# Multi-Armed Bandit

Given unknown probability distributions $R_1, \ldots, R_K$ with means $\mu_1, \ldots, \mu_K$ $\qquad \mu^* = \max_{1 \leq i \leq k} \mu_i$

Choose $i_1, i_2, i_3, \ldots$ to maximize total payout

Regret = difference between what you got and best possible expectation

$$\rho_T = T \cdot \mu^* - \sum_{t=1}^{T} \hat{P}_t \qquad \hat{P}_t = \text{reward at time } t$$

"optimal" means $P\left( \lim_{T \to \infty} \frac{\rho_T}{T} = 0 \right) = 1$

Ex:

| Arm 1 | | Arm 2 | | Arm 3 | |
|---|---|---|---|---|---|
| Prob | payout | prob | payout | prob | payout |
| $\frac{1}{3}$ | 2 | $\frac{1}{4}$ | 3 | $\frac{1}{100}$ | 200 |
| $\frac{2}{3}$ | 0 | $\frac{1}{4}$ | $\frac{1}{2}$ | $\frac{99}{100}$ | 0 |
| | | $\frac{1}{2}$ | 0 | | |
| $\mu_1 = \frac{2}{3}$ | | $\mu_2 = \frac{7}{8}$ | | $\mu_3 = 2 = \mu^*$ | |

uniform rotation: cycle through each arm

$$1 \ 2 \ 3 \quad 1 \ 2 \ 3 \quad \overbrace{1 \ 2 \ 3}, \ldots$$

expected regret = $\frac{4}{3} + \frac{9}{8} + 0 = \frac{59}{24}$ over 3 plays

$$\lim_{T \to \infty} \frac{\rho_T}{T} = \frac{59}{72} \quad \left( = \frac{59}{24} \cdot \frac{1}{3} \right)$$

greedy : play each once, then highest observed payout forever

$$1 \ 2 \ 3 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ \ldots \qquad \lim_{T \to \infty} \frac{\rho_T}{T} = \frac{4}{3} \quad \text{prob} > 0$$

$$1 \ 2 \ 3 \ 2 \ 2 \ 2 \ 2 \ 2 \ \ldots \qquad \lim_{T \to \infty} \frac{\rho_T}{T} = \frac{9}{8} \quad \text{prob} > 0$$

$$1 \ 2 \ 3 \ 3 \ 3 \ 3 \ \ldots \qquad \lim_{T \to \infty} \frac{\rho_T}{T} = 0 \quad \text{prob} < 1$$

$\varepsilon$-greedy: play 1 round, then play best observed reward w/ prob $1 - \varepsilon$ random arm w/ prob $\varepsilon$

$$\lim_{T \to \infty} \frac{\rho_T}{T} = \underbrace{\frac{\varepsilon}{k}}_{>0} \cdot \underbrace{\sum_{i=1}^{k} (\mu^* - \mu_i)}_{\substack{>0 \text{ for} \\ \text{non-optimal}}} > \underset{\substack{\text{(assuming not} \\ \text{all optimal)}}}{0}$$

tunable parameter
total plays
$\varepsilon$

zero regret

zero regret

Choose arm $j$ that maximizes $\overline{r_j}$ + $\sqrt{\dfrac{2 \cdot \ln T}{n_j}}$

UCB
(upper confidence bound)

tunable parameter

total plays

avg observed reward for j

# times arm $j$ played

exploitation          exploration

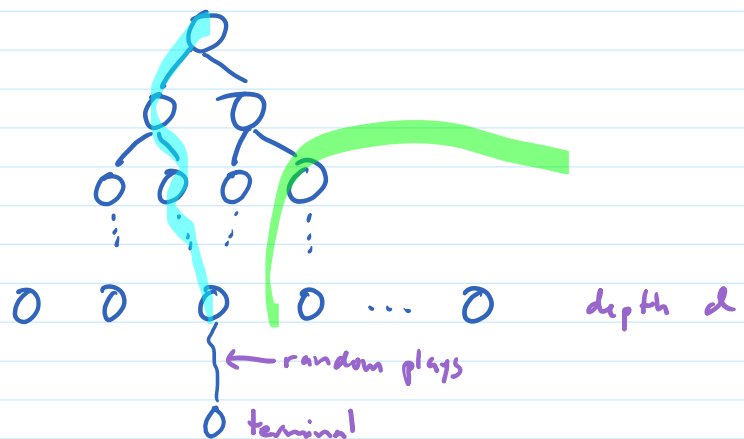$$P\left(\lim_{T \to \infty} \frac{P_T}{T} = 0\right) = 1$$

