

Lecture 11 - Feb 17, 2023

- Dynamic Programming

- Policy Iteration
 - Value Iteration
- } Vector-Form

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

- Approximate Dynamic Programming

- Asynchronous DP
- Generalized Policy Iteration

Exam → Tuesday, Feb 21

HW2 → Due Feb 18

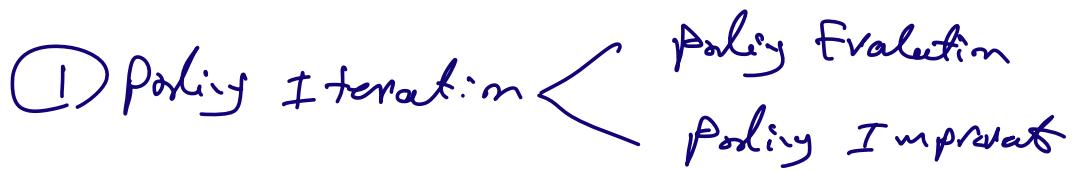
Project 2 → Due March 3

TA's office hour:

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

Fridays, 2pm-3pm (virtual)

Review:



$$\pi^0 \xrightarrow{\text{PE}} V_{\pi^0} \xrightarrow{\text{PI}} \pi^1 \xrightarrow{\text{PE}} V_{\pi^1} - - -$$

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

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

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

② Value Iteration

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

$$V_0 - - - V_T V_{T+1} = V^*$$

$$\pi^*(s) = \underset{a \in A}{\text{argmax}} \sum_{s'} p(s'|s, a) [R(s, a, s) + \delta V^*(s)]$$

Matrix-form DP:

$$M(a) = \begin{bmatrix} & \overset{P}{\swarrow} & & & & \overset{P(s'|s, a)}{\searrow} \\ & s & \downarrow & \downarrow & \downarrow & \\ & & i & j & & \\ & & & & \vdots & \\ & & & & & N \times N \end{bmatrix}$$

$\boxed{p(s'=s^i | s=s^i, a_i)}$

A	B
---	---

$M(a') = \begin{bmatrix} A & B \\ B & A \end{bmatrix}$

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

$R_{SS}^{a^1} = \begin{bmatrix} R(A, a^1, A) & R(A, a^1, B) \\ R(B, a^1, A) & R(B, a^1, B) \end{bmatrix}$

$R_{SS}^{a^2} = \begin{bmatrix} R(A, a^2, A) & R(A, a^2, B) \\ R(B, a^2, A) & R(B, a^2, B) \end{bmatrix}$

$R_{SS}^a = \begin{bmatrix} & \overset{A}{\swarrow} & & & & \overset{B}{\searrow} \\ & s & \downarrow & \downarrow & \downarrow & \\ & & i & j & & \\ & & & & \vdots & \\ & & & & & N \times N \end{bmatrix}$

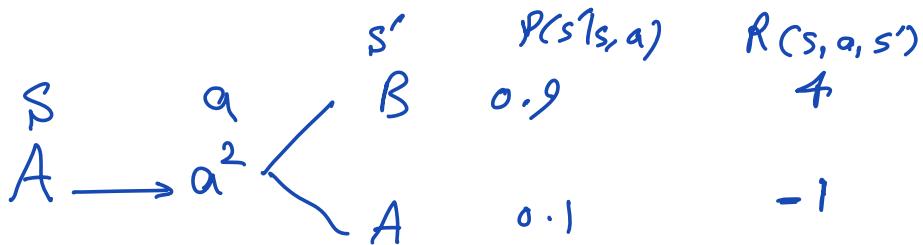
$\boxed{R(s=s^i, a, s=s^j)}$

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

$$M(a \in U) = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & & \\ 0 & 0 & 0 & 0 & 1 & - & - & - & - \end{pmatrix}$$

$$R(s, a, s')$$

$$R(s, a) = \sum_{s'} p(s'|s, a) R(s, a, s')$$



$$R(A, a^2) = 0.9 \times (-4) + 0.1 \times (-1) = -3.5$$

$$R_S^a = \begin{bmatrix} R(s', a) \\ \vdots \\ R(s^N, a) \end{bmatrix} = \left(M(a) \odot R_{ss'}^a \right) \frac{1}{N \times 1}$$

