

STOCHASTIC PROCESSES WITH EXPECTED STOPPING TIME

KRISHNENDU CHATTERJEE ^a AND LAURENT DOYEN ^b

^a IST Austria

e-mail address: krishnendu.chatterjee@ist.ac.at

^b CNRS & LMF, ENS Paris-Saclay, France

e-mail address: doyen@lsv.fr

ABSTRACT. Markov chains are the de facto finite-state model for stochastic dynamical systems, and Markov decision processes (MDPs) extend Markov chains by incorporating non-deterministic behaviors. Given an MDP and rewards on states, a classical optimization criterion is the maximal expected total reward where the MDP stops after T steps, which can be computed by a simple dynamic programming algorithm. We consider a natural generalization of the problem where the stopping times can be chosen according to a probability distribution, such that the expected stopping time is T , to optimize the expected total reward. Quite surprisingly we establish inter-reducibility of the expected stopping-time problem for Markov chains with the Positivity problem (which is related to the well-known Skolem problem), for which establishing either decidability or undecidability would be a major breakthrough. Given the hardness of the exact problem, we consider the approximate version of the problem: we show that it can be solved in exponential time for Markov chains and in exponential space for MDPs.

1. INTRODUCTION

Stochastic models and optimization. The de facto model for stochastic dynamical systems is finite-state Markov chains [FV97, Gal13, KSK66], with several application domains [BK08]. In modeling optimization problems, rewards are associated with states of the Markov chain, and the optimization criterion is formalized as the expected total reward provided that the Markov chain is stopped after T steps [PT87, FV97]. The extension of Markov chains to allow non-deterministic behavior gives rise to Markov decision processes (MDPs), and the optimization criterion is to maximize, over all non-deterministic choices, the expected total reward for T steps. This notion of optimization for fixed time is called *finite-horizon planning*, which has many applications in logic and verification [EMSS92, BCC⁺03] and control problems in artificial intelligence and robotics [NR10, Chapter 10, Chapter 25], [OR94, Chapter 6].

Optimization with expected stopping time. In the most basic case the stopping time for collecting rewards in the stochastic model is a fixed constant T . A natural generalization

Key words and phrases: Markov chain, stopping time, expected reward, Skolem problem, approximation.

* A preliminary version of this paper appeared in the *Proceedings of the 36th Annual Symposium on Logic in Computer Science (LICS)*, IEEE Computer Society Press, 2021 [CD21].

is to consider that the stochastic model can be stopped at a random time such that the expectation of the stopping time is T . We consider the problem of optimizing (maximizing/minimizing) the expected total reward, when the stopping-time probability distribution can be chosen arbitrarily such that the expected stopping time is T . In other words, we consider stochastic models of Markov chains/MDPs with total reward, and instead of fixed stopping time T , we consider expected stopping time T .

Example and motivation. Consider a classical example where a robot explores a region for natural resources (e.g., the well-studied RockSample problem in AI literature [SS04]), and the exploration of the robot is modeled as a Markov chain. The success of the exploration is characterized by the expected total reward, and the stopping time T denotes the expected duration of the exploration. The expected stopping-time problem asks to choose the probability distribution of the exploration duration to optimize the collected reward, satisfying the average exploration time. A classical stopping-time distribution is the geometric distribution where the stochastic model is stopped at every instant with probability λ , called *discount factor*, which entails that the expected stopping time is $T = 1/\lambda$ [FV97]. The discount-factor model makes an assumption on the shape of the stopping-time distribution, whereas in realistic scenarios the distribution is not precisely known, or time-varying discount factors are considered [DB12]. When the discount factors are not known, then robust solutions require the worst-case choice of the factors. Thus in many examples realistic modeling requires complex stopping-time distributions, and if the precise parameters are unknown, then a robust analysis requires to consider the worst-case stopping-time distribution. Hence, when the stopping-time distribution is important yet unknown, a conservative estimate (i.e., lower bound) of the optimal value is obtained using the worst-case choices. Thus we consider problems that represent robust extensions of the classical finite-horizon planning.

Previous and our results. For fixed stopping time T , the expected total reward for Markov chains and MDPs can be computed via a simple dynamic programming (or backward induction) approach [Put94, Chapter 4], [FV97, How60, BKN⁺19]. The optimization problem for Markov chains and MDPs with expected stopping time has not been considered in the literature (to the best of our knowledge). Our main results are as follows:

- In contrast to the simple algorithm for fixed stopping time T , we show that quite surprisingly the expected stopping-time problem is *Positivity*-hard. The Positivity problem is known to be at least as hard as the well-known Skolem problem, whose decidability has been open for more than eight decades [OW14]. Moreover, we establish inter-reducibility between the expected stopping-time problem and the Positivity problem, and thus show that for a simple variant (adding expectation to stopping time) of the classical Markov chain problem, establishing either decidability or undecidability would be a major breakthrough.
- We then consider approximating the optimal expected total reward under the constraint that the expected stopping time is T , and show that for every additive absolute error $\varepsilon > 0$, the approximation can be achieved in time logarithmic in $1/\varepsilon$ and exponential in the size of the Markov chain.
- For MDPs we show that infinite-memory strategies are required. While the expected stopping-time problem is Positivity-hard for MDPs (since Markov chains are a special case), we show that the approximation problem can be solved in exponential space in the size of the MDP and logarithm of $1/\varepsilon$.

Comparison with related work. The optimization problem with fixed expected stopping time has been considered for the simple model of graphs [CD19], which is a model without stochastic aspects. The graph problem can be solved in polynomial time [CD19], while in sharp contrast, we show that the problem is Positivity-hard for Markov chains.

Remark 1.1. The expected stopping-time problem for Markov chains has a similar flavor as probabilistic automata (or blind MDPs) [Rab63]. In probabilistic automata a word (or letter sequence) must be provided without the information about how the probabilistic automaton executes. Similarly, for the expected stopping-time problem for Markov chains the probability distribution for stopping times must be chosen without knowing the execution of the Markov chain (in contrast to stopping criteria based on current state or accumulated reward, which rely on knowing the execution of the Markov chain). For probabilistic automata, even for basic reachability, all problems related to approximation are undecidable [MHC03]. In contrast, we show that while the exact problem for expected stopping time in Markov chains is Positivity-hard, the approximation problem can be solved in exponential time.

2. PRELIMINARIES

A *stopping-time distribution* (or simply, a distribution) is a function $\delta : \mathbb{N} \rightarrow [0, 1]$ such that $\sum_{t \in \mathbb{N}} \delta(t) = 1$. The support of δ is $\text{Supp}(\delta) = \{t \in \mathbb{N} \mid \delta(t) \neq 0\}$. We denote by Δ the set of all stopping-time distributions, and by Δ^\dagger the set of all distributions δ with $|\text{Supp}(\delta)| \leq 2$, called the *bi-Dirac* distributions.

The *expected utility* of a sequence $u = u_0, u_1, \dots$ of real numbers under a distribution δ is $\mathbb{E}_\delta(u) = \sum_{t \in \mathbb{N}} u_t \cdot \delta(t)$. In particular, the expected utility of the sequence $0, 1, 2, 3, \dots$ of all natural numbers is called the *expected time* (of distribution δ), denoted by \mathbb{E}_δ .

We recall the definition of the Positivity problem and of the related Skolem problem. In the sequel, we denote by M_{ij}^t the (i, j) entry of the t -th power of matrix M (we should write it as $(M^t)_{i,j}$, but use this simpler notation when no ambiguity can arise).

Positivity problem [OW14, AAOW15]. Given a square integer matrix M , decide whether there exists an integer $t \geq 1$ such that $M_{1,2}^t > 0$.

Skolem problem [OW14, AAOW15]. Given a square integer matrix M , decide whether there exists an integer $t \geq 1$ such that $M_{1,2}^t = 0$.

The decidability of the Positivity and Skolem problems is a longstanding open question [OW14], and there is a reduction from the Skolem problem to the Positivity problem that increases the matrix dimension quadratically [HHH06, OW14].

3. MARKOV CHAINS

We present the basic definitions related to Markov chains and the decision problems for the optimal total reward with expected stopping time.

3.1. Definitions. A *Markov chain* is a tuple $\langle M, \mu, w \rangle$ consisting of:

- an $n \times n$ stochastic matrix M (in which all entries M_{ij} are nonnegative rationals¹, and the sum $\sum_j M_{ij}$ of the elements in each row i is 1),
- an initial distribution $\mu \in ([0, 1] \cap \mathbb{Q})^n$ (viewed as $1 \times n$ row vector, and such that $\sum_i \mu_i = 1$), and
- a vector $w \in \mathbb{Q}^n$ of weights (or rewards).

