

Lecture 17 - March 17, 2023

- Temporal Difference Learning

- TD(0)

- SARSA

↓ • Q-Learning

- On-Policy vs. Off-Policy

- Expected SARSA

- Double Q-Learning

- Multi-Step Bootstrapping

- SARSA-Lambda

- Actor-Critic Method

HW3 → Due March 17

Project 3 → Due April 14

HW4 → Due March 31

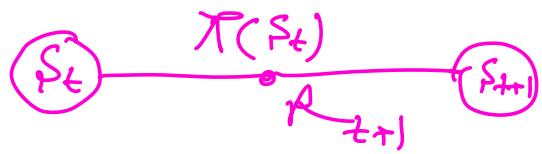
TA's office hour:

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

Fridays, 2pm-3pm (virtual)

Review

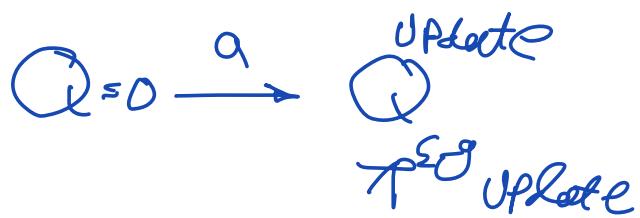
$$TD \Rightarrow V_{\pi}(s_t) = V_{\pi}(s_t) + \alpha \left[\frac{TD \text{ Target}}{R_{t+1} + \gamma V_{\pi}(s_{t+1}) - V_{\pi}(s_t)} \right]$$



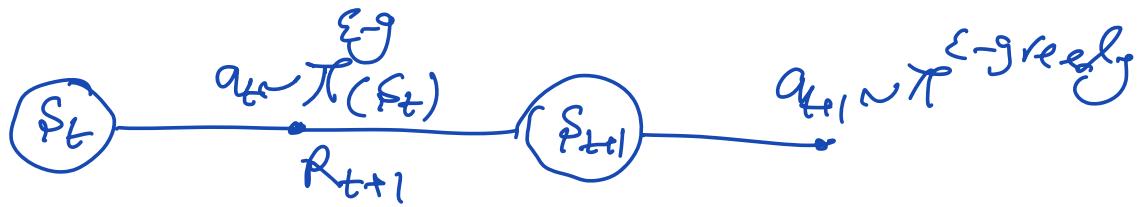
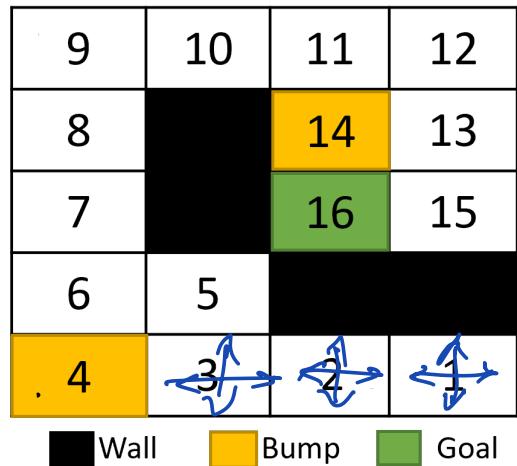
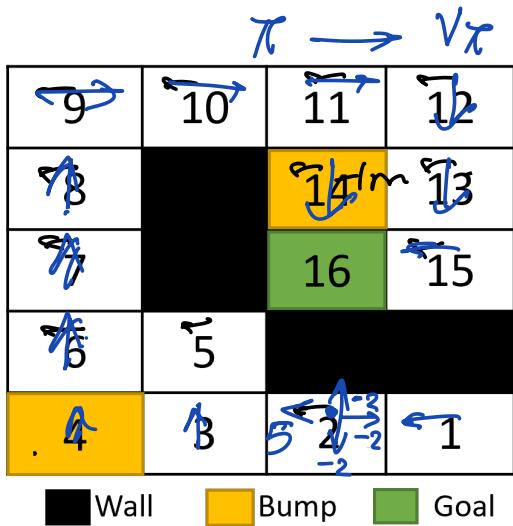
$\pi^{\text{Deterministic}} \Rightarrow \pi^{\text{Stochastic}} := \pi^{\text{argmax}}$

