# Generalized Policy Iteration (Page1-10)
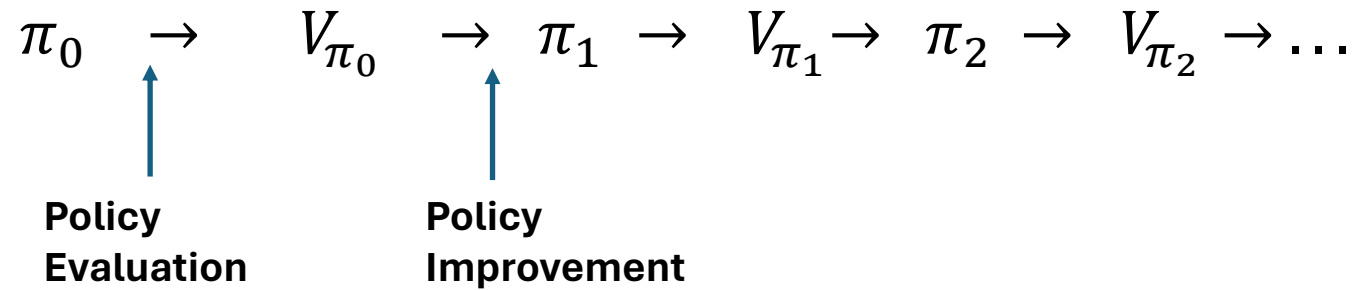# N-Step TD method (Page 11)
# Off-Policy MC method (Page 12-14)

Prof. Subrahmanya Swamy

# Challenges of Policy Iteration in Model-Free Context

**Policy Iteration**

$$\pi_0 \quad \rightarrow \quad V_{\pi_0} \quad \rightarrow \quad \pi_1 \quad \rightarrow \quad V_{\pi_1} \rightarrow \quad \pi_2 \quad \rightarrow \quad V_{\pi_2} \rightarrow \dots$$

Policy Evaluation

Policy Improvement

Repeat till
$$\pi_{k+1} = \pi_k$$
$$\Downarrow$$
$$\pi_k = \pi_*$$

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

- Policy Improvement: $\pi_{i+1}(s) := argmax_a \; R_s^a + \sum_{s'} P_{ss'}^a V_{\pi_i}(s')$

1. Policy Evaluation requires model dynamics

**Solution:**
- Don't use Iterative Policy Evaluation to estimate $V_\pi$
- Instead use MC/TD methods to estimate $V_\pi$

# Challenges of Policy Iteration in Model-Free Context

**Policy Iteration**

$$\pi_0 \overset{PE}{\to} V_{\pi_0} \overset{PI}{\to} \pi_1 \overset{PE}{\to} V_{\pi_1} \overset{PI}{\to} \pi_2 \overset{PE}{\to} V_{\pi_2} \to \dots$$

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

## 2. Policy Improvement requires model dynamics

**Solution:**

$$\textbf{PI}\quad \pi_{i+1}(s) := argmax_a\ R_s^a + \sum_{s'} P_{ss'}^a V_{\pi_i}(s')$$

$$= argmax_a\ Q_{\pi_i}(s, a)$$

- If $Q_{\pi_i}$ is known, model dynamics not required for PI
- Hence, estimate $Q_\pi$ instead of $V_\pi$ in the PE step

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

# MC method to estimate $V_\pi$

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

- Generate multiple episodes starting from $s$
  - Episode 1: $S_0 = s, \ A_0 \sim \pi, \ R_1, \ S_1, \ A_1 \sim \pi, \ R_2, \ S_2, \ \ldots, \ S_T$
  - Episode 2: $S_0 = s, \ A_0 \sim \pi, \ R_1, \ S_1, \ A_1 \sim \pi, \ R_2, \ S_2, \ \ldots, \ S_T$
  - …
  - …

- Compute sample returns of each episode from state $s$
  - $G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \ldots$

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

# MC method to estimate $Q_\pi$

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

- Generate multiple episodes starting from $(s, a)$
  - Episode 1: $S_0 = s$, $A_0 = a$, $R_1$, $S_1$, $A_1 \sim \pi$, $R_2$, $S_2$, $\ldots$, $S_T$
  - Episode 2: $S_0 = s$, $A_0 = a$, $R_1$, $S_1$, $A_1 \sim \pi$, $R_2$, $S_2$, $\ldots$, $S_T$
  - $\ldots$
  - $\ldots$

- Compute sample returns of those episodes starting from $(s, a)$
  - $G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \ldots$

