

Week 5

Policy Iteration (Page 1-9)

Value Iteration (Page 10-18)

Model-Free RL (Page 20-31)

Prof. Subrahmanya Swamy

Iterative Policy Evaluation

- ▶ How to find V_π of a given policy π ?
- ▶ Iteratively apply the BE equation

$$V_{k+1}(s) = R_s^\pi + \sum_{s'} P_{ss'}^\pi V_k(s')$$

V_1 s' V_0

V_0
 V_1

Repeat till $V_{k+1} = V_k$ $\Rightarrow V_k = V_\pi$

Grid Example: Policy Evaluation

| | |
|---|----------|
| A | B |
| C | G |

- **Deterministic** state transitions
- $R_t = -1$ on all transitions
- Terminal state value $V_\pi(G) = 0$
- Discount factor $\gamma = 1$
- **Uniform** Random Policy π

V_π

Uniform Policy Dynamics:

$$\begin{cases} P_{A,A}^\pi = \frac{1}{2}, & P_{A,B}^\pi = \frac{1}{4}, & P_{A,C}^\pi = \frac{1}{4} \\ P_{B,A}^\pi = \frac{1}{4}, & P_{B,B}^\pi = \frac{1}{2}, & P_{B,G}^\pi = \frac{1}{4} \\ P_{C,A}^\pi = \frac{1}{4}, & P_{C,G}^\pi = \frac{1}{4}, & P_{C,C}^\pi = \frac{1}{2} \end{cases}$$

Iterative Policy Evaluation: $V_{k+1}(s) = R_s^\pi + \sum_{s'} P_{ss'}^\pi V_k(s')$

$$\begin{aligned} \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} &\rightarrow \begin{cases} V_1(A) = -1 + \frac{1}{2}V_0(A) + \frac{1}{4}V_0(B) + \frac{1}{4}V_0(C) \\ V_1(B) = -1 + \frac{1}{4}V_0(A) + \frac{1}{2}V_0(B) + \frac{1}{4}V_0(G) \\ V_1(C) = -1 + \frac{1}{4}V_0(A) + \frac{1}{4}V_0(G) + \frac{1}{2}V_0(C) \end{cases} \\ &\rightarrow \begin{bmatrix} -1 \\ -1 \\ -1 \end{bmatrix} \xrightarrow[k=1]{\text{update}} \begin{bmatrix} -2 \\ -1.75 \\ -1.75 \end{bmatrix} \xrightarrow[k=2]{\text{update}} \dots \end{aligned}$$

