

Generalized Policy Iteration (Page 1-12)

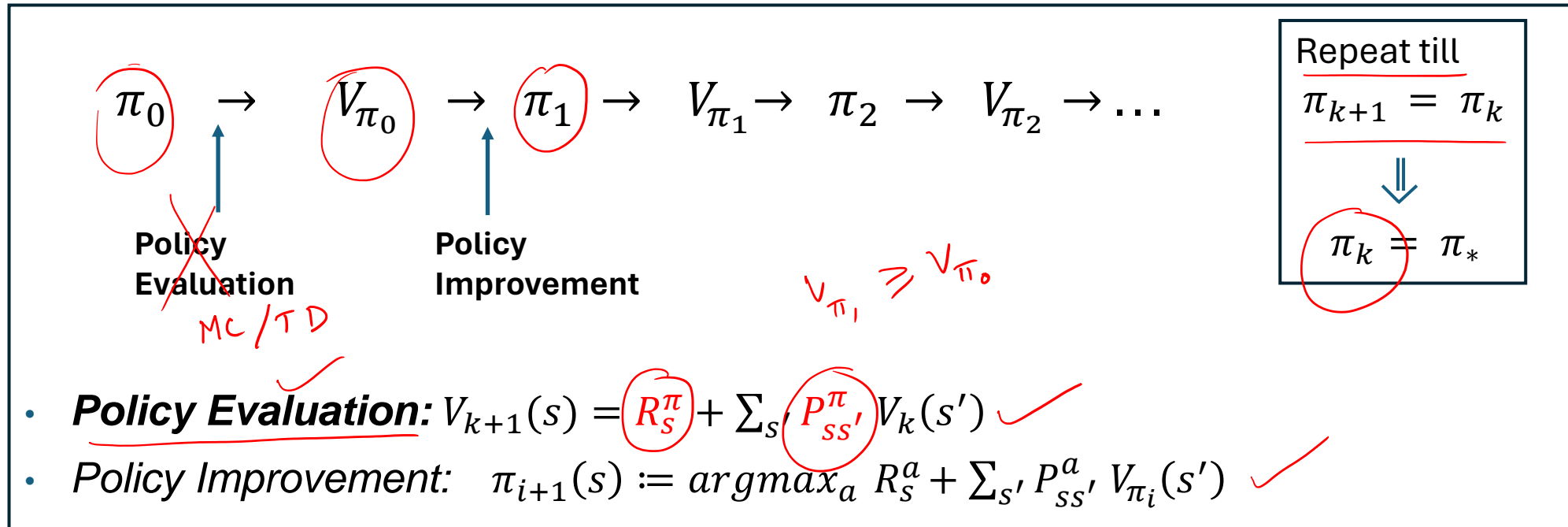
N-Step TD method (Page 13)

Off-Policy MC method (Page 14-16)

Prof. Subrahmanya Swamy

Challenges of Policy Iteration in Model-Free Context

Policy Iteration



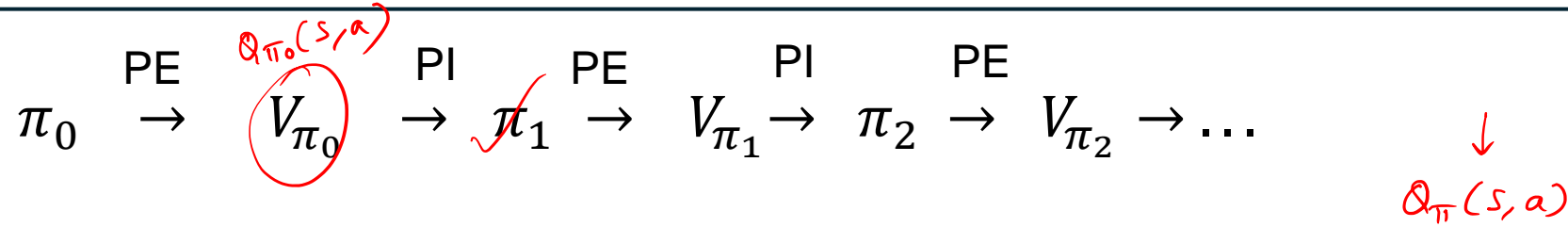
1. Policy Evaluation requires model dynamics

Solution:

- Don't use Iterative ~~Policy Evaluation~~ to estimate V_π
- Instead use MC/TD methods to estimate V_π

Challenges of Policy Iteration in Model-Free Context

Policy Iteration



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

2. Policy Improvement requires model dynamics $\operatorname{argmax}_a \pi_i(s) = \uparrow Q_{\pi_i}(s,a)$

Solution:

$$\text{PI} \quad \pi_{i+1}(s) := \operatorname{argmax}_a \left(R_s^a + \sum_{s'} P_{ss'}^a V_{\pi_i}(s') \right)$$

$$\Rightarrow \operatorname{argmax}_a (Q_{\pi_i}(s,a))$$

- If Q_{π_i} is known, model dynamics not required for PI
- Hence, estimate Q_π instead of V_π in the PE step

How to estimate $Q_{\pi}(s, a)$?

