



HAL
open science

Reachability in MDPs: Refining Convergence of Value Iteration

Serge Haddad, Benjamin Monmege

► **To cite this version:**

Serge Haddad, Benjamin Monmege. Reachability in MDPs: Refining Convergence of Value Iteration. 8th International Workshop on Reachability Problems (RP'14), Sep 2014, Oxford, United Kingdom. pp.125-137, 10.1007/978-3-319-11439-2_10 . hal-01091122v2

HAL Id: hal-01091122

<https://hal.science/hal-01091122v2>

Submitted on 12 Feb 2016

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Copyright

Reachability in MDPs: Refining Convergence of Value Iteration^{*}

Serge Haddad¹ and Benjamin Monmege²

¹ LSV, ENS Cachan, CNRS & Inria, France

`serge.haddad@lsv.ens-cachan.fr`

² Université libre de Bruxelles, Belgium

`benjamin.monmege@ulb.ac.be`

Abstract. Markov Decision Processes (MDP) are a widely used model including both non-deterministic and probabilistic choices. Minimal and maximal probabilities to reach a target set of states, with respect to a policy resolving non-determinism, may be computed by several methods including value iteration. This algorithm, easy to implement and efficient in terms of space complexity, consists in iteratively finding the probabilities of paths of increasing length. However, it raises three issues: (1) defining a stopping criterion ensuring a bound on the approximation, (2) analyzing the rate of convergence, and (3) specifying an additional procedure to obtain the exact values once a sufficient number of iterations has been performed. The first two issues are still open and for the third one a “crude” upper bound on the number of iterations has been proposed. Based on a graph analysis and transformation of MDPs, we address these problems. First we introduce an *interval iteration algorithm*, for which the stopping criterion is straightforward. Then we exhibit convergence rate. Finally we significantly improve the bound on the number of iterations required to get the exact values.

1 Introduction

Markov Decision Processes (MDP) are a commonly used formalism for modelling systems that use both probabilistic and non-deterministic behaviors. These are generalizations of discrete-time Markov chains for which non-determinism is forbidden. MDPs have acquired an even greater gain of interest since the development of quantitative verification of systems, which in particular may take into account probabilistic aspects (see [1] for a deep study of model checking techniques, in particular for probabilistic systems). Automated verification techniques have been extensively studied to handle such probabilistic models, leading to various tools like the PRISM probabilistic model checker [9].

Value iteration for reachability problems. In the tutorial paper [5], the authors cover some of the algorithms for the model-checking of MDPs and Markov

^{*} The research leading to these results has received funding from the European Union Seventh Framework Programme (FP7/2007-2013) under Grant Agreement n601148 (CASSTING).

chains. The first simple, yet intriguing, problem lies in the computation of minimum and maximum probabilities to reach a target set of states of an MDP. Exact polynomial time methods, like linear programming, exist to compute those probabilities, but they seem unable to scale on large systems. Nonetheless, they are based on the fact that these probabilities are indeed fixpoints of some operators. Usually, numerical approximate methods are rather used in practice, the most used one being *value iteration*. The algorithm consists in asymptotically reaching the previous fixpoints by iterating the operators. However, it raises three issues. Since the algorithm must terminate after a finite number of iterations one has to define a stopping criterion ensuring a bound on the difference between the computed and the exact values. From a theoretical point of view, establishing the rate of convergence with respect to the parameters of the MDP (number of states, smallest positive transition probability, etc.) helps to estimate the complexity of value iteration. Sometimes for further application the exact values and/or the optimal policy are required. This is generally done by performing an additional rounding procedure once a sufficient number of iterations has been performed. The first two issues are still open and for the third one a “crude” upper bound on the number of iterations has been proposed [3, Sec 3.5].

Our contributions. Generally the numerical computations of optimal reachability probabilities are preceded by a qualitative analysis that computes the sets of states for which this probability is 0 or 1 and performs an appropriate transformation of the MDP. We adopt here an alternative approach based on the maximal end component (MEC) decomposition of an MDP (that can be computed in polynomial time [4]). We show that for an MDP featuring a particular MEC decomposition, some safety probability is null with an additional convergence rate with respect to the length of the run. Then we design a min- (respectively, max-) reduction that ensures this feature while preserving the minimal (respectively, maximal) reachability probabilities. In both cases, we establish that the reachability probabilities are unique fixed points of some operator.

So we iterate these operators starting from the maximal and the minimal possible vectors. These iterations naturally yield an *interval iteration algorithm* for which the stopping criterion is straightforward since, at any step, the two current vectors constitute a framing of the reachability probabilities. Similar computations of parallel under- and over-approximations have been used in [7], in order to detect steady-state on-the-fly during the transient analysis of continuous-time Markov chains. In [8], under- and over-approximations of reachability probabilities in MDPs are obtained by substituting to the MDP a stochastic game. Combining it with a CEGAR-based procedure leads to an iterative procedure with approximations converging to the exact values. However the speed of convergence is only studied from an experimental point of view. Afterwards, we provide probabilistic interpretations for the adjacent sequences of the interval iteration algorithm. Combining such an interpretation with the safety convergence rate of the reduced MDP allows us to exhibit a convergence rate for interval iteration algorithm. At last, exploiting this convergence rate, we significantly im-

prove the bound on the number of iterations required to get the exact values by a rounding procedure.

Related work. Interestingly, our approach has been realized in parallel of Brázdil et al [2] that solves a different problem with similar ideas. There, authors use some machine learning algorithm, namely real-time dynamic programming, in order to avoid to apply the full operator at each step of the value iteration, but rather to partially apply it based on some statistical test. Using the same idea of lower and upper approximations, they prove that their algorithm *almost surely* converges towards the optimal probability, in case of MDPs without non-trivial end components. In the presence of non-trivial end components, rather than computing in advance a simplified equivalent MDP as we do, they rather compute the simplification on-the-fly. It allows them to also obtain results in the case where the MDP is not explicitly given. However, no analysis of the speed of convergence of their algorithm is provided, nor are given explicit stopping criteria enabling an exact computation of values.

Outline. Section 2 introduces MDPs and the reachability/safety problems. It also includes MEC decomposition, dedicated MDP transformations and characterization of minimal and maximal reachability probabilities as unique fixed points of operators. Section 3 presents our main contributions: the interval iteration algorithm, the analysis of the convergence rate and a better bound for the number of iterations required for obtaining the exact values by rounding. Due to space constraints, a complete version, with full proofs, can be found in [6].

