

Sparse Tree Search Optimality Guarantees in POMDPs with Continuous Observation Spaces

Tianci Li

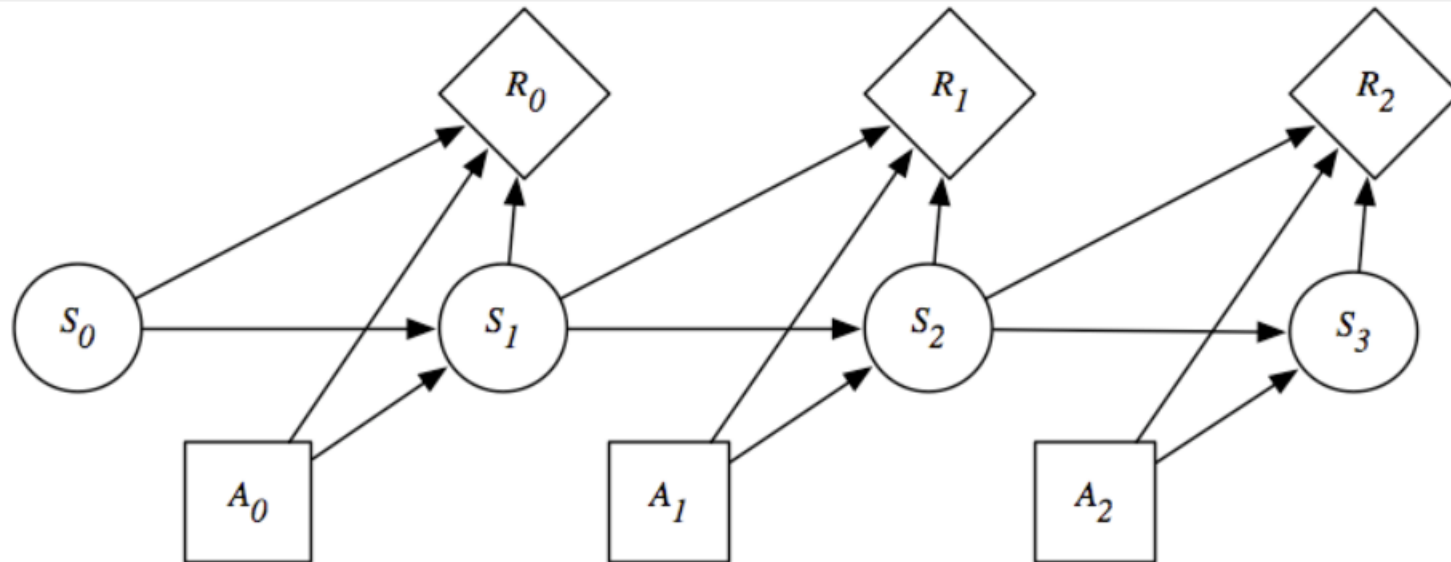
2021/8/27



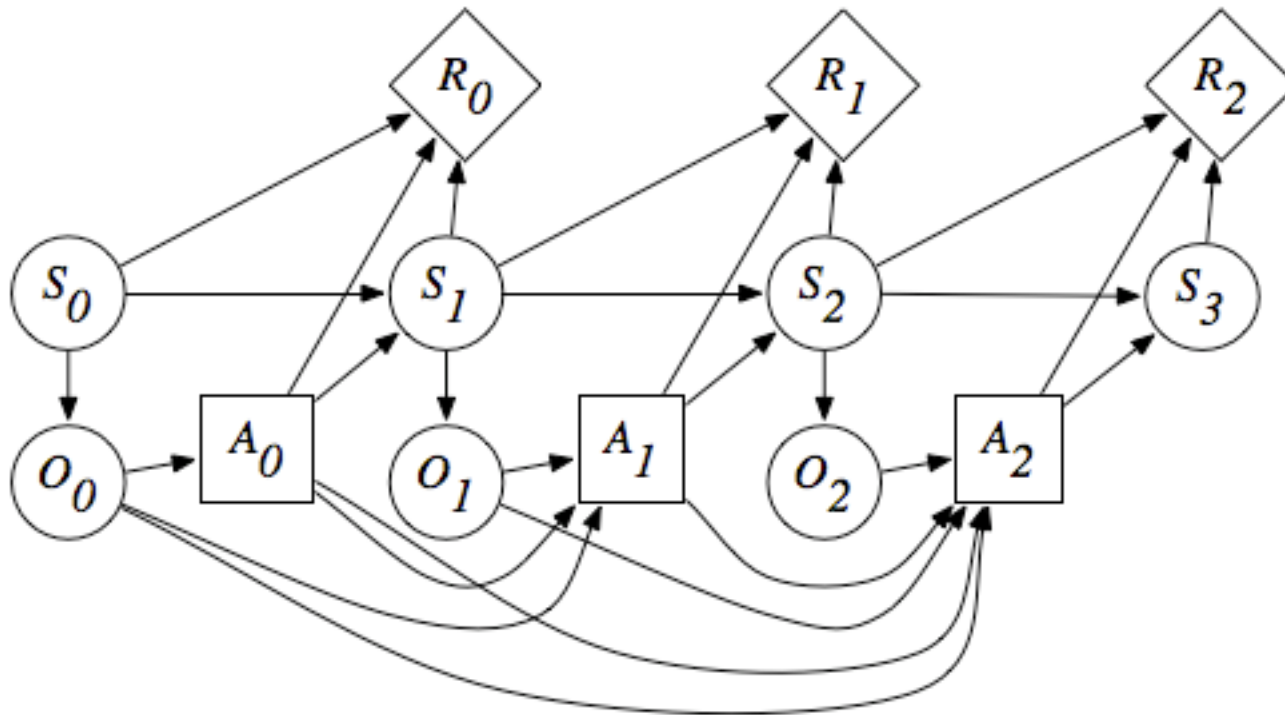
Outline

- Introduction
 - Definition of POMDP
 - Sparse Tree Search Algorithms in POMDP
 - Monte-Carlo Tree Search
 - Particle Filtering
 - Progressive widening
 - Importance Sampling
- Sparse Tree Search Optimality Guarantees

Introduction



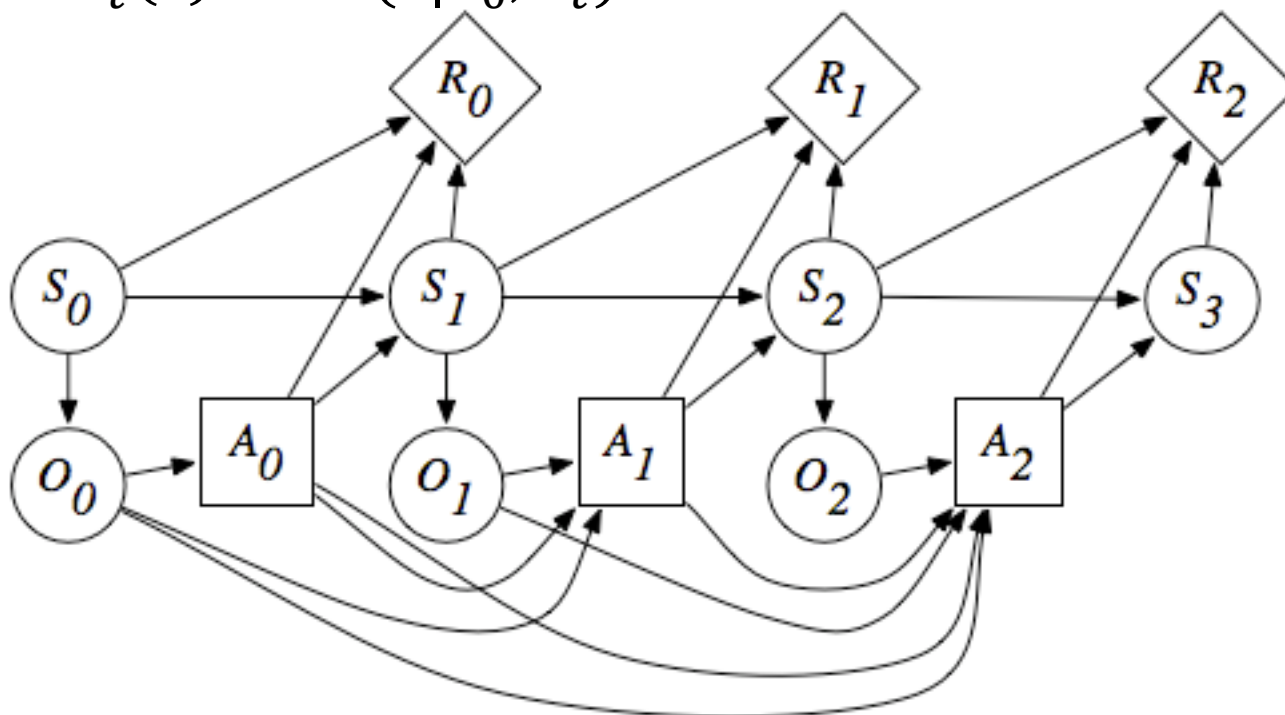
Introduction



Introduction

历史信息 : $h_t = \{a_0, o_1, a_2, \dots, a_{t-1}, o_t\}$

信念状态 : $b_t(s) = Pr(s|b_0, h_t)$



POMDP \rightarrow Belief MDP

Definition of POMDP

A POMDP can formally be described as a 7-tuple $P = (\mathcal{S}, \mathcal{A}, T, R, \Omega, O, \gamma)$,

- $\mathcal{S} = \{s_1, s_2, \dots, s_n\}$ is a set of states,
- $\mathcal{A} = \{a_1, a_2, \dots, a_m\}$ is a set of actions,
- T is a set of conditional transition probabilities $T(s' | s, a)$ for the state transition $s \rightarrow s'$.
- $R: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the reward function,
- $\Omega = \{o_1, o_2, \dots, o_k\}$ is a set of observations,
- O is a set of observation probabilities $O(o|s', a)$
- $\gamma \in [0,1]$ is the discount factor.

Some concepts and equations

历史信息： $h_t = \{a_0, o_1, a_2, \dots, a_{t-1}, o_t\}$

信念状态： $b_t(s) = \Pr(s|b_0, h_t)$

贝叶斯更新： $b_t(s') = \Pr(s'|a_{t-1}, o_t, b_{t-1})$

多变量贝叶斯公式：

$$\Pr(A|B, C) = \frac{\Pr(B|A, C) \Pr(A|C)}{\Pr(B|C)}$$



$O(|S|^2)$

$$\begin{aligned} &= \frac{\Pr(o_t|s', a_{t-1}, b_{t-1}) \Pr(s'|a_{t-1}, b_{t-1})}{\Pr(o_t|a_{t-1}, b_{t-1})} \\ &= \frac{O(o_t|s', a_{t-1}) \sum_{s \in S} T(s'|s, a_{t-1}) b_{t-1}(s)}{\sum_{s' \in S} O(o_t|s', a_{t-1}) \sum_{s \in S} T(s'|s, a_{t-1}) b_{t-1}(s)} \end{aligned}$$

计算成本高！使用粒子滤波

Some concepts and equations

累计奖励 : $G_t = R_t + \gamma R_{t+1} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$

策略 : $\forall s \in S, \forall a \in A, \pi(a|s) = \Pr(a|s)$

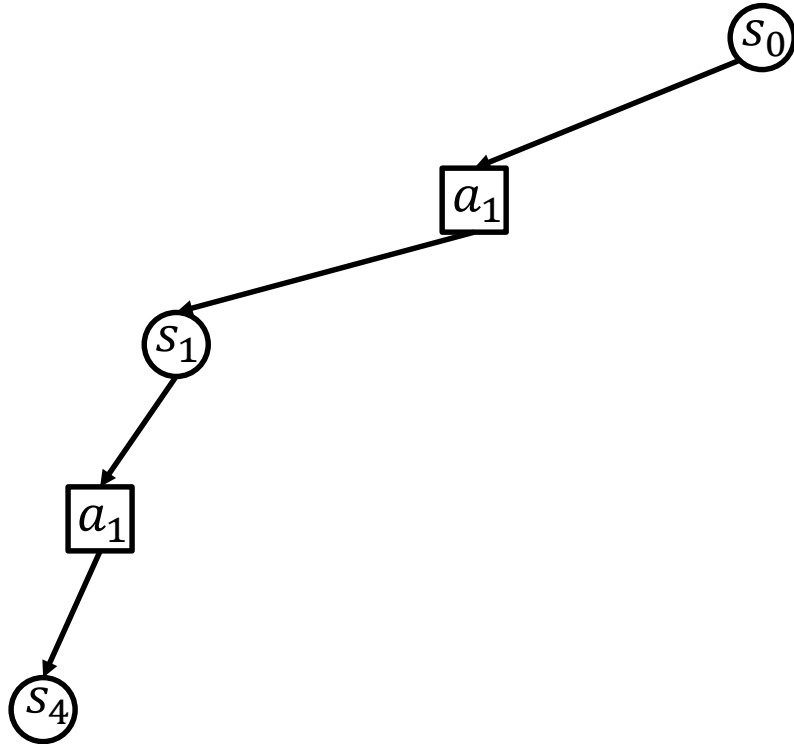
状态值函数 : $V_{\pi}(s) = E[G_t | s_t = s] = E[R_t + \gamma V_{\pi}(s_{t+1}) | s_t = s]$ (贝尔曼方程)

动作值函数 : $Q_{\pi}(s, a) = R(s, a) + \gamma \sum_{s' \in S} P(s' | s, a) V_{\pi}(s')$

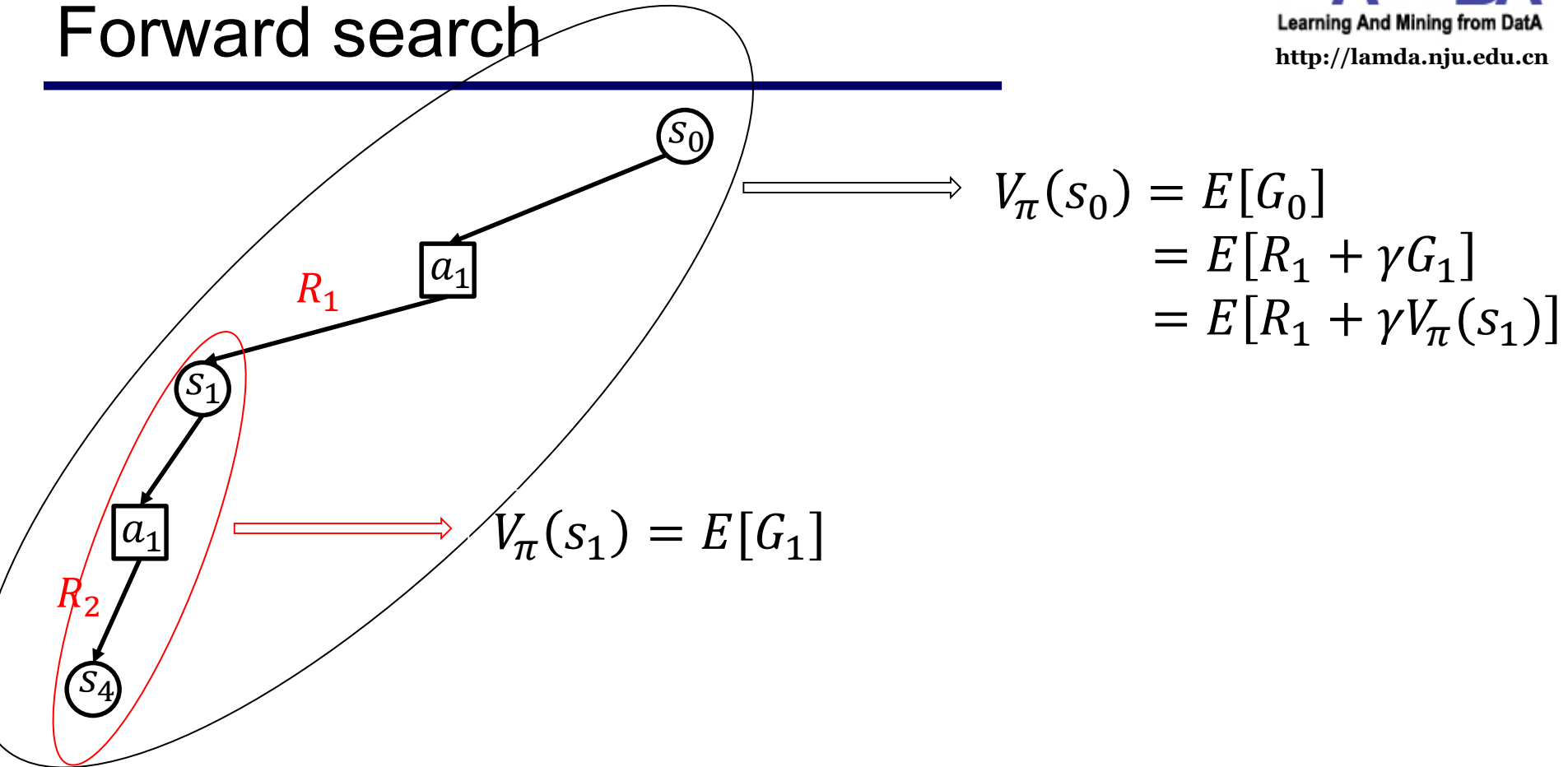
贝尔曼最优方程 : $V^*(s) = \max_a Q(s, a)$

最优策略 : $\pi^*(a|s) = \begin{cases} 1, & \text{if } a \in \underset{a}{\operatorname{argmax}} Q(s, a) \\ 0, & \text{otherwise} \end{cases}$