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

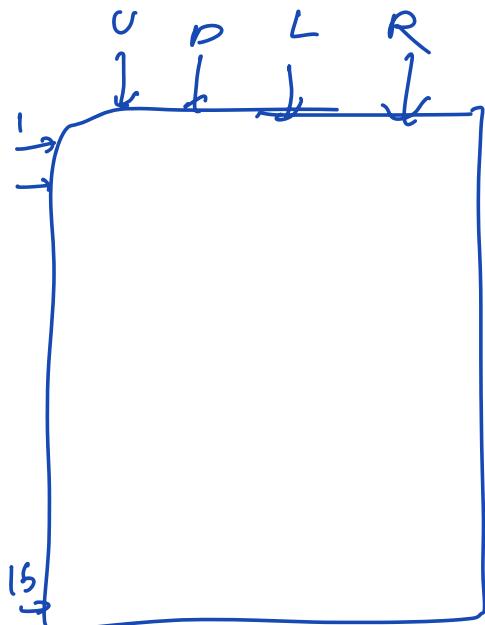
$$\pi^{\text{eg}}(s) = \begin{cases} \underset{a \in A}{\text{argmax}} Q(s, a) & 1 - \epsilon \\ \text{Random}\{a^1, \dots, a^L\} & \epsilon \end{cases}$$



Tabular Learning $\Rightarrow \dot{Q}$



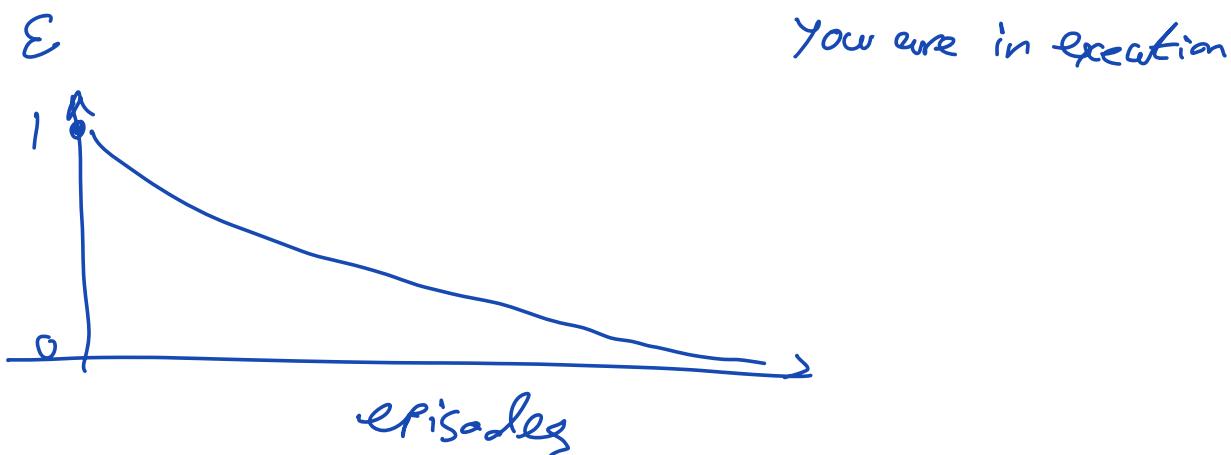
$$Q(s_t, a_t) = Q(s_t, a_t) + \alpha [R_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)]$$



$$\pi^*(s) = \arg \max_{a \in A} Q_{\pi^*}(s, a)$$

The last Q obtained by SARSA is Q^* for the last π^* .

$$\pi^*(s) = \underset{a \in A}{\operatorname{argmax}} Q^{\text{SARSA}}(s, a) \leftarrow \begin{array}{l} \text{You need to follow} \\ \text{this when learning} \\ \text{has finished and} \\ \text{you are in execution} \end{array}$$



$$\epsilon \xrightarrow{\text{Slowly}} 0 \quad \text{SARSA} \Rightarrow Q^{\text{SARSA}} = \pi^*$$