We also view μ and w as functions $V \rightarrow \mathbb{Q}$ where $V = \{1, 2, \dots, n\}$ is the set of vertices of the Markov chain. We often abbreviate Markov chains as M , when μ and w are clear from the context. We denote by $\|w\| = \max_{v \in V} |w(v)|$ the largest absolute value in w .

A Markov chain induces a probability measure on sequences of vertices of a fixed length, namely $\mathbb{P}(v_0 v_1 \dots v_k) = \mu(v_0) \cdot \prod_{i=0}^{k-1} M_{v_i, v_{i+1}}$. Analogously, we denote by $\mathbb{E}(f)$ the expected value of the function $f : V^* \rightarrow \mathbb{Q}$ defined over finite sequences of vertices.

Given a stopping-time distribution $\delta : \mathbb{N} \rightarrow [0, 1]$, let N_δ be a random variable whose distribution is δ . We are interested in computing the *optimal* (worst-case) expected value (or simply the value) of Markov chains with expected stopping time T , defined by:

$$\begin{aligned} \text{val}(M, T) &= \inf_{\substack{\delta \in \Delta \\ \mathbb{E}_\delta = T}} \mathbb{E} \left[\sum_{i=0}^{N_\delta} w(v_i) \right] = \inf_{\substack{\delta \in \Delta \\ \mathbb{E}_\delta = T}} \mathbb{E} \left[\sum_{i=0}^{N_\delta} \mu \cdot M^i \cdot w^\top \right] \\ &= \inf_{\substack{\delta \in \Delta \\ \mathbb{E}_\delta = T}} \sum_{t=0}^{\infty} \delta(t) \cdot u_t, \end{aligned}$$

where w^\top is the transpose of w , and u is the sequence of utilities defined by $u_t = \sum_{i=0}^t \mu \cdot M^i \cdot w^\top$ for all $t \geq 0$. With this definition in mind, we also denote the optimal expected value of a Markov chain M by $\text{val}(u, T)$. The best-case expected value, defined using sup instead of inf in the above definition, can be computed as the opposite of the worst-case expected value for the Markov chain with all weights multiplied by -1 .

Exact value problem with expected stopping time. Given a Markov chain $\langle M, \mu, w \rangle$, a rational stopping time T , and a rational threshold θ , decide whether the optimal expected value of M with expected stopping time T is below θ , i.e., whether $\text{val}(M, T) < \theta$.

Approximation of the value with expected stopping time. We also consider an approximate version of the exact value problem, where the goal is to compute, given $\varepsilon > 0$, a value v_ε such that $|\text{val}(M, T) - v_\varepsilon| \leq \varepsilon$. We say that v_ε is an *approximation with additive error* ε of the optimal value.

3.2. Hardness of the exact value problem. This section is devoted to the proof of the following result, which establishes the inter-reducibility of the exact value problem, the Positivity problem, and the Markov Reachability problem (defined in Section 3.2.3).

Theorem 3.1. *The Positivity problem, the inequality variant of the Markov Reachability problem, and the exact value problem with expected stopping time are inter-reducible.*

¹For decidability and complexity results, we assume the numbers are rationals encoded as two binary numbers.

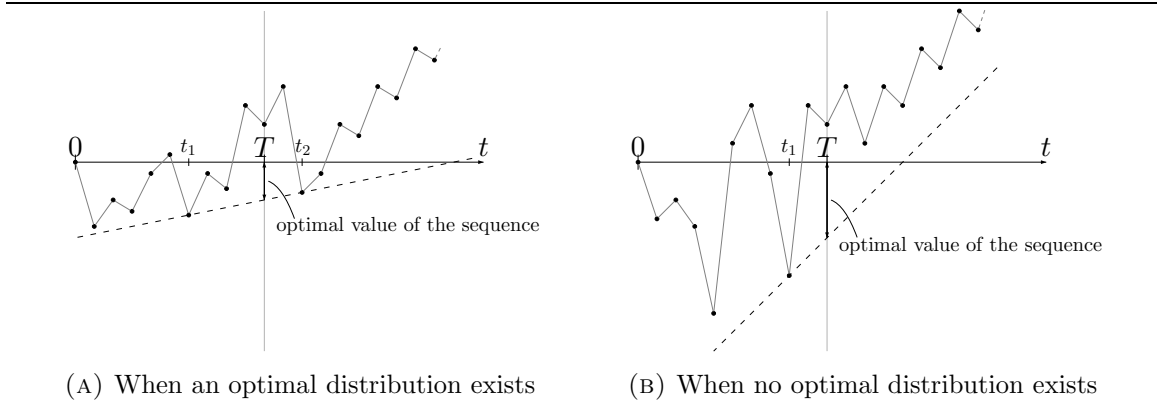


FIGURE 1. Geometric interpretation of the value of a sequence of utilities.

The decidability status of the Positivity problem is a longstanding open question, although decidability is known for dimension $n \leq 5$ [OW14, Section 4]. Therefore, constructing an algorithm to compute the exact value of a Markov chain with expected stopping time T would require the significant advances in number theory that are necessary to solve the Positivity problem [OW14, Section 5].

We also show the converse reduction from the exact value problem to the Positivity problem. Hence proving the undecidability of the exact value problem would also be a major breakthrough, as it would entail the undecidability of the Positivity problem.

The proof of Theorem 3.1 is presented in the rest of this section.

3.2.1. Geometric interpretation. A geometric interpretation for (arbitrary) sequences of real numbers and expected stopping-time was developed in previous work [CD19]. We recall the main result in this section. The rest of our technical results is independent from [CD19] (see also Comparison with related work in Section 1).

It is known that bi-Dirac distributions are sufficient for optimal expected value, namely for all sequences $u = u_0, u_1, \dots$ of utilities, for all time bounds T , the following holds [CD19]:

$$\inf\{\mathbb{E}_\delta(u) \mid \delta \in \Delta \wedge \mathbb{E}_\delta = T\} = \inf\{\mathbb{E}_\delta(u) \mid \delta \in \Delta^\uparrow \wedge \mathbb{E}_\delta = T\}.$$

Moreover the value of the expected utility of the sequence u under a bi-Dirac distribution with support $\{t_1, t_2\}$ (where $t_1 < T < t_2$) and expected time T is given by

$$u_{t_1} + \frac{T - t_1}{t_2 - t_1} \cdot (u_{t_2} - u_{t_1}). \quad (3.1)$$

As illustrated in Figure 1a, this value is obtained as the intersection of the vertical axis at T and the line that connects the two points (t_1, u_{t_1}) and (t_2, u_{t_2}) . Intuitively, the optimal value of a sequence of utilities is obtained by choosing the two points t_1 and t_2 such that the connecting line intersects the vertical axis at T as low as possible.

It is always possible to fix a value of t_1 such that it is sufficient to consider bi-Dirac distributions with support containing t_1 to compute the optimal value (because $t_1 \leq T$ is to be chosen among a finite set of points), but the optimal value of t_2 may not exist, as in Figure 1b. In that case, the value of the sequence of utilities is obtained as $t_2 \rightarrow \infty$.

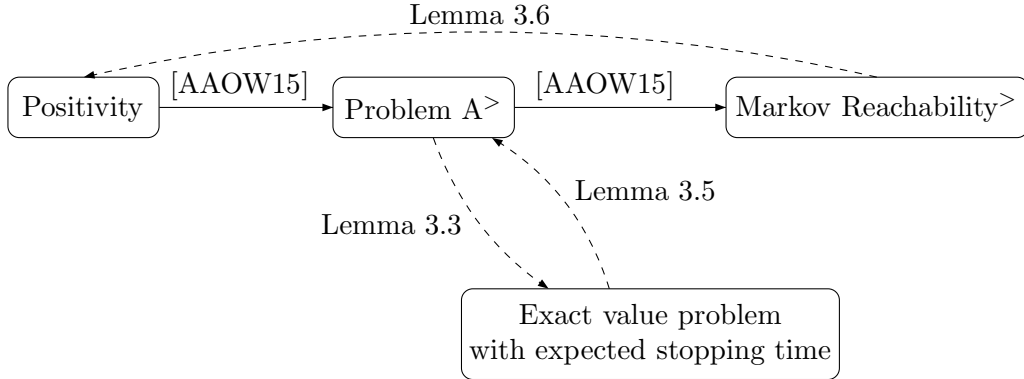


FIGURE 2. Known reductions (solid lines), and reductions established in this paper (dashed lines).

