Ask a Question

Mentor

Progress

NPTEL » Reinforcement Learning

Unit 9 - Week 7

How does an NPTEL online

Course outline

course work?

Week 0

Week 1

Week 2

Week 3

Week 4

Week 5

Week 6

Week 7

Eligibility Traces

Traces

Week 8

Week 9

Week 10

Week 11

Week 12

DOWNLOAD VIDEOS

Assignment Solutions

NPTEL Resources

Backward View of Eligibility

Thompson Sampling Recap

Reinforcement Learning:

Week 7 Feedback form

Eligibility Trace Control

Quiz : Assignment 7

Assignment 7 The due date for submitting this assignment has passed. Due on 2020-03-18, 23:59 IST. As per our records you have not submitted this assignment. 1) Which of the following is equal to $G_t^{(2)}$ 1 point $R_{t+1} + \gamma V(s_{t+1})$ $R_{t+1} + \gamma R_{t+2} + \gamma^2 V(s_{t+2})$ $R_{t+1} + \gamma R_{t+2} + \gamma^2 V(s_{t+3})$ None of the above No, the answer is incorrect. Score: 0 Accepted Answers: $R_{t+1} + \gamma R_{t+2} + \gamma^2 V(s_{t+2})$ Assertion: It is possible to use the forward view of eligibility traces based target for updates before the end of an episode. 1 point **Reason:** G_t^{λ} is the weighed average of terms of the form G_t^i , where each G_t^i is computable after a finite number of steps i, making G_t^{λ} also computable after a finite number of steps in all cases Both Assertion and Reason are true, and Reason is correct explanation for Assertion Both Assertion and Reason are true, but Reason is not correct explanation for assertion. Assertion is true, Reason is false Both Assertion and Reason are false No, the answer is incorrect. Score: 0 Accepted Answers: Both Assertion and Reason are false 3) Consider the current state is s and the action recommended by the policy, a1, is executed. Which of the following is the correct reasoning behind 1 point setting $\forall a \neq a_1, E_t(s, a) = 0$? Rewards obtained by playing a_1 in s should not be attributed to actions other than a_1 played when in state s previously It assumed that the time steps between reaching s are large enough to decay the eligibility trace to 0 Both (a) and (b) None of the above No, the answer is incorrect. Score: 0 Accepted Answers: Rewards obtained by playing a_1 in s should not be attributed to actions other than a_1 played when in state s previously 4) Assertion: When using an ϵ -greedy exploration strategy, and Watkins $Q(\lambda)$, the epsilon value must be kept low 1 point **Reason:**Traces will become too short if a high value of ϵ is used, negating many of the advantages of using eligibility traces. Both Assertion and Reason are true, and Reason is correct explanation for Assertion Both Assertion and Reason are true, but Reason is not correct explanation for assertion. Assertion is true, Reason is fals Both Assertion and Reason are false No, the answer is incorrect. Score: 0 Accepted Answers: Both Assertion and Reason are true, and Reason is correct explanation for Assertion 5) Recollect estimation policy(π) and behavior policy(μ) in importance sampling. 1 point Assertion: It is NOT necessary for behavior policy of an off-policy method to have non-zero probability of selecting all actions. Reason: If probability of certain actions in estiamation policy is zero, then probability of selecting corresponding actions in behavior policy can be zero, and it would not cause trouble in learning procedure Assertion and Reason are both true and Reason is a correct explanation of Assertion Assertion and Reason are both true and Reason is not a correct explanation of Assertion Assertion is true but Reason is false Assertion and Reason are both false No, the answer is incorrect. Score: 0 Accepted Answers: Assertion and Reason are both true and Reason is a correct explanation of Assertion 6) Which of the following is a way to use replacing traces in SARSA(λ) such that we only update the estimate for the action, a, taken in the state, s, in a 1 point trajectory at timestep t? $E_t(s, a) = \begin{cases} \gamma \lambda E_{t-1}(s, a) + 1 & \text{; if } s_t = s, a_t = a \\ (\gamma \lambda) E_{t-1}(s) & \text{; otherwise} \end{cases}$ $E_t(s, a) = \begin{cases} 1 & \text{; if } s_t = s, a_t = a \\ (\gamma \lambda) E_{t-1}(s) & \text{; otherwise} \end{cases}$ $E_t(s, a) = \begin{cases} \gamma \lambda E_{t-1}(s, a) + 1 & \text{; if } s_t = s, a_t = a \\ 0 & \text{; if } s_t = s, a_t \neq a \\ (\gamma \lambda) E_{t-1}(s) & \text{; otherwise} \end{cases}$ $E_t(s, a) = \begin{cases} 1 & \text{; if } s_t = s, a_t = a \\ 0 & \text{; if } s_t = s, a_t \neq a \\ (\gamma \lambda) E_{t-1}(s) & \text{; otherwise} \end{cases}$ No, the answer is incorrect. Score: 0 Accepted Answers: $E_t(s, a) = \begin{cases} 1 & \text{; if } s_t = s, a_t = a \\ 0 & \text{; if } s_t = s, a_t \neq a \\ (\gamma \lambda) E_{t-1}(s) & \text{; otherwise} \end{cases}$