Initial
estimate

Update at $k=0$

Repeat till

$$V_{k+1} = V_k$$

$$\Rightarrow V_k = V_\pi$$

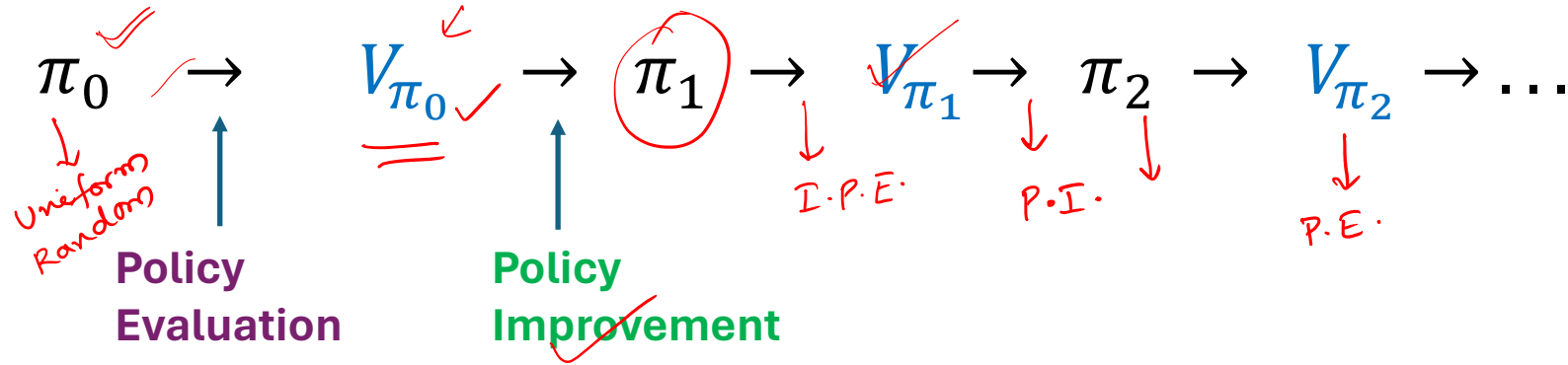
Policy Evaluation gives V_{π}

How to find Optimal Policy V^* ?
 π^*



Policy Iteration

Policy Iteration Algorithm



Repeat till

$$\pi_{k+1} = \pi_k$$

↓

$$\pi_k = \pi_*$$

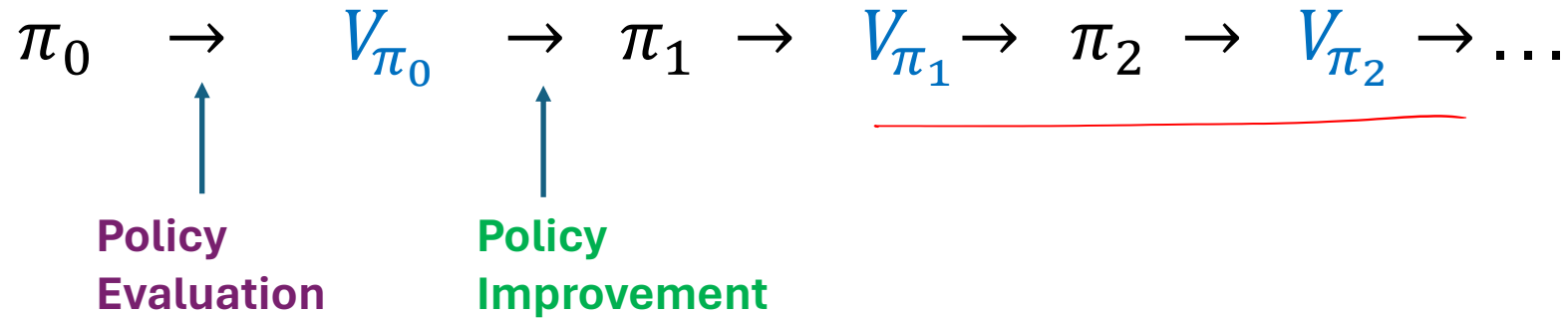
► Policy Evaluation: Iteratively apply BE equation $V_{k+1}(s) = R_s^\pi + \sum_{s'} P_{ss'}^\pi V_k(s')$

► Policy Improvement: $\pi_{i+1}(s) := \underset{a}{\operatorname{argmax}} R_s^a + \sum_{s'} P_{ss'}^a V_{\pi_i}(s')$

optimal substructure ✓

$$\pi(s) = \left\{ \begin{array}{l} a_1 \\ a_2 \\ a_3 \end{array} \right\} \left\{ \begin{array}{l} R_s^a + \sum_{s'} P_{ss'}^a V_{\pi_0}(s') \\ \end{array} \right\}$$

Policy Iteration Algorithm



Repeat till

$$\pi_{k+1} = \pi_k$$

⇓

$$\pi_k = \pi_*$$

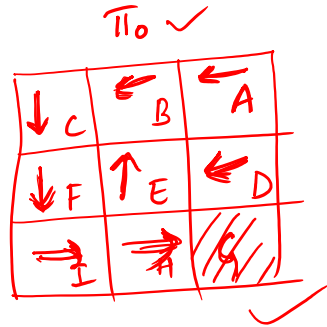
- **Policy Evaluation:** Iteratively apply BE equation $V_{k+1}(s) = R_s^\pi + \sum_{s'} P_{ss'}^\pi V_k(s')$
- **Policy Improvement:** $\pi_{i+1}(s) := \operatorname{argmax}_a R_s^a + \sum_{s'} P_{ss'}^a V_{\pi_i}(s')$
- $V_{\pi_{i+1}} \geq V_{\pi_i}$ due to Policy Improvement Theorem

Grid Example: Policy Iteration

| | |
|---|----------|
| A | B |
| C | G |

Deterministic Transitions

-1
 $\gamma=1$

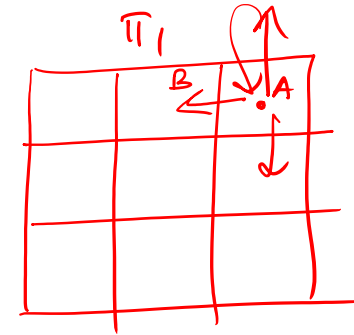


P.E.

V_{π_0}

| | | |
|----|----|----|
| -4 | -5 | -6 |
| -3 | -6 | -7 |
| -2 | -1 | 0 |

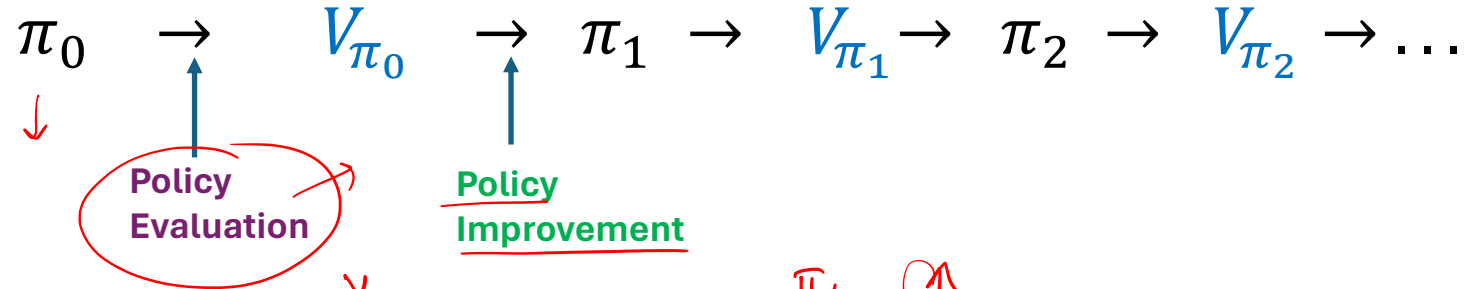
→



$$\pi_1(s) = \underset{a}{\operatorname{argmax}} \left\{ R_s^a + \gamma \sum_{s'} P_{ss'}^a V_{\pi_0}(s') \right\}$$