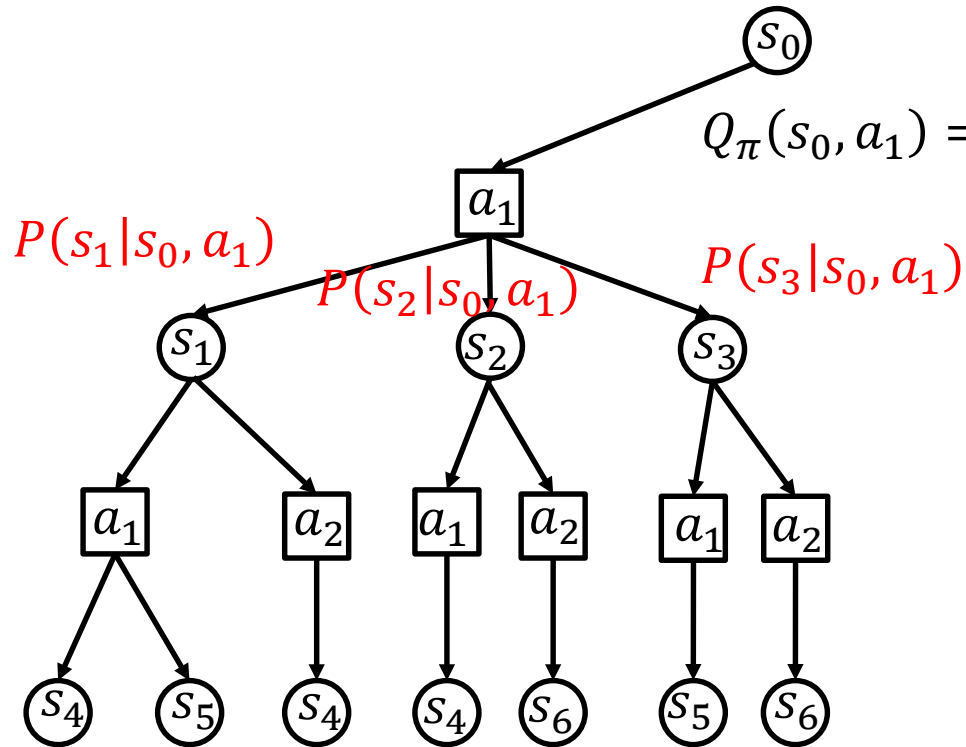
Forward search



Forward search

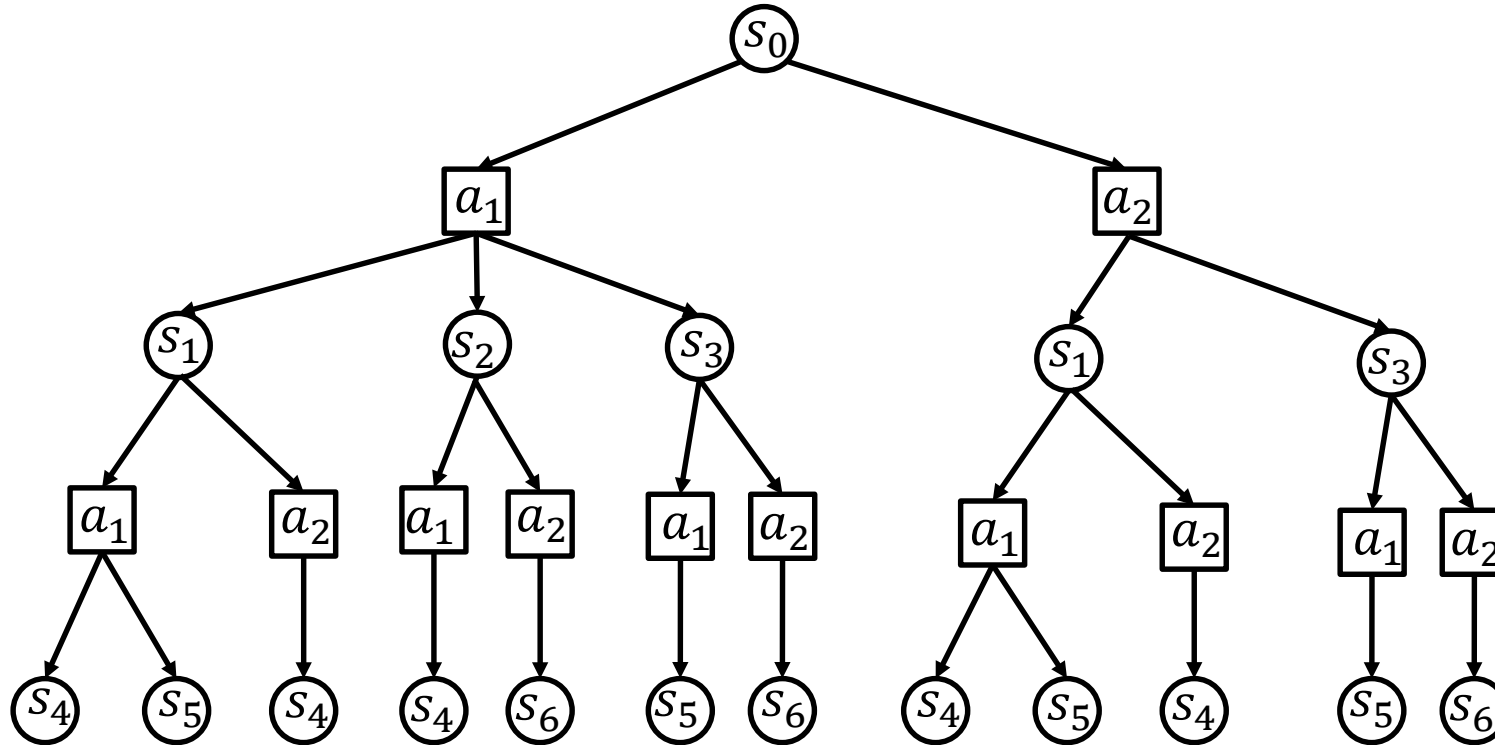


Forward search

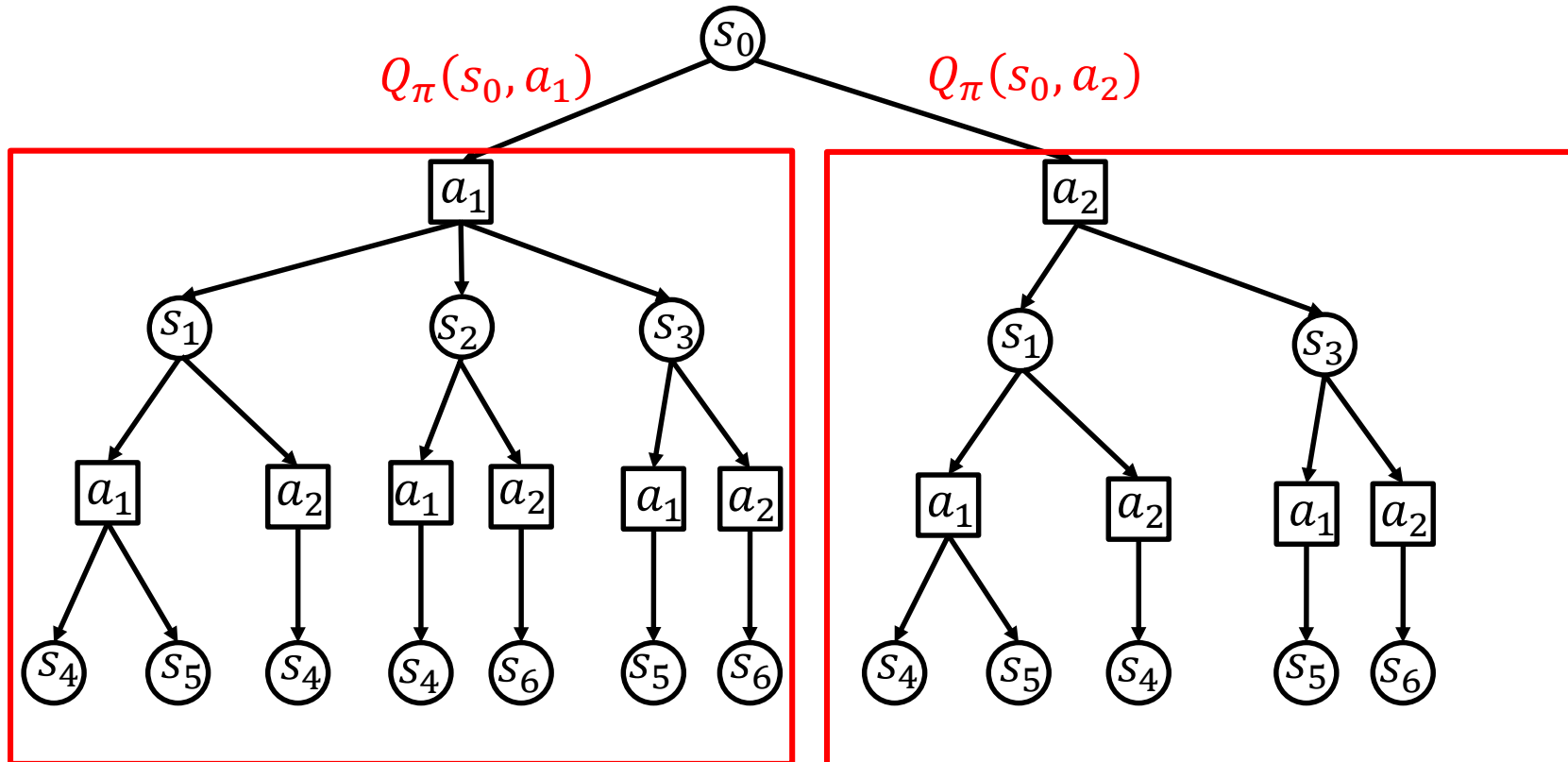


$$Q_{\pi}(s_0, a_1) = R(s_0, a_1) + \gamma \sum_{s' \in S} P(s' | s_0, a_1) V_{\pi}(s')$$

Forward search

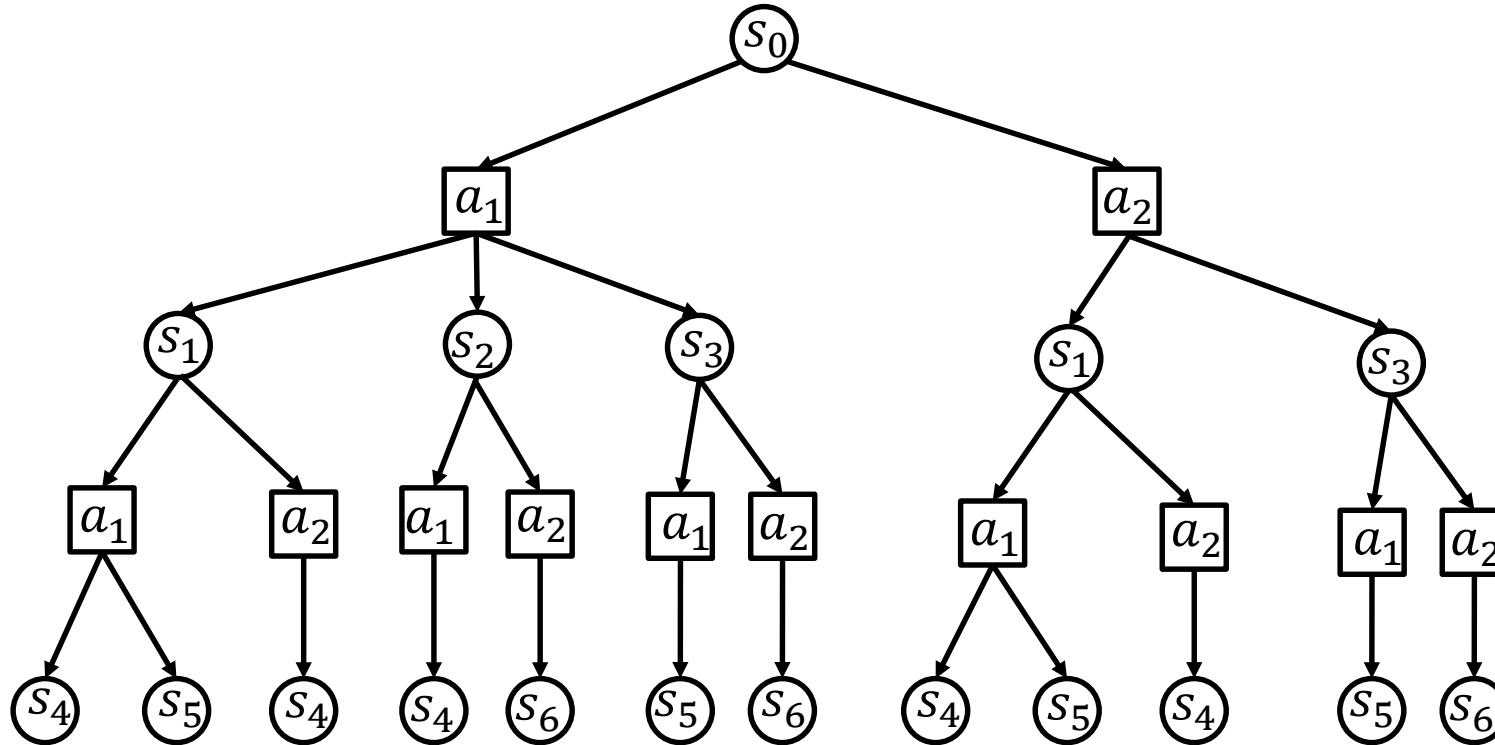


Forward search

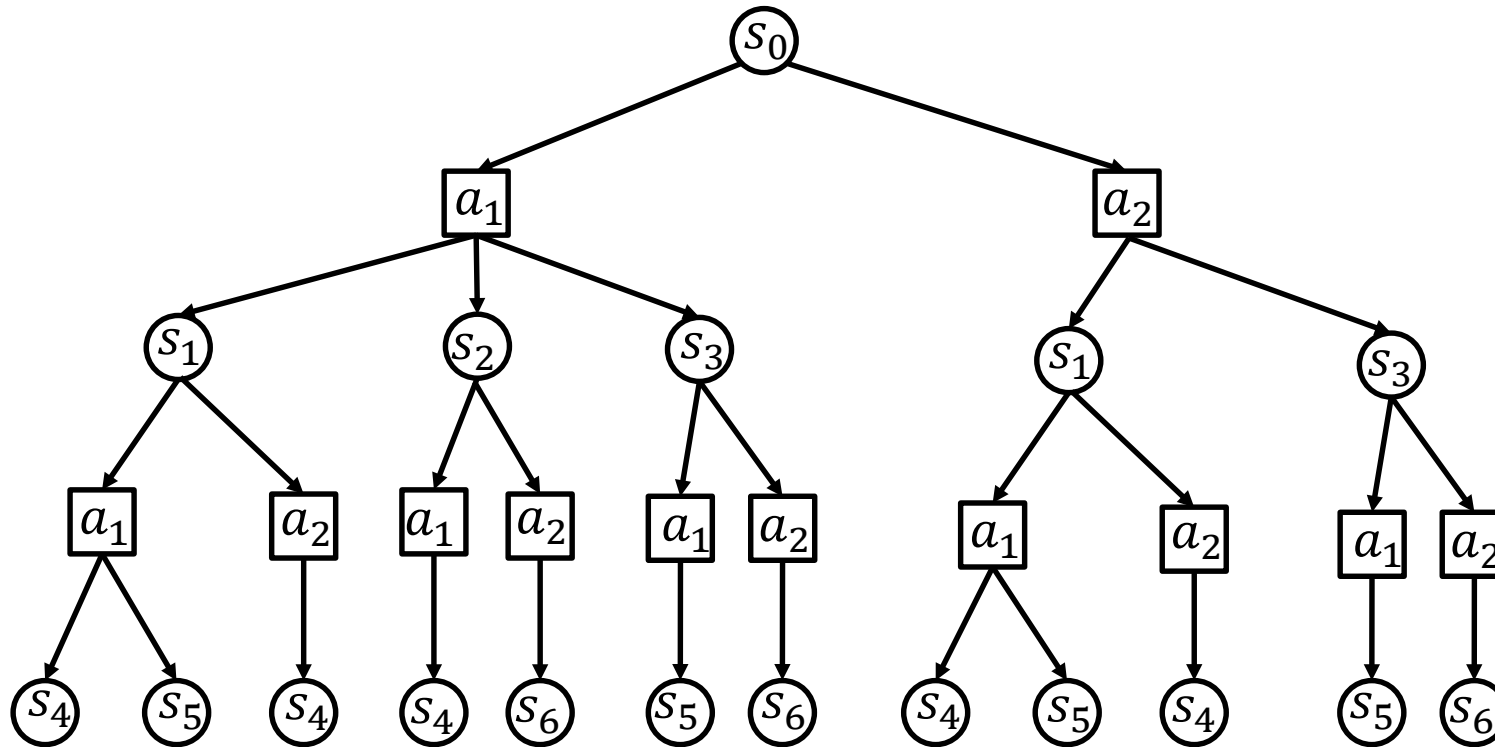


$$V_\pi(s_0) = \sum_{a \in A} \pi(a|s_0) Q_\pi(s_0, a) \leq \max_a Q(s_0, a) \quad V^*(s_0) = \max_a Q^*(s_0, a)$$