Component wise
 Hadamard product

$$\begin{bmatrix} p(s'|s', a) & \dots & p(s^N|s', a) \end{bmatrix} \odot \begin{bmatrix} R(s', a, s') & \dots & R(s', a, s^N) \end{bmatrix}$$

$$\begin{bmatrix} p(s'|s', a) R(s', a, s') & \dots & p(s^N|s', a) R(s', a, s^N) \end{bmatrix}$$

$$M(\alpha^1) = \begin{bmatrix} 0.9 & 0.1 \\ 0.1 & 0.9 \end{bmatrix} \quad R_{SS'}^{\alpha^1} = \begin{bmatrix} 0 & 5 \\ 0 & 5 \end{bmatrix}$$

$$M(\alpha^2) = \begin{bmatrix} 0.1 & 0.9 \\ 0.9 & 0.1 \end{bmatrix} \quad R_{SS'}^{\alpha^2} = \begin{bmatrix} -1 & 4 \\ -1 & 4 \end{bmatrix}$$

$$R_S^{\alpha^1} = \begin{bmatrix} R(A, \alpha^1) \\ R(B, \alpha^1) \end{bmatrix} = (M(\alpha^1) \oplus R_{SS'}^{\alpha^1}) \underbrace{\begin{pmatrix} 1 \\ 1 \end{pmatrix}}_{\text{2x1}} = \left(\begin{bmatrix} 0.9 & 0.1 \\ 0.1 & 0.9 \end{bmatrix} \oplus \begin{bmatrix} 0 & 5 \\ 0 & 5 \end{bmatrix} \right) \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

$$\begin{bmatrix} 0 & 0.5 \\ 0 & 4.5 \end{bmatrix} = \begin{bmatrix} 0.5 \\ 4.5 \end{bmatrix}$$

$$R_S^{\alpha^2} = \begin{bmatrix} R(A, \alpha^2) \\ R(B, \alpha^2) \end{bmatrix} = (M(\alpha^2) \oplus R_{SS'}^{\alpha^2}) \underbrace{\begin{pmatrix} 1 \\ 1 \end{pmatrix}}_{\text{2x1}} = \left(\begin{bmatrix} 0.1 & 0.9 \\ 0.9 & 0.1 \end{bmatrix} \oplus \begin{bmatrix} -1 & 4 \\ -1 & 4 \end{bmatrix} \right) \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

$$\begin{bmatrix} -0.1 & 3.6 \\ -0.9 & 0.4 \end{bmatrix} = \begin{bmatrix} 3.5 \\ -0.5 \end{bmatrix}$$

Value Iteration :

$$V_{k+1} = \max_{a \in A} R_s^a + \gamma M(a) V_k$$

(D)

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

$$= \max_{a \in A} \underbrace{\sum_{s'} P(s'|s, a) R(s, a, s')}_{R(s, a)} + \gamma \sum_{s'} P(s'|s, a) V_k(s')$$

$$\underbrace{P(s'=s'|s, a) V_k(s')}_{\text{red}} + \underbrace{P(s'=s^2|s, a) V_k(s^2)}_{\text{red}} + \dots + \underbrace{P(s=s^N|s, a) V_k(s^N)}_{\text{red}}$$

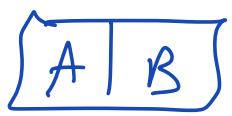
$$M(a) V_k$$

$$V_0 \rightarrow V_1 \rightarrow V_2 \rightarrow \dots \rightarrow V_T$$

$$\lim_{n \rightarrow \infty} \|V_T - V_{T-1}\| < \theta \rightarrow V^*$$

$$\pi^*(s) = \arg\max_{a \in A} \sum_{s'} p(s'|s, a) [r(s, a, s') + \gamma V^*(s')]$$

$$\begin{bmatrix} \pi^*(s') \\ \vdots \\ \pi^*(s^n) \end{bmatrix} = \overbrace{\pi^* = \arg\max_{a \in A} R_s^a + \gamma M(a) V^*}^{\text{(II)}}$$



$$M(a_1) = \begin{bmatrix} 0.9 & 0.1 \\ 0.1 & 0.9 \end{bmatrix} \quad M(a_2) = \begin{bmatrix} 0.1 & 0.9 \\ 0.9 & 0.1 \end{bmatrix}$$

$$R_{ssr}^{a_1} = \begin{bmatrix} 0 & 5 \\ 0 & 5 \end{bmatrix} \quad R_{ssr}^{a_2} = \begin{bmatrix} -1 & 4 \\ -1 & 4 \end{bmatrix}$$

$$R_s^{a_1} = \begin{bmatrix} 0.5 \\ 4.5 \end{bmatrix} \quad R_s^{a_2} = \begin{bmatrix} 3.5 \\ -0.5 \end{bmatrix}$$

