

Lec-9, 1T567, 24-25

Without proof - there exists at least one policy that is better than or equal to all other policies. This is the optimal policy (or policies)  $\triangleq \pi^*$

↳ more than one are possible, but they share the same  $V_{\pi^*}(s)$ , called the optimal state-value function  $V^*(s)$

$$V^*(s) \triangleq \max_{\pi} V_{\pi}(s), \forall s \in \mathcal{S}$$

optimal policies also share the same optimal action-value function  $\triangleq Q^*$

$$\begin{aligned} \pi_1 &\geq \pi_2 \\ V_{\pi_1}(s) &\geq V_{\pi_2}(s) \\ &\forall s \in \mathcal{S} \end{aligned}$$

For optimal policies  $\pi_1^*$  &  $\pi_2^*$ ,

$$V_{\pi_1^*}(s) = V_{\pi_2^*}(s)$$

$$q_{\pi}(s, a) \triangleq \max_{\pi} q_{\pi}(s, a) \quad \forall s \in \mathcal{S} \text{ \& } \forall a \in \mathcal{A}(s)$$

Q.  $V_{\pi}(s) \stackrel{?}{=} \max_{a \in \mathcal{A}(s)} q_{\pi_{\pi}}(s, a) \text{ or } q_{\pi}(s, a)$

$$V_{\pi}(s) = \sum_{a \in \mathcal{A}(s)} \pi(a|s) q_{\pi}(s, a) \leq \max_{a \in \mathcal{A}(s)} q_{\pi}(s, a)$$

ex<sub>2</sub> -  $c_1 x_1 + c_2 x_2 + c_3 x_3 \leq \max \{x_1, x_2, x_3\} \quad (V_{\pi}(s))$

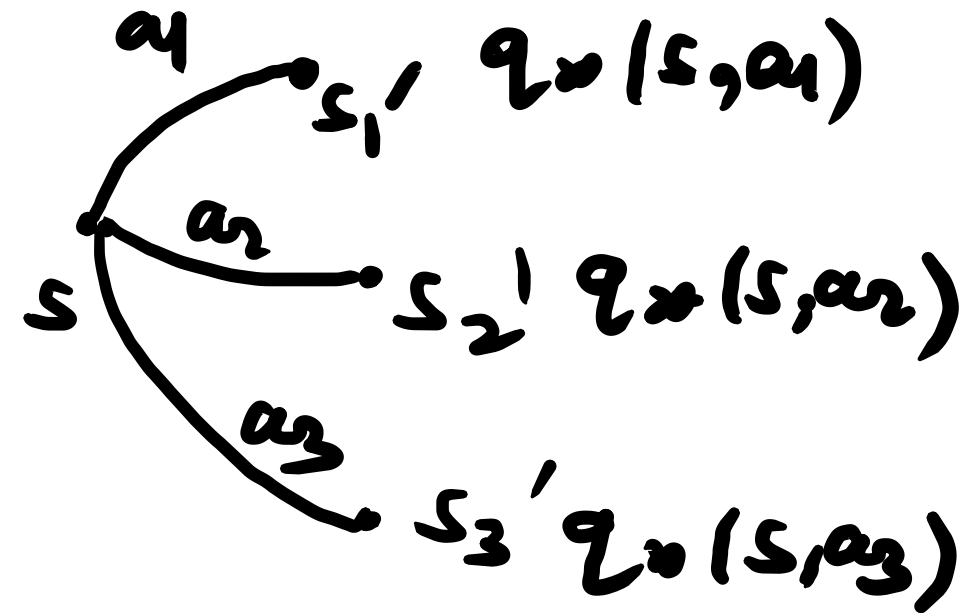
where  $c_1, c_2, c_3 \geq 0$

&  $\sum_{i=1}^3 c_i = 1$

$$0.2(2) + 0.3(1) + 0.5(3) <$$

$$0.2(3) + 0.3(3) + 0.5(3) = 3$$

ex.



$$q_{\pi}(s, a) = E[G_t | S_t = s, A_t = a]$$

$$= E[R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a] - \textcircled{5}$$

$$\stackrel{?}{=} E[R_{t+1} + \gamma V_{\pi}(S_{t+1}) | S_t = s, A_t = a]$$

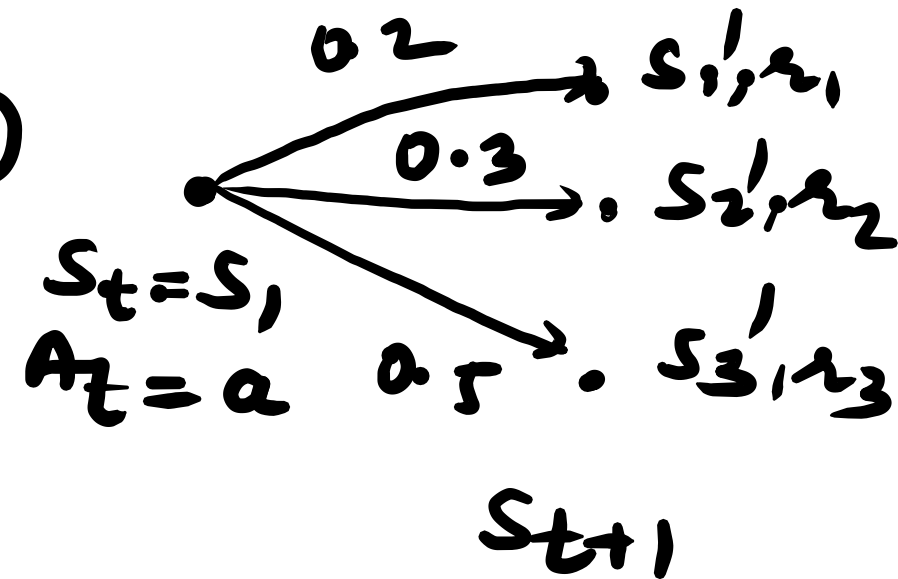
$$\gamma E_{\pi}[G_{t+1} | S_t = s, A_t = a]$$

$$\stackrel{(2)}{=} \gamma \sum_{S_{t+1}} V_{\pi}(S_{t+1}) P(S_{t+1} | S_t = s, A_t = a)$$

$$\stackrel{(1)}{=} \gamma \sum_{S_{t+1}} E_{\pi}[G_{t+1} | S_t = s, A_t = a, S_{t+1}] P[S_{t+1} | S_t = s, A_t = a]$$

$$(3) \quad \gamma \sum_{S_{t+1}} \sum_{\mathbf{r}} V_{\pi}(S_{t+1}) P(S_{t+1}, \mathbf{r} | S_t = s, A_t = a)$$

$$V_{\pi}(s) = E[G_t | S_t = s]$$



$$p(s', \mathbf{r} | s, a)$$

$$E_{\pi^*} [R_{t+1} | S_t = s, A_t = a] = \sum r P(r | s, a)$$

$$= \sum_r \sum_{s'} r P(s', r | s, a)$$

Going back to  
eq. (5)

$$E_{\pi^*} [R_{t+1} + \gamma V_{t+1}(S_{t+1}) | S_t = s, A_t = a]$$

$$= \sum_{S_{t+1}} \sum_r (r + \gamma V_{t+1}(S_{t+1})) P(S_{t+1}, r | s, a)$$

$$E[f(x) | y]$$

$$\int_x f(x) P(x | y)$$