

On the Hardness of the MDP Safety Problem

Ashwin Abraham

Indian Institute of Technology Bombay, India
ashwinabraham@cse.iitb.ac.in

Abstract. We show that the initialized safety problem for MDPs is co-NP hard when the safe set is affine. We then describe an algorithm to solve the initialized safety problem in MDPs, where the safe set is a finite union of affine sets. We show that this version of the problem is Σ_2^P -hard.

1 Introduction

The terminology we use to describe MDPs has been introduced in [1].

1.1 The Initialized Safety Problem

Given an MDP $\mathcal{M} = (S, Act, \delta)$ and a safe set $\mathcal{H} \subseteq \Delta(S)$, a distribution $\mu_0 \in \Delta(S)$ is said to be \mathcal{H} -safe iff there exists a policy over π such that $\forall i \geq 0, \mathcal{M}^\pi(\mu_0, i) \in \mathcal{H}$. Such a policy π is said to be a safe policy.

Definition 1 (Distributional Policies). Any policy π is said to be distributional if it is a function from $\Delta(S)$ to Π_1 , where Π_1 is the set of one step strategies. These distributional strategies induce a sequence of distributions μ_1 given by $\mu_{i+1} = \mathcal{M}^{\pi(\mu_i)}(\mu_i)$.

Theorem 1. If an initial distribution μ_0 is \mathcal{H} -safe, then it is \mathcal{H} -safe under some distributional policy $\pi : \Delta(S) \rightarrow \Pi_1$

Proof. Can be found in [1]. □

1.2 Inductive Invariants

Definition 2 (Inductive Invariants). A set $\mathcal{I} \subseteq \mathcal{H}$ for an initialized safety problem is said to be an inductive invariant for μ_0 under policy π iff

1. $\mu_0 \in \mathcal{I}$
2. $\forall \mu, \mu \in \mathcal{I} \implies \mathcal{M}^\pi(\mu) \in \mathcal{I}$

Theorem 2. Inductive Invariants are sound and complete certificates for safety.

Proof. Inductive Invariants are sound, as if an inductive invariant exists, we can show by induction on i , that $\forall i, \mathcal{M}^\pi(\mu_0, i) \in \mathcal{I} \subseteq \mathcal{H}$, which means that the MDP is clearly safe.

Now, if the MDP is safe then the set $\{\mathcal{M}^\pi(\mu_0, i) : i \in \mathbb{N}\}$ is clearly an inductive invariant, and hence inductive invariants are complete as well. □

Theorem 3 (Affine Inductive Invariants). *If the set μ_0 is safe for an affine safe set \mathcal{H} under a memoryless policy π , and if the set of obtained distributions $\{\mathcal{M}^\pi(\mu_0, i) : i \in \mathbb{N}\}$ is finite, then there exists an affine inductive invariant.*

Proof. Consider the convex hull of the set $\{\mathcal{M}^\pi(\mu_0, i) : i \in \mathbb{N}\}$

$$\mathcal{I} = \left\{ \sum_{i \in S} \alpha_i \mathcal{M}^\pi(\mu_0, i) : S \subseteq \mathbb{N} \wedge |S| < \infty \wedge \sum_{i \in S} \alpha_i = 1 \right\}$$

Clearly we have $\mu_0 \in \mathcal{I}$ and since we have $\{\mathcal{M}^\pi(\mu_0, i) : i \in \mathbb{N}\} \subseteq \mathcal{H}$ and since \mathcal{H} is convex, $\mathcal{I} \subseteq \mathcal{H}$.

Now, since \mathcal{M}^π is a linear transformation, we have:

$$\mathcal{M}^\pi \left(\sum_{i \in S} \alpha_i \mathcal{M}^\pi(\mu_0, i) \right) = \sum_{i \in S} \alpha_i \mathcal{M}^\pi(\mu_0, i+1) \in \mathcal{I}$$

From which we can conclude that \mathcal{I} is an inductive invariant.

Now, since $\{\mathcal{M}^\pi(\mu_0, i) : i \in \mathbb{N}\}$ is finite, its convex hull, \mathcal{I} , is affine.

Therefore, \mathcal{I} is an affine inductive invariant. \square

1.3 Affine Inductive Invariant Synthesis

An algorithm was presented in [1] to synthesize affine inductive invariants for memoryless policies. The algorithm proceeds with the following steps:

1. *Setting up Templates.* We set up templates for our (memoryless) policy π and our inductive invariant \mathcal{I} . Our template for our policy is $\pi(a|s) = p_{sa}$. This adds constraints $0 \leq p_{sa} \leq 1$ and $\sum_{a \in \text{Act}(s)} p_{sa} = 1$. Our template for \mathcal{I} is $\mathcal{I}(x) = (x \in \Delta(S)) \wedge \bigwedge_{i=0}^N (b_i + \sum_{j=0}^n a_{ij} x_j \geq 0)$. Our affine inductive invariant is $\{x : \mathcal{I}(x)\}$, and the first clause ensures that $\mathcal{I} \subseteq \Delta(S)$.
2. *Constraint Collection.* We then add the constraints $\mathcal{I}(\mu_0)$ (ie $\mu_0 \in \mathcal{I}$) and the inductive condition $\forall x, \mathcal{I}(x) \implies \mathcal{I}(\text{next}(x))$, where $\text{next}(x)$ is the distribution obtained after one step from initial distribution x and with policy π . The last constraint to be collected is $\forall x, \mathcal{I}(x) \implies \mathcal{H}(x)$.
3. *Quantifier Elimination.* The current set of constraints contains quantifier alternation, which makes it hard to solve. We remove the quantifier alternation by noticing that we can remove the universal quantifiers in the following way. The universal quantifiers appear in the constraints $\forall x, \mathcal{I}(x) \implies \mathcal{I}(\text{next}(x))$ and $\forall x, \mathcal{I}(x) \implies \mathcal{H}(x)$. Since \mathcal{I} and \mathcal{H} are both affine sets, ie they are formed by the conjunction of a set of affine (non-strict) inequalities, we can write these constraints as $\bigwedge_{C \in \mathcal{I}_{\text{next}}} (\forall x, \mathcal{I}(x) \implies C)$ and $\bigwedge_{C \in \mathcal{H}} (\forall x, \mathcal{I}(x) \implies C)$. Now, we can apply *Farkas' Lemma* on these equations to remove the universal quantification over x and get a constraint on the coefficients, and the new variables introduced. Note that these constraints will be polynomials in the coefficients and template variables.

4. *Constraint Solving.* The constraints finally obtained are purely existentially qualified, and therefore can be solved within the existential theory of the reals, $\exists\mathbb{R}$.

This algorithm is parameterized by N , the size of the encoding of \mathcal{I} (the number of faces of \mathcal{I}). If the algorithm fails for one particular value of N , we move on to a larger value. Note that this means that the algorithm will not terminate if there exists no safe memoryless policy with an affine inductive invariant witnessing it.

Theorem 4 (Soundness and Completeness). *If this algorithm returns a memoryless policy π and an affine inductive invariant \mathcal{I} , then π is \mathcal{H} -safe and \mathcal{I} is an affine inductive invariant witnessing the safety (Soundness).*

If there exists a safe memoryless policy π and an affine inductive invariant witnessing the safety of π , then this algorithm will eventually terminate on some safe memoryless policy and affine inductive invariant. There exists a minimum value $N^ \in \mathbb{N}$ which is the smallest value of N required by the algorithm to compute a safe memoryless policy and invariant (Completeness).*

Proof. Soundness follows directly from the soundness of Farkas' Lemma. For completeness, note that if there exists a safe memoryless policy π and an affine inductive invariant witnessing it, then there exists a safe memoryless policy π and affine inductive invariant with the minimum number of faces. Let this number be N^* . For $N = N^*$, due to the completeness of Farkas' Lemma, our algorithm will return the safe memoryless policy and affine inductive invariant, and for $N < N^*$, no safe memoryless policy and affine inductive invariant will be found, due to the minimality of N^* and the soundness of the algorithm. \square

Since the size of the constraints generated by our algorithm is polynomial in the size of the MDP, the encoding of \mathcal{H} and N , this algorithm has a runtime in $PSPACE$ in terms of the size of the MDP, the encoding of \mathcal{H} and N^* , for the MDPs where a safe memoryless policy and affine inductive invariant exist. Note that since this algorithm only terminates on instances where a safe memoryless policy and affine inductive invariant exist, this complexity bound may not hold for the decision problem of checking if an affine inductive invariant exists - we cannot even conclude from this that it is decidable.

2 co-NP Hardness of the initialized safety problem

We prove that the initialized safety problem is co-NP hard, even for Markov Chains and even when the safe set \mathcal{H} is restricted to an affine set.

We do this by reducing the co-NP complete problem of checking the validity of a 3-DNF to the initialized safety problem.

Consider a 3-DNF $\varphi = \bigvee_{i=0}^{k-1} C_i$ where each C_i is a conjunction of at most 3 literals. Let the variables of this DNF be x_1, \dots, x_n and let p_i denote the i^{th}

prime. Without loss of generality, we assume that no clause C_i is a contradiction, ie no clause contains a variable and its negation.

For each clause C , let $N(C)$ denote $\prod_{x_i \in C} p_i \cdot \prod_{\neg x_i \in C} p_i$, ie $N(C)$ denotes the products of all the primes corresponding to the variables in C - note that this product has at most 3 terms.

Define a Markov Chain \mathcal{M} whose set of states $S = \{(i, j) : 0 \leq i < k, 0 \leq j < N(C_i)\}$ and whose transition function is given by $\delta((j, p)|(i, q)) = 1$ if $i = j$ and $p \equiv q + 1 \pmod{N(C_i)}$ and 0 otherwise. We define an initial distribution over the states as μ_0 such that $\mu_0(i, 0) = \frac{1}{k}$ and $\mu_0(i, j) = 0$ for $j \neq 0$.