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

$$V_1 = \max_{a \in A} R_s^a + \gamma M(a) V_0$$

$$= \max \left\{ \underbrace{R_s^{a_1} + \gamma M(a_1) V_0}_{a_1}, \underbrace{R_s^{a_2} + \gamma M(a_2) V_0}_{a_2} \right\}$$

$$= \max \left\{ \begin{bmatrix} 0.5 \\ 4.5 \end{bmatrix} + 0.9 \begin{bmatrix} 0.9 & 0.1 \\ 0.1 & 0.9 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 3.5 \\ -0.5 \end{bmatrix} + 0.9 \begin{bmatrix} 0.1 & 0.9 \\ 0.9 & 0.1 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \end{bmatrix} \right\}$$

Row-wise

$$= \max \left\{ \begin{bmatrix} 0.5 \\ 4.5 \end{bmatrix}, \begin{bmatrix} 3.5 \\ -0.5 \end{bmatrix} \right\} = \begin{bmatrix} 3.5 \\ 4.5 \end{bmatrix}$$

$$V_2 = \max_{a \in A} R_S^a + \gamma M(a) V_1$$

$$= \max \left\{ R_S^{a_1} + \gamma M(a_1) V_1, R_S^{a_2} + \gamma M(a_2) V_1 \right\}$$

a_1 a_2

$$= \max \left\{ \begin{bmatrix} 0.5 \\ 4.5 \end{bmatrix} + 0.9 \begin{bmatrix} 0.9 & 0.1 \\ 0.1 & 0.9 \end{bmatrix} \begin{bmatrix} 8.5 \\ 4.5 \end{bmatrix}, \begin{bmatrix} 3.5 \\ -0.5 \end{bmatrix} + 0.9 \begin{bmatrix} 0.1 & 0.9 \\ 0.9 & 0.1 \end{bmatrix} \begin{bmatrix} 3.5 \\ 4.5 \end{bmatrix} \right\}$$

$$= \max \left\{ \begin{bmatrix} 3.74 \\ 8.46 \end{bmatrix}, \begin{bmatrix} 7.46 \\ 2.74 \end{bmatrix} \right\} = \begin{bmatrix} 7.46 \\ 8.46 \end{bmatrix}$$

$$\|V_2 - V_1\|_{\infty} < \theta$$

$$\left\| \begin{bmatrix} 7.46 \\ 8.46 \end{bmatrix} - \begin{bmatrix} 3.5 \\ 4.5 \end{bmatrix} \right\| < 0.001 \times \downarrow$$

V_3

V_4

$$V_{1000} \approx V_{1001} = \begin{bmatrix} 43.1 \\ 44.1 \end{bmatrix} = V^*$$

$$\pi^* = \begin{bmatrix} \pi^*(A) \\ \pi^*(B) \end{bmatrix} = \arg\max \mathcal{R}_S^a + \gamma M(a) V^*$$

$$\pi^* = \arg\max_{a \in A} \left\{ \mathcal{R}_S^{a1} + \gamma M(a1) V^*, \mathcal{R}_S^{a2} + \gamma M(a2) V^* \right\}$$

$$= \arg\max_{a \in A} \left\{ \begin{bmatrix} 0.5 \\ 4.5 \end{bmatrix} + 0.9 \begin{bmatrix} 0.9 & 0.1 \\ 0.1 & 0.9 \end{bmatrix} \begin{bmatrix} 43.1 \\ 44.1 \end{bmatrix}, \begin{bmatrix} 3.5 \\ -0.5 \end{bmatrix} + 0.9 \begin{bmatrix} 0.1 & 0.9 \\ 0.9 & 0.1 \end{bmatrix} \begin{bmatrix} 43.1 \\ 44.1 \end{bmatrix} \right\}$$

$$= \arg\max_{\substack{\uparrow \\ \text{row-wise}}} \left\{ \underbrace{\begin{bmatrix} 39.38 \\ 44.1 \end{bmatrix}}_{a1}, \underbrace{\begin{bmatrix} 43.1 \\ 38.38 \end{bmatrix}}_{a2} \right\} = \begin{bmatrix} a2 \\ a1 \end{bmatrix} \checkmark$$

Value Iteration - Matrix Form

Transition matrices $M(a)$, $a \in A$, Reward R_S^a , Number of states N .