V_{π} $\overset{\pi}{\text{MC}}$ / TD

MC method to estimate $\underline{V_\pi}$

- $V_\pi(s) = \mathbb{E}_\pi[G_t \mid S_t = s]$ ✓ $\rightarrow Q_\pi(s, a) = \mathbb{E}_\pi[G_t \mid S_t = s, A_t = a]$

- Generate multiple episodes starting from s
 - Episode 1: $S_0 = s, A_0 \sim \pi, R_1, S_1, A_1 \sim \pi, R_2, S_2, \dots, S_T$ $q^{(1)}$
 - Episode 2: $S_0 = s, A_0 \sim \pi, R_1, S_1, A_1 \sim \pi, R_2, S_2, \dots, S_T$ $q^{(2)}$
 - ... \vdots
 - ... $q^{(N)}$
- Compute sample returns of each episode from state s
 - $G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$

- $V_\pi(s) \approx$ sample avg of the returns

MC method to estimate Q_π

- $Q_\pi(s, a) = \mathbb{E}_\pi[G_t \mid S_t = s, A_t = a]$

TD Q_π ?

(s, a)

- Generate multiple episodes starting from (s, a)

- Episode 1: $(S_0 = \underline{s}, A_0 = \underline{a}) R_1, S_1, A_1 \sim \pi, R_2, S_2, \dots, S_T \rightarrow G^{(1)}$
 - Episode 2: $(S_0 = \underline{s}, A_0 = \underline{a}) R_1, S_1, A_1 \sim \pi, R_2, S_2, \dots, S_T \rightarrow G^{(2)}$
 - ...
 - ...
- $G^{(N)}$

- Compute sample returns of those episodes starting from (s, a)

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

- $Q_\pi(s, a) \approx$ sample avg of the returns

TD Method for Q_π

$$E[R_{t+1} | s_t = s, A_t = a]$$

$$Q_\pi(s, a) = R_s^a + \gamma \sum_{s'} P_{ss'}^a \underbrace{V_\pi(s')}_{(s', a')}$$

$$V_\pi(s') = \left(\sum_{a'} Q_\pi(s', a') \pi(a' | s') \right)$$

$$Q_\pi(s, a) = R_s^a + \gamma \sum_{s'} P_{ss'}^a \left(\sum_{a'} Q_\pi(s', a') \pi(a' | s') \right)$$

$$= R_s^a + \gamma \sum_{s'} \sum_{a'} P_{ss'}^a Q_\pi(s', a') \pi(a' | s')$$

$$= E[R_{t+1} | s_t = s, A_t = a]$$

$$+ \gamma E[Q_\pi(s', a')]$$

$s', a' \sim \underbrace{P_{ss'}^a}, \underbrace{\pi(a' | s')}$

$$Q_\pi(s, a) = E[$$

TD Method for Q_π (SARSA)

$$\underline{V_\pi(s)} = E_\pi [\underline{G_t} | s_t = s]$$

$$\begin{aligned} Q_\pi(s, a) &= \mathbb{E}[R_{t+1} | S_t = s, A_t = a] + \gamma \mathbb{E}[V_\pi(S_{t+1})] \\ &= \mathbb{E}[\underline{R_{t+1}} | S_t = s, A_t = a] + \gamma \mathbb{E}[\mathbb{E}[\underline{Q_\pi(S_{t+1}, A_{t+1})}]] \\ \underline{Q_\pi(s_t, a_t)} &\approx \underline{R_{t+1} + \gamma Q_\pi(S_{t+1}, A_{t+1})} \end{aligned}$$

