

Problem 1

(a) Proof of Eq (1): $V_*(s) = \max_a Q_*(s, a)$

$$\textcircled{1} V_*(s) = \max_{\pi} V^{\pi}(s)$$

$$= \max_{\pi} \left(\sum_a \pi(a|s) Q^{\pi}(s, a) \right)$$

↙ 對所有 Q 值取平均 \leq 取 max over a

$$\Rightarrow \sum_a \pi(a|s) Q^{\pi}(s, a) \leq \max_a Q^{\pi}(s, a)$$

$$\leq \max_{\pi} \max_a Q^{\pi}(s, a) = \max_a \max_{\pi} Q^{\pi}(s, a) = \max_a Q_*(s, a)$$

② 對於比較 policy 的好壞，定義：

$$\pi \geq \pi' \text{ if } V^{\pi}(s) \geq V^{\pi'}(s), \forall s$$

對於 optimal policy π_* ，定義：

$$\pi_* \geq \pi, \text{ for all } \pi$$

從①，得到 $V_*(s) \leq \max_a Q_*(s, a)$

若 $V_*(s) < \max_a Q_*(s, a)$ 成立，表示我們可以找到一個比 optimal policy 更好的 policy，為在 state s 時執行 action a 使 $V_*(s) < \max_a Q_*(s, a)$ 成立。但這個 policy 不應該存在，因此 $V_*(s) < \max_a Q_*(s, a)$ 不成立。

所以最後可得 $V_*(s) = \max_a Q_*(s, a)$ 。

Proof of Eq (2): $Q_*(s, a) = R_s^a + \gamma \sum_{s'} P_{ss'}^a V_*(s')$

$$Q_*(s, a) = \max_{\pi} Q^{\pi}(s, a)$$

$$= \max_{\pi} (R_s^a + \gamma \sum_{s'} P_{ss'}^a V^{\pi}(s'))$$

$$= R_s^a + \gamma \sum_{s'} P_{ss'}^a \cdot \max_{\pi} V^{\pi}(s')$$

$$= R_s^a + \gamma \sum_{s'} P_{ss'}^a V_*(s')$$

$$(V_*(s) = \max_{\pi} V^{\pi}(s))$$

Problem 1

$$\begin{aligned}
 (b) \quad \|T^*(Q) - T^*(Q')\|_{\infty} &= \max_{s,a} |T^*(Q)(s,a) - T^*(Q')(s,a)| \\
 &= \max_{s,a} \left| R_s^a + \gamma \sum_{s'} p_{ss'}^a \max_{a'} Q(s',a') - R_s^a - \gamma \sum_{s'} p_{ss'}^a \max_{a'} Q'(s',a') \right| \\
 &= \max_{s,a} \left| \gamma \sum_{s'} p_{ss'}^a \max_{a'} (Q(s',a') - Q'(s',a')) \right| \\
 &= \gamma \max_{s,a} \max_{a'} \left| \sum_{s'} p_{ss'}^a (Q(s',a') - Q'(s',a')) \right| \\
 &\leq \gamma \|Q - Q'\|_{\infty}
 \end{aligned}$$

$\therefore T^*$ is a γ -contraction operator in terms of ∞ -norm.

Problem 2

$$L(\pi) = \sum_{a \in A} (\pi(a|s) Q_{\Omega}^{\pi_k}(s, a) - \pi(a|s) \log \pi(a|s)) - \mu (\sum_{a \in A} \pi(a|s) - 1)$$

and the optimal satisfies $\frac{\partial L(\pi)}{\partial \pi(a|s)} = 0$ for every $a \in A$

假設 action 是離散的，則有 $a_1, a_2, a_3 \dots a_n \in A$ ，

且 optimal 在 $\frac{\partial L(\pi)}{\partial \pi(a_1|s)} = \frac{\partial L(\pi)}{\partial \pi(a_2|s)} = \dots = \frac{\partial L(\pi)}{\partial \pi(a_n|s)} = 0$

將 $L(\pi)$ 的 sigma 寫開。 $L(\pi)$ 為以下式子的總和：

$$\begin{cases} \pi(a_1|s) Q_{\Omega}^{\pi_k}(s, a_1) - \pi(a_1|s) \log \pi(a_1|s) - \mu \pi(a_1|s) \\ \pi(a_2|s) Q_{\Omega}^{\pi_k}(s, a_2) - \pi(a_2|s) \log \pi(a_2|s) - \mu \pi(a_2|s) \\ \vdots \\ \pi(a_n|s) Q_{\Omega}^{\pi_k}(s, a_n) - \pi(a_n|s) \log \pi(a_n|s) - \mu \pi(a_n|s) - \mu \end{cases}$$

