Mentor

Course outline

course work?

Week 0

Week 1

Week 2

Week 3

Week 4

Week 5

Week 6

Week 7

Week 8

Function Approximation

Linear Parameterization

Eligibility Traces

LSTD and LSTDQ

LSPI and Fitted Q

Week 9

Week 10

Week 11

Week 12

DOWNLOAD VIDEOS

Assignment Solutions

NPTEL Resources

Quiz : Assignment 8

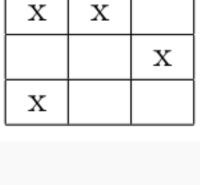
Reinforcement Learning: Week 8 Feedback form

NPTEL » Reinforcement Learning

## Unit 10 - Week 8 Assignment 8 How does an NPTEL online The due date for submitting this assignment has passed. Due on 2020-03-25, 23:59 IST. As per our records you have not submitted this assignment. 1) In which of the following cases would you expect to find that defining the loss of a function appoximator as $\sum_{s \in S} (\hat{V}(s) - V(s))^2$ leads to poor performance. Consider 'relevant' states to be those which are visited frequently when executing near optimal policies Small state space, large percentage of relevant states Large state space, large percentage of relevant states. Large state space, small percentage of relevant states None of the above No, the answer is incorrect. Score: 0 Accepted Answers: Large state space, small percentage of relevant states Assertion: It is not possible to use look-up table based methods to solve continuous state space or continuous action space problems. (Assume) discretization of continuous space is not allowed) Reason: A table which stores the Q value estimate for every state action pair can be built, however the look up time is too large for practical purposes. 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 State Aggregation Methods No, the answer is incorrect. Function Approximation and Score: 0 Accepted Answers: Assertion is true, Reason is false 3) Which of the following is the correct update equation for eligibility traces with linear function approximator, $V(s_t) = w^{\top} \phi(s_t)$ ? $\vec{e}_t = \gamma \lambda \vec{e}_{t-1} + w$ $\vec{e}_t = \gamma \lambda \vec{e}_{t-1} + \phi(s_t)$ $\vec{e}_t = \gamma \lambda \vec{e}_{t-1} + \phi(s_{t-1})$ None of the above No, the answer is incorrect. Score: 0 Accepted Answers: $\vec{e}_t = \gamma \lambda \vec{e}_{t-1} + \phi(s_t)$ Assertion: To solve the given optimization problem for some states with linear function approximator, $\pi_{t+1}(s) = \operatorname{argmax} \ \hat{Q}^{\kappa_t}(s, a)$ we formulate a classification problem for discrete action space. Reason: The given problem is equivalent to solving: $\pi_{t+1} = \operatorname{argmax} \Phi \hat{\Theta}^{\pi_t}$ which can be solved with categorical output, a and elements of $\Phi(s)$ as input 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 true, and Reason is correct explanation for Assertion 5) Which of the following is true about the LSTD and LSTDQ algorithm?

1 point 1 point 1 point 1 point 1 point LSTD estimates matrices  $A = \mathbb{E}[x_t(x_t - \gamma x_{t+1})^{\top}]$  and  $b = \mathbb{E}[R_{t+1}x_t]$  and computes  $\mathbf{w}_{TD} = A^{-1}b$ LSTDQ approximates functions which are linear in feature vector Both LSTD and LSTDQ can reuse samples Both LSTD and LSTDQ are linear function approximation methods No, the answer is incorrect. Score: 0 Accepted Answers:

LSTD estimates matrices  $A = \mathbb{E}[x_t(x_t - \gamma x_{t+1})^{\mathsf{T}}]$  and  $b = \mathbb{E}[R_{t+1}x_t]$  and computes  $\mathbf{w}_{TD} = A^{-1}b$ LSTDQ approximates functions which are linear in feature vector Both LSTD and LSTDQ are linear function approximation methods 6) Consider the following small part of a gridworld. Say blocks with x in it are inaccessible, and empty blocks are accessible. And we want to do linear 1 point parameterization for states by denoting information of neighbouring states. Say the parameterization for state is as follows. '0' shows that the state is accessible, and '1' shows that it is inaccessible. For a state's vector, we define it as at vector of status' of its neigbours in the order (top, right, bottom, left). For the following sub-gridworld, what will be the state vector for the middle state?



(1,1,0,0)(0,1,1,0)No, the answer is incorrect.

 $\bigcirc$  (0,0,1,1)

(1,0,0,1)

Score: 0 Accepted Answers:

(1,1,0,0)

states and for learning, will assume that states within a cluster will have same value function. Assume that clusters are defined in such a way that all states belong to some cluster. Say true/false for the following statement. Statement: A single state can belong to more than one clusters.

Suppose we have a grid-world problem where the number of states is very large(say in millions). So, we define in some way clusters of neighbouring 1 point

True False

No, the answer is incorrect. Score: 0

Accepted Answers: True

Tile coding is a method of state aggregation for gridworld problems. Consider the following statements.

i the number of indicators for each state is equal to number of tilings ii Number of tilings for gridworld must always be 2 iii Tile coding is also a form of Coarse coding Say which of the above statements are true

iii only ○ i, iii

i only ○ i, ii, iii No, the answer is incorrect. Score: 0

i, iii

Accepted Answers:

Note the samples are  $D = \{s_i, a_i, s_i', r_i\}$ , when a policy  $\phi$  is followed.

9) Which of the following are the correct values for A in LSTDQ method.

1 point

1 point

1 point

$$1/L \ \Sigma_{i=1}^L \left[ \phi(s_i)(\phi(s_i) - \gamma \phi(s_i'))^T \right]$$

$$1/L \ \Sigma_{i=1}^L \left[ \phi(s_i, a_i)(\phi(s_i, a_i) - \gamma \phi(s_i'))^T \right]$$

$$1/L \ \Sigma_{i=1}^L \left[ \phi(s_i, a_i)(\phi(s_i, a_i) - \gamma \phi(s_i', \pi(s_i')))^T \right]$$
No, the answer is incorrect.

 $1/L \sum_{i=1}^{L} [\phi(s_i, a_i)(\phi(s_i, a_i) - \gamma \phi(s_i', a_i))^T]$ 

Accepted Answers:  $1/L \sum_{i=1}^{L} [\phi(s_i, a_i)(\phi(s_i, a_i) - \gamma \phi(s_i', \pi(s_i')))^T]$ 

No, the answer is incorrect.

Score: 0

Score: 0

10) In the LSTD method if number of samples  $L \to \infty$  will  $\Phi \theta^{\pi}$  always tend to the true returns of the underlying MDP, for the policy pi?

yes O no

Accepted Answers: