MDP Algorithms

Thomas Keller

University of Basel

June 20, 2018

MC Methods

Outline of this lecture

- Markov decision processes
- Planning via determinization
- Monte-Carlo methods
- Monte-Carlo Tree Search
- Heuristic Search
- Trial-based Heuristic Tree Search

Please ask questions at any time!

Introduction $\circ \bullet$

MDPs 000000000

Determinization

MC Methods

Motivation



MDPs 000000000

Determinization

MC Methods

Motivation



MDPs 000000000

Determinization

MC Methods

Motivation









MDPs •oooooooo Determinization

MC Methods

Markov decision process

Definition (Markov decision process)

A (finite-horizon) Markov decision process (MDP) is a 6-tuple $\mathcal{M} = \langle S, \mathcal{A}, \mathcal{T}, \mathcal{R}, s_0, H \rangle$, where

- \blacksquare S is a finite set of states
- A is a finite set of actions
- $\blacksquare \mathcal{T}: S \times A \times S \mapsto [0, 1]$ is the transition function
- $\blacksquare \ \mathcal{R}: S \times A \mapsto \mathbb{R} \text{ is the reward function}$
- $s_0 \in S$ is the initial state
- \blacksquare $H \in \mathbb{N}$ is the horizon

MDPs ooooooooo Determinization

MC Methods

Markov decision process: example



MDPs are acyclic in finite-horizon setting

Policy

Definition (Policy)

A partial policy for a Markov decision process $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, s_0, H \rangle$ is a mapping $\pi : \mathcal{S} \times \mathcal{A} \mapsto [0, 1] \cup \{\bot\}$. A partial policy π is executable in \mathcal{M} if $\pi(s, a) \neq \bot$ for all states *s* and actions *a* that can be reached by π from s_0 .

 $\pi(s, a)$ gives probability that action *a* is executed in state *s* under application of policy π

MDPs 000000000

Determinization

MC Methods

State- and Action Value

Definition (State- and Action Values)

The state-value $V^{\pi}(s)$ of a state $s \in S$ under policy π is

$$V^{\pi}(s):=egin{cases} 0 & ext{if }s ext{ is termina}\ Q^{\pi}(s,\pi(s)) & ext{otherwise,} \end{cases}$$

where $Q^{\pi}(s, a)$ is the action-value of s and action $a \in \mathcal{A}$ under π

$$Q^{\pi}(s,a) := R(s,a) + \sum_{s' \in \mathcal{S}} \left(T(s,a,s') \cdot V^{\pi}(s') \right)$$

for all state-action pairs (s, a).

	MDPs 0000●0000	Determinization ০০০০০০০০০	MC Methods
Optimal Policy			

- policy π is optimal if $V^{\pi}(s)$ and $Q^{\pi}(s, a)$ are maximal among all π (in the following, π^* , $V^*(s)$ and $Q^*(s, a)$)
- compute $V^{\star}(s)$ and $Q^{\star}(s, a)$ as

$$V^{\star}(s) := \begin{cases} 0 & \text{if } s \text{ is terminal} \\ \max_{a \in \mathcal{A}} Q^{\star}(s, a) & \text{otherwise} \end{cases}$$
$$Q^{\star}(s, a) := R(s, a) + \sum_{s' \in \mathcal{S}} \left(T(s, a, s') \cdot V^{\star}(s') \right).$$

start with terminal states and proceed backwards in DAG until initial state

MDPs

Determinization

MC Methods

Optimal solution



State- and action-values describe expected reward

MDPs 0000000000 MC Methods

Lecture goals

How can we act well despite the complexity of the problem?

- Which algorithms for probabilistic planning are there?
- Which are their strengths and weaknesses?
- What do they have in common, and what are the differences?

