

Lecture 16 - March 14, 2023

- Temporal Difference Learning

- TD(0)
- SARSA
- Q-Learning
- On-Policy Vs. Off-Policy

Project 3 is posted → Due April 14

HW3 → Due March 17

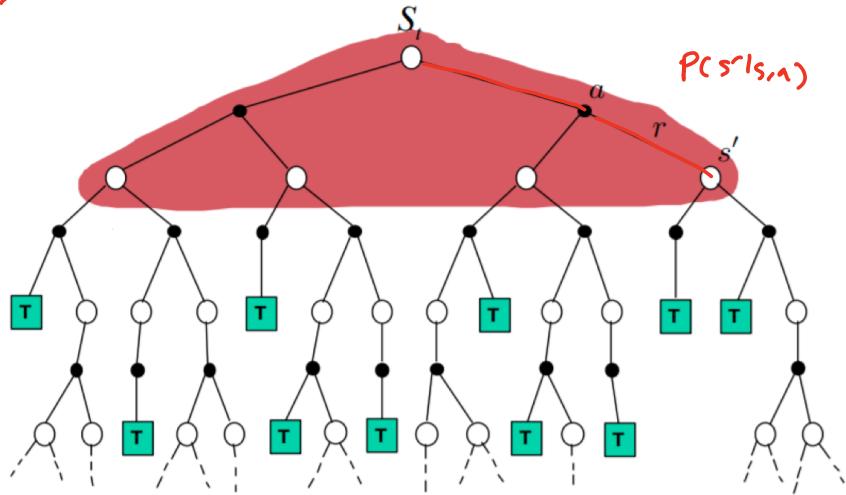
TA's office hour:

Wednesdays, 2pm-3pm (in-person)

Fridays, 2pm-3pm (virtual)

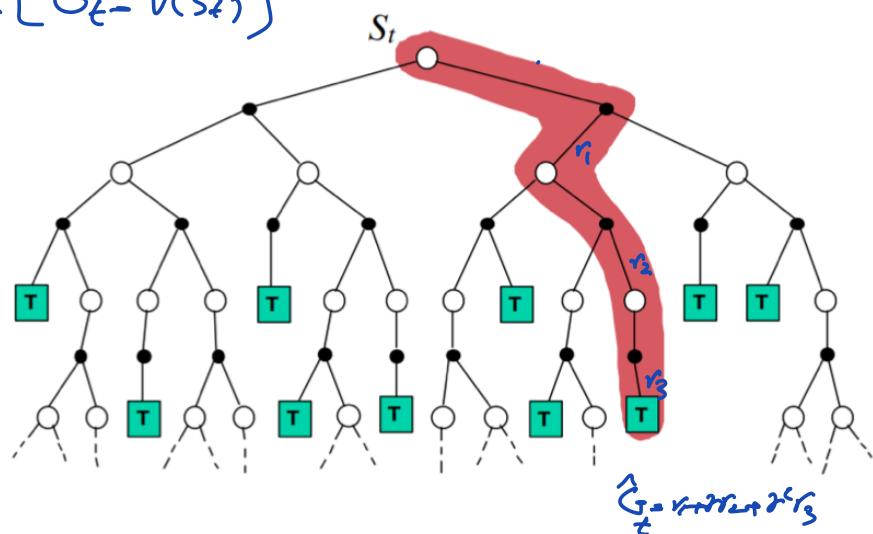
$$DP: V_{k+1}(s) = \max_{a \in A} \sum_{s'} P(s'|s, a) [R + \gamma V_k(s')]$$

Model of system



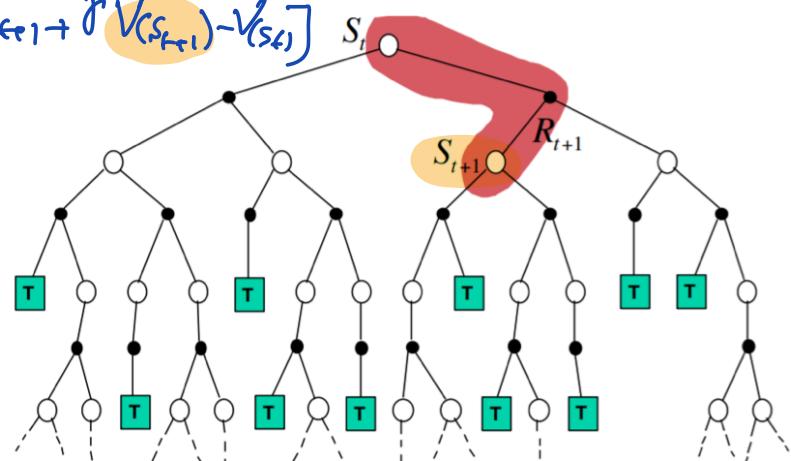
$$MC: V(s_t) = V(s_t) + \alpha [\hat{G}_t - V(s_t)]$$

