

MDP Algorithms

Part II

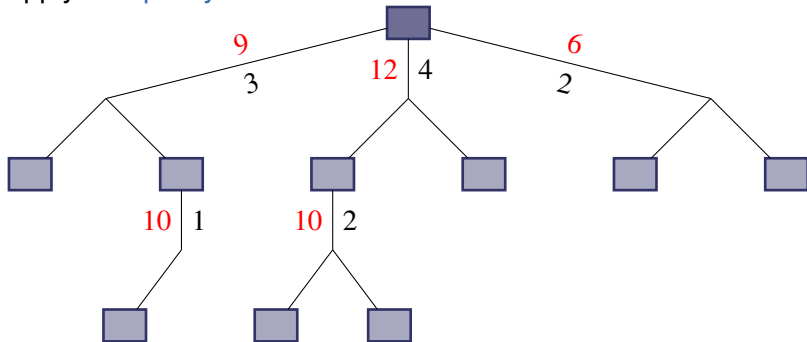
Thomas Keller

University of Basel

June 20, 2018

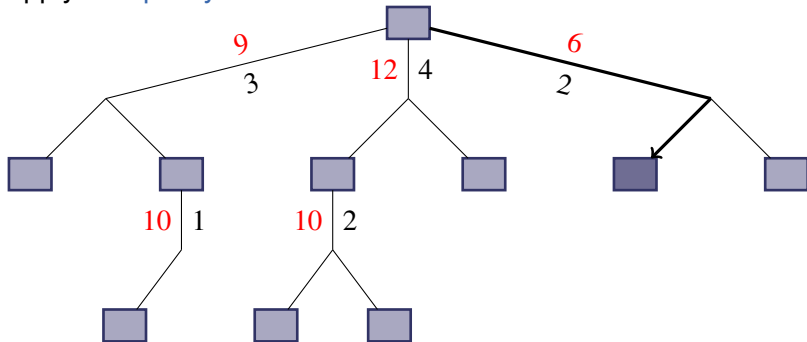
Selection Phase

apply **tree policy** to traverse tree



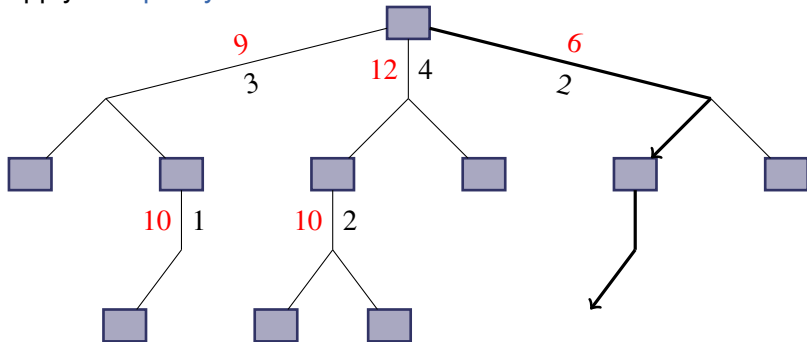
Selection Phase

apply **tree policy** to traverse tree



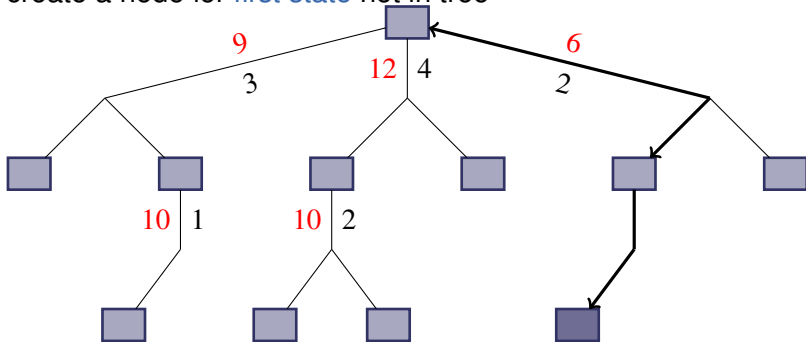
Selection Phase

apply **tree policy** to traverse tree



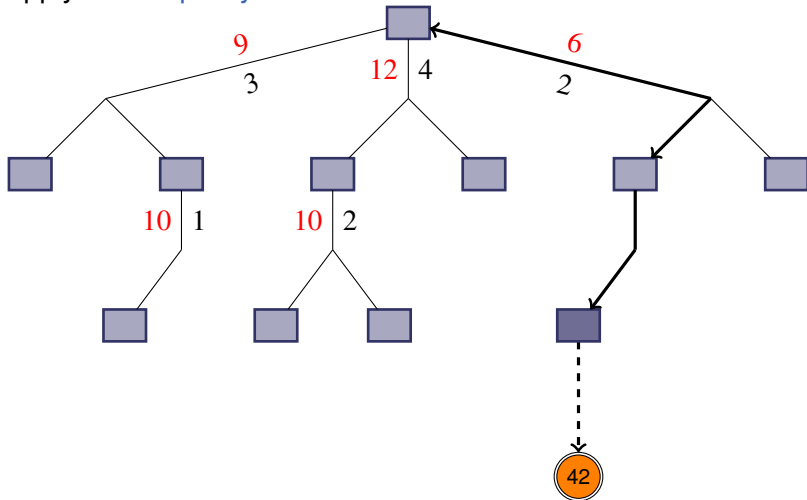
Expansion Phase

create a node for **first state** not in tree



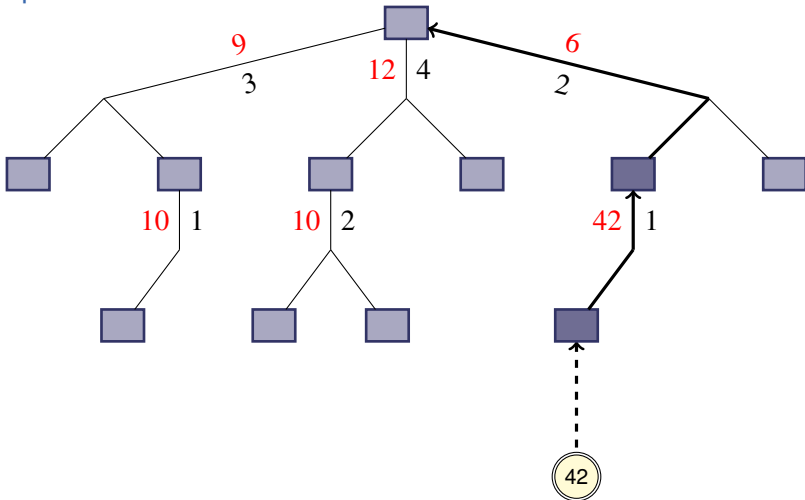
Simulation Phase

apply **default policy** until terminal state is reached



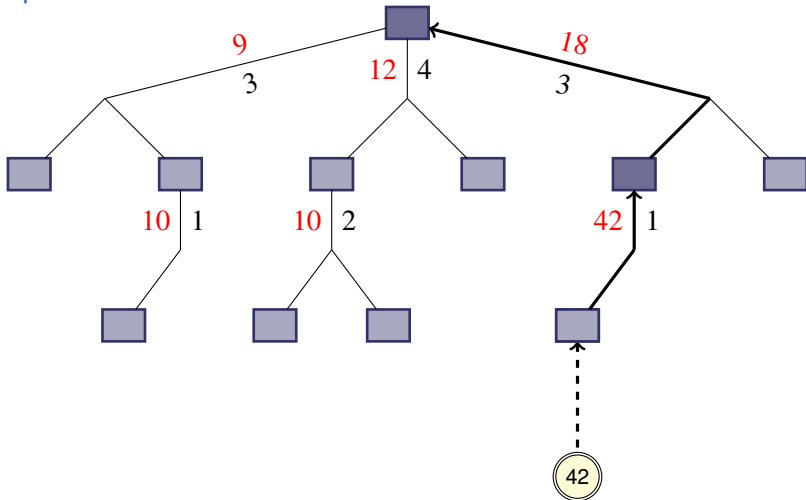
Backpropagation Phase

update visited nodes



Backpropagation Phase

update visited nodes



MCTS Configurations

- MCTS is a **framework** of algorithms
 - concrete algorithms are specified in terms of:
 - tree policy
 - default policy
 - for many probabilistic planning problems, there is a **well-suited MCTS configuration**
 - **However:** there are also many poorly suited MCTS configurations
