

RL: Lecture 19

Harvey Mudd College

April 8, 2020

Neil Rhodes

Course Logistics

MCTS HW 9 is out

~~Quiz 9 postponed~~

MCTS Programming Assignment 6
Midterm 2

Monday's lecture

Expected updates > sample updates?

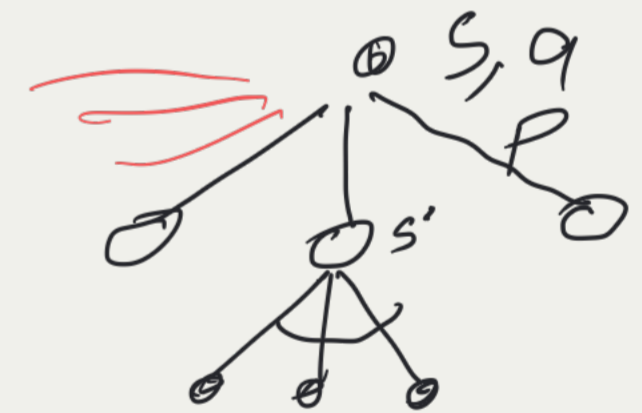
Sample update (less computation):

$$Q(s, a) = Q(s, a) + \alpha [R + \gamma \max_{a'} Q(S', a') - Q(s, a)]$$

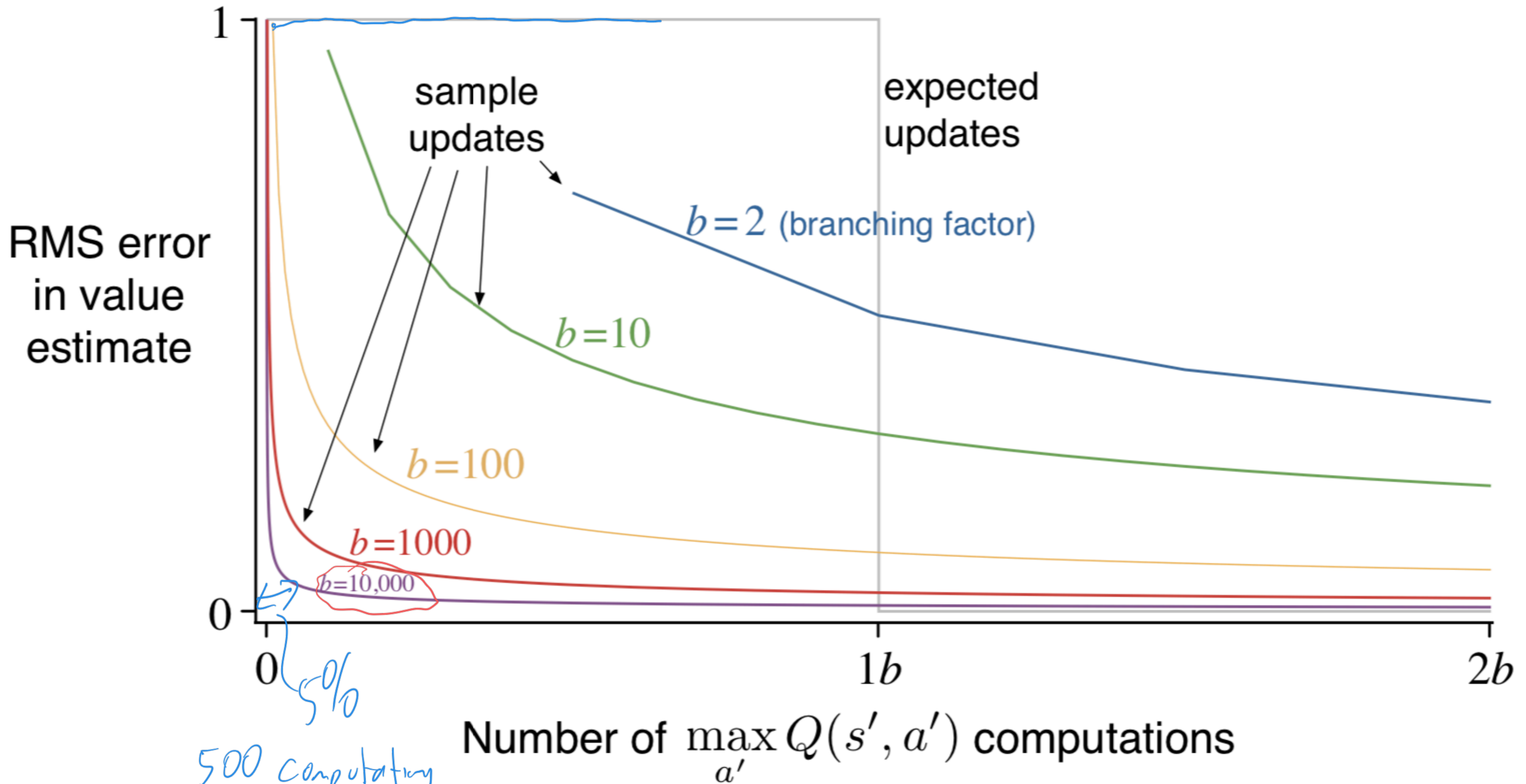
Expected update (more accurate):

