# Lecture 10 - Feb 14, 2023

- Dynamic Programming

    - Policy Iteration
    - Value Iteration     } Vector-Form

    - Policy Iteration
    - Value Iteration     } Matrix-Form

- Approximate Dynamic Programming
    - Asynchronous DP
    - Generalized Policy Iteration

Exam 1 ⟶ Tuesday, Feb 21

HW 2 ⟶ Due Feb 17

Project 2 ⟶ Due March 3

TA's office hour:    Wendsdays, 2pm-3pm (in-person)
                     Fridays, 2pm-3pm (virtual)

# Dynamic Programming

Bellman Eq $\longrightarrow$

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

Bellman Optimally Eq $\longrightarrow$

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

## Approach 1: Policy Iteration

Policy Evaluation (PE)

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

Policy Improvement

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

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

$$\pi_T = \pi_{T-1} = \pi^*$$

# Approach 2: Value Iteration (VI)

Value Iteration Backup (VIB)

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

for all $s \in S$

$$V_0 = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ \vdots \end{bmatrix} \longrightarrow V_1 = \begin{bmatrix} \\ \end{bmatrix} \longrightarrow V_2 \longrightarrow \cdots V_T$$

$$\max \|V_T - V_{T-1}\| < \theta \longrightarrow V_T = V^*$$

$$\pi^*(s) = \text{argmax}_a Q^*(s,a)$$

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

$$V_0 \longrightarrow V_1 \longrightarrow V_2 - - - - \quad V_T = V^* \longrightarrow \pi^*$$

## Proof.

$$V^*_{(\hat{s})} = \max_{a \in A} \sum_{s'} P(s' | \hat{s}', a) \left[ R(s', a, s) + \gamma V^*(s) \right]$$

$$\vdots$$

$$V^*_{(\hat{s})N} = \max_{a \in A} \sum_{s'} P(s' | \hat{s}', a) \left[ R(s', a, s) + \gamma V^*(s) \right]$$

$$V^* = TV^*$$

$$U_0 \xrightarrow{TU_0} U_1 \xrightarrow{TU_1} U_2 \; - \; - \; - \; \cdot \quad |U_T - U_{T-1}| \leq \epsilon$$

V and U are two random vectors

$$|TV_{(s)} - TU_{(s)}| \leq \gamma \| V - U \|_\infty$$

$$|TV_{(s)} - TU_{(s)}| = \left| \max_{a \in A} \sum_{s'} P(s' | s, a) \left[ R(s, a, s') + \gamma V(s) \right] \right.$$

$$\underbrace{\qquad\qquad\qquad\qquad}_{f(a)}$$

$$\left. - \max_{s'} \sum_{s'} P(s' | s, a) \left[ R(s, a, s') + \gamma U(s) \right] \right|$$

$$\underbrace{\qquad\qquad\qquad\qquad}_{g(a)}$$

**Lemma**

$$\left| \max_a f(a) - \max_a g(a) \right| \leq \max_a | f(a) - g(a) |$$

$$\leq \max_{a \in A} \left| \gamma \sum_{s'} P(s' | s, a) \left[ V(s') - U(s') \right] \right|$$

$\hookrightarrow$ max happens at $a^s$

$$= \gamma \left| \sum_{s'} P(s' | s, a^s) \left[ V(s') - U(s') \right] \right|$$

$$\leq \gamma \| V - U \|_\infty$$

Example:

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

$$M(a^1) = \begin{array}{c} A \\ B \end{array} \begin{bmatrix} \overset{A}{0.9} & \overset{B}{0.1} \\ 0.1 & 0.9 \end{bmatrix}$$

Reward $\begin{cases} +5 & \text{ending at } B \\ -1 & \text{for } a^2 \end{cases}$

$$M(a^2) = \begin{array}{c} A \\ B \end{array} \begin{bmatrix} \overset{A}{0.1} & \overset{B}{0.9} \\ 0.9 & 0.1 \end{bmatrix}$$

