

Lecture 15 - March 2, 2023

- Temporal Difference Learning

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

HW3 → Due March 17

Project 2 → Due March 7

TA's office hour:

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

Fridays, 2pm-3pm (virtual)

Known MDP ^k	Finite S A	Dynamic Programming: Policy Iteration (PI) Dynamic Programming: Value Iteration (VI)
	Large S A	APPROXIMATE Dynamic Programming (ADP)
Unknown MDP ^u	Finite S & A	Monte Carlo Methods (MC)
		Temporal Difference (TD) Learning: Q-Learning
		Temporal Difference (TD) Learning: SARSA
		Temporal Difference (TD) Learning: Double Q-Learning
		Temporal Difference (TD) Learning: SARSA (λ)
		Temporal Difference (TD) Learning: Actor Critic
Large S & Finite A	Batch Learning	Least Squares Policy Iteration (LSPI)
		Neural Fitted Q Iteration (NFQI)
	Interventive Learning	Deep Reinforcement Learning (DRL): Deep Q Network (DQN)
		Deep Reinforcement Learning (DRL): Double DQN
		Deep Reinforcement Learning (DRL): Dueling DQN
		Deep Reinforcement Learning (DRL): Prioritized DQN
Large S Continuous A	Policy Gradient (PG)	Policy Gradient (PG): REINFORCE
		Policy Gradient (PG): REINFORCE with Baseline
		Policy Gradient (PG): One-Step Actor Critic
		Policy Gradient (PG): Deep Deterministic Policy Gradient (DDPG)

Temporal Difference Learning (TD Learning)

MDP: $(S, A, \cancel{P(s'|s,a)}, \cancel{R(s,a,s')})$

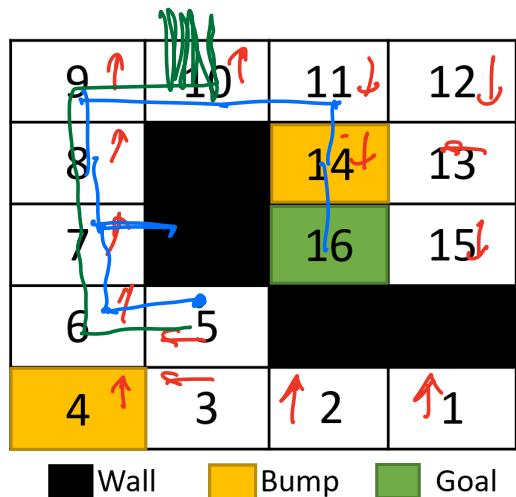
Known Transition \Rightarrow DP

Unknown Trans \Rightarrow MC

$$\begin{aligned}
 V_{\pi}(s) &= E[G_t \mid s_t = s, \pi] \approx \frac{1}{N} \sum_{i=1}^N G_t^i \quad) \text{MC} \\
 &= E[R_{t+1} + \gamma V_{\pi}(s_{t+1}) \mid s_t = s, \pi] \\
 &= \sum_{s'} P(s' | s, \pi(s)) [R_{t+1} + \gamma V_{\pi}(s')] \quad) \text{DP}
 \end{aligned}$$

$$V_{\pi}(s) = \frac{1}{N} \sum_{i=1}^N G_t^i$$

$$\begin{aligned}
 V_{\pi}(s) &= 0.9 [-1 + \gamma V_{\pi}(s)] \\
 &\quad + \frac{0.1}{3} [-4 + \gamma V_{\pi}(s)] + \frac{0.1}{3} [-4 + \gamma V_{\pi}(s)] \\
 &\quad = \frac{0.1}{3} [-1 + \gamma V_{\pi}(s)]
 \end{aligned}$$



Monte Carlo

$$V(s) = V(s) + \alpha [G_t^n - V(s)]$$

$$V(s) = \frac{G_t^1 + G_t^2 + \dots + G_t^n}{n}$$
$$G_t^1 = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$$