$$\pi_1(A) = \underset{a}{\operatorname{argmax}} \left\{ \begin{array}{l} \text{UP: } -1 + V_{\pi_0}(A) \\ \text{Left: } -1 + V_{\pi_0}(B) \\ \text{Down: } -1 + V_{\pi_0}(D) \\ \text{Right: } -1 + V_{\pi_0}(A) \end{array} \right\} = \underset{a}{\operatorname{argmax}} \left\{ \begin{array}{l} \text{U: } -1 - 6 \\ \text{✓ Left: } -1 - 5 \\ \text{D: } -1 - 7 \\ \text{R: } -1 - 6 \end{array} \right\}$$

$\Rightarrow \pi_1(A) = \text{Left}$

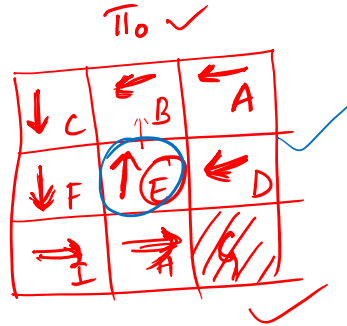


Grid Example: Policy Iteration

| | |
|---|----------|
| A | B |
| C | G |

Deterministic Transitions

-1
 $\gamma=1$

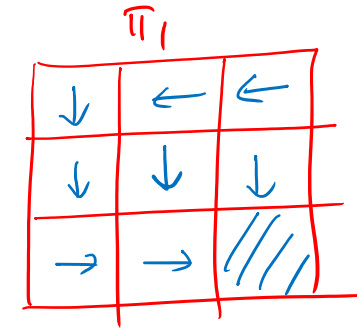


P.E.

V_{π_0}

| | | |
|----|----|----|
| -4 | -5 | -6 |
| -3 | -6 | -7 |
| -2 | -1 | 0 |

→

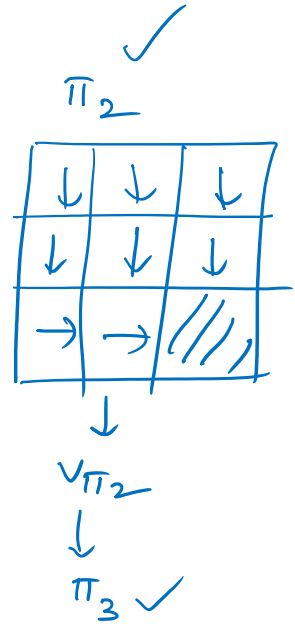


P.E.

V_{π_1}

| | | |
|----|----|----|
| -4 | -5 | -6 |
| -3 | -2 | -1 |
| -2 | -1 | 0 |

→



$$\pi_i(s) = \underset{a}{\operatorname{argmax}} \left\{ R_s^a + \gamma \sum_{s'} P_{ss'}^a V_{\pi_0}(s') \right\}$$

$$\pi_1(E) = \underset{a}{\operatorname{argmax}} \left\{ \begin{array}{l} L: -1 + V_{\pi_0}(F) \\ R: -1 + V_{\pi_0}(D) \\ U: -1 + V_{\pi_0}(B) \\ D: -1 + V_{\pi_0}(H) \end{array} \right\} = \underset{a}{\operatorname{argmax}} \left\{ \begin{array}{l} L: -1 - 3 \\ R: -1 - 7 \\ U: -1 - 5 \\ D: -1 - 1 \end{array} \right\}$$

Exercise

9

- Take π_0 as a uniform random policy
- Apply Policy iteration algorithm
- Show the sequence of policies that we get

$\pi_0 \rightarrow \pi_1 \rightarrow \dots$

| | |
|---|----------|
| A | B |
| C | G |

Deterministic Transitions

Value Iteration

Value Iteration

Bellman Optimality (BO)

$$V^*(s) = \max_a R_s^a + \sum_{s'} P_{ss'}^a V^*(s') \quad \text{(Optimal Substructure)}$$

Value Iteration

Iteratively apply BO equation till $V_{k+1} = V_k \Rightarrow V_k = V^*$

$$V_{k+1}(s) = \max_a R_s^a + \sum_{s'} P_{ss'}^a V_k(s')$$

$$\textcircled{V_0} \xrightarrow[\text{update}]{B_0} V_1 \xrightarrow[\text{update}]{B_0} V_2 \dots$$

Optimal Policy from V^*

$$\pi^*(s) = \arg \max_a R_s^a + \sum_{s'} P_{ss'}^a \underline{V^*(s')}$$

Handwritten red annotations: a checkmark under $\pi^*(s)$, a circled a under $\arg \max$, and a_2 and a_3 written below the circled a . A red arrow points down to the $V^*(s')$ term, which is underlined.

