

Lecture 7 - Feb 3, 2023

- Reinforcement Learning Preliminaries
 - State, Action, Reward, Policy
 - Returns and Expected Returns
 - State Value Function
 - State-Action Value Function
- Bellman Equation and Optimality
- Dynamic Programming
 - Policy Iteration
 - Value Iteration

Project 1 → Due Feb 7

TA's office hour:

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

$$MDP(S, A, R, T, \gamma)$$

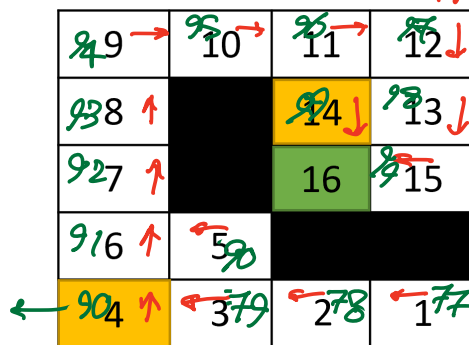
$$P(s'|s, a)$$

$$P(a) = \mathcal{M}(a) = \begin{bmatrix} \downarrow & \downarrow & \downarrow & \downarrow & \downarrow \\ \rightarrow & & & & \\ \rightarrow & \square & & & \\ \rightarrow & & & & \end{bmatrix}$$

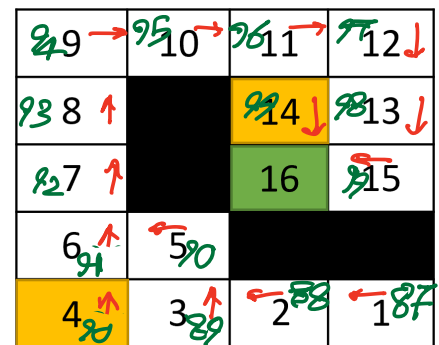
$$\pi = \begin{bmatrix} \pi(s^1) \\ \pi(s^2) \\ \vdots \\ \pi(s^N) \end{bmatrix} = \begin{bmatrix} a^1 \\ a^2 \\ a^1 \\ \vdots \\ a^1 \end{bmatrix}$$

$|A| \ll |S|$ diff Policy

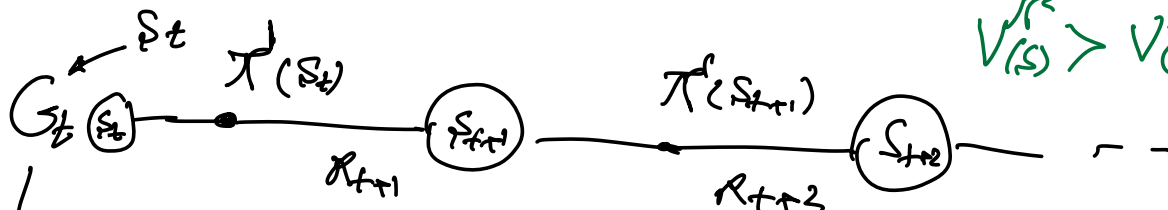
$\sqrt{\pi^2}$



Wall Bump Goal



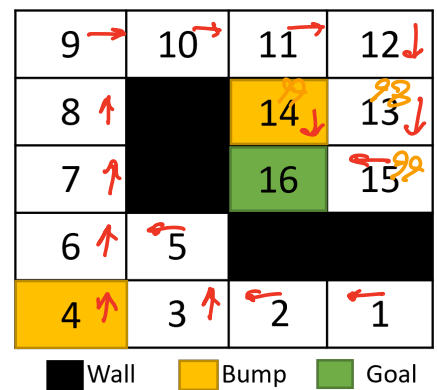
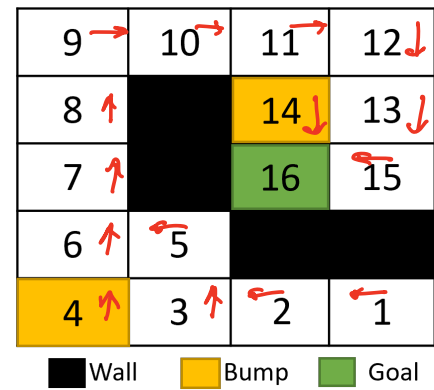
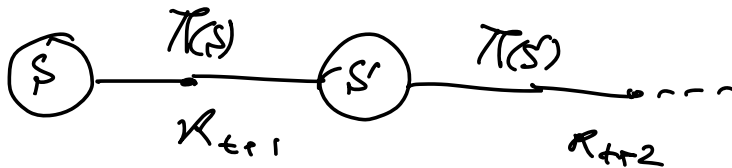
Wall Bump Goal



