NPTEL » Reinforcement Learning

Unit 5 - Week 3

How does an NPTEL online

Course outline

course work?

Week 0

Week 1

Week 2

Week 3

Policy Search

REINFORCE

MDPs

Week 4

Week 5

Week 6

Week 7

Week 8

Week 9

Week 10

Week 11

Week 12

DOWNLOAD VIDEOS

Assignment Solutions

 R_{n-1}

 R_n

Score: 0

Score: 0

Score: 0

No, the answer is incorrect.

Accepted Answers:

Both (a) and (b)

Score: 0

two

 R_{n-1} R_n

NPTEL Resources

Contextual Bandits

Full RL Introduction

O Quiz: Assignment 3

Reinforcement Learning:

Week 3 Feedback form

O Returns, Value Functions and

	Assignment 3	
As	e due date for submitting this assignment has passed. Due on 2020-02-19, 23:59 per our records you have not submitted this assignment.)
1)	Consider the following policy-search algorithm for a multi-armed binary bandit:	1
	$\forall a, \ \pi_{t+1}(a) = \pi_t(a)(1-\alpha) + \alpha(1_{a=a_t}R_t + (1-1_{a=a_t})(1-R_t))$ where $1_{a_t=a}$ is 1 if $a=a_t$ and 0 otherwise. Which of the following is true for the above algorithm?	
	t is L_{R-I} algorithm	
	t is $L_{R-\epsilon P}$ algorithm	
	It would work well if the best arm had probability of 0.9 of resulting in +1 reward and the next best arm had probability of 0.5 of resulting in +1 reward.	ď
	It would work well if the best arm had probability of 0.3 of resulting in +1 reward and the worst arm had probability of 0.25 of resulting in +1 reward and the worst arm had probability of 0.25 of resulting in +1 reward and the worst arm had probability of 0.25 of resulting in +1 reward and the worst arm had probability of 0.25 of resulting in +1 reward and the worst arm had probability of 0.25 of resulting in +1 reward and the worst arm had probability of 0.25 of resulting in +1 reward and the worst arm had probability of 0.25 of resulting in +1 reward and the worst arm had probability of 0.25 of resulting in +1 reward and the worst arm had probability of 0.25 of resulting in +1 reward and the worst arm had probability of 0.25 of resulting in +1 reward and the worst arm had probability of 0.25 of resulting in +1 reward and the worst arm had probability of 0.25 of resulting in +1 reward and the worst arm had probability of 0.25 of resulting in +1 reward and the worst arm had probability of 0.25 of resulting in +1 reward and the worst arm had probability of 0.25 of resulting in +1 reward and the worst arm had probability of 0.25 of resulting in +1 reward and the worst arm had probability of 0.25 of resulting in +1 reward and the worst arm had probability of 0.25 of resulting in +1 reward and the worst arm had probability of 0.25 of resulting in +1 reward and the worst arm had probability of 0.25 of resulting in +1 reward and the worst arm had probability of 0.25 of resulting in +1 reward and the worst arm had probability of 0.25 of resulting in +1 reward and the worst arm had probability of 0.25 of resulting in +1 reward and the worst arm had probability of 0.25 of resulting in +1 reward and the worst arm had probability of 0.25 of resulting in +1 reward and the worst arm had probability of 0.25 of resulting in +1 reward and the worst arm had probability of 0.25 of resulting in +1 reward and 0.25 of resulting in +1 reward and 0.25 of resulting in +1 reward are worst arm had 0.25 of resulting in +1 reward and 0.2	
Sc	o, the answer is incorrect. ore: 0	
It v	cepted Answers: vould work well if the best arm had probability of 0.9 of resulting in +1 reward and the next best arm had obability of 0.5 of resulting in +1 reward	
	Assertion: Contextual bandits can be modeled as a full reinforcement learning problem. Reason: We can define an MDP with n states where n is the number of bandits and each of them having k actions corresponding to the arms in each	n I
with	single action leading to termination of episode, giving a reward which corresponds to the state and action(bandit and arm)	
	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 and Reason is false Both Assertion and Reason are false	
	the answer is incorrect.	
Sc	ore: 0	
	cepted Answers: sertion and Reason are both true and Reason is a correct explanation of Assertion	
3)	In a full RL setting is it possible for an agent to land up in the same state after performing an action	
	Yes	
	No , the answer is incorrect.	
Sc	ore: 0	
Ac Yes	cepted Answers:	
4)	Late sequence for some full DL problem we are esting according to a policy of At some time two are in a state Curbon we took action A. After four	
4) time	Lets assume for some full RL problem we are acting according to a policy π . At some time t we are in a state S where we took action A_1 . After few steps at time t' the same state S was reached where we performed an action $A_2 (\neq A_1)$. What can you say about the policy π .	
time		1
time	steps at time t' the same state S was reached where we performed an action $A_2 (\neq A_1)$. What can you say about the policy π . It is not a Stationary policy It may be Stationary policy	
time	steps at time t' the same state S was reached where we performed an action $A_2 (\neq A_1)$. What can you say about the policy π . It is not a Stationary policy It may be Stationary policy It is a Non-Stationary policy	
time	steps at time t' the same state S was reached where we performed an action $A_2(\neq A_1)$. What can you say about the policy π . It is not a Stationary policy It may be Stationary policy It is a Non-Stationary policy	
time	steps at time t' the same state S was reached where we performed an action $A_2 (\neq A_1)$. What can you say about the policy π . It is not a Stationary policy It may be Stationary policy It may be a Non-Stationary policy It may be a Non-Stationary policy the answer is incorrect.	