 - and for every MCTS configuration that works well in one problem, there is **another problem** where it **performs poorly**
- ⇒ Are there properties we can analyse theoretically?

Asymptotic Optimality

Asymptotic Optimality

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

Asymptotic Optimality

Asymptotic Optimality

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

Note: there are MCTS instantiations that act optimally although the values do not converge in this way (e.g., if all $Q_k^{MCTS}(n)$ converge to $\ell \cdot Q^*(n.state)$ for a constant $\ell > 0$)

Asymptotic Optimality

A tree policy is **asymptotically optimal** if

- it **explores forever**:

- every position is **expanded eventually** and **visited infinitely often**

- and it is **greedy in the limit**:

- the probability that an optimal move action selected converges to 1

⇒ in the limit, backups based on iterations where only an **optimal policy** is followed dominate suboptimal backups

Tree Policy: Objective

tree policies have two contradictory objectives:

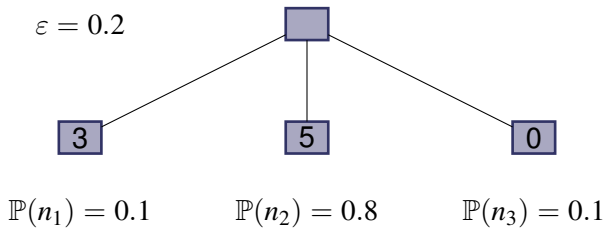
- **explore** parts of the game tree that have not been investigated thoroughly
- **exploit** knowledge about good moves to focus search on promising areas

central challenge: **balance** exploration and exploitation

ϵ -greedy: Idea

- tree policy with constant parameter ϵ
- with probability $1 - \epsilon$, pick a **greedy** move (i.e., one that leads to a successor node with the best action-value estimate)
- otherwise, pick a non-greedy successor **uniformly at random**

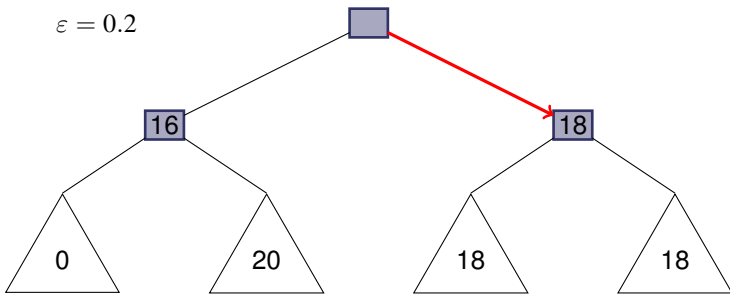
ϵ -greedy: Example



ϵ -greedy: Asymptotic Optimality

Asymptotic Optimality of ϵ -greedy

- explores forever
- not greedy in the limit
- ↪ not asymptotically optimal



ϵ -greedy: Asymptotic Optimality

Asymptotic Optimality of ϵ -greedy

- explores forever
- not greedy in the limit

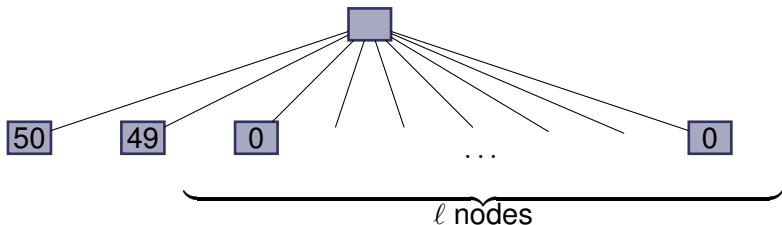
⇒ **not asymptotically optimal**

asymptotically optimal variant: use **decaying** ϵ , e.g. $\epsilon = \frac{1}{k}$