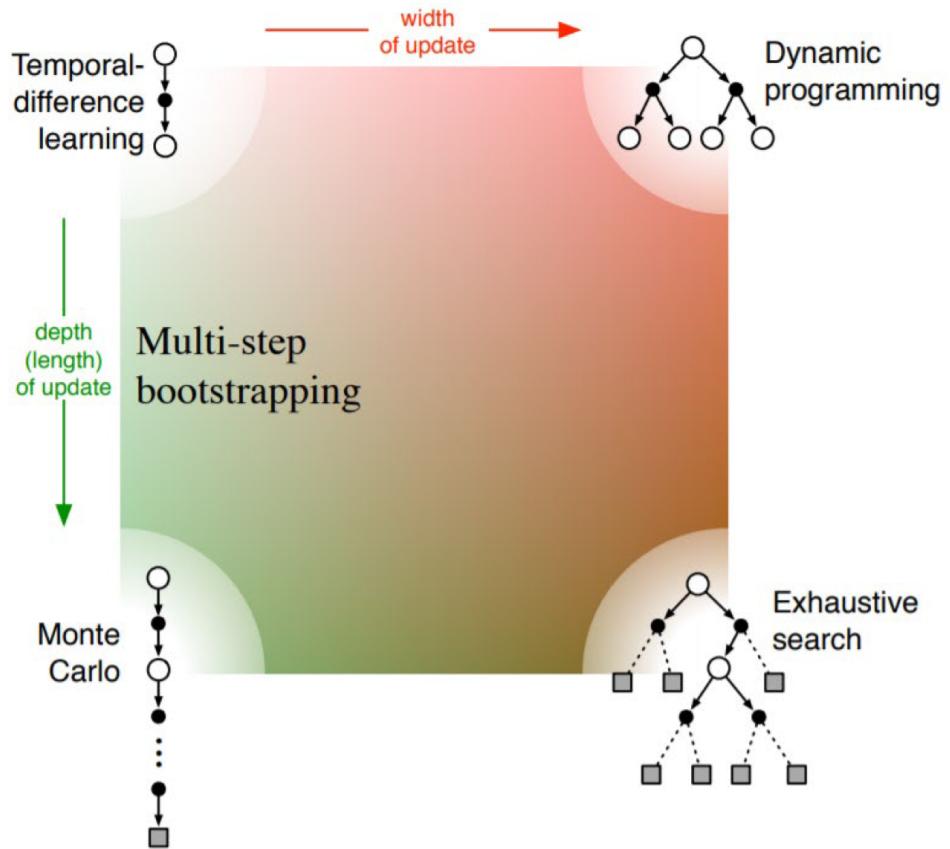
Model ~~X~~
delay ✓
Interactive ~~X~~



$$TD: V(s_t) = V(s_t) + \alpha [R_{t+1} + \gamma V(s_{t+1}) - V(s_t)]$$

Model ~~X~~
Delay one step
Interactive ✓





TD(0)