Implementation Details

- **Issue:** Takes a very long time to see $V_{k+1} \overset{\downarrow}{\underset{\uparrow}{=}} V_k$ $V_{k+1} \approx V_k$
- **Solution:** Stop when $\|V_{k+1} - V_k\|$ is small *enough*.
 $\uparrow \quad \uparrow \quad \infty \quad \underline{\epsilon = 0.01}$
- Typically, max-norm is used

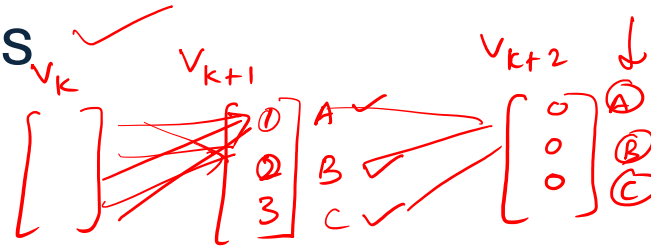
$$\begin{matrix} V_2 & V_1 \\ \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} & \begin{bmatrix} 3 \\ 4 \\ 7 \end{bmatrix} \end{matrix}$$

$$\|V_2 - V_1\|_{\infty} = \max \{ |1-3|, |2-4|, |3-7| \}$$

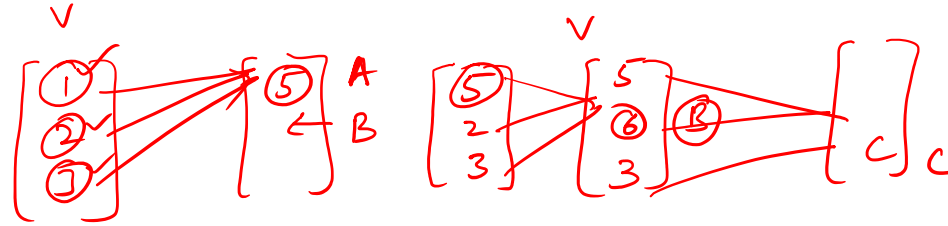
$$= 4$$

Implementation Details

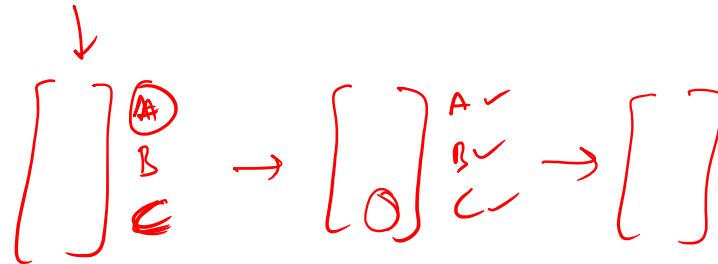
- Synchronous updates



- In-place updates



- Asynchronous updates

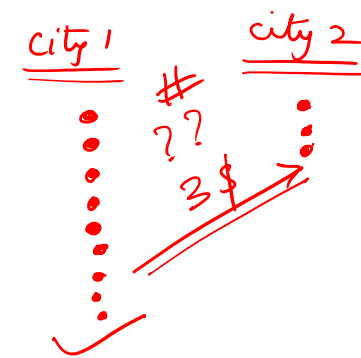


$$\Rightarrow |V_{k+1}(A) - V_k(A)| = 10 \quad V_{k+2}$$

$$\Rightarrow |V_{k+1}(B) - V_k(B)| = 0.01 \quad \checkmark$$

Car rental Example

- A car rental company operates in two cities
- Customers arrive at these cities and rent a car for \$10. If a customer arrives when cars are unavailable: **Business Loss**
- #car requests and returns are Poisson random variables
- At most, 20 cars can be parked at each location
- Upto 5 Cars be transferred overnight between cities at 3\$