$$Q(s, a) = \sum_{s', r} p(s', r | s, a) \alpha [r + \gamma \max_{a'} Q(s', a')]$$

could be expensive



Expected updates > sample updates?



- 1) Sample updates can
- 1) more quickly update $Q(s, a)$ allowing other updates to be more accurate
 - 2) approximate min of values in time it takes to do on expected update

Trajectory Sampling

alternative
to uniform
or prioritized sweeping

Sampling:

- Uniform

Dyna-Q⁺

choose a state-action
pair @ random
DO one-step update

vs.

- According to on-policy distribution

trajectory sampling

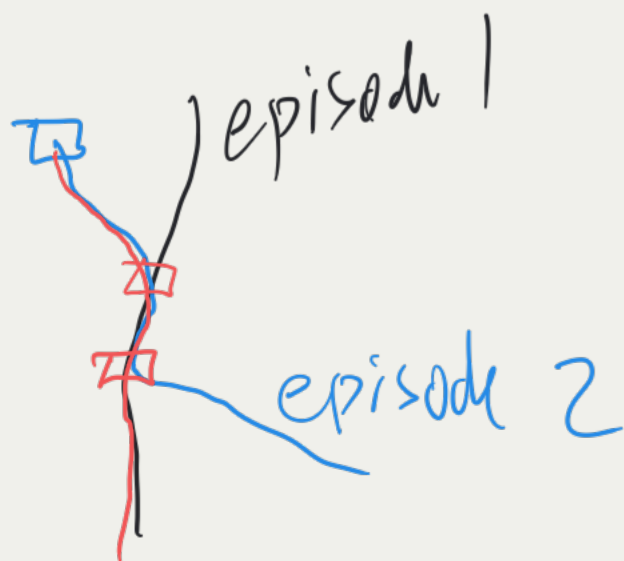
start w/ a start state

keep sampling until we run out
of time or complete an episode

episode seq of state/action
from environment

trajectory seq of state actions
from model

complete an episode



planning



When to plan

- In the background *Always be planning*
- At decision time
when a decision is to be made, do some planning

What is a rollout algorithm?

term comes from backgammon

these days
we don't physically
roll the dice

1. Begin with current state
2. Simulate trajectories starting with that state (following a rollout policy)
3. Choose action with highest estimated value