SARSA

$$Q_{\pi}(s, a) = E[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | s_t=s, a_t=a, q_{t+1:\infty} \sim \pi]$$

Bellman $\rightsquigarrow = \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V_{\pi}(s')]$

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

Sample at this SARSA Target

$$\underbrace{Q(s_t, a_t)}_{\text{new}} = \underbrace{Q(s_t, a_t)}_{\text{old}} + \alpha \left[R_{t+1} + \gamma \underbrace{Q(s_{t+1}, a_{t+1})}_{\text{old}} - \underbrace{Q(s_t, a_t)}_{\text{old}} \right]$$

Q-Learning

$$Q_{\pi^*}(s, a) = E[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | s_t=s, a_t=a, q_{t+1:\infty} \sim \pi^*]$$

$$= \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V_{\pi^*}(s')]$$

$$= \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma Q_{\pi^*}(s', \pi^*(s'))]$$

$$= \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma \max_{a'} Q_{\pi^*}(s', a')]$$

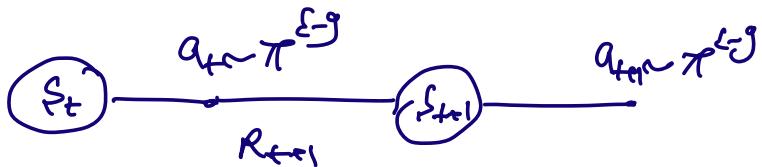
Q-Learning Target

Q-Learning

$$Q(S_t, a_t) = Q(S_t, a_t) + \alpha [R_{t+1} + \gamma \max_{a \in A} Q(S_{t+1}, a) - Q(S_t, a_t)]$$

SARSA

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



Q-Learning

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

 Repeat (for each step of episode):

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

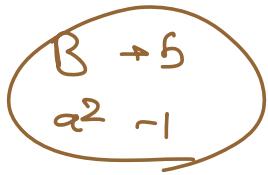
 Take action A , observe R, S'

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

$S \leftarrow S'$;

 until S is terminal

$$\boxed{A \mid B}$$



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

$$M(a^2) = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

Q-learning ↓

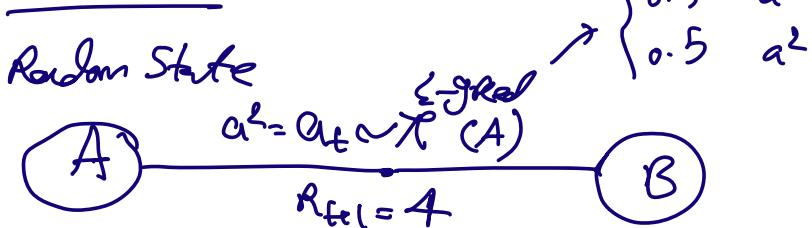
$$\pi^{(0)}(A) = \begin{cases} \text{argmax } Q(A, a) & 1-\epsilon \\ \text{Random} & \epsilon \end{cases}$$

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

Closest Deterministic Policy to $\pi^{(0)}$

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

Episode 1

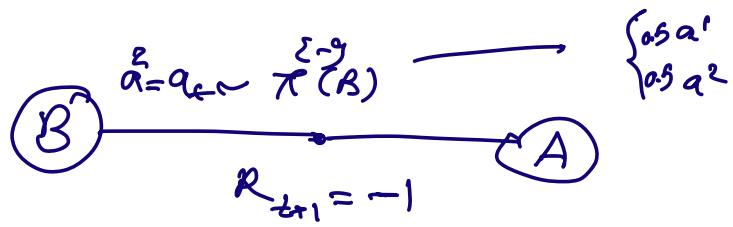


$$Q(A, a^2) = \overbrace{\underbrace{Q(A, a^2)}_{0}} + \alpha \left[\frac{R_{t+1}}{0.5} + \gamma \max_{a \in A} \underbrace{Q(B, a)}_{\overbrace{2}} - \overbrace{Q(A, a^2)}_{0} \right] = 2$$

