

Lecture 8 - Feb 7, 2023

- Dynamic Programming

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

Project 1 → Due Feb 8

HW2 is posted → Due Feb 17

TA's office hour:

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

Fridays, 2pm-3pm (virtual)

Bellman Eq.

$$\begin{aligned}
 V_{\pi}(s) &= E[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | s_t = s, \pi] \\
 &= E[R_{t+1} + \gamma V_{\pi}(s') | s_t = s, \pi] \\
 &= \sum_{s'} P(s' | s, \pi(s)) [R(s, \pi(s), s') + \gamma V_{\pi}(s')]
 \end{aligned}$$

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

Wall
 Bump
 Goal

$$\begin{aligned}
 Q_{\pi}(s, a) &= E[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | s_t = s, a_t = a, \pi] \\
 &= E[R_{t+1} + \gamma V_{\pi}(s') | s_t = s, a_t = a, \pi] \\
 &= \sum_{s'} P(s' | s, a) [R(s, a, s') + \gamma V_{\pi}(s')]
 \end{aligned}$$

Deterministic

$$\begin{aligned}
 V_{\pi}(15) &= \sum_{s'} P(s' | s, \pi(s)) [R + \gamma V_{\pi}(s')] \\
 &= -1 [R(15, L, \text{Goal}) + \gamma \underbrace{V_{\pi}(\text{Goal})}_0] \\
 &= 99
 \end{aligned}$$

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

Wall
 Bump
 Goal

$$V_{\pi}(13) = \sum_{s'} P(s' | 13, D) [R + \gamma V_{\pi}(s')]$$

$$= 1 \times [\underbrace{R(13, D, 15)}_{-1} + \gamma \underbrace{V_{\pi}(15)}_{99}] = 98$$

$$Q_{\pi}(15, U) = -1 + \gamma \underbrace{V_{\pi}(13)}_{98} = 97$$

$$Q_{\pi}(15, D) = -1 + \gamma \underbrace{V_{\pi}(15)}_{99} = 98$$

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

$$Q_{\pi}(15, L) = 99 + \gamma \underbrace{V_{\pi}(\text{Goal})}_0 = 99$$

Bellman Eq \Rightarrow allows us to evaluate $V_{\pi}(s)$

$$\pi_{(s)}^* = \operatorname{argmax}_{\pi \in \Pi} V_{\pi}(s) \quad \text{for all } s \in \mathcal{S}$$

$$\pi_{(s)}^* = \operatorname{argmax}_{\pi \in \Pi} Q_{\pi}(s, a) \quad \text{for all } s, a$$

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

■ Wall ■ Bump ■ Goal

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

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

Bellman Optimality Eq.

$$\pi^*: V_{\pi^*}(s) \geq V_{\pi}(s), \text{ for all } s, \pi \in \Pi$$

$$V_{\pi^*}(s) = V^*(s) = \max_{a \in A} E[\overset{R(s,a,s')}{\underset{\uparrow}{R_{t+1}}} + \gamma V^*(s') | S_t = s, a_t = a]$$

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

$$V^*(s) = \max_{a \in A} Q^*(s, a)$$

$$V_{\pi}(s) \leq Q_{\pi}(s, a) \quad \times \text{ Not true}$$

$$\underline{V_{\pi}(s)} \leq \max_{a \in A} Q_{\pi}(s, a)$$

$$Q_{\pi}(s, \pi(s))$$

$$Q_{\pi}(s, L) \rightarrow V_{\pi}(s)$$

U
D
R

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

Wall
 Bump
 Goal

Bellman Opt: $V_{\pi^*}(s) = \max_{a \in A} Q^*(s, a)$

For π^* $V_{\pi^*}(s) < \max_{a \in A} Q_{\pi^*}(s, a) = Q^*(s, \underline{a_s})$