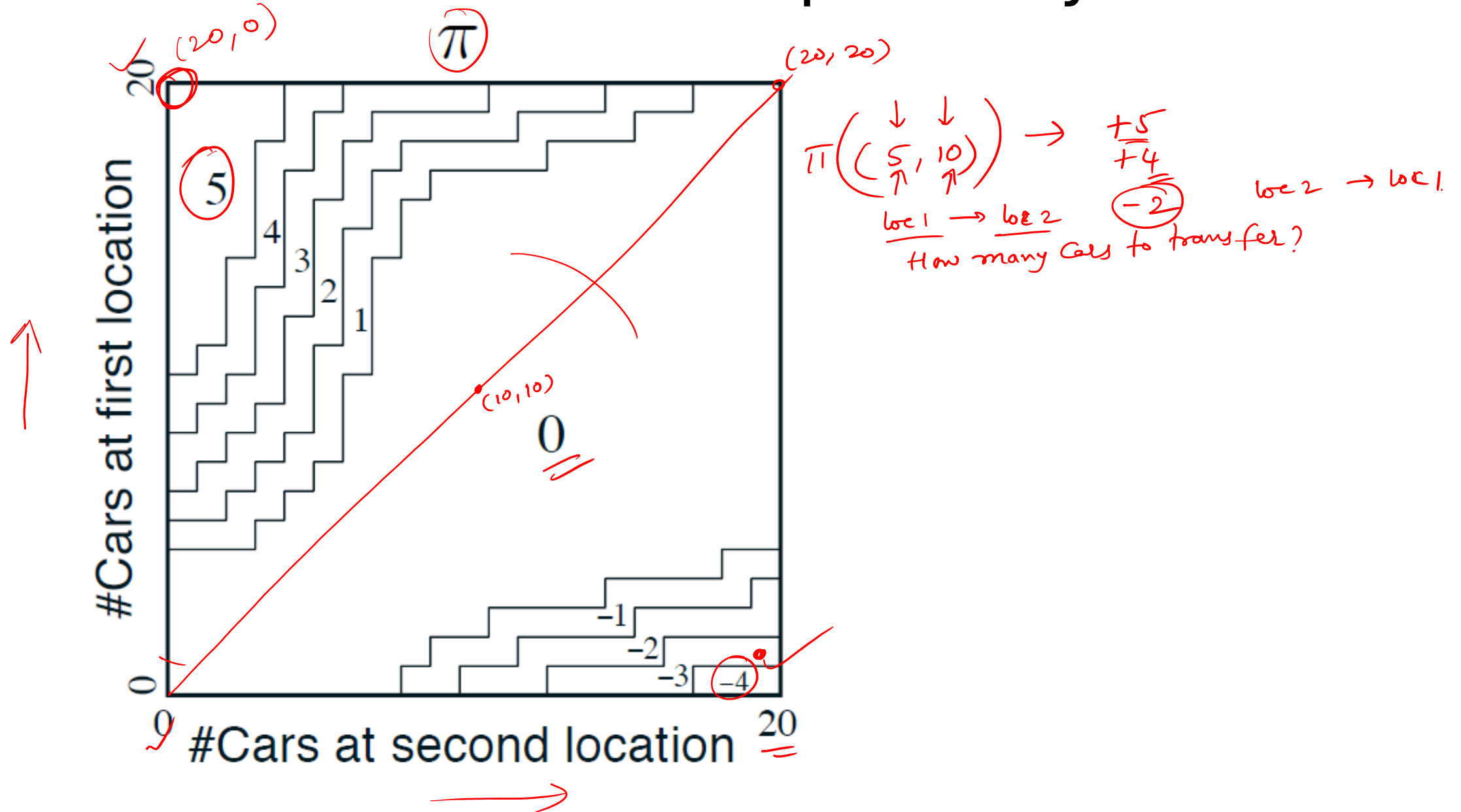
MDP

Car Rental Example

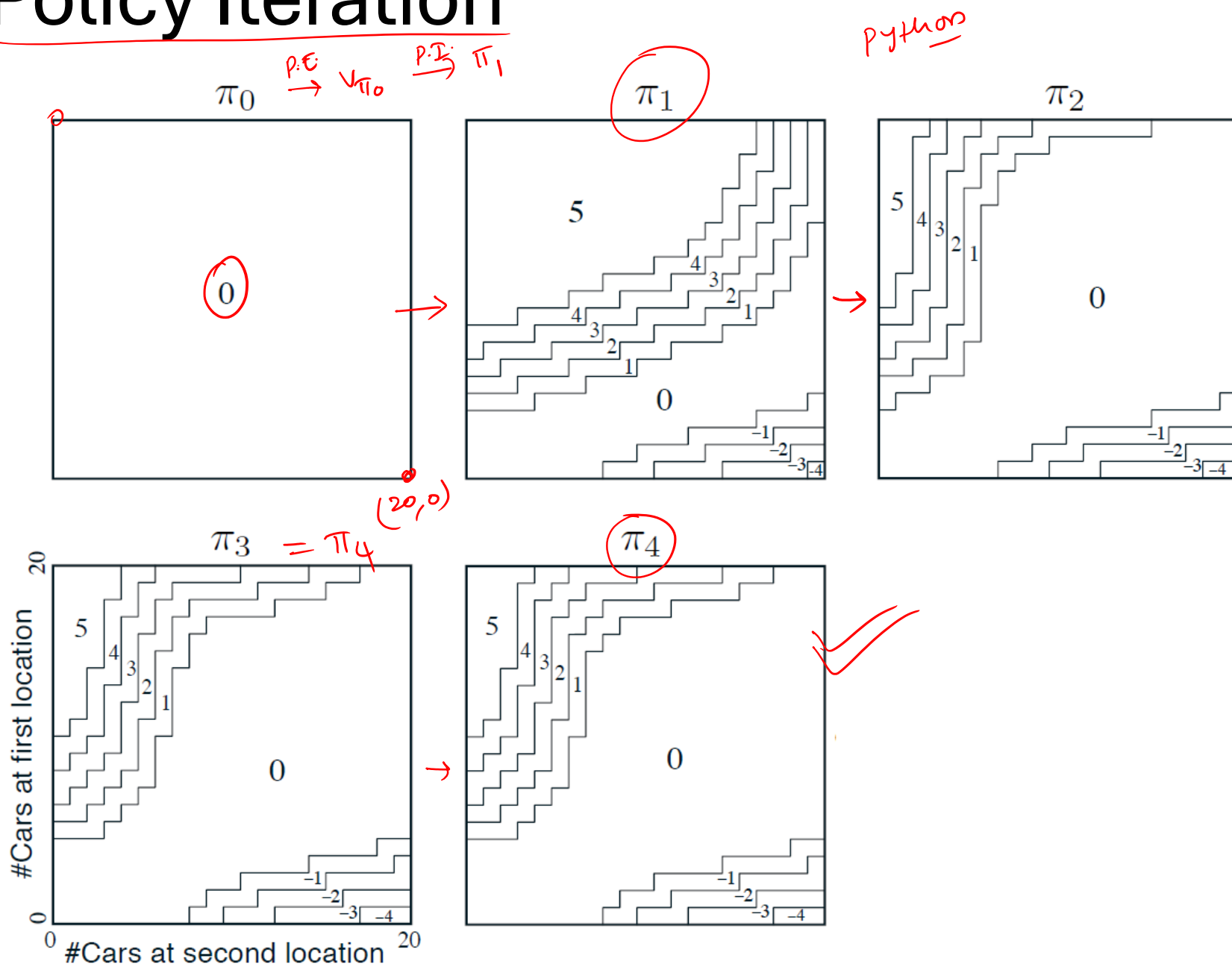
"Sutton & Belto"



