Lecture 6 - Jan 31, 2023

- Reinforcement Learning Preliminaries
  - State, Action, Reward, Policy
  - Returns and Expected Returns
  - State Value Function
  - State-Action Value Function
  - Bellman Equation and optimality

Project 1 → Due Feb 7

TA's office hour:    Wendsdays, 12pm-1pm (in-person)
                     Fridays, 12pm-1pm (virtual)

$$MDP(\ S,\ A,\ R,\ P\ )$$

imediate $\longrightarrow R(s, a, s')$ \qquad $P(s' \mid s, a)$

Expeted \qquad $R(s, a) = \sum_{s'} P(s' \mid s, a)\, R(s, a, s')$

Task $\Big\langle$ Episodic

\qquad\qquad\qquad Continuing

Return: Accumulated Reward

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

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

⬛ Wall \quad 🟧 Bump \quad 🟩 Goal

$\gamma \leq 1$

$r \leq 0.9$

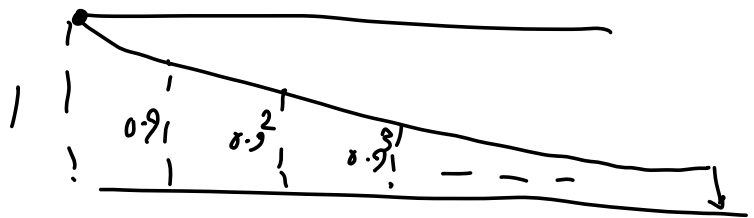$-1 \quad -1 \quad -1 \quad -1 \quad -1) \quad -1, \sim \sim$

$(8.9)^{20}$

$99$

Continuing Task
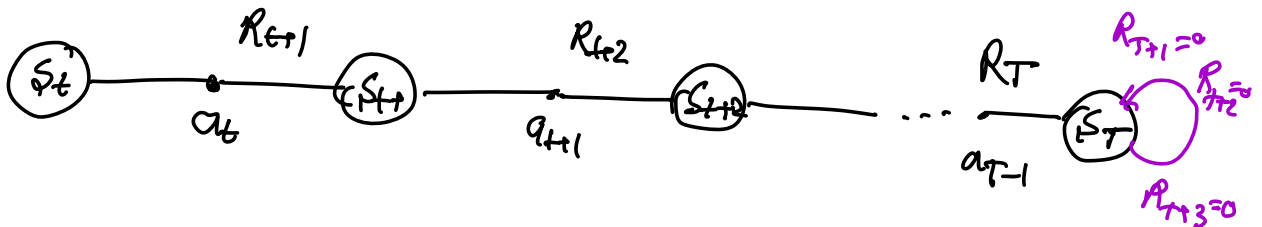
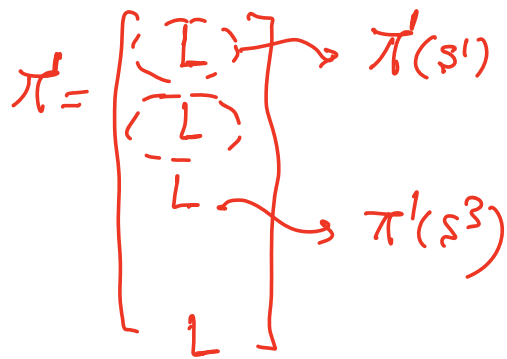$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots$$

Discount Factor

$$0 < \gamma < 1$$



Episodic = Continuing Tasks



$$G_t = R_{t+1} + \gamma R_{t+2} + \cdots + \gamma^T R_T + \gamma^{T+1} R_{t+1} + \gamma^{T+2} R_{t+2} + \cdots$$

