

Lecture 9 - Feb 10, 2023

- Dynamic Programming

- ↓
- Policy Iteration
 - Value Iteration
- } Vector-Form
-
- Policy Iteration
 - Value Iteration
- } Matrix-Form

HW2 → Due Feb 17

Project 2 will be assigned → Due March 3

TA's office hour:

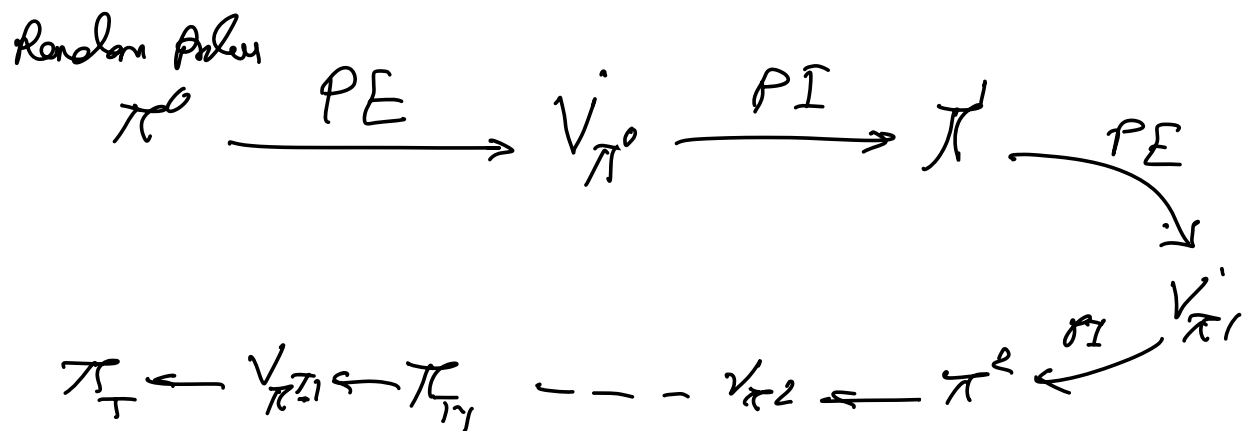
- Wednesdays, 2pm-3pm (in-person)
- Fridays, 2pm-3pm (virtual)

Policy Iteration

$$\pi^* = \arg \max_{\pi \in \Pi} V_{\pi}(s) \text{ for all } s$$

- ① Policy Evaluation (PE)
- ② Policy Improvement (PI)

Random Policy



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

PE: computes V_{π} for any given π

→ Bellman Eq

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

① Solving Bellman Eqs $\rightarrow \begin{cases} N \text{ Equations} \\ N \text{ Variables} \end{cases}$

(Small state space)

② $V_0 \leftarrow \text{random}$

$$V_0 = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}_{N \times 1}$$

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

$$\begin{bmatrix} V_0 \\ 0 \\ 0 \\ 0 \\ \vdots \end{bmatrix} \rightarrow \begin{bmatrix} V_1 \\ 0.1 \\ 0.2 \\ 0.3 \\ \vdots \end{bmatrix} \dots \begin{bmatrix} V_t \\ \vdots \end{bmatrix} \begin{bmatrix} V_{t+1} \\ \vdots \end{bmatrix}$$

$$\max |V_{t+1} - V_t| = \|V_{t+1} - V_t\|_{\infty} < \theta$$

Bellman Eq \leftarrow vector form

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

$$\vdots$$

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

$$\hookrightarrow \overline{T}^{\pi}(V)$$

$$\overline{T}^{\pi}(V)(s') = \sum_{s'} P(s'|s', \pi(s')) [R(s', \pi(s'), s) + \gamma V(s')]$$

$$\vdots$$

$$\overline{T}^{\pi}(V)(s^N) = \sum_{s'} P(s'|s^N, \pi(s^N)) [R(s^N, \pi(s^N), s) + \gamma V(s')]$$

$$V = \overline{T}^{\pi} V$$

\hookrightarrow Contraction-mapping Theorem

$$\|\overline{T}^{\pi}(V) - \overline{T}^{\pi}(U)\|_{\infty} \leq \gamma \|V - U\|_{\infty}$$

