# Foundations of Artificial Intelligence
## 45. Monte-Carlo Tree Search: Advanced Topics

Thomas Keller

Universität Basel

May 30, 2016

## Board Games: Overview

chapter overview:

- 41. Introduction and State of the Art
- 42. Minimax Search and Evaluation Functions
- 43. Alpha-Beta Search
- 44. Monte-Carlo Tree Search: Introduction
- 45. Monte-Carlo Tree Search: Advanced Topics
- 46. AlphaGo and Outlook

Optimality
●○○○○○

Tree Policy
○○○○○○○○○○○○○○○○

Other Techniques
○○○○○

Summary
○○

# Optimality of MCTS

# Reminder: Monte-Carlo Tree Search

- as long as time allows, perform iterations
  - selection: traverse tree
  - expansion: grow tree
  - simulation: play game to final position
  - backpropagation: update utility estimates
- execute move with highest utility estimate

Optimality
○○●○○○

Tree Policy
○○○○○○○○○○○○○○○○

Other Techniques
○○○○○

Summary
○○

## Monte-Carlo Tree Search: Updated Pseudo-Code

---

**function** visit_node(*tree*, *n*)

**if** is_final(*n*.state):
    **return** $u$(*n*.state)
$s = tree$.get_unvisited_successor(*n*)
**if** $s \neq$ **none**:
    $n' = tree$.add_child_node(*n*, *s*)
    *utility* = apply_default_policy()
    backup(*n'*, *utility*)
**else**:
    $n' =$ apply_tree_policy(*n*)
    *utility* = visit_node(*tree*, *n'*)
backup(*n*, *utility*)
**return** *utility*

---

## Optimality

complete "minimax tree" computes optimal utility values $Q^*$

Optimality
○○○○●○

Tree Policy
○○○○○○○○○○○○○○○○○

Other Techniques
○○○○○

Summary
○○

## Asymptotic Optimality

### Asymptotically Optimality

An MCTS algorithm is asymptotically optimal if $\hat{Q}^k(n)$ converges to $Q^*(n)$ for all $n \in \text{succ}(n_0)$ with $k \to \infty$.

## Asymptotic Optimality

### Asymptotically Optimality

An MCTS algorithm is asymptotically optimal if $\hat{Q}^k(n)$ converges to $Q^*(n)$ for all $n \in \text{succ}(n_0)$ with $k \to \infty$.

Note: there are MCTS instantiations that play optimally even though the values do not converge in this way
(e.g., if all $\hat{Q}^k(n)$ converge to $l \cdot Q^*(n)$ for a constant $l > 0$)

# Asymptotic Optimality

For a tree policy to be asymptotically optimal, it is required that it

- explores forever:
  - every position is expanded eventually and visited infinitely often (given that the game tree is finite)
  - after a finite number of iterations, only true utility values are used in backups
- is greedy in the limit:
  - the probability that the optimal move is selected converges to 1
  - in the limit, backups based on iterations where only an optimal policy is followed dominate suboptimal backups

Optimality
○○○○○○

Tree Policy
●○○○○○○○○○○○○○○○○

Other Techniques
○○○○○

Summary
○○

# Tree Policy

Optimality
○○○○○○

Tree Policy
○●○○○○○○○○○○○○○○○

Other Techniques
○○○○○

Summary
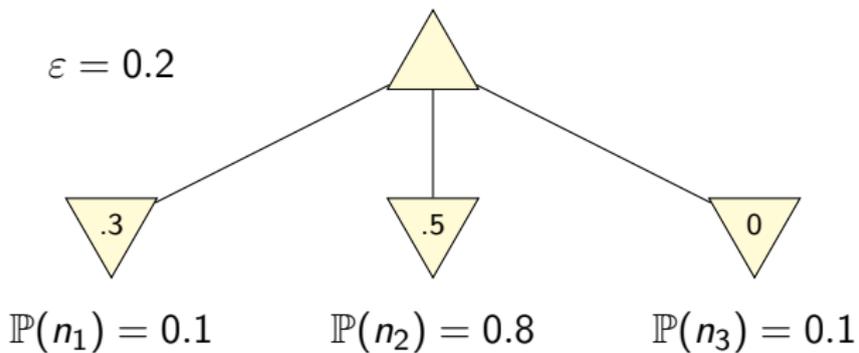○○

## Objective

tree policies have two contradictory objectives:

- explore parts of the game tree that have not been investigated thorroughly
- exploit knowledge about good moves to focus search on promising areas

central challenge: balance exploration and exploitation

Optimality
○○○○○○

Tree Policy
○○●○○○○○○○○○○○○○

Other Techniques
○○○○○

Summary
○○

# $\varepsilon$-greedy: Idea

- tree policy with constant parameter $\varepsilon$
- with probability $1 - \varepsilon$, pick the greedy move (i.e., the one that leads to the successor node with the highest utility estimate)
- otherwise, pick a non-greedy successor uniformly at random

Optimality
oooooo

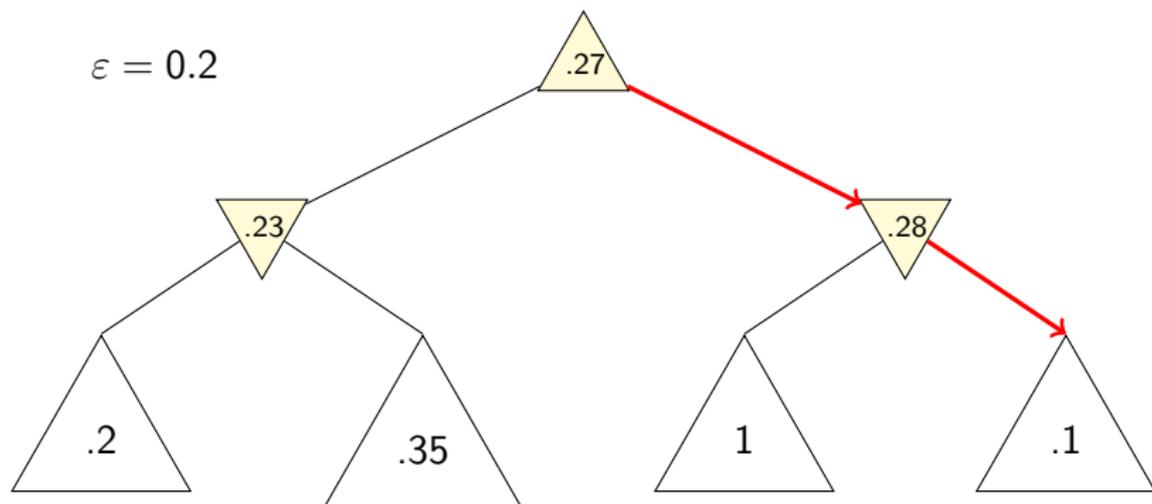Tree Policy
oooooooooooooooo

Other Techniques
ooooo

Summary
oo

# $\varepsilon$-greedy: Example



$\varepsilon = 0.2$

$\mathbb{P}(n_1) = 0.1 \qquad \mathbb{P}(n_2) = 0.8 \qquad \mathbb{P}(n_3) = 0.1$

Optimality
○○○○○○

Tree Policy
○○○○●○○○○○○○○○○○○

Other Techniques
○○○○○

Summary
○○

# $\varepsilon$-greedy: Asymptotic Optimality

## Asymptotic Optimality of $\varepsilon$-greedy

- explores forever
- not greedy in the limit
- $\Rightarrow$ not asymptotically optimal



$\varepsilon = 0.2$

Optimality
○○○○○○

Tree Policy
○○○○●○○○○○○○○○○

Other Techniques
○○○○○

Summary
○○

# $\varepsilon$-greedy: Asymptotic Optimality

### Asymptotic Optimality of $\varepsilon$-greedy

- explores forever
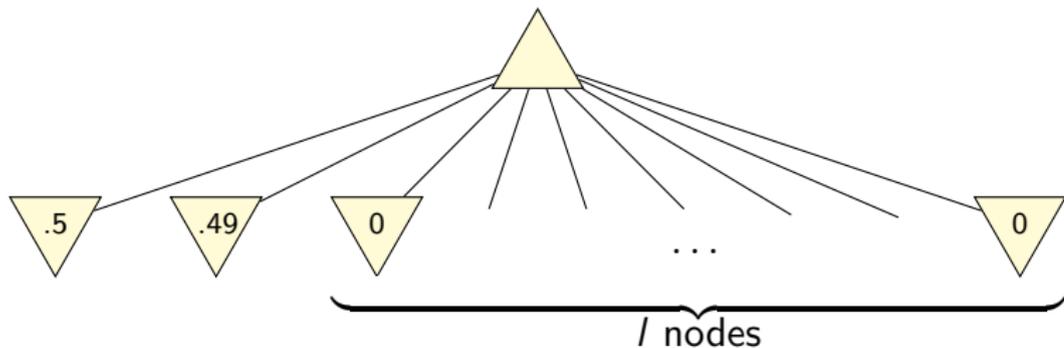- not greedy in the limit
- $\Rightarrow$ not asymptotically optimal

asymptotically optimal variants:

- use decaying $\varepsilon$, e.g. $\varepsilon = \frac{1}{k}$
- use minimax backups

Optimality
oooooo

Tree Policy
ooooo●oooooooooo

Other Techniques
ooooo

Summary
oo

## $\varepsilon$-greedy: Weakness

Problem:
when $\varepsilon$-greedy explores, all non-greedy moves are treated equally



$l$ nodes

e.g., $\varepsilon = 0.2, l = 9$: $\mathbb{P}(n_1) = 0.8$, $\mathbb{P}(n_2) = 0.02$

Optimality
oooooo

Tree Policy
ooooooo●oooooooo

Other Techniques
ooooo

Summary
oo

## Softmax: Idea

- tree policy with constant parameter $\tau$
- select moves proportionally to their utility estimate
- Boltzmann exploration selects moves proportionally to
  $$\mathbb{P}(n) \propto e^{\frac{\hat{Q}(n)}{\tau}}$$

Optimality
○○○○○○

Tree Policy
○○○○○○○●○○○○○○○

Other Techniques
○○○○○

Summary
○○

## Softmax: Example



e.g., $\tau = 0.1, l = 9$: $\mathbb{P}(n_1) \approx 0.51$, $\mathbb{P}(n_2) \approx 0.46$

Optimality
○○○○○○

Tree Policy
○○○○○○○○○●○○○○○○○

Other Techniques
○○○○○

Summary
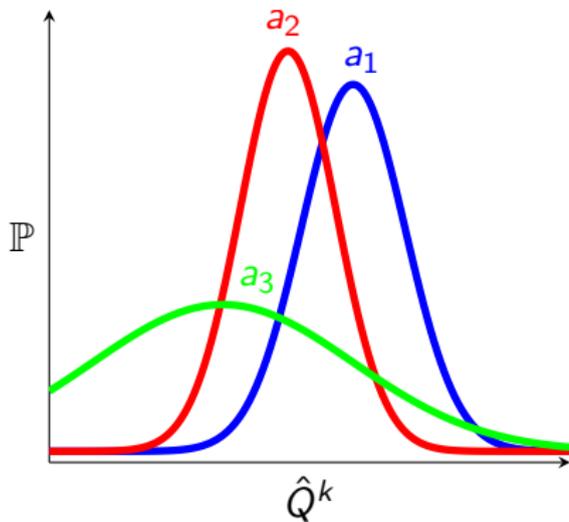○○

# Boltzmann Exploration: Asymptotic Optimality

## Asymptotic Optimality of Boltzmann Exploration

- explores forever
- not greedy in the limit
  (probabilities converge to positive constant)
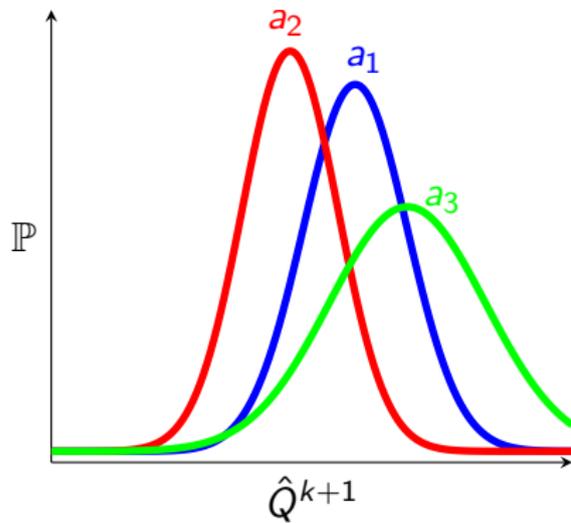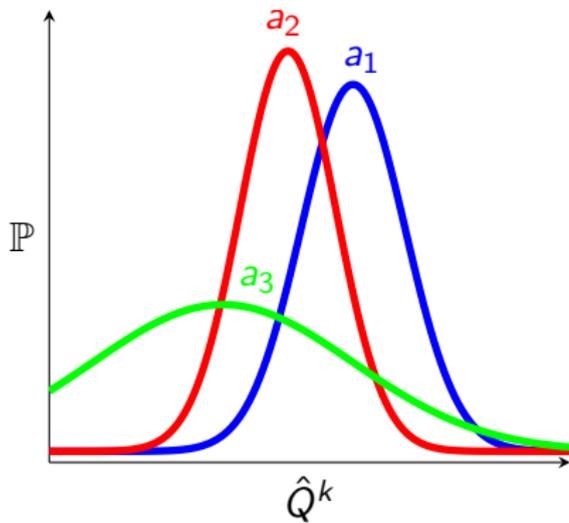- $\Rightarrow$ not asymptotically optimal

