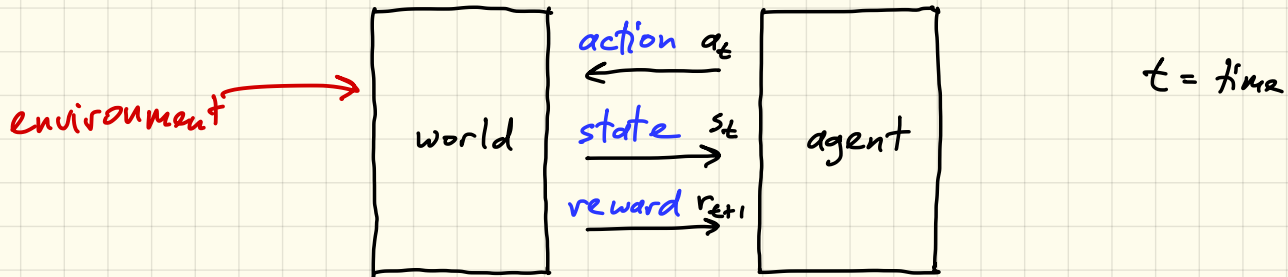


MDP: Markov Decision Process (SAB ch 3)

- stochastic, discrete state, discrete action, state feedback



3 system spaces: state $s_t \in S$, action $a_t \in A$, reward $r_t \in \mathbb{R}$

3 system functions: model $p(s_{t+1} | s_t, a_t) = \text{probability of } s_{t+1} \text{ given } s_t, a_t$

policy $\pi(a_t | s_t) = \text{probability of } a_t \text{ given } s_t$

reward $r_{t+1} = R(s_t, a_t) \leftarrow \text{not stochastic (if so, } E[R])$

reinforcement learning : $a_t = \pi_{\theta}(s_t)$ $(s_{t+1}, r_t) = p(s_t, a_t)$

policy (arrow from π_{θ} to a_t)

recursive (arrow from p to s_{t+1})

evaluation (arrow from p to r_t)

$$\underset{\theta}{\text{maximize}} \sum_t r_t$$

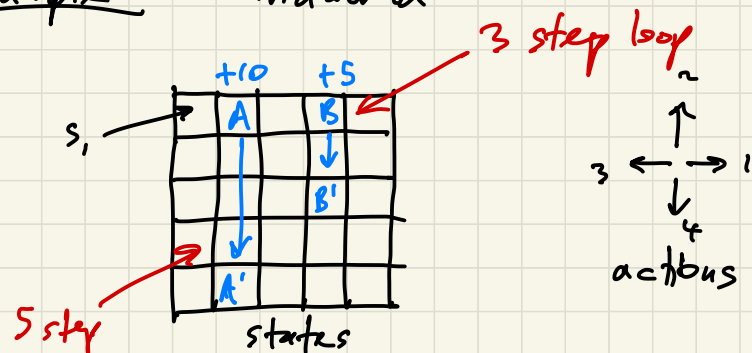
model-based RL : algorithms can directly call $p()$

model-free RL : only have access to $p()$ via trajectories

↳ pros : - can do it on reality
- only need simulator

can we
call $\nabla p()$ (arrow from $\nabla p()$ to $p()$)

Example : Gridworld



5 step loop

$s \in \{1, \dots, 25\}$

(0 to 24 in Python)

$a \in \{1, \dots, 4\}$

reward: $r_t = -1$ except
at A, B

If the system is time-dependent, we can add time as a state.

Value functions

$$V_{\pi}(s) = E \left[\sum_{t=0}^{\infty} \gamma^t r_{t+1} \mid \pi, s_0 = s \right]$$

$$Q_{\pi}(s, a) = E \left[\sum_{t=0}^{\infty} \gamma^t r_{t+1} \mid \pi, s_0 = s, A_0 = a \right]$$

$$V_{\pi}(s) = E_{a \sim \pi(\cdot|s)} [Q_{\pi}(s, a)] = \sum_a \pi(a|s) Q_{\pi}(s, a)$$

$$Q_{\pi}(s, a) = R(s, a) + \gamma E_{s' \sim p(\cdot|s, a)} [V_{\pi}(s')]$$

$$V_{\pi}(s) = E_{a \sim \pi(\cdot|s)} [R(s, a) + \gamma E_{s' \sim p(\cdot|s, a)} [V_{\pi}(s')]]$$