Forward search



Forward search



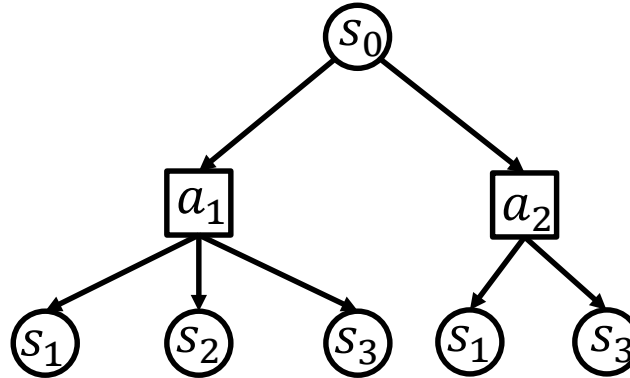
When the search space is large

Forward search ✗

Monte-Carlo Tree Search ✓

Monte-Carlo Tree Search

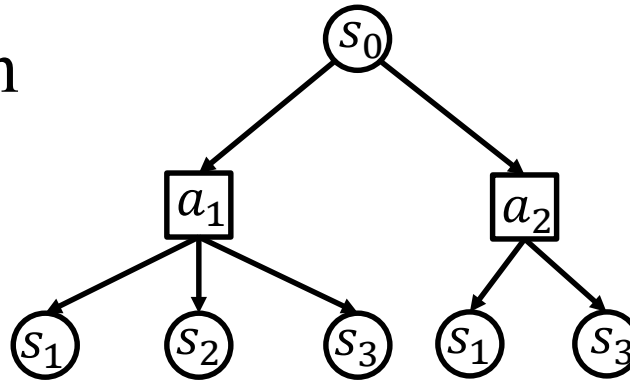
Each node saves the **N value**(visited counts) and **V value**(Node value)



Monte-Carlo Tree Search

Each node saves the **N value**(visited counts) and **V value**(Node value)

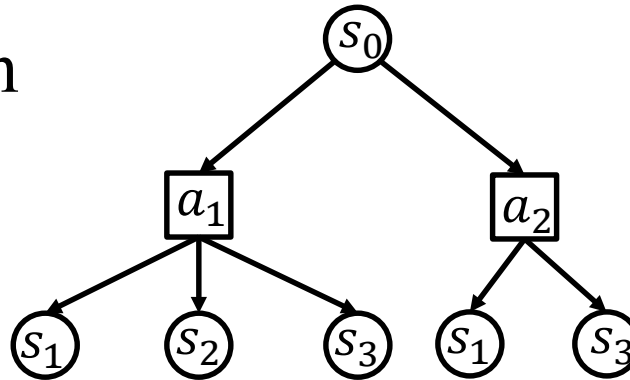
1. Selection



Monte-Carlo Tree Search

Each node saves the **N value**(visited counts) and **V value**(Node value)

1. Selection



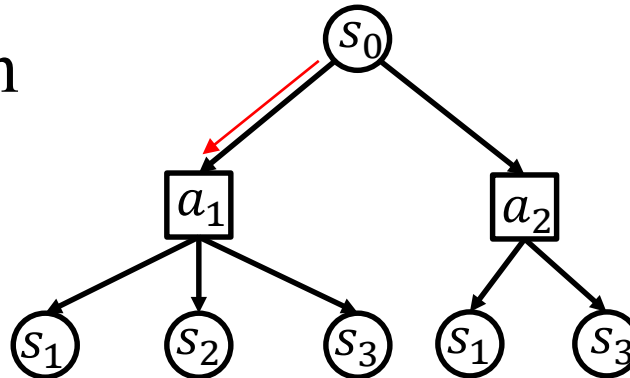
A tree policy that is used while within the search tree. **Greedy tree policy ; UCT**

$$\text{UCT: } A = \underset{a}{\operatorname{argmax}} \left[V(sa) + c \sqrt{\frac{\log N(s)}{N(sa)}} \right]$$

Monte-Carlo Tree Search

Each node saves the **N value**(visited counts) and **V value**(Node value)

1. Selection



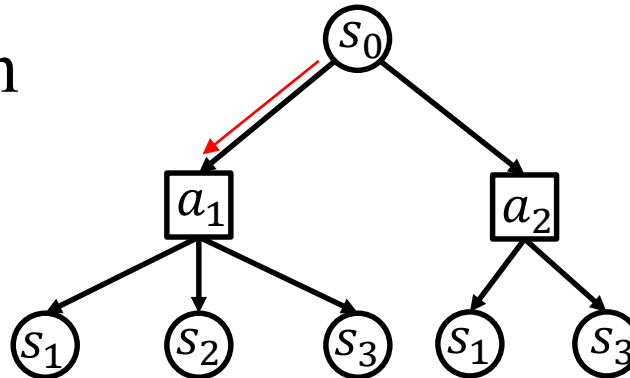
A tree policy that is used while within the search tree. **Greedy tree policy ; UCT**

$$\text{UCT: } A = \underset{a}{\operatorname{argmax}} \left[V(sa) + c \sqrt{\frac{\log N(s)}{N(sa)}} \right]$$

Monte-Carlo Tree Search

Each node saves the **N value**(visited counts) and **V value**(Node value)

1. Selection



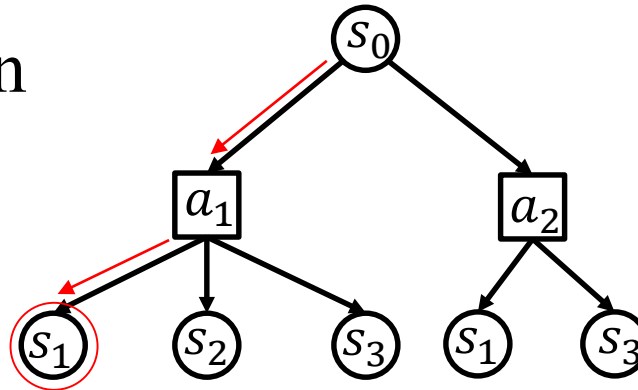
Generative model to generate sequences of states and rewards

$$G(s_t, a_t) = (R_{t+1}, s_{t+1})$$

Monte-Carlo Tree Search

Each node saves the **N value**(visited counts) and **V value**(Node value)

1. Selection

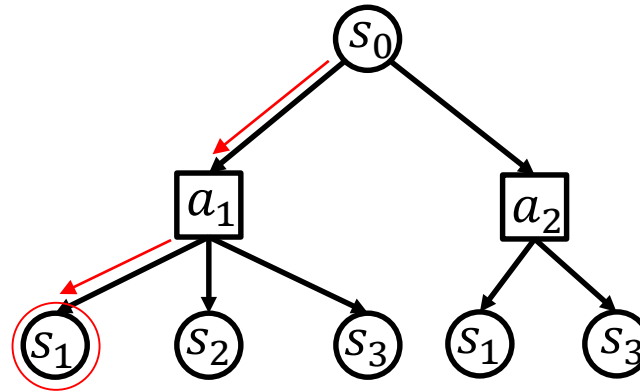


Generative model to generate sequences of states and rewards

$$G(s_t, a_t) = (R_{t+1}, s_{t+1})$$

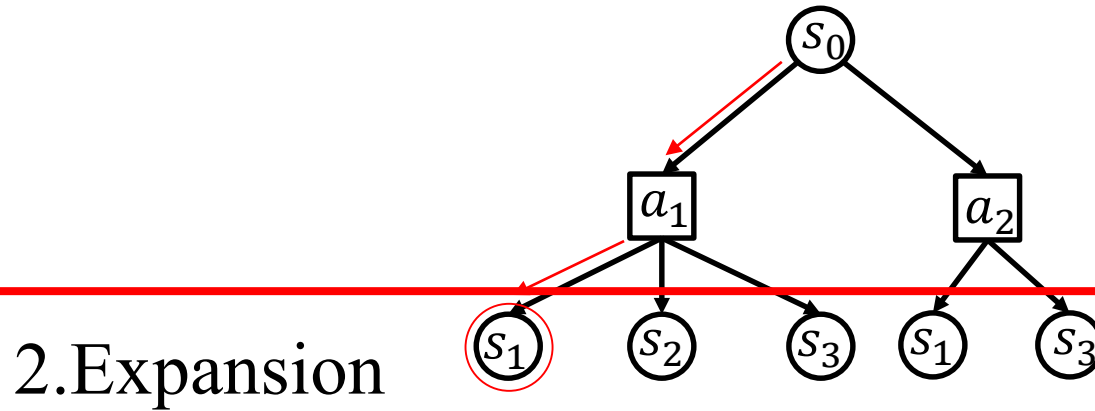
Monte-Carlo Tree Search

Each node saves the **N value**(visited counts) and **V value**(Node value)



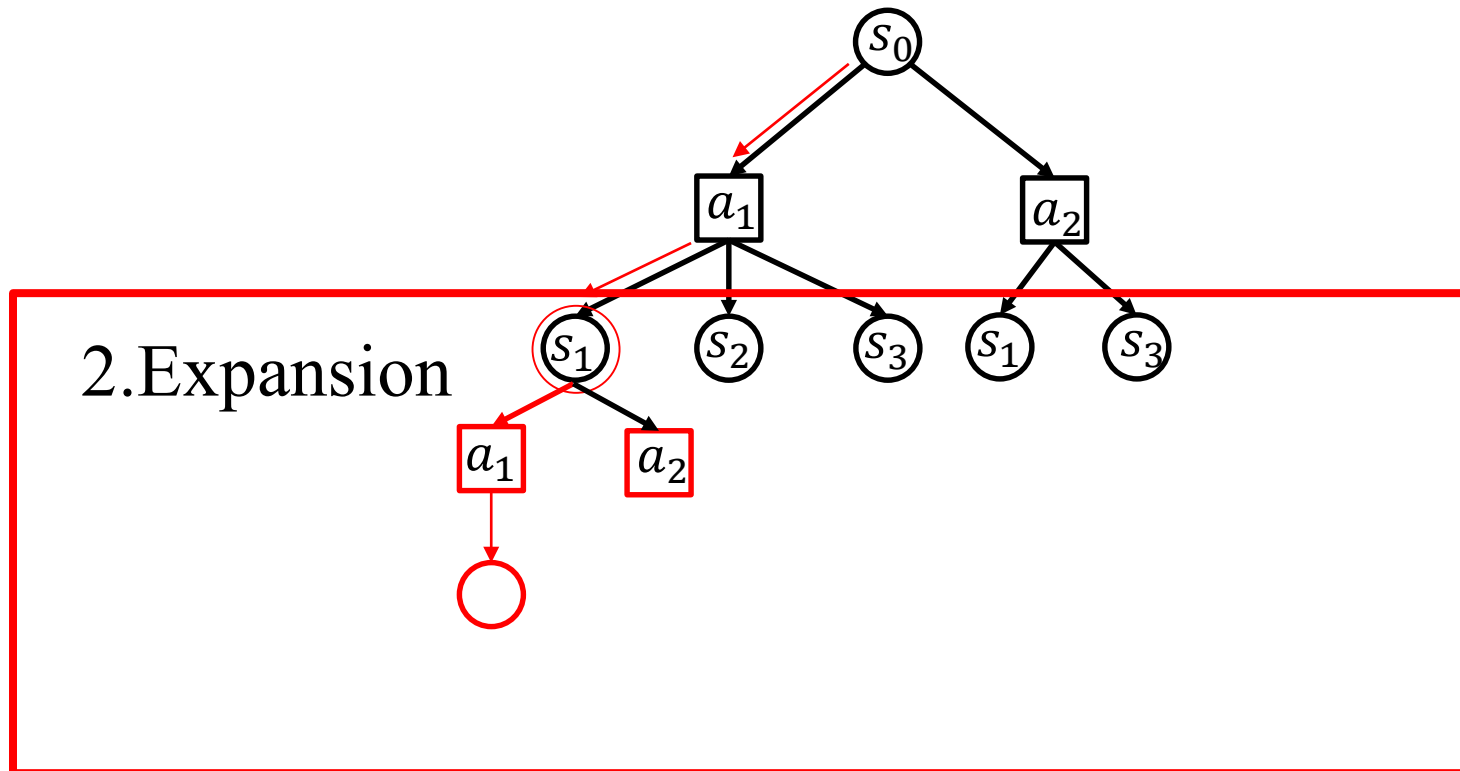
Monte-Carlo Tree Search

Each node saves the **N value**(visited counts) and **V value**(Node value)



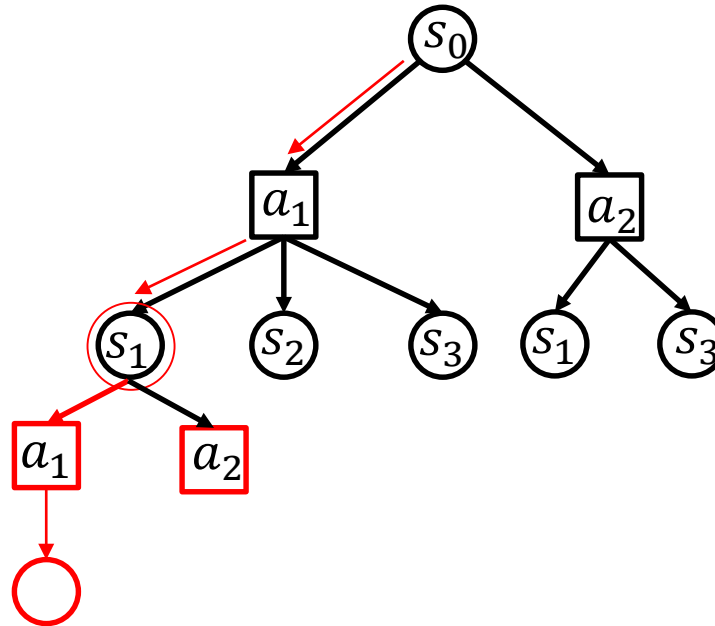
Monte-Carlo Tree Search

Each node saves the **N value**(visited counts) and **V value**(Node value)



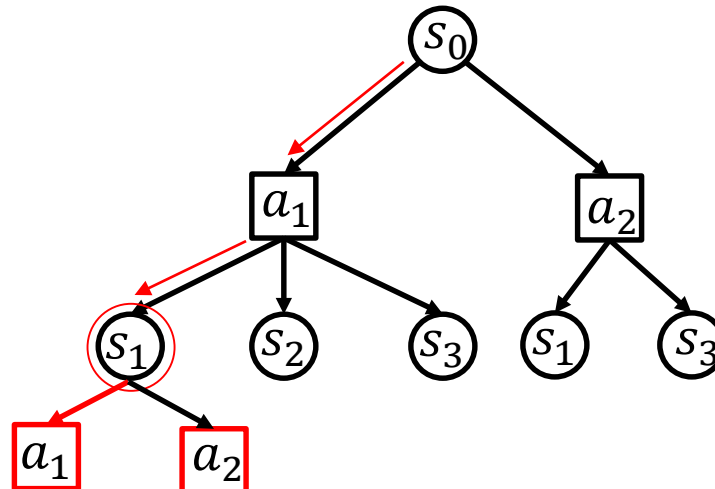
Monte-Carlo Tree Search

Each node saves the **N value**(visited counts) and **V value**(Node value)



Monte-Carlo Tree Search

Each node saves the **N value**(visited counts) and **V value**(Node value)



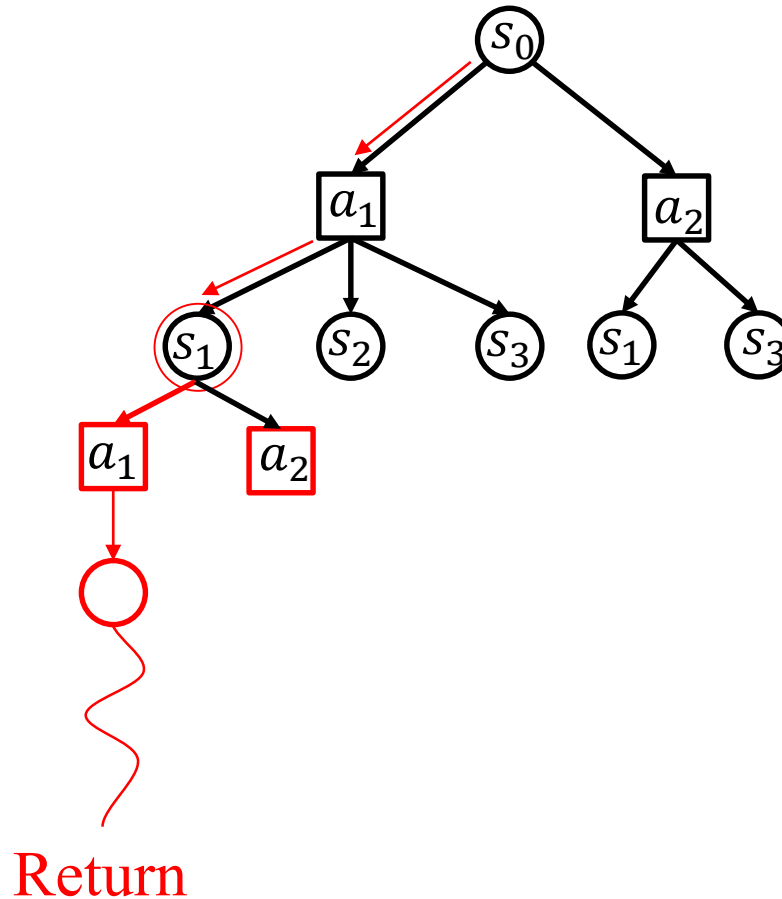
3. Simulation

Return

A rollout policy that is used once simulations leave the scope of the search tree