$$\sqrt{\pi^2} > \sqrt{\pi^1}$$

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

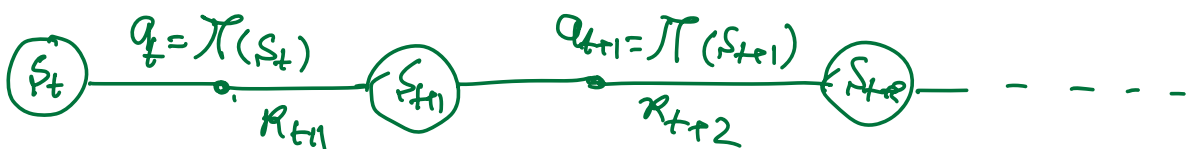
$$V_{\pi}^{\pi}(s) = E[G_t \mid s_t = s, \pi]$$



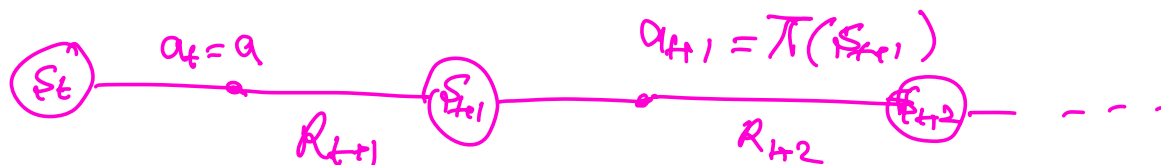
$\pi \geq \pi'$ if and only if

$V_{\pi}(s) \geq V_{\pi'}(s)$ for all $s \in \mathcal{S}$

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



$$Q_{\pi}(s, a) = E[G_t \mid S_t = s, a_t = a, \pi]$$



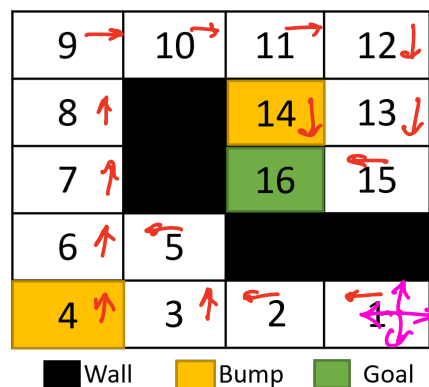
$$Q_{\pi}(S=15, a=U) = 97$$

$$15 \xrightarrow{U -1} 13 \xrightarrow{D -1} 15 \xrightarrow{R 99} 16$$

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

$$Q_{\pi}(15, D) = 98$$

$$Q_{\pi}(15, L) = 99$$



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

$$V_{\pi}(s) = E [G_t \mid s_t = s, a_t = \pi(s), a_{t+1:\infty} \sim \pi]$$

$$Q_{\pi}(s, a) = E [G_t \mid s_t = s, a_t = a, a_{t+1:\infty} \sim \pi]$$

V_{π^1}		V_{π^2}	
10	_____	10	
9	_____	10	
5	_____	9	\rightarrow It is not comparable
6	_____	4	
7	_____	7	
11	_____	12	

$$\Pi = \{ \pi^1, \pi^2, \dots, \pi^{4^{15}} \}$$

π^* optimal policy

$$V^{\pi^*}(s) \geq V^{\pi}(s) \quad \text{for all } s \quad \pi \in \Pi$$

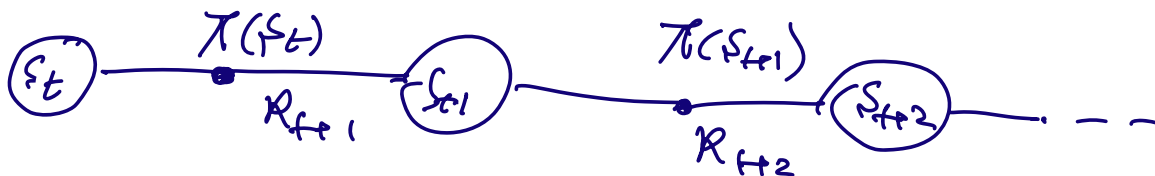
$$V_{\pi^*}(s) = V^*(s) = \max_{\pi \in \Pi} V_{\pi}(s), \text{ for all } s \in \mathcal{S}$$

$$Q_{\pi^*}(s, a) = Q^*(s, a) = \max_{\pi \in \Pi} Q_{\pi}(s, a) \text{ for all } s, a$$