$\gamma = 0.9$

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

$$V_0 = \begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} V_0(A) \\ V_0(B) \end{bmatrix} \xrightarrow{\text{VIB}} V_1 = \begin{bmatrix} V_1(A) \\ V_1(B) \end{bmatrix}$$

$$V_1(A) = \max_{a \in \{a^1, a^2\}} \sum_{s'} P(s' \mid A, a) \left[ R(A, a, s') + \gamma V_0(s') \right]$$

$$= \max \begin{cases} \underbrace{P(A \mid A, a^1)}_{0.9} \overbrace{\left[ \underbrace{R(A, a^1, A)}_{0} + \gamma \underbrace{V_0(A)}_{} \right]}^{0.5} + \underbrace{P(B \mid A, a^1)}_{0.1} \left[ \underbrace{R(A, a^1, B)}_{5} + \gamma \underbrace{V_0(B)}_{0} \right] & \xleftarrow{a^1} \end{cases}$$

$$\underbrace{P(B \mid A, a^2)}_{0.9} \left[ \underbrace{R(A, a^2, B)}_{4} + \gamma \overset{0}{V_0(B)} \right] + \underbrace{P(A \mid A, a^2)}_{0.1} \left[ \underbrace{R(A, a^2, A)}_{-1} + \gamma \overset{0}{V_0(A)} \right] \xleftarrow{a^2}$$

$$\underbrace{\qquad \qquad \qquad \qquad 3.6 \qquad \qquad \qquad \qquad}$$

$$= 3.5$$

$$V_1 = \begin{bmatrix} V_1(A) \\ V_1(B) \end{bmatrix} = \begin{bmatrix} 3.5 \\ \end{bmatrix}$$

$$V_1(B) = \max_{a \in A} \sum_{s'} P(s' | B, a) [R(B, a, s') + \gamma V_0(s')]$$

$$= \max \left\{ \underbrace{P(\cdot B | B, a')}_{0.9} [\underbrace{R(B, a', B)}_{5} + \overset{\overset{0}{\frown}}{\gamma V_0(B)}] + \underbrace{P(A | B, a')}_{0.1} [\overset{0}{\overbrace{R(B, a', A)}} + \gamma V_0(A)] \right.$$

$$\overbrace{\qquad\qquad\qquad 4.5 \qquad\qquad\qquad}^{} \qquad a'$$

$$\underbrace{P(A | B, a^2)}_{0.9} [\underbrace{R(B, a^2, A)}_{-1} + \gamma \overset{0}{\overbrace{V_0(A)}}] + \underbrace{P(B | B, a^2)}_{0.1} [\underbrace{R(B, a^2, B)}_{4} + \gamma \overset{0}{\overbrace{V_0(B)}}] \right\}$$

$$\underbrace{\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad}_{-0.5} \qquad a_2$$

$$= 4.5$$

$$V_1 = \begin{bmatrix} 3.5 \\ 4.5 \end{bmatrix} \longrightarrow V_2 \longrightarrow V_3 \longrightarrow V_\top$$

$$V^* = V_{100} = \begin{bmatrix} 43.1 \\ 44.1 \end{bmatrix}$$

$$\max |V_{100} - V_{99}| < \overbrace{\gamma}^{0.01} \checkmark\checkmark\checkmark \quad V_{100} = V^*$$

$$\pi^* = \begin{bmatrix} \pi^*(A) \\ \pi^*(B) \end{bmatrix} \longleftarrow V^*$$

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

$$V^* = \begin{bmatrix} 43.1 \\ 44.1 \end{bmatrix}$$

$$\pi^*(A) = \operatorname*{argmax}_{a \in \{a^1, a^2\}} \left\{ \overbrace{\underbrace{P(A|A,a^1)}_{0.9}\left[\underbrace{R(A,a^1,A)}_{0} + \gamma \underbrace{V^*(A)}_{43.1}\right] + \underbrace{P(B|A,a^1)}_{0.1}\left[\underbrace{R(A,a^1,B)}_{5} + \gamma \underbrace{V^*(B)}_{44.1}\right]}^{39.38 \quad a^1} \right.$$

$$\left. \overbrace{\underbrace{P(B|A,a^2)}_{0.9}\left[\underbrace{R(A,a^2,B)}_{4} + \gamma \underbrace{V^*(B)}_{44.1}\right] + \underbrace{P(A|A,a^2)}_{0.1}\left[\underbrace{R(A,a^2,A)}_{-1} + \gamma \underbrace{V^*(A)}_{43.1}\right]}^{43.1 \quad a^2} \right\} = a^2$$