start w/ rollout policy

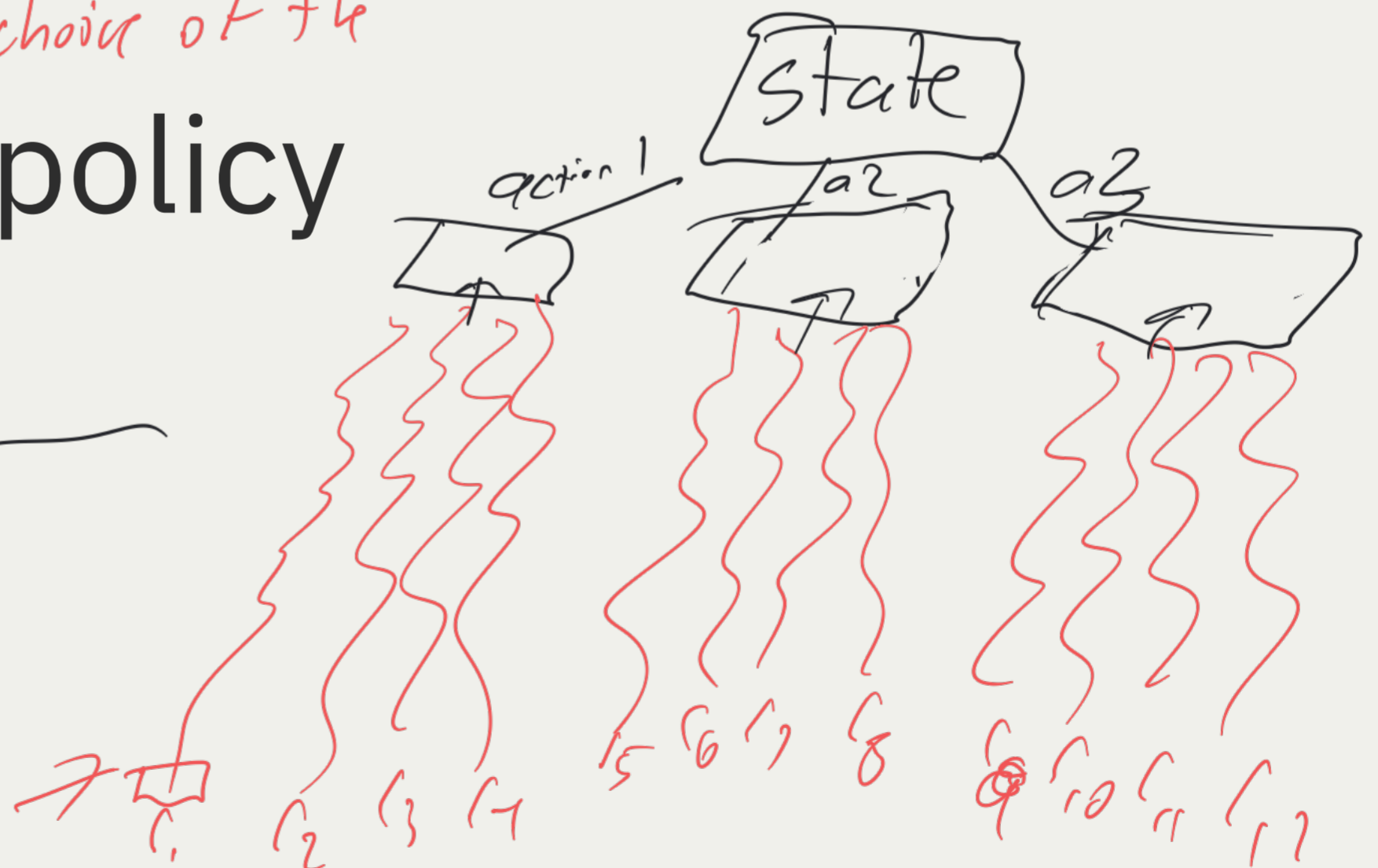
A rollout algorithm ^{current action choice of the}

