$$\pi'(s) \doteq \arg\max_a q_\pi(s, a)$$

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

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

$$V_\pi(s) = V_{\pi'}(s) \quad \forall s \in S \text{ or } S^+ \quad \Rightarrow$$

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

$$= \arg\max_a q_{\pi'}(s, a)$$

$$V_{\pi'}(s) = q_{\pi'}(s, \pi'(s)) = \max_a q_{\pi'}(s, a)$$

$\Rightarrow$ $V_{\pi'}(s) = \max\limits_{a} E[R_{t+1} + \gamma V_{\pi'}(S_{t+1}) | S_t = s, A_t = a]$

this is BOE, implying $\pi'(\cdot)$ as one of the optimal policies & $V_{\pi'}(s)$ as the optimal state-value function

$\hookrightarrow$ $V_*(s)$

$\pi_* = \{\pi_1, \pi_2\}$

$\hookrightarrow$ set of optimal policies

Value iteration:- ( PE $\to$ PI $\to$ PI$_{t_2}$ )

$\hookrightarrow$ prediction

Policy itr involves policy evaluation which itself is iterative & requires multiple sweeps through the state space.

Is it possible to truncate PE?
without losing the convergence guarantee of Policy itr

what about if PE is stopped after just one sweep?

(one update of each state)

This is called as Value iteration (V Itr) → combination of PI & truncated PE steps

$$V_{k+1}(s) \triangleq \max_a E\left[R_{t+1} + \gamma V_k(S_{t+1}) \mid S_t = s, A_t = a\right]$$

$$\forall s \in S$$

---

$$\pi_0 \xrightarrow{E} V_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} V_{\pi_1} \cdots\cdots \pi_* \to V_{\pi_*}$$

V Itr

optimal policy $\pi_*$

for any policy $\pi$, $V_\pi(s)$ is called the state-value func. if it satisfies BE.

BE

BOE