

Lecture 5 - Jan 27, 2023

- Reinforcement Learning Preliminaries

I

- State, Action, Reward, Policy
- Returns and Expected Returns
- State Value Function
- State-Action Value Function
- Bellman Equation and Optimality

HW1 → Due Jan 28

Project 1 → Due Feb 7

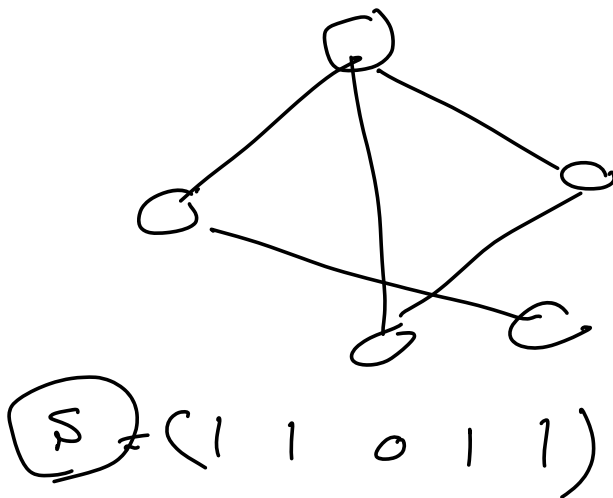
TA's office hour:

Wednesdays, 12pm - 1pm (in-person)

Fridays, 12pm - 1pm (virtual)

Markov Decision Process (MDP)

$\langle \underline{S}, \underline{A}, \underline{T}, R \rangle$



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

Wall Bump Goal

→ Transition Probability
 $T: S \times A \times S$

$A = \{U, D, L, R\}$

$P(\underline{S'} | \underline{S}, a)$
 next correct

$S = \textcircled{3}, a = \textcircled{L} \rightarrow \begin{cases} S' = 4 & \text{w.p. } 0.9 \\ S' = 5 & \text{w.p. } \frac{0.1}{3} \\ S' = 2 & \text{w.p. } \frac{0.1}{3} \\ S' = 3 & \text{w.p. } \frac{0.1}{3} \end{cases}$

$$P(S' \mid S=3, a=L) = \begin{cases} 0.9 \\ 0.1/3 \\ 0 \end{cases}$$

$$S'=4$$

$$S'=5 \text{ or } 2 \text{ or } 3$$

$$S'=15$$

$$R: S \times A \times S$$

$$R(S, a, S')$$

$$R(S=3, a=L, S'=5) = -1$$

$$R(S=3, a=L, S'=4) = -10 - 1 = -11$$

$$R(S=3, a=U, S'=15) = -1$$

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

Wall
 Bump
 Goal

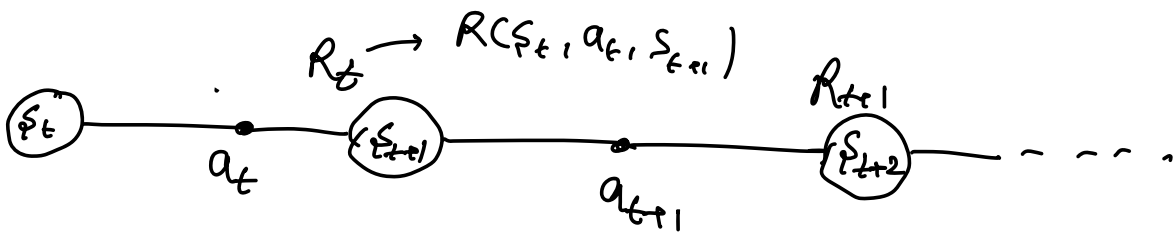
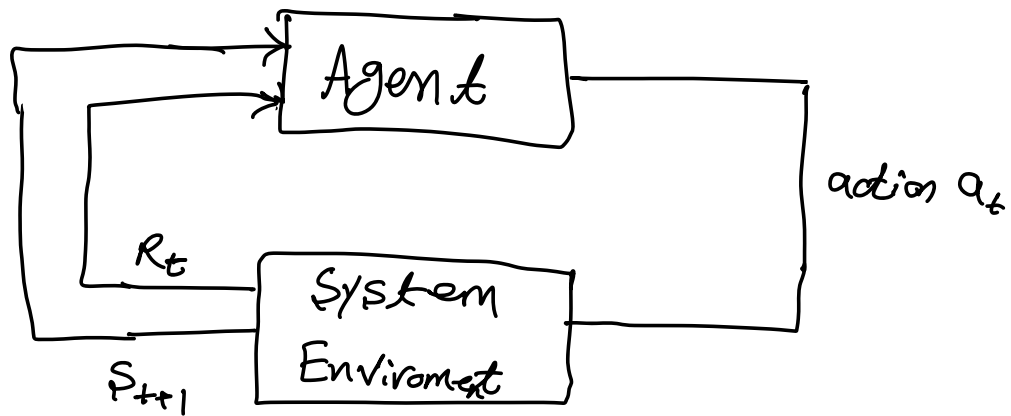
$$\begin{cases} \text{Goal: } 100 \\ \text{Bump: } -10 \\ \text{movement: } -1 \end{cases}$$