Bellman Equation

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

Policy π



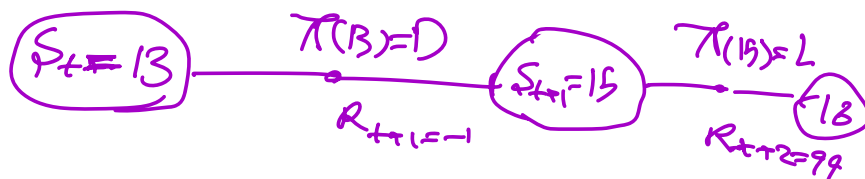
$$V_{\pi}(s) = E[G_t \mid S_t = s, \pi]$$

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

$$= E[R_{t+1} + \gamma \underbrace{(R_{t+2} + \gamma R_{t+3} + \dots)}_{G_{t+1}} \mid S_t = s, \pi]$$

$$V_{\pi}(s) = E[R_{t+1} + \gamma \underbrace{G_{t+1}}_{V_{\pi}(S_{t+1})} \mid S_t = s, \pi]$$

$$V_{\pi}(s) = E[R_{t+1} + \gamma V_{\pi}(S_{t+1}) \mid S_t = s, \pi] \quad \leftarrow \text{Bellman Eq.}$$



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

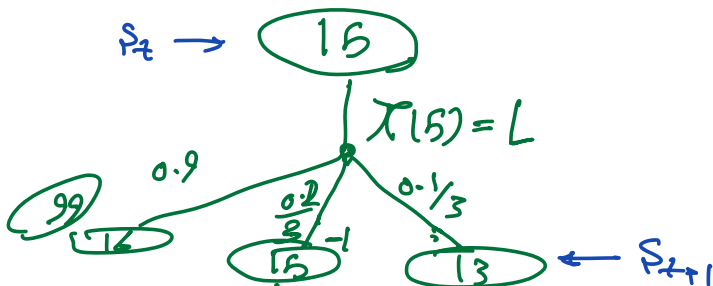
Wall
 Bump
 Goal

$$V_{\pi}(15) = -1 + 100 = 99$$

$$V_{\pi}(13) = E[R_{t+1} + \gamma V_{\pi}(S_{t+1}) | S_t = 13, \pi]$$

$$= -1 + \gamma V_{\pi}(15)$$

Stochastic Environment

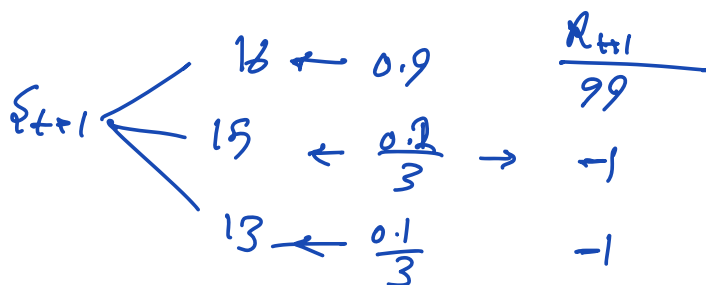


9 →	10 →	11 →	12 ↓
8 ↑		14 ↓	13 ↓
7 ↑		16	15 ←
6 ↑	5 ←		
4 ↑	3 ↑	2 ←	1 ←

Wall
 Bump
 Goal

$$V_{\pi}(s) = E[R_{t+1} + \gamma V_{\pi}(S_{t+1}) | S_t = s, \pi]$$

$$= P(s' | S_t = s, a_t = \pi(s_t))$$



$$E[G_t | S_t = s, \pi]$$

$$= \sum p^{s'aj} G_t^{s'aj}$$

$$V_{\pi}(s) = E [R_{t+1} + \gamma V_{\pi}(s_{t+1}) \mid s_t = s, \pi]$$

$$V_{\pi}(s) = \sum_{s'} P(s' \mid s, \pi(s)) [R(s, \pi(s), s') + \gamma V_{\pi}(s')]$$

$$V_{\pi}(15) = \underbrace{P(s'=16 \mid s=15, \pi(15)=L)}_{0.9} [\underbrace{R(15, L, 16)}_{99} + \gamma \underbrace{V_{\pi}(16)}_0]$$

$$+ \underbrace{P(s'=15 \mid s=15, \pi(15)=L)}_{\frac{0.2}{3}} [\underbrace{R(15, L, 15)}_{-1} + \gamma V_{\pi}(15)]$$

$$+ \underbrace{P(s'=13 \mid s=15, \pi(15)=L)}_{\frac{0.1}{3}} [\underbrace{R(15, L, 13)}_{-1} + \gamma V_{\pi}(13)]$$

