

Lecture 13 - Feb 24, 2023

- Dynamic programming

- Policy Iteration
- Value Iteration



- Approximate Dynamic Programming

- Asynchronous DP

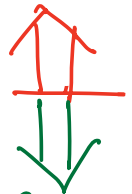
- * In-Place DP
- * Prioritized Sweeping
- * Real-Time DP

- Generalized Policy Iteration

- Monte-Carlo Methods

- First-Visit MC
- Online MC

Model-Based



Model-Free

HW3 is posted → Due March 17

Project 2 → Due March 3

TA's office hour:

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

Fridays, 2pm-3pm (virtual)

Approximate Dynamic Programming

Matrix-Form \times

$$S = \{10,000,000\}$$

$$\underbrace{M(a)}_{\begin{matrix} A \\ B \end{matrix}} \begin{matrix} A \\ B \end{matrix} \left[\begin{matrix} A \\ B \end{matrix} \right]$$

$$\underset{10,000,000}{V_k(s)} = \max_{a \in A} \sum_{s'} P(s'|s,a) \left[R(s,a,s') + \delta \underset{10,000,000}{V_{k-1}(s')} \right]$$

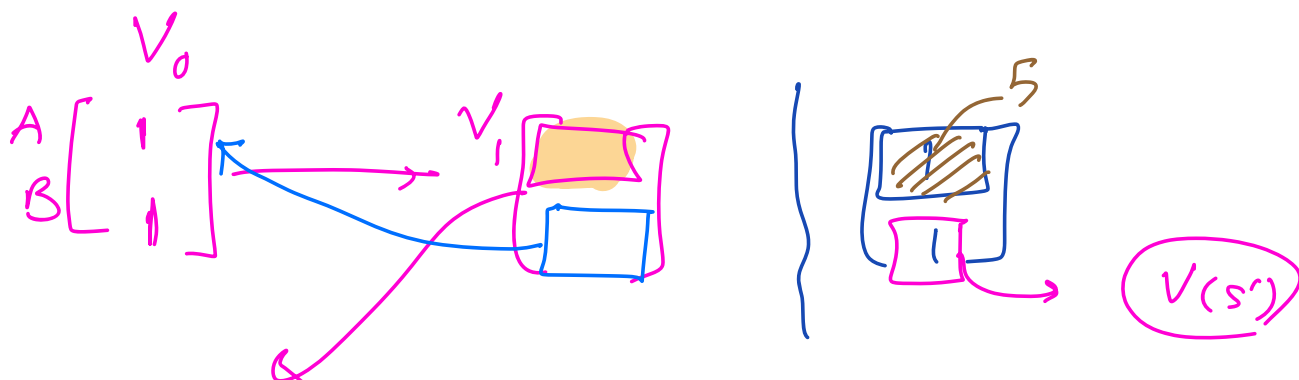
$$V_0 = \begin{bmatrix} 0 \\ c \\ \vdots \\ 0 \end{bmatrix} \xrightarrow{VIB} V_1 = \begin{bmatrix} \vdots \\ 0 \end{bmatrix}$$

10,000,000

In-place Dynamic Programming (Space-Efficient DP)

$$\begin{matrix} V_{old} \\ \uparrow \\ |S| \end{matrix} \quad \begin{matrix} V_{new} \\ \uparrow \\ |S| \end{matrix} \rightarrow \text{a single } V$$

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



$$V_1(s) = \max \sum P(s'|s, a) [R + \gamma V_0(s')]$$

$$V_1(s) = \max \sum P(s'|s, a) [R + \gamma V_0(s')]$$

$$V^*(s) = \arg \max_a \sum P(s'|s, a) [R + \gamma V^*(s')]$$

Approximate DP \Leftarrow Warren Powell

$$\pi \xrightarrow{PE} V^*$$

Prioritized Sweeping

Do I need to Backup for all $s \in \mathcal{S}$.

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

$$|V_{k+1}(s) - V_k(s)|$$

$\Delta \leftarrow$ priority

$$\Delta = \lfloor \cdot \rfloor$$

$$\frac{1}{|\mathcal{S}|}$$

$$s_k \sim \Delta$$

$$V^{\text{new}}(s_k) = \max_{a \in A} \sum_{s'} P(s' | s_k, a) [R + \gamma V^{\text{old}}(s')]$$

$$\Delta(s_k) = |V^{\text{new}}(s_k) - V^{\text{old}}(s_k)|$$

$$\Delta(s) = \Delta(s) + \gamma \max_{a \in A} P(s_k | s, a) \Delta(s_k) \leftarrow \begin{matrix} \text{for all } s \\ \text{that lead} \\ \text{to } s_k \end{matrix}$$

$$\Delta = [1 \ 1 \ 1 \ 1 \ \dots \ 1] \quad \Delta(s_{13})$$

$$S_K \sim [t_S, t_S - \dots]$$

$$S_k = 14$$

$$V(s'_k) = \max_{a \in A} \sum p(c|s'_k, a) [R + \gamma V(s')]^{\text{old}}$$

$$\max \left\{ \overbrace{-1+r_0}^{\text{up}}, \overbrace{99+r_0}^{\text{Down}}, \underbrace{-11}_L, \underbrace{-1}_R \right\} = 99$$

$$\Delta \overset{S_k}{(14)} = \left[\underset{99}{V^{new}(14)} - \underset{0}{V^{old}(14)} \right] = 99$$