Given such a value of t_1 , let $\nu = \inf_{t_2 \geq T} \frac{u_{t_2} - u_{t_1}}{t_2 - t_1}$, and Lemma 3.2 shows that $u_t \geq f_u(t)$, for all $t \geq 0$ where $f_u(t) = u_{t_1} + (t - t_1) \cdot \nu$. The optimal expected utility is

$$\begin{aligned}
 \text{val}(u, T) &= \min_{0 \leq t_1 \leq T} \inf_{t_2 \geq T} u_{t_1} + \frac{T - t_1}{t_2 - t_1} \cdot (u_{t_2} - u_{t_1}) \\
 &= \min_{0 \leq t_1 \leq T} u_{t_1} + (T - t_1) \cdot \nu \\
 &= f_u(T),
 \end{aligned}$$

hence $f_u(T)$ is the optimal value.

Lemma 3.2 (Geometric interpretation [CD19]). *For all sequences u of utilities:*

- if $u_t \geq a \cdot t + b$ for all $t \geq 0$, then the optimal value of the sequence u is at least $a \cdot T + b$;
- we have $u_t \geq f_u(t)$ for all $t \geq 0$, and the optimal expected value of u is $f_u(T)$.

It follows from Lemma 3.2 that the optimal value of the sequence u is the largest possible value at T of a line that lies below u : $\text{val}(u, T) = \sup\{f(T) \mid \exists a, b \cdot \forall t : f(t) = a \cdot t + b \leq u_t\}$.

3.2.2. Reduction of the Positivity problem to the exact value problem. It is known that the Positivity problem can be reduced to the inequality variant of [AAOW15, Problem A], defined below as $A^>$. A subsequent reduction of $A^>$ to the exact value problem with expected stopping time establishes one direction of Theorem 3.1. We present such a reduction in the proof of Lemma 3.3 (see also Figure 2).

Problem $A^=$ [AAOW15]. Given a $n \times n$ aperiodic² stochastic matrix M with rational entries, an initial distribution $\mu = (1, 0, \dots, 0)$, and a vector $z \in \{0, 1, 2\}^n$, decide whether there exists an integer $t \geq 1$ such that $\mu \cdot M^t \cdot z^\top = 1$.

²Although in the original formulation of Problem A, the stochastic matrix M need not be aperiodic, the reduction of the Positivity problem to Problem A produces stochastic matrices that define aperiodic Markov chains (even ergodic unichains) [AAOW15].

Problem A[>] [AAOW15]. Given a $n \times n$ aperiodic stochastic matrix M with rational entries, an initial distribution $\mu = (1, 0, \dots, 0)$, and a vector $z \in \{0, 1, 2\}^n$, decide whether there exists an integer $t \geq 1$ such that $\mu \cdot M^t \cdot z^\top > 1$.

Problems A⁼ and A[>] are difficult to solve only in the case where $\mu \cdot M^t \cdot z^\top$ converges to 1 as $t \rightarrow \infty$. Otherwise, an argument based on the definition of convergence to a limit shows that the problems are decidable [KA04, Theorem 1]. Note that $\lim_{t \rightarrow \infty} \mu \cdot M^t$ exists since M is aperiodic, and the limit is the steady-state vector π , which is algorithmically computable. Hence we can assume that the instances of Problem A[>] are such that

$$\pi \cdot z^\top = 1. \quad (3.2)$$

Moreover, without loss of generality, we can modify M such that there is no incoming transition to the initial vertex 1 (remember that $\mu(1) = 1$) by creating a copy of the initial vertex, and redirecting the transitions to 1 towards the copy vertex. Thus we require the matrix M in Problem A[>] to define a Markov chain consisting of an initial vertex 1 with no incoming transition. This may however increase the dimension of the matrix by 1.

Lemma 3.3. *Problem A[>] can be reduced to the exact value problem with expected stopping time.*

Corollary 3.4. *The Positivity problem can be reduced to the exact value problem with expected stopping time.*

The proof of Lemma 3.3 is organized as follows: we first recall basic results from the theory of Markov chains, then present a reduction of Problem A[>] to the exact value problem with expected stopping time, and establish its correctness.

Basic results. First we show that, given an aperiodic Markov chain $\langle M, \mu, w \rangle$ that has a single recurrent class, there exist vectors x, y such that the expected utility after t steps tends to $\mu \cdot x^\top \cdot t + \mu \cdot y^\top$ as $t \rightarrow \infty$, formally:

$$\lim_{t \rightarrow \infty} \left| \sum_{i=0}^{t-1} \mu \cdot M^i \cdot w^\top - \mu \cdot (x^\top \cdot t + y^\top) \right| = 0. \quad (3.3)$$

The vector x is called the *gain* per time unit, and y is the *relative-gain* vector. They can be computed by solving the following equations (following [Gal13, Section 4.5]):

$$\begin{cases} x_i = \pi \cdot w^\top \text{ for all vertices } i \in V \\ y^\top = M \cdot (y - x)^\top + w^\top \\ \pi \cdot y^\top = 0 \end{cases} \quad (3.4)$$

The number $g = \pi \cdot w^\top$ is the gain per time unit. Note that $x = g \cdot e$ where $e = (1, 1, \dots, 1)$, and that $M \cdot x^\top = x^\top = g \cdot e^\top$ because M is stochastic and the sum of the elements in each of its row is 1. It follows that Equation (3.4) can be written as $y^\top = M \cdot y^\top + w^\top - x^\top$, and by $t - 1$ successive substitutions of y^\top , we get

$$\begin{aligned} y^\top &= M^t \cdot y^\top + \sum_{i=0}^{t-1} M^i \cdot w^\top - \sum_{i=0}^{t-1} M^i \cdot x^\top \\ &= M^t \cdot y^\top + \sum_{i=0}^{t-1} M^i \cdot w^\top - t \cdot x^\top \end{aligned}$$

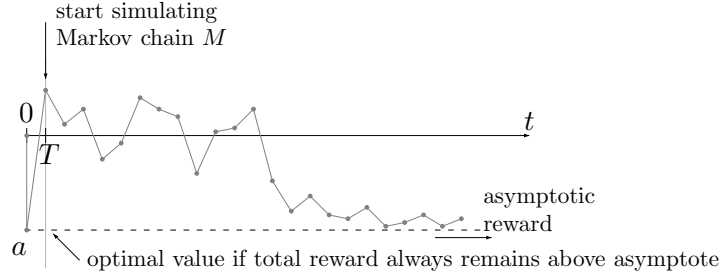


FIGURE 3. Reduction of the Positivity problem to the exact value problem.

Then, the rate of convergence of the expected utility evaluates as follows, for all $t \geq 1$:

$$\begin{aligned}
 & \sum_{i=0}^{t-1} \mu \cdot M^i \cdot w^\top - \mu \cdot (x^\top \cdot t + y^\top) \\
 &= \sum_{i=0}^{t-1} \mu \cdot M^i \cdot w^\top - \mu \cdot x^\top \cdot t \\
 & \quad - \mu \cdot M^t \cdot y^\top - \mu \cdot \sum_{i=0}^{t-1} M^i \cdot w^\top + t \cdot \mu \cdot x^\top \\
 &= -\mu \cdot M^t \cdot y^\top
 \end{aligned} \tag{3.5}$$

which tends to $-\pi \cdot y^\top = 0$ as $t \rightarrow \infty$, establishing (3.3).

In the case of an aperiodic Markov chain with multiple recurrent classes, the gain and relative gain satisfying Equation (3.3) can be computed as the linear combination of the vectors x, y obtained for each recurrent class, where the coefficient in the linear combination is the mass of probability that reaches (in the limit) the recurrent class from the initial distribution μ .

Reduction. The reduction from Problem $A^>$ to the exact value problem with expected stopping time is as follows. Given an instance (M, μ, z) of Problem $A^>$, we construct an instance of the exact value problem in two stages. First, let $w^\top = z^\top - M \cdot z^\top$ be a reward vector defining a Markov chain $\langle M, \mu, w \rangle$. We explain later why w is defined in this way.