For any clause C , we define the set $f(C) = \bigcap_{x_i \in C} \{x : 0 \leq x < N(C) \wedge x \equiv 0 \pmod{p_i}\} \cap \bigcap_{\neg x_i \in C} \{x : 0 \leq x < N(C) \wedge x \not\equiv 0 \pmod{p_i}\}$. By the Chinese Remainder Theorem, $f(C)$ is non-empty, for each clause C . We define the set $T \subseteq S$ as the set $\{(i, j) : 0 \leq i < k \wedge j \in f(C_i)\}$.

We define the safe-set $\mathcal{H} = \{\mu \in \Delta(S) : \sum_{s \in T} \mu(s) \geq \frac{1}{k}\}$

Theorem 5 (The Reduction). μ_0 is \mathcal{H} -safe if and only if φ is valid.

Proof. At timestep t , the probability distribution on this Markov Chain, μ_t , will satisfy $\mu_t(i, j) = \frac{1}{k} \mathbb{I}(t \equiv j \pmod{N(C_i)})$. Note that this Markov Chain is periodic with period $\prod_{i=1}^n p_i$. Now, the safety condition $\forall t \geq 0, \sum_{s \in T} \mu_t(s) \geq \frac{1}{k}$ becomes $\forall t \geq 0, \sum_{i=0}^{k-1} \sum_{j \in f(C_i)} \mathbb{I}(t \equiv j \pmod{N(C_i)}) \geq 1$, which is equivalent to

$$\bigvee_{i=0}^{k-1} \bigvee_{j \in f(C_i)} (t \equiv j \pmod{N(C_i)}), \forall t \geq 0$$

Now, $\bigvee_{j \in f(C)} (t \equiv j \pmod{N(C_i)})$ is equivalent to

$$\bigwedge_{x_i \in C} (t \equiv 0 \pmod{p_i}) \wedge \bigwedge_{\neg x_i \in C} (t \not\equiv 0 \pmod{p_i})$$

and so our safety condition becomes

$$\bigvee_{i=0}^{k-1} \left[\bigwedge_{x_j \in C_i} (t \equiv 0 \pmod{p_j}) \wedge \bigwedge_{\neg x_j \in C_i} (t \not\equiv 0 \pmod{p_j}) \right], \forall t \geq 0$$

Define the Boolean variables $x_i(t) = (t \equiv 0 \pmod{p_i})$. In terms of these new variables, our safety condition becomes

$$\varphi(x_1(t), \dots, x_n(t)), \forall t \geq 0$$

Now, by the Chinese Remainder Theorem, for every choice of $0 \leq a_i < p_i$, the system of equations $t \equiv a_i \pmod{p_i}$ for $i = 1, \dots, n$ has a unique solution modulo

$\prod_{i=1}^n p_i$. Therefore, for each $\alpha_i \in \{0, 1\}$, the system of equations $x_i(t) = \alpha_i$ for $i = 1, \dots, n$ has a solution in $\{0, \dots, \prod_{i=1}^n p_i - 1\}$. Therefore, our safety condition is equivalent to

$$\forall x_1, \dots, x_n \varphi(x_1, \dots, x_n)$$

ie, \mathcal{M} is safe if and only if φ is valid. \square

Theorem 6 (co-NP hardness). *The initialized safety problem is co-NP hard.*

Proof. The sizes of the Markov Chain and safe set generated during our reduction are polynomial in the size of the DNF. To see this, note that there are no actions, and the number of states of the Markov Chain is given by $\sum_{i=0}^{k-1} N(C_i) \leq kp_n^3$. Since $p_n \in \Theta(n \log n)$, we have that the number of states is in $O(kn^3 \log^3(n))$. Therefore, our reduction is a polynomial many-one reduction from 3-VALIDITY (which is co-NP complete) to the initialized safety of Markov Chains, which therefore, must be co-NP hard. \square

Theorem 7. *If the constructed Markov Chain is safe, then there exists a polynomial sized affine inductive invariant certifying it's safety.*

Proof. The constructed Markov chain is periodic with period $\prod_{i=1}^n p_i$, and therefore, the set of obtained distributions $\{\mathcal{M}(\mu_0, i) : i \in \mathbb{N}\}$ is finite. Hence, by Theorem 3, there exists an affine inductive invariant certifying its safety. Note however, that this invariant may not be polynomial sized. \square

3 Σ_2^P Hardness of the deterministic initialized safety problem

We show that the problem of checking if there exists a deterministic safe policy for a given MDP is Σ_2^P hard. We do this by reducing Σ_2^P complete problem of finding the truth value of $\exists x_1 \dots x_m \forall y_1 \dots y_n \varphi(x_1, \dots, x_m, y_1, \dots, y_n)$ to the deterministic initialized safety problem.

References

1. Akshay, S., Chatterjee, K., Meggendorfer, T., Žikelić, Đ.: MDPs as Distribution Transformers: Affine Invariant Synthesis for Safety Objectives. Computer Aided Verification, 35th International Conference, CAV 2023, Paris, France