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

\vdots

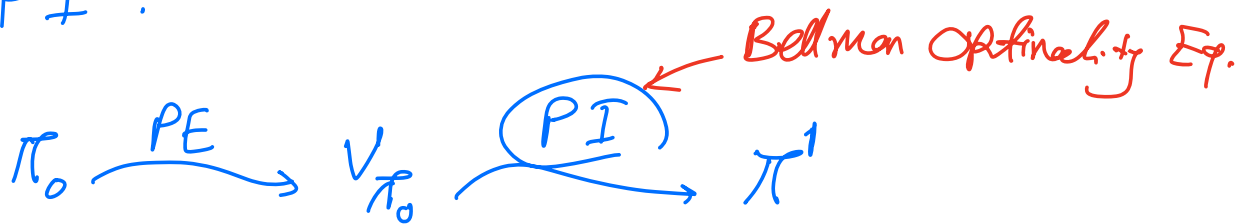
$$\sum_{s'} P(s'|s, \pi(s)) [R(s, \pi(s), s') + \gamma U(s')]$$

\vdots

$$\|\overline{T}^{\pi}(V) - \overline{T}^{\pi}(U)\|_{\infty} = \left\| \gamma \sum_{s'} P(s'|s, \pi(s)) [V(s') - U(s')] \right\|_{\infty}$$

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

PI :

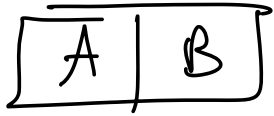


$$V^*(s) = \max_{a \in A} \underbrace{\sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V^*(s')]}_{Q^*(s, a)}$$

Policy Improvement

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

Example:



$$M(a^1) = \begin{matrix} & \begin{matrix} A & B \end{matrix} \\ \begin{matrix} A \\ B \end{matrix} & \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \end{matrix}$$

$$R_{\text{reward}} \rightarrow \begin{matrix} & B \\ +5 & \\ -1 & a^2 \end{matrix}$$

$$M(a^2) = \begin{matrix} & \begin{matrix} A & B \end{matrix} \\ \begin{matrix} A \\ B \end{matrix} & \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \end{matrix}$$

$$\Pi = \left\{ \pi^1 = \begin{bmatrix} a^1 \\ a^1 \end{bmatrix}, \pi^2 = \begin{bmatrix} a^1 \\ a^2 \end{bmatrix}, \pi^3 = \begin{bmatrix} a^2 \\ a^1 \end{bmatrix}, \pi^4 = \begin{bmatrix} a^2 \\ a^2 \end{bmatrix} \right\}$$

Policy Iteration:

$$\text{Random Policy } \pi^0 = \begin{bmatrix} a^1 \\ a^2 \end{bmatrix}, \quad \theta = 0.1, \quad \gamma = 0.9$$

$$\pi^0 \xrightarrow{\text{PE}} \underbrace{V^{\pi^0}} \longrightarrow \pi^1 \longrightarrow \underbrace{V^{\pi^1}} \longrightarrow \dots$$

$$V_0 = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \leadsto V_{k+1}(s) = \sum_{s'} P(s'|s, \pi^0(s)) [R + \gamma V_k(s')] \\ \downarrow \\ \underline{V_1 = \begin{bmatrix} \\ \end{bmatrix}}$$

$$V_1(A) = \sum_{s'} P(s'|A, \pi^0(A)=a^1) [R(A, a^1, s') + \gamma V_0(s')] \\ = P(A|A, \pi^0(A)=a^1) [R(A, a^1, A) + \gamma V_0(A)] = 0$$

$$\underbrace{\quad}_1 \quad \underbrace{\quad}_0 \quad \underbrace{\quad}_0$$

$$\begin{aligned} V_1(B) &= \sum_{s'} P(s' | B, \pi^0(B) = a^2) [R(B, a^2, s') + \gamma V_0(s')] \\ &= \underbrace{P(A | B, \pi^0(B) = a^2)}_1 \left[\underbrace{R(B, a^2, A)}_{-1} + \gamma \underbrace{V_0(A)}_0 \right] = -1 \end{aligned}$$

$$V_1 = \begin{bmatrix} 0 \\ -1 \end{bmatrix} \quad \max_{\substack{\downarrow \\ [-1]}} \| \underset{\substack{\downarrow \\ [-1]}}{V_1} - \underset{\substack{\downarrow \\ [-1]}}{V_0} \| = 1 < \theta^{0.1} \quad \times \quad \downarrow \text{Continue}$$

$$\begin{aligned} V_2(A) &= \sum_{s'} P(s' | A, \pi^0(A) = a^1) [R(A, a^1, s') + \gamma V_1(s')] \\ &= \underbrace{P(A | A, a^1)}_1 \left[\underbrace{R(A, a^1, A)}_0 + \gamma \underbrace{V_1(A)}_0 \right] = 0 \end{aligned}$$

