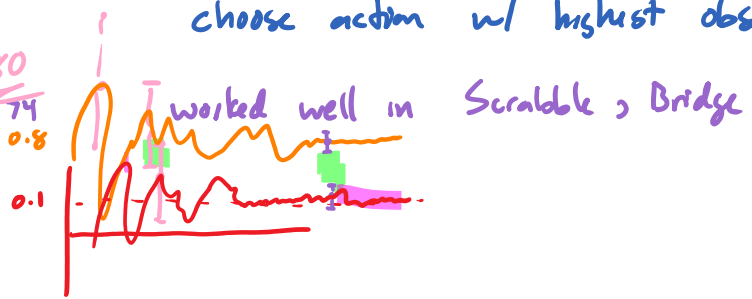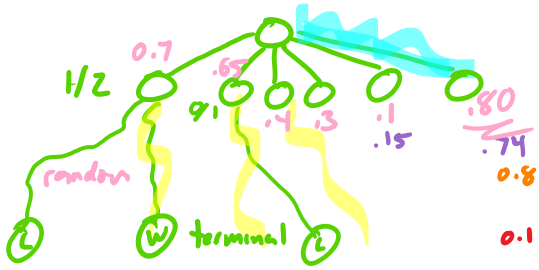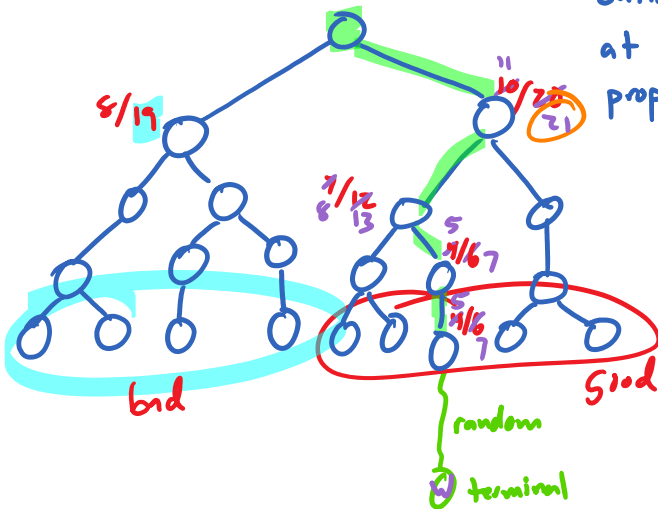**Monte Carlo Techniques**

Flat Monte Carlo : for each action — simulate to terminal using random play
sample many times
choose action w/ highest observed avg

worked well in Scrabble, Bridge

Combine with UCB: choose action max $\bar{r}_j + \sqrt{\dfrac{2 \ln T}{n_j}}$

Combine with tree search: (Flat UCB)

build tree to fixed depth
traverse to a leaf using UCB on children at each level
at leaf, do random playout
propagate result back up, accumulating stats in nodes

5/19

10/20

11

21

7/12
8
13

5
7/16
7

5
7/16
7

bad

Good

random

terminal

Grow Tree asymmetrically : Monte Carlo Tree Search

Multi-Armed Bandit

Given unknown probability distributions $R_1, ..., R_K$
     with means $\mu_1, ..., \mu_K$ $\quad \mu^* = \max_i \mu_i$

Choose indices $i_1, i_2, ....$ to optimize payout

Regret = difference between obtained reward and best possible expected reward

$$\rho_T = T \cdot \mu^* - \sum_{t=1}^{T} \hat{r}_t \qquad \hat{r}_t = \text{reward you got at time } t \text{ (from playing } i_t)$$

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

Ex:

| Arm 1 | | Arm 2 | | Arm 3 | |
|---|---|---|---|---|---|
| prob | payout | prob | payout | prob | payout |
| $\frac{1}{3}$ | 2 | $\frac{1}{4}$ | 5 | $\frac{1}{100}$ | 200 |
| $\frac{2}{3}$ | 0 | $\frac{1}{4} - \frac{1}{100}$ | $\frac{1}{2}$ | $\frac{99}{100}$ | 0 |
| | | $\frac{1}{2}$ | 0 | | |

