

$E[X|Y] \stackrel{?}{=} f(Y)$; $E_2[E_1[X|Y]]$, the first expectation is
 over x , for a given y . the
 $\int x p(x|y) dx$
 or $\sum x p(x|y)$
 Second one is over y .

Substituting $E_{\pi}[R_{t+2} + \gamma V_{\pi}(S_{t+2}) | S_{t+1}, A_{t+1} = \pi'(S_{t+1})]$

as $E_{\pi}[R_{t+2} + \gamma V_{\pi}(S_{t+2}) | S_{t+1}] \rightarrow$ follows from

Going back to the last inequality
 on slide 3 of lec 13, we get

top of slide
 no. 2 of lec 13.

$$V_{\pi}(s) \leq \mathbb{E}_{\pi} [R_{t+1} + \gamma \mathbb{E}_{\pi} [R_{t+2} + \gamma V_{\pi}(s_{t+2}) \mid s_{t+1}] \mid s_t = s]$$

$$= E_{\pi} [R_{t+1} + \gamma R_{t+2} + \gamma^2 \underset{\gamma}{V_{\pi}(s_{t+2})} \mid s_t = s]$$

To show $\gamma E_{\pi'} \left[\overbrace{E_{\pi'} [R_{t+2} + \gamma V_{\pi}(s_{t+2}) | s_{t+1}]}^x | s_t = s \right] - \textcircled{1}$

$$= \gamma E_{\pi} [R_{t+2} + \gamma V_{\pi}(S_{t+2}) | S_t = s] \quad - (2)$$

$$\textcircled{1} \rightarrow E_{\pi} \left\{ \left[\sum_{r', s''} \gamma p(r', s'' | s') \right] \middle| s_t = s \right\} = \sum_{s'} \sum_{r', s''} \gamma p(r', s'' | s') \cdot p(s' | s)$$

$\rightarrow \sum_{r', s''} \gamma p(r', s'' | s)$
 \parallel
 $\frac{\sum_{s'} p(r', s'' | s') p(s' | s)}{p(s' | s)}$

$$p(r', s'' | s', s) = p(r', s'' | s') \because \text{MDP}$$

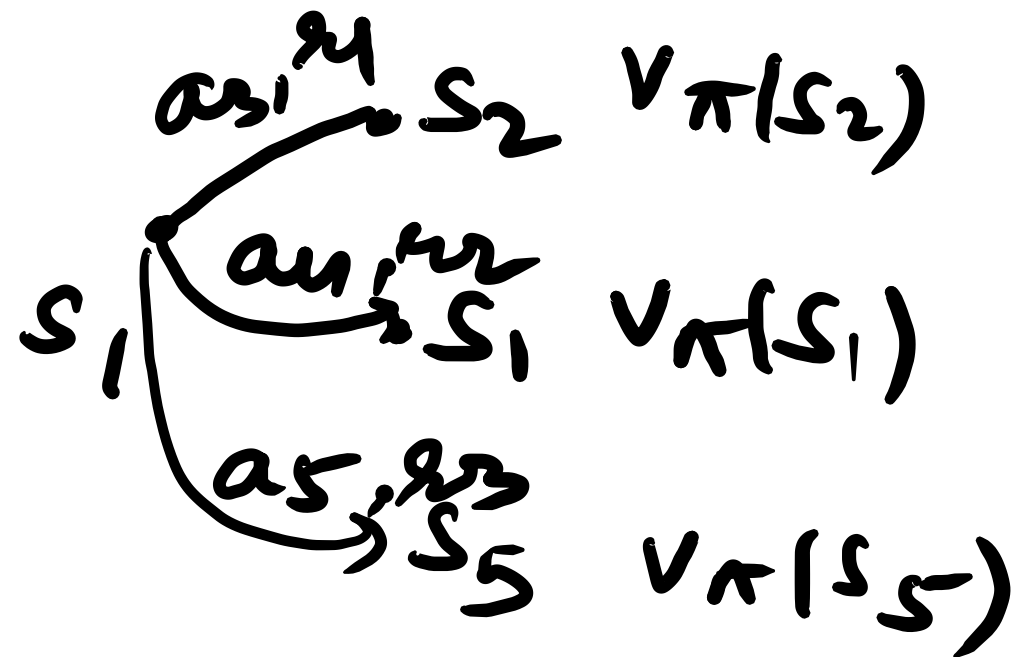
In all states, given $V_\pi(s)$, i.e., expected return following policy π .

$$q_\pi(s, a) \quad \forall a \in A(s)$$

& over all s .

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

$s \backslash a$	a_1	a_2	a_3	a_4	a_5
s_1	-	-	✓	✓	✓
\vdots					
s_n	✓	✓	✓	-	-



So far, we have seen how given a policy & its value function we can easily evaluate a change in the policy at a single state to a particular action.

We can apply changes at all states & to all possible actions, selecting at each state the action that appears best according to $q_{\pi}(s, a)$.

Consider the new 'greedy' policy π' , given

$$\text{by } \pi'(s) \triangleq \arg \max_a q_{\pi}(s, a)$$

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

$$A = \{1, -1, 3, 9\}$$

$$\arg \max A = 4$$

$$\max A = 9$$

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

The process of making a new policy that improve on an original policy, by making it greedy w.r.t the value function of the original policy is called policy improvement (PI).

Suppose the new greedy policy π' is as good as, but not better than the old policy π . Then

$$V_{\pi} = V_{\pi'} \quad \forall s \in \mathcal{S}^+$$