

HW3

111652040. 吳子倫 誌

Problem 1

$$\begin{aligned}
 (i) \quad J_{\pi_{\theta_1}}(\bar{\pi}_{\theta_1}) &= J(\bar{\pi}_{\theta_1}) + \underbrace{\sum_{s \in S} d_{\mu}^{\bar{\pi}_{\theta_1}}(s) \cdot \sum_{a \in A} \bar{\pi}_{\theta_1}(a|s) \cdot A^{\bar{\pi}_{\theta_1}}(s, a)} \\
 &= J(\bar{\pi}_{\theta_1}) + \sum_{s \in S} d_{\mu}^{\bar{\pi}_{\theta_1}}(s) \cdot 0 \\
 &\quad (\text{by lemma 1}) \\
 &= J(\bar{\pi}_{\theta_1}) \quad \#
 \end{aligned}$$

$$\begin{aligned}
 (ii) \quad \nabla_{\theta} J_{\pi_{\theta_1}}(\bar{\pi}_{\theta})|_{\theta=\theta_1} &= \nabla_{\theta} J(\bar{\pi}_{\theta_1})|_{\theta=\theta_1} + \nabla_{\theta} \left(\sum_{s \in S} d_{\mu}^{\bar{\pi}_{\theta_1}}(s) \cdot \sum_{a \in A} \bar{\pi}_{\theta}(a|s) \cdot A^{\bar{\pi}_{\theta_1}}(s, a) \right)|_{\theta=\theta_1} \\
 &\quad (\text{by lemma 2}) \\
 &= \nabla_{\theta} J(\bar{\pi}_{\theta_1})|_{\theta=\theta_1} + \nabla_{\theta} \left(\sum_{s \in S} d_{\mu}^{\bar{\pi}_{\theta_1}}(s) \cdot \sum_{a \in A} \bar{\pi}_{\theta}(a|s) \cdot A^{\bar{\pi}_{\theta_1}}(s, a) \right)|_{\theta=\theta_1} \\
 &= \nabla_{\theta} \left(J(\bar{\pi}_{\theta_1}) + \sum_{s \in S} d_{\mu}^{\bar{\pi}_{\theta_1}}(s) \cdot \sum_{a \in A} \bar{\pi}_{\theta}(a|s) \cdot A^{\bar{\pi}_{\theta_1}}(s, a) \right)|_{\theta=\theta_1} \\
 &\quad (\text{by performance difference lemma}) \\
 &= \nabla_{\theta} (J(\bar{\pi}_{\theta}))|_{\theta=\theta_1}
 \end{aligned}$$

$$\text{lemma 1. } \sum_{a \in A} \bar{\pi}_{\theta_1}(a|s) \cdot A^{\bar{\pi}_{\theta_1}}(s, a) = 0$$

$$\begin{aligned}
 &\sum_{a \in A} \bar{\pi}_{\theta_1}(a|s) \cdot A^{\bar{\pi}_{\theta_1}}(s, a) \\
 &= \sum_{a \in A} \bar{\pi}_{\theta_1}(a|s) \cdot (\underbrace{Q^{\bar{\pi}_{\theta_1}}(s, a) - V^{\bar{\pi}_{\theta_1}}(s)}_{=0}) \\
 &= \underbrace{\sum_{a \in A} \bar{\pi}_{\theta_1}(a|s) \cdot Q^{\bar{\pi}_{\theta_1}}(s, a)}_{=0} - \left(\sum_{a \in A} \bar{\pi}_{\theta_1}(a|s) \right) \cdot V^{\bar{\pi}_{\theta_1}}(s)
 \end{aligned}$$

(by bellman equation)

$$= \underbrace{V^{\bar{\pi}_{\theta_1}}(s)}_{=0} - (1 \cdot V^{\bar{\pi}_{\theta_1}}(s)) = 0 \quad \square$$

$$\text{lemma 1. } \nabla_{\theta} \left(\sum_{s \in S} d_{\mu}^{\pi_{\theta}}(s) \cdot \sum_{a \in A} \pi_{\theta}(a|s) \cdot A^{\pi_{\theta}}(s, a) \right) \Big|_{\theta = \theta_1}$$

$$= \nabla_{\theta} \left(\sum_{s \in S} d_{\mu}^{\pi_{\theta_1}}(s) \cdot \sum_{a \in A} \pi_{\theta}(a|s) \cdot A^{\pi_{\theta_1}}(s, a) \right) \Big|_{\theta = \theta_1}$$