We proceed to the second stage of the construction, and define the instance of the exact value problem, namely the Markov chain $\langle M', \mu', w' \rangle$, the expected time T , and the threshold θ . The key idea of the construction is illustrated in Figure 3. Given the Markov chain $\langle M, \mu, w \rangle$, we can compute its asymptotic expected utility, shown as the dashed line $\mu \cdot (x^\top \cdot t + y^\top)$ in Figure 3 which also plots the sequence u_{t-1} for $t \geq 1$. Note that by Equation (3.3) we have $\lim_{t \rightarrow \infty} |u_{t-1} - \mu \cdot (x^\top \cdot t + y^\top)| = 0$ and by Equations (3.4) $x = \pi \cdot z^\top - \pi \cdot M \cdot z^\top = 0$.

We construct an instance of the exact value problem in such a way that, if the utility of M always remains above its asymptote, then the optimal value is the value of the asymptote at time T , and otherwise, the optimal value is strictly smaller. We achieve this by having an initial vertex with weight a such that, if the Markov chain $\langle M, \mu, w \rangle$ is executed (simulated)

after the initial vertex, then the weight a lies exactly on the asymptote of $\langle M, \mu, w \rangle$ (see Figure 3 and the geometric interpretation in Section 3.2.1). Since we simulate $\langle M, \mu, w \rangle$ after one time step, the value of a is chosen such that the point $(0, a)$ belongs to the line $\mu \cdot (x^\top \cdot t + y^\top)$. Since $\mu = (1, 0, \dots, 0)$ in Problem A[>], we have $a = y(1)$. To recover the original behavior of the Markov chain $\langle M, \mu, w \rangle$, we subtract a from the weight of the initial vertex of M , thus $w'(1) = w(1) - a$. As we assumed that the initial vertex in M has no incoming transition, it is never re-visited later. We take $T = 1$ and the value of the asymptote at time T is $\mu \cdot (x^\top + y^\top) = x(1) + y(1) = a$, which we define as the threshold θ of the exact value problem, thus $\theta = a$.

Formally, the instance of the exact value problem is defined as follows:

$$w' = \begin{pmatrix} a \\ w(1) - a \\ w(2) \\ \vdots \\ w(n) \end{pmatrix}, \quad M' = \begin{pmatrix} 0 & \mu \\ 0 & M \end{pmatrix}, \quad \begin{array}{l} \mu' = (1, 0, \dots, 0), \\ T = 1, \\ \theta = a \end{array}$$

where $a = y(1)$ and y is the relative-gain vector of the Markov chain $\langle M, \mu, w \rangle$. Note that the initial vertex of M has no incoming transition (in M), and thus the sequence of expected utilities in M' indeed simulates the sequence of expected utilities in M , and the asymptotic expected utilities as well as the steady-state vectors of $\langle M, \mu, w \rangle$ and $\langle M', \mu', w' \rangle$ coincide.

Correctness of the reduction. To establish the correctness of the reduction, we show the following equivalences:

- (1) the optimal expected value of M' with expected stopping time T is smaller than θ (i.e., the answer to the exact value problem is YES) if and only if the utility sequence of M eventually drops below its asymptote;
- (2) the utility sequence of M eventually drops below its asymptote if and only if $\mu \cdot M^t \cdot z^\top > 1$ for some $t \geq 1$ (i.e., the answer to Problem A[>] is YES).

To show the first equivalence, consider the first direction and assume that the value of M' is smaller than θ . Given that the line $\mu \cdot (x^\top \cdot t + y^\top)$ has value θ at $t = T$, it follows from Lemma 3.2 that the utility sequence of M' does not always remain above that line, and thus the utility sequence of M eventually drops below its asymptote.

Now consider the second direction of the first equivalence and assume that the utility sequence of M eventually drops below its asymptote. Then the utility sequence of M' drops below the line $\mu \cdot (x^\top \cdot t + y^\top)$, say at time $t_2 \geq 1$. We construct a distribution δ with $\mathbb{E}_\delta = T$ such that the value of the expected reward under δ is less than $\mu \cdot (x^\top \cdot T + y^\top) = \theta$ (which implies that the optimal value, obtained as the infimum over all distributions, is also below θ).

We consider two cases: (1) if $t_2 = 1$ (i.e., $t_2 = T$), consider the distribution δ such that $\delta(t_2) = 1$ (note that $\mathbb{E}_\delta = T$) and the result follows immediately; (2) otherwise, $t_2 > 1$ and consider the bi-Dirac distribution with support $\{t_1, t_2\}$ where $t_1 = 0$. Note that $t_1 < T < t_2$ and the value of the expected reward under this distribution is given by the value at time T of the line connecting the point (t_1, a) and a point below the asymptote (at t_2), see Equation (3.1). This value is below the value θ of the asymptote at time T since (t_1, a) is on the asymptote, and the other point (at t_2) is strictly below the asymptote.

To show the second equivalence, note that by Equation (3.5) the utility sequence of M eventually drops below its asymptote if and only if $-\mu \cdot M^t \cdot y^\top < 0$ for some $t \geq 1$. Hence

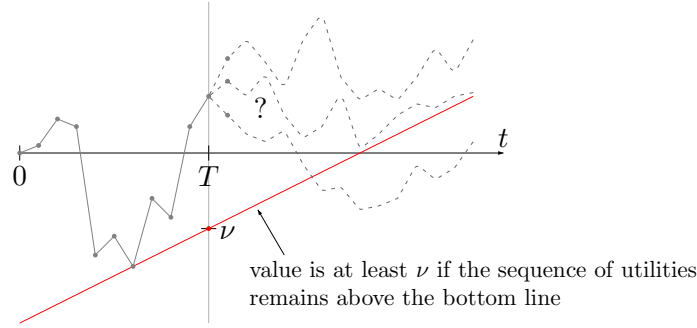


FIGURE 4. Solving the exact value problem using an oracle for Problem $A^>$.

we can establish the second equivalence by showing that $-\mu \cdot M^t \cdot y^\top < 0$ if and only if $\mu \cdot M^t \cdot z^\top > 1$. This is where the value of w is important. The result holds if $y = z - e$, and we just need to show that $y = z - e$ satisfies Equations (3.4), namely that

$$\begin{aligned} (z - e)^\top &= M \cdot (z - e - x)^\top + w^\top \\ \pi \cdot (z - e)^\top &= 0 \end{aligned}$$

that is

$$\begin{aligned} (z - e)^\top &= M \cdot z^\top - e^\top - x^\top + z^\top - M \cdot z^\top \\ \pi \cdot z^\top - \pi \cdot e^\top &= 0 \end{aligned}$$

which hold since $x = 0$ and $\pi \cdot z^\top = 1 = \pi \cdot e^\top$ (Equation (3.2)). This concludes the proof of Lemma 3.3.

Using the reduction of the Positivity problem to Problem $A^>$ [AAOW15], we obtain Corollary 3.4, showing that a decidability result for the exact value problem would imply the decidability of the Positivity problem, which is a longstanding open question.

3.2.3. Reduction of the exact value problem to the Positivity problem. We present the converse reduction of Section 3.2.2, showing that to potentially prove the exact value problem is undecidable would require such a proof for the Positivity problem as well. We sketch the reduction by showing how the exact value problem can be solved using an oracle for Problem $A^>$, illustrated in Figure 4, and then present a reduction of Problem $A^>$ to the Positivity problem.

Lemma 3.5. *The exact value problem with expected stopping time can be reduced to Problem $A^>$.*

Proof. Given a Markov chain $\langle M, \mu, w \rangle$ with expected stopping time T and threshold θ , we solve the exact value problem using an oracle for Problem $A^>$ as follows. First, if $u_T < \theta$ then the answer to the exact value problem is YES. Otherwise, we compute the value of utilities $u_t = \sum_{i=0}^t \mu \cdot M^i \cdot w^\top$ for all $0 \leq t < T$, and let $b = \max_{0 \leq t < T} \frac{u_t - \theta}{t - T}$. Consider the *bottom line* of equation $b \cdot (t - T) + \theta$ and observe that $u_t \geq b \cdot (t - T) + \theta$ for all $0 \leq t \leq T$ (see Figure 4). By the geometric interpretation lemma (Lemma 3.2), it suffices to determine

whether the sequence of utilities ever drops below the bottom line to answer the exact value problem.

For simplicity of presentation, we assume that the recurrent classes of the Markov chain are all aperiodic. The case of periodic recurrent classes (say C_1, C_2, \dots, C_l with respective period d_1, \dots, d_l) can be solved analogously by applying the procedure to each initial distribution $\mu, \mu \cdot M, \mu \cdot M^2, \dots, \mu \cdot M^{d-1}$, and the transition matrix M^d where d is the period of the Markov chain, i.e., $d = \text{lcm}\{d_1, \dots, d_l\}$.

