

# Adaptive Multistage Sampling Algorithm: The Origins of Monte Carlo Tree Search

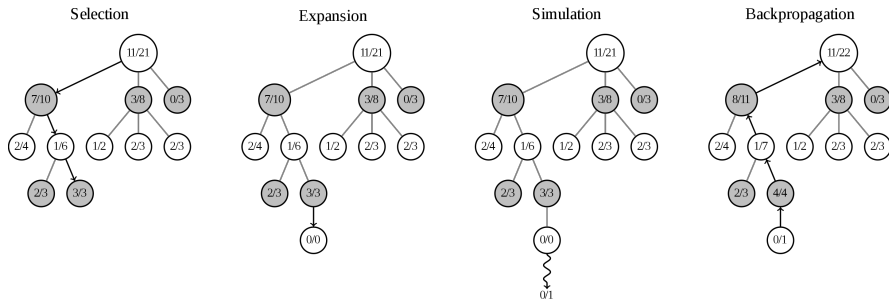
Ashwin Rao

ICME, Stanford University

December 9, 2019

# Monte Carlo Tree Search (MCTS)

- MCTS was popularized a few years ago by [Deep Mind's AlphaGo](#)
- It is a simulation-based method to identify the best action in a state
- MCTS term was first introduced by [Remi Coulom](#) for game trees
- Each round of MCTS consists of four steps:
  - Selection: Successively select children from root R to leaf L
  - Expansion: Create node C as a new child of L
  - Simulation: Complete a random playout from C
  - Backpropagation: Use result of playout to update nodes from C to R



# The Selection Step in MCTS

- Selection involves picking a child with “most promise”
- This means prioritizing children with higher success estimates
- For estimate confidence, we need sufficient playouts under *each* child
- This is our usual *Explore v/s Exploit* dilemma (Multi-armed Bandit)
- Explore v/s Exploit formula for games first due to [Kocsis-Szepesvari](#)
- Formula called *Upper Confidence Bound 1 for Trees* (abbrev. UCT)
- Most current MCTS Algorithms are based on some variant of UCT
- UCT is based on UCB1 formula of [Auer, Cesa-Bianchi, Fischer](#)
- However, MCTS and UCT concepts first appeared in the [Adaptive Multistage Sampling algorithm of Chang, Fu, Hu, Marcus](#)
- Adaptive Multistage Sampling (AMS) is a generic simulation-based algorithm to solve a finite-horizon Markov Decision Process (MDP)
- AMS can be considered as the “spiritual origin” of MCTS/UCT
- Hence, this lecture is dedicated to AMS

# The Setting for the AMS Algorithm

- MDP with finite number of time steps  $t = 0, 1, \dots, T$
- State denoted  $s_t \in \mathcal{S}$ , where  $\mathcal{S}$  is very large
- Action denoted  $a_t \in \mathcal{A}$ , where  $\mathcal{A}$  is fairly small
- Reward  $r_t \in \mathbb{R}$ , with  $\mathbb{E}[r_t | (s_t, a_t)]$  provided as a function  $R(s_t, a_t)$
- Next time step's state  $s_{t+1}$  can be generated by invoking a random sampling function  $SF(s_t, a_t)$ , i.e.,  $s_{t+1} = SF(s_t, a_t)(\cdot)$
- Discount factor denoted as  $\gamma$ , and  $r_T = 0$
- The problem is to calculate the Optimal Value function  $V_t^*(s_t)$
- Unlike tabular backward induction where state transition probabilities are given, here only a sampling function (for next state) is given
- Armed with the sampling function, can we do better than backward induction for the case where  $\mathcal{S}$  is very large and  $\mathcal{A}$  is small?