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

# TD Method for $Q_\pi$ (SARSA)

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

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

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

**S A R S A**

# $Q_\pi$ Estimation: Challenges

- Observation: Only Deterministic policies are encountered in Policy Iteration



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

$$\pi_1(E) = \operatorname{argmax} \begin{cases} L: & -1 + V_{\pi_0}(F) \\ R: & -1 + V_{\pi_0}(D) \\ U: & -1 + V_{\pi_0}(B) \\ D: & -1 + V_{\pi_0}(H) \end{cases} = \operatorname{argmax}_a \begin{cases} L: -1 - 3 \\ R: -1 - 7 \\ U: -1 - 5 \\ D: -1 - 1 \end{cases}$$
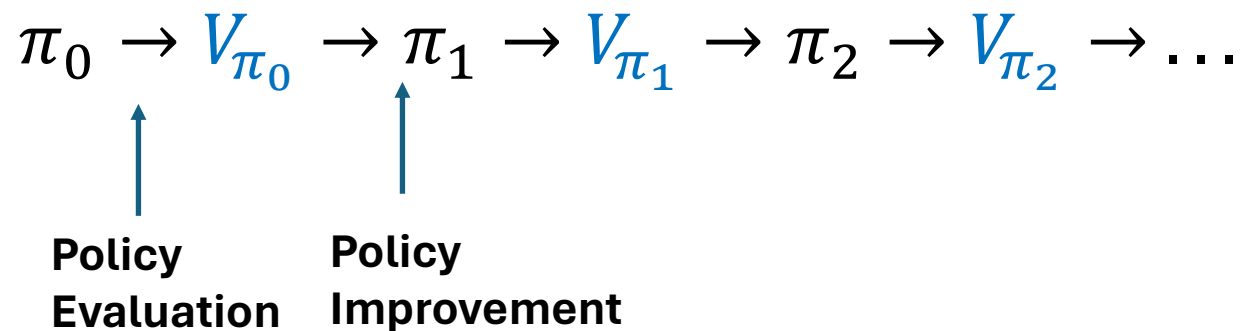
# Issue with Deterministic Policy

- Consider 3 states $A, B, C$

- 4 actions in each state: Left, Right, Up, Down

- Consider a deterministic policy
$$\pi(A) = Right$$
$$\pi(B) = Right$$
$$\pi(C) = Left$$

- Sample Episode

$A, Left, -1, B, Right, -1, A, Right, -1, C, Left, -1, \ldots$

# Generalized Policy Iteration (GPI)

**Policy Iteration**

$$\pi_0 \rightarrow V_{\pi_0} \rightarrow \pi_1 \rightarrow V_{\pi_1} \rightarrow \pi_2 \rightarrow V_{\pi_2} \rightarrow \dots$$

**Policy Evaluation**

**Policy Improvement**

**GPI**

$$\pi_0 \rightarrow Q_{\pi_0} \rightarrow \pi_1 \rightarrow Q_{\pi_1} \rightarrow \pi_2 \rightarrow Q_{\pi_2} \rightarrow \dots$$

**Approx. Policy Evaluation**

**$\epsilon-$greedy Policy Improvement**

$$\Pi_0 \xrightarrow[MC/TD]{} Q_{\Pi_0}$$

$$\rightarrow G_t = R_{t+1} + \gamma G_{t+1}$$

$\downarrow$ Immediate reward

$\downarrow$ remaining return

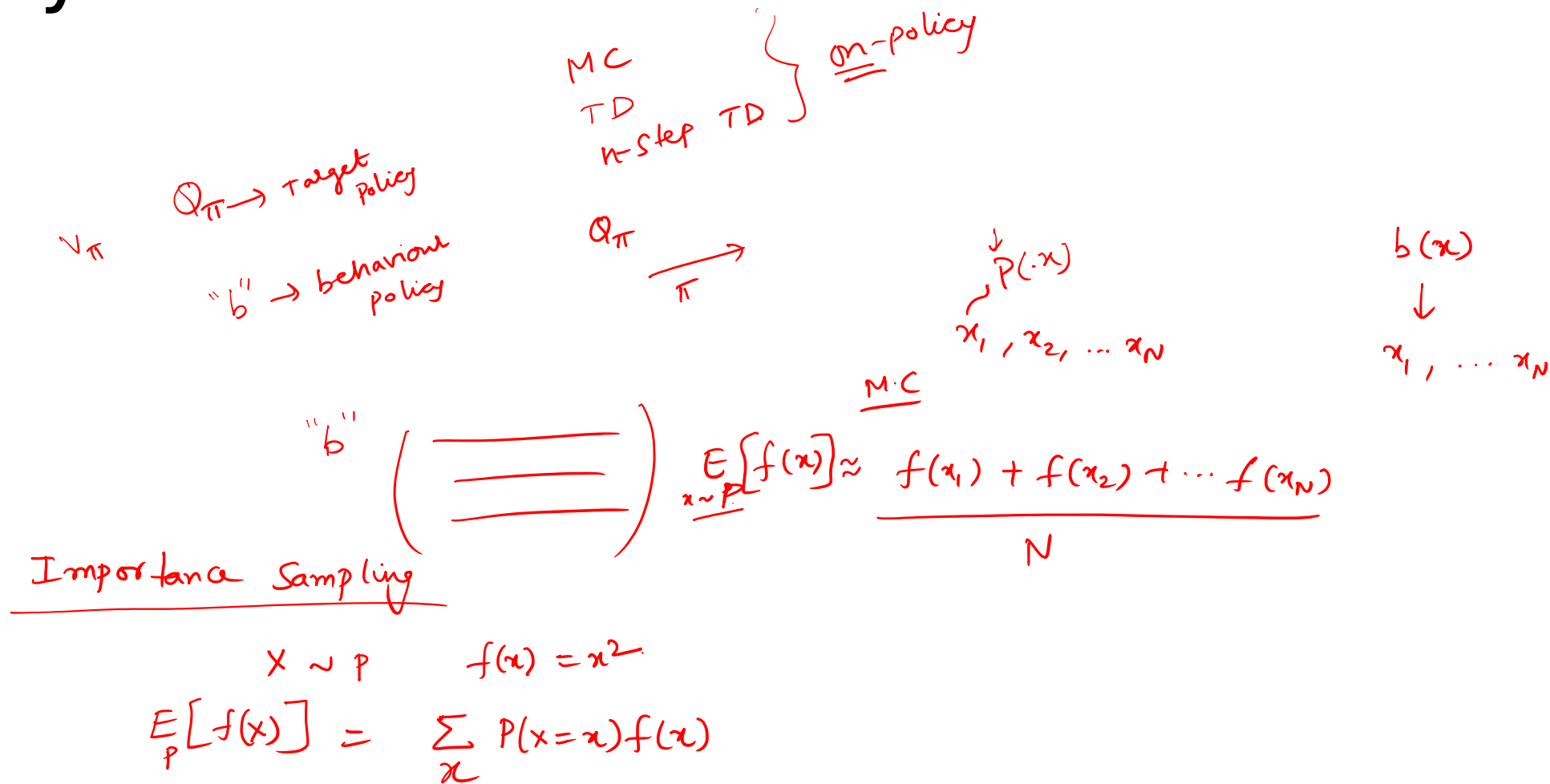$$Q_\pi(s,a) \approx R_{t+1} + \gamma Q_\pi(S_{t+1}, A_{t+1})$$