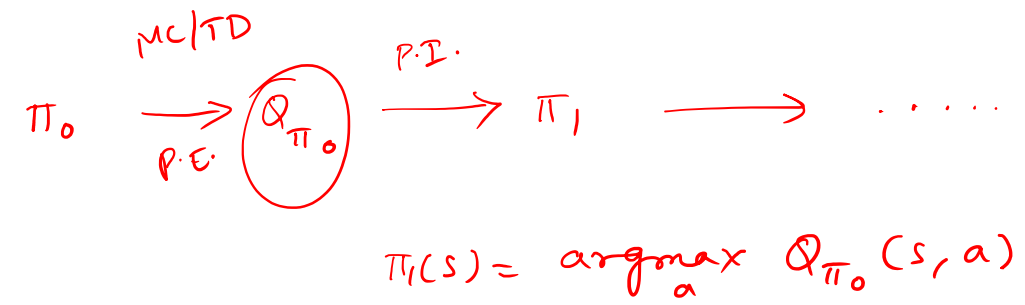
$$V_{\text{new}}(s) = V_{\text{old}}(s) + \alpha [\underline{G_t} - V_{\text{old}}(s)]$$

$$Q_{\text{new}}(S_t, A_t) = Q_{\text{old}}(S_t, A_t) + \underline{\alpha} (\underline{R_{t+1} + \gamma Q_{\text{old}}(S_{t+1}, A_{t+1})} - \underline{Q_{\text{old}}(S_t, A_t)})$$

SARSA

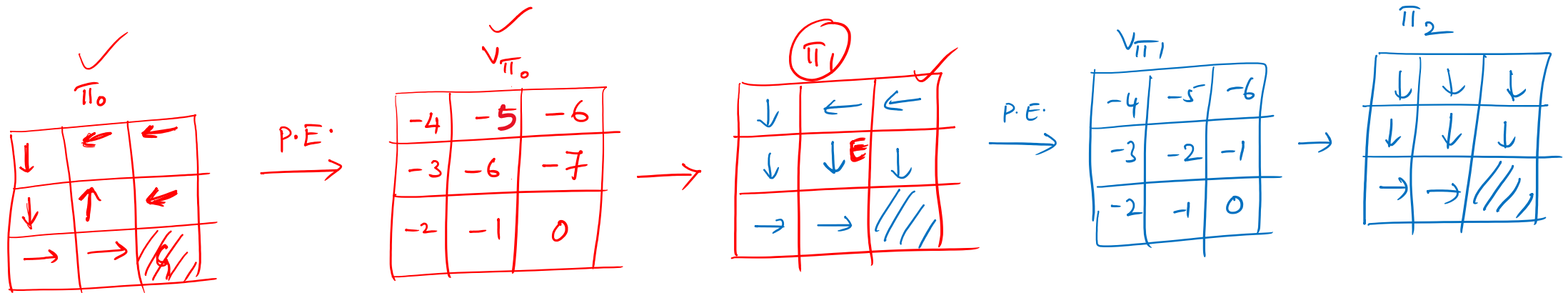
$S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1}$

Q_π
MC/TD



Q_π Estimation: Challenges

- **Observation:** Only Deterministic policies are encountered in Policy Iteration



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

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

Issue with Deterministic Policy

- Consider 3 states A, B, C ✓
- 4 actions in each state: Left, Right, Up, Down
- Consider a deterministic policy
 - $\pi(A) = \text{Right}$ ✓
 - $\pi(B) = \text{Right}$ ✓
 - $\pi(C) = \text{Left}$ ✓
- Sample Episode

ϵ -greedy policies

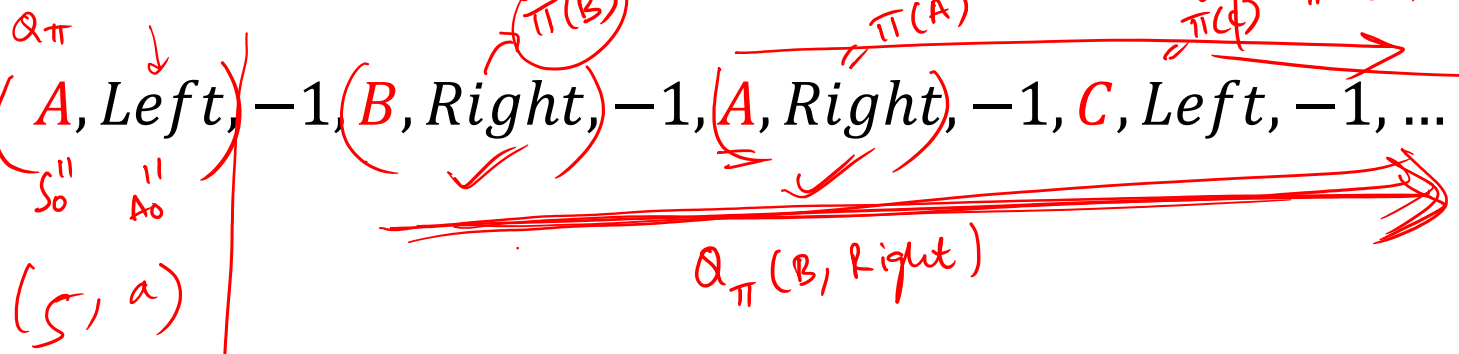
$\rightarrow \pi_t(s) = \underset{a}{\operatorname{argmax}} Q_{\pi_0}(s, a)^{1-\epsilon}$