$$P(S_{t+1}=4 \mid S_t=3, a=L, S_{t-1}=2, a_{t-1}=L)$$

$$= P(S_{t+1}=4 \mid S_t=3, a=L) = 0.9$$

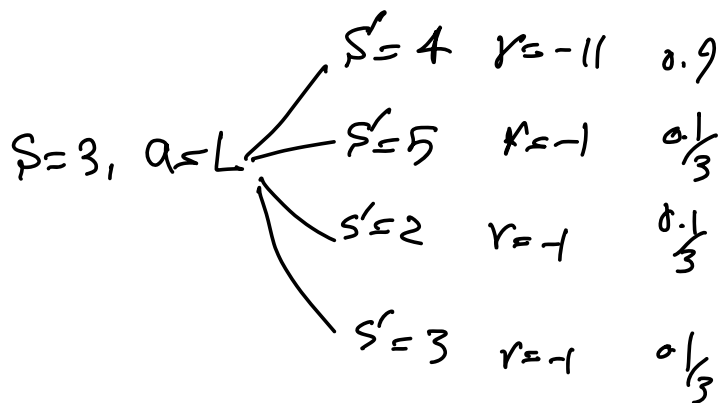
9	10	11	12
8		14	13
7		16	15
6	5		
4	3	2	1

Wall
 Bump
 Goal



MDP(S, A, T, R)

$P(S', r | S, a)$



$$P(S'=4, r=-1 | S=3, a=L) = 0$$

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

Wall
 Bump
 Goal

Types of Rewards

Type I: $R(s, a, s')$ Immediate Reward

Type II: $R(s, a)$ Expected Immediate Reward

$$R(\bar{s}^3, a=L, s') = \begin{cases} -11 & s'=4 & 0.9 \\ -1 & s'=5 & 0.1/3 \\ -1 & s'=2 & 0.1/3 \\ -1 & s'=3 & 0.1/3 \end{cases}$$

$$\begin{aligned} R(s=3, a=L) &= -11 \times 0.9 + (-1) \times \frac{0.1}{3} \\ &\quad + (-1) \times \frac{0.1}{3} + (-1) \times \frac{0.1}{3} \\ &= -10 \end{aligned}$$

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

Wall
 Bump
 Goal

$$P(s', r | s, a)$$

$$P(A, B)$$

$$P(A) = \sum_{b \in B} P(A, B=b)$$

$$P(s' | s, a) = \sum_r P(s', r | s, a)$$

\Rightarrow Marginalization

$$P(r | s, a) = \sum_{s'} P(s', r | s, a)$$

$$R(s, a) = E [R(s, a, s') | s_t = s, a_t = a]$$

$$= \sum_{s'} P(s' | s, a) R(s, a, s')$$

$$R(s=15, a=0) =$$

$$0.9 \times \frac{s'=13}{-1} + 0.1 \times \frac{s'=15}{-1} + 0.1 \times \frac{s'=15}{-1} + 0.1 \times \frac{s'=16}{99}$$

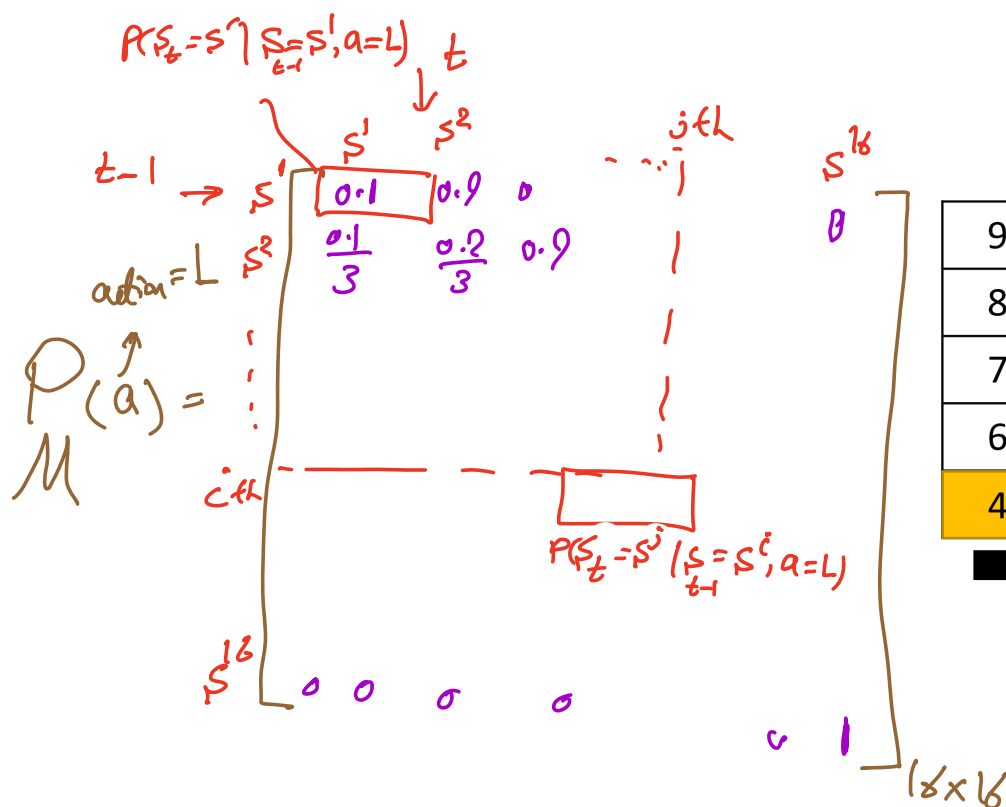
9	10	11	12
8		14	13
7		16	15
6	5		
4	3	2	1