$$\begin{aligned} V_2(B) &= \sum_{s'} P(s' | B, \pi^0(B) = a^2) [R(B, a^2, s') + \gamma V_1(s')] \\ &= \underbrace{P(A | B, a^2)}_1 \left[\underbrace{R(B, a^2, A)}_{-1} + \gamma \underbrace{V_1(A)}_0 \right] = -1 \end{aligned}$$

$$V_2 = \begin{bmatrix} 0 \\ -1 \end{bmatrix} \quad \max_{\substack{\downarrow \\ [-1]}} \| \underset{\substack{\downarrow \\ [-1]}}{V_2} - \underset{\substack{\downarrow \\ [-1]}}{V_1} \| = 0 < \frac{0.1}{\theta} \quad \checkmark \quad \text{stop}$$

$$\underbrace{V_2 = V_{\pi^0}} \quad V_{\pi^0} = \begin{bmatrix} V_{\pi^0}(A) \\ V_{\pi^0}(B) \end{bmatrix} = \begin{bmatrix} 0 \\ -1 \end{bmatrix}$$

$$PI: \pi'(s) = \underset{a \in A}{\operatorname{argmax}} \sum_{s'} P(s' | s, a) [R(s, a, s') + \gamma V_{\pi}(s')] \quad \text{for all } s$$

$$\pi^0 = \begin{bmatrix} a^1 \\ a^2 \end{bmatrix} \quad V_{\pi^0} = \begin{bmatrix} 0 \\ -1 \end{bmatrix}$$

$$\pi^1(A) = \underset{a \in \{a^1, a^2\}}{\operatorname{argmax}} \left\{ \sum_{s'} P(s' | A, a) [R(A, a, s') + \gamma V_{\pi^0}(s')] \right\}$$

$$= \overbrace{P(A | A, a^1)}^1 [\overbrace{R(A, a^1, A)}^0 + \gamma \overbrace{V_{\pi^0}(A)}^0]$$

$$\pi^1(A) = \underset{a^1}{\operatorname{argmax}} \left\{ \underbrace{\sum_{s'} P(s' | A, a^1) [R(A, a^1, s') + \gamma V_{\pi^0}(s')]}_{a^1} \right\} = 0$$

$$\overbrace{P(B | A, a^2)}^1 [\overbrace{R(A, a^2, B)}^4 + \gamma \overbrace{V_{\pi^0}(B)}^{\frac{-1}{\gamma}}]$$

$$\underbrace{\sum_{s'} P(s' | A, a^2) [R(A, a^2, s') + \gamma V_{\pi^0}(s')]}_{a^2} = 3 \cdot 1 = a^2$$

$$\pi^1(B) = \operatorname{argmax}_{a \in A} \left\{ \sum_{s'} P(s' | B, a) [R + \gamma V_{\pi^0}(s')] \right\}$$

$$\pi^1(B) = \operatorname{argmax} \left\{ \underbrace{P(B | B, a^1) [R(B, a^1, B) + \gamma V_{\pi^0}(B)]}_{a^1} = 4.1 \right\}$$

$$\underbrace{P(A | B, a^2) [R(B, a^2, A) + \gamma V_{\pi^0}(A)]}_{a^2} = a^1$$

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

$$\text{PE} \rightarrow V_{\pi^1}$$

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

Bedman Eqs

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

$$V_{\pi^1}(s) = \sum_{s'} P(s' | s, \pi^1(s)) [R(s, \pi^1(s), s') + \gamma V_{\pi^1}(s')] \quad \text{for all } s$$

$\left. \begin{array}{l} N \text{ equations} \\ N \text{ variables} \end{array} \right\} \text{linear}$

$$\begin{aligned} V_{\pi^1}(A) &= \sum_{s'} P(s' | A, \pi^1(A) = a^2) [R(A, a^2, s') + \gamma V_{\pi^1}(s')] \\ &= \underbrace{P(B | A, a^2)}_1 \left[\underbrace{R(A, a^2, B)}_4 + \gamma \underbrace{V_{\pi^1}(B)} \right] \end{aligned}$$

$$V_{\pi^1}(A) - 0.9 V_{\pi^1}(B) = 4 \quad (I)$$

$$\begin{aligned} V_{\pi^1}(B) &= \sum_{s'} P(s' | B, \pi^1(B) = a^1) [R(B, a^1, s') + \gamma V_{\pi^1}(s')] \\ &= \underbrace{P(B | B, a^1)}_1 \left[\underbrace{R(B, a^1, B)}_5 + \underbrace{\gamma}_{0.9} \underbrace{V_{\pi^1}(B)} \right] \end{aligned}$$

$$0.1 V_{\pi^1}(B) = 5 \rightarrow V_{\pi^1}(B) = 50$$

$\rightarrow (I) \quad V_{\pi^1}(A) = 49$

$$\pi^0 = \begin{bmatrix} a^1 \\ a^2 \end{bmatrix} \rightsquigarrow V_{\pi^0} = \begin{bmatrix} 0 \\ -1 \end{bmatrix} \rightsquigarrow \pi^1 = \begin{bmatrix} a^2 \\ a^1 \end{bmatrix} \quad V_{\pi^1} = \begin{bmatrix} 49 \\ 50 \end{bmatrix}$$