asymptotically optimal variants:

- use decaying $\tau$
- use minimax backups

Optimality
oooooo

Tree Policy
oooooooooo●ooooooo

Other Techniques
ooooo

Summary
oo

# Boltzmann Exploration: Weakness

Optimality
○○○○○○

Tree Policy
○○○○○○○○○○●○○○○○○

Other Techniques
○○○○○

Summary
○○

# Boltzmann Exploration: Weakness

## Upper Confidence Bounds: Idea

balance exploration and exploitation by preferring moves that

- have been successful in earlier iterations (exploit)
- have been selected rarely (explore)

# Upper Confidence Bounds: Idea

## Upper Confidence Bounds

- select successor $n'$ of $n$ that maximizes $\hat{Q}(n') + \hat{U}(n')$
- based on utility estimate $\hat{Q}(n')$
- and a bonus term $\hat{U}(n')$
- select $\hat{U}(n')$ such that $Q^*(n') \leq \hat{Q}(n') + \hat{U}(n')$ with high probability
- $\hat{Q}(n') + \hat{U}(n')$ is an upper confidence bound on $Q^*(n')$ under the collected information

Optimality
oooooo

Tree Policy
oooooooooooooo●ooo

Other Techniques
ooooo

Summary
oo

## Upper Confidence Bounds: UCB1

- use $\hat{U}(n') = \sqrt{\frac{2 \cdot \ln N(n)}{N(n')}}$ as bonus term
- bonus term is derived from Chernoff-Hoeffding bound:
    - gives the probability that a sampled value (here: $\hat{Q}(n')$)
    - is far from its true expected value (here: $Q^*(n')$)
    - in dependence of the number of samples (here: $(N(n'))$
- picks the optimal move exponentially more often

Optimality
oooooo

Tree Policy
ooooooooooooo●oo

Other Techniques
ooooo

Summary
oo

# Upper Confidence Bounds: Asymptotic Optimality

## Asymptotic Optimality of UCB1

- explores forever
- greedy in the limit
- $\Rightarrow$ asymptotically optimal

Optimality
○○○○○○

Tree Policy
○○○○○○○○○○○○○○○●○○

Other Techniques
○○○○○

Summary
○○

# Upper Confidence Bounds: Asymptotic Optimality

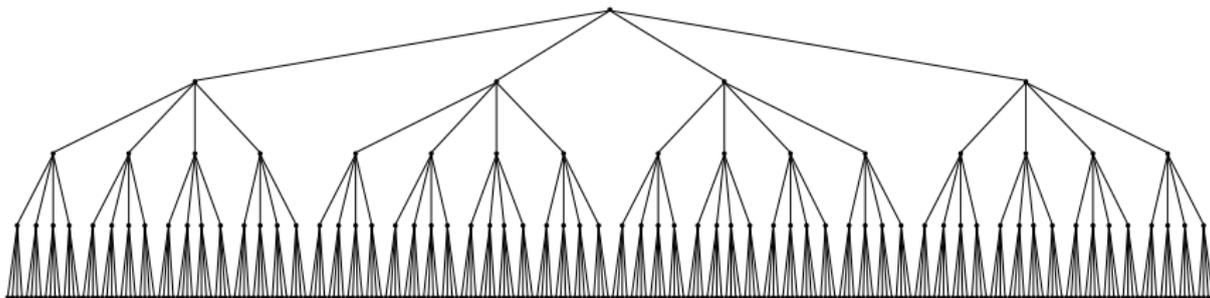## Asymptotic Optimality of UCB1

- explores forever
- greedy in the limit
- $\Rightarrow$ asymptotically optimal

However:

- no theoretical justification to use UCB1 in trees or planning scenarios
- development of tree policies active research topic
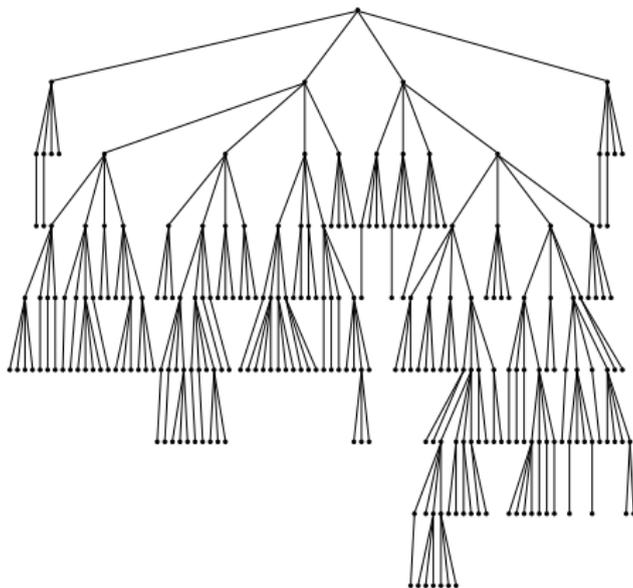
# Tree Policy: Asymmetric Game Tree

full tree up to depth 4

## Tree Policy: Asymmetric Game Tree

UCT tree (equal number of search nodes)

Optimality
oooooo

Tree Policy
oooooooooooooooo

Other Techniques
●oooo

Summary
oo

# Other Techniques

Optimality
○○○○○○

Tree Policy
○○○○○○○○○○○○○○○○

Other Techniques
○●○○○○

Summary
○○

## Default Policy: Instantiations

default: Monte-Carlo Random Walk

- in each state, select a legal move uniformly at random
- very cheap to compute
- uninformed
- usually not sufficient for good results

Optimality
○○○○○○

Tree Policy
○○○○○○○○○○○○○○○

Other Techniques
○●○○○

Summary
○○

# Default Policy: Instantiations

default: Monte-Carlo Random Walk

- in each state, select a legal move uniformly at random
- very cheap to compute
- uninformed
- usually not sufficient for good results

only significant alternative: domain-dependent default policy

- hand-crafted
- offline learned function

Optimality
oooooo

Tree Policy
ooooooooooooooooo

Other Techniques
oo●oo

Summary
oo

# Default Policy: Alternative

- default policy simulates a game to obtain utility estimate
- ⇒ default policy must be evaluated in many positions
- if default policy is expensive to compute, simulations are expensive
- solution: replace default policy with heuristic that computes a utility estimate directly

Optimality
○○○○○○

Tree Policy
○○○○○○○○○○○○○○○○

Other Techniques
○○○●○

Summary
○○

## Other MCTS Enhancements

there are many other techniques to increase information gain from iterations, e.g.,

- All Moves As First
- Rapid Action Value Estimate
- Move-Average Sampling Techique
- and many more

Literature: A Survey of Monte Carlo Tree Search Methods
Browne et. al., 2012

Optimality
oooooo

Tree Policy
oooooooooooooooo

Other Techniques
ooooo●

Summary
oo

## Expansion

- to proceed deeper into the tree, each node must be visited at least once for each legal move
- $\Rightarrow$ deep lookaheads not possible
- rather than add a single node, expand encountered leaf node and add all successors
  - allows deep lookaheads
  - needs more memory
  - needs initial utility estimate for all children
- alternative solution without significant memory overhead:
  - ignore restriction that unvisited successors must be created
  - move annotations to parent node

Optimality
oooooo

Tree Policy
oooooooooooooooo

Other Techniques
ooooo

Summary
●o

# Summary

Optimality
○○○○○○

Tree Policy
○○○○○○○○○○○○○○○○

Other Techniques
○○○○○

Summary
○●

## Summary

- tree policy is crucial for MCTS
  - $\epsilon$-greedy favors the greedy move and treats all other equally
  - Boltzmann exploration selects moves proportionally to their utility estimates
  - UCB1 favors moves that were successful in the past or have been explored rarely
- there are applications for each where they perform best
- good default policies are domain-dependent and hand-crafted or learned offline
- using heuristics instead of a default policy often pays off