$$\pi' = \begin{bmatrix} L \\ L \\ L \\ L \end{bmatrix}$$

$\pi'(s^1)$

$\pi'(s^3)$

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

■ Wall   ▢ Bump   ▢ Goal

Expected Return = Expected Accumulated Reward

$$V(s_t) = E\left[ G_t \mid S_t = 15, \pi' \right]$$

State-Value function

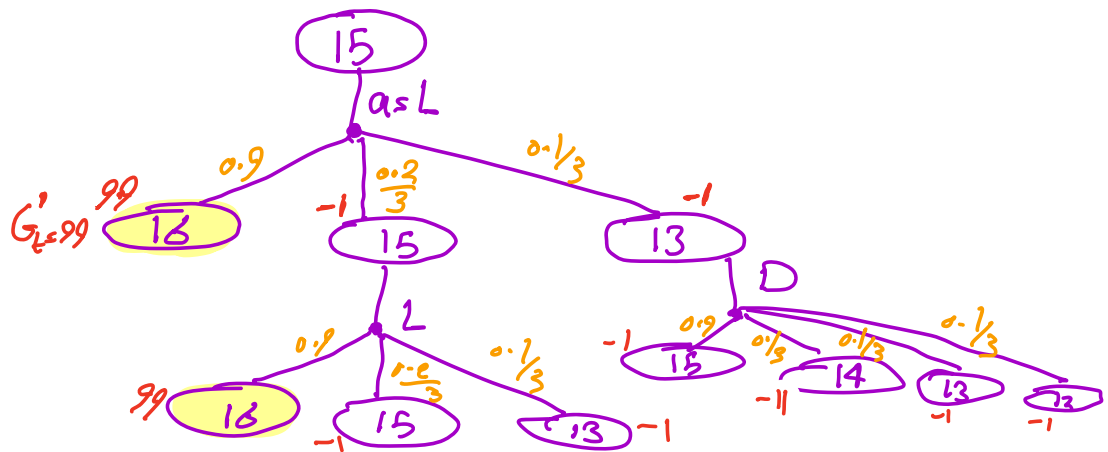$$V^{\pi'}(15) = E\left[ G_t \mid S_t = 15, \pi' \right]$$

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

$$P(s'|s=15, a=L) = \begin{cases} 0.9 & s' = 16 \\ 2\frac{0.1}{3} & s' = 15 \\ \frac{0.1}{3} & s' = 13 \end{cases}$$
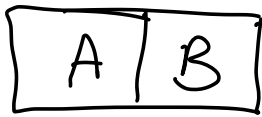
$\pi'$

| 9 → | 10 → | 11 ↗ | 12 ↓ |
|---|---|---|---|
| 8 ↑ | ■ | 14 ↓ | 13 ↓ |
| 7 ↑ | ■ | 16 | 15 ↩ |
| 6 ↑ ↩ 5 | | ■ | ■ |
| 4 ↑ | 3 ↑ ↩ | 2 ↩ | 1 ↩ |

■ Wall   ▮ Bump   ▮ Goal



$$V^{\pi'}(15) = E\left[ G_t \mid S_t = 15, \pi' \right]$$

$$= \sum_{\text{trajectories}} P_t^{j} \, G_t^{j}$$

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

$$S = \{A, B\}$$

$$\text{Reward} \to \begin{cases} -1 \text{ if } \overbrace{\text{switch}}^{a^2} \\ +5 \cdot \text{ be in } B \end{cases}$$

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

$$\underset{\text{keep your location}}{\swarrow} \quad \underset{\text{switch}}{\searrow}$$

$R(S, a, S')$

$R(A, a^1, A) = 0$

$R(A, a^1, B) = 5$

$R(B, a^1, A) = 0$

$\vdots$

$R(B, a^2, B) = 4$

$$P(a^1) = M(a^1) = \begin{array}{c} \\ A \to \\ B \to \end{array}\begin{bmatrix} A \overset{k}{\downarrow} & B \downarrow \\ 1 & 0 \\ \boxed{0} & 1 \end{bmatrix} \begin{array}{c} k-1 \end{array}$$

$$P(S_k = A \mid S_{k-1} = B, \ a = a^1)$$
$$\underset{k-1}{}$$

$$P(a^2) = M(a^2) = \begin{array}{c} A \\ B \end{array}\begin{bmatrix} A & k & B \\ 0 & 1 \\ 1 & 0 \end{bmatrix}$$

$$\pi^1 = \begin{bmatrix} a^1 \\ a^1 \end{bmatrix} \quad \pi^2 = \begin{bmatrix} a^1 \\ a^2 \end{bmatrix} \quad \pi^3 = \begin{bmatrix} a^2 \\ a^1 \end{bmatrix} \quad \pi^4 = \begin{bmatrix} a^2 \\ a^2 \end{bmatrix}$$

$$\pi' = \begin{bmatrix} \pi'(A) \\ \pi'(B) \end{bmatrix} = \begin{bmatrix} a' \\ a' \end{bmatrix} \quad \leftarrow \text{ Stay at your state} \qquad \gamma = 0.9$$
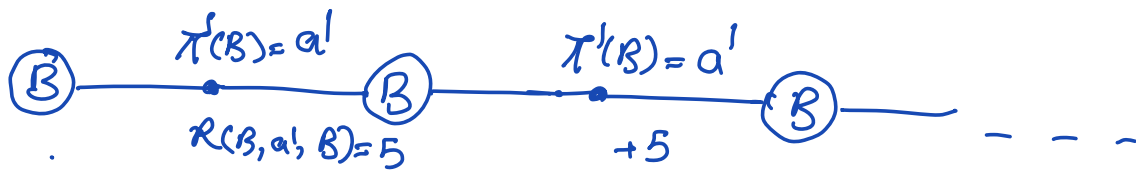
$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots$$

