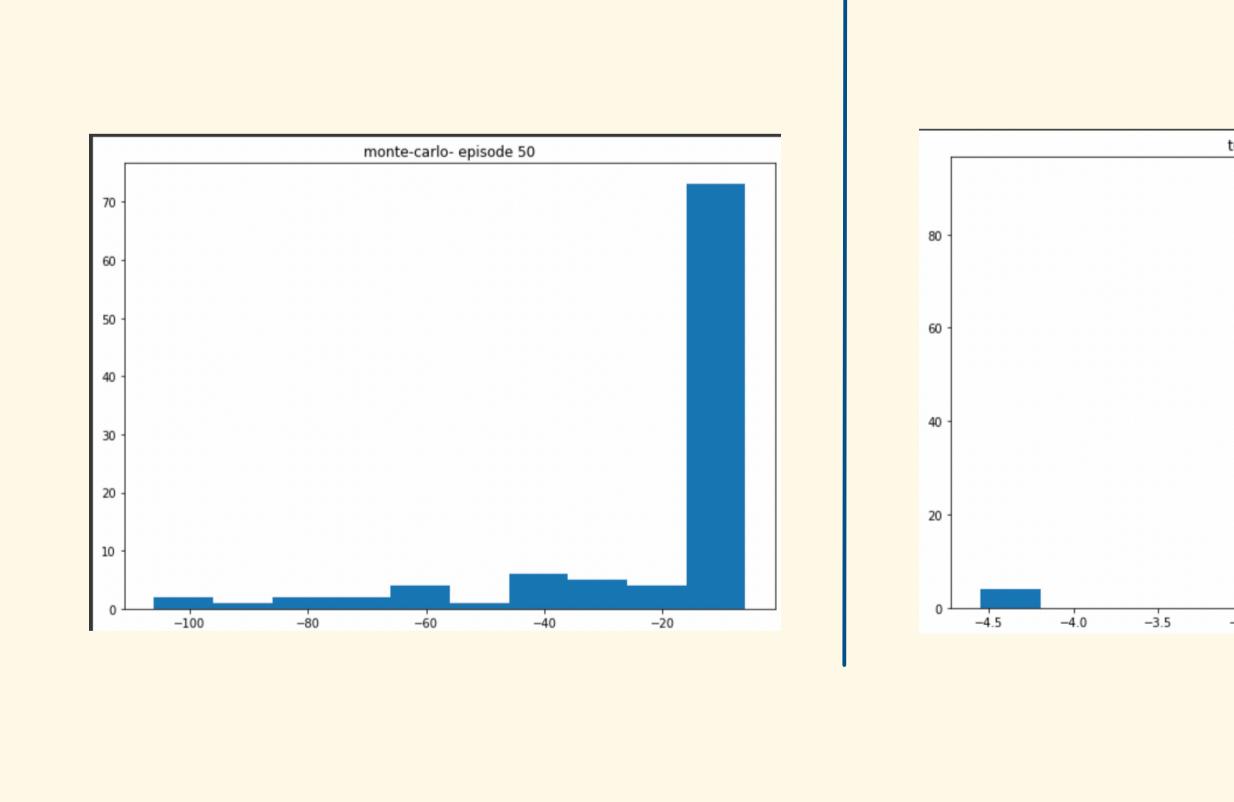
With small α values it's more likely that it converges to the optimal value as we can observe TD methods performing better for the small α , & large α TD would have more oscillation around true values.

Even with wider range of α the MC would become any better, as we update MC only after the episode termination

