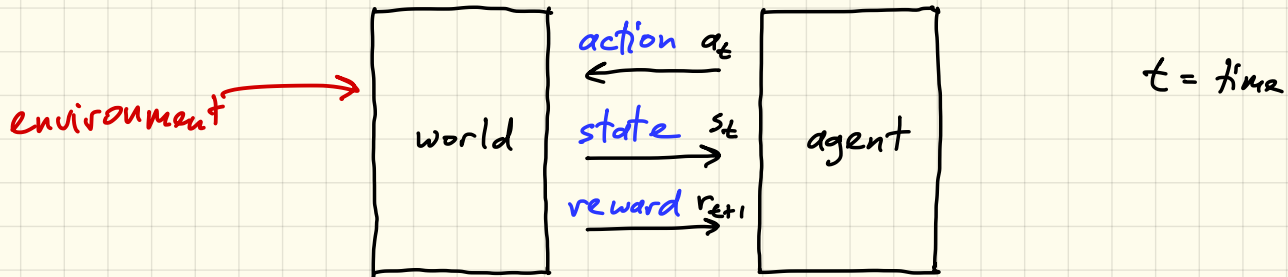# MDP: Markov Decision Process   (S&B Ch 3)

- stochastic, discrete state, discrete action, state feedback



$t$ = time

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

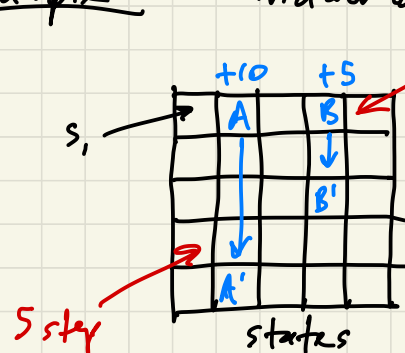3 system functions:    model   $p(s_{t+1} \mid s_t, a_t)$ = probability of $s_{t+1}$ given $s_t, a_t$

policy   $\pi(a_t \mid s_t)$ = probability of $a_t$ given $s_t$

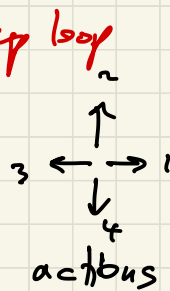reward   $r_{t+1} = R(s_t, a_t)$ ← not stochastic (if so, take E[·])

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

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

# Example : Gridworld



**3 step loop**

+10 +5

A B
B'
A'

**5 step loop**

$s_1$

states

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

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

actions

$$\begin{array}{ccc} s_{t+1} & s_t & a_t \end{array}$$

$$p = 25 \times 25 \times 4$$

$$R = 25 \times 4$$

$$s_{t+1} \sim p(\cdot \mid s_t, a_t)$$

$$p[:, s_t, a_t] = [0.1 \ 0.3 \ \ldots]_{25}$$

$$V = \begin{bmatrix} \\ \\ \\ \\ \end{bmatrix} \begin{matrix} 25 \\ \\ \\ \\ 1 \end{matrix}$$

$$Q = \begin{bmatrix} \\ \\ \\ \\ \end{bmatrix} \begin{matrix} 25 \\ \\ \\ 4 \end{matrix}$$

$$\pi = \begin{bmatrix} & \\ & \end{bmatrix} 4 \\ 25$$

# Policy evaluation (S&B 4.1) : prediction $\quad \pi \longrightarrow V \quad$ (know $p$)

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

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

$(I - \gamma A) V = b$

— model-based

Iterative solution :

$$v_{k+1}(s) = E_{a \sim \pi(\cdot | s)} \left[ R(s,a) + \gamma E_{s' \sim p(\cdot | a,s)} \left[ v_k(s') \right] \right] \qquad \forall_s$$

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

## Policy improvement theorem  Given $\pi, \pi'$ deterministic policies

If $Q_\pi(s, \pi'(s)) \geq V_\pi(s)$ $\forall s$ <span style="color:red">one step improvement</span>

then $V_{\pi'}(s) \geq V_\pi(s)$ $\forall s$ <span style="color:red">$\Longrightarrow$ all-steps improvement</span>

## Greedy policy:  $\pi'(s) = \arg\max_a Q_\pi(s, a)$

## Policy iteration (S&B 4.3)

$$\pi_0 \xrightarrow{E} V_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} V_{\pi_1} \xrightarrow{I} \pi_2 \to \dots \to \pi^* \to V^*$$

## Value iteration (S&B 4.4)

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

<u>Model-free methods</u> → only access trajectories

<u>Temporal-difference</u> TD(0)  (Sa B 6.1)    $\pi \to V$

$$(s_0, a_0, r_1, s_1, a_1, r_2, s_2, a_2, \ldots \ldots) \leftarrow \text{trajectory}$$

$$(s_t, a_t, r_{t+1}, s_{t+1}) \leftarrow \text{piece of trajectory}$$

$$V_\pi(s_t) = E_{a_t \sim \pi(\cdot | s_t)} E_{s_{t+1} \sim p(\cdot | s_t, a_t)} \left[ r_{t+1} + \gamma V_\pi(s_{t+1}) \right]$$

model-based method, $E \to \sum_{a_t}, \sum_{s_{t+1}}$

model-free method, sample $a_t, s_{t+1}$

$$V_\pi(s_t) = r_{t+1} + \gamma V_\pi(s_{t+1}) \quad \text{where} \quad (s_t, a_t, r_{t+1}, s_{t+1}) \text{ is a sample}$$

