

Lecture 2a: Generative Models - Monte Carlo Tree Search

Qi Zhang

Last Updated: September 2025

1 Online Planning with a Generative Model

Planning algorithms like value iteration (VI) and policy iteration (PI) solve Bellman (optimality) equations with the full knowledge of the MDP (e.g., the exact distribution over next states given any state-action pair). We here consider a more relaxed setting, where transition and reward functions (P, R) of the underlying (infinite-horizon) MDP $M = (\mathcal{S}, \mathcal{A}, P, R, \gamma)$ are unknown, but we can query them for any state-action pair to get samples, $s' \sim P(s, a)$, $r \sim R(s, a)$, or $s', r \sim M(s, a)$ as a shorthand. In this setting, we often say we have a *simulator* or a *generative model*. In this note, we restrictively consider MDPs where every trajectory ends in a terminal state after a finite number of transitions so we can safely choose $\gamma = 1$, examples including board games like chess and many others.

Moreover, VI and PI compute an optimal policy purely offline for all possible states and then look up the policy to select actions to take. While ensuring optimality, such *offline planning* methods are not feasible for large state spaces. In this note, we are instead interested in the *online planning* paradigm where planning and action execution is interleaved: at the current state, a certain computation budget is allocated to try to figure out a good action to take; this process is repeated after taking that action and transiting to a next state. This is much closer to how humans play games like chess.

2 Monte Carlo Tree Search - The Algorithmic Template

Monte Carlo Tree Search (MCTS) is a algorithmic technique to figure out a good action to take in current state by progressively building up a search tree where a parent-child connection represents a transition and nodes in the tree accumulates statistics from multiple sampled episodes. More specifically, initializing the tree with the only node being the current state as the root, MCTS repeatedly performs the following procedures as depicted in Figure 1 and outlined in 1 and 2:

1. Selection: From the root, navigate down the tree with transitions sampled from the generative model and some action-selection rule until a new transition occurs from a leaf node.
2. Expansion: Create a new leaf by attaching that new transition to the tree.
3. Simulation: Starting from the new leaf's state, simulate a trajectory by the generative model and some action-selection rule until reaching a terminal state.
4. Backpropagation: Using the rewards from the root to the terminal state, modify the value estimates along the path.

Each node τ tracks N_τ , the total number of sampled paths that visited this node's state, and G_τ^{total} (i.e., $\tau.\text{total_reward}$ in Algorithms 1 and 2), the total reward accumulated from that node to the terminal states over all N paths. To select a good action from the root node, we hope the empirical mean of $G_{\tau,a}^{\text{total}}/N_{\tau,a}$ well approximates Q-value of a reasonably good policy, where $G_{\tau,a}^{\text{total}}$ is the total

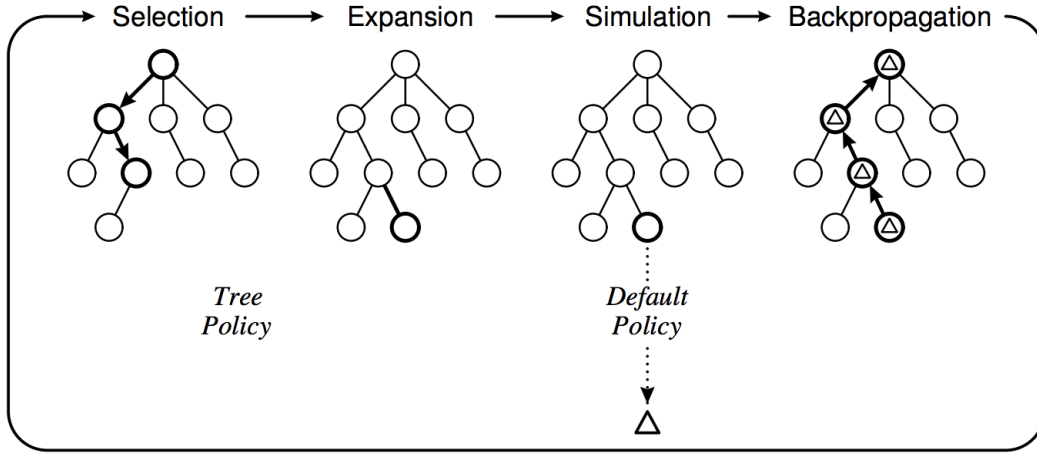


Figure 1: Illustration of an MCTS iteration (credit to [1]).

