

Reinforcement Learning

HW2

309505018
EP/EE22

Problem 1

Recall that $L_{\pi_\theta}(\pi_\theta) = \eta(\pi_\theta) + \sum_s d_\mu(s) \sum_a \pi_\theta(a|s) A^{\pi_\theta}(s, a)$

Then we have

$$\nabla_\theta L_{\pi_\theta}(\pi_\theta) = \sum_s d_\mu(s) \sum_a \nabla_\theta \pi_\theta(a|s) A^{\pi_\theta}(s, a)$$

On the other hand, note that $\eta(\pi_\theta) = \eta(\pi_{\theta_1}) + \sum_s d_\mu(s) \sum_a \pi_\theta(a|s) A^{\pi_\theta}(s, a)$

$$\nabla_\theta \eta(\pi_\theta) = \sum_s \nabla_\theta (d_\mu^\pi(s)) \cdot \sum_a \pi_\theta(a|s) A^{\pi_\theta}(s, a)$$

$$= \sum_s \left\{ (\nabla_\theta d_\mu(s)) \left(\sum_a \pi_{\theta_1}(a|s) A^{\pi_{\theta_1}}(s, a) \right) + d_\mu(s) \sum_a \nabla_\theta \pi_{\theta_1}(a|s) A^{\pi_{\theta_1}}(s, a) \right\}$$

Then,

$$\begin{aligned} \frac{\nabla_\theta \eta(\pi_\theta)}{\theta=\theta_1} &= \sum_s \left\{ (\nabla_\theta d_\mu^\pi(s)) \left(\sum_a \pi_{\theta_1} A^{\pi_{\theta_1}}(s, a) \right) + d_\mu(s) \sum_a \nabla_\theta \pi_{\theta_1}(a|s) A^{\pi_{\theta_1}}(s, a) \right\} \\ &= \sum_s d_\mu(s) \sum_a \nabla_\theta \pi_{\theta_1}(a|s) A^{\pi_{\theta_1}}(s, a) \\ &= \nabla_\theta L_{\pi_{\theta_1}}(\pi_\theta) \end{aligned}$$

Problem 2

$$D(\lambda) = \min \left\{ -(\nabla_{\theta} L_K(\theta)|_{\theta=\theta_K})^T (\theta - \theta_K) + \lambda \left(\frac{1}{2} (\theta - \theta_K)^T H (\theta - \theta_K) - \delta \right) \right\}$$

H is a positive matrix and hence $L(\theta, \lambda)$ is strictly convex

$$\nabla_{\theta} L(\theta, \lambda) = -(\nabla_{\theta} L_K(\theta)|_{\theta=\theta_K}) + \lambda H(\theta - \theta_K)$$

since $L(\theta, \lambda)$ is strictly convex, then a point θ^*, λ^* the global minimum if and only if $\nabla_{\theta} L(\theta, \lambda)|_{\theta=\theta^*, \lambda=\lambda^*} = 0$

$$\Rightarrow \nabla_{\theta} L(\theta^*, \lambda^*) = -(\nabla_{\theta} L_K(\theta)|_{\theta=\theta_K}) + \lambda^* H(\theta^* - \theta_K) = 0$$

$$\Leftrightarrow \theta^* - \theta_K = \frac{1}{\lambda^*} H^T (\nabla_{\theta} L_K(\theta)|_{\theta=\theta_K})$$

Hence,

$$D(\lambda) = L(\theta^*, \lambda) = \frac{1}{\lambda^*} (\nabla_{\theta} L_K(\theta)|_{\theta=\theta_K})^T H^{-1} (\nabla_{\theta} L_K(\theta)|_{\theta=\theta_K})$$

$$+ \lambda^* \left(\frac{1}{2\lambda^*} (\nabla_{\theta} L_K(\theta)|_{\theta=\theta_K})^T H^{-1} H \cdot H^T (\nabla_{\theta} L_K(\theta)|_{\theta=\theta_K}) - \delta \right)$$

$$= \frac{-1}{2\lambda^*} \left((\nabla_{\theta} L_K(\theta)|_{\theta=\theta_K})^T H^{-1} (\nabla_{\theta} L_K(\theta)|_{\theta=\theta_K}) \right) - \lambda^* \delta$$

(a)

$$\lambda^* = \sqrt{\frac{(\nabla_{\theta} L_K(\theta)|_{\theta=\theta_K})^T H^{-1} (\nabla_{\theta} L_K(\theta)|_{\theta=\theta_K})}{2\delta}}$$

(b)

$$\alpha = \frac{1}{\lambda^*}$$