Example 4.2: Jack's Car Rental Jack manages two locations for a nationwide car rental company. Each day, some number of customers arrive at each location to rent cars. If Jack has a car available, he rents it out and is credited \$10 by the national company. If he is out of cars at that location, then the business is lost. Cars become available for renting the day after they are returned. To help ensure that cars are available where they are needed, Jack can move them between the two locations overnight, at a cost of \$2 per car moved. We assume that the number of cars requested and returned at each location are Poisson random variables, meaning that the probability that the number is n is $\frac{\lambda^n}{n!} e^{-\lambda}$, where λ is the expected number. Suppose λ is 3 and 4 for rental requests at the first and second locations and 3 and 2 for returns. To simplify the problem slightly, we assume that there can be no more than 20 cars at each location (any additional cars are returned to the nationwide company, and thus disappear from the problem) and a maximum of five cars can be moved from one location to the other in one night. We take the discount rate to be $\gamma = 0.9$ and formulate this as a continuing finite MDP, where the time steps are days, the state is the number of cars at each location at the end of the day, and the actions are the net numbers of cars moved between the two locations overnight. Figure 4.2 shows the sequence of policies found by policy iteration starting from the policy that never moves any cars.



Policy Iteration



Model-Free RL

Prof. Subrahmanya Swamy

How to find V_π for a given π ?

Policy Evaluation

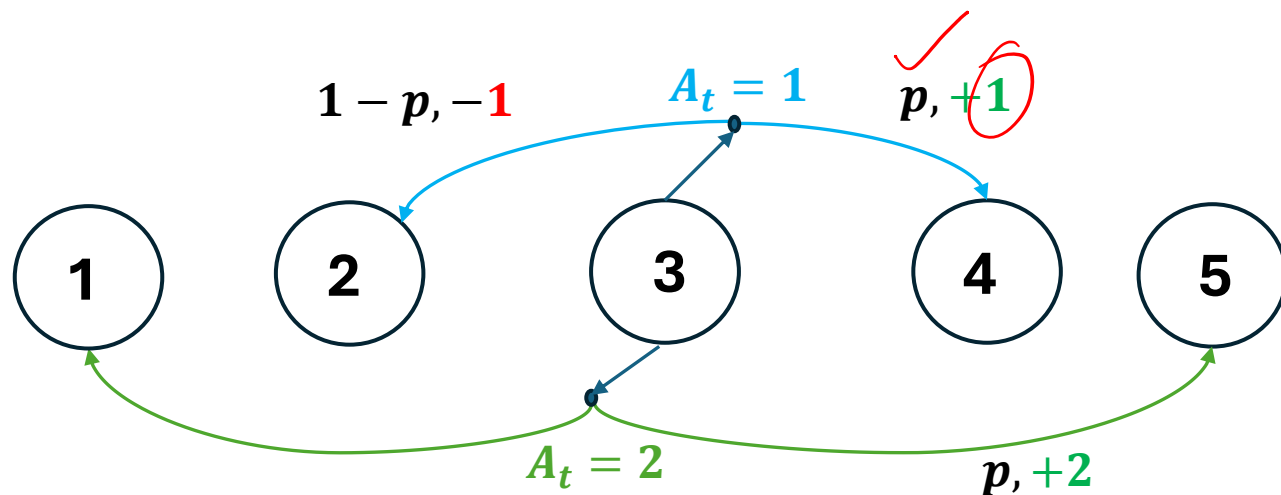
$$V_{k+1}(s) = R_s^\pi + \sum_{s'} P_{ss'}^\pi V_k(s') \quad \checkmark$$