MC Methods

Anytime Planning: Plan-Execute-Monitor Cycle

- size of executable policy exponential in horizon
- compact representation of executable policy required to describe solution ⇒ not part of this lecture
- instead: perform plan-execute-monitor cycle:
 - plan action a for the current state s₀
 - execute a
 - update s_0 and H in \mathcal{M}
 - repeat until H = 0

MC Methods

Anytime Planning: Plan-Execute-Monitor Cycle

Advantages and disadvantages of plan-execute-monitor:

- avoid loss of precision that often comes with compact description of executable policy
- do not waste time with planning for states that are never reached during execution
- poor decisions can be avoided by spending more time with deliberation before execution



- replace probabilistic actions with deterministic ones
- leads to classical planning problem
- (often) determinization can be solved in practice even if MDP cannot

How do we come up with a determinization?

MDPs 000000000

Determinization

MC Methods

Determinization: Example



MDPs 000000000

Determinization

MC Methods

Determinization: Single-outcome determinization

Remove all outcomes but one



MDPs 000000000 Determinization

MC Methods

All-outcomes determinization: Example



MDPs 000000000

Determinization

MC Methods

Determinization: All-outcomes determinization

Generate one action per outcome



MDPs 000000000

Determinization

MC Methods

Single-outcome Determinization: Limitations



MDPs 000000000

Determinization

MC Methods

Single-outcome Determinization: Limitations

Important parts of the MDP can become unreachable



MDPs 000000000 Determinization

MC Methods

All-outcomes Determinization: Limitations

All-outcomes determinization is too optimistic



MC Methods

Determinizations in Probabilistic Planning

in combination with classical planner in plan-execute-monitor cycle approach

- well-suited if uncertainty has certain form (e.g., actions can fail or succeed)
- well-suited if information on probabilities noisy (e.g., path planning for robots in uncertain terrain)

domain-independent implementation: FF-Replan (Yoon, Fern & Givan)

MC Methods

Determinizations in Probabilistic Planning

as subsolver of a more complex system:

- FPG (Buffet and Aberdeen) learns a policy utilizing FF-Replan
- RFF (Teichteil-Königsbuch, Infantes & Kuter) extends determinization-based plan to policy

PROST-2011 (Keller & Eyerich) and PROST-2014 (Keller & Geißer) use a determinization-based iterative deepening search as heuristic

MC Methods •••••••••

Monte-Carlo Tree Search: Brief History

- Starting in the 1930s: first researchers experiment with Monte-Carlo methods
- 1998: Ginsberg's GIB player, based on Hindsight Optimization, achieves strong performance playing Bridge
- 2002: Kearns et al. propose Sparse Sampling
- 2002: Auer et al. present UCB1 action selection for multi-armed bandits
- 2006: Coulom coins the term Monte-Carlo Tree Search (MCTS)
- 2006: Kocsis and Szepesvári combine UCB1 and MCTS into the most famous MCTS variant, UCT
- 2007-2016: Constant progress in MCTS-based Go player lead to AlphaGo's defeat of 9-dan Go player Lee Sedol

MDPs 000000000 Determinization

MC Methods

Monte-Carlo Methods: Idea

- subsume a broad family of algorithms
- decisions are based on random samples
- results of samples are aggregated by computing the average
- apart from these points, algorithms differ significantly

Hindsight Optimization: Idea

- perform samples as long as resources (deliberation time, memory) allow:
 - sample a determinization of the MDP
 - solve the determinization for each applicable action
 - update average reward $\hat{Q}^{HOP}(s_0, a)$ for each action *a* with action-value estimate of *a* in the sample
- execute the action with the highest average reward

MDPs 000000000 Determinization

MC Methods

Hindsight Optimization: Example



South to play, three tricks to win, trump suit ♣

MDPs 00000000000 Determinization

MC Methods

Hindsight Optimization: Example



South to play, three tricks to win, trump suit ♣

MDPs 000000000 Determinization

MC Methods

Hindsight Optimization: Example

Å







South to play, three tricks to win, trump suit \$

MDPs 00000000000 Determinization

MC Methods

Hindsight Optimization: Example



South to play, three tricks to win, trump suit ♣

Determinizatior

MC Methods

Hindsight Optimization: Evaluation

- HOP well-suited for some problems
- must be possible to solve sampled MDP efficiently:
 - domain-dependent knowledge (various games and MDPs)
 - classical planner (FF-Hindsight)
 - LP solver (Issakkimuthu et al., ICAPS 2015)
- What about optimality with unbounded resources?