Using Equations (3.4), we can compute the *asymptote* of the sequence of utilities as $\mu \cdot x^\top \cdot (t+1) + \mu \cdot y^\top$. Comparing the slope of the bottom line and the slope of the asymptote, we have the following cases:

- if $b < \mu \cdot x^\top$, then from some point on the utilities always remain above the bottom line. Such a point can be computed using the convergence rate of Markov chains (Appendix B). Then the answer to the exact value problem is NO if $u_t \geq b \cdot (t - T) + \theta$ for all (the finitely many) values of t up to that point. Otherwise the answer is YES.
- if $b > \mu \cdot x^\top$, then eventually the utility gets below the bottom line, thus the answer to the exact value problem is YES.
- if $b = \mu \cdot x^\top$, then either the bottom line is different from the asymptote, and we can reuse one of the above cases, or the bottom line is equal to the asymptote, and the condition for the sequence of utilities to eventually drop below the bottom line is that $-\mu \cdot M^t \cdot y^\top < 0$ for some $t \geq 1$ (Equation (3.5)), which is equivalent to $\mu \cdot M^t \cdot (e + y)^\top > 1$ for some $t \geq 1$ (where $e = (1, 1, \dots, 1)$). We cannot immediately use an oracle for Problem $A^>$ with $z = e + y$ because Problem $A^>$ requires $z \in \{0, 1, 2\}^n$. However, we can obtain a valid instance of Problem $A^>$ as follows. First we apply a scaling factor to y to ensure that its components are in the interval $[-1, 1]$. Scaling the vector y does not affect the answer to the original question (whether $-\mu \cdot M^t \cdot y^\top < 0$). Then we construct a Markov chain $\langle M', \mu', w' \rangle$ as follows: for each vertex $v \in V$ of the Markov chain M , we create a copy v^0 and $w'(v^0) = 0$. Define $w'(v) = 1$ if $y(v) \geq 0$ and $w'(v) = -1$ if $y(v) < 0$, where $y(v)$ is the component in y corresponding to vertex v . For each transition $v \xrightarrow{p} u$ (i.e., $M_{v,u} = p$), we have the transitions $v \xrightarrow{p \cdot |y(u)|} u$ and $v \xrightarrow{p \cdot (1 - |y(u)|)} u^0$ in M' . The outgoing transitions from the copy v^0 are the same as from v . Note that the weights in w' are in the set $\{-1, 0, 1\}$, and therefore the vector $z = e + w'$ is in $\{0, 1, 2\}^{2n}$. We now call the oracle for Problem $A^>$ with the Markov chain M' and $z = e + w'$ which gives the answer to the exact value problem (in this way we transferred the value of y into the transition probabilities of M' , and note that the dimension of M' is twice the dimension of M). \square

To obtain the inter-reducibility result of Theorem 3.1, we need to show that Problem $A^>$ can be reduced to the Positivity problem, which we establish by showing that the inequality version of the Markov reachability problem (defined below) can be reduced to the Positivity problem, as it is known that Problem $A^>$ can be reduced to the inequality variant of the Markov reachability problem [AAOW15] (see also Figure 2). This is a straightforward result established in Lemma 3.6.

Markov reachability[>] problem [AAOW15]. Given a square stochastic matrix M with rational entries and a rational number $r > 0$, decide whether there exists an integer $t \geq 1$ such that $M_{1,2}^t > r$.

Markov reachability⁼ problem [AAOW15]. Given a square stochastic matrix M with rational entries and a rational number $r > 0$, decide whether there exists an integer $t \geq 1$ such that $M_{1,2}^t = r$.

Lemma 3.6. *The Markov reachability[>] problem can be reduced to the Positivity problem.*

Proof. The reduction of the Markov reachability[>] problem to the Positivity problem is as follows. Given a $n \times n$ stochastic matrix M and a rational number $r > 0$, define the $(n+3) \times (n+3)$ matrix P as follows, where $0_{n \times 1}$ is the zero column-vector of length n , and $e_1 = (1, 0, 0, \dots, 0)$ and $e_2 = (0, 1, 0, \dots, 0)$ (so that $e_1 \cdot M \cdot e_2^\top = M_{1,2}$):

$$P = \begin{pmatrix} M & 0_{n \times 1} & 0_{n \times 1} & e_2^\top \\ e_1 \cdot M & 0 & -r & 0 \\ 0_{1 \times n} & 0 & 1 & 1 \\ 0_{1 \times n} & 0 & 0 & 0 \end{pmatrix}$$

All entries of P are rational numbers, hence there exists $d \in \mathbb{N}$ such that $\tilde{P} = d \cdot P$ is an integer matrix. It is easy to show by induction that for all $t \geq 2$:

$$\tilde{P}^t = d^t \cdot \begin{pmatrix} M^t & 0_{n \times 1} & 0_{n \times 1} & M^{t-1} \cdot e_2^\top \\ e_1 \cdot M^t & 0 & -r & e_1 \cdot M^{t-1} \cdot e_2^\top - r \\ 0_{1 \times n} & 0 & 1 & 1 \\ 0_{1 \times n} & 0 & 0 & 0 \end{pmatrix}$$

It follows that $\tilde{P}_{n+1, n+3}^t = d^t \cdot (e_1 \cdot M^{t-1} \cdot e_2^\top - r) = d^t \cdot (M_{1,2}^{t-1} - r)$, and therefore $\tilde{P}_{n+1, n+3}^t > 0$ if and only if $M_{1,2}^{t-1} > r$ (for all $t \geq 2$). Since $\tilde{P}_{n+1, n+3} = 0$, we conclude that there exists an integer $t \geq 1$ such that $\tilde{P}_{n+1, n+3}^t > 0$ if and only if there exists an integer $t \geq 1$ such that $M_{1,2}^t > r$, which concludes the reduction by rearranging the order of the rows and columns of \tilde{P} , and the desired result follows. \square

The results of Lemma 3.3, 3.5 and 3.6 establish Theorem 3.1. The reduction in Lemma 3.6 can easily be adapted to show that the Markov reachability⁼ problem can be reduced to the Skolem problem and thus these problems are inter-reducible with Problem A⁼ (the adaptation is to replace $\tilde{P}_{n+1, n+3}$ by an arbitrary nonzero value, which has no effect on the value of the powers of \tilde{P}).

Theorem 3.7. *The Skolem problem, Problem A⁼, and the Markov reachability⁼ problem are inter-reducible.*

3.3. Approximation of the optimal value. We can compute an approximation of the optimal value with additive error by considering an approximation u' of the exact sequence u of expected utilities of the Markov chain as follows: for a large number of time steps, let the approximate sequence u' be equal to u , and then from some point on it switches to the value of the limit (asymptotic, and possibly periodic) sequence of expected utilities at the steady-state distribution(s). By taking the switching point large enough, the approximation sequence u' can be made arbitrarily close to the exact sequence u . We show that the value of the sequences u' approximates arbitrarily closely the (exact) optimal value of u .

By the results of Section 3.2.1, the optimal expected value of any sequence u' of utilities is given by the expression

$$val(u', T) = \min_{0 \leq t_1 \leq T} \inf_{t_2 \geq T} \frac{u'_{t_1}(t_2 - T) + u'_{t_2}(T - t_1)}{t_2 - t_1}. \quad (3.6)$$

We can effectively compute the value of $val(u', T)$ when u' is an ultimately periodic sequence, i.e. $u' = A.C^\omega$ where A, C are finite sequences (with C nonempty): we show in Lemma 3.8 that the infinite range of t_2 in the expression (3.6) can be replaced by a finite range, because the optimal value is obtained either by taking t_2 before the first repetition of the cycle C , or by taking $t_2 \rightarrow \infty$ (i.e., if repeating the cycle once improves the value, then repeating the cycle infinitely often improves the value even more). Let S_A and S_C be the sum of the weights in A and C respectively, let $M_C = \frac{S_C}{|C|}$ be the average weight of the cycle C .