$$\nabla_{\theta} \left(\sum_{s \in S} d_{\mu}^{\pi_{\theta}}(s) \cdot \sum_{a \in A} \pi_{\theta}(a|s) \cdot A^{\pi_{\theta}}(s, a) \right) \Big|_{\theta = \theta_1}$$

$$= \sum_{s \in S} \nabla_{\theta} \left(d_{\mu}^{\pi_{\theta}}(s) \cdot \sum_{a \in A} \pi_{\theta}(a|s) \cdot A^{\pi_{\theta}}(s, a) \right) \Big|_{\theta = \theta_1}$$

$$= \sum_{s \in S} \left(\nabla_{\theta} \left(d_{\mu}^{\pi_{\theta}}(s) \right) \Big|_{\theta = \theta_1} \cdot \underbrace{\sum_{a \in A} \pi_{\theta_1}(a|s) \cdot A^{\pi_{\theta_1}}(s, a)}_{= 0} \right.$$

$$\left. + d_{\mu}^{\pi_{\theta_1}}(s) \cdot \nabla_{\theta} \left(\sum_{a \in A} \pi_{\theta}(a|s) \cdot A^{\pi_{\theta_1}}(s, a) \right) \Big|_{\theta = \theta_1} \right)$$

by lemma 1)

$$= \sum_{s \in S} d_{\mu}^{\pi_{\theta_1}}(s) \cdot \nabla_{\theta} \left(\sum_{a \in A} \pi_{\theta}(a|s) \cdot A^{\pi_{\theta_1}}(s, a) \right) \Big|_{\theta = \theta_1}$$

$$= \nabla_{\theta} \left(\sum_{s \in S} d_{\mu}^{\pi_{\theta_1}}(s) \cdot \sum_{a \in A} \pi_{\theta}(a|s) \cdot A^{\pi_{\theta_1}}(s, a) \right) \Big|_{\theta = \theta_1} \quad \square$$

Problem 2

(a)

① Find $\min_{\theta \in \mathbb{R}^d} L(\theta, \lambda)$

(take derivative)

$$\therefore \nabla_{\theta} L(\theta, \lambda) = -(\nabla_{\theta} L_{\theta_K}(\theta))|_{\theta=\theta_K} + \lambda \cdot (I - \theta_K \cdot (\theta - \theta_K)) = 0$$

$$\therefore \theta - \theta_K = \frac{1}{\lambda} \cdot (I - \theta_K)^{-1} \cdot \nabla_{\theta} L_{\theta_K}(\theta)|_{\theta=\theta_K}$$

$$\begin{aligned} \Rightarrow V(\lambda) &= -(\nabla_{\theta} L_{\theta_K}(\theta)|_{\theta=\theta_K})^T \cdot \frac{1}{\lambda} \cdot (I - \theta_K)^{-1} \cdot (\nabla_{\theta} L_{\theta_K}(\theta)|_{\theta=\theta_K}) - \lambda \cdot \delta \\ &+ \lambda \cdot \frac{1}{2} \cdot (\nabla_{\theta} L_{\theta_K}(\theta)|_{\theta=\theta_K})^T \cdot (I - \theta_K)^{-1} \cdot I \cdot (I - \theta_K)^{-1} \cdot (\nabla_{\theta} L_{\theta_K}(\theta)|_{\theta=\theta_K}) \\ &= \frac{-1}{2 \cdot \lambda} (\nabla_{\theta} L_{\theta_K}(\theta)|_{\theta=\theta_K})^T \cdot (I - \theta_K)^{-1} \cdot (\nabla_{\theta} L_{\theta_K}(\theta)|_{\theta=\theta_K}) - \lambda \cdot \delta \end{aligned}$$

② Find $\text{Max}_{\lambda \geq 0} V(\lambda)$

$$\therefore \nabla_{\lambda} V(\lambda) = \frac{1}{2 \cdot \lambda^2} \cdot (\nabla_{\theta} L_{\theta_K}(\theta)|_{\theta=\theta_K})^T \cdot (I - \theta_K)^{-1} \cdot (\nabla_{\theta} L_{\theta_K}(\theta)|_{\theta=\theta_K}) - \delta = 0$$

$$\therefore \lambda = \sqrt{\frac{1}{2 \cdot \delta} \cdot (\nabla_{\theta} L_{\theta_K}(\theta)|_{\theta=\theta_K})^T \cdot (I - \theta_K)^{-1} \cdot (\nabla_{\theta} L_{\theta_K}(\theta)|_{\theta=\theta_K})} \quad (\text{take } \lambda > 0)$$