$$V_{\pi^*}(s) < V_{\pi^1}(s) = Q_{\pi^*}(s, a^s)$$

$$a_s \neq \pi^*(s)$$

↳

$$\pi^1(s) = a_s$$

$$\pi^1(s) = \pi(s)$$

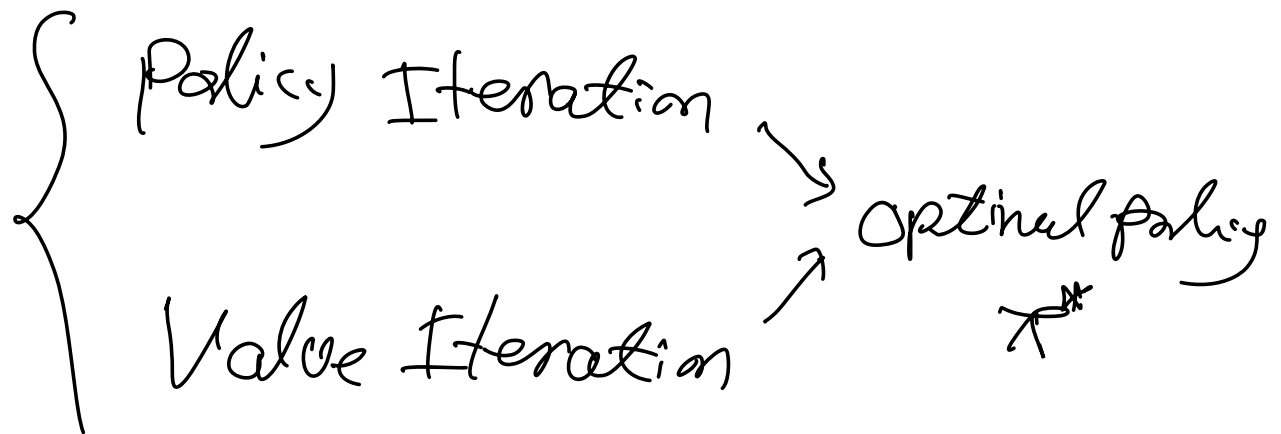
$$\begin{aligned}
 Q^*(s, a) &= E \left[\overset{R(s, a, s')}{\overset{\uparrow}{R_{t+1}}} + \gamma \underbrace{V^*(s')}_{\max_{a' \in A} Q^*(s', a')} \mid s_t = s, a_t = a \right] \\
 &= E \left[R_{t+1} + \gamma \max_{a' \in A} Q^*(s', a') \mid s_t = s, a_t = a \right]
 \end{aligned}$$

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

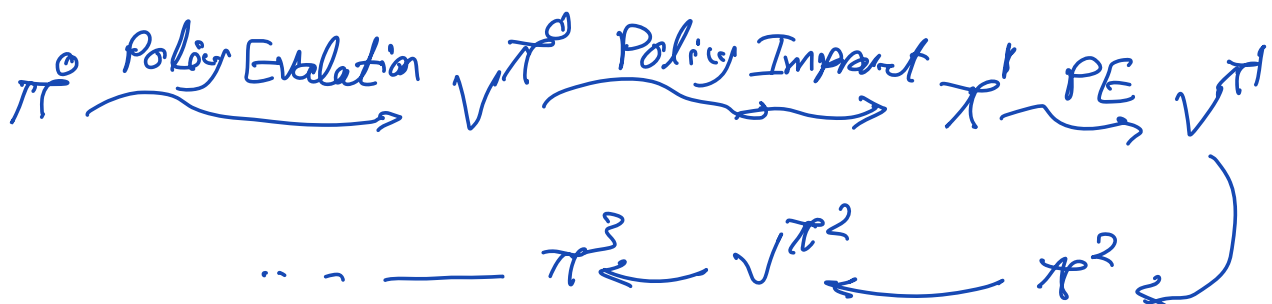
9 →	10 →	11 →	↓ 12
8 ↑		14 ↓	↓ 13
7 ↑		16	← 15
6 ↑	5		
4 ↑	3 ↑	← 2	← 1

Wall
 Bump
 Goal

Dynamic Programming (DP)



- Policy Iteration



π^T and π^{T+1} are the same $\Rightarrow \pi^T = \pi^{T+1} = \pi^*$

Policy Evaluation Step is Policy Iteration Algorithm

① Exact solution according to Bellman Eq.

A	B
---	---

Reward $\begin{cases} +5 & \text{to B} \\ -1 & \text{if } a^2 \end{cases}$

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

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

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

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

Bellman Eq.

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

$\hookrightarrow N$ States
 N Equation

N Variables $V_{\pi}(s^1) \dots V_{\pi}(s^N)$

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

$$V_{\pi}(A) = \underbrace{P(A | A, \overbrace{\pi(A)}^{a^1})}_{0.9} [\underbrace{R(A, a^1, A)}_0 + \overbrace{\gamma V_{\pi}(A)}^{0.9}]$$

$$+ \underbrace{P(B | A, \overbrace{\pi(A)}^{a^1})}_{0.1} [\underbrace{R(A, a^1, B)}_{+5} + \gamma V_{\pi}(B)]$$

$$V_{\pi}(A) = 0.81 V_{\pi}(A) + 0.5 + 0.09 V_{\pi}(B)$$

$$V_{\pi}(B) = \underbrace{P(A | B, \overbrace{\pi(B)}^{a^2})}_{0.9} [\underbrace{R(B, a^2, A)}_{-1} + \gamma V_{\pi}(A)]$$

$$+ \underbrace{P(B | B, \overbrace{\pi(B)}^{a^2})}_{0.1} [\underbrace{R(B, a^2, B)}_4 + \gamma V_{\pi}(B)]$$

