Mentor

Unit 4 - Week 2

How does an NPTEL online

Concentration Bounds

UCB 1 Theorem

Median Elimination

Thompson Sampling

O Quiz : Assignment 2

 Reinforcement Learning: Week 2 Feedback form

PAC Bounds

Week 3

Week 4

Week 5

Week 6

Week 7

Week 8

Week 9

Week 10

Week 11

Week 12

DOWNLOAD VIDEOS

Assignment Solutions

NPTEL Resources

Course outline

course work?

Week 0

Week 1

Week 2

OUCB 1

NPTEL » Reinforcement Learning

Assignment 2 The due date for submitting this assignment has passed. As per our records you have not submitted this assignment.	Due on 2020-02-12, 23:59	IST.
Which of the following happens when using the UCB algorithm? Assume infinite time case.		1 po
Action with highest Q value will be chosen in every iteration O A certain action may never be chosen again after a certain amount of time		
The confidence interval of a certain action decreases in size every time that action is not chosen		
The confidence interval of a certain action increases in size every time that action is chosen		
No, the answer is incorrect. Score: 0		
Accepted Answers: A certain action may never be chosen again after a certain amount of time		
2) In UCB, the term $\sqrt{\frac{2 \ln(n)}{n_i}}$ is added to each arm's Q value and the arm with the highest value after adding the	2 terms is chosen. Which one of the	1 po
Illowing would be an effect of adding $\sqrt{\frac{\ln(n)}{n_j}}$ instead of $\sqrt{\frac{2\ln(n)}{n_j}}$?	2 terms is chosen. Which one of the	
Sub-optimal arms would be chosen less frequently		
Sub-optimal arms would be chosen more frequently		
Makes no change on the frequency of picking sub-optimal arms		
Can't say any of the above		
No, the answer is incorrect. Score: 0		
Accepted Answers: Sub-optimal arms would be chosen less frequently		
3) In a 4-arm bandit problem, after executing 100 iterations of UCB1 algorithm, the estimates of Q values are $Q_{100}(1) = 1.73$, $Q_{100}(2) = 0.83$, $Q_{100}(3) = 1.19$, $Q_{100}(4) = 1.25$ and the number of times each of them are san $Q_{100}(1) = 1.25$, $Q_{100}(2) = 0.83$, $Q_{100}(3) = 0.83$, $Q_{100}(3) = 0.83$, $Q_{100}(4) = 0.83$,	npled are-	1 po
O Arm 1		
O Arm 2		
○ Arm 3 ○ Arm 4		
No, the answer is incorrect.		
Score: 0 Accepted Answers:		
Arm 1		
Now, the epsilon is halved keeping delta unchanged. How many samples would be needed to re-run naive (ϵ, δ) – 400 \odot 800 \odot 1600 \odot 100	PAC algorithm?	
No, the answer is incorrect. Score: 0 Accepted Answers: 800		
 Assertion: The confidence bound of each arm in UCB1 algorithm cannot increase with iterations. 		1 poi
Reason: The n_j term in denominator ensures that the confidence bound remains the same for unselected arm	s and decreases for selected 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		
No, the answer is incorrect.		
Score: 0 Accepted Answers:		
Both Assertion and Reason are false		
6) In UCB 1 Theorem proof when $l = \lceil 8 \ln m/\Delta_i^2 \rceil$, $q_*(a^*) - q_*(i) - 2C_{m,l} = 0$. Which of the following correct $q_*(a^*) - q_*(i) - 2C_{m,T_l(s_i)}$ being greater than 0 for the same l . Note: $C_{m,l} = \sqrt{2 \ln(m)/l}$, $q_*(a^*) - q_*(i) = \Delta_i$	ly explains	1 po
$T_i(s_i) < l \text{ thus } C_{m,T_i(s_i)} < C_{m,l}$		
$T_i(s_i) > m \text{ thus } C_{m,T_i(s_i)} < C_{m,l}.$		
$C_{m,T_{l}(s_{l})} = l \text{ thus } C_{m,T_{l}(s_{l})} < C_{m,l}$		
$T_i(s_i) < l \text{ thus } C_{m,T_i(s_i)} > C_{m,l}$ No, the answer is incorrect.		
Score: 0		
Accepted Answers: $T_i(s_i) > l \text{ thus } C_{m,T_i(s_i)} < C_{m,l}$		
7) Assertion: If a sub optimal arm is replaced by some other arm whose expected reward is farther away from the gret($8\sum_{i\neq a^*} \ln n/\Delta_i + (1+\pi^2/3)\sum_{j=1}^K \Delta_i$) for UCB 1 algorithm will decrease(for $n>3$) Reason: As the reward difference increases the arm will be selected more often	e optimal arm then the bound on	1 po
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		
O Both Assertion and Reason are false		
No, the answer is incorrect. Score: 0		
Accepted Answers: Assertion is true and Reason is false		
	ę	
8) In median elimination method of (ϵ, δ) PAC bounds, what should be the update rule for ϵ and δ if $\epsilon_1 = \frac{\epsilon}{3}$ and ϵ	$\delta = \frac{\delta}{4}$?	1 po

 $\epsilon_{l+1}=rac{\epsilon_l}{2}$ and $\delta_{l+1}=rac{\delta_l}{2}$ $\epsilon_{l+1}=rac{3\epsilon_l}{4}$ and $\delta_{l+1}=rac{\delta_l}{2}$ $\epsilon_{l+1}=rac{2\epsilon_l}{3}$ and $\delta_{l+1}=rac{3\delta_l}{4}$

1 point

1 point

No, the answer is incorrect. Score: 0

 $\epsilon_{l+1}=rac{2\epsilon_l}{3}$ and $\delta_{l+1}=rac{\delta_l}{2}$

Accepted Answers: $\epsilon_{l+1}=rac{2\epsilon_l}{3}$ and $\delta_{l+1}=rac{3\delta_l}{4}$

Follow the notations as

9) In median elimination method for (ϵ, δ) PAC bounds, we claim that for every phase l, for an event E_l , $Pr[E_l] > 1 - \delta_l$

 S_l -- is the set of arms remaining in the l^{th} phase A -- is the maximum of rewards of true best arm in S_l , i.e. in l^{th} phase B -- is the maximum of rewards of true best arm in S_{l+1} , i.e. in $l+1^{th}$ phase Which of the following is E_l ?

 $A+\epsilon_l \leq B$

 $A \leq B + \epsilon_l$ $A \geq B$ $A \geq B + \epsilon_l$ No, the answer is incorrect. Score: 0

 $A \leq B + \epsilon_l$

Accepted Answers:

It is shown that Thompson sampling generally gives better regret bound than UCB In Thompson sampling, eliminating arms is necessary for getting better complexity

After each sampling of an arm, the distributions for the arm becomes thinner(considering bell shaped reward distributions for arms)

All are true No, the answer is incorrect.

10) Which of the following statements is NOT true about Thompson Sampling or Posterior Sam- pling?

Score: 0 Accepted Answers: In Thompson sampling, eliminating arms is necessary for getting better complexity