$$\pi^*(B) = \operatorname*{argmax}_{a \in \{a^1, a^2\}} \left\{ \overbrace{\underbrace{P(B|B,a^1)}_{0.9}\left[\underbrace{R(B,a^1,B)}_{5} + \gamma \underbrace{V^*(B)}_{44.1}\right] + \underbrace{P(A|B,a^1)}_{0.1}\left[\underbrace{R(B,a^1,A)}_{0} + \gamma \underbrace{V^*(A)}_{43.1}\right]}^{44.1 \quad a^1} \right.$$

$$\left. \overbrace{\underbrace{P(A|B,a^2)}_{0.9}\left[\underbrace{R(B,a^2,A)}_{-1} + \gamma \underbrace{V^*(A)}_{43.1}\right] + \underbrace{P(B|B,a^2)}_{0.1}\left[\underbrace{R(B,a^2,B)}_{4} + \gamma \underbrace{V^*(B)}_{44.1}\right]}^{38.38} \right\} = a^1$$

## Value Iteration

Initialize array $V$ arbitrarily (e.g., $V(s) = 0$ for all $s \in \mathcal{S}^+$)

Repeat
$\quad \Delta \leftarrow 0$
$\quad$ For each $s \in \mathcal{S}$:
$\quad\quad v \leftarrow V(s)$
$\quad\quad V(s) \leftarrow \max_a \sum_{s',r} p(s',r|s,a)\left[r + \gamma V(s')\right]$
$\quad\quad \Delta \leftarrow \max(\Delta, |v - V(s)|)$
until $\Delta < \theta$ (a small positive number)

Output a deterministic policy, $\pi$, such that
$\quad \pi(s) = \arg\max_a \sum_{s',r} p(s',r|s,a)\left[r + \gamma V(s')\right]$

$V_0$

| | | | |
|---|---|---|---|
| 9 $^0$ | 10 $^0$ | 11 $^0$ | 12 $^0$ |
| 8 $^0$ | ■ | 14 $^{0}$ | 13 $^0$ |
| 7 $^0$ | ■ | 16 $^0$ | 15 $^0$ |
| 6 $^0$ | 5 $^0$ | ■ | ■ |
| 4 $^0$ | 3 $^0$ | 2 $^0$ | 1 $^0$ |

■ Wall   ▧ Bump   ▨ Goal

VIB $\rightarrow$

$V_1$

| | | | |
|---|---|---|---|
| 9 $^{-1}$ | 10 $^{-1}$ | 11 $^{-1}$ | 12 $^{-1}$ |
| 8 $^{-1}$ | ■ | 14 $^{99}$ | 13 $^{-1}$ |
| 7 $^{-1}$ | ■ | 16 | 15 $^{99}$ |
| 6 $^{-1}$ | 5 $^{-1}$ | ■ | ■ |
| 4 $^{-1}$ | 3 $^{-1}$ | 2 $^{-1}$ | 1 $^{-1}$ |

■ Wall   ▧ Bump   ▨ Goal

$$V_1(s) = \max_{a \in A} \sum_{s'} P(s'|s,a)\left[R(s,a,s') + \gamma V_0(s')\right]$$

$s=15 \rightarrow V_1(15) = \max_{a \in \{U\,L\,D\,R\}}$

$\overbrace{P(13|15,U)}^{1}\left[\overbrace{R(15,U,13)}^{-1} + \gamma \overbrace{V_0(13)}^{0}\right] = -1$   U

$\overbrace{P(Goal|15,L)}^{1}\left[\overbrace{R(15,L,Goal)}^{99} + \gamma \overbrace{V_0(Goal)}^{0}\right] = 99$   L

$\overbrace{P(15|15,D)}^{1}\left[\overbrace{R(15,D,15)}^{-1} + \gamma \overbrace{V_0(15)}^{0}\right] = -1$   D