Temporal Difference Learning

$$G_t^1 = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$$

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

$$\begin{matrix} s_0 \\ | \\ R_{t+1} \\ | \\ s_{t+1} \end{matrix}$$

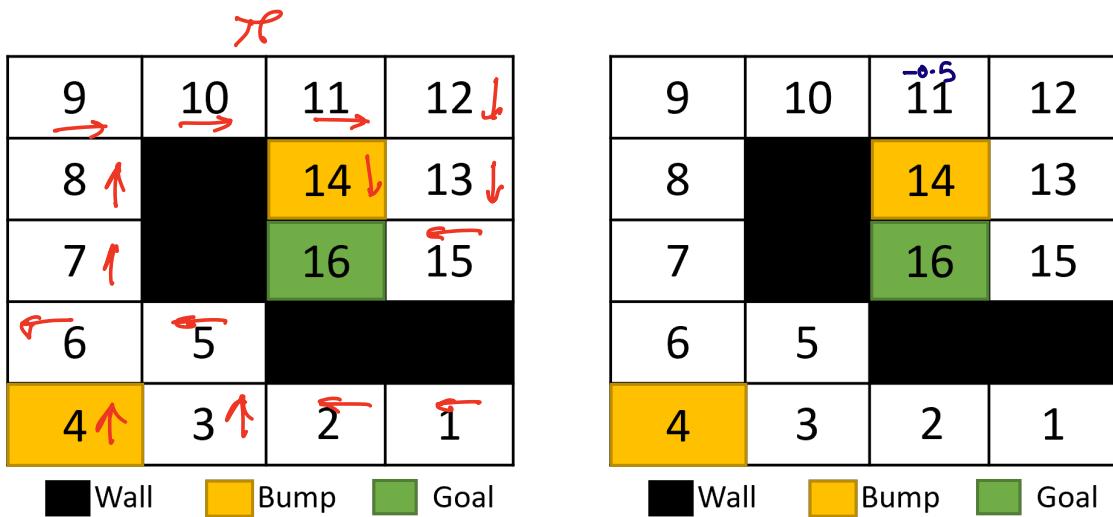
$$V^\pi(s) = E[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \mid s_f = s, \pi]$$

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

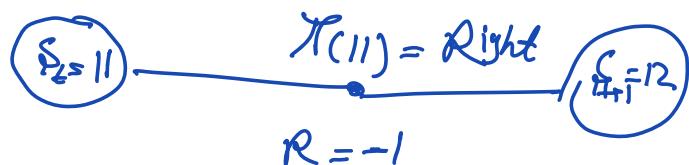
$R_{t+1} + \gamma V(s_{t+1})$: TD Target

$R_{t+1} + \gamma V(s_{t+1}) - V(s_t)$: TD Error

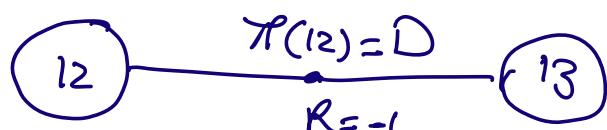
$$\left\{ \begin{array}{l} V_T(s) = \sum_s P(s' | s, \pi(s)) [R + \gamma V_T(s')] \rightarrow DP \\ V(s) = V(s) + \alpha [G_t - V(s)] \rightarrow MC \\ V(s_t) = V(s_t) + \alpha [R_{t+1} + \gamma V(s_{t+1}) - V(s_t)] \downarrow TD \end{array} \right.$$