Flat Monte Carlo: for each action   simulate to terminal using random play

choose action w/ highest observed average

Scrabble, Bridge

Combine with UCB:   choose action to max   $\bar{r_j} \pm \sqrt{\frac{2 \ln T}{n_j}}$

Combine with tree search:   build tree to depth $d$
  (Flat UCB)   traverse to leaf using UCB at each level
  at leaf, play randomly to end
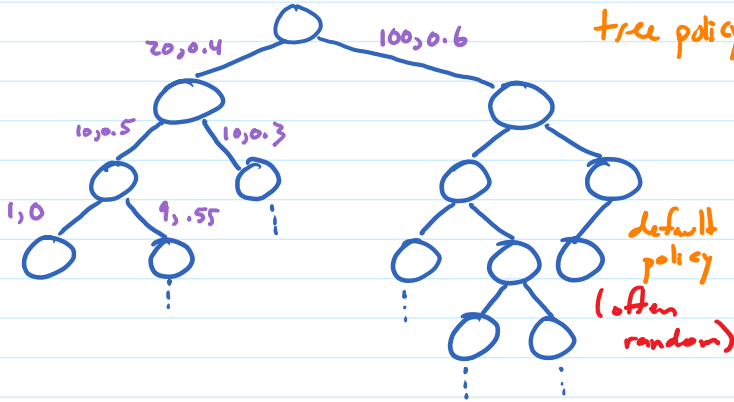  propagate stats back along path



← random plays

terminal

Grow Tree asymmetrically: Monte Carlo Tree Search

↳ explore good parts deeper than bad parts

depth $d$

# Monte Carlo Tree Search

while time left

tree policy traverse tree root ⟶ expandable node
or terminal node

(some missing children ↓)

add child

default policy
(often random)
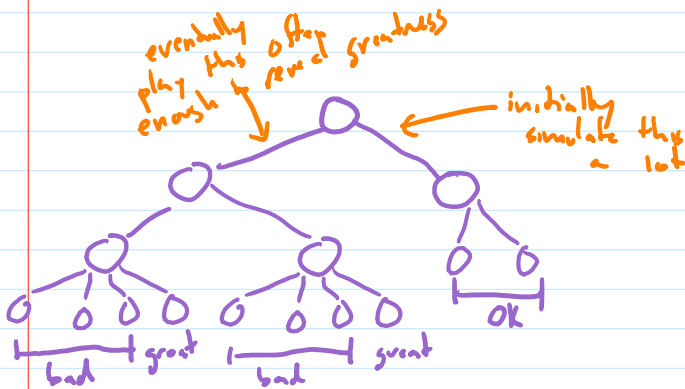playout from child

backpropagate result

return best child of root
↓
highest observed reward
(or played most often)
(or keep going until same)

eventually play this often enough → reveal greatness

initially simulate this a lot

ok

bad | great | bad | great

## UCT = MCTS + UCB as tree policy

advantages: convergent — converges to minimax

anytime — can suggest move after any iteration
(compare to iterative deepening)

no domain knowledge — no heuristic

easily parallelized

leaf parallel — multiple parallel playouts from leaf

root parallel — separate tree on each CPU
combine results in end

tree parallel — parallel traversals on same tree
(but needs locking)
if playout time >> traversal time
waits for locks insignif

disadvantages: no domain knowledge

some games not amenable

Tree vs DAG

default policy: random

move-averaged sampling technique (MAST)

PAST (predicate averaged sampling techniques)

|  | | $n$ | $r$ |
|---|---|---|---|
| $P_1$ : | A has more in store | 100 | 0.52 |
| $P_2$ : | B has more in store | 50 | 0.4 |
| $P_3$ : | A has no pit w/ >2 seeds | 50 | 0.38 |
| $P_4$ | ~$P_3$ | 100 | 0.55 |
| $P_5$ : | A has empty pits | 50 | 0.7 |
| $P_6$ : | ~$P_5$ | 100 | 0.3 |

(10)  (2)  (4)  (2)  (1)  (5)  (1)  (12)
      (3)  (1)  (◯)  (1)  (2)  (4)

MCTS tree data ⟶ playability of game?