π^2 PI

$$\pi^2(A) = \underset{a \in \{a^1, a^2\}}{\operatorname{argmax}} \left\{ \sum_{s'} P(s' | A, a) [R(A, a, s') + \gamma V_{\pi^1}(s')] \right\}$$

$$= \overbrace{P(A | A, a^1)}^1 [\overbrace{R(A, a^1, A)}^0 + \gamma \overbrace{V_{\pi^1}(A)}^{49}] = 44$$

$$\pi^2(A) = \underset{a \in \{a^1, a^2\}}{\operatorname{argmax}} \left\{ \sum_{s'} P(s' | A, a) [R(A, a, s') + \gamma V_{\pi^1}(s')] \right\} = 0$$

$$\overbrace{P(B | A, a^2)}^1 [\overbrace{R(A, a^2, B)}^4 + \gamma \overbrace{V_{\pi^1}(B)}^{50}]$$

$$\underbrace{\sum_{s'} P(s' | A, a^2) [R(A, a^2, s') + \gamma V_{\pi^1}(s')]}_{Q2} = 49 = a^2$$

$$\pi^2(B) = \operatorname{argmax}_{a \in A} \left\{ \sum_{s'} P(s' | B, a) [R + \gamma V_{\pi^1}(s')] \right\}$$

$$\pi^2(B) = \operatorname{argmax}_x \left\{ \underbrace{P(B | B, a^1)}_1 \left[\underbrace{R(B, a^1, B)}_5 + \gamma \underbrace{V_{\pi^1}(B)}_{50} \right] \right\} = 50$$

a^1

$$\underbrace{P(A | B, a^2)}_1 \left[\underbrace{R(B, a^2, A)}_{-1} + \gamma \underbrace{V_{\pi^1}(A)}_{49} \right] = 44$$

a^2

$= a^1$

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

$$\pi^1 = \pi^2 \checkmark \checkmark \checkmark = \pi^*$$

Policy Iteration (using iterative policy evaluation) for estimating $\pi \approx \pi_*$

1. Initialization

$V(s) \in \mathbb{R}$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in \mathcal{S}$

2. Policy Evaluation

Loop:

$\Delta \leftarrow 0$

Loop for each $s \in \mathcal{S}$:

$v \leftarrow V(s)$

$V(s) \leftarrow \sum_{s',r} p(s',r|s,\pi(s)) [r + \gamma V(s')]$

$\Delta \leftarrow \max(\Delta, |v - V(s)|)$

until $\Delta < \theta$ (a small positive number determining the accuracy of estimation)

3. Policy Improvement

policy-stable \leftarrow *true*

For each $s \in \mathcal{S}$:

old-action $\leftarrow \pi(s)$

$\pi(s) \leftarrow \operatorname{argmax}_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$

If *old-action* $\neq \pi(s)$, then *policy-stable* \leftarrow *false*

If *policy-stable*, then stop and return $V \approx v_*$ and $\pi \approx \pi_*$; else go to 2

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

π' is better than π

$$V_{\pi'} \geq V_{\pi} \Rightarrow V_{\pi'}(s) \geq V_{\pi}(s) \text{ for all } s$$

$$V_{\pi}(s) \leq Q_{\pi}(s, \pi'(s))$$

$$= E \left[\overbrace{R_{t+1} + \gamma V_{\pi}(s_{t+1})}^{Q_{\pi}(s_{t+1}, \pi'(s_{t+1}))} \mid s_t = s, a_t = \pi'(s), a_{t+1: \infty} \sim \pi \right]$$

$$\leq E \left[R_{t+1} + \gamma Q_{\pi}(s_{t+1}, \pi'(s_{t+1})) \mid s_t = s, a_t = \pi'(s) \right]$$

$$= E \left[R_{t+1} + \gamma \left(\overbrace{R_{t+2} + \gamma V_{\pi}(s_{t+2})}^{\pi(s_{t+1})} \right) \mid s_t = s, a_t = \pi'(s), a_{t+1} = \pi'(s_{t+1}) \right]$$

<

$$= E \left[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \mid s_t = s, a_t = \pi'(s), a_{t+1} = \pi'(s_{t+1}), \dots \right]$$

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