$$V_{\pi}(B) = -0.9 + 0.81 V_{\pi}(A) + 0.4 + 0.09 V_{\pi}(B)$$

$$0.19 V_{\pi}(A) - 0.09 V_{\pi}(B) = 0.5$$

$$-0.81 V_{\pi}(A) + 0.91 V_{\pi}(B) = -0.5$$

$$\underbrace{\begin{bmatrix} 0.19 & -0.09 \\ -0.81 & 0.91 \end{bmatrix}}_A \begin{bmatrix} V_{\pi}(A) \\ V_{\pi}(B) \end{bmatrix} = \begin{bmatrix} 0.5 \\ -0.5 \end{bmatrix}$$

$$\begin{bmatrix} V_{\pi}(A) \\ V_{\pi}(B) \end{bmatrix} = A^{-1} \begin{bmatrix} 0.5 \\ -0.5 \end{bmatrix}$$

② Approximate Policy Evaluation

→ Start with arbitrary $V_0(s)$ for all $s \in \mathcal{S}$

$$V_{k+1}(s) = \sum_{s'} P(s' | s, \pi(s)) [R(s, \pi(s), s') + \gamma V_k(s')] \quad \text{for } s \in \mathcal{S} \quad (I)$$

$$\|V_{t+1} - V_t\|_{\infty} = \max_{s \in \mathcal{S}} |V_{t+1}(s) - V_t(s)| < \theta$$

$$\hookrightarrow V^{\pi} = V_{t+1}$$

↑
threshold

$$V_0 = \begin{bmatrix} 0 \\ c \\ \vdots \\ 0 \end{bmatrix} \xrightarrow{(I)} V_1 = \begin{bmatrix} 0 \\ 0 \\ c \\ \vdots \\ c \end{bmatrix} \xrightarrow{(I)} V_2 = \begin{bmatrix} 1 \\ \vdots \\ \vdots \\ \vdots \end{bmatrix} \dots$$

$$\pi^2 = \begin{bmatrix} a^1 \\ a^2 \end{bmatrix} \xleftarrow{\text{Randomly}} \xrightarrow{\text{Policy Evaluation}} V_{\pi^2}$$

$$V_0 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

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

$$= 0.9 \left[\underbrace{R(A, a^1, A)}_0 + \gamma \underbrace{V_0(A)}_0 \right] + 0.1 \left[\underbrace{R(A, a^1, B)}_5 + \gamma \underbrace{V_0(B)}_0 \right] = 0.5$$

$$V_1(B) = 0.9 \left[\underbrace{R(B, a^1, A)}_{-1} + \gamma \underbrace{V_0(A)}_0 \right] + 0.1 \left[\underbrace{R(B, a^1, B)}_{-4} + \gamma \underbrace{V_0(B)}_0 \right] = -0.5$$

$$V_1 = \begin{bmatrix} 0.5 \\ -0.5 \end{bmatrix}$$

$$\|V_1 - V_0\|_{\infty} = \max_i |V_1(i) - V_0(i)| = 0.5$$

$$\gamma = 0.9$$

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

$$V_2(A) = 0.9 \left[\underbrace{R(A, a^1, A)}_0 + \underbrace{\delta V_1(A)}_{0.5} \right] + 0.1 \left[\underbrace{R(A, a^1, B)}_5 + \underbrace{\delta V_1(B)}_{-0.5} \right] = 0.86$$

$$V_2(B) = 0.9 \left[\underbrace{R(B, a^2, A)}_{-1} + \underbrace{\delta V_1(A)}_{0.5} \right] + 0.1 \left[\underbrace{R(B, a^2, B)}_4 + \underbrace{\delta V_1(B)}_{-0.5} \right] = -0.14$$

$$V_2 = \begin{bmatrix} 0.86 \\ -0.14 \end{bmatrix}$$

$$\max_{\substack{b \\ \begin{bmatrix} 0.5 \\ -0.5 \end{bmatrix}}} |V_2 - V_1| = 0.36 < \underline{0.5} \quad \chi$$

$$V_3$$

$$|V_{99} - V_{100}| < \epsilon \quad \checkmark \quad V_{100} = V^* = \begin{bmatrix} 3.73 \\ 2.82 \end{bmatrix}$$

$$|T(u) - T(v)|_{\infty} < \gamma |u - v|_{\infty}$$

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

(u)

(v)

)

Policy Evaluation

Input π , the policy to be evaluated.

Initialization: a small threshold $\theta > 0$, $V(s) = 0$ for $s \in \mathcal{S}$.

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)

$V_s^\pi = V(s)$ for $s \in \mathcal{S}$



Policy Improvement

π' that is better than π

$$\pi'(s) = \arg \max_{a \in A} Q^\pi(s, a) \text{ for all } s$$

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