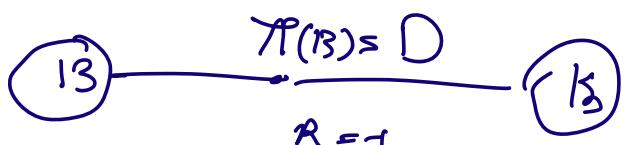
Random State



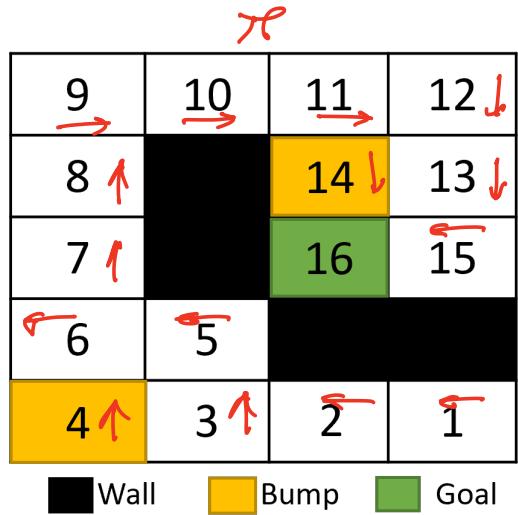
$$V(11) = \frac{V(11)}{\sigma} + \alpha \left[\frac{R}{\sigma \cdot 5} + \gamma \frac{V(12)}{\sigma} - \frac{V(11)}{\sigma} \right] = -0.5$$



$$V(12) = \frac{V(12)}{\sigma} + \alpha \left[\frac{R}{\sigma \cdot 5} + \gamma \frac{V(13)}{\sigma} - \frac{V(12)}{\sigma} \right] = -0.5$$

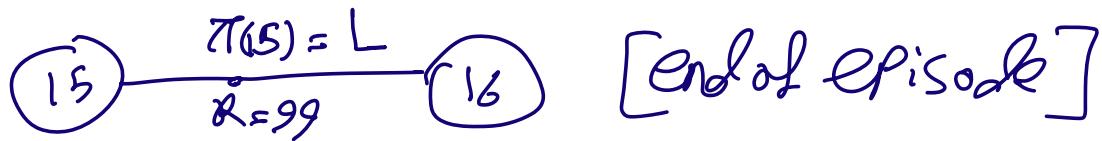


$$V(13) = V(13) + \alpha [R + \gamma V(15) - V(13)] = 0.5$$



9	10	11 -0.5	12 -0.5
8		14 1.0	13 -0.5
7		16	15 1.0
6	5		
4	3	2	1

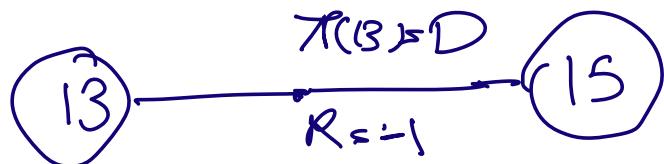
■ Wall ■ Bump ■ Goal



$$V(5) = \frac{V(15)}{0} + \frac{\alpha}{0.5} \left[\frac{R}{99} + \gamma V(\text{Goal}) - \frac{V(15)}{0} \right] = 49.5$$

Episode 2

Random state



$$V(13) = \frac{V(13)}{-0.5} + \frac{\alpha}{0.5} \left[\frac{R}{-1} + \gamma \frac{V(15)}{49.5} - \frac{V(13)}{-0.5} \right] = 24$$

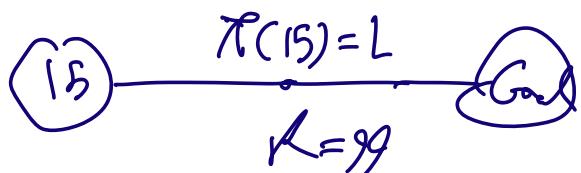
π

9	10	11	12 ↓
8 ↑		14 ↓	13 ↓
7 ↑		16	15 ←
6 ←	5 ←		
4 ↑	3 ↑	2 ←	1 ←

■ Wall ■ Bump ■ Goal

9	10	11	12
8		14	13
7		16	15
6	5		
4	3	2	1

■ Wall ■ Bump ■ Goal

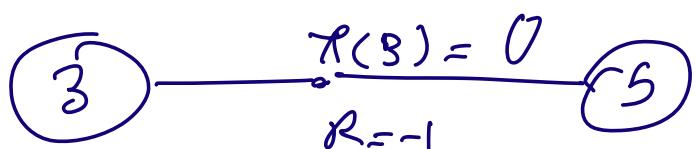


[Terminal End of episode]

$$V(15) = \frac{V(15)}{49.5} + \alpha \left[\frac{R}{0.5} + \gamma \frac{V(\text{Goal}) - V(15)}{0} \right]$$

$$= 74.25$$

episode 3



$$V(3) = -0.5$$

Tabular TD(0) for estimating v_π

Input: the policy π to be evaluated

Initialize $V(s)$ arbitrarily (e.g., $V(s) = 0$, for all $s \in S^+$)

