

Lecture 14 - Feb 28, 2023

- Monte-Carlo Methods
  - First-Visit MC
  - Online MC
- Temporal Difference Learning
  - TD(0)
  - SARSA
  - Q-Learning
  - On-Policy Vs. Off-Policy

HW3 is posted → Due March 17

Project 2 → Due March 5

TA's office hour:

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

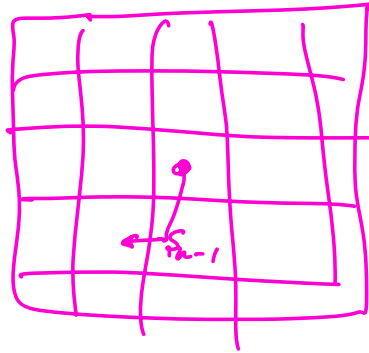
Fridays, 2pm-3pm (virtual)

$$S' = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix} \leftarrow \sum_i \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

$$S^1 = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

نہی نہی =  $\Sigma_k$

$$(U(a))_{ij} = \left( \begin{array}{c} \text{---} \text{---} \boxed{\cdot} \end{array} \right)_{16 \times 16}$$



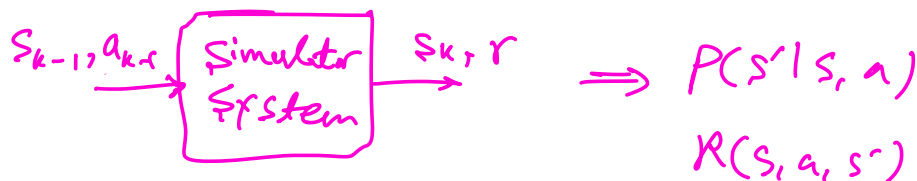
# Monte-Carlo Methods:

MDP( $S, A, \gamma, r$ )

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

Wall
  Bump
  Goal

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



$$\pi^*(s) \leftarrow a$$

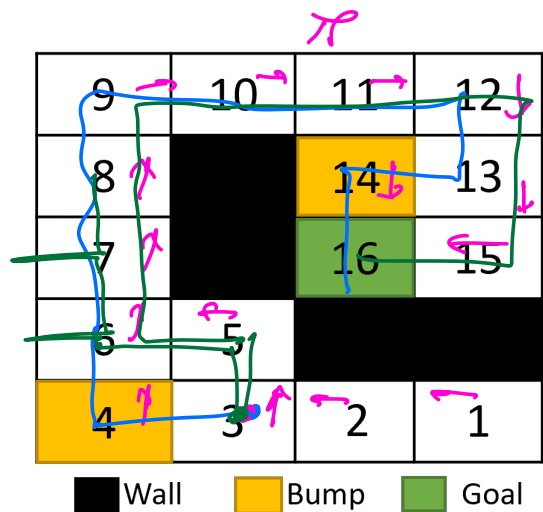
$$V(s) = E \left[ \overbrace{R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots}^{G_t} \mid s_t = s, \pi \right]$$

$\pi^* \quad V_{\pi^*}(s) \geq V_{\pi}(s) \text{ for all } \pi \quad \checkmark \checkmark$

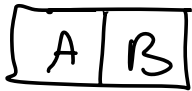
$$\begin{aligned}
 V^\pi(s) &= E \left[ R_{t+1} + \gamma \underbrace{V^\pi(s_{t+1})}_{\text{Bellman Eq}} \mid s_t = s, \pi \right] \\
 &= \sum_{s'} \underbrace{P(s' | s, \pi(s))}_{\text{Pr}} [R + \gamma V^\pi(s')]
 \end{aligned}$$

$s_t = 3$

$-1, -1, -1, \dots, -1, 99$   $G_t$



$$\begin{aligned}
 V^\pi(s) &= E [ G_t \mid s_t = s, \pi ] \\
 &\approx \frac{1}{N} \sum_{i=1}^N G_t^i
 \end{aligned}$$