doesn't work, because one sample is not enough

$$V_\pi(s_t) \sim \underbrace{r_{t+1} + \gamma V_\pi(s_{t+1})}_{\text{target}}$$

<span style="color:red">↰ want to average these</span>

<span style="color:red">— use incremental average</span>

$$V_\pi(s_t) = (1-\alpha) V_\pi(s_t) + \alpha \underbrace{\left( r_{t+1} + \gamma V_\pi(s_{t+1}) \right)}_{\text{target}}$$

$$\alpha = \text{learning rate} \in (0, 1)$$

$$V_\pi(s_t) = V_\pi(s_t) + \alpha \underbrace{\left( r_{t+1} + \gamma V_\pi(s_{t+1}) - V_\pi(s_t) \right)}_{\delta_t = \text{increment}}$$

where $a_t \sim \pi(\cdot \mid s_t)$ and $s_{t+1} \sim p(\cdot \mid s_t, a_t)$
<span style="color:red">↰ sampled from</span>

Q: how do we improve our policy?

Without access to $p$, we can't convert $V \to \pi$

But, we can convert $Q \to \pi$     $\pi(s) = a = \underset{a'}{\arg\max} \; Q(s, a')$

    without knowing $p$

We are going to work with $Q$, not $V$.

# SARSA  $(S\&B\ 6.4)$

Use traj segments $(s_t,\ a_t,\ r_{t+1},\ s_{t+1},\ a_{t+1})$

$\begin{array}{cccccc} S & A & R & S & A \end{array}$

$Q \rightarrow$ better $Q$

$\pi \xcancel{\rightleftharpoons} Q$

$$Q(s_t, a_t) = r_{t+1} + \gamma E_{s_{t+1}\sim p(\cdot | s_t, a_t)}\left[ E_{a_{t+1}\sim \pi(\cdot | s_{t+1})}\left[ Q(s_{t+1}, a_{t+1})\right]\right]$$

simplified:

$$Q(s_t, a_t) = r_{t+1} + \gamma E\left[ Q(s_{t+1}, a_{t+1})\right]$$

$$= \underbrace{E\left[ r_{t+1} + \gamma Q(s_{t+1}, a_{t+1})\right]}_{\text{target}}$$

"on-policy" learning

$$Q(s_t, a_t) \leftarrow (1-\alpha)\, Q(s_t, a_t) + \alpha \underbrace{\left( r_{t+1} + \gamma Q(s_{t+1}, a_{t+1})\right)}_{\text{target}}$$

targets should be sampled from E distribution

$$= Q(s_t, a_t) + \alpha \underbrace{\left( r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)\right)}_{\delta_t = \text{increment}}$$

Given $(s_t, a_t, r_{t+1}, s_{t+1}, a_{t+1})$

Update $\boxed{Q(s_t, a_t)} \leftarrow Q(s_t, a_t) + \alpha \, \delta_t$

$(s_0, a_0, r_1, s_1, a_1)$

$(12, 3, 2, 2, 1)$

$Q = \begin{bmatrix} 5 \\ \cdot \\ ? \\ 4 \end{bmatrix}$ 25

$0 \rightarrow 0.2$

$$\delta_t = r_{t+1} + \gamma \, Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)$$
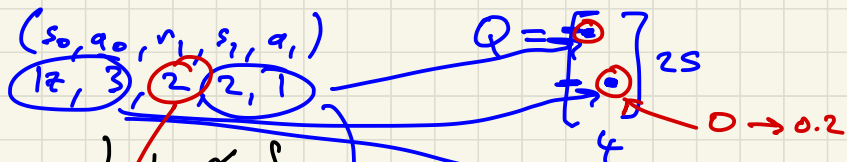
Choose actions:

greedy: $a_t = \arg\max_a Q(s_t, a) \quad \leftarrow$ best choice based on current Q function (no exploration) (all exploitation)

$\varepsilon$-greedy: $a_t = \begin{cases} \text{random } a & \text{with probability } \varepsilon \\ \arg\max_a Q(s_t, a) & \text{otherwise} \end{cases}$

$\varepsilon$-greedy is hopefully close to optimal when Q is close to optimal

exploration/ exploitation tradeoff param.

**Thm** Sarsa converges to the optimal Q if all
(s, a) pairs are visited infinitely often
and the policy converges to greedy. (e.g. $\varepsilon = \frac{1}{t}$)

Sarsa will converge to the Q for whichever policy we use.

If we use $\varepsilon$-greedy $\longrightarrow$ converge to a good Q

This is "on-policy" learning, where we need to choose good
actions.

This means that aggressive exploration hurts learning.

$\longrightarrow$ instead, we can use Q-learning, which
doesn't' have this problem (Q-learning is off-policy)

## Q-learning   (S&B 6.5)

Use   $(s_t, a_t, r_{t+1}, s_{t+1})$

Won't assume that are chosen from $Q$

$$Q(s_t, a_t) = r_{t+1} + E_{s_{t+1} \sim p(\cdot \mid s_t, a_t)} \left[ E_{a_{t+1} \sim \pi(\cdot \mid s_{t+1})} \left[ Q(s_{t+1}, a_{t+1}) \right] \right]$$

next time.