2 Reachability problems for Markov decision processes

2.1 Problem specification

We mainly follow the notations of [5]. We denote by $Dist(S)$ the set of *distributions* over a finite set S , i.e., every mapping $p: S \rightarrow [0, 1]$ from S to the set $[0, 1]$ such that $\sum_{s \in S} p(s) = 1$. The support of a distribution p , denoted by $Supp(p)$, is the subset of S defined by $Supp(p) = \{s \in S \mid p(s) > 0\}$.

A *Markov Decision Process* (MDP) is a tuple $\mathcal{M} = (S, \alpha_{\mathcal{M}}, \delta_{\mathcal{M}})$ where S is a finite set of states; $\alpha_{\mathcal{M}} = \bigcup_{s \in S} A(s)$ where every $A(s)$ is a non empty finite set of actions with $A(s) \cap A(s') = \emptyset$ for all $s \neq s'$; and $\delta_{\mathcal{M}}: S \times \alpha_{\mathcal{M}} \rightarrow Dist(S)$ is a partial probabilistic transition function defined for (s, a) if and only if $a \in A(s)$.

The dynamic of the system is defined as follows. Given a current state s , an action $a \in A(s)$ is chosen non deterministically. The next state is then randomly selected, using the corresponding distribution $\delta_{\mathcal{M}}(s, a)$, i.e., the probability that a transition to s' occurs equals $\delta_{\mathcal{M}}(s, a)(s')$. In a more suggestive way, one denotes $\delta_{\mathcal{M}}(s, a)(s')$ by $\delta_{\mathcal{M}}(s'|s, a)$ and $\sum_{s' \in S'} \delta_{\mathcal{M}}(s'|s, a)$ by $\delta_{\mathcal{M}}(S'|s, a)$.

More formally, an *infinite path* through an MDP is a sequence $\pi = s_0 \xrightarrow{a_0} s_1 \xrightarrow{a_1} \dots$ where $s_i \in S$, $a_i \in A(s_i)$ and $\delta_{\mathcal{M}}(s_{i+1}|s_i, a_i) > 0$ for all $i \in \mathbb{N}$: in the following, state s_i is denoted by $\pi(i)$. For every $i \in \mathbb{N}$, $\pi_{\uparrow i}$ denotes the suffix of π starting in s_i , i.e., $\pi_{\uparrow i} = s_i \xrightarrow{a_i} s_{i+1} \dots$. A *finite path* $\rho = s_0 \xrightarrow{a_0} s_1 \xrightarrow{a_1} \dots \xrightarrow{a_{n-1}} s_n$ is a prefix of an infinite path ending in a state s_n , denoted by $\text{last}(\rho)$. We

denote by $\text{Path}_{\mathcal{M},s}$ (respectively, $\text{FPath}_{\mathcal{M},s}$) the set of infinite paths (respectively, finite paths) starting in state s , whereas $\text{Path}_{\mathcal{M}}$ (respectively, $\text{FPath}_{\mathcal{M}}$) denotes the set of all infinite paths (respectively, finite paths).

To associate a probability space with an MDP, we need to eliminate the non-determinism of the behaviour. This is done by introducing policies (also called schedulers or strategies). A *policy* of an MDP $\mathcal{M} = (S, \alpha_{\mathcal{M}}, \delta_{\mathcal{M}})$ is a function $\sigma: \text{FPath}_{\mathcal{M}} \rightarrow \text{Dist}(\alpha_{\mathcal{M}})$ such that $\sigma(\rho)(a) > 0$ only if $a \in A(\text{last}(\rho))$. One denotes $\sigma(\rho)(a)$ by $\sigma(a|\rho)$. We denote by $\text{Pol}_{\mathcal{M}}$ the set of all policies of \mathcal{M} . A policy σ is *deterministic* when $\sigma(\rho)$ is a Dirac distribution for every $\rho \in \text{FPath}_{\mathcal{M}}$ (in that case, $\sigma(\rho)$ denotes the action $a \in A(\text{last}(\rho))$ associated to probability one); it is *stationary* (also called memoryless) if $\sigma(\rho)$ only depends on $\text{last}(\rho)$.

A policy σ and an initial state $s \in S$ yields a discrete-time Markov chain \mathcal{M}_s^σ (see [5, Definition 10]), whose states are the finite paths of $\text{FPath}_{\mathcal{M},s}$. The probability measure $Pr_{\mathcal{M}^\sigma, s}$ over paths of the Markov chain starting in s (with basic cylinders being generated by finite paths) defines a probability measure $Pr_{\mathcal{M}, s}^\sigma$ over $\text{Path}_{\mathcal{M},s}$, capturing the behavior of \mathcal{M} from state s under policy σ . Let $\rho_n = s_0 \xrightarrow{a_0} s_1 \xrightarrow{a_1} \dots \xrightarrow{a_{n-1}} s_n$ and $\rho_{n+1} = s_0 \xrightarrow{a_0} s_1 \xrightarrow{a_1} \dots \xrightarrow{a_n} s_{n+1}$, the probability measure is inductively defined by

$$Pr_{\mathcal{M}, s_0}^\sigma(\rho_{n+1}) = Pr_{\mathcal{M}, s_0}^\sigma(\rho_n) \sum_{a \in A(s_n)} \sigma(a|\rho_n) \delta_{\mathcal{M}}(s_{n+1}|s_n, a).$$

One specifies properties on infinite paths as follows. Given a subset $S' \subseteq S$ of states and $\pi = s_0 \xrightarrow{a_0} s_1 \xrightarrow{a_1} \dots \in \text{Path}_{\mathcal{M}}$, $\pi \models S'$ iff $s_0 \in S'$. The atomic proposition $\{s\}$ is more concisely denoted by s . One also uses Boolean operators \neg , \wedge and \vee for building formulas. We finally use temporal operators **F** (for *Finally*) and **G** (for *Globally*). For a property φ , we let $\pi \models \text{F}\varphi$ if there exists $i \in \mathbb{N}$ such that the suffix $\pi_{\uparrow i}$ of π verifies $\pi_{\uparrow i} \models \varphi$. The dual operator **G** is defined by $\text{G}\varphi \equiv \neg \text{F}\neg\varphi$. One also considers restricted scopes of these operators: $\pi \models \text{F}^{\leq n}\varphi$ if there exists $0 \leq i \leq n$ such that $\pi_{\uparrow i} \models \varphi$, and $\text{G}^{\leq n}\varphi \equiv \neg \text{F}^{\leq n}\neg\varphi$. Given a property φ on infinite paths one denotes $Pr_{\mathcal{M}, s}^\sigma(\{\pi \in \text{Path}_{\mathcal{M}, s} \mid \pi \models \varphi\})$ more concisely by $Pr_{\mathcal{M}, s}^\sigma(\varphi)$.