Repeat (for each episode):

Initialize S

Repeat (for each step of episode):

$A \leftarrow$ action given by π for S

Take action A , observe R, S'

$V(S) \leftarrow V(S) + \alpha [R + \gamma V(S') - V(S)]$

$S \leftarrow S'$

until S is terminal

$$\pi \rightarrow V^\pi \Rightarrow \underbrace{\pi' > \pi}_{\text{.}}$$

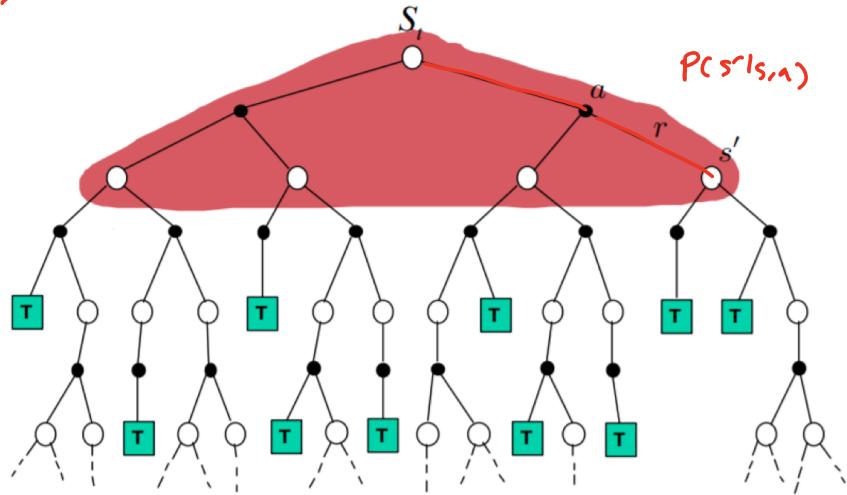
$$\pi'(s) = \arg \max_{a \in A} \sum_{s'} p(s'|s, a) [R + \gamma V_\pi(s')] \quad \begin{array}{l} \text{Unknown} \\ \hline \end{array}$$

$Q_\pi(s, a)$

$$\pi \rightarrow Q^\pi(s, a) \rightarrow \pi'$$

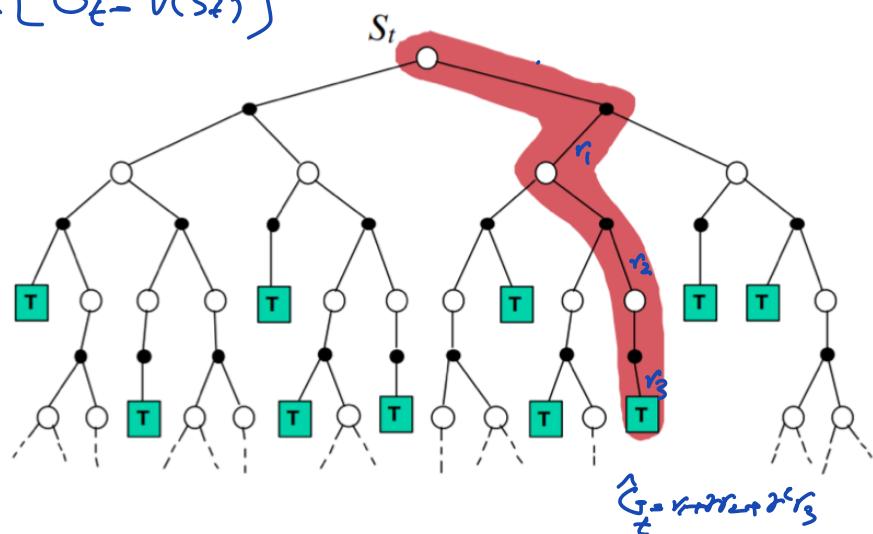
$$DP: V_{k+1}(s) = \max_{a \in A} \sum_{s'} P(s'|s, a) [R + \gamma V_k(s')]$$