$$V_{\pi'}(s) = E[G_t \mid S_t = s, \pi']$$

$S_0$
$(A') \xrightarrow{\pi'(A) = a'} (A) \xrightarrow{\pi'(A) = a'} (A) \; - - - -$
$\quad\quad R(A,a',A) = 0 \quad , \quad R(A,a',A) = 0$

$$\underline{V_{\pi'}(A)} = 0 + \gamma 0 + \gamma^2 0 + \cdots = 0$$

Expected Accumulated rewards Starting from A and following Policy $\pi'$

$(B) \xrightarrow{\pi'(B) = a'} (B) \xrightarrow{\pi'(B) = a'} (B) \; - - -$
$\quad\quad R(B,a',B) = 5 \quad\quad +5$

$$V_{\pi'}(B) = E[R_{t+1} + \gamma R_{t+2} + \cdots \mid S_t = B, \pi']$$
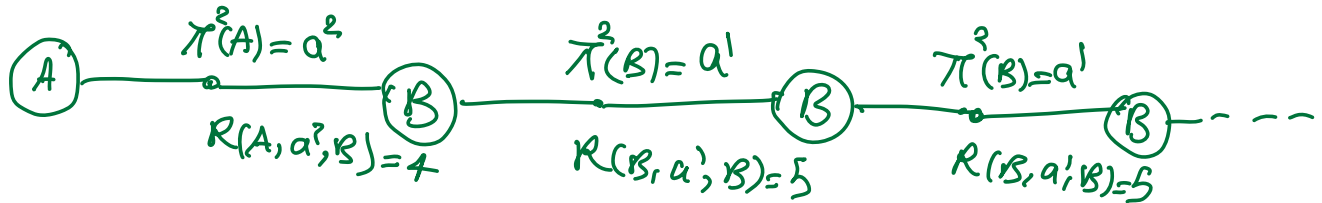
$$\gamma = 0.9$$

$$= 5 + \gamma 5 + \gamma^2 5 + \gamma^3 5 + \cdots$$

$$= 5(1 + \gamma + \gamma^2 + \gamma^3 + \cdots)$$

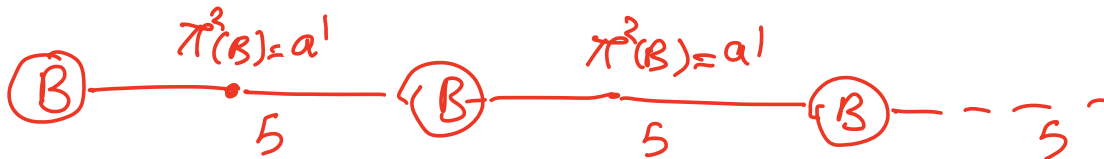$$= 5 \; \underline{\frac{1}{1-\gamma}} = \frac{5}{1-0.9} = 50$$

$$1 + x + x^2 + \cdots = \frac{1}{1-x}$$
$$-1 < x < 1$$

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



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

$R(A,a^2,B)=4$  $R(B,a^1,B)=5$  $R(B,a^1,B)=5$

$$V_{\pi^2}(A) = 4 + \gamma 5 + \gamma^2 5 + \gamma^3 5 + \cdots$$

$$= 4 + \gamma 5 \underbrace{\left( 1 + \gamma + \gamma^2 + \cdots \right)}_{\frac{1}{1-\gamma}} = 49$$



$(B)$ —— $\pi^2(B)=a^1$ —— $(B)$ —— $\pi^2(B)=a^1$ —— $(B)$ - - -

   $5$                    $5$              $5$

$$V_{\pi^2}(B) = 5 + \gamma 5 + \gamma^2 5 + \cdots = 5 \left( 1 + \gamma + \gamma^2 + \cdots \right) = 50$$
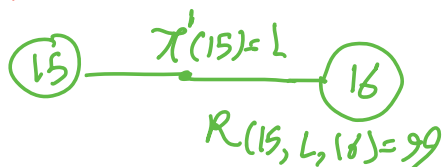
$$\pi' = \begin{bmatrix} \pi'(1) \\ \vdots \\ \vdots \\ \vdots \\ \pi'(15) \end{bmatrix} = \begin{bmatrix} L \\ L \\ U \\ \vdots \\ L \end{bmatrix}$$
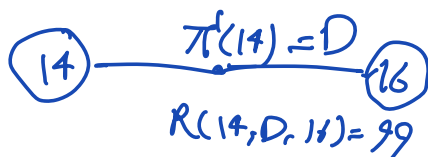
Deterministic