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

Closest Deterministic Policy to $\pi^{(0)}$

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

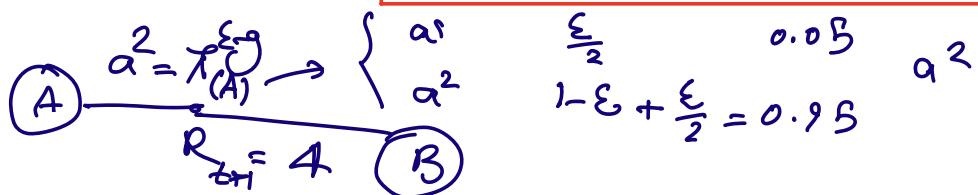


$$Q(B, a^2) = \underbrace{Q(B, a^2)}_0 + \alpha \left[\underbrace{R_{t+1}}_{-1} + \gamma \max_{a \in A} \frac{Q(A, a) - Q(B, a^2)}{2} \right] = 0.4$$

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

Closest Deterministic Policy to $\pi^{\epsilon-g}$

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

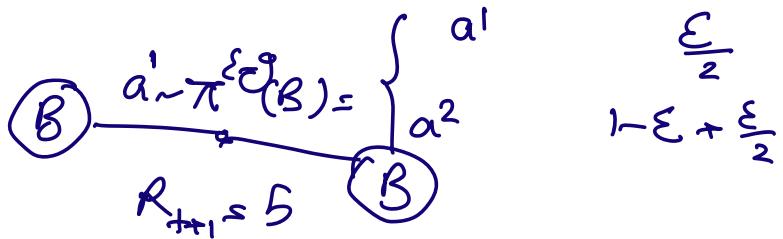


$$Q(A, a^2) = \underbrace{Q(A, a^2)}_{2} + \alpha \left[\underbrace{R_{t+1}}_{4} + \gamma \max_{a \in A} \frac{Q(B, a) - Q(A, a^2)}{0.4} \right] = 3.18$$

$$\begin{bmatrix} Q(A, a^1) \\ Q(B, a^1) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \begin{bmatrix} Q(A, a^2) \\ Q(B, a^2) \end{bmatrix} = \begin{bmatrix} 3.18 \\ 0.4 \end{bmatrix}$$

Closest Deterministic Policy to $\pi^{\epsilon-g}$

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



$$Q(B, a^1) = \underbrace{Q(B, a^1)}_0 + \gamma [R_{t+1} + \delta \max_{a \in A} \underbrace{Q(B, a)}_{0.4} - \underbrace{Q(B, a^1)}_0] = 2.68$$

$$\begin{bmatrix} Q(A, a^1) \\ Q(B, a^1) \end{bmatrix} = \begin{bmatrix} 0 \\ 2.68 \end{bmatrix} \quad \begin{bmatrix} Q(A, a^2) \\ Q(B, a^2) \end{bmatrix} = \begin{bmatrix} 3.18 \\ 0.4 \end{bmatrix}$$