Given a subset of target states T , *reachability* properties are specified by FT and *safety* properties by $\text{G}\neg T$. Our main goal is to compute the infimum and supremum reachability and safety probabilities, with respect to the policies, i.e., for $\varphi \in \{\text{FT}, \text{G}\neg T\}$: $Pr_{\mathcal{M}, s}^{\min}(\varphi) = \inf_{\sigma \in \text{Pol}_{\mathcal{M}}} Pr_{\mathcal{M}, s}^\sigma(\varphi)$ and $Pr_{\mathcal{M}, s}^{\max}(\varphi) = \sup_{\sigma \in \text{Pol}_{\mathcal{M}}} Pr_{\mathcal{M}, s}^\sigma(\varphi)$. Since $Pr_{\mathcal{M}, s}^\sigma(\text{G}\neg T) = 1 - Pr_{\mathcal{M}, s}^\sigma(\text{FT})$, one immediately gets: $Pr_{\mathcal{M}, s}^{\max}(\text{G}\neg T) = 1 - Pr_{\mathcal{M}, s}^{\min}(\text{FT})$, and $Pr_{\mathcal{M}, s}^{\min}(\text{G}\neg T) = 1 - Pr_{\mathcal{M}, s}^{\max}(\text{FT})$.

Thus we focus on reachability problems and without loss of generality, all the states of T may be merged in a single state called s_+ with $A(s_+) = \{\text{loop}_+\}$ such that $\delta_{\mathcal{M}}(s_+|s_+, \text{loop}_+) = 1$. In the sequel, the vector $(Pr_{\mathcal{M}, s}^\sigma(\varphi))_{s \in S}$ (respectively, $(Pr_{\mathcal{M}, s}^{\min}(\varphi))_{s \in S}$ and $(Pr_{\mathcal{M}, s}^{\max}(\varphi))_{s \in S}$) of probabilities will be denoted by $Pr_{\mathcal{M}}^\sigma(\varphi)$ (respectively, $Pr_{\mathcal{M}}^{\min}(\varphi)$ and $Pr_{\mathcal{M}}^{\max}(\varphi)$).

2.2 MEC decomposition and transient behaviour

In our approach, we first reduce an MDP by a qualitative analysis based on *end components*. We adopt here a slightly different definition of the usual one by allowing trivial end components (see later on). Preliminarily, the *graph* of an MDP \mathcal{M} is defined as follows: the set of its vertices is S and there is an edge from s to s' if there is some $a \in A(s)$ with $\delta_{\mathcal{M}}(s'|s, a) > 0$.

Definition 1 (end component). *Let $\mathcal{M} = (S, \alpha_{\mathcal{M}}, \delta_{\mathcal{M}})$. Then (S', α') with $\emptyset \neq S' \subseteq S$ and $\alpha' \subseteq \bigcup_{s \in S'} A(s)$ is an end component if (i) for all $s \in S'$ and $a \in A(s) \cap \alpha'$, $\text{Supp}(\delta_{\mathcal{M}}(s, a)) \subseteq S'$; (ii) the graph of (S', α') is strongly connected.*

Given two end components, one says that (S', α') is smaller than (S'', α'') , denoted by $(S', \alpha') \preceq (S'', \alpha'')$, if $S' \subseteq S''$ and $\alpha' \subseteq \alpha''$. Given some state s , there is a minimal end component containing s namely $(\{s\}, \emptyset)$. Such end components are called *trivial* end components. The union of two end components that share a state is also an end component. Hence, *maximal* end components (MEC) do not share states and cover all states of S . Furthermore, we consider *bottom* MEC (BMEC): a MEC (S', α') is a BMEC if $\alpha' = \bigcup_{s \in S'} A(s)$. For instance $(\{s_+\}, \{\text{loop}_+\})$ is a BMEC. Every MDP contains at least one BMEC.

Fig. 1-(a) shows the decomposition in MEC of an MDP. There are two BMECs $(\{s_+\}, \{\text{loop}_+\})$ and $(\{b, b'\}, \{d, e\})$, one trivial MEC $(\{t\}, \emptyset)$ and another MEC $(\{s, s'\}, \{a, c\})$.

The set of MECs of an MDP can be computed in polynomial time (see for instance [4]). It defines a partition of $S = \bigsqcup_{i=k}^K S_k \uplus \bigsqcup_{\ell=1}^L \{t_\ell\} \uplus \bigsqcup_{m=0}^M B_m$ where $\{t_\ell\}$ is the set of states of a trivial MEC, B_m is the set of states a BMEC and S_k 's are the set of states of the other MECs. By convention, $B_0 = \{s_+\}$. The next proposition is the key ingredient of our approach.

Proposition 2. *Let \mathcal{M} be an MDP such that its MEC decomposition only contains trivial MECs and BMECs, i.e., $S = \bigsqcup_{\ell=1}^L \{t_\ell\} \uplus \bigsqcup_{m=0}^M B_m$. Then:*

1. *There is a partition $S = \bigsqcup_{0 \leq i \leq I} G_i$ such that $G_0 = \bigsqcup_{m=0}^M B_m$ and for all $1 \leq i \leq I$, $s \in G_i$ and $a \in A(s)$, there is $s' \in \bigcup_{j < i} G_j$ such that $\delta_{\mathcal{M}}(s'|s, a) > 0$.*
2. *Let η be the smallest positive probability occurring in the distributions of \mathcal{M} . Then for all $n \in \mathbb{N}$, and for all $s \in S$, $Pr_{\mathcal{M}, s}^{\max}(\mathbf{G}^{\leq n I} \neg G_0) \leq (1 - \eta^I)^n$.*
3. *For all $s \in S$, $Pr_{\mathcal{M}, s}^{\max}(\mathbf{G} \neg G_0) = 0$.*

Proof. (Sketch) 1. One builds the partition of S by induction. We first let $G_0 = \bigsqcup_{m=0}^M B_m$. Then, assuming that G_0, \dots, G_i have been defined, we let $G_{i+1} = \{s \in S \setminus \bigcup_{j \leq i} G_j \mid \forall a \in A(s) \exists s' \in \bigcup_{j \leq i} G_j \delta_{\mathcal{M}}(s'|s, a) > 0\}$. The construction stops when some G_i is empty. If G_I is the last non-empty set, it can easily be checked that $S = \bigcup_{i \leq I} G_i$.

