Short-Sighted Stochastic Shortest Path Problems

Felipe Trevizan and Manuela Veloso

School of Computer Science Carnegie Mellon University June 29, 2012



Announcement

 Wrong version of the paper in the proceedings in the *thumbdrive* and the cd-rom:



RTDP (Bonet and Geffner 2003), resulting in optimal algorithms with convergence bounds. Due to the pruning in

• The right version is in the online proceedings

model in which solutions computed based on it can be ex-

ecuted for at least t steps as a closed form solution. Using

short-sighted SSPs, we present a novel probabilistic planner

Motivation

- Two classes of solutions to probabilistic planning problems:
 - Complete policy (a.k.a. universal plan):
 - Maps every state to an action
 - Never fails, i.e., no need to replan
 - Optimal
 - Doesn't scale up

Motivation

- Two classes of solutions to probabilistic planning problems:
 - Complete policy (a.k.a. universal plan):
 - Maps every state to an action
 - Never fails, i.e., no need to replan
 - Optimal
 - Doesn't scale up
 - Partial policy:
 - Maps some states to an action
 - Can fail, i.e., reaches an unpredicted state and replan from there
 - Non-optimal
 - Scales up

Contributions

- A framework that offers a **new trade-off** between complete and partial policies:
 - A new model: short-sighted Stochastic Shortest
 Path Problems
 - A new **planner**: short-sighted probabilistic planner



Model problems as Stochastic Shortest Path Problems



- Model problems as Stochastic Shortest Path Problems
- Generate short-sighted subproblems



- Model problems as Stochastic Shortest Path Problems
- Generate short-sighted subproblems
- Solve the subproblems and execute this solution



- Model problems as Stochastic Shortest Path Problems
- Generate short-sighted subproblems
- Solve the subproblems and execute this solution



- Model problems as Stochastic Shortest Path Problems
- Generate short-sighted subproblems
- Solve the subproblems and execute this solution



- Model problems as Stochastic Shortest Path Problems
- Generate short-sighted subproblems
- Solve the subproblems and execute this solution



• Analyze single and multiple execution cases

- Model problems as Stochastic Shortest Path Problems
- Generate short-sighted subproblems
- Solve the subproblems and execute this solution



Analyze single and multiple execution cases

Stochastic Shortest Path Problems (SSPs)



Optimal policies

- An **optimal** policy π^* minimizes the expected cost to reach a goal state from s₀
- The **minimum** expected cost to reach a goal state from a state s is:

$$V^*(s) = \begin{cases} 0 & \text{if } s \in \mathsf{G} \\ \min_{a \in \mathsf{A}} \sum_{s' \in \mathsf{S}} P(s'|s, a) [C(s, a, s') + V^*(s')] & \text{otherwise} \end{cases}$$

Short-Sighted SSPs: Idea

- Manage uncertainty by:
 - Considering the uncertainty structure in the neighborhood of the current state; and
 - Adding artificial goals to heuristically approximate the pruned states.



Short-Sighted SSPs: Idea

- Manage uncertainty by:
 - Considering the uncertainty structure in the neighborhood of the current state; and
 - Adding artificial goals to heuristically approximate the pruned states.



• $\delta(s, s')$: minimum number of actions to reach s' from s

- Given: an SSP <S,s₀,G,A,P,C>,
 - $s \in \mathsf{S}$
 - t > 0 and
 - a heuristic function H

the (s,t)-short-sighted SSP is <S',s,G',A,P,C'>:

- Given: an SSP <S,s₀,G,A,P,C>,
 - $s \in \mathsf{S}$
 - t > 0 and
 - a heuristic function H

the (s,t)-short-sighted SSP is <S',s,G',A,P,C'>:

States reachable using **up to t** actions

•
$$S' = \{s' \in S | \delta(s, s') \leq t\}$$

- Given: an SSP <S,s₀,G,A,P,C>,
 - $s \in \mathsf{S}$
 - t > 0 and
 - a heuristic function H

the (s,t)-short-sighted SSP is <S',s,G',A,P,C'>:

States reachable using **up to t** actions

• $\mathsf{S}' = \{s' \in \mathsf{S} | \delta(s, s') \leq t\}$

Artificial goal: states reachable using **exactly t** actions

• $\mathsf{G}' = \{s' \in \mathsf{S} | \delta(s, s') = t\} \cup (\mathsf{G} \cap \mathsf{S}')$

- Given: an SSP <S,s₀,G,A,P,C>,
 - $s \in S$
 - t > 0 and
 - a heuristic function H

the (s,t)-short-sighted SSP is <S',s,G',A,P,C'>:

•
$$S' = \{s' \in S | \delta(s, s') \leq t\}$$

• $G' = \{s' \in S | \delta(s, s') = t\} \cup (G \cap S')$
• $C'(s, a, s') = \begin{cases} C(s, a, s') + H(s') & \text{if } s' \in G' \\ C(s, a, s') & \text{otherwise} \end{cases}$

If s' is an artificial goal, then its cost is incremented by its heuristic value

Short-Sighted SSPs: Examples



Short-Sighted SSPs: Examples





Short-Sighted SSPs: Examples



(s₀,1)-short-sighted:



 $(s_0, 2)$ -short-sighted:



22







Key difference: Short-sighted SSPs preserve the **action structure**, e.g., self-loop actions and loops of actions



Key difference: Short-sighted SSPs preserve the **action structure**, e.g., self-loop actions and loops of actions

Theorem: the optimal value-function for an (s,t)-short-sighted SSP is at least as good as the t-look-ahead value of s, i.e.,

$$L_t(s_0) \le \hat{V}_t^*(s_0) \le V^*(s_0)$$

Short-Sighted Probabilistic Planner (SSiPP)



Since short-sighted SSPs are much smaller than the original problem. we can compute a complete policy for them.

SSiPP and replanning

• **Theorem**: at least t actions are executed in the environment before replanning is needed



• The policy $\frac{s_0}{a_0} = \frac{s_3}{s_3}$ can be executed for **more** than 2

timesteps:

$$s_0 \xrightarrow{a_0} s_3 \xrightarrow{a_0} s_3 \xrightarrow{a_0} s_3 \xrightarrow{a_0} \cdots$$

- Two scenarios:
 - the same problem is solved more than once; or
 - simulation is allowed for a given amount of time

- Two scenarios:
 - the same problem is solved more than once; or
 simulation is allowed for a given amount of time
- **Theorem**: SSiPP is **asymptotically optimal**, i.e., if the same problem is solved sufficiently many times, then the optimal policy is found.

- Two scenarios:
 - the same problem is solved more than once; or
 simulation is allowed for a given amount of time
- Theorem: SSiPP is asymptotically optimal, i.e., if the same problem is solved sufficiently many times, then the optimal policy is found.



- Two scenarios:
 - the same problem is solved more than once; or
 simulation is allowed for a given amount of time
- Theorem: SSiPP is asymptotically optimal, i.e., if the same problem is solved sufficiently many times, then the optimal policy is found.





- Two scenarios:
 - the same problem is solved more than once; or
 simulation is allowed for a given amount of time
- Theorem: SSiPP is asymptotically optimal, i.e., if the same problem is solved sufficiently many times, then the optimal policy is found.





- Two scenarios:
 - the same problem is solved more than once; or
 simulation is allowed for a given amount of time
- Theorem: SSiPP is asymptotically optimal, i.e., if the same problem is solved sufficiently many times, then the optimal policy is found.



2nd run of the same problem



- Two scenarios:
 - the same problem is solved more than once; or
 simulation is allowed for a given amount of time
- Theorem: SSiPP is asymptotically optimal, i.e., if the same problem is solved sufficiently many times, then the optimal policy is found.