Closest Deterministic Policy to $\pi^{\epsilon-g}$

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

π_{ter}

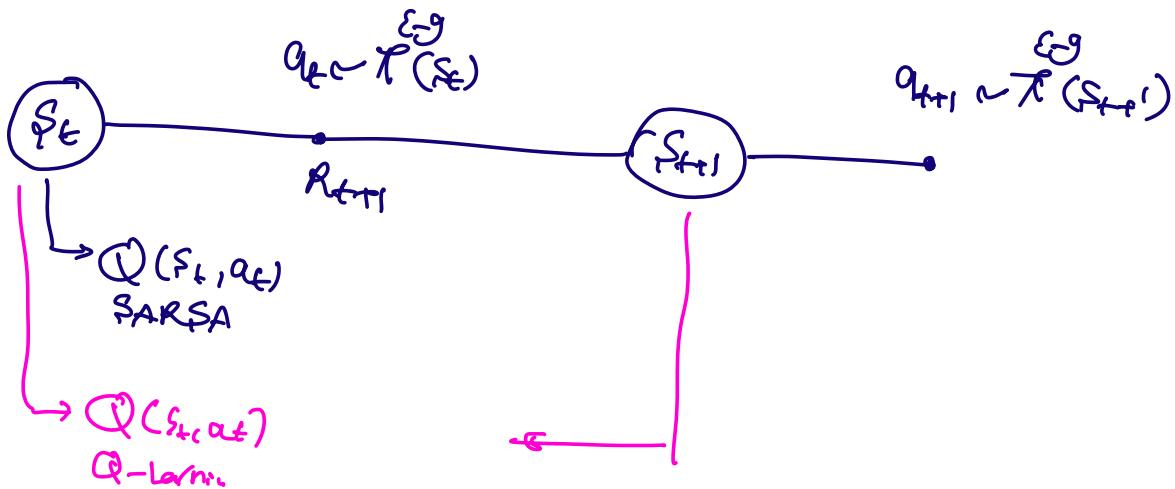
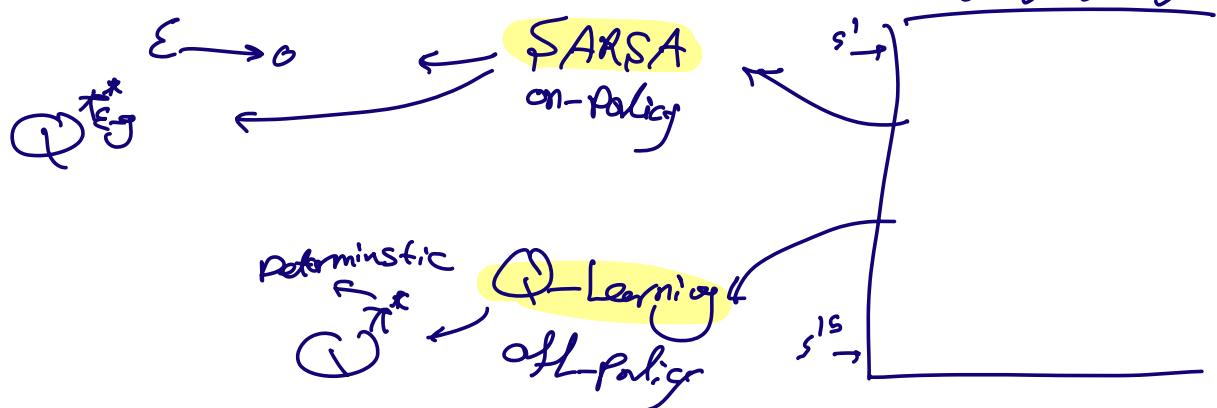
$$\pi(s) = \arg \max_{a \in A} Q(s, a)$$

Off-Policy vs On-Policy

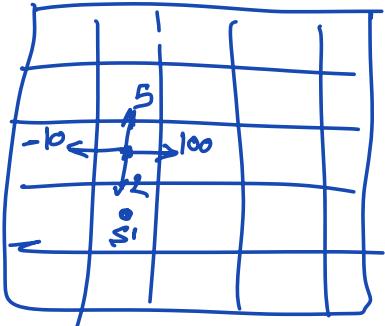
$$a_t \sim \pi^{\text{eg}}_t : \begin{cases} \text{greedy} & 1-\epsilon \\ \text{random} & \epsilon \end{cases} = \begin{cases} \arg\max Q(s_t, a) & 1-\epsilon \\ \text{random} & \epsilon \end{cases}$$

^{last}

$$\pi^*(s) = \arg\max_{a \in A} Q(s, a)$$



$$Q(S_t, a_t) = Q(S_t, a_t) + \alpha [R_{t+1} + \gamma \max_{a'} Q(S_{t+1}, a') - Q(S_t, a_t)]$$



$$T^{Eg}(s_{2+1}) = \begin{cases} R & 1 - \frac{\epsilon}{4} + \frac{\epsilon_1}{4} \\ L & \frac{\epsilon}{4} \\ D & \frac{\epsilon_1}{4} \\ U & \frac{\epsilon_1}{4} \end{cases}$$

$$Q(s', u) = Q(s', u) + \alpha \left[R + \gamma \underbrace{- Q(s', u)}_{\text{Q-loss} \max\{100, 5, -10, 2\}} \right]$$

$$(1 - \varepsilon + \frac{\varepsilon}{4})(100) + \frac{\varepsilon}{4}(-10) + \frac{\varepsilon}{4}(5) + \frac{\varepsilon}{4}(2)$$