2. One observes that the path property $\mathbf{G}^{\leq n} \neg G_0$ only depends on the prefix of length n . So there is only a finite number of policies up to n and we denote σ_n the policy that achieves $Pr_{\mathcal{M}, s}^{\max}(\mathbf{G}^{\leq n} \neg G_0)$. Observe also that after a path of length $k < n$ leading to state $s \notin G_0$, policy σ_n may behave as policy σ_{n-k} starting in s . The property may then be shown by using the fact that for all

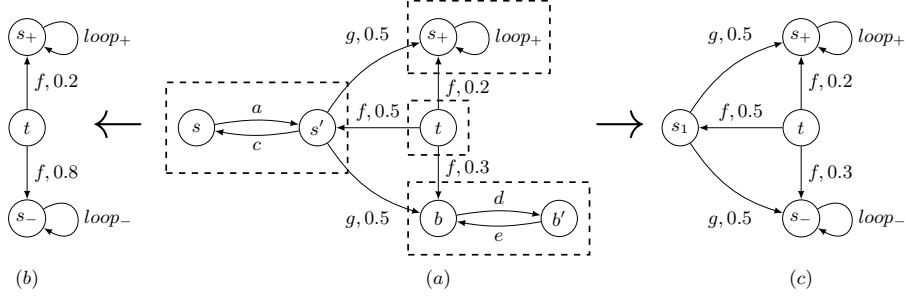


Fig. 1. (a) An MDP and its MEC decomposition, (b) its min-reduction, and (c) its max-reduction

state s and policy σ , there is a path of length at most I in \mathcal{M}^σ from s to ρ with $\text{last}(\rho) \in G_0$, showing that $\Pr_{\mathcal{M},s}^\sigma(\mathbf{G}^{\leq I} \neg G_0) \leq (1 - \eta^I)$.

3. The last assertion is a straightforward consequence of the previous one. \square

This proposition shows the interest of eliminating MECs that are neither trivial ones nor BMECs. In the following, we consider the partition $S = \bigsqcup_{i=k}^K S_k \sqcup \bigsqcup_{\ell=1}^L \{t_\ell\} \sqcup \bigsqcup_{m=0}^M B_m$ where $\{t_\ell\}$'s are trivial MECs, B_m 's are BMECs and S_k 's are all the other MECs. A quotienting of an MDP has been introduced in [4, Algorithm 3.3] in order to decrease the complexity of the computation for reachability properties. We now introduce two variants of reductions for MDPs depending on the kind of probabilities we want to compute.

2.3 Characterization of minimal reachability probabilities

The reduction in the case of minimal reachability probabilities consists in merging all non-trivial MECs different from $(\{s_+\}, \{loop_+\})$ into a fresh state s_- : all these states merged into s_- will have a zero minimal reachability probability.

Definition 3 (min-reduction). Let \mathcal{M} be an MDP with the partition of $S = \bigsqcup_{i=k}^K S_k \sqcup \bigsqcup_{\ell=1}^L \{t_\ell\} \sqcup \bigsqcup_{m=0}^M B_m$. We define $\mathcal{M}^\bullet = (S^\bullet, \alpha_{\mathcal{M}^\bullet}, \delta_{\mathcal{M}^\bullet})$ by:

- $S^\bullet = \{s_-, s_+, t_1, \dots, t_L\}$, and for all $s \in S$, s^\bullet is defined by: (1) $s^\bullet = t_\ell$ if $s = t_\ell$, (2) $s^\bullet = s_+$ if $s = s_+$, and (3) $s^\bullet = s_-$ otherwise.
- $A^\bullet(s_-) = \{loop_-\}$, $A^\bullet(s_+) = \{loop_+\}$ and for all $1 \leq \ell \leq L$, $A^\bullet(t_\ell) = A(t_\ell)$.
- For all $1 \leq \ell, \ell' \leq L$, $a \in A^\bullet(t_\ell)$,

$$\begin{aligned} \delta_{\mathcal{M}^\bullet}(s_-|t_\ell, a) &= \delta_{\mathcal{M}}(\bigsqcup_{i=k}^K S_k \sqcup \bigsqcup_{m=1}^M B_m|t_\ell, a), \\ \delta_{\mathcal{M}^\bullet}(s_+|t_\ell, a) &= \delta_{\mathcal{M}}(s_+|t_\ell, a), \quad \delta_{\mathcal{M}^\bullet}(t_{\ell'}|t_\ell, a) = \delta_{\mathcal{M}}(t_{\ell'}|t_\ell, a), \\ \delta_{\mathcal{M}^\bullet}(s_+|s_+, loop_+) &= \delta_{\mathcal{M}}(s_-|s_-, loop_-) = 1. \end{aligned}$$

An MDP \mathcal{M} is called *min-reduced* if $\mathcal{M} = \mathcal{N}^\bullet$ for some MDP \mathcal{N} . The min-reduction of an MDP is illustrated in Fig. 1-(b). The single trivial MEC $(\{t\}, \emptyset)$ is preserved while MECs $(\{b, b'\}, \{d, e\})$ and $(\{s, s'\}, \{a, c\})$ are merged in s_- .

Proposition 4. *Let \mathcal{M} be an MDP and \mathcal{M}^\bullet be its min-reduced MDP. Then for all $s \in S$, $Pr_{\mathcal{M},s}^{\min}(\mathbf{F} s_+) = Pr_{\mathcal{M}^\bullet,s}^{\min}(\mathbf{F} s_+)$.*

We now establish another property of the min-reduced MDP that allows us to use Proposition 2.