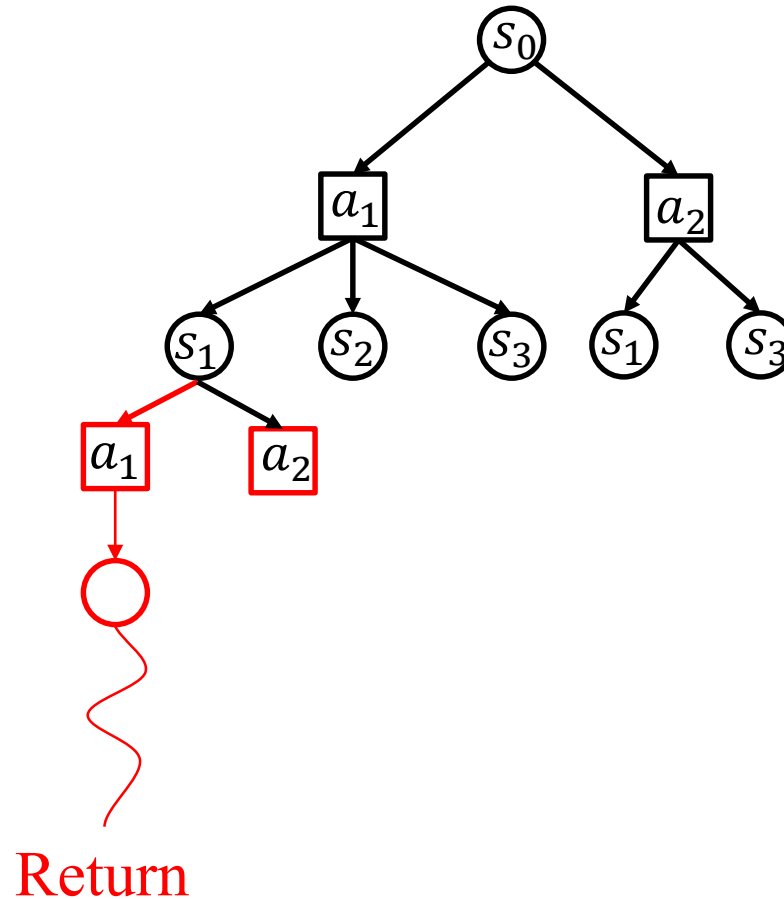
Monte-Carlo Tree Search

Each node saves the **N value**(visited counts) and **V value**(Node value)



Monte-Carlo Tree Search

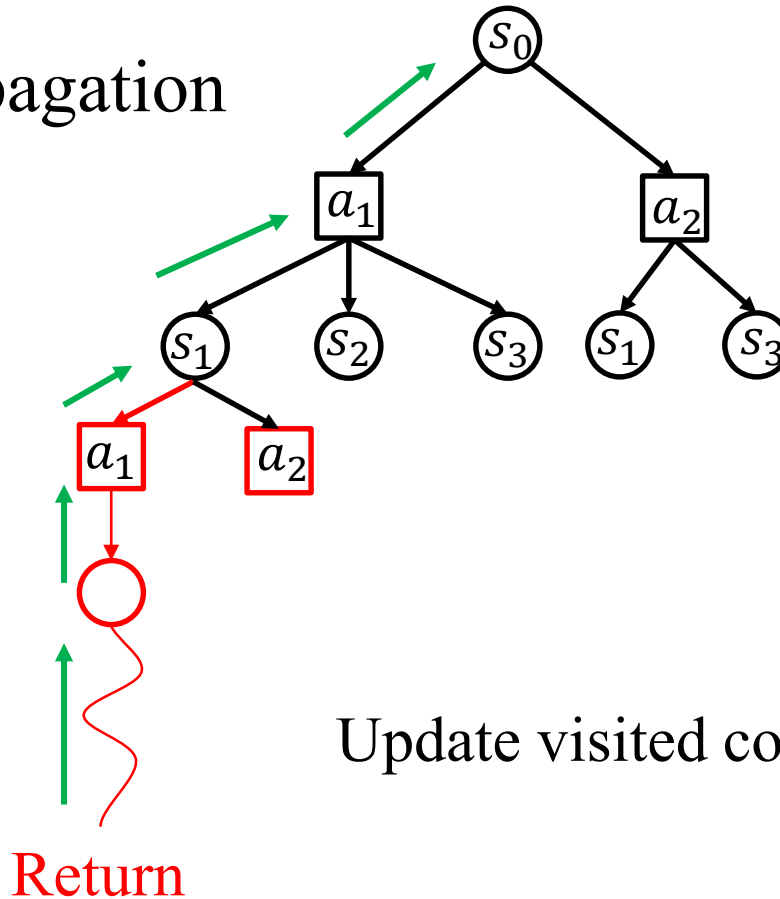
Each node saves the **N value**(visited counts) and **V value**(Node value)



Monte-Carlo Tree Search

Each node saves the **N value(visited counts)** and **V value(Node value)**

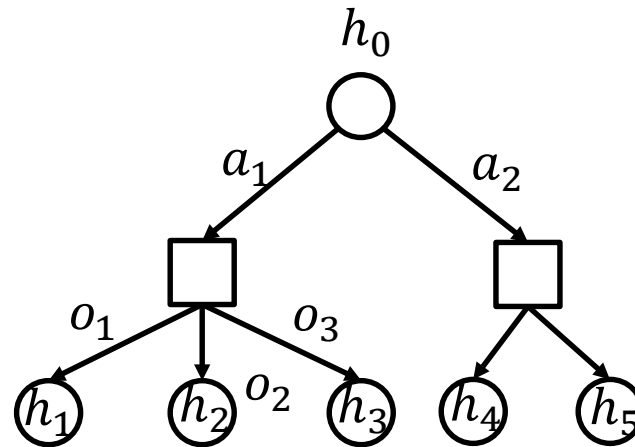
4. Backpropagation



Update visited counts and node value

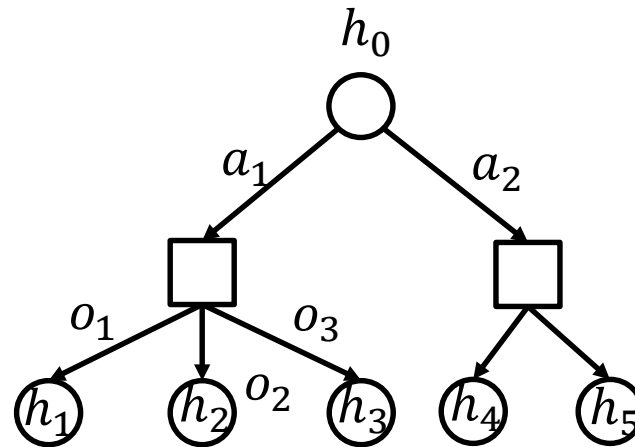
Monte-Carlo Tree Search

In POMCP:



Monte-Carlo Tree Search

In POMCP:

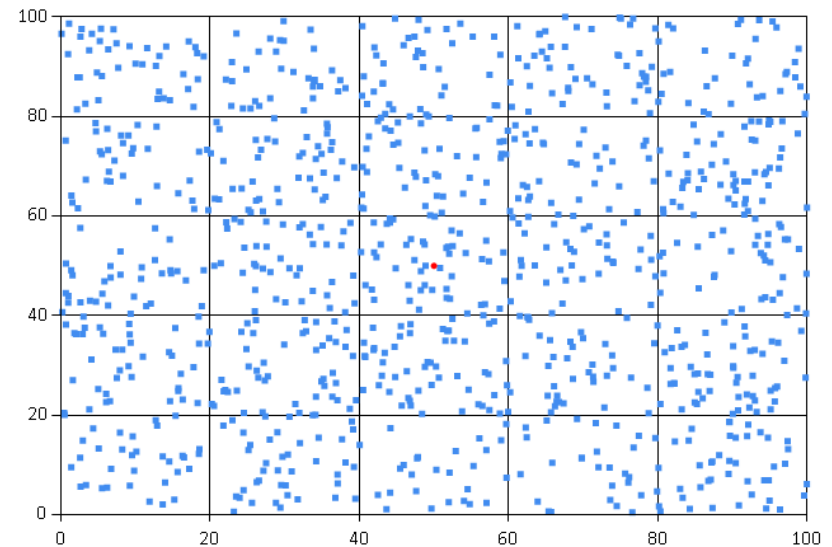
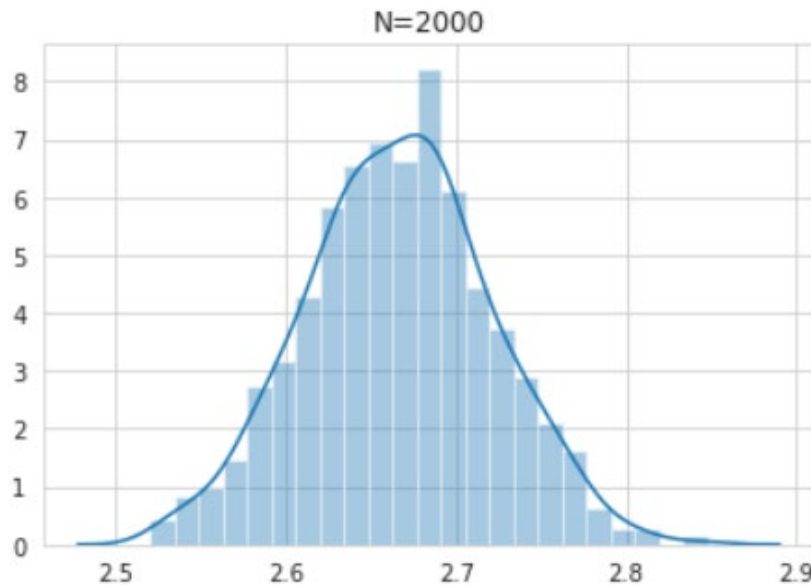


PO-UCT:
$$A = \underset{a}{\operatorname{argmax}} \left[V(ha) + c \sqrt{\frac{\log N(h)}{N(ha)}} \right]$$

Generative model:
$$G(s_t, a_t) = (R_{t+1}, o_{t+1}, s_{t+1})$$

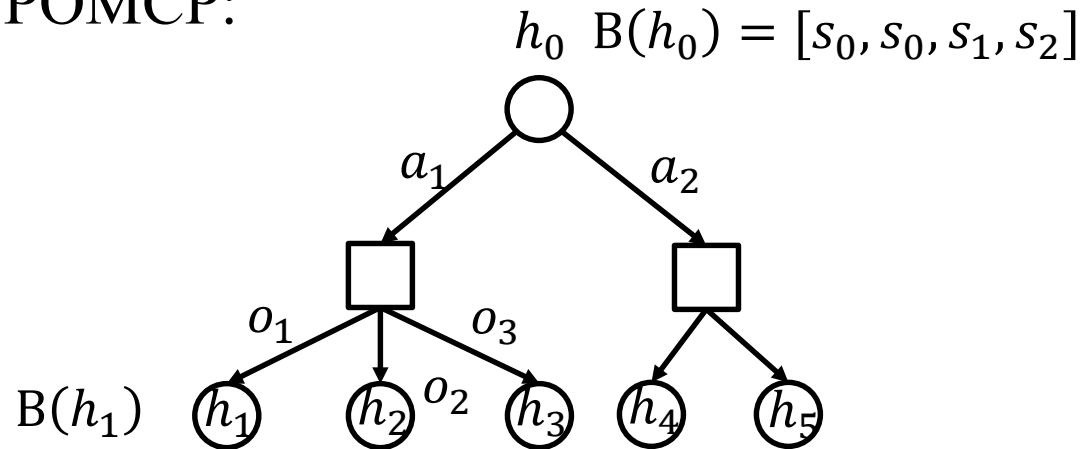
Particle Filtering

Using **a set of particles** to represent the posterior distribution.



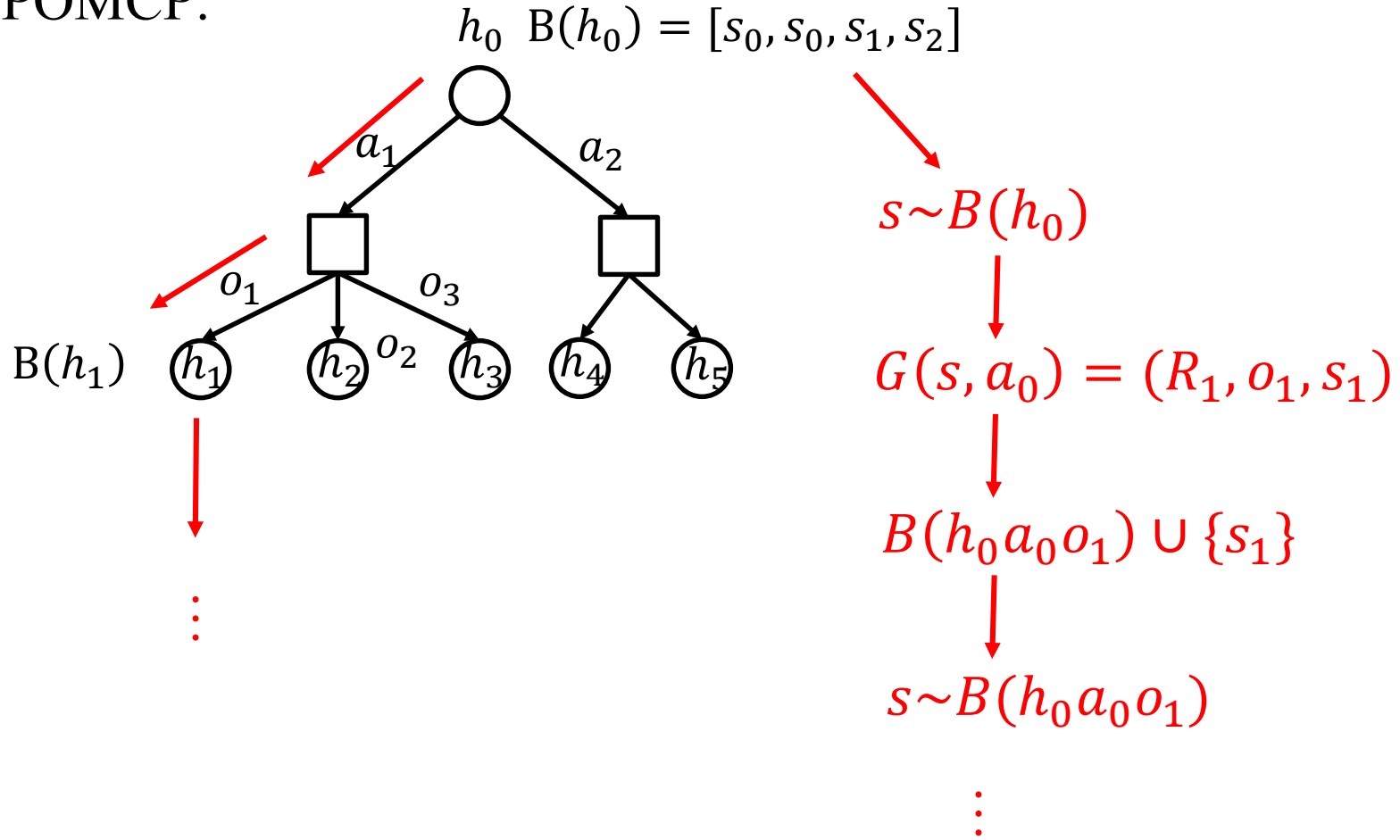
Particle Filtering

In POMCP:



Particle Filtering

In POMCP:



Partially Observable Monte-Carlo

Algorithm 1 Partially Observable Monte-Carlo Planning

```

procedure SEARCH( $h$ )
  repeat
    if  $h = \text{empty}$  then
       $s \sim \mathcal{I}$ 
    else
       $s \sim B(h)$ 
    end if
    SIMULATE( $s, h, 0$ )
  until TIMEOUT()
  return  $\underset{b}{\operatorname{argmax}} V(hb)$ 
end procedure

```

```

procedure ROLLOUT( $s, h, \text{depth}$ )
  if  $\gamma^{\text{depth}} < \epsilon$  then
    return 0
  end if
   $a \sim \pi_{\text{rollout}}(h, \cdot)$ 
   $(s', o, r) \sim \mathcal{G}(s, a)$ 
  return  $r + \gamma \cdot \text{ROLLOUT}(s', hao, \text{depth}+1)$ 
end procedure

```

```

procedure SIMULATE( $s, h, \text{depth}$ )
  if  $\gamma^{\text{depth}} < \epsilon$  then
    return 0
  end if
  if  $h \notin T$  then
    for all  $a \in \mathcal{A}$  do
       $T(ha) \leftarrow (N_{\text{init}}(ha), V_{\text{init}}(ha), \emptyset)$ 
    end for
    return ROLLOUT( $s, h, \text{depth}$ )
  end if
   $a \leftarrow \underset{b}{\operatorname{argmax}} V(hb) + c \sqrt{\frac{\log N(h)}{N(hb)}}$ 
   $(s', o, r) \sim \mathcal{G}(s, a)$ 
   $R \leftarrow r + \gamma \cdot \text{SIMULATE}(s', hao, \text{depth} + 1)$ 
   $B(h) \leftarrow B(h) \cup \{s\}$ 
   $N(h) \leftarrow N(h) + 1$ 
   $N(ha) \leftarrow N(ha) + 1$ 
   $V(ha) \leftarrow V(ha) + \frac{R - V(ha)}{N(ha)}$ 
  return  $R$ 
end procedure