time No Sc Ac	steps at time t' the same state S was reached where we performed an action $A_2 (\neq A_1)$. What can you say about the policy π . It is not a Stationary policy It may be Stationary policy It is a Non-Stationary policy It may be a Non-Stationary policy It may be a Non-Stationary policy	
No Sc Ac	steps at time t' the same state S was reached where we performed an action $A_2 (\neq A_1)$. What can you say about the policy π . It is not a Stationary policy It may be Stationary policy It is a Non-Stationary policy It may be answer is incorrect. Orce: 0 Cepted Answers:	
No Sc Ac	steps at time t' the same state S was reached where we performed an action $A_2 (\neq A_1)$. What can you say about the policy π . It is not a Stationary policy It may be Stationary policy It may be a Non-Stationary policy It may be a Non-Stationary policy It may be a Non-Stationary policy It may be Stationary policy It stationary policy	
No Sc Ac It r. It r. 5)	steps at time t' the same state S was reached where we performed an action $A_2 (\neq A_1)$. What can you say about the policy π . It is not a Stationary policy It may be Stationary policy It may be a Non-Stationary policy It may be a Non-Stationary policy It may be a Non-Stationary policy It may be a Stationary policy It may be a Non-Stationary policy	
No Sc Ac It r. It r. 5)	steps at time t' the same state S was reached where we performed an action $A_2 (\neq A_1)$. What can you say about the policy π . It is not a Stationary policy It may be Stationary policy It is a Non-Stationary policy It may be a Non-Stationary policy The answer is incorrect. The answers: The answers: The answers: The answers: The answers: The answers: The answer is incorrect. The answers: The answers	
No Sc Ac It r. It r. 5)	steps at time t' the same state S was reached where we performed an action $A_2 (\neq A_1)$. What can you say about the policy π . It is not a Stationary policy It may be Stationary policy It may be a Non-Stationary policy It may be a Non-Stationary policy It may be a Non-Stationary policy It may be π stationary policy It may be π stationary policy It may be π stationary policy It may be a Non-Stationary policy	
No Sc Ac It r. It r. 5)	steps at time t' the same state S was reached where we performed an action $A_2 (\neq A_1)$. What can you say about the policy π . It is not a Stationary policy It may be Stationary policy It is a Non-Stationary policy It may be a Non-Stationary policy It may be a Non-Stationary policy It may be Stationary policy It may be a Non-Stationary policy	
No Sc Ac It r. It r. 5)	steps at time t' the same state S was reached where we performed an action $A_2(\neq A_1)$. What can you say about the policy π . It is not a Stationary policy It may be Stationary policy It is a Non-Stationary policy It may be a Non-Stationary policy It may be a Non-Stationary policy It may be a Stationary policy It may be a Non-Stationary policy It may be	
No Sc Ac It r. It r. S)	steps at time t' the same state S was reached where we performed an action $A_2 (\neq A_1)$. What can you say about the policy π . It is not a Stationary policy It may be Stationary policy It may be a Non-Stationary policy It may be a Stationary policy It may be a Non-Stationary policy It	
No Sc Ac At At	steps at time t' the same state S was reached where we performed an action $A_2 (\neq A_1)$. What can you say about the policy π . It is not a Stationary policy It may be Stationary policy It may be a Non-Stationary policy It may be Stationary policy It may be a Non-Stationary policy It may be a Non-Station	
No Sc Ac At At	steps at time t' the same state S was reached where we performed an action $A_2 (\neq A_1)$. What can you say about the policy π . It is not a Stationary policy It may be Stationary policy It may be a Non-Stationary policy It may be a Non-Stationary policy It may be a Non-Stationary policy It may be a Stationary policy It may be a Non-Stationary policy It may be a Non-Stationar	
No Sc Ac It r. It r. 5) No Sc Ac At Ze 6)	steps at time t' the same state S was reached where we performed an action $A_2 (\neq A_1)$. What can you say about the policy π . It is not a Stationary policy It may be Stationary policy It may be Stationary policy It may be a Non-Stationary policy Which of the following statements are true. In RL problems We assume that the agent determines the reward based on the current state and action Our main aim is to get a net positive reward At one time step we can perform only one action Zero rewards may be possible It he answer is incorrect. One: 0 coepted Answers: one time step we can perform only one action ro rewards may be possible Stochastic gradient ascent/descent update occurs in the right direction at every step	;
No Sc Ac It r. It r. 5) No Sc Ac At Ze 6)	steps at time t' the same state S was reached where we performed an action $A_2(\neq A_1)$. What can you say about the policy π . It is not a Stationary policy It may be Stationary policy It may be a Non-Stationary policy It may be a Non-Stationary policy It may be a Non-Stationary policy It may be stationary policy It may be a Non-Stationary policy It may be a Non-Stati	;
No Sc Ac It	steps at time t' the same state S was reached where we performed an action $A_2 (\neq A_1)$. What can you say about the policy π . It is not a Stationary policy It may be Stationary policy It may be Stationary policy It may be a Non-Stationary policy It may be Stationary policy It may be A Non-Stationary policy It may be Stationary policy It	;