Wall Bump Goal

0.9
0.1
0.1
0.1

Transition Probability

- element-wise $P(s'|s, a)$
- Matrix form
 - Transition Matrix



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

Wall (Black) Bump (Yellow) Goal (Green)

Policy: $\begin{cases} \text{Deterministic} \iff \text{Most Common} \\ \text{Stochastic} \end{cases}$

$$\pi: S \rightarrow A$$

$$\Pi = \{ \pi^1, \pi^2, \dots, \pi^{15} \}$$

π

9 \rightarrow	10 \rightarrow	11 \downarrow	12 \downarrow
8 \uparrow		14 \downarrow	13 \downarrow
7 \uparrow		16 \leftarrow	15 \leftarrow
6 \uparrow	5 \leftarrow		
4 \rightarrow	3 \leftarrow	2 \leftarrow	1 \leftarrow

Wall
 Bump
 Goal

$$\pi = \begin{bmatrix} \pi(s^1) \\ \pi(s^2) \\ \vdots \\ \pi(s^{16}) \end{bmatrix} = \begin{bmatrix} L \\ L \\ \vdots \\ L \\ \textcircled{D} \end{bmatrix}$$

$$\pi(s) = U$$

$\pi(a|s) = \begin{cases} a=L & 0.4 \\ a=U & 0.2 \\ a=D & 0.4 \\ a=R & 0 \end{cases}$

\uparrow
stochastic

$\pi(s)$ action that is chosen in state s under Policy π
 $\pi(s) = a$

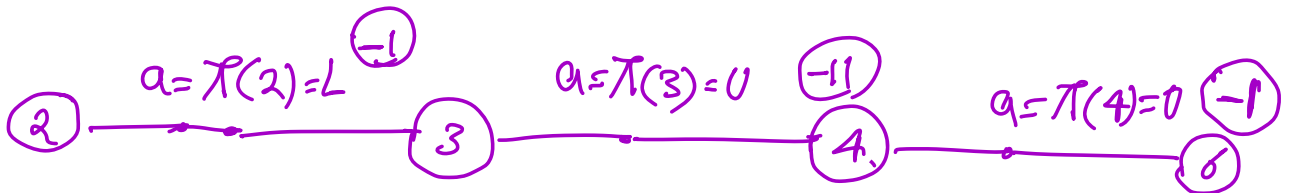
$$\underset{a_{0:T-1}}{\operatorname{argmax}} E \left[\sum_{t=0}^T R(s_t, a_t, s_{t+1}) \mid s_0 = 2, a_{0:T-1} \right]$$

\Downarrow Policy

9 \rightarrow	10 \rightarrow	11 \rightarrow	12 \downarrow
8 \uparrow		14 \downarrow	13 \downarrow
7 \uparrow		16	15 \leftarrow
6 \uparrow	5 \leftarrow		
4 \uparrow	3 \uparrow	2 \leftarrow	1 \leftarrow

Wall
 Bump
 Goal

$$\pi^* = \underset{\pi \in \Pi}{\operatorname{argmax}} E \left[\sum_{t=0}^T R(s_t, a_t, s_{t+1}) \mid s_0 = 2, a_{0:T} \sim \pi \right]$$



Return

$$\pi \rightarrow R_t \rightarrow R_{t+1} \rightarrow R_{t+2} + \dots$$

(s_t, r_t)

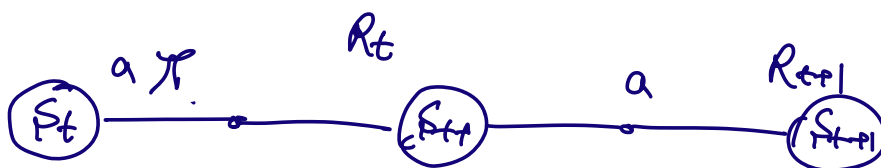
$$G_t =$$

s_{t+2}

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

Wall
 Bump
 Goal

- Total Reward G_t : sum of all future rewards in the episode
- Discounted Reward G_t : sum of all future discounted rewards
- Average Reward G_t : Avg Reward per time step





Episodic Tasks:

"Break into episodes"

Total Reward

$$G_t = R_{t+1} + R_{t+2} + \dots + R_T$$

Terminal Point

$$E[G_t]$$

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

Wall
 Bump
 Goal

Episodic