100 $1 - \varepsilon + \frac{\varepsilon}{4}$
 -10 $\frac{\varepsilon}{4}$
 5 $\frac{\varepsilon}{4}$
 2 $\frac{\varepsilon}{4}$

$$\text{SARSA} \quad T(s) = \text{argmax}_a Q(s, a) \quad \neq \quad \text{Q-Lear} \quad T(s) = \text{argmax}_a Q(s, a)$$

If you are using off-policy \Rightarrow you can explore as you want
Batch Learning

For On-Policy $\rightarrow \epsilon$ should be small or approach to small value

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-Learning

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

 Repeat (for each step of episode):

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

 Take action A , observe R, S'

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

$$S \leftarrow S';$$

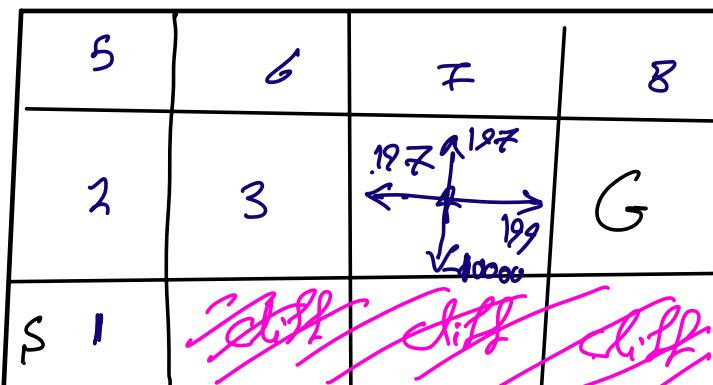
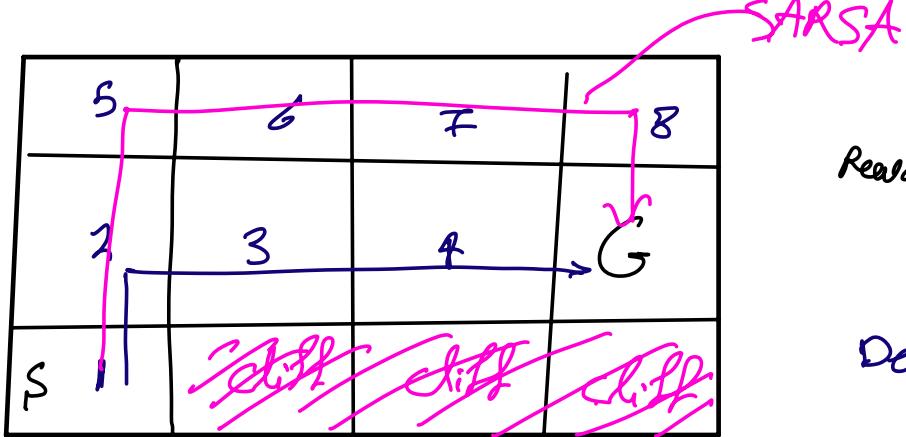
 until S is terminal

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

SARSA

Q Learning

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



$$R = \pi_{(3)}^{\{g\}} = \begin{cases} R & 0.8 + 0.025 \\ \text{oth} & 0.025 \end{cases}$$

(3) → (4)

$$Q(3, R) = Q(3, R) + \alpha [R + \gamma \max_a Q(4, a) - Q(3, R)]$$

α learning

$$3 \xrightarrow{R = \pi_{(3)}^{\{g\}}} 4 \xrightarrow{D = \pi_{(4)}^{\{g\}}}$$

$$\pi^{(4)} \leftarrow \begin{cases} R & 0.925 \rightarrow 199 \\ L & 0.025 \\ U & 0.025 \\ D & 0.025 \end{cases}$$

SARSA

$$Q(s, a) = Q(s, a) + \alpha [R_{-1} + \gamma Q(s_1, \underbrace{\pi(a_1)}_{\text{argmax}}) - Q(s, a)]$$