

Lec-13, ITS27, 24-25

Policy improvement theorem:-

Let  $\pi$  &  $\pi'$  be pair of **deterministic** policies s.t.  $\forall s \in \mathcal{S}$

$$Q_{\pi}(s, \pi'(s)) \geq V_{\pi}(s) \quad \text{--- (1)}$$

then  $\pi'$  must be as good as or better than  $\pi$ , i.e., it must obtain  $\geq$  expected return from all states  $s \in \mathcal{S}$

$$V_{\pi'}(s) \geq V_{\pi}(s) \quad \text{--- (2)}$$

now, if there is strict inequality at any state in (1)  $\Rightarrow$  " in (2) for that state.

$$\begin{aligned} \pi(s) &= a_1 \\ a_2 &\neq \pi(s) \end{aligned}$$

$$\text{if } Q_{\pi}(s, a_2) > V_{\pi}(s)$$

then following  $a_2$  is  $s$  always should be beneficial

We know that  $\forall s \in \mathcal{S}$

$$V_{\pi}(s) \leq Q_{\pi}(s, \pi'(s))$$

$$= E[\underbrace{R_{t+1} + \gamma V_{\pi}(S_{t+1})}_X \mid S_t = s, A_t = \pi'(s)]$$

$$= E_{\pi'}[R_{t+1} + \gamma V_{\pi}(S_{t+1}) \mid S_t = s]$$

R.H.S  $\sum_{s'} \sum_a x p(s', a \mid s)$

$$\sum_{s'} \sum_a \sum_{\pi'} \pi'(a \mid s) p(s', a \mid s, a) x$$

because  $\pi'$  is det.  $\downarrow$

$$\sum_{s'} \sum_a p(s', a \mid s, \pi'(s)) x$$

$$x, y, z = x + y$$

$$E[z]$$

$$z_1 = x_1 + y_1$$

$$\sum_{x,y} (x+y) p(x,y)$$

$$p(x \mid y) = \sum_z p(x, z \mid y)$$

$$= \sum_z p(z \mid y) p(x \mid y, z)$$

$$E[x \mid y, z]$$

$$= \sum x p(x \mid y, z)$$

$$p(z \mid y)$$

$$p(x \mid y, z)$$

$$p(x, z \mid y) =$$

$$p(z \mid y) p(x \mid y, z)$$

$$\begin{aligned}
 & E_{\pi'} [R_{t+1} + \gamma V_{\pi}(S_{t+1}) \mid S_t = s] \leq E_{\pi'} [R_{t+1} + \gamma q_{\pi}(S_{t+1}, \\
 & \quad \pi'(S_{t+1})) \mid S_t = s] \\
 & \left( \because V_{\pi}(s) \leq q_{\pi}(s, \pi'(s)) \quad \forall s \in \mathcal{S} \right. \\
 & \quad \left. \rightarrow \right. \\
 & = E_{\pi'} [R_{t+1} + \gamma E_{\pi'} [R_{t+2} + \gamma V_{\pi}(S_{t+2}) \mid S_{t+1}, A_{t+1} = \pi'(S_{t+1})] \mid \\
 & \quad S_t = s]
 \end{aligned}$$