Lemma 5. *Let \mathcal{M}^\bullet be the min-reduced MDP of an MDP \mathcal{M} . Then every state $s \in S^\bullet \setminus \{s_-, s_+\}$ is a trivial MEC.*

In order to characterize $Pr_{\mathcal{M}}^\sigma(\mathbf{F} s_+)$ with a fixpoint equation, we define the set of S -vectors as $\mathcal{V} = \{x = (x_s)_{s \in S} \mid \forall s \in S \setminus \{s_-, s_+\} 0 \leq x_s \leq 1 \wedge x_{s_+} = 1 \wedge x_{s_-} = 0\}$. We also introduce the operator $f_{\min}: \mathcal{V} \rightarrow \mathcal{V}$ by letting for all $x \in \mathcal{V}$: $f_{\min}(x)_s = \min_{a \in A(s)} \sum_{s' \in S} \delta_{\mathcal{M}}(s'|s, a)x_{s'}$ for every $s \in S \setminus \{s_-, s_+\}$, $f_{\min}(x)_{s_-} = 0$ and $f_{\min}(x)_{s_+} = 1$.

We claim that there is a single fixed point of f_{\min} . In order to establish that claim, given a stationary deterministic strategy σ , we introduce the operator $f_\sigma: \mathcal{V} \rightarrow \mathcal{V}$ defined for all $x \in \mathcal{V}$ by: $f_\sigma(x)_s = \sum_{s' \in S} \delta_{\mathcal{M}}(s'|s, \sigma(s))x_{s'}$ for every $s \in S \setminus \{s_-, s_+\}$, $f_\sigma(x)_{s_-} = 0$ and $f_\sigma(x)_{s_+} = 1$.

Proposition 6. *Let \mathcal{M} be a min-reduced MDP. $Pr_{\mathcal{M}}^\sigma(\mathbf{F} s_+)$ is the unique fixed point of f_σ . $Pr_{\mathcal{M}}^{\min}(\mathbf{F} s_+)$ is the unique fixed point of f_{\min} and it is obtained by a stationary deterministic policy.*

2.4 Characterization of maximal reachability probabilities

The reduction for maximal reachability probabilities is more complex. Indeed, we cannot merge any non-trivial MEC different from $(\{s_+\}, \{loop_+\})$ into the state s_- anymore, since some of these states may have a non-zero maximal reachability probability. Hence, we consider a fresh state s_k for each MEC S_k and simply merge all BMECs B_m 's different from $(\{s_+\}, \{loop_+\})$ into state s_- .

Definition 7 (max-reduction). *Let \mathcal{M} be a MDP with the partition of $S = \bigsqcup_{i=k}^K S_k \sqcup \bigsqcup_{\ell=1}^L \{t_\ell\} \sqcup \bigsqcup_{m=0}^M B_m$. Then the max-reduced $\mathcal{M}^\bullet = (S^\bullet, \alpha_{\mathcal{M}^\bullet}, \delta_{\mathcal{M}^\bullet})$ is defined by:*

- $S^\bullet = \{s_-, s_+, t_1, \dots, t_L, s_1, \dots, s_K\}$. For all $s \in S$, one defines s^\bullet by: (1) $s^\bullet = t_\ell$ if $s = t_\ell$, (2) $s^\bullet = s_+$ if $s = s_+$, (3) $s^\bullet = s_k$ if $s \in S_k$, and (4) $s^\bullet = s_-$ otherwise.
- $A^\bullet(s_-) = \{loop_-\}$, $A^\bullet(s_+) = \{loop_+\}$ for all $1 \leq \ell \leq L$, $A^\bullet(t_\ell) = A(t_\ell)$, and for all $1 \leq k \leq K$, $A^\bullet(s_k) = \{a \mid \exists s \in S_k a \in A(s) \wedge \text{Supp}(\delta_{\mathcal{M}}(s, a)) \not\subseteq S_k\}$.
- For all $1 \leq \ell, \ell' \leq L$, $a \in A^\bullet(t_\ell)$, $1 \leq k, k' \leq K$, $b \in A^\bullet(s_k) \cap A_s$ with $s \in S_k$,

$$\begin{aligned} \delta_{\mathcal{M}^\bullet}(s_-|t_\ell, a) &= \delta_{\mathcal{M}}(\bigsqcup_{m=1}^M B_m|t_\ell, a), & \delta_{\mathcal{M}^\bullet}(s_+|t_\ell, a) &= \delta_{\mathcal{M}}(s_+|t_\ell, a), \\ \delta_{\mathcal{M}^\bullet}(t_{\ell'}|t_\ell, a) &= \delta_{\mathcal{M}}(t_{\ell'}|t_\ell, a), & \delta_{\mathcal{M}^\bullet}(s_k|t_\ell, a) &= \delta_{\mathcal{M}}(S_k|t_\ell, a), \\ \delta_{\mathcal{M}^\bullet}(s_-|s_k, b) &= \delta_{\mathcal{M}}(\bigsqcup_{m=1}^M B_m|s, b), & \delta_{\mathcal{M}^\bullet}(s_+|s_k, b) &= \delta_{\mathcal{M}}(s_+|s, b), \\ \delta_{\mathcal{M}^\bullet}(t_\ell|s_k, b) &= \delta_{\mathcal{M}}(t_\ell|s, b), & \delta_{\mathcal{M}^\bullet}(s_{k'}|s_k, b) &= \delta_{\mathcal{M}}(S_{k'}|s, b), \\ \delta_{\mathcal{M}^\bullet}(s_+|s_+, loop_+) &= \delta_{\mathcal{M}}(s_-|s_-, loop_-) = 1. \end{aligned}$$

Observe that \mathcal{M}^\bullet is indeed an MDP since $A^\bullet(s_k)$ cannot be empty (otherwise S_k would be BMEC). Fig. 1-(c) illustrates the max-reduction of an MDP. The single trivial MEC $(\{t\}, \emptyset)$ is preserved while MEC $(\{b, b'\}, \{d, e\})$ is merged in s_- . The MEC $(\{s, s'\}, \{a, c\})$ is now merged into s_1 with only action g preserved.

The following propositions are similar to Proposition 4 and Lemma 5 for the min-reductions.