Bellman Eq

$$V_{\pi}(s) = E [R_{t+1} + \gamma V_{\pi}(s') \mid s_t = s, \pi]$$

$$= \sum_{s'} P(s' \mid s, \pi(s)) [\underbrace{R_{t+1}}_{R(s, \pi(s), s')} + \gamma V_{\pi}(s')]$$

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

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

$$P(S' | S, a) \leftarrow \text{stochastic}$$

$$P(S' = D | S = B, a = 1) = 0.9$$

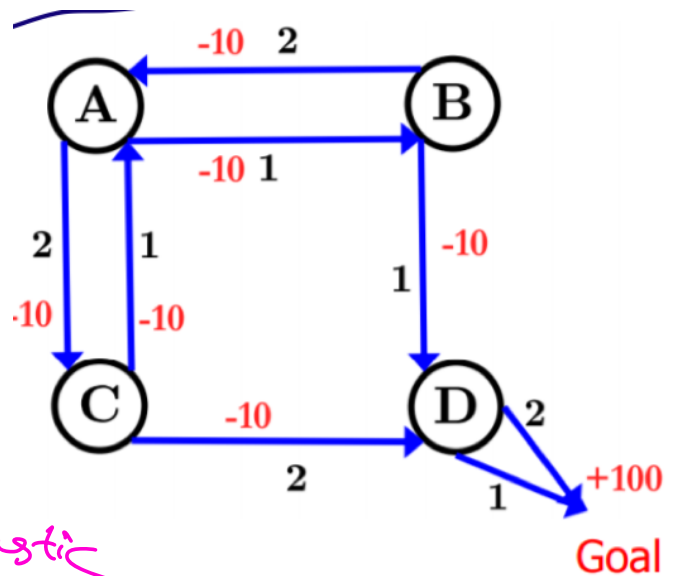
$$P(S' = A | S = B, a = 1) = 0.1$$

$$\pi = \begin{bmatrix} \pi(A) \\ \pi(B) \\ \pi(C) \\ \pi(D) \end{bmatrix}$$

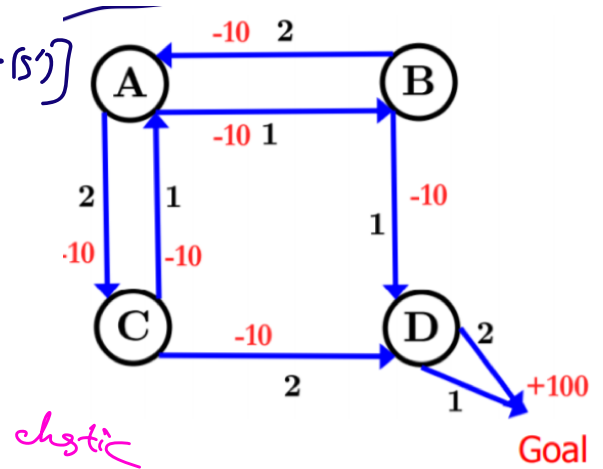
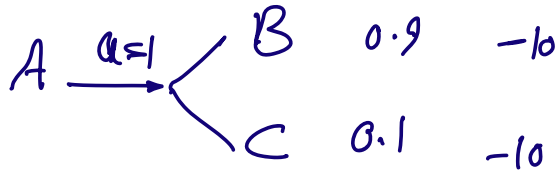
$$\begin{matrix} 1 & 1 & 1 & ? \\ 1 & ? & 1 & ? \\ 1 & ? & 2 & ? \\ 1 & 1 & 2 & 2 \end{matrix} \quad 2^4 = 16$$