$$R_s^\pi = \sum_a \pi(a | s) R_{sa}^a$$

$$P_{ss'}^\pi = \sum_a \pi(a | s) P_{ss'}^a$$

$V_{k+1} = V_k$

Requires complete knowledge of the environment: MDP model



Model-Free RL: Unknown $R_s^a, P_{ss'}^a$

| Task | Model Available | Model <u>Unknown</u> |
|-----------------------------|-------------------------------------|----------------------|
| Policy Evaluation V_π ✓ | Iterative Policy Evaluation ✓ | ?? |
| Optimal Policy π^* ✓ | Policy Iteration, Value Iteration ✓ | ?? |

✓ Monte-Carlo (MC)
Temporal Difference (TD)

Learn through real-time interaction with environment



Monte-Carlo (MC) method to estimate V_π

 π

Prediction

$$\bullet \quad \underline{V_\pi(s)} = \mathbb{E}_\pi[\underline{G_t} \mid S_t = s]$$

$$\mu(a) = E[R^a]$$

$$R_1, R_2, \dots, R_N$$

$$\mu(a) \approx \frac{R_1 + R_2 + \dots + R_N}{N}$$

- Interact with the environment and generate multiple episodes of data

$$\begin{aligned} \bullet \text{ Episode 1: } S_0 = \underline{s}, A_0 \sim \underline{\pi}, R_1, S_1, A_1 \sim \underline{\pi}, R_2, S_2, \dots, S_T \rightarrow G^{(1)} \\ \bullet \text{ Episode 2: } S_0 = \underline{s}, A_0 \sim \underline{\pi}, R_1, S_1, A_1 \sim \underline{\pi}, R_2, S_2, \dots, S_T \rightarrow G^{(2)} \\ \bullet \dots \\ \bullet \dots \end{aligned}$$

$G^{(N)}$

- Compute sample returns of each episode from state s

$$\bullet \quad G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$$

$$\bullet \quad \underline{V_\pi(s)} \approx \text{sample avg of the returns} \checkmark$$

Uniform Random Policy

Sample Episode

MC - Every-visit

MC First-Visit

- Two states: $\{A, B\}$ ✓
- Observed episodes:
 - $A, 1, B, -2, B, 4, A, 0, B, -2 \rightarrow$ Terminated ✓
 - $B, -1, B, 3, A, 2, B, 0, A, -3 \rightarrow$ Terminated ✓
- Observed returns
 - Episode 1: Return from first-visit of state A : $1 - 2 + 4 + 0 - 2 = 1$ ✓
 - Episode 1: Return from first-visit of state B : $-2 + 4 + 0 - 2 = 0$ ✓
 - Episode 2: Return from first-visit of state A : $2 + 0 - 3 = -1$
 - Episode 2: Return from first-visit of state B : $-1 + 3 + 2 + 0 - 3 = 1$ ✓
- MC estimates (average of observed returns):
 - $V(A) \approx \frac{1}{2}(1 - 1) = 0$ ✓
 - $V(B) \approx \frac{1}{2}(0 + 1) = \frac{1}{2}$ ✓

MC Every-Visit

- Two states: $\{A, B\}$

- Observed episodes:

- $A, 1, B, -2, B, 4, A, 0, B, -2 \rightarrow$ Terminated
- $B, -1, B, 3, A, 2, B, 0, A, -3 \rightarrow$ Terminated

- Observed returns

- Episode 1: Returns from all-visits of state A:
 - Episode 1: Return from all-visits of state B:
 - Episode 2: Return from all-visits of state A:
 - Episode 2: Return from all-visits of state B:
- 2 return samples
3 return samples for B.

- MC estimates (average of observed returns):

- $V(A) \approx$
- $V(B) \approx$

MC: Disadvantages

- Requires

G_t

$V(\cdot)$ →

online fashion

- Can be obtained only after the (episode terminates) – **Offline method**

- Not suitable for Continuing MDPs

episodic
Continuing ✓ S_T

- Alternate approach: **Temporal Difference (TD) methods!**

MC Incremental form

TD

$\checkmark G^{(1)}, \checkmark G^{(2)}, \dots \checkmark G^{(N)}$

Estimate after observing returns from N episode: $\underline{\underline{V_N(s)}}$ $\underline{\underline{V_N(s)}} = \frac{G^{(1)} + G^{(2)} + \dots + G^{(N)}}{N}$

Observed return in next episode: $\underline{\underline{G^{(N+1)}}}$ $\underline{\underline{V_{N+1}(s)}} = \frac{G^{(1)} + G^{(2)} + \dots + G^{(N)} + G^{(N+1)}}{N+1}$

What is $V_{N+1}(s)$?

