

$$I. V^*(s) \stackrel{\text{def}}{=} \max_{\pi} V_{\pi}(s) = \max_{\pi} \sum_{a \in A(s)} \pi(a|s) Q_{\pi}(s, a)$$

$$\text{Prove: } V^*(s) = \max_{\pi} \max_{a \in A(s)} Q_{\pi}(s, a)$$

$$\begin{aligned} &= \max_{\pi} \max_a \mathbb{E}_{\pi} [R_t | s_t = s, a_t = a] \\ &= \max_{\pi} \max_a \mathbb{E}_{\pi} \left[\sum_{k=0}^{+\infty} \gamma^k r_{t+k+1} | s_t = s, a_t = a \right] \\ &= \max_{\pi} \max_a \mathbb{E}_{\pi} \left[r_{t+1} + \gamma \sum_{k=0}^{+\infty} \gamma^k r_{t+k+2} | s_t = s, a_t = a \right] \\ &= \max_{a \in A(s)} \max_{\pi} \mathbb{E}_{\pi} [\dots] \\ &= \max_a \left[\mathbb{E}_{\pi} \{ r_{t+1} | s_t = s, a_t = a \} + \gamma \max_{\pi} \mathbb{E}_{\pi} \{ r_{t+1} | s_{t+1} = s, a_{t+1} = a \} \right] \\ &= \max_a \mathbb{E} \{ r_{t+1} + \gamma V^*(s_{t+1}) | s_t = s, a_t = a \} \end{aligned}$$

The same logic applies to derivation of Q^* :

$$Q^*(s, a) = \mathbb{E} \{ r_{t+1} + \gamma \max_{a'} Q^*(s_{t+1}, a') | s_t = s, a_t = a \}$$

II. Let π and π' be two policies such that

$$\forall s \in S \quad Q_{\pi}(s, \pi'(s)) \geq V_{\pi}(s)$$

$$\text{Then } V_{\pi'}(s) \geq V_{\pi}(s) \quad (\Leftrightarrow \pi' \geq \pi)$$

$$\begin{aligned} \triangleright V_{\pi}(s) &\leq Q_{\pi}(s, \pi'(s)) = \mathbb{E}_{\pi'} [r_{t+1} + \gamma V_{\pi}(s_{t+1}) | s_t = s] \\ &\leq \mathbb{E}_{\pi'} [r_{t+1} + \gamma Q_{\pi}(s_{t+1}, \pi'(s_{t+1})) | s_t = s] \\ &= \mathbb{E}_{\pi'} [r_{t+1} + \gamma \mathbb{E}_{\pi'} \{ r_{t+2} + \gamma V_{\pi}(s_{t+2}) | s_{t+1} = s_{t+1} \} | s_t = s] \\ &= \mathbb{E}_{\pi} [r_{t+1} + \gamma r_{t+2} + \gamma^2 V_{\pi}(s_{t+2}) | s_t = s] \\ &\leq \mathbb{E}_{\pi} [r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \gamma^3 V_{\pi}(s_{t+3}) | s_t = s] \\ &\vdots \\ &\leq V_{\pi'}(s) \end{aligned} \quad \square$$

In Policy Improvement:

$$\pi'(s) = \operatorname{argmax}_a Q_{\pi}(s, a)$$

$$Q_{\pi}(s, \pi'(s)) \geq Q_{\pi}(s, a) \quad \forall a \in A$$

$$V_{\pi}(s) = \mathbb{E}_{a \sim \pi(\cdot|s)} Q_{\pi}(s, a) \leq Q_{\pi}(s, \pi'(s)) \Rightarrow \pi' \geq \pi$$

$$\text{When } \pi' = \pi \Rightarrow V_{\pi'} = V_{\pi}$$

then it's optimal, since it satisfies:

$$V_{\pi'}(s) = \max_a \sum_{r, s'} p(r, s' | s, a) [r + \gamma V_{\pi'}(s')]$$