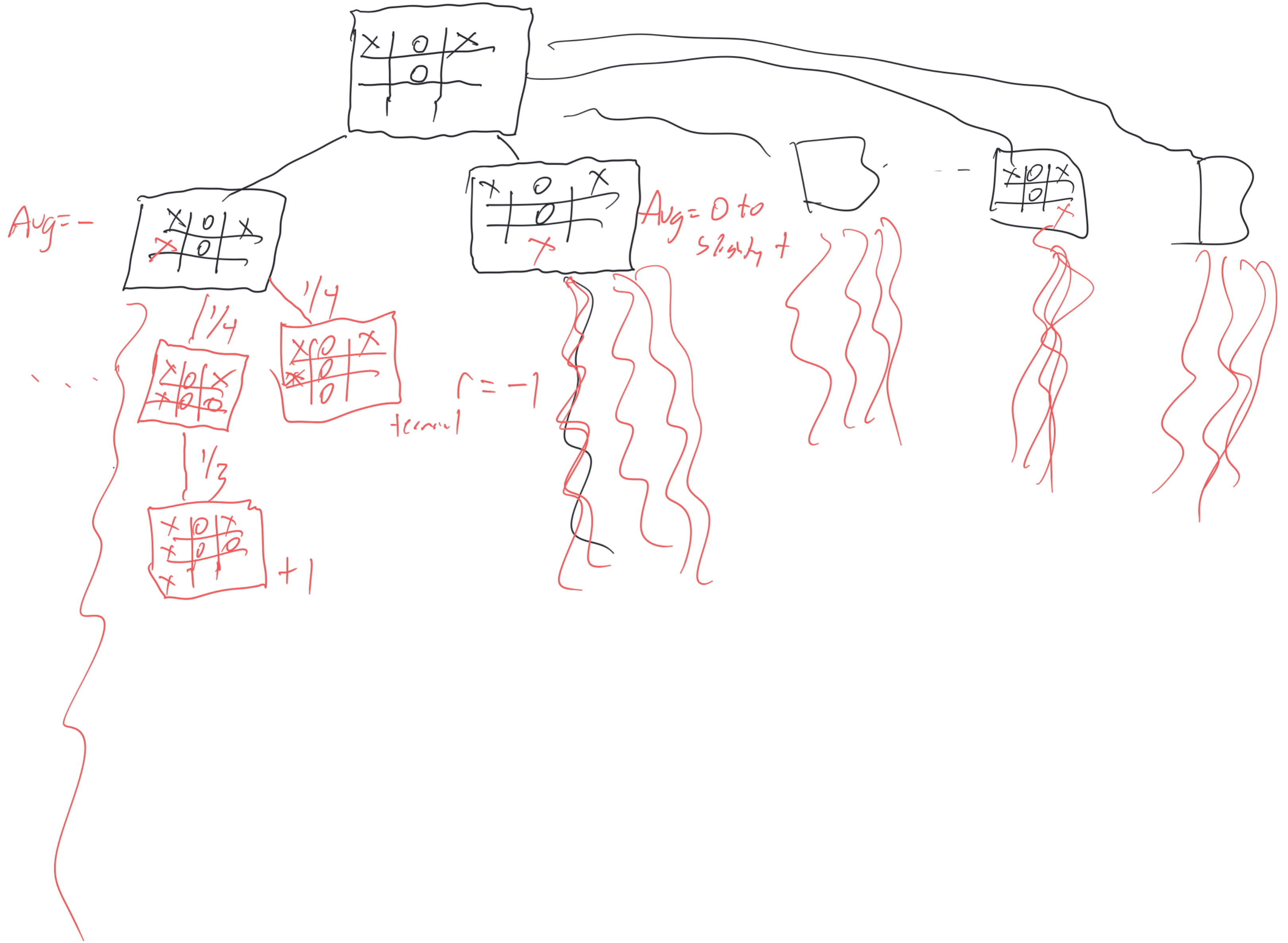
Improves the rollout policy

$$\pi(\text{start state}) \rightarrow a_3$$

often, the rollout policy
is uniform random