MDPs 000000000 Determinization

MC Methods

Hindsight Optimization: Optimality



MDPs 000000000 Determinization

MC Methods

Hindsight Optimization: Optimality



MDPs 000000000 Determinization

MC Methods

Hindsight Optimization: Optimality



Hindsight Optimization: Evaluation

- HOP well-suited for some problems,
- must be possible to solve sampled MDP efficiently:
 - domain-dependent knowledge (various games and MDPs)
 - classical planner (FF-Hindsight, Yoon et. al, AAAI 2008)
 - LP solver (Issakkimuthu et al., ICAPS 2015)
- What about optimality in the limit?
 in many problems not optimal due to assumption of clairvoyance

MC Methods

Hindsight Optimization: Clairvoyance

Idea of Hindsight Optimization (Repetition):

- perform samples as long as resources (deliberation time, memory) allow:
 - sample a determinization of the MDP
 - solve the determinization for each applicable action
 - update average reward $\hat{Q}^{HOP}(s_0, a)$ for each action *a* with action-value estimate of *a* in the sample
- execute the action with the highest average reward

MC Methods

Hindsight Optimization: Clairvoyance

Idea of Hindsight Optimization (Repetition):

- perform samples as long as resources (deliberation time, memory) allow:
 - sample a determinization of the MDP
 - solve the determinization for each applicable action
 - update average reward $\hat{Q}^{HOP}(s_0, a)$ for each action *a* with action-value estimate of *a* in the sample
- execute the action with the highest average reward

MDPs 000000000 Determinization

MC Methods

Policy Simulation: Idea

- separate computation of policy and its evaluation
- by simulating the execution of a policy
- any base policy can be used

MDPs 000000000 Determinization

MC Methods

Optimistic Policy Simulation

- perform samples as long as resources (deliberation time, memory) allow:
 - sample a determinization of the MDP
 - compute a policy by solving the determinization for each applicable action
 - simulate the policy
 - update average reward $\hat{Q}^{OPS}(s_0, a)$ for each action *a* with reward in the simulation
- execute the action with the highest average reward

MDPs 000000000 Determinization

MC Methods

Optimistic Policy Simulation: Example



 $s_1 \rightarrow s_3 \text{ in sample} \qquad s_1 \rightarrow s_3 \text{ simulation} \\ s_1 \rightarrow s_4 \text{ in sample} \qquad s_1 \rightarrow s_4 \text{ in in simulation}$

MDPs 000000000

Determinization

MC Methods

Optimistic Policy Simulation: Example



 $s_1 \rightarrow s_3 \text{ in sample} \qquad s_1 \rightarrow s_3 \text{ simulation} \\ s_1 \rightarrow s_4 \text{ in sample} \qquad s_1 \rightarrow s_4 \text{ in in simulation}$

MC Methods

Optimistic Policy Simulation: Evaluation

- Problem: suboptimal base policy is static
- no mechansim to overcome weakness of base policy
- $\blacksquare \Rightarrow$ repeated suboptimal decisions in simulation affect Policy Simulation

Sparse Sampling: Idea

- rather than restrict the simulated actions (as in policy simulation), restrict the simulated outcomes
- search tree creation: sample a constant number of outcomes according to their probability in each state and ignore the rest
- near-optimal: utility of resulting policy close to utility of optimal policy
- runtime independent from the number of states

MDPs 000000000 Determinization

MC Methods

Sparse Sampling: Search Tree



MDPs 000000000 Determinization

MC Methods

Sparse Sampling: Search Tree



MDPs 000000000 Determinization

MC Methods

Sparse Sampling: Problems

- independent from number of states, but still exponential in horizon
- constant that gives the number of outcomes large for good bounds on near-optimality
- search time difficult to predict
- tree is symmetric ⇒ resources are wasted in non-promising parts of the tree

Summary

presented several algorithms for probabilistic planning:

- Determinizations
- Hindsight Optimization
- Policy Sampling
- Sparse Sampling
- there are applications for all where they perform well
- but all are suboptimal even if provided with unlimited resources