$\pi \longrightarrow V^\pi$

$$V(s_t) = V(s) + \alpha [R_{t+1} + \gamma V(s_{t+1}) - V(s_t)]$$

$\pi(s_t)$

$$R(s_t, \pi(s_t), s_{t+1})$$

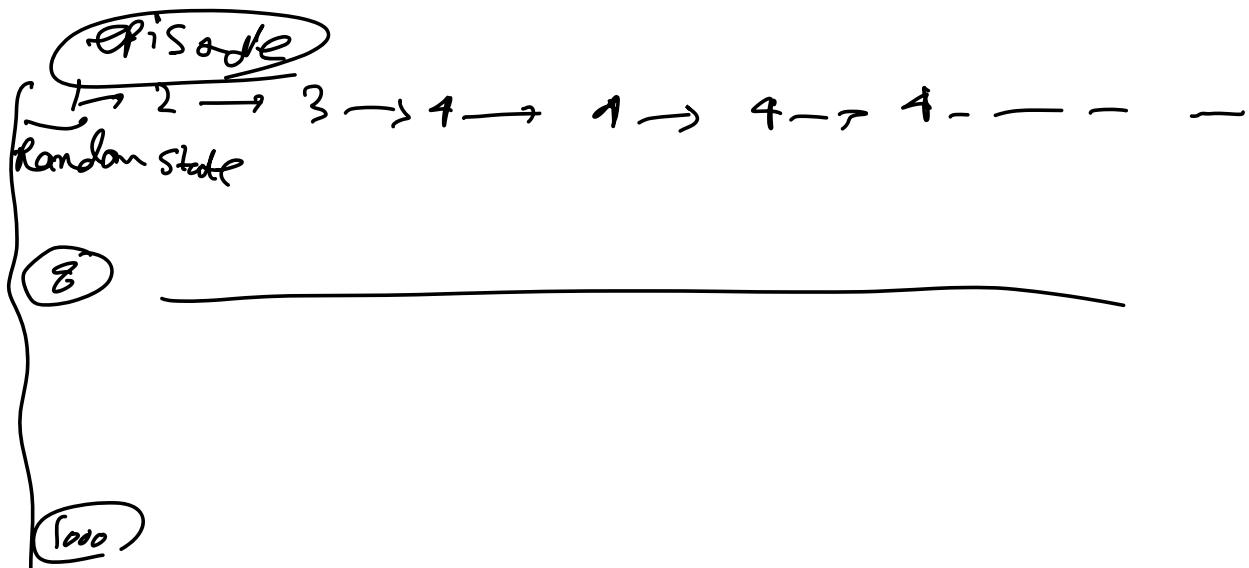
9	10	11	12
8		14	13
7		16	15
6	5		
4	-5.5	-5.5	-0.5

■ Wall
 ■ Bump
 ■ Goal

$\sqrt{\tau}$

$$V(1) = \frac{r(1)}{\gamma} + \alpha \left[-1 + \gamma V(2) - \frac{V(1)}{\gamma} \right] = -0.5$$

$$V(4) = \frac{r(4)}{\gamma} + \alpha \left[-11 + \gamma V(4) - \frac{V(4)}{\gamma} \right] = -9.8$$



$$\pi \rightarrow V^\pi \quad TD(0)$$

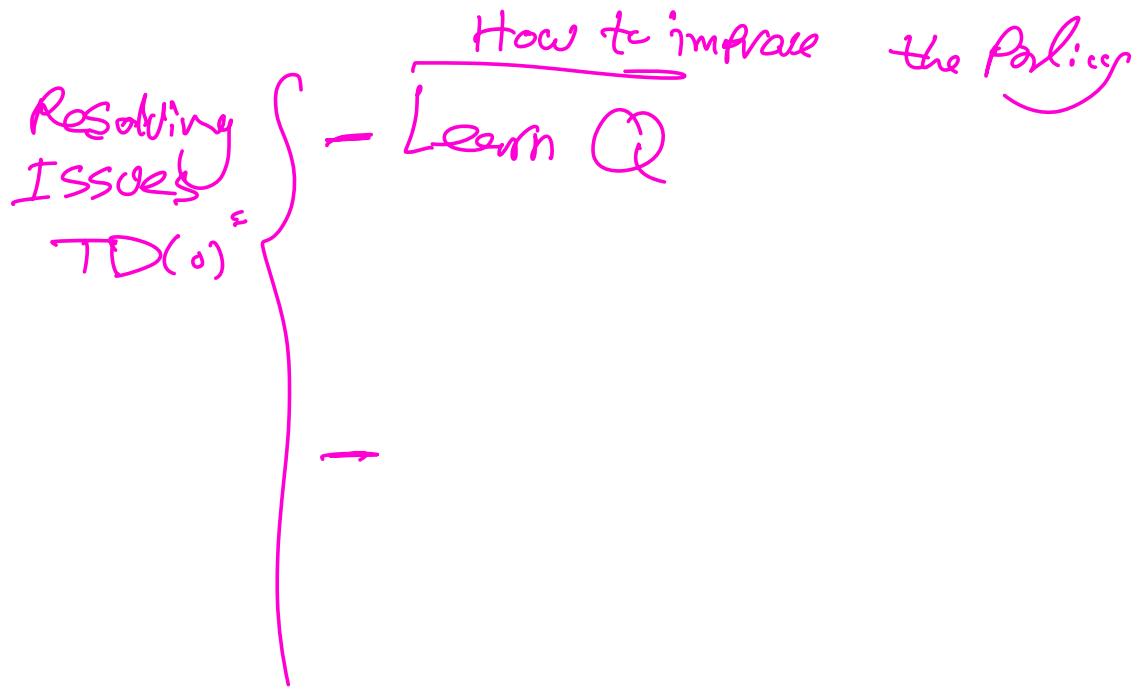
$$\pi'(s) = \operatorname{argmax}_{s'} p(s'|s, a) [R + \gamma V_\pi(s')]$$