$$V_{N+1}(s) = \frac{NV_N(s) + G^{(N+1)}}{N+1}$$

$$V_{N+1}(s) = \frac{V_N(s)(N+1) - V_N(s) + G^{(N+1)}}{N+1}$$

$$\rightarrow V_{N+1}(s) = \underset{\uparrow}{V_N(s)} + \frac{1}{N+1} \left[\underset{\downarrow}{G^{(N+1)}} - V_N(s) \right]$$

MC Incremental form

Update: $V_{N+1}(s) = V_N(s) + \frac{1}{N+1} (G_t^{(N+1)} - V_N(s)) =$

$$V_{new}(s) = V_{old}(s) + \alpha (G_t - V_{old}(s))$$

New Estimate Old Estimate Error between current sample and old estimate

Different Sample estimates \Rightarrow Different algorithms MC

TD

Temporal Difference (TD) method for V_π

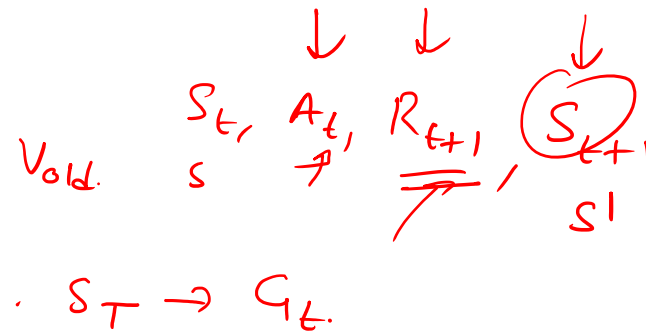
$$\begin{aligned}
 V_\pi(s) &= \mathbb{E}_\pi[\overset{\downarrow}{G_t} \mid \overset{\downarrow}{S_t} = s] \\
 &= \mathbb{E}_\pi[R_{t+1} + \gamma \overset{\downarrow}{G_{t+1}} \mid S_t = s] \\
 &= \mathbb{E}_\pi[R_{t+1} \mid S_t = s] + \gamma \mathbb{E}_\pi[G_{t+1} \mid S_t = s]
 \end{aligned}$$

MC update: $V_{new}(s) = V_{old}(s) + \alpha (\overset{\downarrow}{G_t} - V_{old}(s))$

TD update: $V_{new}(s) = V_{old}(s) + \alpha (\overset{\downarrow}{??} - V_{old}(s))$

Temporal Difference (TD) method for V_π

$$\begin{aligned}
 V_\pi(s) &= \mathbb{E}_\pi[G_t | S_t = s] \\
 &= \mathbb{E}_\pi[R_{t+1} + \gamma G_{t+1} | S_t = s] \\
 &= \mathbb{E}_\pi[R_{t+1} | S_t = s] + \gamma \mathbb{E}_\pi[G_{t+1} | S_t = s] \\
 &\approx R_{t+1} + \gamma V_\pi(S_{t+1})
 \end{aligned}$$



MC update: $V_{new}(s) = V_{old}(s) + \alpha (G_t - V_{old}(s))$

TD update: $V_{new}(s) = V_{old}(s) + \alpha (R_{t+1} + \gamma V_{old}(s') - V_{old}(s))$

MC

TD

V_π

"online"

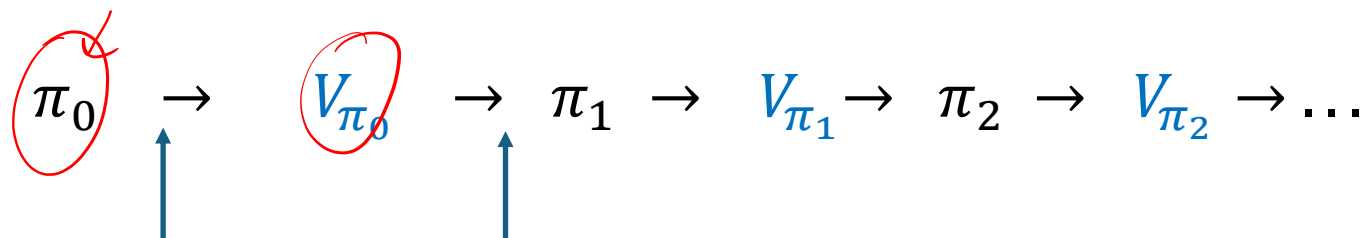
"continuing"

π^*

How to find π^* in model-free setting?

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

Policy Iteration



Policy
Evaluation

TD MC

Policy
Improvement

MC
TD

Repeat till

$$\pi_{k+1} = \pi_k$$



$$\pi_k = \pi_*$$

- Policy Evaluation: Iteratively apply BE equation $V_{k+1}(s) = R_s^\pi + \sum_{s'} P_{ss'}^\pi V_k(s')$
- Policy Improvement: $\pi_{i+1}(s) := \operatorname{argmax}_a (R_s^a + \sum_{s'} P_{ss'}^a V_{\pi_i}(s'))$

Can we use Policy Iteration in a model-free setting?

- Replace PE with MC/TD-based V_π estimation
- What about Policy Improvement?