terminal
state

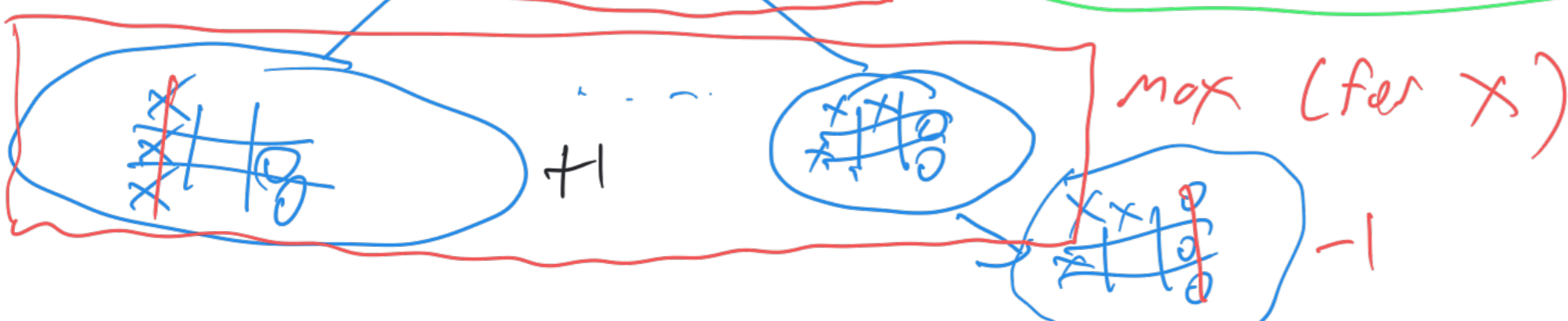
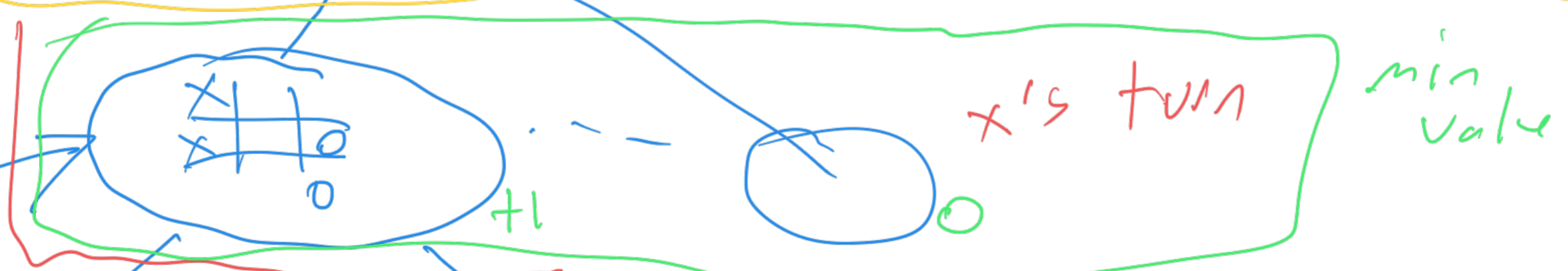
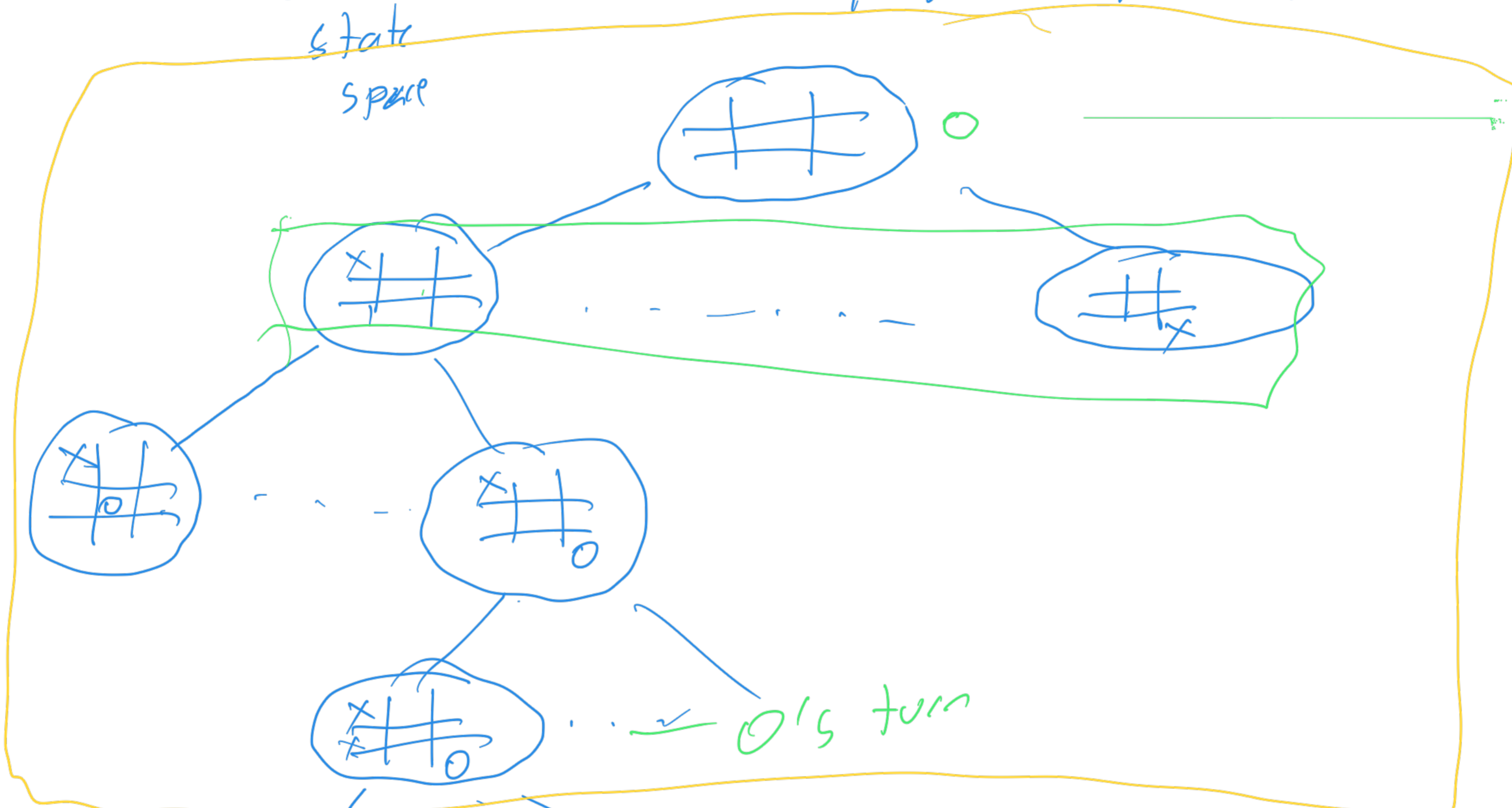




