

Week-12 RL

1 Monte Carlo Tree Search

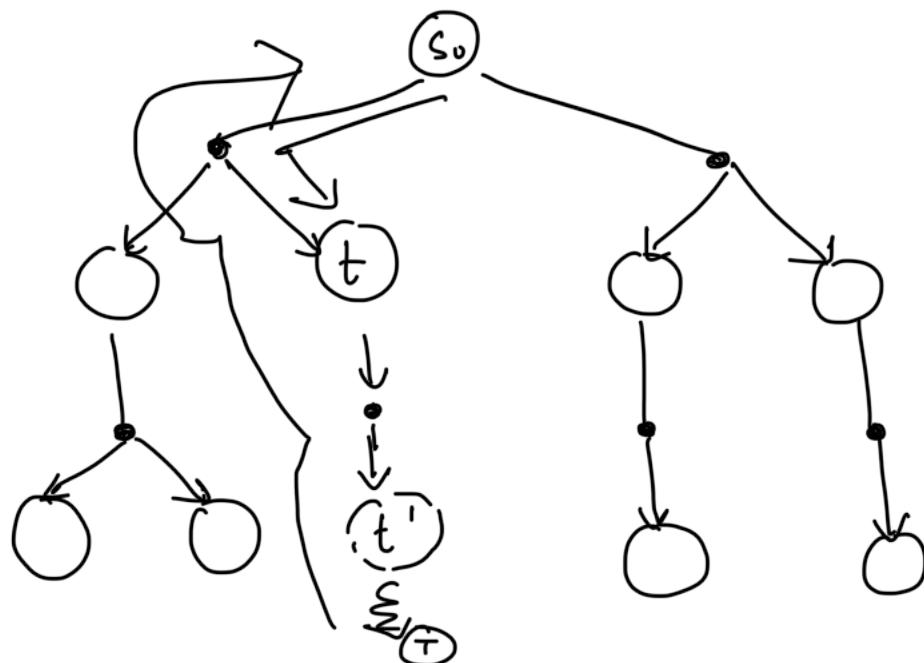
- Concept of roll outs

$$s_t \quad \pi(s_t) \quad \sim \pi(s_t, \cdot)$$

- Online planning \rightarrow planning is done immediately before executing an actⁿ. Once an actⁿ / seq. of acts is executed, we start planning from the new state. So, planning & execⁿ are interleaved.

1. for each state s , the set of acts $A(s)$ are partially evaluated.
2. $Q(s, a)$ is calculated by avg. the expected reward of traj. over s .
3. actⁿ chosen is $\text{argmax}_a Q(s, a)$.

- consider MDP as a tree -



1. select - start at the root node & successively select a child until you hit a node that isn't fully expanded.
 2. expansion - expand the children of the selected node by choosing an ac" with select" policy & creating new nodes using the ac" outcomes.
 3. Simula" - choose one of the new nodes & perform a random simula" of MDP to the terminating state with rollout policy.
 4. backup - given reward r at the terminating state, backup & save the reward to calculate the value $Q(s, a)$ at each state on the path.
- MCTS Search alg.

$\text{func}^{\sim}(\text{MCTS})(s_0)$

while budget do

selected node $\leftarrow \text{select}(s_0)$

child $\leftarrow \text{expand}(\text{selected node})$

$G \leftarrow \text{simulate}(\text{child})$

$\text{Backup}(\text{child}, G)$

return $\text{argmax}_a(Q(s_0, a))$

- properties of MCTS -
 1. heuristic - no domain specific knowledge.
 2. asym. - tree expansion in terms of finding the new nodes is asym.
 3. anytime - it backups the outcome of each simulaⁿ immediately which ensures all-values are up-to-date with every iteratⁿ of algs. So, we can return back to any stage at anytime.

- UCT (Upper confidence for trees)

$$\text{UCT} = \text{UCB}_1 + \text{MCTS}$$

$$\overbrace{\left(\text{argmax}_{a \in A(s)} \left[Q(s, a) + 2c \sqrt{\frac{2 \ln N(s)}{N(s, a)}} \right] \right)}$$

$N(s) \rightarrow$ no. of times a state node was visited.

$N(s, a) \rightarrow$ no. of times 'a' has been selected from this node.

$c > 0 \rightarrow$ exploration constant

- applications -

1. AlphaGo 0
2. Chess, Shogi
3. Real time strategy games
4. planning & scheduling