reward for all children of the node connected by a branch for action a (respectively for $N_{\tau,a}$). In the typical case where the state space is too huge to be fully explored, Selection (line 4 in Algorithm 2) should strike a balance between exploration of unknown paths and exploitation of historically promising paths. Expansion grows the tree, often by one leaf per iteration. We hope that each node on average gets a good number of sampled episodes to accumulate its statistics, so it is important not to grow the tree too fast. Selection and Expansion together are sometimes referred to as Tree Policy. Simulation aims to choose actions quickly, usually by some heuristics or just random action selection. Various algorithmic choices for these procedures give rise to specific instantiations. Sections ?? and 3 respectively discuss two well-known instantiations from this template, UCT and AlphaZero.

Implementation details

As above, the algorithmic description of MCTS is often based on the data structure of *state node*: a node is associated with a state s along paths in the tree that visit this node, with its children indexed by (a, s') for a transition of (s, a, s') . In hw2.ipynb of our Homework 2, we adopt an alternative implementation that explicitly distinguishes between the state node, referred to as the *V-node*, and the action node that is referred to as the *Q-node*. In a path, V-nodes are still associated with the states therein, while Q-nodes are associated with the actions. Therefore,

- The two types of nodes are interleaved depth-wise: a V-node has only Q-nodes (but not V-nodes) as its children and a Q-node has only V-nodes (but not Q-nodes) as its children;
- The root node is always a V-node of current state.
- The expansion happens upon an unseen transition (s, a, s') starting from a leaf V-node of state s and grows the tree with a new leaf V-node of state s' and a new Q-node of action a as its parent if a has not been tried, so the leaves are always V-nodes.
- Statistics such as $N_{\tau,a}$, $G_{\tau,a}^{\text{total}}$, and step reward are accumulated at the corresponding Q-nodes, and V-node can thus access its statistics by aggregating its Q-node children. For example, $N_{\tau} = \sum_a N_{\tau,a}$ where N_{τ} is the number of visits for a V-node τ , which equals to the summation over its Q-node children.

Algorithm 1 Monte Carlo Tree Search

```

1: procedure MONTECARLOTREESearch( $\tau_0$ )                                ▷ Root node  $\tau_0$ 
2:   repeat
3:      $\tau \leftarrow \text{TreePolicy}(\tau_0)$                                 ▷ New leaf  $\tau$ , total reward  $G_i$  cumulated from root  $\tau_0$ 
4:      $s \leftarrow \tau.\text{state}$ 
5:      $G \leftarrow \text{DefaultPolicy}(s)$                                 ▷  $G$  is total reward cumulated from state  $s$ 
6:      $\text{Backup}(\tau, G)$                                             ▷ Update node statistics along the path
7:   until  $\text{Timeout}()$ 
8:   return  $\text{BestAction}(\tau_0)$                                 ▷ Pick the approximately optimal root-level action
9: end procedure
10:
11: procedure INITNODE( $s$ )
12:    $\tau.\text{parent} \leftarrow \text{null}$ 
13:    $\tau.\text{state} \leftarrow s$ 
14:    $\tau.\text{count} \leftarrow 0$ 
15:    $\tau.\text{total\_reward} \leftarrow 0$                                 ▷ Cumulated over timesteps and episodes
16:    $\tau.\text{children}[a][s'] \leftarrow \text{null}$ , for all  $a, s'$ 
17:   return  $\tau$ 
18: end procedure

