

RL HW0, 311555031, 黃伯金

1. (a) 假設  $s_0 \in S$ , 且  $a_* = \operatorname{argmax}_{a \in A(s_0)} q_*(s_0, a)$ , 故針對所有 policy  $\pi$ :

$$V_\pi(s_0) = \sum_{a \in A(s_0)} \pi(a|s_0) * q_\pi(s_0, a)$$

$$\leq \sum_{a \in A(s_0)} \pi(a|s_0) * q_*(s_0, a) \quad (q_*(s, a) \geq q_\pi(s, a))$$

$$\leq \sum_{a \in A(s_0)} \pi(a|s_0) * q_*(s_0, a_*) \quad \text{在 optimal action 下得 } q_*(s, a_*)$$

$$= q_*(s_0, a_*) * \sum_{a \in A} \pi(a|s_0) \quad (\sum \pi(a|s_0) = 1)$$

$$= q_*(s_0, a_*) = \max_{a \in A(s_0)} q_*(s_0, a)$$

故針對所有 policy  $\pi$ , 可得  $V_*(s_0) = \max_{\pi} V_\pi(s_0) \leq \max_{a \in A(s_0)} q_*(s_0, a)$

By Contradiction, 先假設  $V_*(s_0) < \max_{a \in A(s_0)} q_*(s_0, a)$

故依據假設可推得存在一 policy  $\phi_t$  使得  $V_{\phi_t}(s_0) = \max_{a \in A(s_0)} q_*(s_0, a) > V_*(s_0)$

但與  $V_*(s_0) = \max_{\pi} V_\pi(s_0)$  的前提產生矛盾!!

故  $V_*(s_0) = \max_{a \in A(s_0)} q_*(s_0, a)$ . For all  $s \in S_a$

代入  $\Delta \pi(s, a) = R_s^a + \gamma \sum_{s'} p_{ss'}^a V_\pi(s')$  可得  $Q_*(s, a) = R_s^a + \gamma \sum_{s'} p_{ss'}^a V_*(s')$

1. (b) For any two action-Value  $Q, \hat{Q}$ ,  $\|T^*(Q) - T^*(\hat{Q})\|_\infty = \max_{(s,a)} (|T^*(Q)(s,a) - T^*(\hat{Q})(s,a)|)$

$$T^*(Q)(s,a) = R_s^a + \gamma \sum_{s'} p_{ss'}^a \left( \max_{a'} Q(s', a') \right)$$

$$\begin{aligned} \text{可得 } \|T^*(Q) - T^*(\hat{Q})\|_\infty &= \cancel{R_s^a} + \gamma \sum_{s'} p_{ss'}^a (\max_{a'} Q(s', a')) - \cancel{R_s^a} - \gamma \sum_{s'} p_{ss'}^a (\max_{a'} \hat{Q}(s', a')) \\ &= \gamma \sum_{s'} p_{ss'}^a \left[ \max_{a'} Q(s', a') - \max_{a'} \hat{Q}(s', a') \right] \end{aligned}$$

故可推得  $T^*$  is a  $\gamma$ -Contraction operator in terms of  $\infty$ -norm.

$$2. L(\pi) = \sum_{a \in A} \left[ \pi(a|s) Q_\pi^k(s,a) - \pi(a|s) (\log \pi(a|s)) \right] - M \left( \sum_{a \in A} \pi(a|s) - 1 \right)$$

$$\frac{\partial L(\pi)}{\partial \pi(a|s)} = Q_\pi^k(s,a) - (\log \pi(a|s) + 1) - M = 0 \quad \text{for all } a \in A$$

$$\text{两边同时} \times \exp \Rightarrow \exp Q_\pi^k(s,a) \times \pi(a|s) \times \frac{1}{e} \times \frac{1}{\exp(M)} = 1$$

$$\text{得到} \quad \frac{\exp Q_\pi^k(s,a)}{e \times \exp(M)} = \pi(a|s) \quad \text{--- ①}$$

$$\text{又 } \sum \pi(a|s) = 1 \quad \text{可得} \quad \sum \frac{\exp Q_\pi^k(s,a)}{e \times \exp(M)} = 1 \quad \text{--- ②}$$

$$\text{由 ①/② 可得} \quad \frac{\exp Q_\pi^k(s,a)}{\sum \exp Q_\pi^k(s,a)} = \pi(a|s) \quad \text{for all } a \in A$$