```

POMCP in Continuous Space

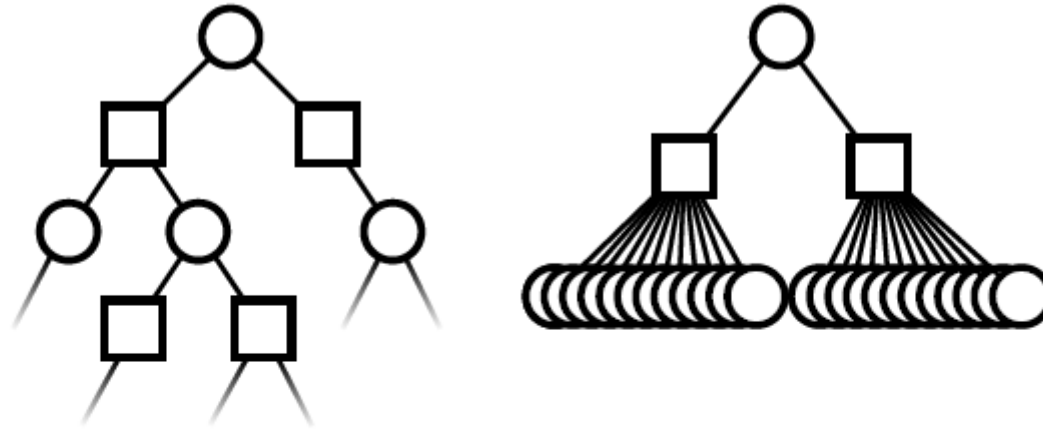
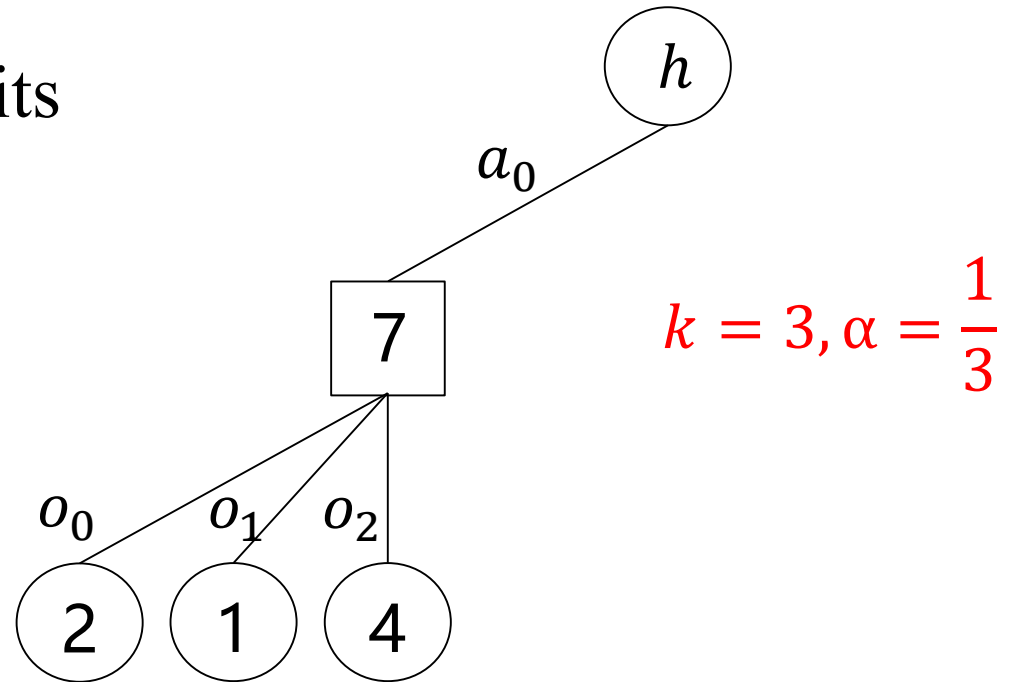


Figure 1: POMCP tree for a discrete POMDP (left), and for a POMDP with a continuous observation space (right). Because the observation space is continuous, each simulation creates a new observation node and the tree cannot extend deeper.

Progressive widening

- Limit the number of child nodes: $\leq k N^\alpha$

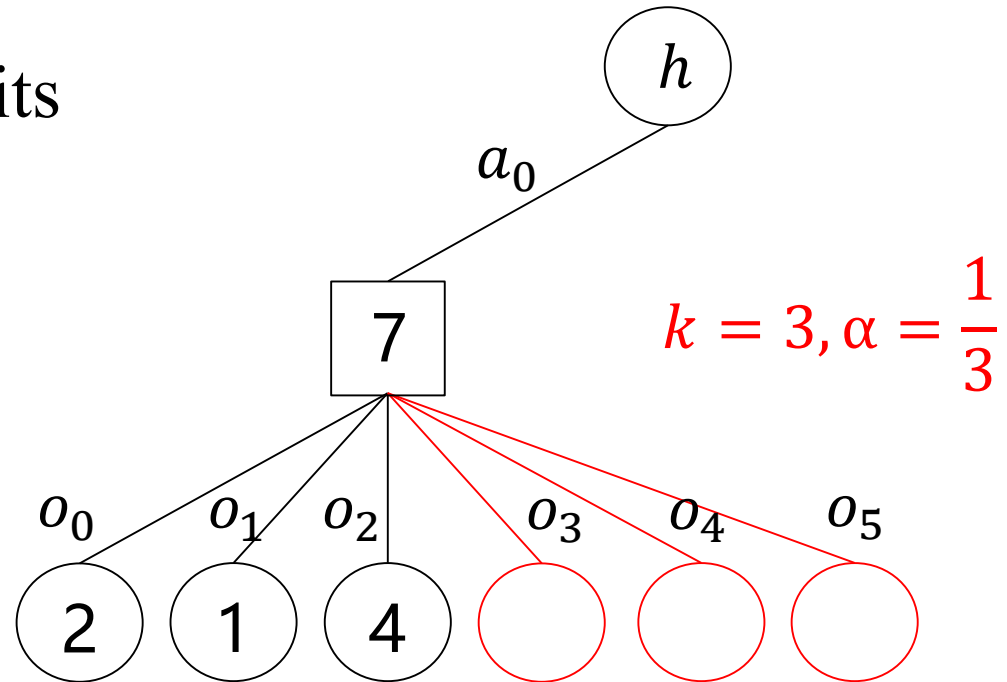
N is the number of visits



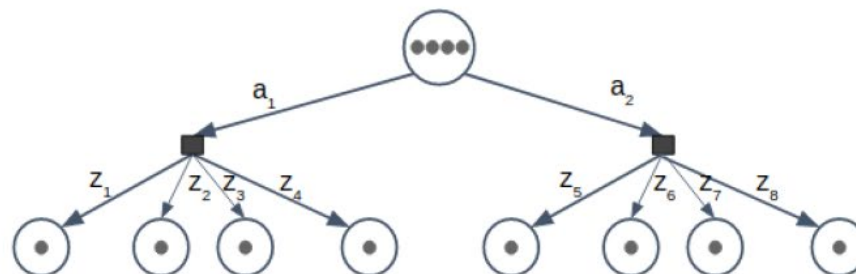
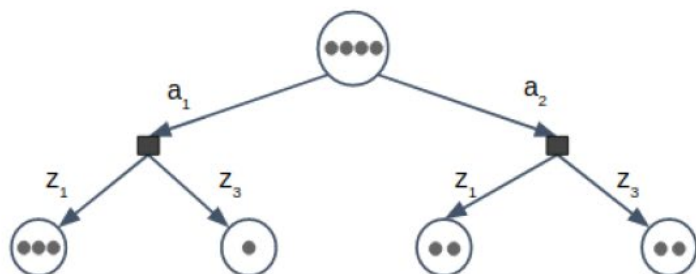
Progressive widening

- Limit the number of child nodes: $\leq k N^\alpha$

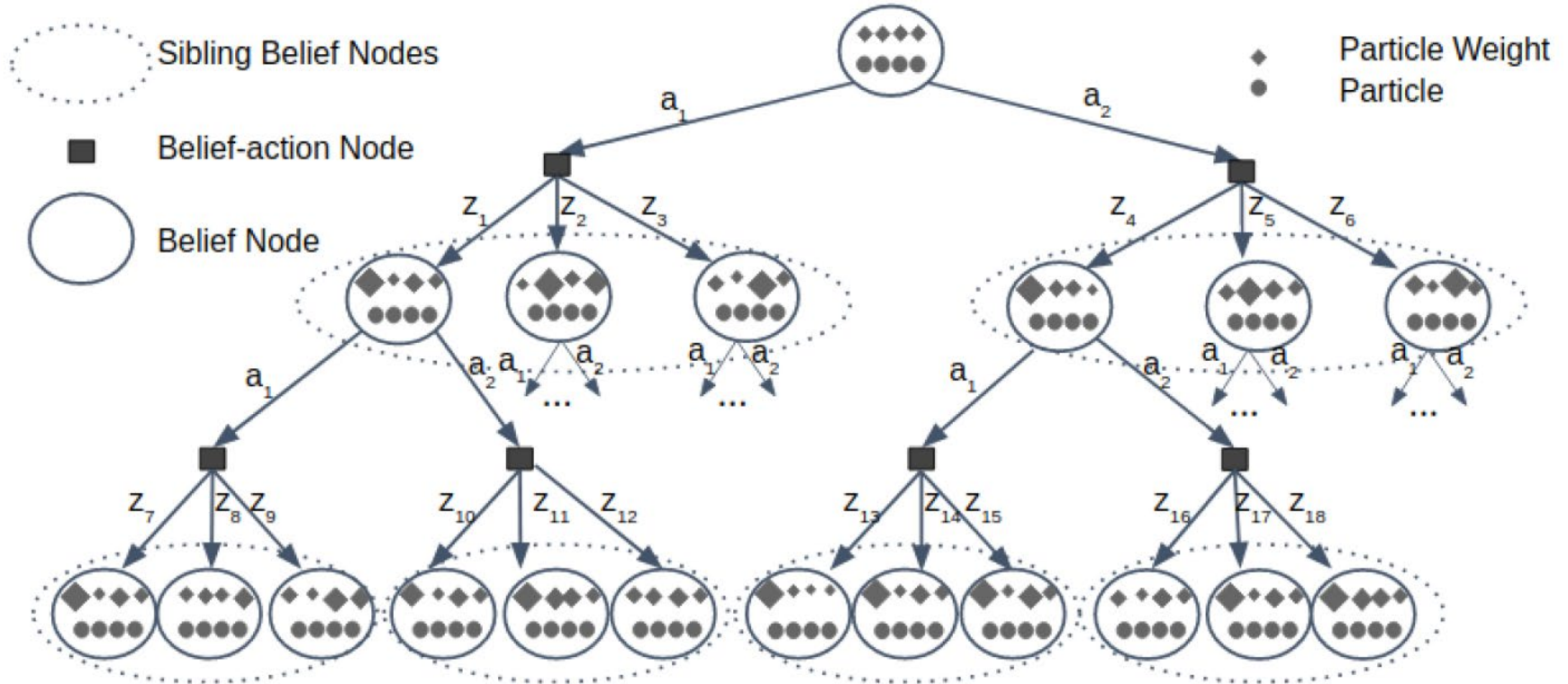
N is the number of visits



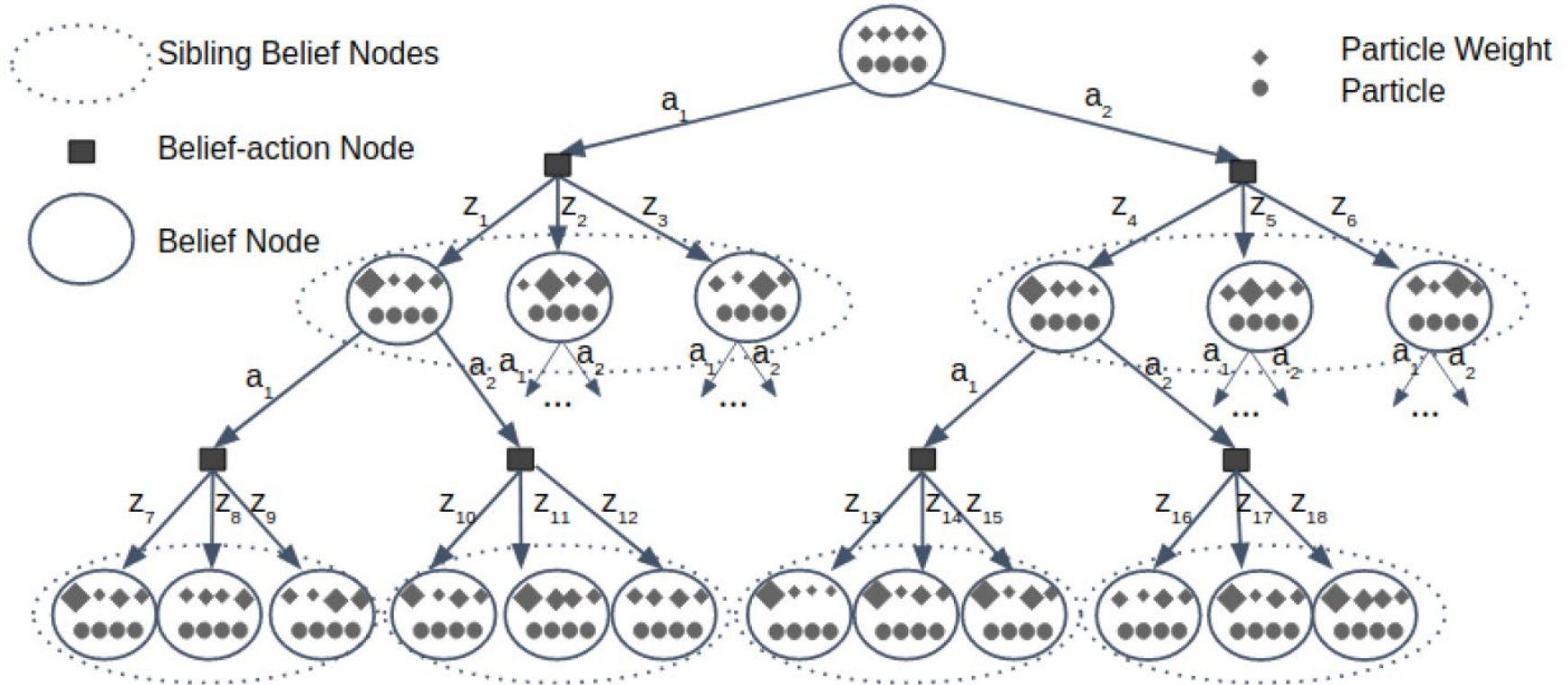
Belief Update



Belief Update



Belief Update



State particles are **shared** among all observation branches and **importance sampling weights** are given

Importance Sampling

Constructs an approximation using samples from a **proposal distribution**, and corrects for the discrepancy between the target and proposal using **(importance) weights**.

Importance Sampling

Constructs an approximation using samples from a **proposal distribution**, and corrects for the discrepancy between the target and proposal using **(importance) weights**.

$$E_p[f(x)] = \int f(x)p(x)dx = \int \frac{f(x)p(x)}{g(x)} g(x)dx = \frac{E_g \left[\frac{f(x)p(x)}{g(x)} \right]}{E_g \left[\frac{p(x)}{g(x)} \right]}$$

Importance Sampling

Constructs an approximation using samples from a **proposal distribution**, and corrects for the discrepancy between the target and proposal using **(importance) weights**.

$$\begin{aligned} E_p[f(x)] &= \int f(x)p(x)dx = \int \frac{f(x)p(x)}{g(x)} g(x)dx = \frac{E_g \left[\frac{f(x)p(x)}{g(x)} \right]}{E_g \left[\frac{p(x)}{g(x)} \right]} \\ &= \frac{E_g \left[\frac{f(x)p(x)}{g(x)} \right]}{E_g \left[\frac{p(x)}{g(x)} \right]} = \frac{\frac{1}{N} \sum_{i=1}^N f(x_i)W(x_i)}{\frac{1}{N} \sum_{j=1}^N W(x_j)} \approx \sum_{i=1}^N \hat{W}(x_i) f(x_i) \end{aligned}$$

Importance Sampling

Constructs an approximation using samples from a **proposal distribution**, and corrects for the discrepancy between the target and proposal using **(importance) weights**.

$$\begin{aligned} E_p[f(x)] &= \int f(x)p(x)dx = \int \frac{f(x)p(x)}{g(x)} g(x)dx = \frac{E_g \left[\frac{f(x)p(x)}{g(x)} \right]}{E_g \left[\frac{p(x)}{g(x)} \right]} \\ &= \frac{E_g \left[\frac{f(x)p(x)}{g(x)} \right]}{E_g \left[\frac{p(x)}{g(x)} \right]} = \frac{\frac{1}{N} \sum_{i=1}^N f(x_i)W(x_i)}{\frac{1}{N} \sum_{j=1}^N W(x_j)} \approx \sum_{i=1}^N \hat{W}(x_i) f(x_i) \end{aligned}$$

$$W(x_i) = \frac{p(x_i)}{q(x_i)}$$

$$\hat{W}(x_i) = \frac{W(x_i)}{\sum_{j=0}^N W(x_j)}$$

Importance Sampling

Constructs an approximation using samples from a **proposal distribution**, and corrects for the discrepancy between the target and proposal using **(importance) weights**.