Proposition 8 ([4, Thm. 3.8]). *Let \mathcal{M} be an MDP and \mathcal{M}^\bullet be its max-reduced MDP. Then for all $s \in S$, $Pr_{\mathcal{M},s}^{\max}(\mathbf{F} s_+) = Pr_{\mathcal{M}^\bullet, s^\bullet}^{\max}(\mathbf{F} s_+)$.*

Lemma 9. *Let \mathcal{M}^\bullet be the max-reduced MDP of an MDP \mathcal{M} . Then every state $s \in S^\bullet \setminus \{s_-, s_+\}$ is a trivial MEC.*

As for minimal reachability probabilities, we introduce operator $f_{\max}: \mathcal{V} \rightarrow \mathcal{V}$ by letting for all $x \in \mathcal{V}$: $f_{\max}(x)_s = \max_{a \in A(s)} \sum_{s' \in S} \delta_{\mathcal{M}}(s, a)(s')x_{s'}$ for all $s \in S \setminus \{s_-, s_+\}$, $f_{\max}(x)_{s_-} = 0$ and $f_{\max}(x)_{s_+} = 1$.

We observe that Lemma 9 combined with Proposition 2 ensures that in a max-reduced MDP \mathcal{M} , for any policy σ , $S \setminus \{s_-, s_+\}$ is a set of transient states of \mathcal{M}^σ . This helps to prove that Proposition 6 also holds for max-reduced MDPs:

Proposition 10. *Let \mathcal{M} be a max-reduced MDP. $Pr_{\mathcal{M}}^\sigma(\mathbf{F} s_+)$ is the unique fixed point of f_σ . $Pr_{\mathcal{M}}^{\max}(\mathbf{F} s_+)$ is the unique fixed point of f_{\max} and it is obtained by a stationary deterministic policy.*

Discussion. Usually, algorithms that compute maximal and minimal reachability probabilities first determine the set of states for which those probabilities are 0 or 1, and merge them in states s_- and s_+ respectively (see for instance [5, Algorithms 1-4]). For the case of minimal reachability probabilities, the MDP obtained after this transformation—which is a quotient of our \mathcal{M}^\bullet —fulfills the hypotheses of Proposition 2 and our further development is still valid.

Unfortunately, it does not hold in the maximal case: for the MDP on the left of Fig. 1-(a), the obtained MDP, that we call \mathcal{M}' , simply merges $\{b, b'\}$ into s_- , without merging $\{s, s'\}$ (since the maximal probability to reach s_+ from s or s' is equal to 0.5, when choosing action b in s'). Moreover, Proposition 10 does not hold either in \mathcal{M}' for maximal probabilities³. In fact, the vector of maximal probabilities in \mathcal{M}' is only the smallest fixed point of f_{\max} . Indeed, the reader can check that the vector which is equal to 0 for s_- , 0.7 for t , and 1 for all the other states is also a fixed point of f_{\max} , whereas the maximal reachability probability to reach s_+ from s or s' is equal to 0.5. Notice that in the max-reduction \mathcal{M}^\bullet of this MDP, the MEC $(\{s, s'\}, \{a, c\})$ is merged into a single state, hence removing this non-unicity problem, as shown in Proposition 10.

While this issue does not preclude the standard computation of the probabilities, the approach we have followed enables us to solve the convergence issues of the previous methods. This is the subject of the next section.

³ This is already observed in [5], but a wrong statement is made in [1, Thm. 10.100].

3 Value iteration for reachability objectives

This section presents the value iteration algorithm used, for example in the PRISM model-checker [9], to compute optimal reachability probabilities of an MDP. After stating convergence issues of this method, we give a new algorithm, called *interval iteration algorithm*, and the strong guarantees that it gives.

3.1 Convergence issues

The idea of the value iteration algorithm is to compute the fixed points of f_{\min} and f_{\max} (more precisely, the smallest fixed points of f_{\min} and f_{\max}) by iterating them on a given initial vector, until a certain convergence criterion is met. More precisely, as recalled in [5], we let $x^{(0)}$ defined by $x_{s_+}^{(0)} = 1$ and $x_s^{(0)} = 0$ for $s \neq s_+$ (observe that $x^{(0)}$ is the minimal vector of \mathcal{V} for the pointwise order), and we then build one of the two sequences $\underline{x} = (\underline{x}^{(n)})_{n \in \mathbb{N}}$ or $\bar{x} = (\bar{x}^{(n)})_{n \in \mathbb{N}}$ defined by

- $\underline{x}^{(0)} = x^{(0)}$ and for all $n \in \mathbb{N}$, $\underline{x}^{(n+1)} = f_{\min}(\underline{x}^{(n)})$;
- $\bar{x}^{(0)} = x^{(0)}$ and for all $n \in \mathbb{N}$, $\bar{x}^{(n+1)} = f_{\max}(\bar{x}^{(n)})$.

Since f_{\min} and f_{\max} are monotonous operators and due to the choice of the initial vector, \underline{x} and \bar{x} are non-decreasing bounded sequences, hence convergent. Let $\underline{x}^{(\infty)}$ and $\bar{x}^{(\infty)}$ be their respective limits. By continuity of f_{\min} and f_{\max} , $\underline{x}^{(\infty)}$ (respectively, $\bar{x}^{(\infty)}$) is a fixed point of f_{\min} (respectively, f_{\max}). Due to Propositions 6 and 10, $\underline{x}^{(\infty)}$ (respectively, $\bar{x}^{(\infty)}$) is the vector $Pr_{\mathcal{M}}^{\min}(\mathbf{F} s^+)$ (respectively, $Pr_{\mathcal{M}}^{\max}(\mathbf{F} s^+)$) of minimal (respectively, maximal) reachability probabilities.

In practice, several stopping criteria can be chosen. In the model-checker PRISM [9], two criteria are implemented. For a vector $x \in \mathcal{V}$, we let $\|x\| = \max_{s \in S} |x_s|$. For $x \in \{\underline{x}, \bar{x}\}$ and a given threshold $\varepsilon > 0$, the *absolute criterion* consists in stopping once $\|x^{(n+1)} - x^{(n)}\| \leq \varepsilon$, whereas the *relative criterion* considers $\max_{s \in S} (x_s^{(n+1)} - x_s^{(n)}) / x_s^{(n)} \leq \varepsilon$. However, as noticed in [5], no guarantees are obtained when using such value iteration algorithms, whatever the stopping criterion. As an example, consider the MDP (indeed the Markov chain) of Fig. 2. It is easy to check that (minimal and maximal) reachability probability of $s^+ = 0$ in state n is $1/2$. However, if ε is chosen as $1/2^n$ (or any value above), the sequence of vectors computed by the value iteration algorithm will be $x^{(0)} = (1, 0, 0, \dots, 0, 0, \dots, 0)$, $x^{(1)} = (1, 1/2, 0, \dots, 0, 0, \dots, 0)$, $x^{(2)} = (1, 1/2, 1/4, \dots, 0, 0, \dots, 0)$, \dots , $x^{(n)} = (1, 1/2, 1/4, \dots, 1/2^n, 0, \dots, 0)$, at which point the absolute stopping criterion is met. Hence, the algorithm outputs $x_n^{(n)} = 1/2^n$ as the reachability probability of $s_+ = \{0\}$ in state n .

