Practice Assignment

$$\pi(a|s) = rac{e^{ heta_{as}}}{\sum_i e^{ heta_{is}}}$$

$$\sum_i e^{ heta_{is}}$$
 $\bigcirc \quad \pi(a|s) = rac{ heta_a}{\sum_i e^{ heta_i}}$.

$$\bigcirc \quad \pi(a|s) = rac{e^{ heta_a}}{\sum_i heta_i}$$









An optimal policy could not be represented by the parameterisation used to represent the policy





3) Which of the following is the correct update rule for θ (representing the policy) with the policy gradient method?

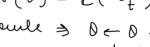


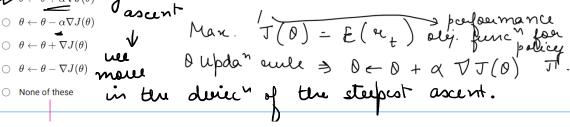
1 point

$$\theta \leftarrow \theta + \alpha \nabla J(\theta)$$
 grading as unt

$$\supset$$
 $heta \leftarrow heta -
abla J(heta)$ moul







4) Which of the following is the correct formulation of cost function for multi arm bandit problem?

1 point

$$J(\theta) = \sum_a q^*(a) \pi_{\theta}(a)$$

$$\bigcirc J(\theta) = \sum_{a} q(a) \pi_{\theta}(a)$$

$$\bigcirc \ \ J(heta) = q^*(a)\pi_ heta(a)$$

$$\bigcirc J(\theta) = \sum_a q^*(a)\pi(a)$$

$$J(0) = E(u_t) = \underbrace{\underbrace{\underbrace{g_{\mu}(a)}}_{a} J_{0}(a)}_{\text{true}}$$

true

conjugate

man of the gaussian.

5) Which of the following is the correct formulation to estimate cost function for multi arm bandit problem from N samples? r_i is the reward received after pulling 1 point an arm at i^{th} timestamp.

$$\hat{\nabla} J(\theta) = rac{1}{N} \sum_{i=1}^{N} r_i rac{
abla \pi_a(\theta)}{\pi_a(\theta)}$$

$$\bigcirc \hat{
abla} J(heta) = rac{1}{N} \sum_{i=1}^N r_i
abla \pi_a(heta)$$

$$\bigcirc \hat{\nabla} J(\theta) = \sum_{i=1}^{N} r_i \frac{\nabla \pi_a(\theta)}{\pi_a(\theta)}$$

$$\bigcirc \quad \hat{
abla} J(heta) = rac{1}{N} \sum_{i=1}^N rac{
abla \pi_a(heta)}{\pi_a(heta)}$$

$$\bigcirc \quad \hat{
abla} J(heta) = rac{1}{N} \sum_{i=1}^N r_i rac{\pi_a(heta)}{
abla \pi_a(heta)}$$

$$\widehat{\nabla} J(0) = \frac{1}{N} \sum_{i=1}^{N} \kappa_i \frac{\nabla J \delta(a_i)}{J \delta(a_i)}$$

kind of Cfd the likelihood impositance sampling.