分別求得對 $\pi(a_1|s), \pi(a_2|s), \dots, \pi(a_n|s)$ 的偏微分並令其 = 0

$$\frac{\partial L(\pi)}{\partial \pi(a_1|s)} = Q_{\Omega}^{\pi_k}(s, a_1) - (\log \pi(a_1|s) + 1) - \mu = 0$$

$$\frac{\partial L(\pi)}{\partial \pi(a_2|s)} = Q_{\Omega}^{\pi_k}(s, a_2) - (\log \pi(a_2|s) + 1) - \mu = 0$$

\vdots

$$\frac{\partial L(\pi)}{\partial \pi(a_n|s)} = Q_{\Omega}^{\pi_k}(s, a_n) - (\log \pi(a_n|s) + 1) - \mu = 0$$

綜合以上關係可得聯立方程式：

解此方程：

$$e^{Q_{\Omega}^{\pi_k}(s, a_1) - \mu - 1} + e^{Q_{\Omega}^{\pi_k}(s, a_2) - \mu - 1} + \dots + e^{Q_{\Omega}^{\pi_k}(s, a_n) - \mu - 1} = 1$$

同乘 $e^{\mu+1}$ $\rightarrow e^{Q_{\Omega}^{\pi_k}(s, a_1)} + e^{Q_{\Omega}^{\pi_k}(s, a_2)} + \dots + e^{Q_{\Omega}^{\pi_k}(s, a_n)} = e^{\mu+1}$

$$\begin{cases} \pi(a_1|s) = e^{Q_{\Omega}^{\pi_k}(s, a_1) - \mu - 1} \\ \pi(a_2|s) = e^{Q_{\Omega}^{\pi_k}(s, a_2) - \mu - 1} \\ \vdots \\ \pi(a_n|s) = e^{Q_{\Omega}^{\pi_k}(s, a_n) - \mu - 1} \\ \pi(a_1|s) + \pi(a_2|s) + \dots + \pi(a_n|s) = 1 \end{cases}$$

改寫原聯立方程式:

$$\begin{cases} \pi(a_1|s) = e^{Q_{\Omega}^{\pi_k}(s, a_1) - \mu - 1} = \frac{e^{Q_{\Omega}^{\pi_k}(s, a_1)}}{e^{\mu+1}} \\ \pi(a_2|s) = e^{Q_{\Omega}^{\pi_k}(s, a_2) - \mu - 1} = \frac{e^{Q_{\Omega}^{\pi_k}(s, a_2)}}{e^{\mu+1}} \\ \vdots \\ \pi(a_n|s) = e^{Q_{\Omega}^{\pi_k}(s, a_n) - \mu - 1} = \frac{e^{Q_{\Omega}^{\pi_k}(s, a_n)}}{e^{\mu+1}} \end{cases}$$

將 $e^{\mu+1} = e^{Q_{\Omega}^{\pi_k}(s, a_1)} + e^{Q_{\Omega}^{\pi_k}(s, a_2)} + \dots + e^{Q_{\Omega}^{\pi_k}(s, a_n)} = e^{\mu+1}$ 代入,

可得:

$$\begin{cases} \pi(a_1|s) = \frac{e^{Q_{\Omega}^{\pi_k}(s, a_1)}}{e^{\mu+1}} = \frac{e^{Q_{\Omega}^{\pi_k}(s, a_1)}}{e^{Q_{\Omega}^{\pi_k}(s, a_1)} + e^{Q_{\Omega}^{\pi_k}(s, a_2)} + \dots + e^{Q_{\Omega}^{\pi_k}(s, a_n)}} \\ \vdots \\ \pi(a_n|s) = \frac{e^{Q_{\Omega}^{\pi_k}(s, a_n)}}{e^{\mu+1}} = \frac{e^{Q_{\Omega}^{\pi_k}(s, a_n)}}{e^{Q_{\Omega}^{\pi_k}(s, a_1)} + e^{Q_{\Omega}^{\pi_k}(s, a_2)} + \dots + e^{Q_{\Omega}^{\pi_k}(s, a_n)}} \end{cases}$$

$$= \frac{e^{Q_{\Omega}^{\pi_k}(s, a)}}{\sum_{a \in A} e^{Q_{\Omega}^{\pi_k}(s, a)}}$$

改寫上式可得知 $L(\pi)$ 在

$$\pi(\cdot|s) = \frac{e^{Q_{\Omega}^{\pi_k}(s, \cdot)}}{\sum_{a \in A} e^{Q_{\Omega}^{\pi_k}(s, a)}} \quad \text{有最大值。}$$