- Reward: $R_S^a = (M(a) \odot R_{SS}^a) \mathbb{1}_{N \times 1}$, for $a \in A$.
 Hadamard Product

- $V_0 = 0_{N \times 1}$, $K=0$

Repeat

- Value Iteration Backup: $V_{K+1} = \max_{a \in A} \left\{ R_S^a + \gamma M(a) V_K \right\}$
 Row-wise

- $K = K + 1$

Until $\max_{i \in \{1, \dots, N\}} |V_K(i) - V_{K-1}(i)| < \theta$

- Optimal State Values: $V^* = V_K$

- Optimal Policy: $\pi^* = \arg \max_{a \in A} R_S^a + \gamma M(a) V^*$
 Row-wise

Policy Iteration:



(PE)

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

$$③ V_{k+1}(s) = \sum_{s'} p(s'|s, \pi(s)) [R(s, \pi(s), s') + \gamma V_k(s')]$$

Matrix-form

$$V_{\pi}(s) = \underbrace{\sum_{s'} p(s'|s, \pi(s)) R(s, \pi(s), s')}_{R(s, \pi(s))} + \gamma \underbrace{\sum_{s'} p(s'|s, \pi(s)) V_{\pi}(s')}_{V_{\pi}(s')}$$

$$p_{s \leftarrow s'}(s, \pi(s)) V_{\pi}(s') + p_{s \leftarrow s^2}(s, \pi(s)) V_{\pi}(s^2) + \dots + p_{s \leftarrow s^n}(s, \pi(s)) V_{\pi}(s^n)$$

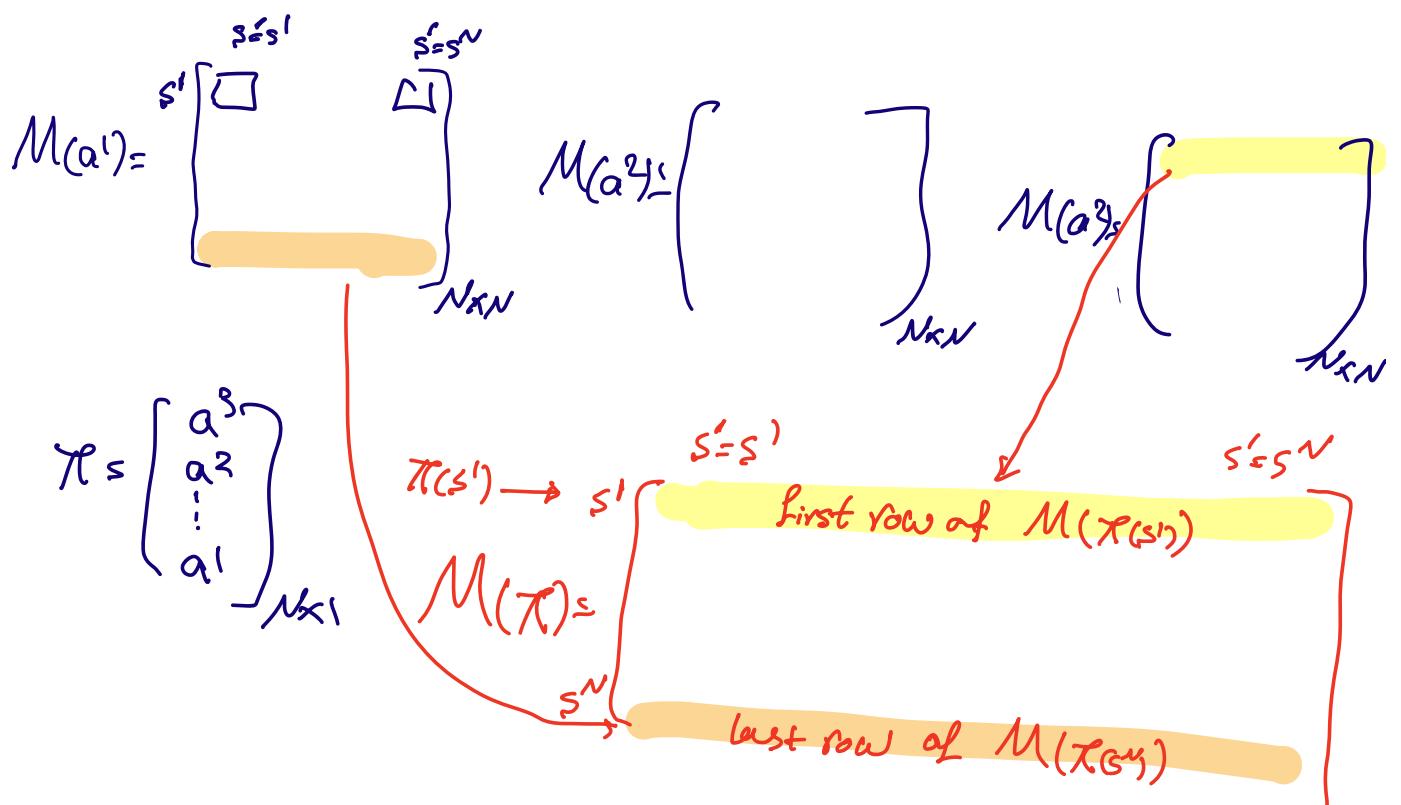
$$s \rightarrow \left[M(\pi(s)) \right]$$