behaviour policy, and $\rho_t = \frac{\pi(s_t, a_t)}{\rho(s_t, a_t)}$? $E_t(s) = \begin{cases} \rho_t \gamma \lambda E_{t-1}(s) + 1 & \text{; if } s_t = s \\ \rho_t \gamma \lambda E_{t-1}(s) & \text{; otherwise} \end{cases}$

7) In off-policy $TD(\lambda)$, which of the following is a correct way to update eligibility trace, if π is the policy whose Q-values we want to estimate, μ is the **1** point

```
E_t(s) = \begin{cases} \rho_t(\gamma \lambda E_{t-1}(s) + 1) & \text{; if } s_t = s \\ \rho_t \gamma \lambda E_{t-1}(s) & \text{; otherwise} \end{cases}
   E_t(s) = \begin{cases} \gamma \lambda E_{t-1}(s) + 1 & \text{; if } s_t = s \\ \rho_t \gamma \lambda E_{t-1}(s) & \text{; otherwise} \end{cases}
   E_t(s) = \begin{cases} \rho_t(\gamma \lambda E_{t-1}(s) + 1) & \text{; if } s_t = s \\ \gamma \lambda E_{t-1}(s) & \text{; otherwise} \end{cases}
 No, the answer is incorrect.
 Score: 0
Accepted Answers:

E_t(s) = \begin{cases} \rho_t(\gamma \lambda E_{t-1}(s) + 1) & \text{; if } s_t = s \\ \rho_t \gamma \lambda E_{t-1}(s) & \text{; otherwise} \end{cases}
```

Yes

O No

Yes

Accepted Answers:

8) If $\lambda = 1$ then we add only G_t at the end of the episode. In MC we do the same thing

No, the answer is incorrect. Score: 0

Consider the following trajectory s₃, s₂, s₁, s₂, s₂, s₄, s₅, s₆

initial value of eligibility is zero for all states $\gamma^3 \lambda^3 (\gamma^2 \lambda^2 + 2)$

What would be the eligibility value $E_8(s_2)$, for state s_2 after the 8th time step if we use accumulating trace. Discount factor $= \gamma$, trace decay parameter $= \lambda$,

1 point

1 point

1 point

$$\gamma^3 \lambda^3 (\gamma^3 \lambda^3 + \gamma \lambda + 1)$$
 $\gamma^3 \lambda^3 (\gamma^3 \lambda^3 + \gamma \lambda + 1)$
 $\gamma^3 \lambda^3 (\gamma^3 \lambda^3 + \gamma \lambda + 1)$

Accepted Answers:
 $\gamma^3 \lambda^3 (\gamma^3 \lambda^3 + \gamma \lambda + 1)$

Score: 0

TD(1) with replacing trace is identical to first visit Monte Carlo algorithm TD(1) with replacing trace is identical to every visit Monte Carlo algorithm

10) Which of the following statements are true?

TD(1) with replacing trace is identical to last visit Monte Carlo algorithm TD(1) with accumulating trace is identical to every visit Monte Carlo algorithm No, the answer is incorrect.

Accepted Answers: TD(1) with replacing trace is identical to first visit Monte Carlo algorithm TD(1) with accumulating trace is identical to every visit Monte Carlo algorithm