Lemma 3.8. *The optimal value of an ultimately periodic sequence $u = A.C^\omega$ is $val(u, T) = \min\{E_1, E_2\}$ where*

$$E_1 = \min_{0 \leq t_1 \leq T} \min_{T \leq t_2 \leq |A| + |C|} \frac{u_{t_1}(t_2 - T) + u_{t_2}(T - t_1)}{t_2 - t_1}, \text{ and}$$

$$E_2 = \min_{0 \leq t_1 \leq T} u_{t_1} + M_C \cdot (T - t_1).$$

If $T \geq |A| + |C|$, then $val(u, T) = \min_{0 \leq t_1 \leq |A| + |C|} u_{t_1} + M_C \cdot (T - t_1)$.

Proof. The expression E_1 is the expression (3.6) where the range of t_2 is the interval $[T, |A| + |C|]$. We now show that the expression E_2 corresponds to $t_2 \geq |A|$ (assuming always $t_2 \geq T$), which covers $t_2 \geq |A| + |C|$. For all $t_2 \geq |A|$, we can express t_2 as $t_0 + k \cdot |C|$ where $|A| \leq t_0 \leq |A| + |C|$ and $k \geq 0$. We have two cases:

- If $T \leq |A| + |C|$, then the expression (3.6) gives $val(u, T) =$

$$\min_{0 \leq t_1 \leq T} \min_{\max(T, |A|) \leq t_0 \leq |A| + |C|} \inf_{k \geq 0} \frac{u_{t_1}(t_0 + k \cdot |C| - T) + (u_{t_0} + k \cdot S_C)(T - t_1)}{t_0 + k \cdot |C| - t_1}$$

where the numerator and denominator of the fraction are linear in k . Such functions $\frac{a \cdot k + b}{c \cdot k + d}$ are monotone (their first derivative has constant sign), and note that the denominator $t_0 + k \cdot |C| - t_1$ is nonzero for all $k \geq 0$. It follows that the infimum over $k \geq 0$ is obtained either for $k = 0$, which is covered by the expression E_1 , or for $k \rightarrow \infty$, which gives the expression $\frac{u_{t_1} \cdot |C| + S_C \cdot (T - t_1)}{|C|} = u_{t_1} + M_C \cdot (T - t_1)$, corresponding to E_2 .

- If $T \geq |A| + |C|$, then the same reasoning as above shows that $val(u, T) = \min_{0 \leq t_1 \leq T} u_{t_1} + M_C \cdot (T - t_1)$, as the range of t_2 in expression E_1 is empty. Now observe that if $t_1 = t_0 + k \cdot |C|$ where $|A| \leq t_0 \leq |A| + |C|$, then $u_{t_1} + M_C \cdot (T - t_1) = u_{t_0} + M_C \cdot (T - t_0) + k \cdot S_C - M_C \cdot k \cdot |C| = u_{t_0} + M_C \cdot (T - t_0)$ and thus we can restrict the range of t_1 to $[0, |A| + |C|]$. The result follows. \square

We show that for a sequence u' of utilities that approximates the sequence u , the value of u' approximates the value of u and the error can be bounded. Precisely, if the weights in a Markov chain are shifted by at most η , then the optimal expected value of the Markov chain with expected stopping time T is shifted by at most $\eta \cdot (T + 1)$. Consider w' such that $|w'(v) - w(v)| \leq \eta$ for all vertices $v \in V$, and consider the sequences u and u' of utilities of a path according to w and w' respectively. Then we have $|u'_t - u_t| \leq (t + 1) \cdot \eta$ for all $t \geq 0$,

and for all distributions δ with $\mathbb{E}_\delta = T$:

$$\begin{aligned} \left| \sum_i \delta(i) \cdot u'_i - \sum_i \delta(i) \cdot u_i \right| &\leq \sum_i \delta(i) \cdot |u'_i - u_i| \\ &\leq \sum_i \delta(i) \cdot (i+1) \cdot \eta \\ &= (T+1) \cdot \eta. \end{aligned}$$

It follows that $|\text{val}(u', T) - \text{val}(u, T)| \leq (T+1) \cdot \eta$, that is the value of the sequence is shifted by at most $(T+1) \cdot \eta$ (it is easy to see that if $\forall \delta : |f(\delta) - g(\delta)| \leq K$, then $|\inf_\delta f - \inf_\delta g| \leq K$).

Lemma 3.9. *Given $\eta \geq 0$ and two sequences u and u' of utilities such that $|u'_t - u_t| \leq (t+1) \cdot \eta$ for all $t \geq 0$, we have $|\text{val}(u', T) - \text{val}(u, T)| \leq (T+1) \cdot \eta$. Analogously, if $u'_t = u_t + (t+1) \cdot \eta$ for all $t \geq 0$, then $\text{val}(u', T) = \text{val}(u, T) + (T+1) \cdot \eta$.*

We recall a result about Markov chains, which states that for Markov chains with only aperiodic recurrent classes, the vector $\mu \cdot M^t$ converges to a steady-state vector π , and the rate of convergence is bounded by an exponential in n [Gal13, Theorem 4.3.7] (see Appendix B for detailed computation). For all $j \in V$:

$$|(\mu \cdot M^t)_j - \pi_j| \leq K_1 \cdot K_2^t$$

where K_1, K_2 are constants with $K_2 < 1$, namely $K_2 = (1 - \alpha^{n^2})^{1/3n^2}$ where α is the smallest non-zero probability in M (i.e., $\alpha = \min\{M_{ij} \mid M_{ij} > 0\}$) and n is the number of vertices of M .

For general Markov chains (with possibly periodic recurrent classes), we adapt the above result as follows. Consider the set \mathcal{T} of transient vertices, each recurrent class C_1, C_2, \dots, C_l with their respective period d_1, d_2, \dots, d_l , and let $d = \text{lcm}\{d_1, \dots, d_l\}$ be their least common multiple. Note that $d_i \leq n$ for all $1 \leq i \leq l$ and d is at most the product of all prime numbers smaller than n , thus at most exponential in n [Erd89]. Then M^d can be viewed as the transition matrix of a Markov chain with aperiodic recurrent classes, and thus $\mu \cdot M^{d \cdot t}$ converges to a steady-state vector π as $t \rightarrow \infty$. Considering a recurrent class C_i , and the vertices $j \in C_i \cup \mathcal{T}$ the rate of convergence can be bounded as follows, where α^{d_i} is a lower bound on the smallest non-zero probability in M^{d_i} :

$$\begin{aligned} |(\mu \cdot M^{d \cdot t})_j - \pi_j| &= |(\mu \cdot (M^{d_i})^{\frac{d \cdot t}{d_i}})_j - \pi_j| \\ &\leq K_1 \cdot (1 - \alpha^{d_i \cdot n^2})^{\frac{d \cdot t}{d_i \cdot 3n^2}} \\ &\leq K_1 \cdot (1 - \alpha^{n^3})^{\frac{t}{3n^2}}, \end{aligned}$$

which is independent of i , and thus holds for all $j \in V$. Let $K_3 = (1 - \alpha^{n^3})^{1/3n^2}$.

It follows that $|\mu \cdot M^{d \cdot t} \cdot w^\top - \pi \cdot w^\top| \leq n \cdot W \cdot K_1 \cdot K_3^t$ where $W = \|w\|$ is the largest absolute weight in w . Then for all $\varepsilon > 0$, for all $t \geq \frac{\ln(\frac{n \cdot W \cdot K_1}{\varepsilon})}{\ln(K_3)} =: B$, we have $|\mu \cdot M^{d \cdot t} \cdot w^\top - \pi \cdot w^\top| \leq \varepsilon$, and by the same reasoning with initial distributions $\mu \cdot M, \mu \cdot M^2, \dots, \mu \cdot M^{d-1}$ we get $|\mu \cdot M^{d \cdot t + k} \cdot w^\top - \pi \cdot M^k \cdot w^\top| \leq \varepsilon$ for all $0 \leq k < d$.