$$V^\pi = R_S^\pi + \gamma M(\pi) V^\pi \Rightarrow \text{Calculate } V^\pi$$

$$(I - \gamma M(\pi)) V^\pi = R_S^\pi$$

Policy Evaluation

$$V^\pi = (I - \gamma M(\pi))^{-1} R_S^\pi$$



$$R_S^\pi = (M(\pi) \odot R_{SS'}^\pi) \mathbf{1}_N$$

$$R_S^{a_1} = \begin{bmatrix} R(s^1, a^1) \\ \vdots \\ R(s^N, a^1) \end{bmatrix}$$

$$R_S^{a_2} = \begin{bmatrix} R(s^1, a^2) \\ \vdots \\ R(s^N, a^2) \end{bmatrix} - R_S^{a_3} = \begin{bmatrix} R(s^1, a^3) \\ R(s^N, a^3) \end{bmatrix}$$

$$R_S^\pi = \begin{bmatrix} R(s^1, \pi(s^1)) \\ \vdots \\ R(s^N, \pi(s^N)) \end{bmatrix}$$

First row of $R_S^{\pi(s^1)}$

Last row of $R_S^{\pi(s^N)}$

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

↓

$$V_{k+1} = R_S^\pi + \gamma M(\pi) V_k$$

Policy Evaluation: Matrix-form

Input: Policy π

- Number of states: N , Transition Matrices: $M(a^1), \dots, M(a^L)$.

- Reward: $R_s^a = (M(a) \odot R_{ss'})^T_{N \times 1}$, for all a .

- $M(\pi) = \begin{bmatrix} & \\ & \vdots & \\ & \end{bmatrix}^{\text{Row } 1 \text{ of } M(\pi(s))} \quad , \quad R_s^\pi = \begin{bmatrix} & \\ & \vdots & \\ & \end{bmatrix}^{\text{Element 1 of } R_s^{\pi(s')}} \quad , \quad \text{Row } N \text{ of } M(\pi(s'))$

- $V^\pi = (I - \gamma M(\pi))^{-1} R_s^\pi$

Output: V^π

Policy Improvement

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

Policy Improvement

$$\pi' = \operatorname{argmax}_{a \in A} R_s^a + \gamma M(a) V^\pi$$

Example:

$$R_s^{a_1} = (M(a_1) \odot R_{ss'})^T$$

$$= \begin{bmatrix} 0.5 \\ 4.5 \end{bmatrix}$$

$$R_s^{a_2} = \begin{bmatrix} 3.5 \\ -0.5 \end{bmatrix}$$

$$M(a_1) = \begin{bmatrix} 0.9 & 0.1 \\ 0.1 & 0.9 \end{bmatrix}$$

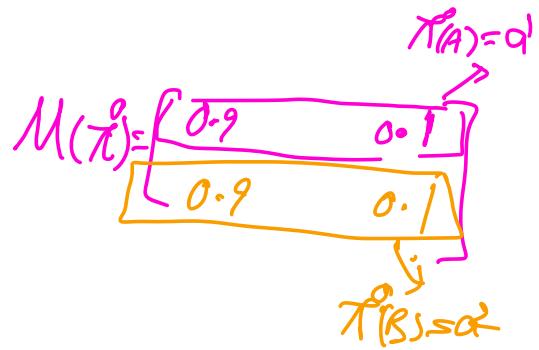
$$M(a_2) = \begin{bmatrix} 0.1 & 0.9 \\ 0.9 & 0.1 \end{bmatrix}$$