Searching a state space

Search state space

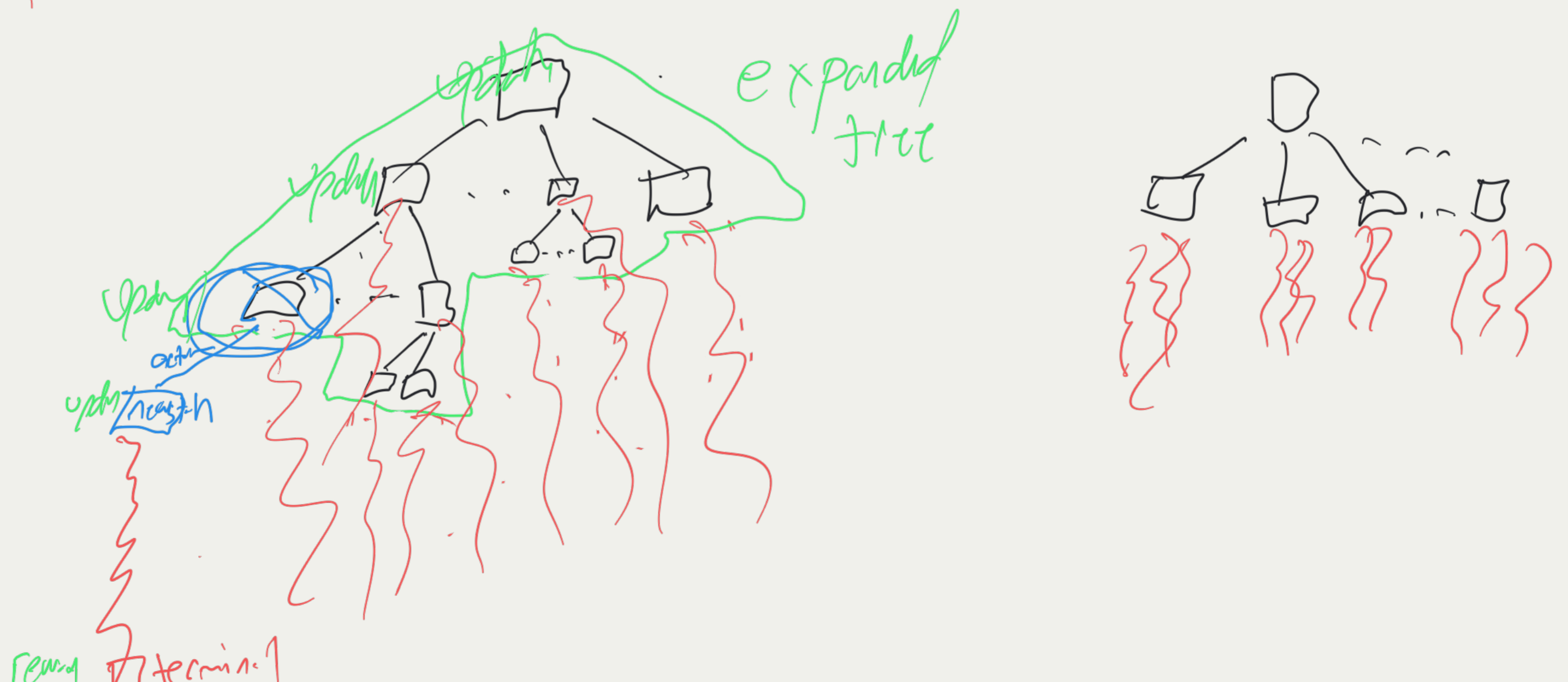
2-player competitive game



Monte Carlo Tree Search (MCTS)

Repeat while time remaining starting with current state:

1. Selection: Select a leaf node in the expanded tree
2. Expansion: Expand a child of the leaf node
3. Simulation: Follow rollout-policy from expanded node to simulate complete episode
4. Backup: Backup action values to nodes in the tree