Example 11. The use of PRISM confirms this phenomenon. On this MDP, choosing $n = 10$ and threshold $\varepsilon = 10^{-3} < 1/2^{10}$, the absolute stopping criterion leads to the probability $9.77 \times 10^{-4} \approx 1/2^{10}$, whereas the relative stopping criterion leads to the probability 0.198. It has to be noticed that the tool indicates that the value iteration has converged, and does not warn the user that a possible problem may have arisen.

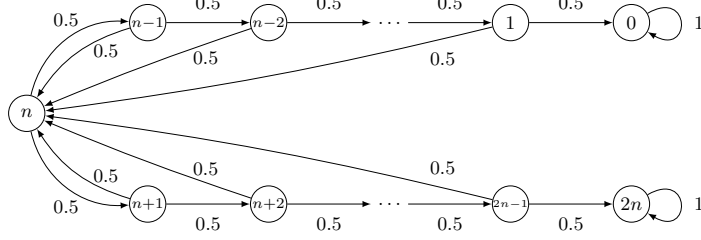


Fig. 2. A Markov chain with problems of convergence in value iteration

Algorithm 1: Interval iteration algorithm for minimum reachability

Input: Min-reduced MDP $\mathcal{M} = (S, \alpha_{\mathcal{M}}, \delta_{\mathcal{M}})$, convergence threshold ε

- 1 $x_{s_+} := 1; x_{s_-} := 0; y_{s_+} := 1; y_{s_-} := 0$
- 2 **foreach** $s \in S \setminus \{s_+, s_-\}$ **do** $x_s := 0; y_s := 1$
- 3 **repeat**
- 4 **foreach** $s \in S \setminus \{s_+, s_-\}$ **do**
- 5 $x'_s := \min_{a \in A(s)} \sum_{s' \in S} \delta_{\mathcal{M}}(s, a)(s') x_{s'}$
- 6 $y'_s := \min_{a \in A(s)} \sum_{s' \in S} \delta_{\mathcal{M}}(s, a)(s') y_{s'}$
- 7 $\delta := \max_{s \in S} (y'_s - x'_s)$
- 8 **foreach** $s \in S \setminus \{s_+, s_-\}$ **do** $x'_s := x_s; y'_s := y_s$
- 9 **until** $\delta \leq \varepsilon$
- 10 **return** $(x_s)_{s \in S}, (y_s)_{s \in S}$

We consider a modification of the algorithm in order to obtain a convergence guarantee when stopping the value iteration algorithm. We provide (1) stopping criteria for approximation and exact computations and, (2) rate of convergence.

3.2 Stopping criterion for ε -approximation

Here, we introduce two other sequences. For that, let vector $y^{(0)}$ be the maximal vector of \mathcal{V} , defined by $y_{s_-}^{(0)} = 0$ and $y_s^{(0)} = 1$ for $s \neq s_-$. We then define inductively the two sequences \underline{y} and \overline{y} of vectors by

- $\underline{y}^{(0)} = y^{(0)}$ and for all $n \in \mathbb{N}$, $\underline{y}^{(n+1)} = f_{\min}(\underline{y}^{(n)})$;
- $\overline{y}^{(0)} = y^{(0)}$ and for all $n \in \mathbb{N}$, $\overline{y}^{(n+1)} = f_{\max}(\overline{y}^{(n)})$.

Because of the new choice for the initial vector, notice that \underline{y} and \overline{y} are non-increasing sequences. Hence, with the same reasoning as above, we know that these sequences converge, and that their limit, denoted by $\underline{y}^{(\infty)}$ and $\overline{y}^{(\infty)}$ respectively, are the minimal (respectively, maximal) reachability probabilities. In particular, notice that \underline{x} and \underline{y} , as well as \overline{x} and \overline{y} , are adjacent sequences, and that $\underline{x}^{(\infty)} = \underline{y}^{(\infty)} = Pr_{\mathcal{M}}^{\min}(\mathbf{F} s^+)$ and $\overline{x}^{(\infty)} = \overline{y}^{(\infty)} = Pr_{\mathcal{M}}^{\max}(\mathbf{F} s^+)$.

Let us first consider a min-reduced MDP \mathcal{M} . Then, our new value iteration algorithm computes both in the same time sequences \underline{x} and \underline{y} and stops as soon

the sequel, we assume that there is at least one transition probability $0 < \delta < 1$ (otherwise the problems are trivial).

Theorem 15. *For a min- or max-reduced MDP \mathcal{M} , and a convergence threshold ε , the interval iteration algorithm converges in at most $I \lceil \frac{\log \varepsilon}{\log(1-\eta^I)} \rceil$ steps, where I and η are introduced in Proposition 2.*

Proof. Let σ be the policy corresponding to the minimal probability of satisfying $\mathbf{G}^{\leq n} \neg s_-$ and σ' be the policy corresponding to the minimal probability of satisfying $\mathbf{F}^{\leq n} s_+$. In particular, notice that $Pr_{\mathcal{M},s}^{\sigma}(\mathbf{G}^{\leq nI} \neg s_-) \leq Pr_{\mathcal{M},s}^{\sigma'}(\mathbf{G}^{\leq nI} \neg s_-)$.