# Outline of AMS Algorithm

- AMS Algorithm is based on a fixed allocation of action selections for each state in each time step
- Denote number of action selections per state in time step  $t$  as  $N_t$
- Denote  $\hat{V}_t^{N_t}(s_t)$  as the AMS Algorithm estimate of  $V_t^*(s_t)$
- Let  $N_t^{s_t, a_t}$  be the number of selections of  $a_t$  for  $s_t$  ( $\sum_{a_t \in \mathcal{A}} N_t^{s_t, a_t} = N_t$ )
- Proportions of  $N_t^{s_t, a_t}$  based on Explore v/s Exploit UCT formula
- For each of the  $N_t^{s_t, a_t}$  selections of  $a_t$ , one next-state  $s_{t+1}$  is sampled
- Each  $s_{t+1} = SF(s_t, a_t)(\cdot)$  sample leads to recursive call  $\hat{V}_{t+1}^{N_{t+1}}(s_{t+1})$
- Optimal Action Value Function  $Q_t^*(s_t, a_t)$  estimated as:

$$\hat{Q}_t^{N_t}(s_t, a_t) = R(s_t, a_t) + \gamma \cdot \frac{\sum_{j=1}^{N_t^{s_t, a_t}} \hat{V}_{t+1}^{N_{t+1}}(SF(s_t, a_t)(\cdot))}{N_t^{s_t, a_t}}$$

- $V_t^*(s_t) = \max_{a_t} Q_t^*(s_t, a_t)$  approximated as:

$$\hat{V}_t^{N_t}(s_t) = \sum_{a_t} \frac{N_t^{s_t, a_t}}{N_t} \cdot \hat{Q}_t^{N_t}(s_t, a_t)$$

# The AMS Algorithm

## Algorithm 0.1: $\text{OPTVF}(t, s, N_t)$

**if**  $t == T$  **return** (0)

**comment:** Initialize VALS and CNTS by selecting each action once

**for**  $a \leftarrow \mathcal{A}$

**do**  $\begin{cases} \text{VALS}[a] \leftarrow \text{OptVF}(t+1, SF(s, a)(), N_{t+1}) \\ \text{CNTS}[a] \leftarrow 1 \end{cases}$

**for**  $i \leftarrow |\mathcal{A}|$  **to**  $N_t - 1$

**do**  $\begin{cases} \text{comment: Pick action based on UCB1 Explore v/s Exploit formula} \\ a^* \leftarrow \operatorname{argmax}_{a \in \mathcal{A}} (R(s, a) + \gamma \cdot \frac{\text{VALS}[a]}{\text{CNTS}[a]} + \sqrt{\frac{2 \ln i}{\text{CNTS}[a]}}) \\ \text{VALS}[a^*] \leftarrow \text{VALS}[a^*] + \text{OptVF}(t+1, SF(s, a^*)(), N_{t+1}) \\ \text{CNTS}[a^*] \leftarrow \text{CNTS}[a^*] + 1 \end{cases}$

**return**  $(\sum_{a \in \mathcal{A}} \frac{\text{CNTS}[a]}{N_t} \cdot (R(s, a) + \gamma \cdot \frac{\text{VALS}[a]}{\text{CNTS}[a]}))$

**comment:**  $N_t$  next-state samplings and  $N_t$  recursive calls to  $\text{OptVF}$

# Running Time, Bias, Convergence and Code

- Let  $N = \max(N_0, N_1, \dots, N_{T-1})$  and assume  $N > |\mathcal{A}|$
- Running time of AMS Algorithm is of the order of  $N^T \cdot |\mathcal{A}|$
- Compare this versus backward induction running time of  $|\mathcal{S}|^2 \cdot |\mathcal{A}| \cdot T$
- So AMS is more efficient when  $\mathcal{S}$  is very large (typical in real-world)
- [AMS paper](#) proves the estimate  $\hat{V}_0^{N_0}(s_0)$  is asymptotically unbiased

$$\lim_{N_0 \rightarrow \infty} \lim_{N_1 \rightarrow \infty} \dots \lim_{N_{T-1} \rightarrow \infty} \mathbb{E}[\hat{V}_0^{N_0}(s_0)] = V_0^*(s_0) \text{ for all } s_0 \in \mathcal{S}$$

- AMS paper also proves that the worst-possible bias is bounded by a quantity that converges to zero at rate  $O(\sum_{t=0}^{T-1} \frac{\ln N_t}{N_t})$

$$0 \leq V_0^*(s_0) - \mathbb{E}[\hat{V}_0^{N_0}(s_0)] \leq O\left(\sum_{t=0}^{T-1} \frac{\ln N_t}{N_t}\right) \text{ for all } s_0 \in \mathcal{S}$$

- Here's some [Python code for the AMS Algorithm](#) you can play with