↖ $\pi(s)$ if π is deterministic

$$V^*(s) = \max_{\pi} V_{\pi}(s)$$

$$Q^*(s, a) = \max_{\pi} Q_{\pi}(s, a)$$

Deterministic policy

$$a = \pi(s)$$

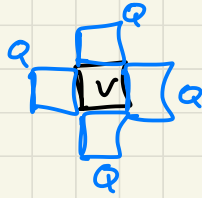
$$\pi(a|s) = \begin{cases} 1 & \text{if } a = a' \\ 0 & \text{otherwise} \end{cases}$$

$$V_{\pi}(s) = Q(s, \pi(s))$$

↙

$a \propto$

$$V^*(s) = \max_a Q^*(s, a)$$



$$V^*(s) = \max_a \left(R(s, a) + \gamma E_{s' \sim p(\cdot | s, a)} [V^*(s')] \right)$$

Bellman
optimality
equation

$$Q^*(s, a) = R(s, a) + \gamma E_{s' \sim p(\cdot | s, a)} \left[\max_{a'} Q^*(s', a') \right]$$

Policy evaluation (S.B 4.1) : prediction $\pi \rightarrow V$ (know p)

$$V_{\pi}(s) = E_{a \sim \pi(\cdot|s)} [R(s, a) + \gamma E_{s' \sim p(\cdot|s, a)} [V_{\pi}(s')]]$$

$$= \sum_a \pi(a|s) \left[R(s, a) + \gamma \sum_{s'} p(s'|a, s) V_{\pi}(s') \right]$$

$s \in \{1, \dots, 25\}$

$a \in \{1, \dots, 4\}$

$$\pi = \begin{bmatrix} \times & \times & \times & \times \\ \pi(16, 2) \\ \bullet \\ \vdots \end{bmatrix}_{25 \times 4}$$

$$V = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}_{25 \times 1}$$

$$V_{\pi}(s) - \gamma \sum_{s'} \underbrace{\left(\sum_a \pi(a|s) p(s'|s, a) \right)}_{A(s, s') \leftarrow 25 \times 25} V_{\pi}(s') = \underbrace{\sum_a \pi(a|s) R(s, a)}_{b(s) \leftarrow 25 \times 1}$$

$$\underbrace{(I - \gamma A)}_{25 \times 25} V = \underbrace{b}_{25 \times 1}$$

$$(I - A)x = b$$

$$\underbrace{[1 \dots 1]}_{x_i} [x] - [A] [x] = [b] \leftarrow \text{ith}$$

$$x_i - \sum_j A_{ij} x_j = b_i$$

$$V(s) - \int A(s, s') V(s') = b(s)$$

$$y = Ax$$

Iterative solution: $v_k(s) \mapsto v_{k+1}(s)$
 \uparrow iteration number

$$v_{k+1}(s) = \mathbb{E}_{a \sim \pi(\cdot|s)} \left[R(s,a) + \gamma \mathbb{E}_{s' \sim p(\cdot|a,s)} [v_k(s')] \right] \quad \forall s$$

$$v_{k+1} = b + \gamma A v_k \quad \left(\text{same as } (I - \gamma A)v = b \right)$$

Can prove this converges if $\gamma < 1$

Policy improvement (S & B 4.2) $\pi \rightarrow \pi'$

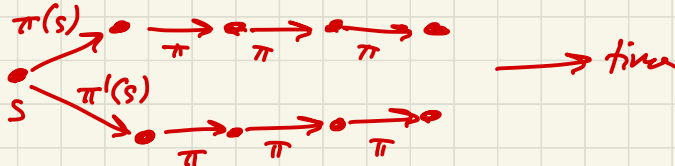
Policy improvement theorem Given π, π' deterministic policies

If $Q_\pi(s, \pi'(s)) \geq V_\pi(s) \quad \forall s$ Better 1-step decisions

then $V_{\pi'}(s) \geq V_\pi(s) \quad \forall s$

⇒ better overall performance
 $Q_\pi(s, \pi(s)) = V_\pi(s)$

Recall



Pf Unroll:

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

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

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

$$\leq E_{\pi'}[R_{t+1} + \gamma Q_\pi(s_{t+1}, \pi'(s_{t+1})) \mid S_t = s]$$

$$= E_{\pi'}[R_{t+1} + \gamma E[R_{t+2} + \gamma V_\pi(s_{t+2}) \mid A_{t+1} = \pi'(s_{t+1})] \mid S_t = s]$$

$$= E_{\pi'}[R_{t+1} + \gamma R_{t+2} + \gamma^2 V_\pi(s_{t+2}) \mid S_t = s] = \dots = V_{\pi'}(s)$$

Greedy policy: $\pi'(s) = \underset{a}{\operatorname{argmax}} Q_{\pi}(s, a)$

$$= \underset{a}{\operatorname{argmax}} E_{s' \sim p(\cdot | s, a)} [R(s, a) + \gamma V_{\pi}(s')]$$

Observe that $Q_{\pi}(s, \pi'(s)) \geq Q_{\pi}(s, a) \quad \forall a$

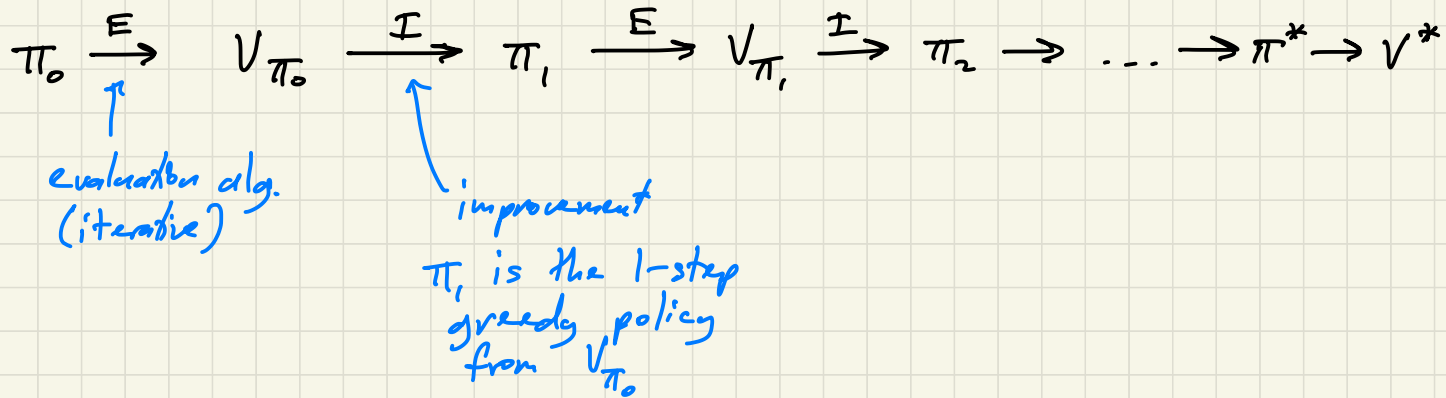
$$\Rightarrow Q_{\pi}(s, \pi'(s)) \geq V_{\pi}(s) \quad (\text{by then } \Rightarrow V_{\pi'}(s) \geq V_{\pi}(s) \text{ for all } s)$$

What if $\pi'(s) = \pi(s) \quad \forall s$?

$$\Rightarrow V_{\pi'}(s) = \max_a E_{s' \sim p(\cdot | s, a)} [R(s, a) + \gamma V_{\pi}(s')] \quad \leftarrow \text{Bellman}$$

Thus if greedy policy makes no improvement, π must already be optimal.

Policy iteration ($S \in \mathcal{B}$ 4.3)



Q: how many iterations of evaluation should we do before improving the policy?

Value iteration (S & B 4.4)

Key idea: do policy improvement but with a single evaluation sweep

$$\begin{aligned} V_{k+1}(s) &= \max_a \left(R(s, a) + \gamma \mathbb{E}_{s' \sim p(\cdot | s, a)} [V_k(s')] \right) \quad \forall s \\ &= \max_a \sum_{s'} \left(R(s, a) + \gamma p(s' | s, a) V_k(s') \right) \end{aligned}$$