```

Algorithm 2 Monte Carlo Tree Search - Procedures

```

1: procedure TREEPOLICY( $\tau_0$ )
2:    $\tau, s \leftarrow \tau_0, \tau_0.\text{state}$ 
3:   while NONTERMINAL( $s$ ) do
4:      $a \leftarrow \text{SELECTACTION}(\tau)$  ▷ Heuristically select an action
5:      $s', r \sim M(s, a)$  ▷ Sample next state and reward from generative model
6:     if  $\tau.\text{children}[a][s'] = \text{null}$  then ▷ Is this the first observation of  $s \xrightarrow{a} s'$ ?
7:        $\tau' \leftarrow \text{INITNODE}(s')$  ▷ Initialize leaf node for state  $s'$ 
8:        $\tau'.\text{parent} \leftarrow \tau.\text{children}[a][s']$  ▷ Attach leaf node to the tree
9:        $\tau.\text{children}[a].\text{reward} \leftarrow r$  ▷ Record  $r = R(s, a)$ 
10:      return  $\tau.\text{children}[a][s']$  ▷ Move on to simulation phase
11:    end if
12:     $\tau.\text{children}[a].\text{reward} \leftarrow r$  ▷ Record  $r = R(s, a)$ 
13:     $\tau \leftarrow \tau.\text{children}[a][s']$ 
14:     $s \leftarrow s'$ 
15:  end while
16:  return  $\tau$ 
17: end procedure
18:
19: procedure DEFAULTPOLICY( $s$ ) ▷ Decision policy for simulation
20:    $G = 0$ 
21:   while NONTERMINAL( $s$ ) do
22:      $a \sim \text{DEFAULTACTIONSELECTION}(s)$  ▷ e.g., randomly select an action from  $\mathcal{A}$ 
23:      $s', r \sim M(s, a)$  ▷ Sample next state and reward from generative model
24:      $G \leftarrow G + r$ 
25:   end while
26:   return  $G$ 
27: end procedure
28:
29: procedure BACKUP( $\tau, G$ ) ▷  $G$  is the total reward cumulated from node  $\tau$ 
30:   repeat
31:      $\tau.\text{total\_reward} \leftarrow \tau.\text{total\_reward} + G$ 
32:      $\tau.\text{count} \leftarrow \tau.\text{count} + 1$ 
33:      $\tau, a \leftarrow \tau.\text{parent}$  ▷  $a$  is the action connecting the two nodes
34:     if  $\tau$  is not null then
35:        $G \leftarrow G + \tau.\text{children}[a].\text{reward}$  ▷ Add step reward along the path
36:     end if
37:   until  $\tau$  is null
38:   return
39: end procedure
40:
41: procedure BESTACTION( $\tau$ ) ▷ Pick the empirically best action at node  $\tau$ 
42:   return  $\arg \max_{a \in \mathcal{A}} \frac{\sum_{\tau' \in \tau.\text{children}[a][\cdot]} \tau'.\text{total\_reward}}{\sum_{\tau' \in \tau.\text{children}[a][\cdot]} \tau'.\text{count}}$ 
43: end procedure

```

3 UCT

UCT [2], Upper Confidence Bounds (UCB) applied to Trees, applies the UCB algorithm [3] to the Selection procedure of MCTS to balance the exploration-exploitation tradeoff. Specifically, UCT's Selection selects the action at each (V-)node τ (line 4 in Algorithm 2) that maximizes the upper confidence bound on the estimated total reward from that node:

$$\arg \max_{a \in \mathcal{A}} \text{UCB}(\tau, a) := \frac{G_{\tau,a}^{\text{total}}}{N_{\tau,a}} + c \sqrt{\frac{\ln(N_s)}{N_{\tau,a}}}$$

where $c > 0$ is a hyperparameter controlling how much exploration is favored over exploitation.

References

- [1] C. B. Browne, E. Powley, D. Whitehouse, S. M. Lucas, P. I. Cowling, P. Rohlfshagen, S. Tavener, D. Perez, S. Samothrakis, and S. Colton. A survey of monte carlo tree search methods. *IEEE Transactions on Computational Intelligence and AI in Games*, 4(1):1–43, March 2012.
- [2] Levente Kocsis and Csaba Szepesvári. Bandit based monte-carlo planning. In *Proceedings of the 17th European Conference on Machine Learning, ECML'06*, pages 282–293, Berlin, Heidelberg, 2006. Springer-Verlag.
- [3] Peter Auer, Nicolò Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine Learning*, 47(2):235–256, May 2002.