r

Q.4

4.a Applied on windy gridworld . No King movement



Around (250 - 250) episodes generated

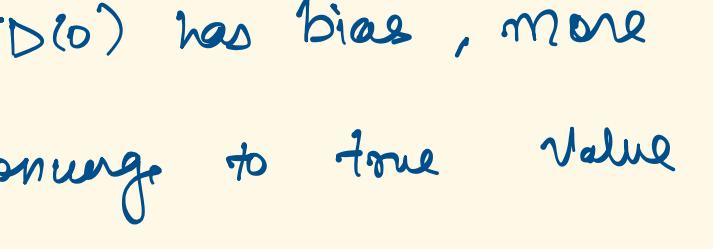
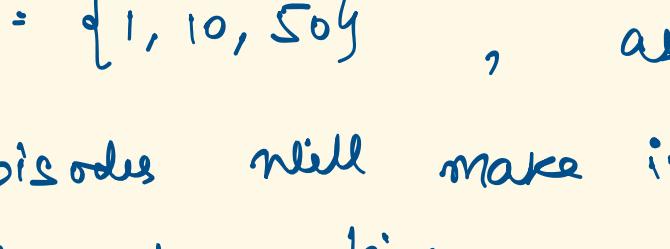
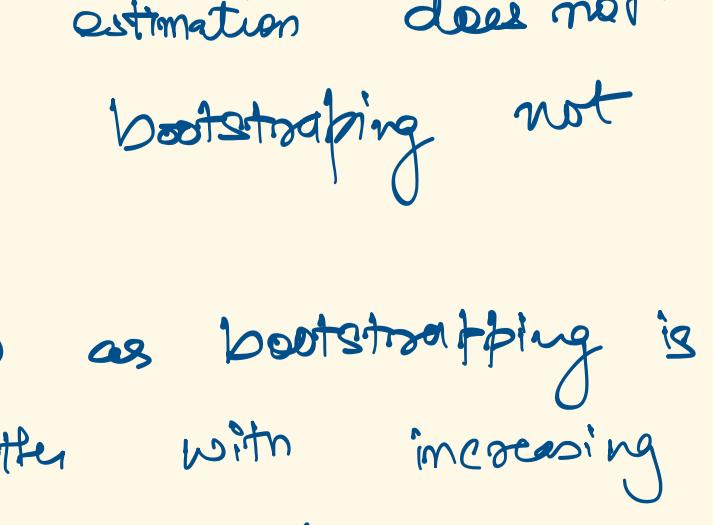
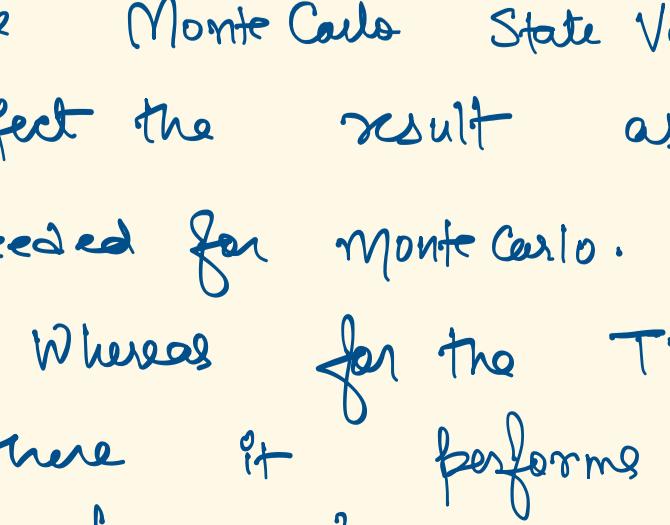
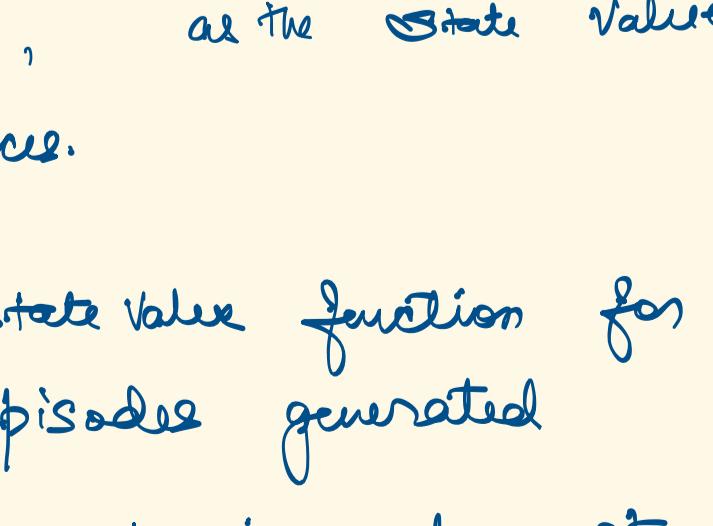
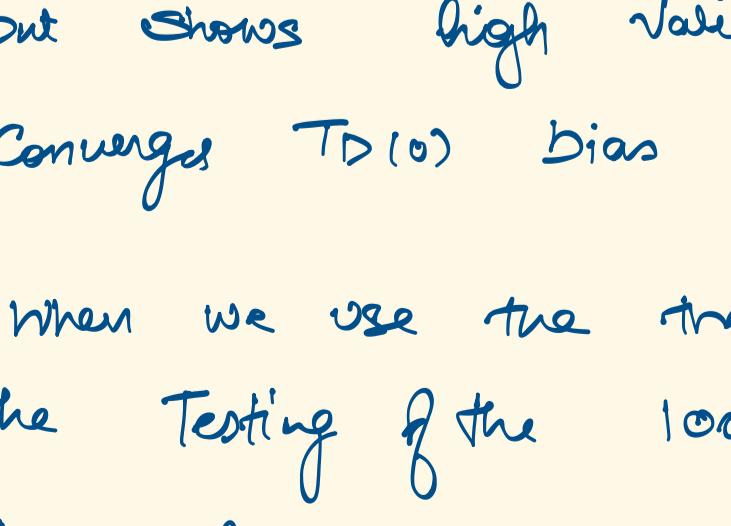
Episodes generated

More than 100

Q.5

Monte Carlo

TD (δ)



Q.6

Monte Carlo give better value estimation with the lower bias but has higher variance

On the other hand, TD(δ) has more stable estimation but shows high variance, as the state value converges TD(δ) bias reduces.

Whereas for the TD(δ) as bootstrapping is there it performs better with increasing $N = 1, 10, 50$, as TD(δ) has bias, more episodes will make it converge to true value & reduce bias.

Plot is better than the 4a but worse than the previous plot with only king's movement

episodes More than 350 less than 450