ϵ

3x4

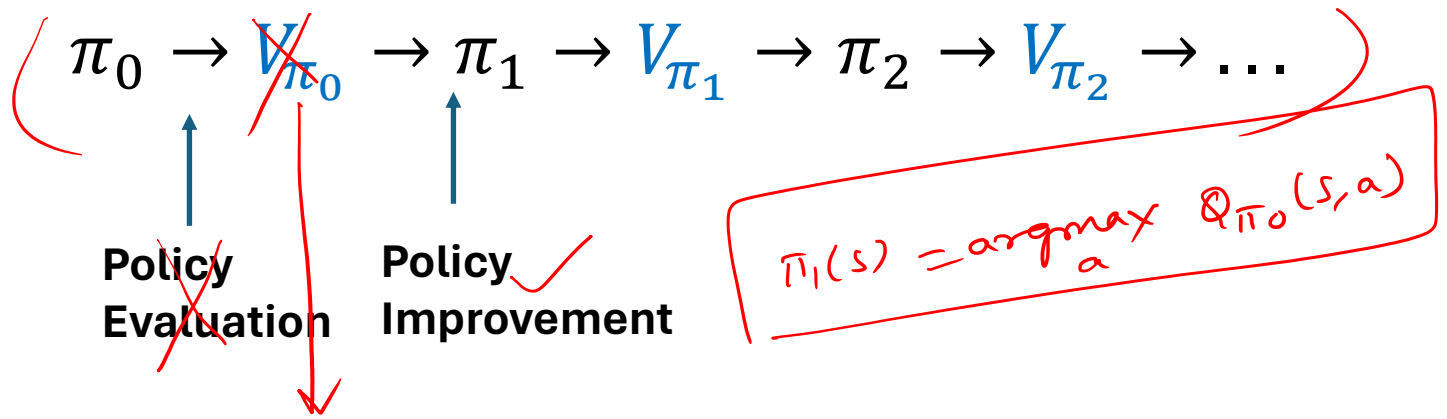
$Q_{\pi}(s, a)$

| $Q_{\pi}(A, \text{left})$ ✓ | $Q(B, L)$ ✓ | $\rightarrow Q(C, L)$ ✓ |
|-------------------------------|-------------------------|-------------------------|
| $\rightarrow Q_{\pi}(A, R)$ ✓ | $\rightarrow Q(B, R)$ ✓ | $Q(C, R)$ |
| $Q_{\pi}(A, \text{up})$ ✓ | $Q(B, \text{up})$ ✓ | $Q(C, U)$ |
| $Q_{\pi}(A, \text{down})$ ✓ | $Q(B, D)$ ✓ | $Q(C, D)$ |