2nd run of the same problem



- RTDP and *anytime* SSiPP:
 - can be seen as asynchronous value iteration
 - differ in the scheduling of Bellman updates (backups)







- RTDP and *anytime* SSiPP:
 - can be seen as asynchronous value iteration
 - differ in the scheduling of Bellman updates (backups)







- RTDP and *anytime* SSiPP:
 - can be seen as asynchronous value iteration
 - differ in the scheduling of Bellman updates (backups)







- RTDP and *anytime* SSiPP:
 - can be seen as asynchronous value iteration
 - differ in the scheduling of Bellman updates (backups)







- RTDP and *anytime* SSiPP:
 - can be seen as asynchronous value iteration
 - differ in the scheduling of Bellman updates (backups)



RTDP



- RTDP and *anytime* SSiPP:
 - can be seen as asynchronous value iteration
 - differ in the scheduling of Bellman updates (backups)



RTDP



- RTDP and *anytime* SSiPP:
 - can be seen as asynchronous value iteration
 - differ in the scheduling of Bellman updates (backups)



RTDP



- RTDP and *anytime* SSiPP:
 - can be seen as asynchronous value iteration
 - differ in the scheduling of Bellman updates (backups)



RTDP



- RTDP and *anytime* SSiPP:
 - can be seen as asynchronous value iteration
 - differ in the scheduling of Bellman updates (backups)



RTDP



- RTDP and *anytime* SSiPP:
 - can be seen as asynchronous value iteration
 - differ in the scheduling of Bellman updates (backups)



RTDP



- RTDP and *anytime* SSiPP:
 - can be seen as asynchronous value iteration
 - differ in the scheduling of Bellman updates (backups)



RTDP



Experiments

- Goal: compare SSiPP against the winners of the previous International Probabilistic Planning Competitions (IPPCs)
- Methodology: (same as IPPC'04 and IPPC'06)
 - For each problem, planners are requested to solve it
 50 times in 20 minutes
 - Learning is allowed between attempts of the same problem
 - The evaluation metric is the number of times the goal is reached

Planners

- We compare the following planners
 - **SSiPP**: using LRTDP as optimal solver
 - LRTDP: 2nd place IPPC'04
 - FF-Replan: 1st place IPPC'04
 - FPG: 1st place IPPC'06
 - **RFF**: 1st place IPPC'08
- Parametrizations for SSiPP and LRTDP:
 - $-t \in \{1, 2, 3, ..., 10\}$
 - H: zero-heuristic, FF+all-outcomes, min-min

Triangle tireworld: results



SSiPP: t = 3, H = FF+all-outcome

SSiPP: t = 8, H = zero-heuristic

LRTDP: t = 3, H = zero-heuristic





Blocks World: result



SSiPP: t = 3, H = FF+all-outcome

SSiPP: t = 2, H = FF+all-outcome

LRTDP: t = 3, H = zero-heuristic

FF-Replan



RFF

Exploding Blocks World: results



SSiPP: t = 3, H = FF+all-outcome

SSiPP: t = 3, H = FF+all-outcome

LRTDP: t = 3, H = zero-heuristic

FF-Replan



RFF

Zeno Travel: results



SSiPP: t = 3, H = FF+all-outcome

SSiPP: t = 2, H = FF+all-outcome

LRTDP: t = 3, H = zero-heuristic





IPPC experiment: summary

| | SSiPP Overall | SSiPP Per Domain |
|--------------------|---------------|------------------|
| Outperforms all | 16.6% | 36.6% |
| Ties with the best | 41.6% | 53.3% |

IPPC experiment: summary

| | SSiPP Overall | SSiPP Per Domain |
|--------------------|---------------|------------------|
| Outperforms all | 16.6% | 36.6% |
| Ties with the best | 41.6% | 53.3% |

- In the considered problems:
 - SSiPP is never the last place in any of the problems
 - LRTDP never outperforms SSiPP

Conclusion and Future Work

- A framework that offers a new trade-off between complete and partial policies
 - Short-sighted SSPs

Conclusion and Future Work

- A framework that offers a new trade-off between complete and partial policies
 - Short-sighted SSPs
 - SSiPP
 - as replanner: no replanning needed for at least t actions
 - as anytime algorithm: same guarantees as RTDP

Conclusion and Future Work

- A framework that offers a new trade-off between complete and partial policies
 - Short-sighted SSPs
 - SSiPP
 - as replanner: no replanning needed for at least t actions
 - as anytime algorithm: same guarantees as RTDP
- Future work:
 - New definitions of short-sighted *spaces* (S')
 - Improve scalability in problems that are **not** probabilistic interesting

Thank you!

Questions?