Model of system



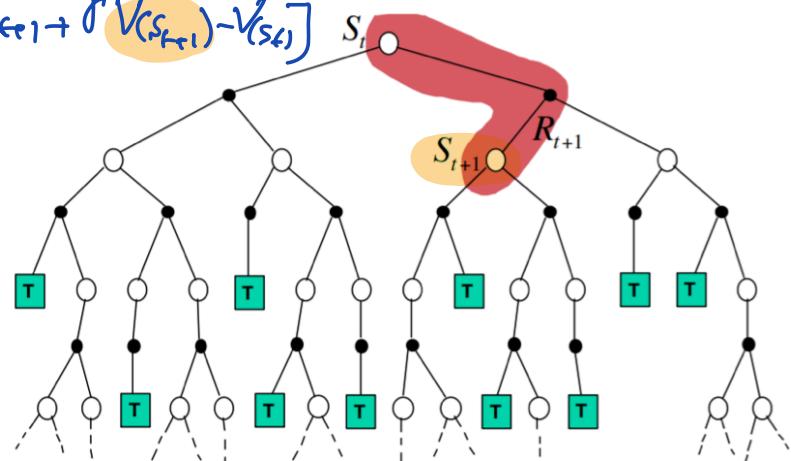
$$MC: V(s_t) = V(s_t) + \alpha [\hat{G}_t - V(s_t)]$$

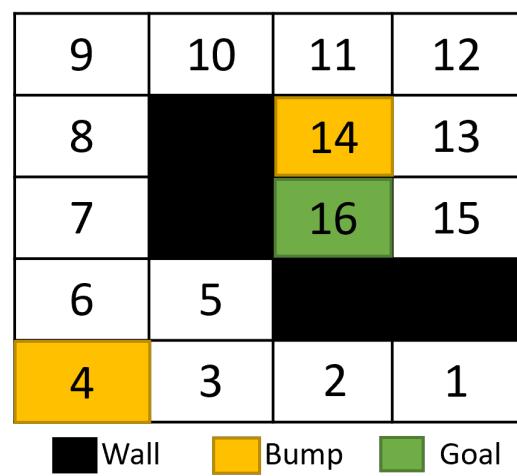
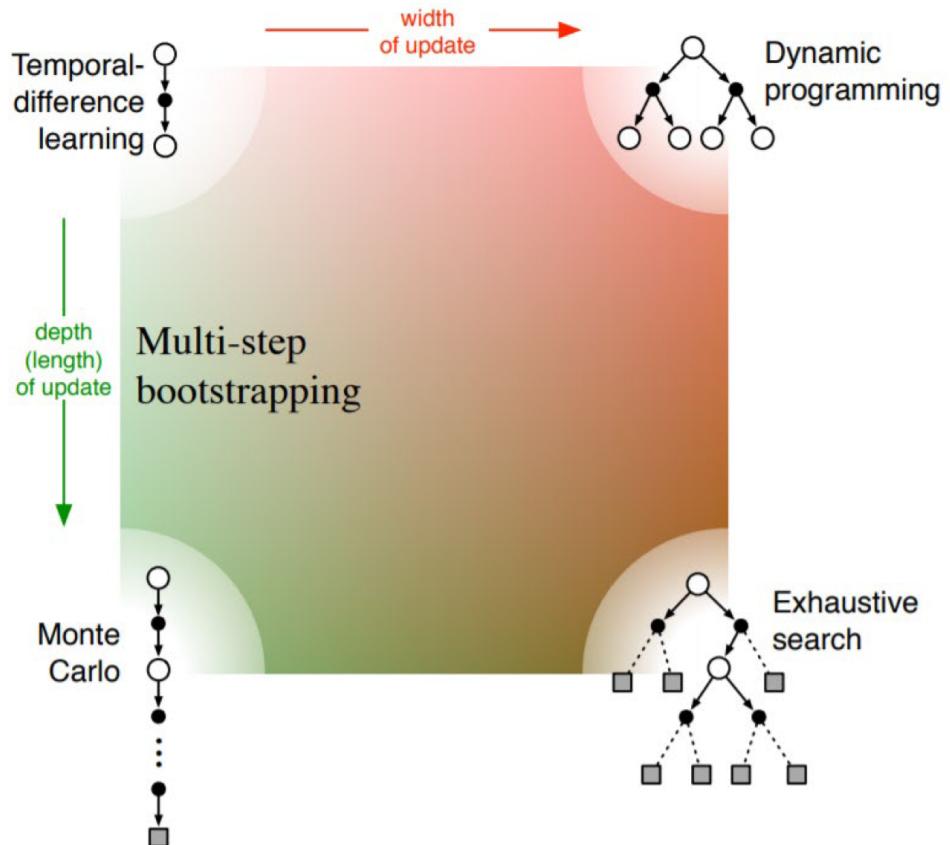
Model ~~X~~
delay ✓
Interactive ~~X~~