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

 Wall
  Bump
  Goal

$$\Delta(11) = \Delta(11) + \gamma \max_{a \in A} \underbrace{P\left(\frac{S_K}{14} \mid \frac{S}{11}, a\right)}_{\text{Down } 1} \underbrace{\Delta(14)}_{\text{eg}} = 100$$

$$\Delta(13) = \underbrace{\Delta(13)}_1 + \max_{\substack{a \in A \\ \text{Left}}} \underbrace{P(14|13,a)}_1 \underbrace{\Delta(14)}_{22} = 100$$

$$\Delta \approx \begin{bmatrix} 1 & 1 & \frac{11}{100} & \frac{12}{1} & \frac{13}{100} & \frac{14}{99} & \frac{15}{1} \end{bmatrix}$$

$$S_k \sim \Delta$$

$$S_k = 13$$

$$V^{new}(13) = \max_{s'} \sum_{s'} P(s' | s_k, a) [R + \gamma V^{old}(s')]$$

$$= \max \{-1, -1, 88, -1\}$$

$$= 88$$

$$\Delta(S_k) = |V^{old}(13) - V^{new}(13)|$$

$$= 88$$

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

Wall
 Bump
 Goal

$$\Delta(12) = \Delta(12) + \gamma \max_{s'} P(13 | 12, a) \Delta(13)$$

$$= 1 + 1 \times 88 = 89$$

$$\Delta(15) = \Delta(15) + \gamma \max_{s'} P(13 | 15, a) \Delta(13) = 89$$

$$\Delta(14) = \Delta(14) + \gamma \max_{s'} P(13 | 14, a) \Delta(13) = 187$$

$$\Delta = \left[\frac{1}{1} \quad \frac{3}{1} \quad \frac{11}{100} \quad \frac{12}{89} \quad \frac{13}{88} \quad \frac{14}{187} \quad \frac{15}{89} \right]$$

$$S_k \sim \Delta$$

$$S_k = 11$$

$$V_{(11)}^{\text{new}} = \max\{-1, -1, -1, 88\} = 88$$

$$\Delta(11) = |0 - 88| = 88$$

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

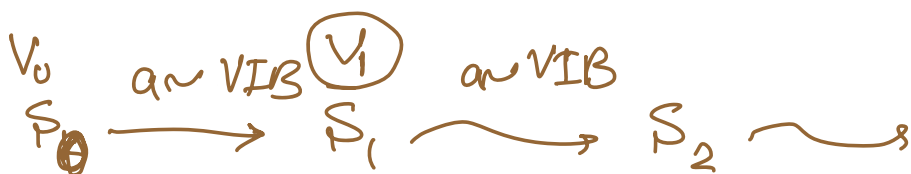
Wall
 Bump
 Goal

$$\Delta_F$$

Real-Time DP

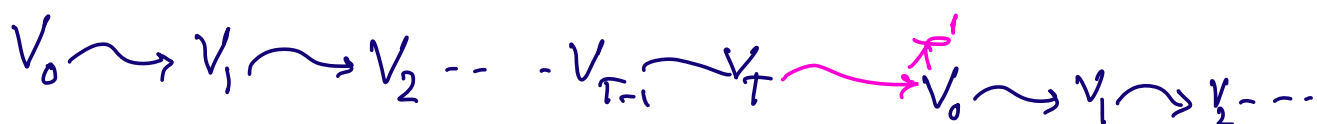
$$V(s_k) = \max_a \sum_{s'} P(s' | s_k, a) [R(s_k, a, s') + \gamma V(s')]$$

$$a_k = \operatorname{argmax}_a \sum_{s'} P(s' | s_k, a) [R + \gamma V(s')]$$



Generalized Policy Iteration (GPI)

$$\pi_0 \xrightarrow{\text{PE}} V_{\pi_0} \xrightarrow{\text{PI}} \pi_1 \xrightarrow{\text{PE}} \pi_2$$



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

$$\pi'(s) = \operatorname{argmax}_{a \in A} \sum_{s'} P(s' | s, a) [R + \gamma V^{\pi_0}(s')]$$

$$V_0 \rightsquigarrow V_1 \rightsquigarrow V_2 \text{ --- }$$

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

one-step PE + PI = Value Iteration

GPI \rightarrow

$$V_0 \rightarrow V_1 \rightarrow V_2 \xrightarrow{\text{PI}} V_3 \rightarrow V_4 \rightarrow V_5 \xrightarrow{\text{PI}}$$

Finger grain switch between PE & PI

Monte-Carlo Methods

$$MDP(S, A, R, \cancel{T})$$

$P(S'|S, a)$

$$S \xrightarrow[a]{a} S'$$

$P(S'|S, a)$

$$V_{k+1}(S) = \max_a \sum_{S'} \cancel{P(S'|S, a)} [R + \gamma V(S')]$$

$$V^\pi(S) = E \left[\underbrace{R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots}_{G_t} \mid S_0 = S, a \sim \pi \right]$$

$\begin{matrix} \nwarrow \pi(S_0) & \nwarrow \pi(S_1) & \nwarrow \pi(S_2) \\ R(S_0, a, S_1) & R(S_1, a, S_2) & R(S_2, \pi(S_2), S_3) \end{matrix}$

$$= E \left[\underbrace{R_{t+1} + \gamma (R_{t+2} + \gamma R_{t+3} + \dots)}_{V^\pi(S')} \mid S_0 = S, a \sim \pi \right]$$

$$= E [R_{t+1} + \gamma V^\pi(S')]$$

$$= \sum_{S'} P(S' \mid \pi(S), S) [R + \gamma V^\pi(S')]$$

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

$$\approx \frac{1}{N} \sum_{i=1}^N G_t^i$$

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

Wall
 Bump
 Goal

$$13 \xrightarrow[\substack{R=-1}]{\pi(13)} 15 \xrightarrow[\substack{99}]{\pi(15)} G \quad G_t^1 = -1 + 99 = 98$$

$$13 \xrightarrow[\substack{-11}]{\pi(13)} 14 \xrightarrow[\substack{99}]{\pi(14)} G \quad G_t^2 = -11 + 99 = 88$$