Importance Sampling

Constructs an approximation using samples from a **proposal distribution**, and corrects for the discrepancy between the target and proposal using **(importance) weights**.

In POMCPOW:

Target Distribution $\Pr(s_{t+1} | b_t, a_t, o_{t+1})$

Proposal Distribution $\Pr(s_{t+1} | b_t, a_t)$

$$W(s_{t+1}) = \frac{p(s_{t+1})}{q(s_{t+1})} = \frac{\Pr(o_t | s', a_{t-1}, b_{t-1})}{\Pr(o_t | a_{t-1}, b_{t-1})}$$

Importance Sampling

Constructs an approximation using samples from a **proposal distribution**, and corrects for the discrepancy between the target and proposal using **(importance) weights**.

In POMCPOW:

$$\begin{aligned} \text{Target Distribution} \quad b_t(s') &= \Pr(s' | a_{t-1}, o_t, b_{t-1}) \\ &= \frac{\Pr(o_t | s', a_{t-1}, b_{t-1}) \Pr(s' | a_{t-1}, b_{t-1})}{\Pr(o_t | a_{t-1}, b_{t-1})} \end{aligned}$$

$$\text{Proposal Distribution} \quad \Pr(s_{t+1} | b_t, a_t)$$

$$W(s_{t+1}) = \frac{p(s_{t+1})}{q(s_{t+1})} = \frac{\Pr(o_t | s', a_{t-1}, b_{t-1})}{\Pr(o_t | a_{t-1}, b_{t-1})}$$

Importance Sampling

Constructs an approximation using samples from a **proposal distribution**, and corrects for the discrepancy between the target and proposal using **(importance) weights**.

In POMCPOW:

Target Distribution $\Pr(s_{t+1} | b_t, a_t, o_{t+1})$

Proposal Distribution $\Pr(s_{t+1} | b_t, a_t)$

$$W(s_{t+1}) = \frac{p(s_{t+1})}{q(s_{t+1})} = \frac{\Pr(o_t | s', a_{t-1}, b_{t-1})}{\Pr(o_t | a_{t-1}, b_{t-1})}$$

Importance Sampling

Constructs an approximation using samples from a **proposal distribution**, and corrects for the discrepancy between the target and proposal using **(importance) weights**.

$$\tilde{w}_{\mathcal{P}/\mathcal{Q}}(x) \equiv \frac{w_{\mathcal{P}/\mathcal{Q}}(x)}{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i)} \quad (\text{SN Importance Weight})$$

$$d_{\alpha}(\mathcal{P}||\mathcal{Q}) \equiv \mathbb{E}_{x \sim \mathcal{Q}}[w_{\mathcal{P}/\mathcal{Q}}(x)^{\alpha}] \quad (\text{Rényi Divergence})$$

$$\tilde{\mu}_{\mathcal{P}/\mathcal{Q}} \equiv \sum_{i=1}^N \tilde{w}_{\mathcal{P}/\mathcal{Q}}(x_i) f(x_i) \quad (\text{SN Estimator})$$

Importance Sampling

Constructs an approximation using samples from a **proposal distribution**, and corrects for the discrepancy between the target and proposal using **(importance) weights**.

Theorem 1

Theorem 1 (SN d_∞ -Concentration Bound). *Let \mathcal{P} and \mathcal{Q} be two probability measures on the measurable space $(\mathcal{X}, \mathcal{F})$ with $\mathcal{P} \ll \mathcal{Q}$ and $d_\infty(\mathcal{P}||\mathcal{Q}) < +\infty$. Let x_1, \dots, x_N be i.i.d.r.v. sampled from \mathcal{Q} , and $f : \mathcal{X} \rightarrow \mathbb{R}$ be a bounded Borel function ($\|f\|_\infty < +\infty$). Then, for any $\lambda > 0$ and N large enough such that $\lambda > \|f\|_\infty d_\infty(\mathcal{P}||\mathcal{Q})/\sqrt{N}$, the following bound holds with probability at least $1 - 3 \exp(-N \cdot t^2(\lambda, N))$:*

$$|\mathbb{E}_{x \sim \mathcal{P}}[f(x)] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}}| \leq \lambda$$

where $t(\lambda, N)$ is defined as:

$$t(\lambda, N) \equiv \frac{\lambda}{\|f\|_\infty d_\infty(\mathcal{P}||\mathcal{Q})} - \frac{1}{\sqrt{N}}$$

Proof 1

Lemma1:

$$\begin{aligned}
 \mathbb{P}(\hat{\mu}_{\mathcal{P}/\mathcal{Q}} - \mathbb{E}_{x \sim P}[f(x)] \geq \lambda) &= \mathbb{P}(\hat{\mu}_{\mathcal{P}/\mathcal{Q}} - \mathbb{E}_{x \sim Q}[W_{\mathcal{P}/\mathcal{Q}}(x)f(x)] \geq \lambda) \\
 &\leq \exp\left(-\frac{2N^2\lambda^2}{\sum_{i=1}^N 2(d_{\infty}(\mathcal{P} \parallel \mathcal{Q}) \parallel f \parallel_{\infty})^2}\right) \\
 &\leq \exp\left(-\frac{N\lambda^2}{d_{\infty}^2(\mathcal{P} \parallel \mathcal{Q}) \parallel f \parallel_{\infty}^2}\right) \equiv \delta \\
 \mathbb{P}(|\hat{\mu}_{\mathcal{P}/\mathcal{Q}} - \mathbb{E}_{x \sim P}[f(x)]| \geq \lambda) &\leq 2\exp\left(-\frac{N\lambda^2}{d_{\infty}^2(\mathcal{P} \parallel \mathcal{Q}) \parallel f \parallel_{\infty}^2}\right) = 2\delta
 \end{aligned}$$

Proof 1

Lemma1:

$$\begin{aligned}
 \mathbb{P}(\hat{\mu}_{\mathcal{P}/Q} - \mathbb{E}_{x \sim P}[f(x)] \geq \lambda) &= \mathbb{P}(\hat{\mu}_{\mathcal{P}/Q} - \mathbb{E}_{x \sim Q}[W_{\mathcal{P}/Q}(x)f(x)] \geq \lambda) \quad \textcircled{1} \\
 &\leq \exp\left(-\frac{2N^2\lambda^2}{\sum_{i=1}^N 2(d_\infty(\mathcal{P} \parallel Q) \parallel f \parallel_\infty)^2}\right) \\
 &\leq \exp\left(-\frac{N\lambda^2}{d_\infty^2(\mathcal{P} \parallel Q) \parallel f \parallel_\infty^2}\right) \equiv \delta \\
 \mathbb{P}(|\hat{\mu}_{\mathcal{P}/Q} - \mathbb{E}_{x \sim P}[f(x)]| \geq \lambda) &\leq 2\exp\left(-\frac{N\lambda^2}{d_\infty^2(\mathcal{P} \parallel Q) \parallel f \parallel_\infty^2}\right) = 2\delta
 \end{aligned}$$

$$\textcircled{1} \quad \mathbb{E}_{x \sim P}[f(x)] = \int f(x)p(x)dx = \int \frac{p(x)f(x)}{q(x)}q(x) = \mathbb{E}_{x \sim Q}[W_{\mathcal{P}/Q}(x)f(x)]$$

Proof 1

Lemma1:

$$\begin{aligned}
 \mathbb{P}(\hat{\mu}_{\mathcal{P}/\mathcal{Q}} - \mathbb{E}_{x \sim P}[f(x)] \geq \lambda) &= \mathbb{P}(\hat{\mu}_{\mathcal{P}/\mathcal{Q}} - \mathbb{E}_{x \sim Q}[W_{\mathcal{P}/\mathcal{Q}}(x)f(x)] \geq \lambda) \\
 &\leq \exp\left(-\frac{2N^2\lambda^2}{\sum_{i=1}^N 2(d_{\infty}(\mathcal{P} \parallel \mathcal{Q}) \parallel f \parallel_{\infty})^2}\right) \\
 &\leq \exp\left(-\frac{N\lambda^2}{d_{\infty}^2(\mathcal{P} \parallel \mathcal{Q}) \parallel f \parallel_{\infty}^2}\right) \equiv \delta \\
 \mathbb{P}(|\hat{\mu}_{\mathcal{P}/\mathcal{Q}} - \mathbb{E}_{x \sim P}[f(x)]| \geq \lambda) &\leq 2\exp\left(-\frac{N\lambda^2}{d_{\infty}^2(\mathcal{P} \parallel \mathcal{Q}) \parallel f \parallel_{\infty}^2}\right) = 2\delta
 \end{aligned}$$

Proof 1

Lemma1:

$$\begin{aligned}
 \mathbb{P}(\hat{\mu}_{\mathcal{P}/\mathcal{Q}} - \mathbb{E}_{x \sim P}[f(x)] \geq \lambda) &= \mathbb{P}(\hat{\mu}_{\mathcal{P}/\mathcal{Q}} - \mathbb{E}_{x \sim Q}[W_{\mathcal{P}/\mathcal{Q}}(x)f(x)] \geq \lambda) \\
 &\leq \exp\left(-\frac{2N^2\lambda^2}{\sum_{i=1}^N 2(d_{\infty}(\mathcal{P} \parallel \mathcal{Q}) \parallel f \parallel_{\infty})^2}\right) \quad \textcircled{2} \\
 &\leq \exp\left(-\frac{N\lambda^2}{d_{\infty}^2(\mathcal{P} \parallel \mathcal{Q}) \parallel f \parallel_{\infty}^2}\right) \equiv \delta \\
 \mathbb{P}(|\hat{\mu}_{\mathcal{P}/\mathcal{Q}} - \mathbb{E}_{x \sim P}[f(x)]| \geq \lambda) &\leq 2\exp\left(-\frac{N\lambda^2}{d_{\infty}^2(\mathcal{P} \parallel \mathcal{Q}) \parallel f \parallel_{\infty}^2}\right) = 2\delta
 \end{aligned}$$

Proof 1

$$\textcircled{2} \quad \hat{\mu}_{\mathcal{P}/\mathcal{Q}} = \frac{1}{N} \sum_{i=1}^N W_{\mathcal{P}/\mathcal{Q}}(x_i) f(x_i) = \frac{1}{N} \sum_{i=1}^N g(x_i) \quad g(x) = W_{\mathcal{P}/\mathcal{Q}}(x) f(x)$$

Proof 1

$$\textcircled{2} \quad \hat{\mu}_{\mathcal{P}/\mathcal{Q}} = \frac{1}{N} \sum_{i=1}^N W_{\mathcal{P}/\mathcal{Q}}(x_i) f(x_i) = \frac{1}{N} \sum_{i=1}^N g(x_i) \quad g(x) = W_{\mathcal{P}/\mathcal{Q}}(x) f(x)$$

$\hat{\mu}_{\mathcal{P}/\mathcal{Q}}$ is unbiased, which means $E[\hat{\mu}_{\mathcal{P}/\mathcal{Q}}] = E[g(x)]$, then:

Proof 1

$$\textcircled{2} \quad \hat{\mu}_{\mathcal{P}/\mathcal{Q}} = \frac{1}{N} \sum_{i=1}^N W_{\mathcal{P}/\mathcal{Q}}(x_i) f(x_i) = \frac{1}{N} \sum_{i=1}^N g(x_i) \quad g(x) = W_{\mathcal{P}/\mathcal{Q}}(x) f(x)$$

$\hat{\mu}_{\mathcal{P}/\mathcal{Q}}$ is unbiased, which means $E[\hat{\mu}_{\mathcal{P}/\mathcal{Q}}] = E[g(x)]$, then:

$$\mathbb{P}(\hat{\mu}_{\mathcal{P}/\mathcal{Q}} - \mathbb{E}_{x \sim \mathcal{Q}}[W_{\mathcal{P}/\mathcal{Q}}(x) f(x)] \geq \lambda) = \mathbb{P}(\hat{\mu}_{\mathcal{P}/\mathcal{Q}} - \mathbb{E}_{x \sim \mathcal{Q}}[\hat{\mu}_{\mathcal{P}/\mathcal{Q}}] \geq \lambda)$$

Proof 1

$$\textcircled{2} \quad \hat{\mu}_{\mathcal{P}/\mathcal{Q}} = \frac{1}{N} \sum_{i=1}^N W_{\mathcal{P}/\mathcal{Q}}(x_i) f(x_i) = \frac{1}{N} \sum_{i=1}^N g(x_i) \quad g(x) = W_{\mathcal{P}/\mathcal{Q}}(x) f(x)$$

$\hat{\mu}_{\mathcal{P}/\mathcal{Q}}$ is unbiased, which means $E[\hat{\mu}_{\mathcal{P}/\mathcal{Q}}] = E[g(x)]$, then:

$$\mathbb{P}(\hat{\mu}_{\mathcal{P}/\mathcal{Q}} - \mathbb{E}_{x \sim \mathcal{Q}}[W_{\mathcal{P}/\mathcal{Q}}(x) f(x)] \geq \lambda) = \mathbb{P}(\hat{\mu}_{\mathcal{P}/\mathcal{Q}} - \mathbb{E}_{x \sim \mathcal{Q}}[\hat{\mu}_{\mathcal{P}/\mathcal{Q}}] \geq \lambda)$$

Hoeffding's inequality:

Let $\{X_1, X_2, \dots, X_N\}$ be independent random variables bounded by the interval $[a, b]$, and $\bar{X} = \frac{1}{N} \sum_{i=1}^N X_i$, then:

Proof 1

$$\textcircled{2} \quad \hat{\mu}_{\mathcal{P}/Q} = \frac{1}{N} \sum_{i=1}^N W_{\mathcal{P}/Q}(x_i) f(x_i) = \frac{1}{N} \sum_{i=1}^N g(x_i) \quad g(x) = W_{\mathcal{P}/Q}(x) f(x)$$

$\hat{\mu}_{\mathcal{P}/Q}$ is unbiased, which means $E[\hat{\mu}_{\mathcal{P}/Q}] = E[g(x)]$, then:

$$\mathbb{P}(\hat{\mu}_{\mathcal{P}/Q} - \mathbb{E}_{x \sim Q}[W_{\mathcal{P}/Q}(x) f(x)] \geq \lambda) = \mathbb{P}(\hat{\mu}_{\mathcal{P}/Q} - \mathbb{E}_{x \sim Q}[\hat{\mu}_{\mathcal{P}/Q}] \geq \lambda)$$

Hoeffding's inequality:

Let $\{X_1, X_2, \dots, X_N\}$ be independent random variables bounded by the interval $[a, b]$, and $\bar{X} = \frac{1}{N} \sum_{i=1}^N X_i$, then:

$$\forall \lambda > 0, \quad \mathbb{P}(\bar{X} - \mathbb{E}_{x \sim Q}[\bar{X}] \geq \lambda) \leq \exp\left(-\frac{2N\lambda^2}{\sum_{i=1}^N (b_i - a_i)^2}\right)$$

Proof 1

$$\textcircled{2} \quad \hat{\mu}_{\mathcal{P}/\mathcal{Q}} = \frac{1}{N} \sum_{i=1}^N W_{\mathcal{P}/\mathcal{Q}}(x_i) f(x_i) = \frac{1}{N} \sum_{i=1}^N g(x_i) \quad g(x) = W_{\mathcal{P}/\mathcal{Q}}(x) f(x)$$

$\hat{\mu}_{\mathcal{P}/\mathcal{Q}}$ is unbiased, which means $E[\hat{\mu}_{\mathcal{P}/\mathcal{Q}}] = E[g(x)]$, then:

$$\mathbb{P}(\hat{\mu}_{\mathcal{P}/\mathcal{Q}} - \mathbb{E}_{x \sim \mathcal{Q}}[W_{\mathcal{P}/\mathcal{Q}}(x) f(x)] \geq \lambda) = \mathbb{P}(\hat{\mu}_{\mathcal{P}/\mathcal{Q}} - \mathbb{E}_{x \sim \mathcal{Q}}[\hat{\mu}_{\mathcal{P}/\mathcal{Q}}] \geq \lambda)$$

$$0 \leq g(x_i) \leq \|g(x)\|_{\infty} = d_{\infty}(\mathcal{P}||\mathcal{Q}) \|f(x)\|_{\infty}$$

$$\forall \lambda > 0, \quad \mathbb{P}(\bar{X} - \mathbb{E}_{x \sim \mathcal{Q}}[\bar{X}] \geq \lambda) \leq \exp\left(-\frac{2N\lambda^2}{\sum_{i=1}^N (b_i - a_i)^2}\right)$$