$$\pi = \begin{bmatrix} a^2 \\ a^1 \end{bmatrix} \Rightarrow V^\pi$$

episode 1

Random

$$\hookrightarrow A, \pi(A) = a^2, A, r = -1$$

$$A, \pi(A) = a^2, B, r = 4$$

$$B, \pi(B) = a^1, B, r = 5$$

$$B, \pi(B) = a^1, B, r = 5$$

$$B, \pi(B) = a^1, A, r = 0$$

$$A, \pi(A) = a^2, B, r = 4$$

$$G(A) = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$$

$$= -1 + \gamma 4 + \gamma^2 5 + \gamma^3 5 + \gamma^4 0 + \gamma^5 4 + \dots$$

$$G(B) = 5 + \gamma 5 + \gamma^2 0 + \gamma^3 4$$

$$V^\pi(A) \approx G(A)$$

$$V^\pi(B) \approx G(B)$$

Unknown

$$M(Q^1) = \begin{bmatrix} 0.9 & 0.1 \\ 0.1 & 0.9 \end{bmatrix}$$

$$M(Q^2) = \begin{bmatrix} 0.1 & 0.9 \\ 0.9 & 0.1 \end{bmatrix}$$

$$A \xleftarrow{5} B$$

$$\xrightarrow{-1} a^2$$

episode 2

$$B, \pi(B) = a^1, B, r = 5$$

$$B, \pi(B) = a^1, A, r = 0$$

$$A, \pi(A) = a^2, B, r = 4$$

First visit of A

$$B, \pi(B) = a^1, B, r = 5$$

$$B, \pi(B) = a^1, B, r = 5$$

$$B, \pi(B) = a^1, B, r = 5$$

$$B, \pi(B) = a^1, B, r = 5$$

$$G(A) = 4 + \gamma 5 + \gamma^2 5 + \gamma^3 5 + \gamma^4 5$$

$$G(B) = 5 + \gamma 0 + \gamma^2 4 + \gamma^3 5 + \dots$$

$$V^\pi(A) = \frac{G(A)^1 + G(A)^2}{2}$$

$$V^\pi(B) = \frac{G(B)^1 + G(B)^2}{2}$$

$$\underbrace{V(s)}_{\text{new estimate}} = \underbrace{V(s)}_{\text{old estimate}} + \alpha [G(s) - V(s)]$$

$$\alpha = \frac{1}{n}$$

No applicability

## First-Visit Monte Carlo Policy Evaluation

Initialize:

$\pi \leftarrow$  policy to be evaluated

$V \leftarrow$  an arbitrary state-value function

$Returns(s) \leftarrow$  an empty list, for all  $s \in \mathcal{S}$

Repeat forever:

Generate an episode using  $\pi$

For each state  $s$  appearing in the episode:

$G \leftarrow$  return following the first occurrence of  $s$

Append  $G$  to  $Returns(s)$

$V(s) \leftarrow \text{average}(Returns(s))$



$$A = \{a^1, a^2\}$$

$$|A|^{|\mathcal{S}|}$$

$\leftarrow$  policies

$$\pi^1 = \begin{bmatrix} a^1 \\ a^1 \end{bmatrix}$$

$$\pi^2 = \begin{bmatrix} a^1 \\ a^2 \end{bmatrix}$$

$$\pi^3 = \begin{bmatrix} a^2 \\ a^1 \end{bmatrix}$$

$$\pi^4 = \begin{bmatrix} a^2 \\ a^2 \end{bmatrix}$$

$$V^{\pi^1} \begin{bmatrix} \bigcirc \\ \bigcirc \end{bmatrix}$$

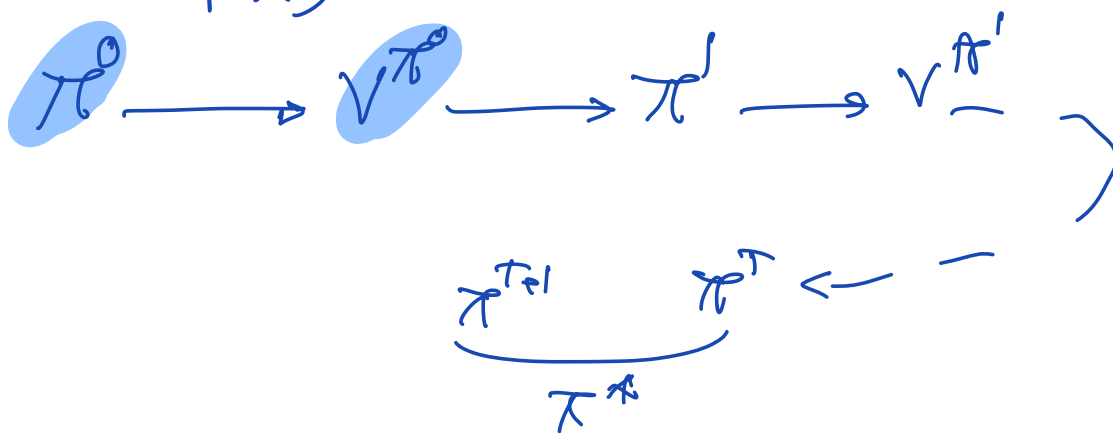
$$V^{\pi^2}$$

$$V^{\pi^3}$$

$$V^{\pi^4}$$

$\pi$

# Policy Iteration



$$\pi(s) = \underset{a \in A}{\operatorname{argmax}} \underbrace{\sum_{s'} P(s'|s,a) [R + \gamma V_{\pi}(s')]}_{\text{Unknown}}$$