  R    $\left.\begin{array}{c} =-1 \\ =99 \end{array}\right\}$

$V_1$

| 9 $^{-1}$ | 10 $^{-1}$ | 11 | 12 $^{-1}$ |
|---|---|---|---|
| 8 $^{-1}$ | Wall | 14 | 13 $^{-1}$ |
| 7 $^{-1}$ | Wall | 16 | 15 $^{99}$ |
| 6 $^{-1}$ | 5 $^{-1}$ | Wall | Wall |
| 4 | 3 $^{-1}$ | 2 | 1 $^{-1}$ |

$V \to -1 + (-1) = -2$
$D \to -1 + 99 = 98$
$R \to -1 + 99 = 98$
$L \to 99$

Wall  Bump  Goal

$V_2$

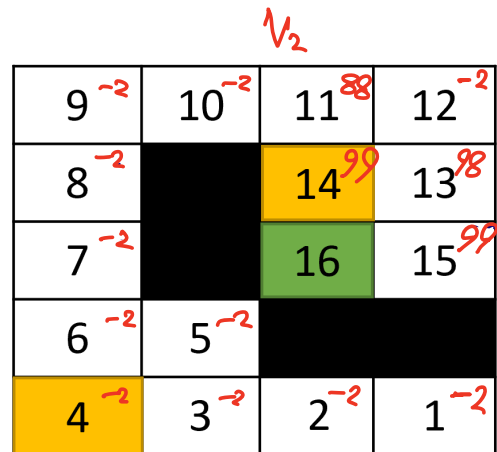| 9 $^{-2}$ | 10 $^{-2}$ | 11 $^{88}$ | 12 $^{-2}$ |
|---|---|---|---|
| 8 $^{-2}$ | Wall | 14 $^{99}$ | 13 $^{98}$ |
| 7 $^{-2}$ | Wall | 16 | 15 $^{99}$ |
| 6 $^{-2}$ | 5 $^{-2}$ | Wall | Wall |
| 4 $^{-2}$ | 3 $^{-2}$ | 2 $^{-2}$ | 1 $^{-2}$ |

Wall  Bump  Goal

$V_2$

| 9 $^{-2}$ | 10 $^{-2}$ | 11 $^{88}$ | 12 $^{-2}$ |
|---|---|---|---|
| 8 $^{-2}$ | Wall | 14 $^{99}$ | 13 $^{98}$ |
| 7 $^{-2}$ | Wall | 16 | 15 $^{99}$ |
| 6 $^{-2}$ | 5 $^{-2}$ | Wall | Wall |
| 4 $^{-2}$ | 3 $^{-2}$ | 2 $^{-2}$ | 1 $^{-2}$ |

Wall  Bump  Goal

$V_3$

| 9 $^{-2}$ | 10 $^{87}$ | 11 $^{88}$ | 12 $^{97}$ |
|---|---|---|---|
| 8 $^{-2}$ | Wall | 14 $^{99}$ | 13 $^{98}$ |
| 7 $^{-2}$ | Wall | 16 | 15 $^{99}$ |
| 6 $^{-2}$ | 5 $^{-2}$ | Wall | Wall |
| 4 $^{-2}$ | 3 $^{-2}$ | 2 $^{-2}$ | 1 $^{-2}$ |

8

Wall  Bump  Goal

$V_4$

| 9 $^{86}$ | 10 $^{87}$ | 11 $^{96}$ | 12 $^{97}$ |
|---|---|---|---|
| 8 $^{-2}$ | Wall | 14 $^{99}$ | 13 $^{98}$ |
| 7 $^{-2}$ | Wall | 16 | 15 $^{99}$ |
| 6 $^{-2}$ | 5 $^{-2}$ | Wall | Wall |
| 4 $^{-2}$ | 3 $^{-2}$ | 2 $^{-2}$ | 1 $^{-2}$ |

8

Wall  Bump  Goal

$V_{13} = V^*$

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

Wall  Bump  Goal

**First grid (top-left):** $V^*$

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

Wall   Bump   Goal

**Second grid (top-right):** $\pi^*$

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

Wall   Bump   Goal

$$\pi^*_{(s)} = \arg\max_a \sum_{s'} P(s' \mid s, a)\left[R \to + \gamma V^*_{(s')}\right]$$

**Third grid (bottom):**

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

Wall   Bump   Goal