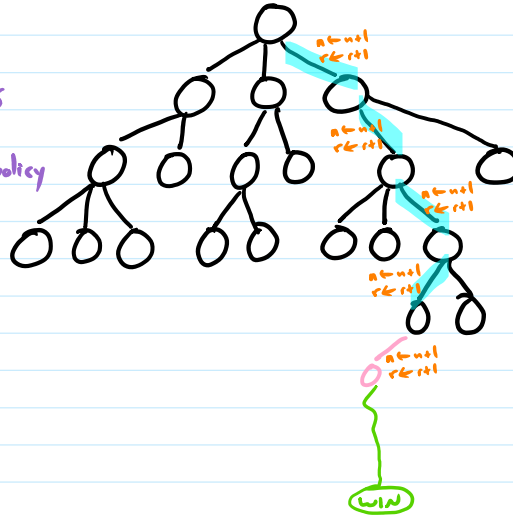


Monte Carlo Tree Search

UCT: UCB for trees
 ||
 MCTS + UCB for tree policy
 converges to minimax!



search tree: nodes = positions
 branches = moves
 children = resulting pos

\bar{r}_j = observed score over all
 playouts after traversing
 tree to child j

n_j = # times we traversed tree
 through child j

T = # times through
 current node

Until out of time

traverse tree root \rightarrow expandable node
 (not all children in tree)

tree policy

expand node (add missing child)

play from new node (playout)

update statistics along path through tree
 # times played
 total score

default policy

↑ random (fast)
 ↓ heuristic (long games hurt)

Advantages: always have move ready (gets better as time ↑)

no domain knowledge except rules

Multi-Armed Bandit

Given unknown probability distributions R_1, \dots, R_k with means μ_1, \dots, μ_k let $\mu^* = \max_i \mu_i$

Choose indices i_1, i_2, \dots to optimize payout

Regret = difference between observed reward and best possible expectation

$$R_T = T \cdot \mu^* - \sum_{t=1}^T \hat{r}_t \quad \hat{r}_t = \text{reward obtained on play } t$$

"optimal" means $P(\lim_{T \rightarrow \infty} \frac{R_T}{T} = 0) = 1$

	$\mu_1 = \frac{2}{3}$ Arm 1	$\mu_2 = \frac{7}{8}$ Arm 2	$\mu_3 = 2$ Arm 3
Ex: prob	$\frac{1}{3}$	$\frac{1}{4}$	$\frac{1}{100}$
payout	2	3	200
	$\frac{2}{3}$ 0	$\frac{1}{4}$ $\frac{1}{2}$	$\frac{59}{100}$ 0
		$\frac{1}{2}$ 0	

uniform rotation: $\frac{4}{3}$ $\frac{7}{8}$ 0 $\frac{4}{3}$ $\frac{7}{8}$ 0 $\frac{4}{3}$ $\frac{7}{8}$ 0 ...

$$\lim_{T \rightarrow \infty} \frac{R_T}{T} = \frac{59}{24}$$

$$\text{expected regret for } T=3n = \frac{\frac{4}{3}n + \frac{7}{8}n + 0n}{3n} = \frac{59}{24}$$

greedy: play each arm once (or n times), then arm with best observed reward forever

w/prob $\frac{1}{100}$	1	2	3	3	3	3	3	3	3	...
expected regret	$\frac{2}{3}$	0	200	0	0	0	0	0	0	...

$$\lim_{T \rightarrow \infty} \frac{R_T}{T} = \frac{\frac{2}{3} + 0}{T} = 0$$

w/prob $\frac{1}{3} \cdot \frac{1}{2} \cdot \frac{99}{100}$	1	2	3	1	1	1	1	1	...
expected regret	2	0	0	$\frac{4}{3}$	$\frac{4}{3}$

$$\lim_{T \rightarrow \infty} \frac{R_T}{T} = \frac{4}{3} \quad \text{so } P(\lim_{T \rightarrow \infty} \frac{R_T}{T} = 0) \neq 1$$

ϵ -greedy: play each once, then play random w/prob ϵ , arm with best avg reward so far otherwise

play best observed machine w/prob $(1-\epsilon) + \frac{\epsilon}{k}$ expected regret 0 if best observed = actual best

each other machine: w/prob $\frac{\epsilon}{k}$ $(\mu^* - \mu_i) > 0$

$$\text{expected regret } (1-\epsilon) + \frac{\epsilon}{k} \cdot 0 + \underbrace{\sum_{i \text{ not optimal}} \frac{\epsilon}{k} \cdot (\mu^* - \mu_i)}_{\text{positive}} > 0$$

zero regret (optimal!)

Choose arm j that maximizes

UCB : upper confidence bound

exploitation term

$$\bar{r}_j +$$

tunable parameter

$$\sqrt{\frac{2 \ln T}{n_j}}$$

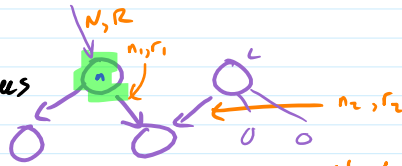
total plays over all arms

number of plays for option j

avg observed reward for arm j

exploration term

But many games are DAGs, not trees



standard UCT uses n_1, r_1, N in UCB formula

but can use n_1, n_2, r_1, r_2 for exploit term
 n_1, N for explore term

(using n_1, n_2, r_1, r_2 in both affects convergence)