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

■ Wall   ▨ Bump   ▨ Goal

$V_{\pi'}(15) = 99$

15 —— $\pi'(15) = L$ —— 16

$R(15, L, 16) = 99$

$V_{\pi'}(14) = 99$

14 —— $\pi'(14) = D$ —— 16

$R(14, D, 16) = 99$

| | | | |
|---|---|---|---|
| 9 94 | 10 95 | 11 96 | 12 97 |
| 8 93 | ■ | 14 99 | 13 98 |
| 7 92 | 16 | 15 99 |
| 6 91 | 5 90 | ■ | ■ |
| 4 90 | 3 89 | 2 88 | 1 87 |

■ Wall   ▨ Bump   ▨ Goal

$V_{\pi'}(13) = -1 + 99 = 98$

13 —— $\pi'(13) = D$ —— 15 —— $\pi'(15)$ —— 16
         $-1$                      $99$

$V_{\pi'}(7) =$

7 —— $\pi'(7) = U$ —— 8 —— $\pi'(8) = 0$ —— 9 —— $\pi'(9) = R$ —— ◯ - - - -
        $-1$                  $-1$

**Grid 1 (top-left)**

| 9 86 | 10 87 | 11 88 | 12 97 |
|---|---|---|---|
| 8 85 | (Wall) | 14 99 | 13 98 |
| 7 . | (Wall) | 16 | 15 99 |
| 6 | 5 | (Wall) | (Wall) |
| 4 | 3 | 2 | 1 |

Wall   Bump   Goal

**Grid 2 (top-right)**

| 9 → | 10 → | 11 ↓ | 12 ↓ |
|---|---|---|---|
| 8 ↑ | (Wall) | 14 ↓ | 13 ↓ |
| 7 ↑ | (Wall) | 16 | 15 ↺ |
| 6 ↑ | 5 ↺ | (Wall) | (Wall) |
| 4 ↑ | 3 ↑ | 2 ↺ | 1 ← |

Wall   Bump   Goal

**Grid 3 (bottom-left)**

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

Wall   Bump   Goal

**Grid 4 (bottom-right)**

| 9 ↓ | 10 ↓ | 11 ↓ | 12 ↓ |
|---|---|---|---|
| 8 ↓ | (Wall) | 14 ↓ | 13 ↓ |
| 7 ↓ | (Wall) | 16 | 15 ↓ |
| 6 ↓ | 5 ↺ | (Wall) | (Wall) |
| 4 ← | 3 ↺ | 2 ↺ | 1 ↺ |

Wall   Bump   Goal

# State-Action Value Function

$$Q_\pi(S, a) = E\left[G_t \mid S_t = S, a_t = a, \pi\right]$$

$$\pi' = \begin{bmatrix} L \\ \tilde{U} \\ | \\ | \\ L \end{bmatrix}$$

| | | | |
|---|---|---|---|
| 9 → | 10 → | 11 ↱ | 12 ↓ |
| 8 ↑ | ■ | 14 ↓ | 13 ↓ |
| 7 ↑ | ■ | 16 | 15 ↰ |
| 6 ↑ | 5 ↰ | ■ | ■ |
| 4 ↑ | 3 ↑ | 2 ↰ | 1 ↰ |

■ Wall   🟧 Bump   🟩 Goal

$$Q_{\pi'}(15, U) =$$

$S_t$  $a_t = U$   $\pi'[13] = D$   $\pi'(15) = L$

(15) —•— (13) —•— (15) —•— (16)

  $-1$      $-1$       $99$

$$Q_{\pi'}(15, U) = -1 -1 + 99 = 97$$

$$Q_{\pi'}(15, R) = -1 + 99 = 98$$

$$(15) \xrightarrow[-1]{a_t = R} (15) \xrightarrow[99]{\pi^1(15) = L} (16)$$

$$Q_{\pi^1}(15, D) = -1 + 99 = 98$$

$$Q_{\pi^1}(15, L) =$$

$$(15) \xrightarrow{a_t = L} (16)$$

$$Q_{\pi}(S, \boxed{\pi(S)}) = V_{\pi}(S)$$

$$\boxed{A \mid B} \qquad \pi^1 = \begin{bmatrix} a1 \\ a1 \end{bmatrix} \qquad \pi^2 = \begin{pmatrix} a1 \\ a2 \end{pmatrix} \pi^3 \qquad \pi^4$$

$$V_{\pi^1}(A) \overset{0}{} \qquad V_{\pi^2}(A) \overset{50}{} \qquad V_{\pi^3}(A) \qquad V_{\pi^4}(A)$$

$$V_{\pi^1}(B)^{49} \qquad V_{\pi^2}(B)^{50} \qquad V_{\pi^3}(B) \qquad V_{\pi^4}(B)$$