Graded Assignment

What are the advantages of policy search methods over other approaches?	1 point	
They can lead to simpler solution description	_	
They offe better convergence as compared to function approximation based methods.	_	
If continuous action setting they work better than value function based approaches	-	
They are robust to partial observability.	n / 1-	
None of the above	0 > 4	< 1
2) Which of the following is the correct way to represent policy for policy search methods?	1 point	
$\pi(a s) = \sum_{a \in \sum_i \rho_i = 1} \operatorname{and} 1 \geq \rho_i \geq 0 \forall i$	addition	should be
$\pi(a s) = \overline{ ho_a}, \sum_i ho_i \geq 1$ and $1 \geq ho_i \geq 0 orall i$	200 %	should be
$\pi(a s) = \underbrace{\sum_{i} \rho_{i} = 1 \text{ and } 1 \geq \rho_{i} \geq 0 \forall i}_{\text{o}} $ $= \underbrace{\pi(a s) = \rho_{a}, \sum_{i} \rho_{i} \geq 1 \text{ and } 1 \geq \rho_{i} \geq 0 \forall i}_{\text{o}} $ $= \underbrace{\pi(a s) = \rho_{a}, \sum_{i} \rho_{i} \geq 1 \text{ and } 1 \geq \rho_{i} \geq 0 \forall i}_{\text{o}} $ $= \underbrace{\pi(a s) = \rho_{a}, \sum_{i} \rho_{i} \leq 1 \text{ and } 1 \geq \rho_{i} \geq 0 \forall i}_{\text{o}} $ $= \underbrace{\pi(a s) = \rho_{a}, \sum_{i} \rho_{i} \leq 1 \text{ and } 1 \geq \rho_{i} \geq 0 \forall i}_{\text{o}} $ $= \underbrace{\pi(a s) = \rho_{a}, \sum_{i} \rho_{i} \leq 1 \text{ and } 1 \geq \rho_{i} \geq 0 \forall i}_{\text{o}} $ $= \underbrace{\pi(a s) = \rho_{a}, \sum_{i} \rho_{i} \leq 1 \text{ and } 1 \geq \rho_{i} \geq 0 \forall i}_{\text{o}} $ $= \underbrace{\pi(a s) = \rho_{a}, \sum_{i} \rho_{i} \leq 1 \text{ and } 1 \geq \rho_{i} \geq 0 \forall i}_{\text{o}} $ $= \underbrace{\pi(a s) = \rho_{a}, \sum_{i} \rho_{i} \leq 1 \text{ and } 1 \geq \rho_{i} \geq 0 \forall i}_{\text{o}} $ $= \underbrace{\pi(a s) = \rho_{a}, \sum_{i} \rho_{i} \leq 1 \text{ and } 1 \geq \rho_{i} \geq 0 \forall i}_{\text{o}} $ $= \underbrace{\pi(a s) = \rho_{a}, \sum_{i} \rho_{i} \leq 1 \text{ and } 1 \geq \rho_{i} \geq 0 \forall i}_{\text{o}} $ $= \underbrace{\pi(a s) = \rho_{a}, \sum_{i} \rho_{i} \leq 1 \text{ and } 1 \geq \rho_{i} \geq 0 \forall i}_{\text{o}} $ $= \underbrace{\pi(a s) = \rho_{a}, \sum_{i} \rho_{i} \leq 1 \text{ and } 1 \geq \rho_{i} \geq 0 \forall i}_{\text{o}} $ $= \underbrace{\pi(a s) = \rho_{a}, \sum_{i} \rho_{i} \leq 1 \text{ and } 1 \geq \rho_{i} \geq 0 \forall i}_{\text{o}} $ $= \underbrace{\pi(a s) = \rho_{a}, \sum_{i} \rho_{i} \leq 1 \text{ and } 1 \geq \rho_{i} \geq 0 \forall i}_{\text{o}} $		accounts be f
$ 0$ $\pi(a s)= ho_a,\sum_i ho_i=1$ and 2 $\rho_i\geq 0 \forall i$	lse	
3 q surays be	ane.	
3) Which of the following is the correct update rule for policy parameter θ if the policy is represented in soft-max fashion for a multi arm bandit?	1 point	
$\bigcirc \Delta heta_i = lpha(r-ar{r})(1-\pi(a_i, heta))$	_	
$\Delta heta_i = lpha(r-ar{r})(-\pi(a_i, heta))$	-	
$\Delta \theta_i = \begin{cases} \alpha(r - \bar{r})(1 - \pi(a_i, \theta)), \text{ if action}} (a_i) \text{ s. chosen.} \\ \alpha(r - \bar{r})(-\pi(a_i, \theta)), \text{ otherwise.} \end{cases}$		
$(\alpha(r-\bar{r})(-\pi(a_i,\theta)), \text{ otherwise.})$		
None of these.	-	
Consider following update rule of policy parameter ($ heta$): $\Delta heta_t = lpha(r_t - b_t) rac{\partial ln \pi_{ heta}(a_t)}{\partial a}$	1 point	
₹		
Choose the correct statements of the following:	-	
Baseling b_t gan help figure out if the reward is high or low.	-	
Baseline b_t can be set as the average of the rewards received.	_	
$\partial ln\pi_0(a)$		
is characteristic eligibility, that decides how much a parameter gets updated.	-	
None of these.		
5) Policy gradient methods can be used for continuous action spaces.	1 point	
False función four continuous actions spaces.	_	
o False function for continuous actions spaces.	_	
U U		