Since $\mathbf{G}^{\leq n} \neg s_- \equiv \mathbf{G}^{\leq n} \neg (s_- \vee s_+) \vee \mathbf{F}^{\leq n} s_+$, with the disjunction being exclusive, we have for all $s \in S$,

$$\begin{aligned} Pr_{\mathcal{M},s}^{\min}(\mathbf{G}^{\leq nI} \neg s_-) - Pr_{\mathcal{M},s}^{\min}(\mathbf{F}^{\leq nI} s_+) &= Pr_{\mathcal{M},s}^{\sigma}(\mathbf{G}^{\leq nI} \neg s_-) - Pr_{\mathcal{M},s}^{\sigma'}(\mathbf{F}^{\leq nI} s_+) \\ &\leq Pr_{\mathcal{M},s}^{\sigma'}(\mathbf{G}^{\leq nI} \neg s_-) - Pr_{\mathcal{M},s}^{\sigma'}(\mathbf{F}^{\leq nI} s_+) = Pr_{\mathcal{M},s}^{\sigma'}(\mathbf{G}^{\leq nI} \neg (s_- \vee s_+)) \leq (1 - \eta^I)^n \end{aligned}$$

due to Proposition 2. It is easy to show by induction that $\underline{x}^{(n)} = Pr_{\mathcal{M}}^{\min}(\mathbf{F}^{\leq n} s_+)$ and $\underline{y}^{(n)} = Pr_{\mathcal{M}}^{\min}(\mathbf{G}^{\leq n} \neg s_-)$. Then, we have $\|\underline{y}^{(nI)} - \underline{x}^{(nI)}\| \leq (1 - \eta^I)^n$. In conclusion, the stopping criterion is met when $(1 - \eta^I)^n \leq \varepsilon$, i.e., after at most $I \lceil \frac{\log \varepsilon}{\log(1-\eta^I)} \rceil$ steps. A similar proof can be made for maximal probabilities. \square

It may also be noticed, from similar arguments, that for all n , $\|\underline{y}^{((n+1)I)} - \underline{x}^{((n+1)I)}\| \leq (1 - \eta^I) \|\underline{y}^{(nI)} - \underline{x}^{(nI)}\|$ (and similarly for the maximum case), implying that the value iteration algorithm has a linear rate of convergence.

Remark 16. One may use this convergence rate to delay the computation of one of the two adjacent sequences of Algorithm 1. Indeed assume that only $x^{(n)}$ is computed until step n . To use the stopping criterion provided by the adjacent sequences, one sets the upper sequence with $y_s^{(n)} = \min(x_s^{(n)} + (1 - \eta^I)^{\lfloor \frac{n}{I} \rfloor}, 1)$ for all $s \notin \{s_-, s_+\}$, $y_{s_+}^{(n)} = 1$, and $y_{s_-}^{(n)} = 0$ and then applies the algorithm. In the favorable cases, this could divide by almost 2 the computation time.

3.4 Stopping criterion for exact computation

In [3], a convergence guarantee was given for MDPs with rational transition probabilities. For such an MDP \mathcal{M} , let d be the largest denominator of transition probabilities (expressed as irreducible fractions), N the number $|S|$ of states, and M the number of transitions with non-zero probabilities. A bound $\gamma = d^{4M}$ was announced so that, after γ^2 iterations, the obtained probabilities lie in intervals that could only contain one possible probability value for the system, permitting to claim for the convergence of the algorithm. So after γ^2 iterations, the actual probability might be computed by considering the rational of the form α/γ closest to the current estimate. However, no proof of this result is given in [3].

Using our simultaneous computation of under- and over-approximations of the probabilities, we provide an alternative stopping criterion for exact computation that moreover exhibits an optimal policy. Its proof is based on the fact that optimal probabilities are rational for which we can control the size of the denominator, and strongly relies on the existence of stationary optimal policies.

Theorem 17. *Let \mathcal{M} be a reduced MDP with rational transition probabilities. Optimal reachability probabilities and optimal policies can be computed by the interval iteration algorithm in at most $\mathcal{O}((1/\eta)^N N^3 \log d)$.*

The theorem also holds for the value iteration algorithm. Observe that our stopping criterion is significantly better than the bound d^{8M} claimed in [3] since $N \leq M$ and $1/\eta \leq d$. Furthermore M may be in $\Omega(N^2)$ even with a single action per state and $1/\eta$ may be significantly smaller than d as for instance in the extreme case $\eta = \frac{1}{2} - \frac{1}{10^n}$ and $d = 10^n$ for some large n .

4 Conclusion

We have provided a framework allowing to guarantee good properties when value iteration algorithm is used to compute optimal reachability probabilities of Markov decision processes. Our study pointed out some difficulties related to non-trivial end components in MDPs, that was not clearly described previously. Moreover, we gave results over the convergence speed, as well as criteria for obtaining exact convergence. As future works, it seems particularly interesting to test this algorithm on real instances, as it is done in [2], where authors moreover apply machine learning techniques.

Acknowledgments. We thank the reviewer that pointed out the similarities between our approach and [2] (to be presented at the next ATVA, in Nov. 2014).

References

1. C. Baier and J.-P. Katoen. *Principles of Model Checking*. MIT Press, 2008.
2. T. Brázdil, K. Chatterjee, M. Chmelík, V. Forejt, J. Křetínský, M. Kwiatkowska, D. Parker, and M. Ujma. Verification of Markov decision processes using learning algorithms. Research Report arXiv:1402.2967, 2014.
3. K. Chatterjee and T. A. Henzinger. Value iteration. In *25 Years of Model Checking, LNCS 5000*, p. 107–138. Springer, 2008.
4. L. de Alfaro. *Formal Verification of Probabilistic Systems*. PhD thesis, Stanford University, 1997.
5. V. Forejt, M. Kwiatkowska, G. Norman, and D. Parker. Automated verification techniques for probabilistic systems. In *SFM'11, LNCS 6659*, p.53–113. Springer, 2011.
6. S. Haddad and B. Monmege. Reachability in MDPs: Refining convergence of value iteration. Technical Report LSV-14-07, LSV, ENS Cachan, 2014. Available at http://www.lsv.ens-cachan.fr/Publis/RAPPORTS_LSV/PDF/rr-lsv-2014-07.pdf.
7. J.-P. Katoen and I. S. Zapreev. Safe on-the-fly steady-state detection for time-bounded reachability. In *QEST'06*, p. 301–310, 2006.
8. M. Kattenbelt and M. Z. Kwiatkowska and G. Norman and D. Parker. A game-based abstraction-refinement framework for Markov decision processes. *Formal Methods in System Design*, vol. 36:3, p. 246–280, 2010
9. M. Kwiatkowska, G. Norman, and D. Parker. PRISM 4.0: Verification of probabilistic real-time systems. In *CAV'11, LNCS 6806*, p. 585–591. Springer, 2011.