Since $\lambda = 0$ boundary value doesn't exist,

$$\lambda^* = \sqrt{\frac{1}{2 \cdot \delta} \cdot (\nabla_{\theta} L_{\theta_K}(\theta)|_{\theta=\theta_K})^T \cdot (I - \theta_K)^{-1} \cdot (\nabla_{\theta} L_{\theta_K}(\theta)|_{\theta=\theta_K})} \quad \#$$

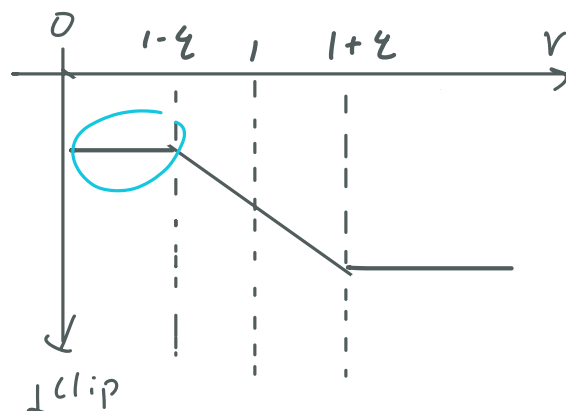
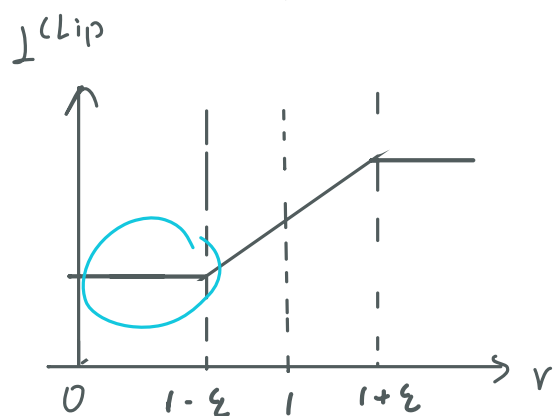
(b)

By (a) ①, $\theta^* - \theta_K = \underbrace{\left(\frac{1}{\lambda^*} \right)}_{\alpha} \cdot (I - \theta_K)^{-1} \cdot \nabla_{\theta} L_{\theta_K}(\theta)|_{\theta=\theta_K}$

$$\Rightarrow \alpha = \sqrt{\frac{2 \cdot \delta}{(\nabla_{\theta} L_{\theta_K}(\theta)|_{\theta=\theta_K})^T \cdot (I - \theta_K)^{-1} \cdot (\nabla_{\theta} L_{\theta_K}(\theta)|_{\theta=\theta_K})}} \quad \#$$

Problem 3

$P_{\tau}(\theta) > 0$	A_{τ}	Return value of min	objective is CLipped	sign of objective	Gradient
$P_{\tau}(\theta) \in [1-\epsilon, 1+\epsilon]$	+	$P_{\tau}(\theta) \cdot A_{\tau}$	No	+	✓
$P_{\tau}(\theta) \in [1-\epsilon, 1+\epsilon]$	-	$P_{\tau}(\theta) \cdot A_{\tau}$	No	-	✓
$P_{\tau}(\theta) < 1-\epsilon$	+	$(1-\epsilon) \cdot A_{\tau}$	yes	+	0
$P_{\tau}(\theta) < 1-\epsilon$	-	$(1-\epsilon) \cdot A_{\tau}$	yes	-	0
$P_{\tau}(\theta) > 1+\epsilon$	+	$(1+\epsilon) \cdot A_{\tau}$	yes	+	0
$P_{\tau}(\theta) > 1+\epsilon$	-	$(1+\epsilon) \cdot A_{\tau}$	yes	-	0



explain

Compared to Figure 1, we observe that $J_{S,A}^{clip}$ give a severe

punishment on loss that is negative,

while $\tilde{J}_{S,A}^{clip}$ is clipped on both sides

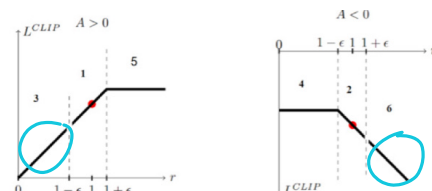


Figure 1: Behavior of the original PPO-clip objective.