Proof 1

Further:

$$\begin{aligned} & \mathbb{P}(|\mathbb{E}_{x \sim \mathcal{P}}[f(x)] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}}| \geq \lambda) \\ & \leq \mathbb{P}(\tilde{\mu}_{\mathcal{P}/\mathcal{Q}} - \mathbb{E}_{x \sim \mathcal{P}}[f(x)] \geq \lambda) + \mathbb{P}(\mathbb{E}_{x \sim \mathcal{P}}[f(x)] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}} \geq \lambda) \\ & \leq \mathbb{P}(\tilde{\mu}_{\mathcal{P}/\mathcal{Q}} - \mathbb{E}_{x \sim \mathcal{Q}}[\tilde{\mu}_{\mathcal{P}/\mathcal{Q}}] \geq \tilde{\lambda}) + \mathbb{P}(\mathbb{E}_{x \sim \mathcal{Q}}[\tilde{\mu}_{\mathcal{P}/\mathcal{Q}}] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}} \geq \tilde{\lambda}) \\ & \leq \tilde{\delta} + \mathbb{P}(\mathbb{E}_{x \sim \mathcal{P}}[f(x)] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}} \geq \lambda) \end{aligned}$$

Proof 1

Further:

$$\begin{aligned} & \mathbb{P}(|\mathbb{E}_{x \sim \mathcal{P}}[f(x)] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}}| \geq \lambda) \\ & \leq \mathbb{P}(\tilde{\mu}_{\mathcal{P}/\mathcal{Q}} - \mathbb{E}_{x \sim \mathcal{P}}[f(x)] \geq \lambda) + \mathbb{P}(\mathbb{E}_{x \sim \mathcal{P}}[f(x)] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}} \geq \lambda) \quad \textcircled{1} \\ & \leq \mathbb{P}(\tilde{\mu}_{\mathcal{P}/\mathcal{Q}} - \mathbb{E}_{x \sim \mathcal{Q}}[\tilde{\mu}_{\mathcal{P}/\mathcal{Q}}] \geq \tilde{\lambda}) + \mathbb{P}(\mathbb{E}_{x \sim \mathcal{Q}}[\tilde{\mu}_{\mathcal{P}/\mathcal{Q}}] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}} \geq \tilde{\lambda}) \\ & \leq \tilde{\delta} + \mathbb{P}(\mathbb{E}_{x \sim \mathcal{P}}[f(x)] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}} \geq \lambda) \end{aligned}$$

$$\textcircled{1} \quad P(A \cup B) = P(A) + P(B) - P(A \cap B) \leq P(A) + P(B)$$

Proof 1

Further:

$$\begin{aligned}
 & \mathbb{P}(|\mathbb{E}_{x \sim \mathcal{P}}[f(x)] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}}| \geq \lambda) \\
 & \leq \mathbb{P}(\tilde{\mu}_{\mathcal{P}/\mathcal{Q}} - \mathbb{E}_{x \sim \mathcal{P}}[f(x)] \geq \lambda) + \mathbb{P}(\mathbb{E}_{x \sim \mathcal{P}}[f(x)] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}} \geq \lambda) \quad \textcircled{1} \\
 & \leq \mathbb{P}(\tilde{\mu}_{\mathcal{P}/\mathcal{Q}} - \mathbb{E}_{x \sim \mathcal{Q}}[\tilde{\mu}_{\mathcal{P}/\mathcal{Q}}] \geq \tilde{\lambda}) + \mathbb{P}(\mathbb{E}_{x \sim \mathcal{Q}}[\tilde{\mu}_{\mathcal{P}/\mathcal{Q}}] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}} \geq \tilde{\lambda}) \quad \textcircled{2} \\
 & \leq \tilde{\delta} + \mathbb{P}(\mathbb{E}_{x \sim \mathcal{P}}[f(x)] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}} \geq \lambda)
 \end{aligned}$$

$$\begin{aligned}
 \textcircled{2} \quad & \mathbb{P}(\tilde{\mu}_{\mathcal{P}/\mathcal{Q}} - \mathbb{E}_{x \sim \mathcal{P}}[f(x)] \geq \lambda) \\
 & = \mathbb{P}(\tilde{\mu}_{\mathcal{P}/\mathcal{Q}} - \mathbb{E}_{x \sim \mathcal{P}}[f(x)] - |\mathbb{E}_{x \sim \mathcal{P}}[f(x)] - \mathbb{E}_{x \sim \mathcal{Q}}[\tilde{\mu}_{\mathcal{P}/\mathcal{Q}}]| \geq \tilde{\lambda}) \\
 & \leq \mathbb{P}(\tilde{\mu}_{\mathcal{P}/\mathcal{Q}} - \mathbb{E}_{x \sim \mathcal{P}}[f(x)] - (\mathbb{E}_{x \sim \mathcal{P}}[f(x)] - \mathbb{E}_{x \sim \mathcal{Q}}[\tilde{\mu}_{\mathcal{P}/\mathcal{Q}}])) \geq \tilde{\lambda})
 \end{aligned}$$

Proof 1

Further:

$$\begin{aligned}
 & \mathbb{P}(|\mathbb{E}_{x \sim \mathcal{P}}[f(x)] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}}| \geq \lambda) \\
 & \leq \mathbb{P}(\tilde{\mu}_{\mathcal{P}/\mathcal{Q}} - \mathbb{E}_{x \sim \mathcal{P}}[f(x)] \geq \lambda) + \mathbb{P}(\mathbb{E}_{x \sim \mathcal{P}}[f(x)] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}} \geq \lambda) \quad \textcircled{1} \\
 & \leq \mathbb{P}(\tilde{\mu}_{\mathcal{P}/\mathcal{Q}} - \mathbb{E}_{x \sim \mathcal{Q}}[\tilde{\mu}_{\mathcal{P}/\mathcal{Q}}] \geq \tilde{\lambda}) + \mathbb{P}(\mathbb{E}_{x \sim \mathcal{Q}}[\tilde{\mu}_{\mathcal{P}/\mathcal{Q}}] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}} \geq \tilde{\lambda}) \quad \textcircled{2} \\
 & \leq \tilde{\delta} + \mathbb{P}(\mathbb{E}_{x \sim \mathcal{P}}[f(x)] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}} \geq \lambda) \quad \textcircled{3}
 \end{aligned}$$

③ Using Lemma 1

$$\begin{aligned}
 \tilde{\lambda} &= \lambda - |\mathbb{E}_{x \sim \mathcal{P}}[f(x)] - \mathbb{E}_{x \sim \mathcal{Q}}[\tilde{\mu}_{\mathcal{P}/\mathcal{Q}}]| \\
 \tilde{\delta} &= \exp\left(-\frac{N\tilde{\lambda}^2}{d_{\infty}^2(\mathcal{P} \parallel \mathcal{Q}) \parallel f \parallel_{\infty}^2}\right)
 \end{aligned}$$

Proof 1

Further:

$$\begin{aligned} & \mathbb{P}(|\mathbb{E}_{x \sim \mathcal{P}}[f(x)] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}}| \geq \lambda) \\ & \leq \mathbb{P}(\tilde{\mu}_{\mathcal{P}/\mathcal{Q}} - \mathbb{E}_{x \sim \mathcal{P}}[f(x)] \geq \lambda) + \mathbb{P}(\mathbb{E}_{x \sim \mathcal{P}}[f(x)] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}} \geq \lambda) \quad \textcircled{1} \\ & \leq \mathbb{P}(\tilde{\mu}_{\mathcal{P}/\mathcal{Q}} - \mathbb{E}_{x \sim \mathcal{Q}}[\tilde{\mu}_{\mathcal{P}/\mathcal{Q}}] \geq \tilde{\lambda}) + \mathbb{P}(\mathbb{E}_{x \sim \mathcal{Q}}[\tilde{\mu}_{\mathcal{P}/\mathcal{Q}}] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}} \geq \tilde{\lambda}) \quad \textcircled{2} \\ & \leq \tilde{\delta} + \mathbb{P}(\mathbb{E}_{x \sim \mathcal{P}}[f(x)] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}} \geq \lambda) \quad \textcircled{3} \end{aligned}$$

Proof 1

Further:

$$\begin{aligned}
 & \mathbb{P}(|\mathbb{E}_{x \sim \mathcal{P}}[f(x)] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}}| \geq \lambda) \\
 & \leq \mathbb{P}(\tilde{\mu}_{\mathcal{P}/\mathcal{Q}} - \mathbb{E}_{x \sim \mathcal{P}}[f(x)] \geq \lambda) + \mathbb{P}(\mathbb{E}_{x \sim \mathcal{P}}[f(x)] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}} \geq \lambda) \quad \textcircled{1} \\
 & \leq \mathbb{P}(\tilde{\mu}_{\mathcal{P}/\mathcal{Q}} - \mathbb{E}_{x \sim \mathcal{Q}}[\tilde{\mu}_{\mathcal{P}/\mathcal{Q}}] \geq \tilde{\lambda}) + \mathbb{P}(\mathbb{E}_{x \sim \mathcal{Q}}[\tilde{\mu}_{\mathcal{P}/\mathcal{Q}}] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}} \geq \tilde{\lambda}) \quad \textcircled{2} \\
 & \leq \tilde{\delta} + \mathbb{P}(\mathbb{E}_{x \sim \mathcal{P}}[f(x)] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}} \geq \lambda) \quad \textcircled{3}
 \end{aligned}$$

$$\begin{aligned}
 \mathbb{P}(\mathbb{E}_{x \sim \mathcal{P}}[f(x)] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}} \geq \lambda) & \leq \mathbb{P}(\mathbb{E}_{x \sim \mathcal{Q}}[\tilde{\mu}_{\mathcal{P}/\mathcal{Q}}] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}} \geq \tilde{\lambda}) \\
 & \leq \mathbb{P}(|\mathbb{E}_{x \sim \mathcal{Q}}[\tilde{\mu}_{\mathcal{P}/\mathcal{Q}}] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}}| \geq \tilde{\lambda}) \leq 2\tilde{\delta}
 \end{aligned}$$

Proof 1

Further:

$$\begin{aligned}
 & \mathbb{P}(|\mathbb{E}_{x \sim \mathcal{P}}[f(x)] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}}| \geq \lambda) \\
 & \leq \mathbb{P}(\tilde{\mu}_{\mathcal{P}/\mathcal{Q}} - \mathbb{E}_{x \sim \mathcal{P}}[f(x)] \geq \lambda) + \mathbb{P}(\mathbb{E}_{x \sim \mathcal{P}}[f(x)] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}} \geq \lambda) \quad \textcircled{1} \\
 & \leq \mathbb{P}(\tilde{\mu}_{\mathcal{P}/\mathcal{Q}} - \mathbb{E}_{x \sim \mathcal{Q}}[\tilde{\mu}_{\mathcal{P}/\mathcal{Q}}] \geq \tilde{\lambda}) + \mathbb{P}(\mathbb{E}_{x \sim \mathcal{Q}}[\tilde{\mu}_{\mathcal{P}/\mathcal{Q}}] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}} \geq \tilde{\lambda}) \quad \textcircled{2} \\
 & \leq \tilde{\delta} + \mathbb{P}(\mathbb{E}_{x \sim \mathcal{P}}[f(x)] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}} \geq \lambda) \quad \textcircled{3}
 \end{aligned}$$

$$\begin{aligned}
 \mathbb{P}(\mathbb{E}_{x \sim \mathcal{P}}[f(x)] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}} \geq \lambda) & \leq \mathbb{P}(\mathbb{E}_{x \sim \mathcal{Q}}[\tilde{\mu}_{\mathcal{P}/\mathcal{Q}}] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}} \geq \tilde{\lambda}) \\
 & \leq \mathbb{P}(|\mathbb{E}_{x \sim \mathcal{Q}}[\tilde{\mu}_{\mathcal{P}/\mathcal{Q}}] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}}| \geq \tilde{\lambda}) \leq 2\tilde{\delta}
 \end{aligned}$$

$$\mathbb{P}(|\mathbb{E}_{x \sim \mathcal{P}}[f(x)] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}}| \geq \lambda) \leq 3\tilde{\delta}$$

Proof 1

$$\begin{aligned}
& |\mathbb{E}_{x \sim \mathcal{P}}[f(x)] - \mathbb{E}_{x \sim \mathcal{Q}}[\tilde{\mu}_{\mathcal{P}/\mathcal{Q}}]| = |\mathbb{E}_{x \sim \mathcal{Q}}[\hat{\mu}_{\mathcal{P}/\mathcal{Q}} - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}}]| \leq \mathbb{E}_{x \sim \mathcal{Q}}[|\hat{\mu}_{\mathcal{P}/\mathcal{Q}} - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}}|] \\
& \leq \mathbb{E}_{x \sim \mathcal{Q}} \left| \frac{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i) f(x_i)}{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i)} - \frac{1}{N} \sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i) f(x_i) \right| \\
& = \mathbb{E}_{x \sim \mathcal{Q}} \left[\left| \frac{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i) f(x_i)}{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i)} \right| \left| 1 - \frac{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i)}{N} \right| \right] \\
& \leq \mathbb{E}_{x \sim \mathcal{Q}} \left[\left(\frac{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i) f(x_i)}{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i)} \right)^2 \right]^{1/2} \mathbb{E}_{x \sim \mathcal{Q}} \left[\left(1 - \frac{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i)}{N} \right)^2 \right]^{1/2} \\
& \leq \|f\|_{\infty} \sqrt{\frac{d_2(\mathcal{P} \parallel \mathcal{Q}) - 1}{N}} \leq \|f\|_{\infty} \frac{d_{\infty}(\mathcal{P} \parallel \mathcal{Q})}{\sqrt{N}}
\end{aligned}$$

Proof 1

$$\begin{aligned}
 |\mathbb{E}_{x \sim \mathcal{P}}[f(x)] - \mathbb{E}_{x \sim \mathcal{Q}}[\tilde{\mu}_{\mathcal{P}/\mathcal{Q}}]| &= |\mathbb{E}_{x \sim \mathcal{Q}}[\hat{\mu}_{\mathcal{P}/\mathcal{Q}} - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}}]| \leq \mathbb{E}_{x \sim \mathcal{Q}}[|\hat{\mu}_{\mathcal{P}/\mathcal{Q}} - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}}|] \\
 &\leq \mathbb{E}_{x \sim \mathcal{Q}} \left| \frac{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i) f(x_i)}{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i)} - \frac{1}{N} \sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i) f(x_i) \right| \\
 &= \mathbb{E}_{x \sim \mathcal{Q}} \left[\left| \frac{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i) f(x_i)}{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i)} \right| \left| 1 - \frac{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i)}{N} \right| \right] \quad \textcircled{1} \\
 &\leq \mathbb{E}_{x \sim \mathcal{Q}} \left[\left(\frac{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i) f(x_i)}{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i)} \right)^2 \right]^{1/2} \mathbb{E}_{x \sim \mathcal{Q}} \left[\left(1 - \frac{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i)}{N} \right)^2 \right]^{1/2} \\
 &\leq \|f\|_{\infty} \sqrt{\frac{d_2(\mathcal{P}||\mathcal{Q}) - 1}{N}} \leq \|f\|_{\infty} \frac{d_{\infty}(\mathcal{P}||\mathcal{Q})}{\sqrt{N}}
 \end{aligned}$$

① Cauchy-Schwarz inequality

$$\left(\sum_{i=1}^n x_i y_i \right)^2 \leq \left(\sum_{i=1}^n x_i^2 \right) \left(\sum_{i=1}^n y_i^2 \right)$$

Proof 1

$$|\mathbb{E}_{x \sim \mathcal{P}}[f(x)] - \mathbb{E}_{x \sim \mathcal{Q}}[\tilde{\mu}_{\mathcal{P}/\mathcal{Q}}]| = |\mathbb{E}_{x \sim \mathcal{Q}}[\hat{\mu}_{\mathcal{P}/\mathcal{Q}} - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}}]| \leq \mathbb{E}_{x \sim \mathcal{Q}}[|\hat{\mu}_{\mathcal{P}/\mathcal{Q}} - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}}|]$$

$$\leq \mathbb{E}_{x \sim \mathcal{Q}} \left| \frac{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i) f(x_i)}{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i)} - \frac{1}{N} \sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i) f(x_i) \right|$$

$$= \mathbb{E}_{x \sim \mathcal{Q}} \left[\left| \frac{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i) f(x_i)}{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i)} \right| \left| 1 - \frac{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i)}{N} \right| \right] \quad \textcircled{1}$$

$$\leq \mathbb{E}_{x \sim \mathcal{Q}} \left[\left(\frac{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i) f(x_i)}{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i)} \right)^2 \right]^{1/2} \mathbb{E}_{x \sim \mathcal{Q}} \left[\left(1 - \frac{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i)}{N} \right)^2 \right]^{1/2} \quad \textcircled{2}$$

$$\leq \|f\|_{\infty} \sqrt{\frac{d_2(\mathcal{P} \parallel \mathcal{Q}) - 1}{N}} \leq \|f\|_{\infty} \frac{d_{\infty}(\mathcal{P} \parallel \mathcal{Q})}{\sqrt{N}}$$