$$\pi = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} \rightarrow$$

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



$$V_{\pi'}(A) = \sum_{s'} P(s' | s=A, \pi'(A)=1) [R + \gamma V_{\pi'}(s')]$$



$$V_{\pi'}(A) = \underbrace{P(s'=B | s=A, \pi'(A)=1)}_{0.9} [-10 + \gamma V_{\pi'}(B)] + \underbrace{P(s'=C | s=A, \pi'(A)=1)}_{0.1} [-10 + \gamma V_{\pi'}(C)]$$

$$V_{\pi'}(B) = 0.9 [-10 + \gamma V_{\pi'}(D)] + 0.1 [-10 + \gamma V_{\pi'}(A)]$$

$$V_{\pi'}(C) = 0.9 [-10 + \gamma V_{\pi'}(A)] + 0.1 [-10 + \gamma V_{\pi'}(D)]$$

$$V_{\pi'}(D) = 100$$

↳ Four Equations + Four Variables

$$V_{\pi'}(A) = x \quad V_{\pi'}(B) = y \quad V_{\pi'}(C) = z \quad V_{\pi'}(D) = d$$

$$x = 0.9 [-10 + 0.9 y] + 0.1 [-10 + 0.9 z]$$

$$y = -9 \quad \quad \quad$$

$$z \quad \quad \quad$$

$$d \quad \quad \quad$$

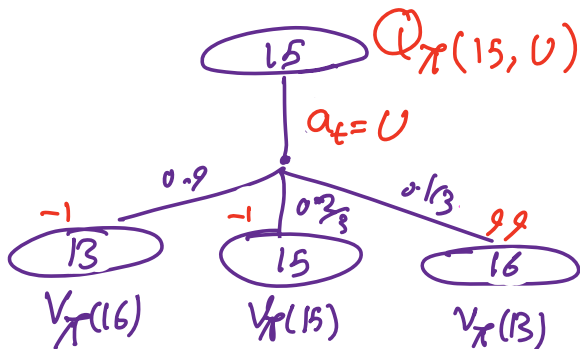
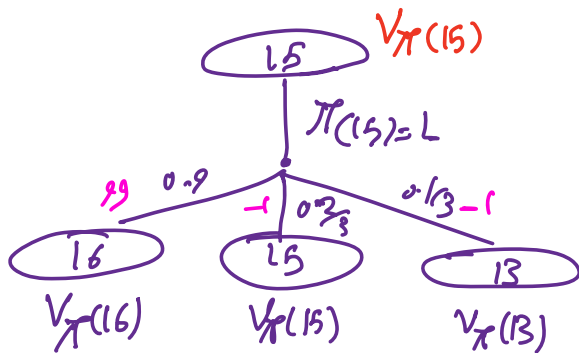
$$\begin{cases} x - 0.81y - 0.09z + 0d = -10 \\ -0.09x + y + 0z - 0.81d = -10 \\ -0.09x + 0y + z - 0.09d = -10 \\ 0x + 0y + 0z + 1d = 100 \end{cases}$$

$$A \begin{bmatrix} x \\ y \\ z \\ d \end{bmatrix} = a \rightarrow \begin{bmatrix} -10 \\ -10 \\ -10 \\ 100 \end{bmatrix} \rightarrow \boxed{A^{-1}a}$$

$$x = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} \rightarrow V_{x1} = \begin{bmatrix} 75 \\ 87 \\ 68 \\ 100 \end{bmatrix}$$

$$\pi^2 = \begin{bmatrix} 2 \\ 2 \\ 2 \\ 2 \end{bmatrix} \rightarrow V_{\pi^2} = \begin{bmatrix} 75 \\ 68 \\ 87 \\ 100 \end{bmatrix}$$

not comparable



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

Wall
 Bump
 Goal

$$Q_{\pi}(s, a) = E[G_t \mid s_t = s, a_t = a, \pi]$$

$$= E \left[\underbrace{R_{t+1}}_a + \underbrace{\gamma}_{\gamma} \underbrace{R_{t+2} + \gamma R_{t+3} + \dots}_{\pi} \mid S_t = s, a_t = a, \pi \right]$$

$$= E \left[R_{t+1} + \gamma \underbrace{(R_{t+2} + \gamma R_{t+3} + \dots)}_{\substack{G_{t+1} \\ \xrightarrow{\pi} \\ V_{\pi}(s_{t+1})}} \mid S_t = s, a_t = a, \pi \right]$$

$$= E \left[R_{t+1} + \gamma V_{\pi}(s_{t+1}) \mid S_t = s, a_t = a, \pi \right]$$

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