$$Q_{\pi}(s, a)$$



$\Rightarrow$  our Goal is to estimate  $Q_{\pi}(s, a)$   
Using MC



$$\boxed{A|B}$$

$$\pi = \begin{bmatrix} a^1 \\ a^2 \end{bmatrix}$$

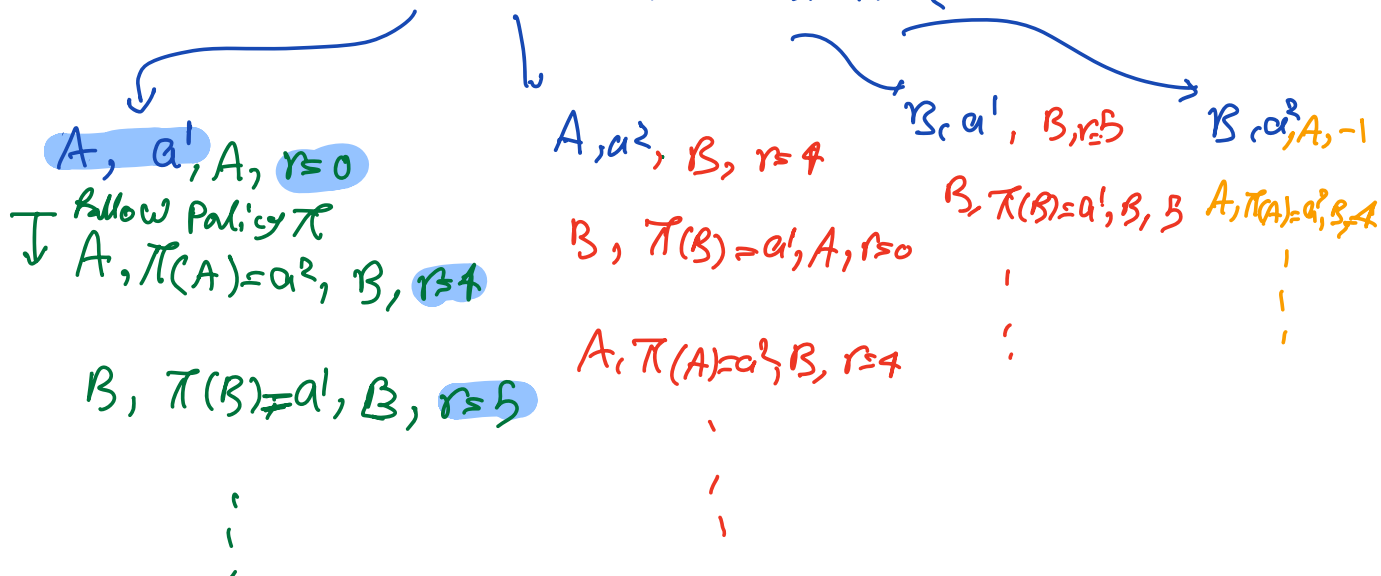
$$Q_{\pi}(A, a^1)$$

$$Q_{\pi}(B, a^1)$$

$$Q_{\pi}(A, a^2)$$

$$Q_{\pi}(B, a^2)$$

Start with all pairs of  $s, a$



$$G_{A, a^1} = 0 + \gamma 4 + \gamma^2 5 + \dots$$

$$Q^{\pi}(A, a^1) \approx G_{A, a^1}$$

$$G_{A, a^2} = 4 + \gamma 0 + \gamma^2 4$$

$$G_{A, a^2}^2 = 4 + \gamma 5 + \gamma^2 \dots$$

$$G_{A, a^2}^3 = 4 + \gamma \dots$$

$$Q^{\pi}(A, a^2) = \frac{G_{A, a^2}^1 + G_{A, a^2}^2 + G_{A, a^2}^3}{3}$$

$$Q^{\pi}(A, a^1)$$

$$Q^{\pi}(B, a^1)$$

$$Q^{\pi}(A, a^2)$$

$$Q^{\pi}(B, a^2)$$

$$\pi^0 \xrightarrow{uc} \bigoplus \pi^i(s, a) \longrightarrow \pi'$$

$$\pi'(s) = \underset{a \in A}{\operatorname{argmax}} Q_{\pi^0}(s, a)$$

$$\underline{\pi'(A)} = \underset{\substack{Q_{\pi}(A, a^1), Q_{\pi}(A, a^2)}}{\operatorname{argmax}} Q_{\pi}(A, a)$$