Policy Improvement

$$\pi'(s) = \operatorname{argmax}_{a \in A} Q_\pi(s; a)$$

$TD(0) \Rightarrow$ Nat application
for Policy optimization

$$\pi \Rightarrow V_\pi(s_t) = V_\pi(s_t) + \alpha [R_{t+1} + \gamma V_\pi(s_{t+1}) - V_\pi(s_t)]$$



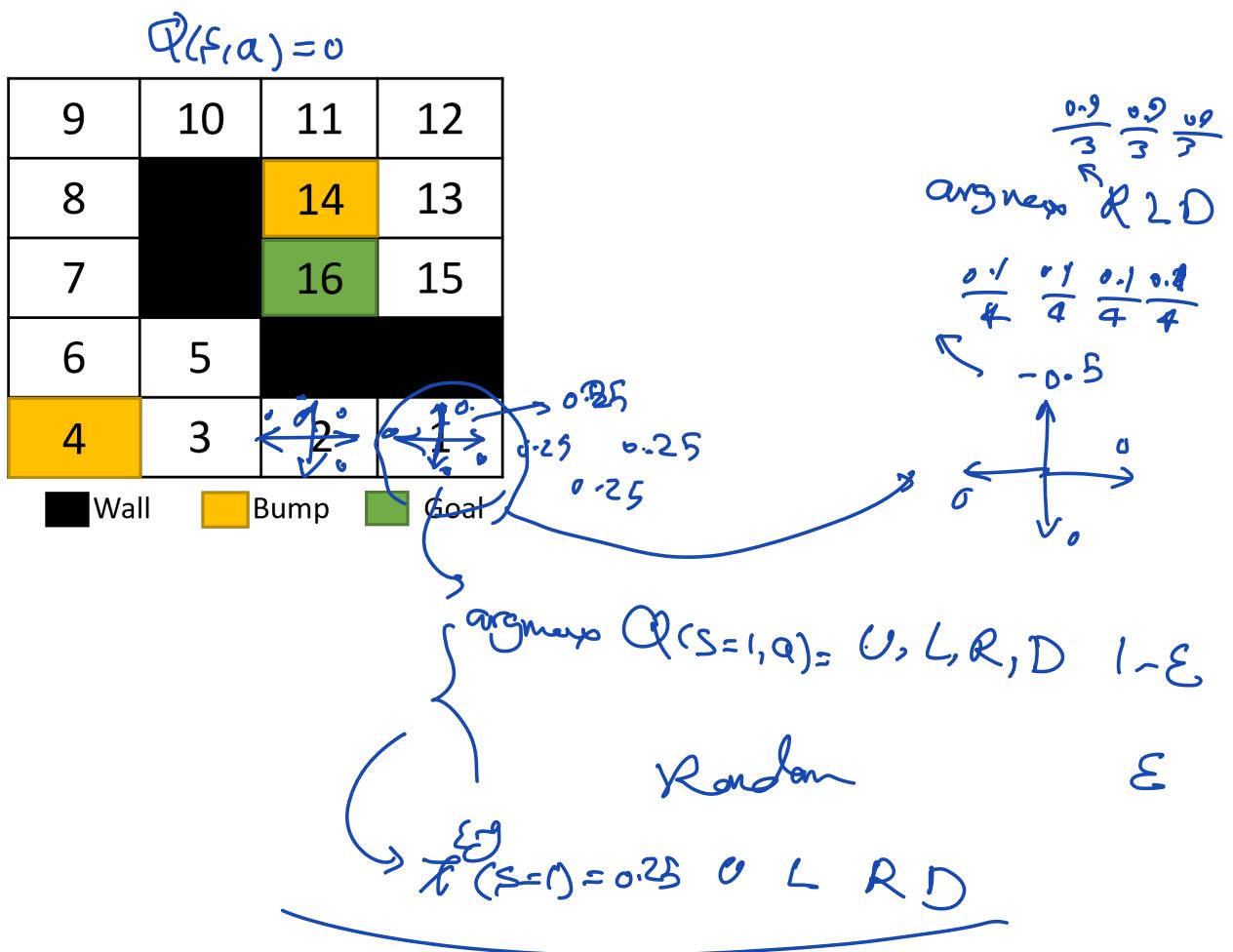
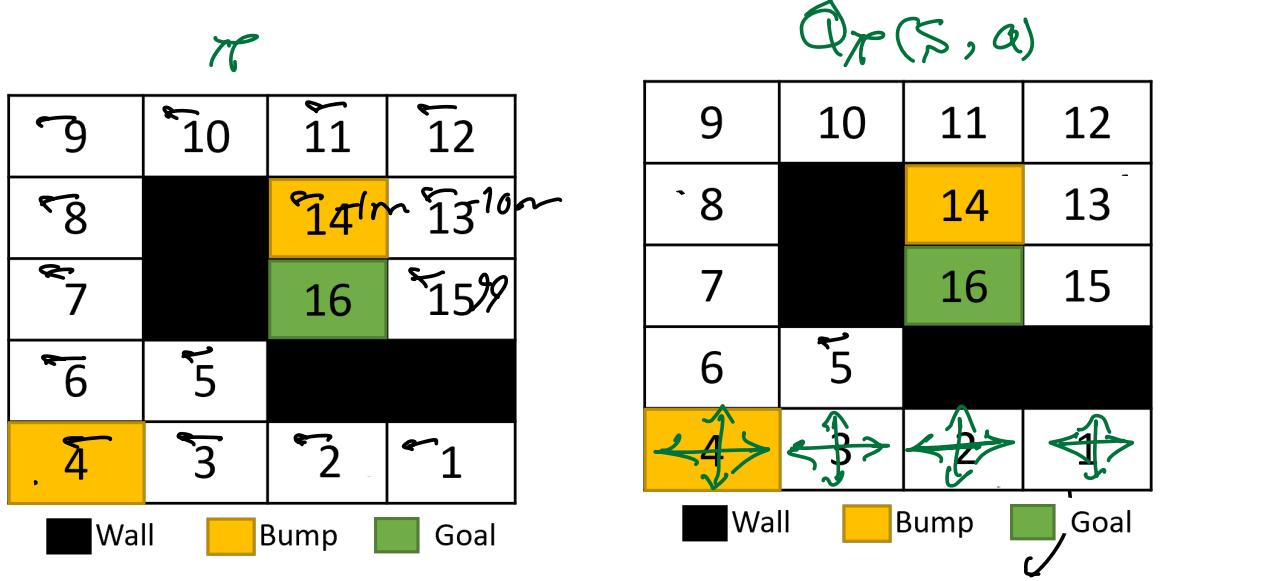
$$V_{\pi}(s_t) = V_{\pi}(s_t) + \alpha [R_{t+1} + \gamma V_{\pi}(s_{t+1}) - V_{\pi}(s_t)] \Rightarrow TD(u)$$

$$Q_{\pi}(s_t, \pi(s_t)) = Q_{\pi}(s_t, \pi(s_t)) + \alpha [R_{t+1} + \gamma Q_{\pi}(s_{t+1}, \pi(s_{t+1})) - Q_{\pi}(s_t, \pi(s_t))]$$

stochastic Policy

$$Q_{\pi}(s_t, \pi(s_t)) = Q_{\pi}(s_t, \pi(s_t)) + \alpha [R_{t+1} + \gamma Q_{\pi}(s_{t+1}, \pi(s_{t+1})) - Q_{\pi}(s_t, \pi(s_t))]$$

