

I would like to use my

2-day extension.

RL HW 3

Q2.2 $d_{\pi}(s) = \sum_{t=0}^{\infty} \gamma^t \underbrace{\text{Pr}(s|s_0, \pi)}$

↓ rewrite

$$v_{\pi}(s) = \left[\sum_a \pi(a|s) \text{Pr}(s|a) \right] R(s, a)$$

$$\therefore v_{\pi}(s) = d_{\pi}(s) R(s) \quad (1)$$

R discount · probability · reward

Policy Gradient Thm:

$$\nabla_{\theta} J(\theta) = \nabla v_{\pi}(s_0) \propto \underbrace{\sum_s \text{Pr}(s)}_{\text{constant}} \sum_a q_{\pi}(s, a) \nabla_{\pi}(a|s, \theta)$$

$$= \nabla d_{\pi}(s_0) R(s_0) \propto \sum_{t=0}^{\infty} \gamma^t \nabla_{\theta} v_{\pi}(s_0) \quad \left(\begin{array}{l} \text{deriv. of} \\ (1) \end{array} \right)$$

$$\nabla_{\pi} E_s [V_{\pi}(s)] = \nabla_{\pi} V_{\pi}(s) \quad \text{where } V_{\pi}(s) = \sum_a q_{\pi}(s, a) v_{\pi}(a|s)$$

from Bellman Eqn.

$$\propto \sum_{t=0}^{\infty} \gamma^t \sum_s \text{Pr}(s) \sum_a q_{\pi}(s, a) \nabla_{\pi}(a|s, \theta)$$

Answers to other problems (incl. 2.1 & discussions)
in colab notebook