Proof 1

$$\begin{aligned}
 & |\mathbb{E}_{x \sim \mathcal{P}}[f(x)] - \mathbb{E}_{x \sim \mathcal{Q}}[\tilde{\mu}_{\mathcal{P}/\mathcal{Q}}]| = |\mathbb{E}_{x \sim \mathcal{Q}}[\hat{\mu}_{\mathcal{P}/\mathcal{Q}} - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}}]| \leq \mathbb{E}_{x \sim \mathcal{Q}}[|\hat{\mu}_{\mathcal{P}/\mathcal{Q}} - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}}|] \\
 & \leq \mathbb{E}_{x \sim \mathcal{Q}} \left| \frac{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i) f(x_i)}{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i)} - \frac{1}{N} \sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i) f(x_i) \right| \\
 & = \mathbb{E}_{x \sim \mathcal{Q}} \left[\left| \frac{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i) f(x_i)}{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i)} \right| \left| 1 - \frac{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i)}{N} \right| \right] \tag{1} \\
 & \leq \mathbb{E}_{x \sim \mathcal{Q}} \left[\left(\frac{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i) f(x_i)}{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i)} \right)^2 \right]^{1/2} \mathbb{E}_{x \sim \mathcal{Q}} \left[\left(1 - \frac{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i)}{N} \right)^2 \right]^{1/2} \tag{2} \\
 & \leq \|f\|_{\infty} \sqrt{\frac{d_2(\mathcal{P} \parallel \mathcal{Q}) - 1}{N}} \leq \|f\|_{\infty} \frac{d_{\infty}(\mathcal{P} \parallel \mathcal{Q})}{\sqrt{N}} \tag{3}
 \end{aligned}$$

Proof 1

$$\begin{aligned}
 & |\mathbb{E}_{x \sim \mathcal{P}}[f(x)] - \mathbb{E}_{x \sim \mathcal{Q}}[\tilde{\mu}_{\mathcal{P}/\mathcal{Q}}]| = |\mathbb{E}_{x \sim \mathcal{Q}}[\hat{\mu}_{\mathcal{P}/\mathcal{Q}} - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}}]| \leq \mathbb{E}_{x \sim \mathcal{Q}}[|\hat{\mu}_{\mathcal{P}/\mathcal{Q}} - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}}|] \\
 & \leq \mathbb{E}_{x \sim \mathcal{Q}} \left| \frac{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i) f(x_i)}{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i)} - \frac{1}{N} \sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i) f(x_i) \right| \\
 & = \mathbb{E}_{x \sim \mathcal{Q}} \left[\left| \frac{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i) f(x_i)}{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i)} \right| \left| 1 - \frac{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i)}{N} \right| \right] \quad \textcircled{1} \\
 & \leq \mathbb{E}_{x \sim \mathcal{Q}} \left[\left(\frac{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i) f(x_i)}{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i)} \right)^2 \right]^{1/2} \mathbb{E}_{x \sim \mathcal{Q}} \left[\left(1 - \frac{\sum_{i=1}^N w_{\mathcal{P}/\mathcal{Q}}(x_i)}{N} \right)^2 \right]^{1/2} \quad \textcircled{2} \\
 & \leq \|f\|_{\infty} \sqrt{\frac{d_2(\mathcal{P}||\mathcal{Q}) - 1}{N}} \leq \|f\|_{\infty} \frac{d_{\infty}(\mathcal{P}||\mathcal{Q})}{\sqrt{N}} \quad \textcircled{3}
 \end{aligned}$$

$$\begin{aligned}
 \textcircled{3} \quad & d_2(\mathcal{P}||\mathcal{Q}) - 1 \leq d_2(\mathcal{P}||\mathcal{Q}) = \mathbb{E}_{x \sim \mathcal{Q}}[W_{\mathcal{P}/\mathcal{Q}}(x)^2] \\
 & \leq W_{\mathcal{P}/\mathcal{Q}}(x)_{\max}^2 = d_{\infty}(\mathcal{P}||\mathcal{Q})^2
 \end{aligned}$$

Proof 1

$$\begin{aligned} \textcircled{2} \quad & \mathbb{E}_{x \sim Q} \left[\left(1 - \frac{\sum_{i=1}^N W_{\mathcal{P}/Q}(x_i)}{N} \right)^2 \right] = \mathbb{E}_{x \sim Q} \left[\sum_{i=1}^N \left(\frac{1 - W_{\mathcal{P}/Q}(x_i)}{N} \right)^2 \right] \\ & \leq \sum_{i=1}^N \mathbb{E}_{x \sim Q} \left[\left(\frac{1 - W_{\mathcal{P}/Q}(x_i)}{N} \right)^2 \right] \leq \sum_{i=1}^N \mathbb{E}_{x \sim Q} \left[\left(\frac{W_{\mathcal{P}/Q}(x_i)}{N} \right)^2 - \left(\frac{1}{N} \right)^2 \right] \\ & = \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{x \sim Q} \left[\left(W_{\mathcal{P}/Q}(x_i) \right)^2 \right] - \frac{1}{N} \leq \frac{d_2(\mathcal{P}||Q) - 1}{N} \end{aligned}$$

Proof 1

Last:

$$\begin{aligned}\tilde{\delta} &\leq \exp \left(-\frac{N(\lambda - \|f\|_{\infty} d_{\infty}(\mathcal{P}||\mathcal{Q})/\sqrt{N})^2}{d_{\infty}^2(\mathcal{P}||\mathcal{Q})\|f\|_{\infty}^2} \right) \\ &= \exp \left(-N \left(\frac{\lambda - \|f\|_{\infty} d_{\infty}(\mathcal{P}||\mathcal{Q})/\sqrt{N}}{\|f\|_{\infty} d_{\infty}(\mathcal{P}||\mathcal{Q})} \right)^2 \right) \\ &\equiv \exp \left(-N \cdot t^2(\lambda, N) \right)\end{aligned}$$

Lemma 1

Lemma 1 (SN Estimator Leaf Node Convergence).
 $\hat{Q}_{D-1}^(\bar{b}_{D-1}, a)$ is an SN estimator of $Q_{D-1}^*(b_{D-1}, a)$, and the following leaf-node concentration bound holds with probability at least $1 - 3 \exp(-C \cdot t_{\max}^2(\lambda, C))$,*

$$|Q_{D-1}^*(b_{D-1}, a) - \hat{Q}_{D-1}^*(\bar{b}_{D-1}, a)| \leq \lambda$$

Proof 2

$\hat{Q}_{D-1}^*(\bar{b}_{D-1}, a)$ is an SN estimator of $Q_{D-1}^*(\bar{b}_{D-1}, a)$:

$$\hat{Q}_{D-1}^*(\bar{b}_{D-1}, a) = \sum_{i=1}^c \tilde{w}_{\mathcal{P}^{D-1}/\mathcal{Q}^{D-1}}(\{s_n\}_i) R(s_{D-1,i}, a)$$

$$Q_{D-1}^*(b_{D-1}, a) = \int_S R(S_{D-1}, a) dS$$

Proof 2

$\hat{Q}_{D-1}^*(\bar{b}_{D-1}, a)$ is **an SN estimator** of $Q_{D-1}^*(\bar{b}_{D-1}, a)$:

Using Theorem 1:

first bound R by $3V_{max}$:

$$V_{max} \equiv \frac{R_{max}}{1-\gamma} \geq R_{max}$$

Proof 2

Theorem 1 (SN d_∞ -Concentration Bound). *Let \mathcal{P} and \mathcal{Q} be two probability measures on the measurable space $(\mathcal{X}, \mathcal{F})$ with $\mathcal{P} \ll \mathcal{Q}$ and $d_\infty(\mathcal{P}||\mathcal{Q}) < +\infty$. Let x_1, \dots, x_N be i.i.d.r.v. sampled from \mathcal{Q} , and $f : \mathcal{X} \rightarrow \mathbb{R}$ be a bounded Borel function ($\|f\|_\infty < +\infty$). Then, for any $\lambda > 0$ and N large enough such that $\lambda > \|f\|_\infty d_\infty(\mathcal{P}||\mathcal{Q})/\sqrt{N}$, the following bound holds with probability at least $1 - 3 \exp(-N \cdot t^2(\lambda, N))$:*

$$|\mathbb{E}_{x \sim \mathcal{P}}[f(x)] - \tilde{\mu}_{\mathcal{P}/\mathcal{Q}}| \leq \lambda \quad (6)$$

where $t(\lambda, N)$ is defined as:

$$t(\lambda, N) \equiv \frac{\lambda}{\|f\|_\infty d_\infty(\mathcal{P}||\mathcal{Q})} - \frac{1}{\sqrt{N}} \quad (7)$$

Proof 2

$\hat{Q}_{D-1}^*(\bar{b}_{D-1}, a)$ is **an SN estimator** of $Q_{D-1}^*(\bar{b}_{D-1}, a)$:

Using Theorem 1:

first bound R by $3V_{max}$:

$$V_{max} \equiv \frac{R_{max}}{1-\gamma} \geq R_{max}$$

Proof 2

$\hat{Q}_{D-1}^*(\bar{b}_{D-1}, a)$ is **an SN estimator** of $Q_{D-1}^*(\bar{b}_{D-1}, a)$:

Using Theorem 1:

first bound R by $3V_{max}$:

$$V_{max} \equiv \frac{R_{max}}{1-\gamma} \geq R_{max}$$

$$t_{D-1}(\lambda, C) = \frac{\lambda}{3V_{max}d_{\infty}(\mathcal{P}^{D-1} \parallel \mathcal{Q}^{D-1})} - \frac{1}{\sqrt{C}} \geq \frac{\lambda}{3V_{max}d_{\infty}^{max}} - \frac{1}{\sqrt{C}} \equiv t_{max}(\lambda, C)$$

$$\forall i < D, \quad d_{\infty}(\mathcal{P}^i \parallel \mathcal{Q}^i) \leq d_{\infty}^{max}$$

$$t_i(\lambda, C) \geq t_{max}(\lambda, C)$$

Lemma 2

Lemma 2 (SN Estimator Step-by-Step Convergence). $\hat{Q}_d^*(\bar{b}_d, a)$ is an SN estimator of $Q_d^*(b_d, a)$, and for all $d = 0, \dots, D - 1$ and a , the following holds with probability at least $1 - 3|A|(3|A|C)^D \exp(-C \cdot t_{\max}^2)$:

$$|Q_d^*(b_d, a) - \hat{Q}_d^*(\bar{b}_d, a)| \leq \alpha_d$$
$$\alpha_d \equiv \lambda + \gamma \alpha_{d+1}; \alpha_{D-1} = \lambda$$

Proof 3

$$|Q_d^*(b_d, a) - \hat{Q}_d^*(\bar{b}_d, a)| \leq \underbrace{\left| \mathbb{E}[R(s_d, a)|b_d] - \frac{\sum_{i=1}^C w_{d,i} r_{d,i}}{\sum_{i=1}^C w_{d,i}} \right|}_{(A)} + \gamma \underbrace{\left| \mathbb{E}[V_{d+1}^*(ba_o)|b_d] - \frac{\sum_{i=1}^C w_{d,i} \hat{V}_{d+1}^*(\overline{b_d a_o_i})}{\sum_{i=1}^C w_{d,i}} \right|}_{(B)}$$

Proof 3

$$|Q_d^*(b_d, a) - \hat{Q}_d^*(\bar{b}_d, a)| \leq \underbrace{\left| \mathbb{E}[R(s_d, a)|b_d] - \frac{\sum_{i=1}^C w_{d,i} r_{d,i}}{\sum_{i=1}^C w_{d,i}} \right|}_{(A)} + \gamma \underbrace{\left| \mathbb{E}[V_{d+1}^*(bao)|b_d] - \frac{\sum_{i=1}^C w_{d,i} \hat{V}_{d+1}^*(\overline{b_d a o_i})}{\sum_{i=1}^C w_{d,i}} \right|}_{(B)}$$

$$|A_1 + B_1 - A_2 - B_2| \leq |A_1 - A_2| + |B_1 - B_2|$$

$$\hat{Q}_d^*(\bar{b}_d, a) = \frac{\sum_{i=1}^C w_{d,i} \left(r_{d,i} + \gamma \bar{V}_{d+1}^*(\overline{b_d a o_i}) \right)}{\sum_{i=1}^C w_{d,i}}$$

$$Q_d^*(b_d, a) = E[R(s_d, a) + \gamma V_{d+1}^*(bao)|b_d]$$

Proof 3

$$|Q_d^*(b_d, a) - \hat{Q}_d^*(\bar{b}_d, a)| \leq \underbrace{\left| \mathbb{E}[R(s_d, a)|b_d] - \frac{\sum_{i=1}^C w_{d,i} r_{d,i}}{\sum_{i=1}^C w_{d,i}} \right|}_{(A)} + \gamma \underbrace{\left| \mathbb{E}[V_{d+1}^*(ba_o)|b_d] - \frac{\sum_{i=1}^C w_{d,i} \hat{V}_{d+1}^*(\bar{b}_d a_{o_i})}{\sum_{i=1}^C w_{d,i}} \right|}_{(B)}$$

For part A: bound R by R_{max} and let $\lambda = \frac{R_{max}}{3V_{max}} \lambda$, using Thero1:

$$A \leq \frac{R_{max}}{3V_{max}} \lambda \quad \Pr \geq 1 - \exp(-t_{max}^2)$$

Proof 3

For part B:

$$\begin{aligned}
 (B) &\leq \underbrace{\left| \mathbb{E}[V_{d+1}^*(bao)|b_d] - \frac{\sum_{i=1}^C w_{d,i} V_{d+1}^*(s_{d,i}, b_d, a)}{\sum_{i=1}^C w_{d,i}} \right|}_{\text{Importance sampling error}} \\
 &+ \underbrace{\left| \frac{\sum_{i=1}^C w_{d,i} V_{d+1}^*(s_{d,i}, b_d, a)}{\sum_{i=1}^C w_{d,i}} - \frac{\sum_{i=1}^C w_{d,i} V_{d+1}^*(b_d a o_i)}{\sum_{i=1}^C w_{d,i}} \right|}_{\text{MC next-step integral approximation error}} \\
 &+ \underbrace{\left| \frac{\sum_{i=1}^C w_{d,i} V_{d+1}^*(b_d a o_i)}{\sum_{i=1}^C w_{d,i}} - \frac{\sum_{i=1}^C w_{d,i} \hat{V}_{d+1}^*(\overline{b_d a o_i})}{\sum_{i=1}^C w_{d,i}} \right|}_{\text{Function estimation error}} \\
 &\leq \frac{1}{3}\lambda + \frac{2}{3\gamma}\lambda + \alpha_{d+1}
 \end{aligned}$$

Proof 3

For Important Sampling error ,

$$\begin{aligned}
 \mathbf{V}_{d+1}^*(s_{d,i}, b_d, a) &\equiv \int_S \int_O V_{d+1}^*(b_d a o) \mathcal{Z}(o|a, s_{d+1}) \mathcal{T}(s_{d+1}|s_{d,i}, a) ds_{d+1} do \\
 \mathbb{E}[V_{d+1}^*(ba o)|b_d] &= \int_S \int_S \int_O V_{d+1}^*(b_d a o) (\mathcal{Z}_{d+1})(\mathcal{T}_{d,d+1}) b_d \cdot ds_{d:d+1} do \\
 &= \int_S \mathbf{V}_{d+1}^*(s_d, b_d, a) b_d \cdot ds_d \\
 &= \frac{\int_{S^{d+1}} \mathbf{V}_{d+1}^*(s_d, b_d, a) (\mathcal{Z}_{1:d})(\mathcal{T}_{1:d}) b_0 ds_{0:d}}{\int_{S^{d+1}} (\mathcal{Z}_{1:d})(\mathcal{T}_{1:d}) b_0 ds_{0:d}}
 \end{aligned}$$

Using SN inequality: $\|V_{d+1}^*(s_d, b_d, a)\|_\infty = V_{max}$

$$\text{IS} \leq \frac{1}{3} \lambda \qquad \text{Pr} \geq 1 - \exp(-t_{\max}^2)$$

Proof 3

For MC error:

the quantity $V_{d+1}^*(b_d a o_i)$ for a given $(s_{d,i}, b_d, a)$ is an **unbiased** 1-sample MC estimate of $\mathbf{V}_{d+1}^*(s_{d,i}, b_d, a)$

$$\Delta_{d+1}(s_{d,i}, b_d, a) \equiv \mathbf{V}_{d+1}^*(s_{d,i}, b_d, a) - V_{d+1}^*(b_d a o_i) \quad E(\Delta_{d+1}) = 0$$

$$\begin{aligned} & \left| \frac{\sum_{i=1}^C w_{d,i} \mathbf{V}_{d+1}^*(s_{d,i}, b_d, a)}{\sum_{i=1}^C w_{d,i}} - \frac{\sum_{i=1}^C w_{d,i} V_{d+1}^*(b_d a o_i)}{\sum_{i=1}^C w_{d,i}} \right| \\ &= \left| \frac{\sum_{i=1}^C w_{d,i} \Delta_{d+1}(s_{d,i}, b_d, a)}{\sum_{i=1}^C w_{d,i}} - 0 \right| \leq \frac{2}{3} \lambda \leq \frac{2}{3\gamma} \lambda \quad \|\Delta_{d+1}\|_{\infty} \leq 2V_{max} \end{aligned}$$

Proof 3

For Function Estimation Error:

the third term is bounded by the **inductive hypothesis**, since each i -th absolute difference of the Q -function and its estimate at step $d + 1$, and furthermore the value function and its estimate at step $d + 1$, are all bounded by α_{d+1} .

Proof 3

For part B:

$$\begin{aligned} |Q_d^*(b_d, a) - \hat{Q}_d^*(\bar{b}_d, a)| &\leq \frac{R_{\max}}{3V_{\max}}\lambda + \gamma\left[\frac{1}{3}\lambda + \frac{2}{3\gamma}\lambda + \alpha_{d+1}\right] \\ &\leq \lambda + \gamma\alpha_{d+1} = \alpha_d \end{aligned}$$

Thanks!