Selection: choose a leaf node in the MCTS tree

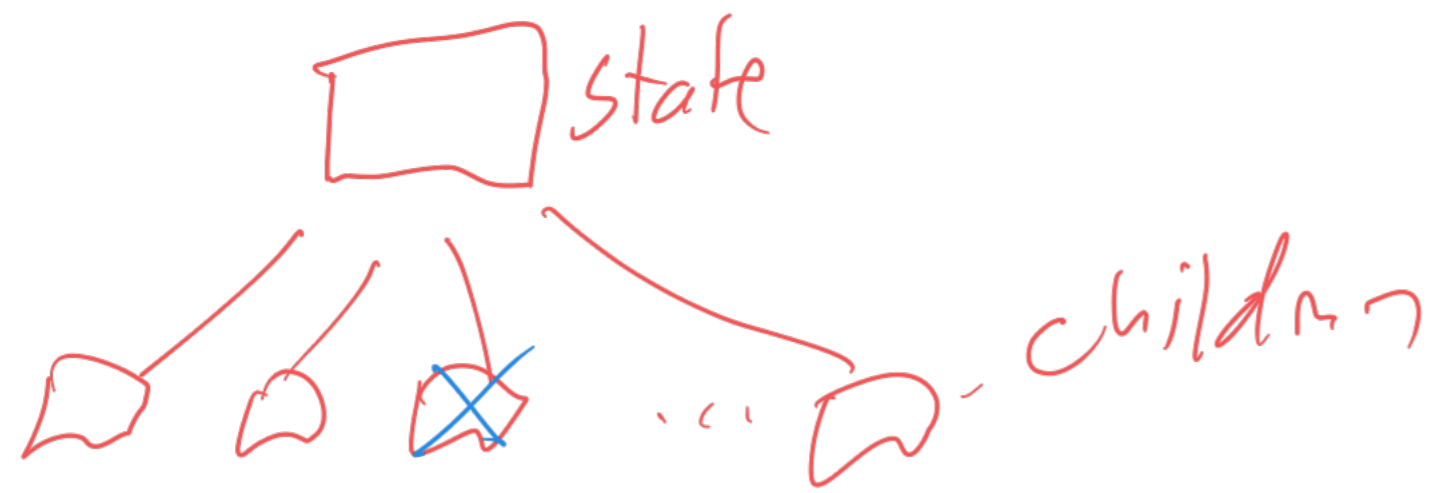
Traverse MCTS tree until reach a node with unexpanded children.

Use Upper Confidence Bound for Trees (UCT) to decide most promising child.

$q(v)$: Total simulation reward for node v
 $n(v)$: Total number of visits (simulation backups) for node v

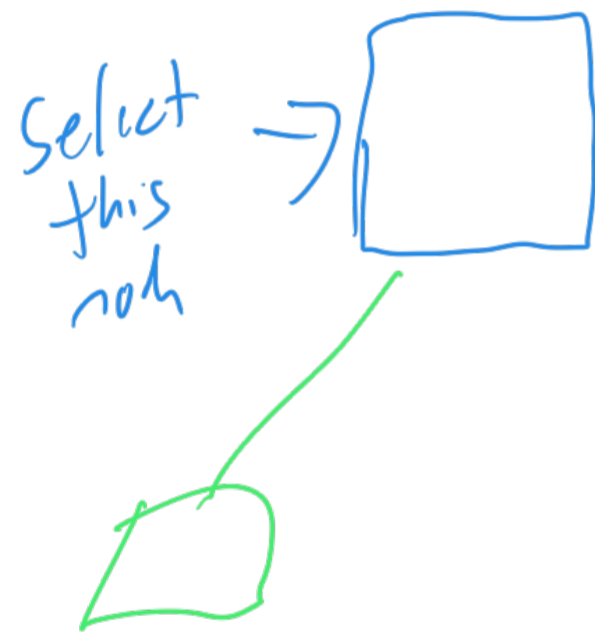
$$UCT(v) = \frac{q(v)}{n(v)} + c \sqrt{\frac{\log n(v.\text{parent})}{n(v)}}$$

Handwritten annotations:
- $q(v)$ is circled in red.
- $n(v)$ is circled in red.
- c is boxed in red.
- The fraction $\frac{\log n(v.\text{parent})}{n(v)}$ is boxed in red.
- $q(v)$ is labeled "Avg reward for v".
- c is labeled "tradeoff".
- The entire fraction is labeled "exploitation".
- "exploitation" is also written below the fraction.
- "max" is written above the text.