Consider the sequence u' defined by

$$u'_t = \begin{cases} u_t & \text{for all } t \leq d \cdot B \\ u_{d \cdot B} + \sum_{k=d \cdot B+2}^t \pi \cdot M^{k \% d} \cdot w^\top & \text{for all } t > d \cdot B \end{cases}$$

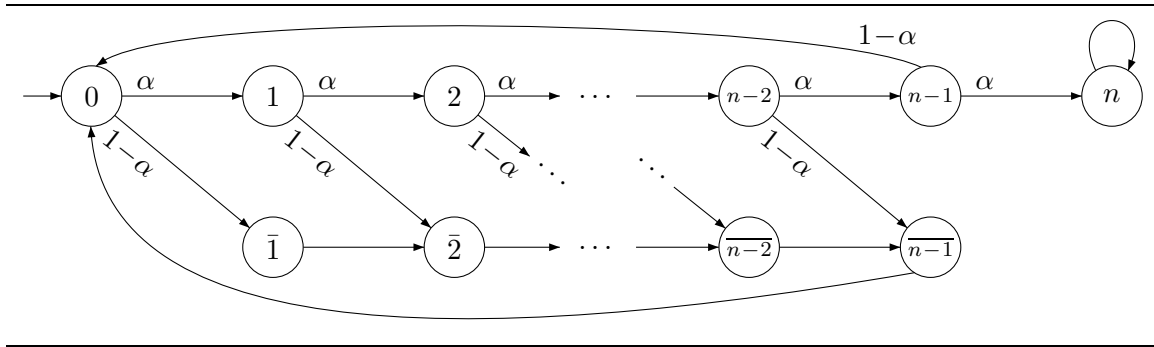


FIGURE 5. A family of Markov chains for Proposition 3.11.

where $k \% d$ is the remainder of the division of k by d . Intuitively, u'_t approximates u_t after time $t = d \cdot B$ by considering the (expected) weight at time t to be given by the limit (expected) weight at the steady-state distribution.

Then $|u'_t - u_t| \leq (t + 1) \cdot \varepsilon$ for all $t \geq 0$, and therefore $|\text{val}(u', T) - \text{val}(u, T)| \leq \varepsilon \cdot (T + 1)$ (by Lemma 3.9). The sequence u' is an ultimately periodic sequence of the form $A.C^\omega$ where $|A| = d \cdot B$ and $|C| = d$. Hence the optimal value of u' is given by Lemma 3.8 and can be obtained by computing the first $d \cdot B + d$ terms of the sequence u' , the steady-state vector π , the number d , and the average weight $M_C = \frac{S_C}{|C|}$ where $S_C = \sum_{i=0}^{d-1} \pi \cdot M^i \cdot w^\top$. This provides a way to compute an approximation with additive error ε of the optimal value of a Markov chain in time $O(P(n) \cdot T \cdot B \cdot d)$ where $P(n)$ is a polynomial in the size of the Markov chain (that accounts for matrix multiplication, steady-state vector computation, etc.). Using the fact that $(1 - \frac{1}{x})^x \in O(1)$, and that $\ln(1 - \frac{1}{x}) \in O(-1/x)$ (by Lemma A.1 in Appendix), we obtain the bounds in Theorem 3.10 in the special cases where α or n is constant.

Theorem 3.10. *The optimal expected value of a Markov chain with expected stopping time T can be computed to an arbitrary level of precision $\varepsilon > 0$, in time*

$$O\left(P(n) \cdot T \cdot \frac{\ln\left(\frac{\varepsilon}{n \cdot W}\right)}{\ln(K_3)} \cdot 2^{O(n)}\right)$$

where $K_3 = (1 - \alpha^{n^3})^{1/3n^2}$ and $P(\cdot)$ is a polynomial.

If α (the smallest non-zero probability) is constant, then the computation time is in

$$O\left(\frac{P(n) \cdot 2^{O(n)}}{\alpha^{n^3}} \cdot T \cdot \ln\left(\frac{n \cdot W}{\varepsilon}\right)\right) \text{ (as } n \rightarrow \infty\text{)}.$$

If n (the number of vertices) is constant, then the computation time is in

$$O\left(\frac{1}{\alpha^{O(1)}} \cdot T \cdot \ln\left(\frac{W}{\varepsilon}\right)\right) \text{ (as } \alpha \rightarrow 0\text{)}.$$

We present a lower bound on the execution time of the approximation algorithm of Theorem 3.10: we show that the algorithm runs in time exponential in the number of vertices of the Markov chains presented in Figure 5. This shows that the complexity analysis of our algorithm cannot be improved to eliminate the exponential dependency in the number of

vertices. However, whether there exists a polynomial-time algorithm for the approximation problem is an open question.

Figure 5 shows a family of Markov chains M with $2n$ vertices and transition probabilities parameterized by $\alpha \leq \frac{1}{2}$, with initial distribution μ that assigns probability 1 to vertex 0. The steady-state distribution assigns probability 1 to vertex n . It is easy to show that from vertex 0 the probability to reach vertex n after $t \cdot n$ steps is $1 - (1 - \alpha^n)^t$. Therefore the distance between the distribution $\mu \cdot M^t$ at time t and the steady-state vector is at least $(1 - \alpha^n)^{\frac{t}{n}}$. It follows that, to ensure that this distance is less than ε , the algorithm needs to compute the sequence of utilities of the Markov chain up to time at least $\frac{n \cdot \ln(\varepsilon)}{\ln(1 - \alpha^n)} \geq \frac{n \cdot \ln(1/\varepsilon)}{\alpha^n}$ (using Lemma A.1 in Appendix to get $\frac{1}{\ln(1 - \frac{1}{x})} \geq -x$ for all $x > 1$).

Proposition 3.11. *There exists a family of aperiodic Markov chains $M(n, \alpha)$ with $2n$ vertices ($n \in \mathbb{N}$) and smallest positive transition probability α ($\alpha \leq \frac{1}{2}$) such that, for the initial distribution $\mu = (1, 0, \dots, 0)$, we have*

$$\max_j |(\mu \cdot M(n, \alpha)^t)_j - \pi_j| \geq (1 - \alpha^n)^{\frac{t}{n}},$$

where π is the steady-state vector of $M(n, \alpha)$, and the computation time of the approximation algorithm (of Theorem 3.10) for $M(n, \alpha)$ is at least

$$\frac{n \cdot \ln(1/\varepsilon)}{\alpha^n}.$$

4. MARKOV DECISION PROCESSES

Markov decision processes (MDPs) extend Markov chains with transition choices determined by control actions. We give the basic definitions of MDPs and of the optimal expected value of an MDP with expected stopping time T .

4.1. Definitions. A *Markov decision process* is a tuple $\mathcal{M} = \langle V, A, \theta, \mu, w \rangle$ consisting of:

- a finite set V of vertices and a finite set A of actions,
- a transition function $\theta : V \times A \rightarrow (V \rightarrow [0, 1])$ such that $\theta(v, a)$ is a probability distribution over V , that is $\sum_{v' \in V} \theta(v, a)(v') = 1$ for all $v \in V$ and $a \in A$.
- $\mu : V \rightarrow [0, 1]$ is an initial distribution and $w : V \rightarrow \mathbb{Q}$ is a vector of weights, as in Markov chains.

Given a vertex $v \in V$ and a set $U \subseteq V$, let $A_U(v)$ be the set of all actions $a \in A$ such that $\text{Supp}(\theta(v, a)) \subseteq U$. A *closed* set in an MDP is a set $U \subseteq V$ such that $A_U(v) \neq \emptyset$ for all $v \in U$. A set $U \subseteq V$ is an *end-component* [dA97, BK08] if (i) U is closed, and (ii) the graph (U, E_U) is strongly connected where $E_U = \{(v, v') \in U \times U \mid \theta(v, a)(v') > 0 \text{ for some } a \in A_U(v)\}$ denote the set of edges given the actions. In the sequel, end-components should be considered *maximal*, that is such that no strict superset is an end-component.

A *strategy* in \mathcal{M} is a function $\sigma : V^+ \rightarrow (A \rightarrow [0, 1])$ such that $\sigma(\rho)$ is a probability distribution over A , for all sequences $\rho \in V^+$. A strategy σ is *pure* if for all $\rho \in V^+$, there exists an action $a \in A$ such that $\sigma(\rho)(a) = 1$; σ is *memoryless* if $\sigma(\rho v) = \sigma(\rho' v)$ for all $\rho, \rho' \in V^*$ and $v \in V$; σ uses *finite memory* if there exists a right congruence \approx over V^+ (i.e., if $\rho \approx \rho'$, then $\rho \cdot v \approx \rho' \cdot v$ for all $\rho, \rho' \in V^+$ and $v \in V$) of finite index such that $\rho \approx \rho'$ implies $\sigma(\rho) = \sigma(\rho')$.