$\pi_{\text{Egocentric}}(s_t) = \begin{cases} \text{argmax } Q(s_t, a) & 1-\epsilon \\ \text{Random} & \epsilon \end{cases} - Q_{\pi}(s_t, \pi(s_t))]$



$$Q_{\pi}(s_t, \pi(s_t)) = Q_{\pi}(s_t, \pi(s_t)) + \alpha [R_{t+1} + \gamma Q_{\pi}(s_{t+1}, \pi(s_{t+1}))]$$

stochastic Policy

$$\pi_{(s_t)} = \begin{cases} \text{argmax } Q(s_t, a) & 1-\epsilon \\ \text{Random} & \epsilon \end{cases}$$

$$I \xrightarrow{\pi^{(1)} \sim V} I \xrightarrow{\pi^{(1)} \sim L}$$

$$Q(1, V) = Q(1, U) + \alpha \left[-1 + \gamma Q(s_{t+1}', L) - Q(1, U) \right]$$

$= \approx 0.5$

SARSA

Initialize $Q(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$, arbitrarily, and $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

 Initialize S

 Choose A from S using policy derived from Q (e.g., ε -greedy)

 Repeat (for each step of episode):

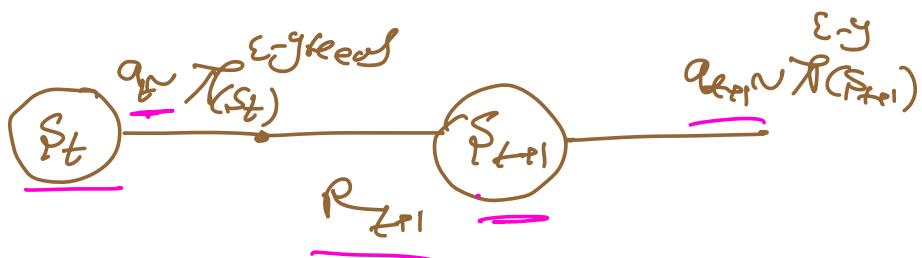
 Take action A , observe R, S'

 Choose A' from S' using policy derived from Q (e.g., ε -greedy)

$$Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma Q(S', A') - Q(S, A)]$$

$$S \leftarrow S'; A \leftarrow A';$$

 until S is terminal



$$Q(s_t, a_t) = Q(s_t, a_t) + \alpha [R_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)]$$

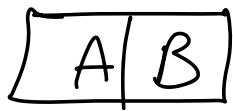
State - Action - Reward - State - Action

(SARSA)

$$\pi(s) = \begin{cases} \underset{a \in A}{\text{argmax}} Q(s, a) & 1 - \varepsilon \\ \text{Random} & \varepsilon \end{cases}$$

Final SARSA: $\pi^*(s) \equiv \underset{a \in A}{\text{argmax}} Q(s, a)$

Example



Intuitively
 $\pi^*(A) = a^2$
 $\pi^*(B) = a^1$

$$M(a^1) = \begin{bmatrix} A & B \\ 1 & 0 \\ B & 0 \\ 0 & 1 \end{bmatrix}$$

$$M(a^2) = \begin{bmatrix} A & B \\ - & 1 \\ B & 1 \\ 1 & - \end{bmatrix}$$

Reward $\rightarrow \begin{cases} 5 & B \\ -1 & a^2 \end{cases}$

SARSA

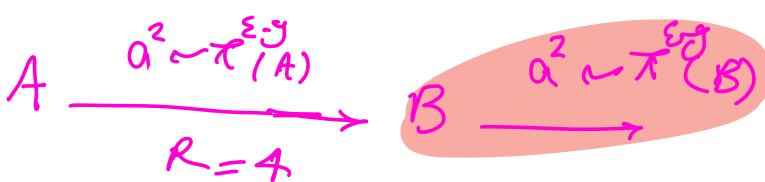
Initial Q

$$\begin{bmatrix} Q(A, a^1) \\ Q(B, a^1) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \begin{bmatrix} Q(A, a^2) \\ Q(B, a^2) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

episode 1

Initial
state
(Random) = A

$$\pi_{\text{eg}}^*(A) = \begin{cases} \underset{a \in A}{\operatorname{argmax}} Q(A, a) & a^1, a^2 \\ \text{Rand} & \sim Q^2 \end{cases}$$



$$\pi_{\text{eg}}^*(B) = \begin{cases} 0.5 & a^1 \sim Q^2 \\ 0.5 & a^2 \end{cases}$$

