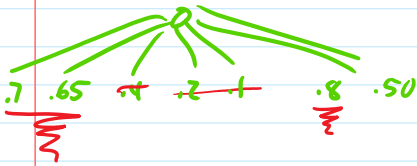


Monte Carlo Techniques

Flat Monte Carlo: for each action simulate to terminal using random play
choose action with highest observed average

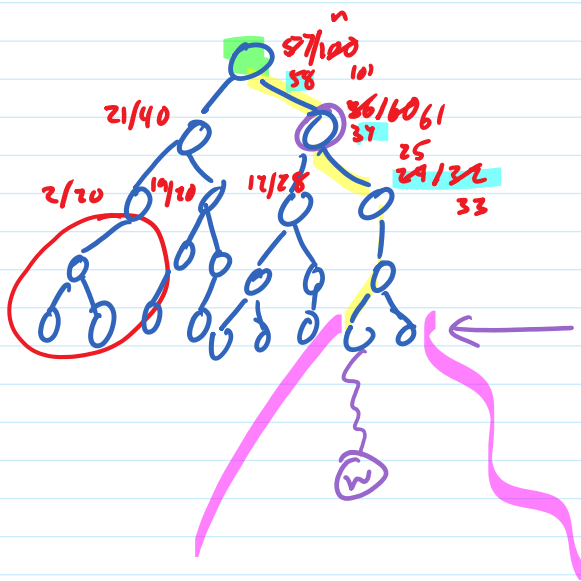
worked well in Scrabble, Bridge



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

Combine with tree search:
(Flat UCB)

build tree to depth d
traverse to leaf using UCB at each level
at leaf do random playout (move randomly until end)
propagate stats back



Grow Tree asymmetrically: Monte Carlo Tree Search

Multi-Armed Bandit

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

Choose indices i_1, i_2, \dots to optimize payout

Regret = difference between reward obtained and best possible expectation

$$P_T = \underbrace{T}_{\text{total plays so far}} \cdot \mu^* - \sum_{t=1}^T \hat{r}_t$$

\hat{r}_t = reward at time t
 (what we got from arm i_t at time t)

"optimal" means "zero regret"

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

Arm 1	Arm 2	Arm 3																				
<table border="1"> <tr><th>Prob</th><th>Payout</th></tr> <tr><td>$\frac{1}{3}$</td><td>2</td></tr> <tr><td>$\frac{2}{3}$</td><td>0</td></tr> </table>	Prob	Payout	$\frac{1}{3}$	2	$\frac{2}{3}$	0	<table border="1"> <tr><th>Prob</th><th>Payout</th></tr> <tr><td>$\frac{1}{4}$</td><td>3</td></tr> <tr><td>$\frac{1}{2}$</td><td>0</td></tr> <tr><td>$\frac{1}{4}$</td><td>0</td></tr> </table>	Prob	Payout	$\frac{1}{4}$	3	$\frac{1}{2}$	0	$\frac{1}{4}$	0	<table border="1"> <tr><th>Prob</th><th>Payout</th></tr> <tr><td>$\frac{1}{100}$</td><td>200</td></tr> <tr><td>$\frac{99}{100}$</td><td>0</td></tr> </table>	Prob	Payout	$\frac{1}{100}$	200	$\frac{99}{100}$	0
Prob	Payout																					
$\frac{1}{3}$	2																					
$\frac{2}{3}$	0																					
Prob	Payout																					
$\frac{1}{4}$	3																					
$\frac{1}{2}$	0																					
$\frac{1}{4}$	0																					
Prob	Payout																					
$\frac{1}{100}$	200																					
$\frac{99}{100}$	0																					
$\mu_1 = \frac{1}{3} \cdot 2 + \frac{2}{3} \cdot 0 = \frac{2}{3}$	$\mu_2 = \frac{1}{4} \cdot 3 + \frac{1}{4} \cdot \frac{1}{2} + \frac{1}{2} \cdot 0 = \frac{7}{8}$	$\mu_3 = \frac{1}{100} \cdot 200 + 0 = 2 = \mu^*$																				

uniform rotation: cycle through each arm 1 2 3 1 2 3 ...