$\mu_1 = \frac{1}{3} \cdot 2 + \frac{2}{3} \cdot 0 = \frac{2}{3}$

$\mu_2 = \frac{1}{4} \cdot 3 + \frac{1}{4} \cdot \frac{1}{2} + \frac{1}{2} = 0 = \frac{7}{8}$

$\mu_3 = 2$

uniform rotation: cycle through each arm $\quad 1 \ 2 \ 3 \ 1 \ 2 \ 3 \ \cdots$

$$\lim_{T \to \infty} \frac{\rho_T}{T} = \lim_{T \to \infty} \frac{T \cdot \mu^* - (\hat{r}_1 + \hat{r}_2 + \hat{r}_3 + \cdots)}{T}$$

$$= \lim_{T \to \infty} \frac{T\mu^* - (\hat{r}_1 + \hat{r}_4 + \hat{r}_7 + \cdots) - (\hat{r}_2 + \cdots) - (\hat{r}_3 + \cdots +)}{T}$$

$$= \mu^* - \left( \frac{1}{3}\mu_1 + \frac{1}{3}\mu_2 + \frac{1}{3}\mu_3 \right)$$
$$\quad\quad\quad\quad < \frac{1}{3}\mu^* + \frac{1}{3}\mu^* + \frac{1}{3}\mu^* = \mu^*$$

$$> 0$$

greedy : play each machine once, then machine w/ highest initial payoff

$\lim_{T \to \infty} \frac{\rho_T}{T} = \mu^* - \mu' \quad$ prob I do this $= \frac{1}{3} \cdot \frac{3}{4} \cdot \frac{99}{100} = \frac{99}{400} > 0$
$\quad\quad\quad 2 - \frac{2}{3} > 0$

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

$1 \ 2 \ 3 \quad 1 \ 1 \ 1 \ 1 \cdots$
$1 \ 2 \ 3 \quad 2 \ 2 \cdots$
$1 \ 2 \ 3 \quad 3 \ 3 \ 3 \cdots$

$1 \ 2 \ 3 \ 1 \ 2 \ 2 \ 2 \ 3 \ 2 \ 2 \ 2 \ 1 \ 2 \ 2 \ 3 \ 3 \ 3 \ 3 \ 1 \ 3 \ 3 \ 3 \ 2 \ 3 \ 3 \ 3 \ 3 \cdots$

$\varepsilon$-greedy : play one round, then $\quad$ play best observed w/ prob $1 - \varepsilon$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ play randomly $\quad$ w/ prob $\varepsilon$

$$\lim_{T \to \infty} \frac{\rho_T}{T} > \underbrace{\frac{\varepsilon}{k}}_{>0} \cdot \underbrace{\frac{(\mu^* - \mu'^n)}{>0}}_{} > 0$$

$\cdots \quad \overline{r}_2 \quad \cdots \quad \overline{r}_5 \quad$

tunable parameter

exploration

total number of plays

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

zero regret

zero regret

Choose arm $j$ that maximizes $\bar{r}_j + \sqrt{\dfrac{2 \ln T}{n_j}}$

exploration

total number of plays

UCB
upper confidence bound

mean observed reward of arm $j$

exploitation

number of times played arm $j$

# Monte Carlo Tree Search

**mCTS+UCB = UCT**

Until out of time

traverse tree    root → leaf    *(for example, UCB)* tree policy

expand    if leaf is expandable, add its children
→ nonterminal and non-zero visits

simulate    play to terminal pos (from arb. selected added child if expanded)
default policy (for ex. random)

update    backpropagate stats along path through tree
from point simulation started → root

Select move from root based on current stats (choose child w/
highest avg or highest visit count)

P1 = max

P2 = min

P1

P2

P1



Oct 28^J 2021 Page 4