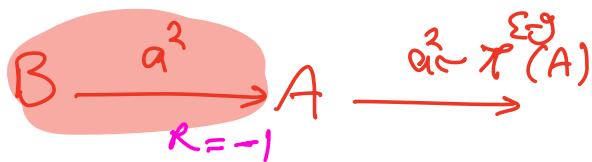
$$Q(A, a^2) = \frac{Q(A, a^2)}{0.5} + \alpha \left[\frac{R}{4} + \gamma \frac{Q(B, a^2)}{0.9} - \frac{Q(A, a^1)}{0} \right] = 2$$

$$\begin{bmatrix} Q(A, a^1) \\ Q(B, a^1) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \begin{bmatrix} Q(A, a^2) \\ Q(B, a^2) \end{bmatrix} = \begin{bmatrix} 2 \\ 0 \end{bmatrix}$$

Closest Deterministic Policy to π^{Greedy}

$$\pi(A) = a^2 \quad \pi(B) = a^1 \text{ or } a^2$$

Previous



$$\pi^{\text{eq}}(A) = \begin{cases} \text{argmax } Q(A, a) = a^2 & 1-\epsilon \\ \text{Random } a^1 \text{ or } a^2 & \epsilon \end{cases} \sim a^2$$

$$Q(B, a^2) = \frac{Q(B, a^2)}{0} + \alpha \left[\frac{R}{-1} + \frac{\gamma}{0.9} \left[Q(A, a^2) - \frac{Q(B, a^2)}{0} \right] \right] = 0.4$$

$$\begin{bmatrix} Q(A, a^1) \\ Q(B, a^1) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \begin{bmatrix} Q(A, a^2) \\ Q(B, a^2) \end{bmatrix} = \begin{bmatrix} 2 \\ 0.4 \end{bmatrix}$$

Closest Deterministic Policy to π^{Greedy}

$$\pi(A) = a^2 \quad \pi(B) = a^2$$

$$A \xrightarrow[\mathcal{R}=4]{a^2} B \xrightarrow{a^1 \sim \pi^{EG}(B)} \left\{ \begin{array}{ll} 0.05 & a^1 \\ 0.95 & a^2 \end{array} \right. \sim a^1$$

$$Q(A, a^2) = \underbrace{Q(A, a^2)}_2 + \alpha \left[\frac{R}{4} + \gamma \underbrace{Q(B, a^1)}_{0.9} - \underbrace{Q(A, a^2)}_2 \right] = 3$$

$$\begin{bmatrix} Q(A, a^1) \\ Q(B, a^1) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \begin{bmatrix} Q(A, a^2) \\ Q(B, a^2) \end{bmatrix} = \begin{bmatrix} 3 \\ 0.4 \end{bmatrix}$$

Closest Deterministic Policy to π^{EG}

$$\pi(A) = a^2 \quad \pi(B) = a^2$$

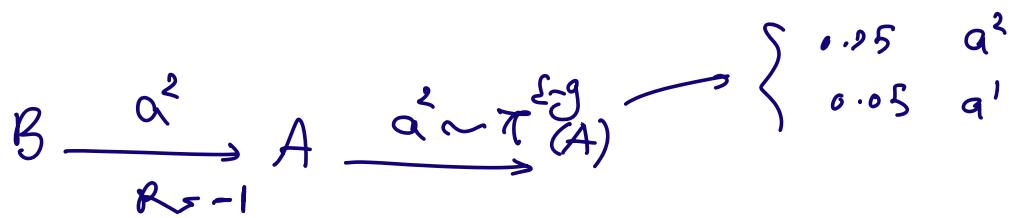
$$B \xrightarrow[\mathcal{R}=5]{a^1} B \xrightarrow{a^2 \sim \pi^{EG}(B)} \pi^{EG}(B) = \begin{cases} a^1 & 0.05 \\ a^2 & 0.95 \end{cases}$$

$$Q(B, a^1) = \underbrace{Q(B, a^1)}_0 + \alpha \left[\frac{R}{5} + \gamma \underbrace{Q(B, a^2)}_{0.9} - \underbrace{Q(B, a^1)}_0 \right] = 2.68$$

$$\begin{bmatrix} Q(A, a^1) \\ Q(B, a^1) \end{bmatrix} = \begin{bmatrix} 0 \\ 2.68 \end{bmatrix} \quad \begin{bmatrix} Q(A, a^2) \\ Q(B, a^2) \end{bmatrix} = \begin{bmatrix} 3 \\ 0.4 \end{bmatrix}$$