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

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

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

$MC$
$TD$
$n\text{-step } TD$ $\Big\}$ $\underline{on}$-policy

$Q_\pi \longrightarrow$ Target policy

$V_\pi$

"$b$" $\longrightarrow$ behaviour policy

$Q_\pi$

$\pi$

$\downarrow$
$P(\cdot x)$

$x_1, x_2, \dots x_N$

$b(x)$

$\downarrow$

$x_1, \dots x_N$

"$b$" $\left( \begin{array}{c} \underline{\phantom{xxxxx}} \\ \underline{\phantom{xxxxx}} \\ \underline{\phantom{xxxxx}} \end{array} \right)$ $\underset{x \sim P}{E}\{f(x)\} \approx$ $\dfrac{f(x_1) + f(x_2) + \cdots f(x_N)}{N}$ $\quad \underline{M.C}$

<u>Importance Sampling</u>

$X \sim P \qquad f(x) = x^2$

$\underset{P}{E}[f(x)] = \sum_x P(x = x) f(x)$

# Importance Sampling

$$E_P\left[f(x)\right] = \sum_x f(x)\, P(x=x)$$

$$= \sum_x f(x)\, b(x=x) \times \frac{P(x=x)}{b(x=x)}$$

$$= \sum_x \left(f(x)\cdot \frac{P(x=x)}{b(x=x)}\right) \underbrace{\phantom{xxx}}_{g(x)} b(x=x)$$

MC  $v_\pi \mid Q_\pi$

$b(\cdot)$

$x_1, \ldots, x_N$

$$\boxed{E_P[f(x)]} = \sum_x g(x)\, b(x=x) = E_b\left[g(x)\right] \approx \frac{g(x_1) + g(x_2) + \cdots g(x_N)}{N}$$

$$= \frac{1}{N}\left(f(x_1)\frac{P(x_1)}{b(x_1)} + f(x_2)\frac{P(x_2)}{b(x_2)} + \cdots\right)$$

Importance Sampling Tech.

MC for $v_\pi$.

$$v_\pi(s) =$$

EP1     $S_0 = s$,   $A_0 \sim \pi$,   $R_1$, $S_1$, $A_1 \sim \pi$, . . . . .    $G^{(1)}$

EP2     $G^{(2)}$

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

behaviour $b$.

EP1     $S_0 = s$,   $A_0 \sim b$,   $R_1$   $S_1$,   $A_1 \sim b$, . . . .    $G^{(1)}$

EP2     $G^{(N)}$

$$v_\pi(s) \approx \frac{1}{N}\left( \left( \frac{G^{(1)} \cdot \frac{P_\pi(\text{EP1})}{}}{P_b(\text{EP1})} \right) + \left( \frac{G^{(2)} \, P_\pi(\text{EP2})}{P_b(\text{EP1})} \right) + \cdots \right)$$

$$\frac{P_\pi(\text{EP1})}{P_b(\text{EP1})} = \frac{\pi(A_0|S_0)\,\pi(A_1|S_1)\,\pi(A_2|S_2)\cdots}{b(A_0|S_0)\,b(A_1|S_1)\,b(A_2|S_2)\cdots}$$