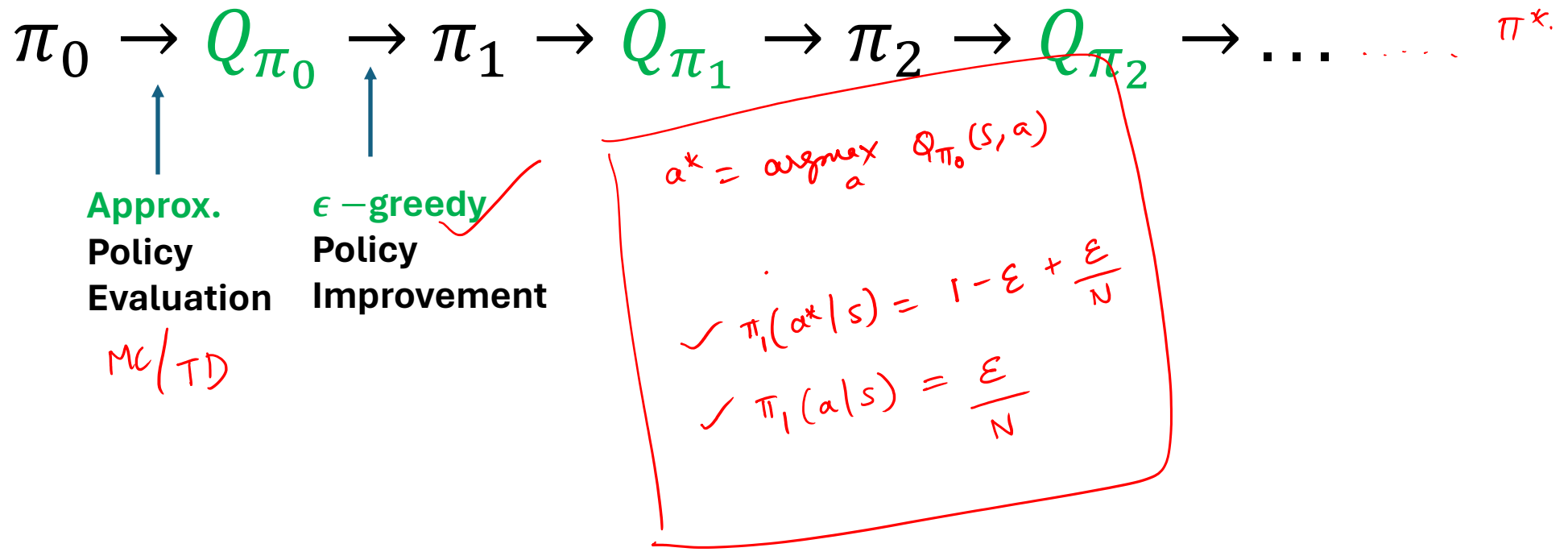


Generalized Policy Iteration (GPI)

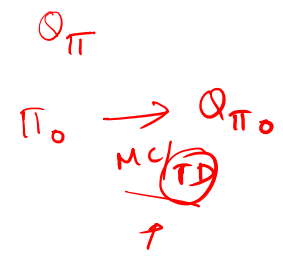
Policy Iteration



GPI



N-Step TD Method



$$\rightarrow G_t = \underbrace{R_{t+1}}_{\text{Immediate reward}} + \gamma \underbrace{G_{t+1}}_{\text{remaining return}}$$

$$Q_\pi(s, a) \approx R_{t+1} + \gamma Q_\pi(s_{t+1}, A_{t+1})$$

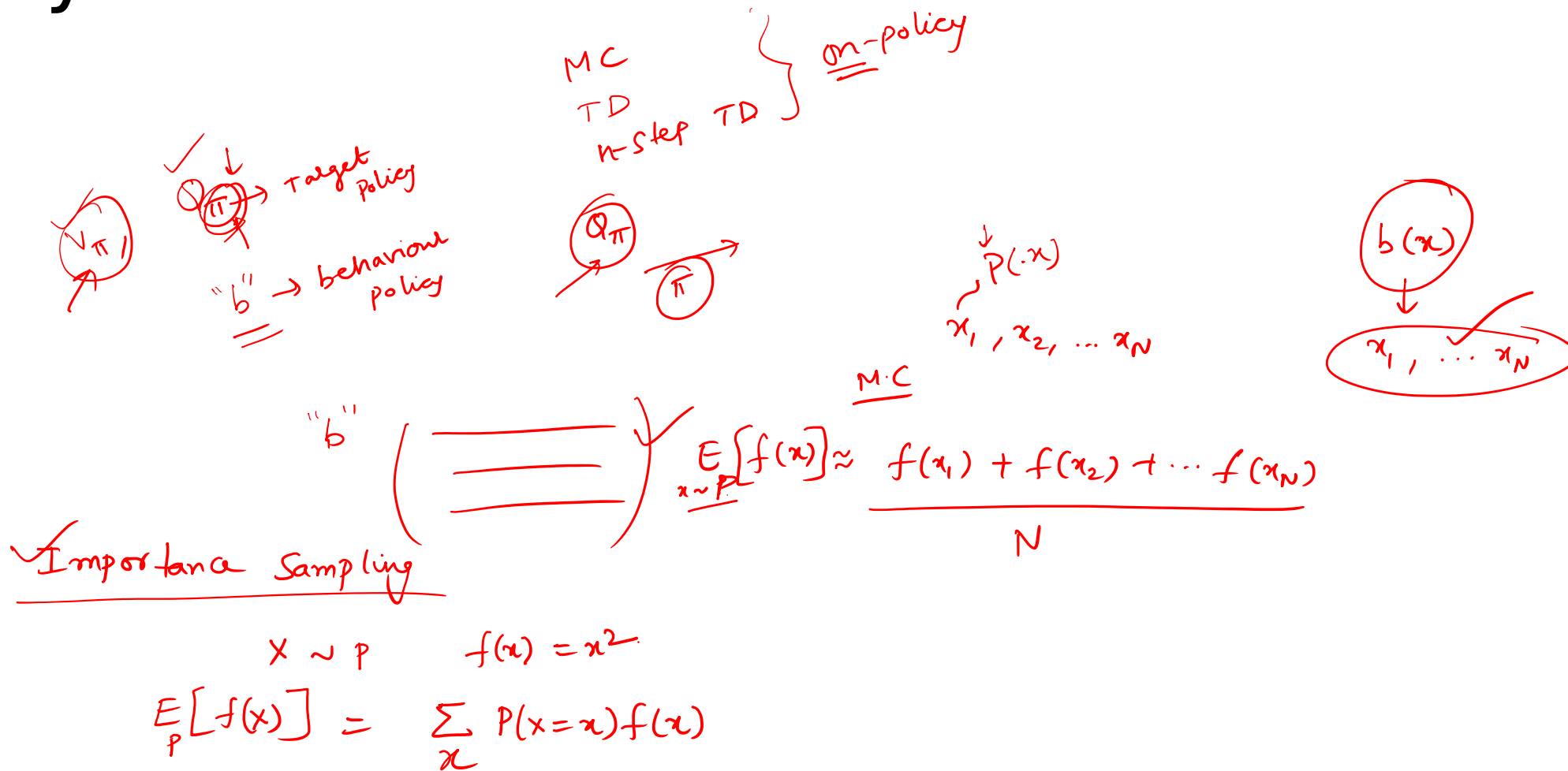
$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$$