$$\lim_{T \rightarrow \infty} \frac{P_T}{T} = \lim_{T \rightarrow \infty} \frac{T \cdot \mu^* - (\hat{r}_1 + \hat{r}_2 + \hat{r}_3 + \dots)}{T}$$

$$\frac{T \cdot \mu^*}{T} - \lim_{T \rightarrow \infty} \frac{(\hat{r}_1 + \hat{r}_2 + \hat{r}_3 + \dots)}{T} = 0$$

$$\mu^* - \lim_{T \rightarrow \infty} \frac{(\hat{r}_1 + \hat{r}_2 + \hat{r}_3 + \dots)}{T} = \frac{\hat{r}_1}{T} + \frac{\hat{r}_2}{T} + \frac{\hat{r}_3}{T} + \dots$$

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

$$\frac{1}{3} (\mu^* - \mu_1) + \frac{1}{3} (\mu^* - \mu_2) > 0$$

greedy: play each once, play arm with highest payout forever

$\frac{1}{3} \cdot \frac{3}{4} \cdot \frac{99}{100} > 0$

1 2 3 | | | | |

$\lim_{T \rightarrow \infty} \frac{P_T}{T} = \mu^* - \mu_1 > 0$

$\frac{1}{3} \cdot \frac{3}{4} \cdot \frac{99}{100} > 0$
 $P(\text{this}) > 0$
 $= \frac{1}{4} \cdot \frac{3}{5} \cdot \frac{99}{100}$

$\lim_{T \rightarrow \infty} \frac{P_T}{T} = \mu^* - \mu_1 > 0$
 $\lim_{T \rightarrow \infty} \frac{P_T}{T} = \mu^* - \mu_2 > 0$
 $\lim_{T \rightarrow \infty} \frac{P_T}{T} = 0$

$P(\lim_{T \rightarrow \infty} \frac{P_T}{T} = 0) = \frac{1}{5} \cdot \frac{99}{100} + \dots > 0$

ϵ -greedy : play one round, play best observed with prob $1 - \epsilon$
 play randomly with prob ϵ

$$\lim_{T \rightarrow \infty} \frac{P_T}{T} > \frac{\epsilon}{K} \cdot (\mu^* - \mu_i) > 0$$

prob of playing non-optimal arm i is at least this

zero regret
Aver

Choose arm j that maximizes

UCB
(upper confidence bound)

$\bar{r}_j + \sqrt{\frac{2 \ln T}{n_j}}$

exploitation
 observed mean reward for arm j

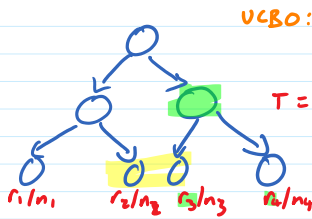
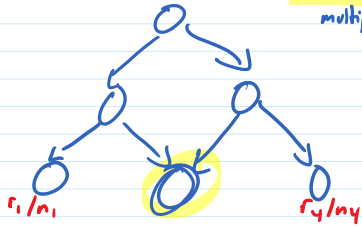
tunable parameter
 total plays
 total plays for arm j
 exploration term

chess has many - MCTS alone doesn't work well

so doesn't → MCTS works well

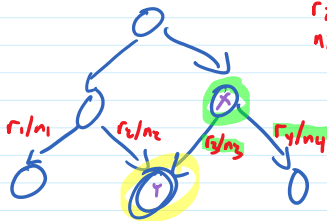
MCTS adapted for games that aren't trees

(some states reachable through multiple paths)



UCBO: ignore multiple paths; duplicate states if reachable along different paths

$$T = n_2 + n_4$$



r_i = total reward after traversing edge
 n_i = total times edge traversed

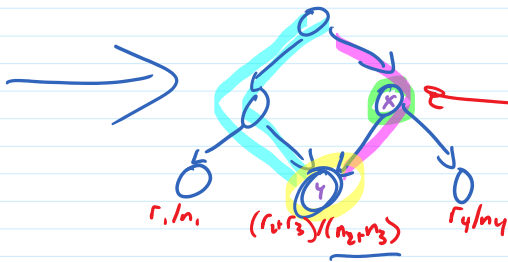
UCB1: ignore other paths (= UCBO but without duplicates)

ex: when evaluating $Y \otimes X$ compute

$$\frac{r_2}{n_2} + \sqrt{\frac{c \ln(n_2 + n_4)}{n_2}}$$

UCB2: combine multiple edges in exploit term

$$\frac{r_2 + r_3}{n_2 + n_3} + \sqrt{\frac{c \ln(n_2 + n_4)}{n_2 + n_3}}$$



no guarantee on convergence
(that I have seen)

$$\frac{r_2 + r_3}{n_2 + n_3} + \sqrt{\frac{c \ln(n_2 + n_4)}{n_2 + n_3}}$$