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

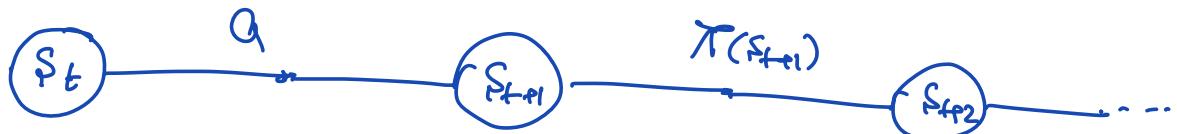
Model ~~X~~
Delay one step
Interactive ✓





$$V_{\pi}(s) \rightarrow \pi'(s) = \sum p(s'|s) \cancel{a}$$

$$Q_{\pi}(s, a) \rightarrow \pi'(s) = \underset{a \in A}{\operatorname{argmax}} Q_{\pi}(s, a)$$



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

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

- {
 1 - How to improve policy
 2 - How to update Q in interactive way

$$\pi(s_t) = \begin{cases} \underset{a \in A}{\operatorname{argmax}} Q(s_t, a) & \text{w.p. } 1-\epsilon \\ \text{Random} & \text{w.p. } \epsilon \end{cases}$$

E-greedy Policy

9	10	11	12
18		14	13
7		16	15
6	5		
4	3	2	1

█ Wall █ Bump █ Goal

deterministic Policy

9	10	11	12
18			13
7		16	15
6	5		
4	3	2	1

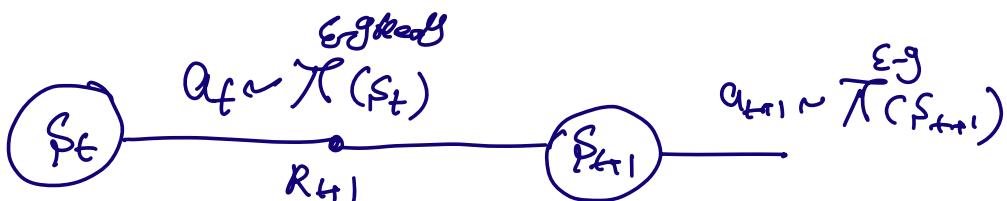
█ Wall █ Bump █ Goal

E-greedy

$$V(s) = E [R_{t+1} + \gamma R_{t+2} + \dots \mid s_t = s, \pi^{\text{E-greedy}}]$$

E-greedy

$$Q(s_t, a_t) = E [R_{t+1} + \gamma R_{t+2} + \dots \mid s_t = s, a_t = a, \pi^{\text{E-greedy}}]$$



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

State – Action – Reward – State – Action

SARSA algorithm

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