$$\boxed{G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 G_{t+2}}$$

2-step TD

$$Q_\pi(s, a) = R_{t+1} + \gamma R_{t+2} + \gamma^2 Q_\pi(s_{t+2}, A_{t+2})$$

Off-Policy MC



Off-Policy MC

$$\begin{aligned}
 E_P[f(x)] &= \sum_x f(x) P(x=x) \\
 &= \sum_x f(x) \cancel{b(x=x)} \times \frac{P(x=x)}{\cancel{b(x=x)}} \\
 &= \sum_x \left(f(x) \cdot \frac{P(x=x)}{\cancel{b(x=x)}} \right) \cancel{b(x=x)}
 \end{aligned}$$

MC $v_\pi | Q_\pi$

$$E_P[f(x)] = \sum_x g(x) \cancel{b(x=x)} = \overset{\downarrow}{E_b[g(x)]} \approx \frac{g(x_1) + g(x_2) + \dots + g(x_N)}{N}$$

E_P

$$\approx \frac{1}{N} \left(f(x_1) \frac{P(x_1)}{\cancel{b(x_1)}} + f(x_2) \frac{P(x_2)}{\cancel{b(x_2)}} + \dots \right)$$

Importance sampling
Tech.

Off-Policy MC MC for V_π .

$$\begin{array}{lcl}
 & V_\pi(s) & \\
 \text{ep1} & S_0 = s, \underline{A_0 \sim \pi}, R_1, S_1, A_1 \sim \pi, \dots & G^{(1)} \\
 \text{ep2} & & G^{(2)}
 \end{array}$$

$$V_\pi(s) \approx \frac{G^{(1)} + G^{(2)} + \dots + G^{(N)}}{N}$$

behavior b.

$$\begin{array}{lcl}
 \text{ep1} & S_0 = s, \underline{A_0 \sim b}, R_1, S_1, A_1 \sim b, \dots & G^{(1)} \\
 \text{ep2} & & \vdots \\
 & & G^{(N)}
 \end{array}$$

$$V_\pi(s) \approx \frac{1}{N} \left(\left(G^{(1)} \cdot \frac{P_\pi(\text{ep1})}{P_b(\text{ep1})} \right) + \left(G^{(2)} \cdot \frac{P_\pi(\text{ep2})}{P_b(\text{ep1})} \right) + \dots \right)$$

$$\frac{P_\pi(\text{ep1})}{P_b(\text{ep1})} = \frac{\pi(A_0|S_0) \pi(A_1|S_1) \pi(A_2|S_2) \dots}{b(A_0|S_0) b(A_1|S_1) b(A_2|S_2) \dots}$$