$$R_{ss'}^{a_1} = \begin{bmatrix} 0 & 5 \\ * & 5 \end{bmatrix}$$

$$R_{ss'}^{a_2} = \begin{bmatrix} -1 & 4 \\ -1 & 4 \end{bmatrix}$$

Policy Iteration:

$$\pi = \begin{bmatrix} a^1 \\ a^2 \end{bmatrix} \xrightarrow{\text{Random}} \pi^0 \xrightarrow{\text{PE}} V^{\pi^0}$$



$$R_S^{\pi^0} = \begin{bmatrix} 0.5 \\ -0.5 \end{bmatrix} \leftarrow R(A, a^1) \quad \leftarrow R(B, a^2)$$

$$V^{\pi^0} = (I - \gamma M(\pi^0))^{-1} R_S^{\pi^0}$$

$$V^{\pi^0} = \left(\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} - 0.9 \begin{bmatrix} 0.9 & 0.1 \\ 0.1 & 0.9 \end{bmatrix} \right)^{-1} \begin{bmatrix} 0.5 \\ -0.5 \end{bmatrix} = \begin{bmatrix} 4.1 \\ 3.1 \end{bmatrix}$$

$$\pi' = \underset{a \in A}{\operatorname{argmax}} \quad R_S^a + \gamma M(a) V^{\pi^0}$$

$$= \underset{a \in A}{\operatorname{argmax}} \left\{ R_S^{a^1} + \gamma M(a^1) V^{\pi^0}, R_S^{a^2} + \gamma M(a^2) V^{\pi^0} \right\}$$

$$= \underset{a \in A}{\operatorname{argmax}} \left\{ \begin{bmatrix} 0.5 \\ 4.5 \end{bmatrix} + 0.9 \begin{bmatrix} 0.9 & 0.1 \\ 0.1 & 0.9 \end{bmatrix} \begin{bmatrix} 4.1 \\ 3.1 \end{bmatrix} \right\}$$

$$, \begin{bmatrix} 3.5 \\ -0.5 \end{bmatrix} + 0.9 \begin{bmatrix} 0.1 & 0.1 \\ 0.9 & 0.1 \end{bmatrix} \begin{bmatrix} 4.1 \\ 5.1 \end{bmatrix} \}$$

$$= \arg \max \left\{ \underbrace{\begin{bmatrix} 4.1 \\ 7.38 \end{bmatrix}}_{\alpha^1}, \underbrace{\begin{bmatrix} 6.38 \\ 3.1 \end{bmatrix}}_{\alpha^2} \right\} = \begin{bmatrix} \alpha^2 \\ \alpha^1 \end{bmatrix}$$

$$\rightarrow M(\chi^2) R_S^{\chi^2} \rightarrow \sqrt{\chi^2}$$

$$\chi^3$$

Policy Iteration: Matrix Form

Initialization: arbitrary Policy π , transition matrices $M(a)$, Reward $R_{SS'}^a$

$$\text{- Reward: } R_S^a = (M(a) \odot R_{SS'}^a) I_{N \times 1} \text{, for } a \in A.$$

Policy Evaluation

$$\begin{aligned} - M(\pi) &= \begin{bmatrix} & \xrightarrow{\text{Row 1 of } M(\pi(s))} \\ & \vdots \\ & \xrightarrow{\text{Row N of } M(\pi(s))} \end{bmatrix}, R_S^\pi = \begin{bmatrix} \xrightarrow{\text{Element 1 of } R_S^{\pi(s')}} \\ \vdots \\ \xrightarrow{\text{Element N of } R_S^{\pi(s')}} \end{bmatrix} \\ - V^\pi &= (I - \gamma M(\pi))^{-1} R_S^\pi \end{aligned}$$

Policy Improvement

$$\pi' = \underset{a \in A}{\operatorname{argmax}} \xrightarrow{\text{Row-wise}} R_S^a + \gamma M(a) V^\pi$$

If $\pi \neq \pi'$, $\pi = \pi'$ and go back to policy Evaluation step.

Else

- optimal policy: $\pi^* = \pi'$
- optimal state values: $V^* = V^\pi$