## Monte Carlo Policy Iteration

Initialize, for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$ :

$Q(s, a) \leftarrow$  arbitrary

$\pi(s) \leftarrow$  arbitrary

$Returns(s, a) \leftarrow$  empty list

Fixed point is optimal  
policy  $\pi^*$

Now proven (almost)

Repeat forever:

Choose  $S_0 \in \mathcal{S}$  and  $A_0 \in \mathcal{A}(S_0)$  s.t. all pairs have probability  $> 0$

Generate an episode starting from  $S_0, A_0$ , following  $\pi$

For each pair  $s, a$  appearing in the episode:

$G \leftarrow$  return following the first occurrence of  $s, a$

Append  $G$  to  $Returns(s, a)$

$Q(s, a) \leftarrow \text{average}(Returns(s, a))$

For each  $s$  in the episode:

$\pi(s) \leftarrow \operatorname{argmax}_a Q(s, a)$

# Online Monte Carlo Algorithm

$$\pi(a|s) \begin{cases} \leftarrow \text{Solomon Policy} \\ \leftarrow \text{Randomness} \end{cases}$$

$$\pi^{\epsilon\text{-greedy}}(s) = \begin{cases} \arg\max_{a \in A} Q(s, a) & 1-\epsilon \\ \text{Random}\{a^1, a^2, \dots, a^L\} & \epsilon \end{cases}$$

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

Wall
  Bump
  Goal

$$a^* = L = \arg\max Q(s, a)$$

$$P(s) = \begin{cases} a^* & 1-\epsilon + \frac{\epsilon}{|A|} \\ a \neq a^* & \frac{\epsilon}{|A|} \end{cases}$$

$$\begin{aligned} \epsilon &= 0.1 \\ |A| &= 4 \\ \pi(s) &= \begin{cases} a^* & 0.925 \\ a \neq a^* & \frac{0.1}{4} = 0.025 \end{cases} \quad a \in A - \{a^*\} \end{aligned}$$

$$a^* = \operatorname{argmax} Q(s, a)$$

$$\pi^{\epsilon}(s) = 0.25 \quad \text{all } a \in A$$

$$\pi^{\epsilon, \text{greedy}} = \begin{cases} \text{greedy} & 1 - \epsilon \\ \text{random} & \epsilon \end{cases}$$

$a^1 a^2 a^3 a^4$   
 $\downarrow$   
 $a^1 a^2 a^3 a^4$

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

Wall
  Bump
  Goal

$$Q(s, a) = 0$$

Random State

$$A, \pi^{\epsilon, \text{greedy}}(A) = a^1, A, r=0 \Rightarrow Q(A, a^1) = 0 + \gamma Q + \gamma^2 S + \gamma^3 C + \dots$$

$$A, \pi^{\epsilon}(A) = a^2, B, r=4 \Rightarrow Q(A, a^2) = \dots$$

$$B, \pi^{\epsilon, \text{greedy}}(B) = a^1, B, r=5 \Rightarrow Q(B, a^1) =$$

$$B, \pi(B) = a^2, A, r=-1 \Rightarrow Q(B, a^2) =$$

⋮

$$\pi^{\text{new } \epsilon\text{-greedy}}(s) = \begin{cases} \operatorname{argmax} Q(s, a) = a_s^* & 1 - \epsilon + \frac{\epsilon}{|A|} \\ a \neq a_s^* & \frac{\epsilon}{|A|} \end{cases}$$

## On-Policy Monte Carlo Control

Initialize, for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$ :

$Q(s, a) \leftarrow$  arbitrary

$Returns(s, a) \leftarrow$  empty list

$\pi(a|s) \leftarrow$  an arbitrary  $\varepsilon$ -soft policy

Repeat forever:

(a) Generate an episode using  $\pi$

(b) For each pair  $s, a$  appearing in the episode:

$G \leftarrow$  return following the first occurrence of  $s, a$

Append  $G$  to  $Returns(s, a)$

$Q(s, a) \leftarrow \text{average}(Returns(s, a))$

(c) For each  $s$  in the episode:

$A^* \leftarrow \arg \max_a Q(s, a)$

For all  $a \in \mathcal{A}(s)$ :

$$\pi(a|s) \leftarrow \begin{cases} 1 - \varepsilon + \varepsilon/|\mathcal{A}(s)| & \text{if } a = A^* \\ \varepsilon/|\mathcal{A}(s)| & \text{if } a \neq A^* \end{cases}$$