

Monte Carlo Tree Search

default policy: random

move-averaged sampling technique (MAST)

maintain rewards for actions taken during play out
bias random moves towards higher reward

	n	F
L	10000	.55
R	9000	.42
	9000	

	n	F
capture	100	.60
move again	110	.53
other	240	.40

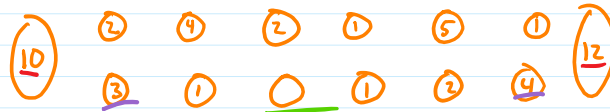
current play out
LLRRRLRRLL wins

predicate average (PAST)

establish predicates
maintain rewards for moves satisfying predicates
bias towards moves to pos that satisfy preds w/ higher rewards

- P₁: A has more in store
- P₂: B has more in store
- P₃: A has no pit w/ > 2 seeds
- P₄: ~P₃
- P₅: A has empty pits
- P₆: ~P₅

P	n	F
P ₁	100+1	0.58
P ₂	50	0.4
P ₃	50	0.34
P ₄	100+1	0.55
P ₅	50+1	0.7
P ₆	100	0.5

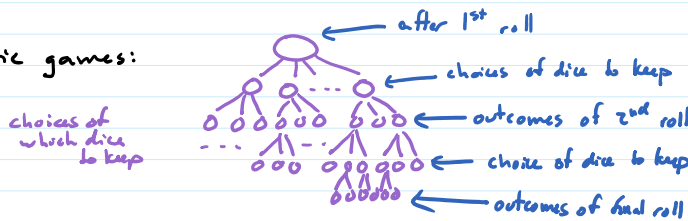


P₁ satisfied
P₄ satisfied
P₅ satisfied

$$\frac{.52 + .55 + .7}{3} = \dots$$

leads to W

stochastic games:



tree policy: UCB at choice nodes
random at random event nodes

imperfect information games:

determinization: pick values for unknown info

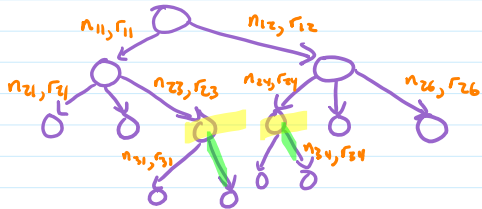
MCTS tree data → playability of game?

variance of score
spread (max-min)
⋮

not too hard/easy
not unbalanced
not drawish
not too long or short

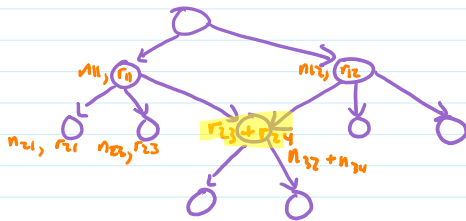
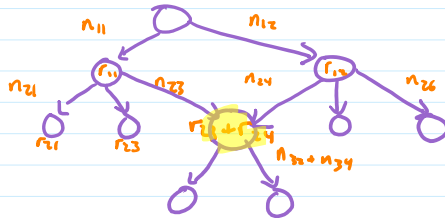
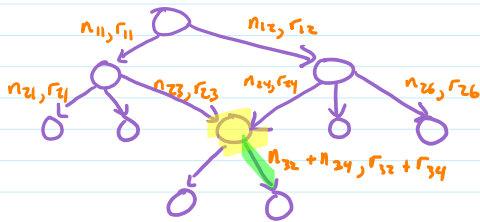
MCDS?

directed acyclic graph (DAG)



n_{ij} = times move played
 r_{ij} = total reward after move

same position after
2 seq of moves



store graph as a map
position → statistics
edges are implicit (given by left-moves
+ result)

Two-player games

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

if r_j always reward for P I

@ P II nodes, modify form