ϵ -greedy: Weakness

Problem:

when ϵ -greedy explores, all non-greedy moves are treated equally



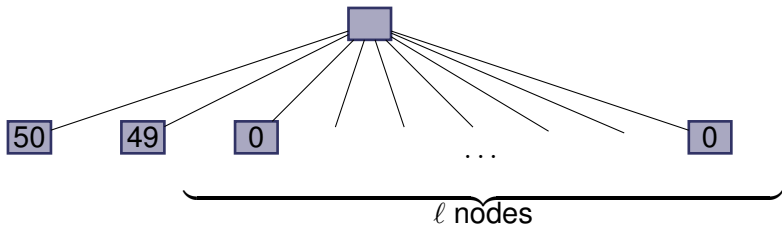
e.g., $\epsilon = 0.2, \ell = 9: \mathbb{P}(n_1) = 0.8, \mathbb{P}(n_2) = 0.02$

Softmax: Idea

- tree policy with constant parameter τ
- select moves **proportionally** to their action-value estimate
- **Boltzmann exploration** selects moves proportionally to

$$\mathbb{P}(n) \propto e^{\frac{\hat{Q}_k^{\text{MCTS}}(n)}{\tau}}$$

Softmax: Example



e.g., $\tau = 10, \ell = 9: \mathbb{P}(n_1) \approx 0.51, \mathbb{P}(n_2) \approx 0.46$

Boltzmann Exploration: Asymptotic Optimality

Asymptotic Optimality of Boltzmann Exploration

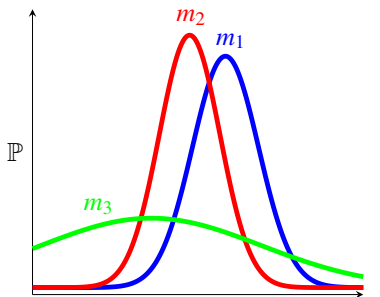
- explores forever
- not greedy in the limit
(probabilities converge to positive constant)

↪ not asymptotically optimal

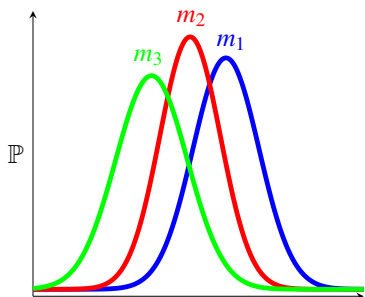
asymptotically optimal variants: use **decaying** τ

careful: τ must not decay faster than logarithmically
(i.e., must have $\tau \geq \frac{\text{const}}{\log k}$) to explore infinitely

Boltzmann Exploration: Weakness



scenario 1: high variance for m_3



scenario 2: low variance for m_3

- Boltzmann exploration only considers **mean** of sampled action-values for the given moves
- as we sample the same node many times, we can also gather information about variance (how **reliable** the information is)
- Boltzmann exploration ignores the variance, treating the two scenarios equally

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 $\hat{Q}_k^{MCTS}(n') + B(n')$
- based on **action-value estimate** $\hat{Q}^{MCTS}(n')$
- and a **bonus term** $B(n')$
- select $B(n')$ such that $Q^*(n'.state) \leq \hat{Q}_k^{MCTS}(n') + B(n')$ with high probability
- idea: $\hat{Q}_k^{MCTS}(n') + B(n')$ is an **upper confidence bound** on $Q^*(n')$ under the collected information

Upper Confidence Bounds: UCB1

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

Upper Confidence Bounds: Asymptotic Optimality

Asymptotic Optimality of UCB1

- explores forever
- greedy in the limit

~> asymptotically optimal

Upper Confidence Bounds: Asymptotic Optimality

Asymptotic Optimality of UCB1

- explores forever
- greedy in the limit
- ⇒ 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

UCT tree (equal number of search nodes)