Closest Deterministic Policy to π^{target}

$$\pi(A) = a^2 \quad \pi(B) = a^1$$

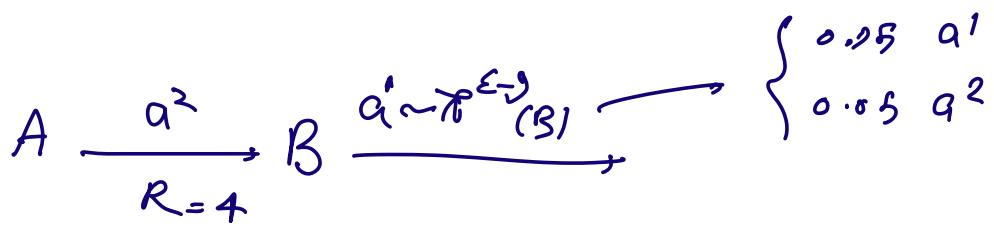


$$Q(B, a^2) = \frac{Q(B, a^2)}{0.4} + \alpha \left[R + \gamma \frac{Q(A, a^2)}{0.5} - \frac{Q(B, a^2)}{0.4} \right] = 1.05$$

$$\begin{bmatrix} Q(A, a^1) \\ Q(B, a^1) \end{bmatrix} = \begin{bmatrix} 0 \\ 2.68 \end{bmatrix} \quad \begin{bmatrix} Q(A, a^2) \\ Q(B, a^2) \end{bmatrix} = \begin{bmatrix} 3 \\ 1.05 \end{bmatrix}$$

Closest Deterministic Policy to π^{target}

$$\pi(A) = a^2 \quad \pi(B) = a^1$$



$$\begin{aligned} Q(A, a^2) &= \frac{Q(A, a^2)}{3} + \alpha \left[\frac{R}{0.5} + \frac{\gamma Q(B, a^1)}{4} - \frac{Q(A, a^2)}{5} \right] \\ &= 4.7 \end{aligned}$$

$$\begin{bmatrix} Q(A, a^1) \\ Q(B, a^1) \end{bmatrix} = \begin{bmatrix} 0 \\ 2.68 \end{bmatrix} \quad \begin{bmatrix} Q(A, a^2) \\ Q(B, a^2) \end{bmatrix} = \begin{bmatrix} 4.7 \\ 1.05 \end{bmatrix}$$

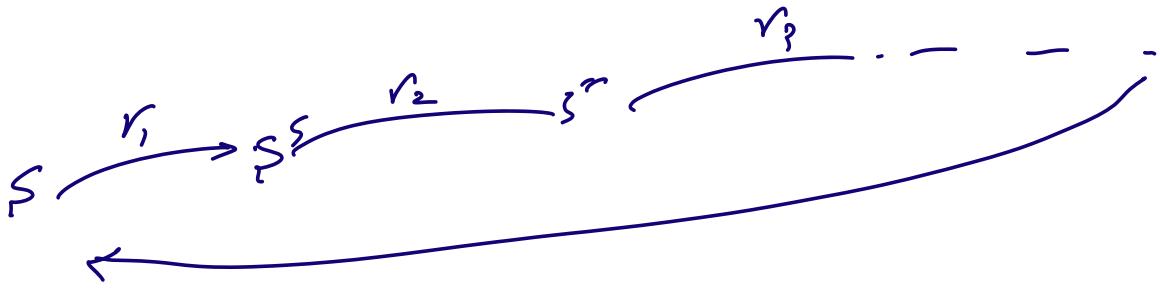
Closest Deterministic Policy to π^{egre}

$$\pi(A) = a^2 \quad \pi(B) = a^1$$

No learning
Execution
 $\xrightarrow{\text{greedy form of } \pi^{\text{egre}}}$

SARSA
 $\pi(s) = \underset{a \in A}{\text{argmax}} Q(s, a)$

SARSA
 $\pi(A) = a^2$
 $\pi(B) = a^1$



$$Q(s, a) = \sum \delta r_c$$