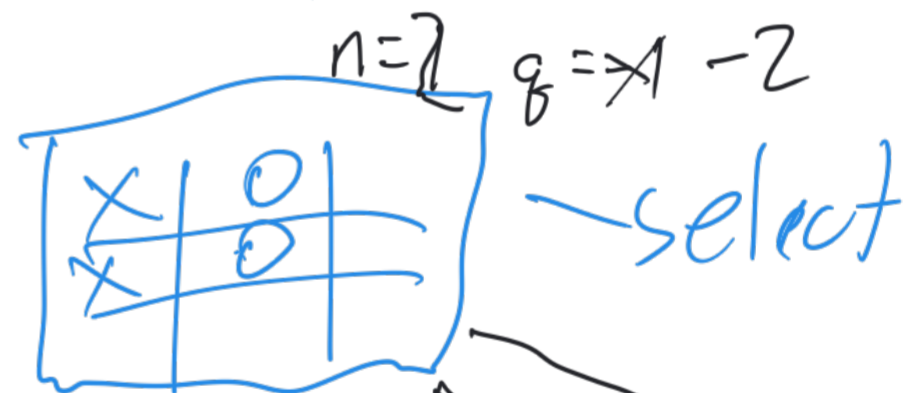


exploitation
 exploration
 use uct

first starting MCTS



expand
 pick random action



expand



$n=1$
 $g=+1$

$n=1$ $g=+1$
 expand

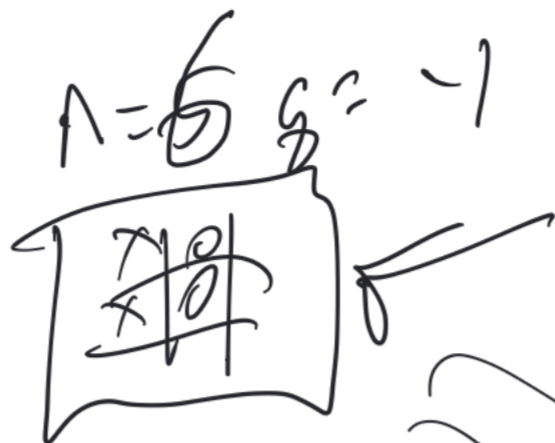


$n=+1$



$+1, -1$ or 0

X	0
X	0



$g=1$

$A=2 \quad g=2$



$A=1 \quad g=0$



$A=2 \quad g=-1$



$A=1 \quad g=1$



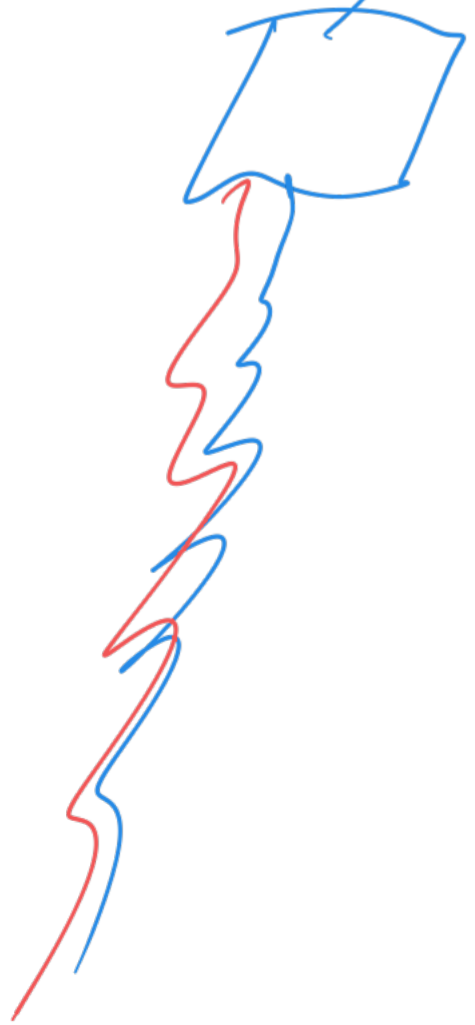
$A=1 \quad g=0$



highest g

because

$g=1$ + expump bonus



Expansion: expand selected node

If selected node is not terminal, choose an untried action and create a new MCTS node for the state that generates

Simulation: rollout starting at expanded node

Monte Carlo simulation using rollout policy until terminal state is reached.

Record total reward.

Backup: backup action values to nodes in the MCTS tree

Update $n(v)$ and $q(v)$ for each node v in the MCTS tree from simulation node up to root.

Note: If two-player competitive game, adjust reward to reflect who made the move.

For example, if reward is +1 (player 1 won):

- For v reached from player 1 move, increment $q(v)$
- For v reached from player 2 move, decrement $q(v)$

Selecting a final action

Probably don't want exploration term in UCT

- Child with highest $\frac{q(v)}{n(v)}$, or *highest avg reward*
- Child of root with highest $N(v)$ — it's the one that was explored the most so must have been most promising overall. *what if many states had high \uparrow one state with low \uparrow but high \uparrow*