Given the initial distribution μ , and a strategy σ , a probability can be assigned to every finite path $\rho = v_0 \cdots v_n$ as follows:

$$\mathbb{P}_\mu^\sigma(v_0 v_1 \dots v_k) = \mu(v_0) \cdot \prod_{i=0}^{k-1} \sum_{a \in A} \sigma(v_0 \cdots v_i)(a) \cdot \theta(v_i, a)(v_{i+1}).$$

Analogously, we denote by $\mathbb{E}_\mu^\sigma(f)$ the expected value of the function $f : V^* \rightarrow \mathbb{Q}$ defined over finite sequences of vertices. Let $u_t = \mathbb{E}_\mu^\sigma(\sum_{i=0}^t w(v_i))$ and define the optimal expected value of \mathcal{M} with expected stopping time $T \in \mathbb{Q}$ as follows:

$$\text{val}(\mathcal{M}, T) = \sup_{\sigma} \inf_{\substack{\delta \in \Delta \\ \mathbb{E}_\delta = T}} \sum_{t=0}^{\infty} \delta(t) \cdot u_t.$$

The strategy σ is ε -optimal if the sequence $u = (u_t)_{t \in \mathbb{N}}$ it induces is such that $\text{val}(u, T) \geq \text{val}(\mathcal{M}, T) - \varepsilon$. For $\varepsilon = 0$, we simply say that σ is *optimal* (instead of 0-optimal).

For an arbitrary strategy σ , with probability 1 the set of states visited infinitely often along an (infinite) path is an end-component [CY95, dA97]. Let the limit-probability of a (maximal) end-component U be the probability that the set of states visited infinitely often along a path is a subset of U . A *limit distribution* under σ is a distribution δ^* such that, for every end-component U , the limit-probability of U is $\sum_{v \in U} \delta^*(v)$.

4.2. Infinite memory is necessary. Since MDPs are an extension of Markov chains, the problem of computing the optimal expected value $\text{val}(\mathcal{M}, T)$ is Positivity-hard (by Corollary 3.4). Another source of hardness for this problem is that infinite memory is required for optimal strategies, as illustrated in the following example.

Example. We show in Figure 6 an MDP where infinite memory is required for optimal expected value. The only strategic choice is in vertex v'_1 (we omit the actions in the figure, and all weights not shown are 0). In particular, the upper part $\{v_1, \dots, v_6\}$ is a Markov chain and after $3k + 2$ steps, the probability mass in v_4 is $p_k = \frac{1}{3} \cdot (1 - \frac{1}{2^{k+1}})$. For instance $p_0 = \frac{1}{6}$. Note that one step before, the probability mass in v_1 is $\frac{1}{3} \cdot \frac{1}{2^k}$.

We claim that the optimal expected value of the MDP is 0, which can be obtained by a strategy σ_{opt} that ensures utility 0 at every step: let m_k be the mass of probability in v'_1 after $3k + 1$ steps (thus $m_0 = \frac{2}{3}$, and m_1, m_2, \dots depend on the strategy). In v'_1 , after $3k + 1$ steps, the strategy σ_{opt} chooses v'_4 with probability α_k such that $m_0 \cdot \alpha_0 = p_0$, thus $\alpha_0 = \frac{1}{4}$, and $m_k \cdot \alpha_k = p_k - p_{k-1}$ for all $k \geq 1$. It is easy to see that $m_k = \frac{1}{3} + \frac{1}{3} \cdot \frac{1}{2^k}$ and $\alpha_k = \frac{1}{2+2^{k+1}}$ ensure this as well as $m_{k+1} = m_k \cdot (1 - \alpha_k)$ for all $k \geq 0$. Therefore the strategy σ_{opt} maintains always the same probability in v'_4 as in v_4 , and the expected total reward is 0 at every step.

It is easy to show that any other strategy (with a different value of some α_k) produces a negative total utility at some time step (either by putting too much probability into v'_4 , and thus too much probability for weight -2 in v'_5 , as compared to the weight 2 in v_5 , or by putting too little probability into v'_4 , and thus too little probability for weight 1 in v'_4 , as compared to the weight -1 in v_4), and that it entails a negative expected value of the MDP.

The strategy σ_{opt} requires infinite memory, since the sequence α_k is strictly decreasing, and the vertex v'_1 is reached after $3k + 1$ steps along a unique path $\rho_k = v_0 v'_1 (v'_2 v'_3 v'_1)^k$. It

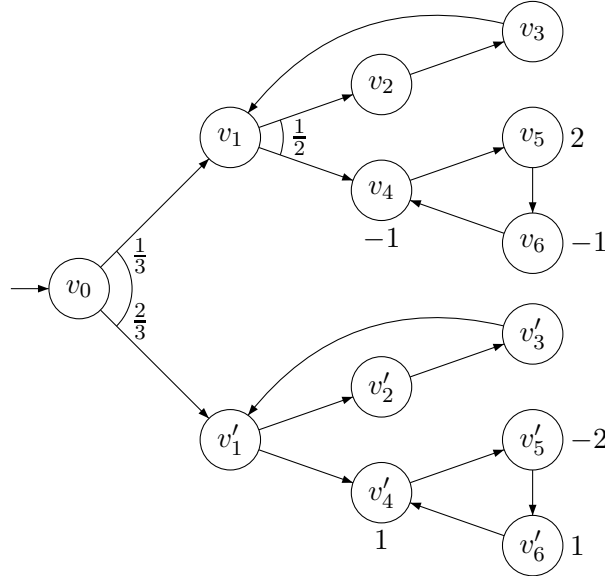


FIGURE 6. An MDP where infinite memory is required for optimal expected value.

follows that for all right congruences \approx over V^+ such that $\rho \approx \rho'$ implies $\sigma_{\text{opt}}(\rho) = \sigma_{\text{opt}}(\rho')$, we have $\rho_k \not\approx \rho_l$ for $k \neq l$ since $\alpha_k \neq \alpha_l$ for $k \neq l$, thus \approx cannot have finite index.

As the above example illustrates, infinite-memory strategies are required in MDPs. The expected stopping-time problem can be formulated as a game between a player that controls the transition choice and the opponent that chooses the stopping times. However, the game is not a perfect-information game as the opponent chooses the stopping times without knowing the execution of the MDP (in particular, the stopping-time distribution cannot be adapted according to the outcome of the probabilistic choices in the MDP). As a consequence, while finite-memory strategies are sufficient in finite-horizon planning (even in perfect-information stochastic games), in contrast we show infinite-memory strategies are required. In general, in imperfect-information probabilistic models such as probabilistic automata [Paz71, Rab63, Rei99], infinite-memory strategies are required [BGB12], and the basic computational problems (such as optimal reachability probability) as well as their *approximation* are undecidable [MHC03]. However, our setting only represents limited imperfect information for the opponent, and we establish in the rest of this section that the approximation problem is decidable.

4.3. Approximation of the optimal value. The problem of computing $\text{val}(\mathcal{M}, T)$ up to an additive error ε can be solved as follows. We show that there exist ε -optimal strategies of a simple form: after some time t^* (that depends on ε), it is sufficient to play a (memoryless) strategy that maximizes the mean-payoff expected reward, defined as follows for a strategy

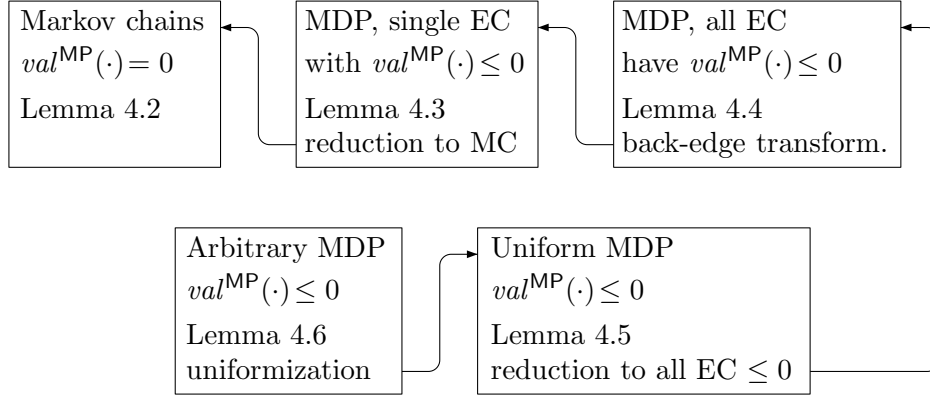


FIGURE 7. Main steps towards the proof that the supremum of total expected reward is bounded in MDPs with mean-payoff value at most 0 (Theorem 4.7).

σ in \mathcal{M} :

$$\text{MP}(\mathcal{M}, \sigma) = \limsup_{t \rightarrow \infty} \frac{1}{t} \sum_{i=0}^{t-1} \mathbb{E}_{\mu}^{\sigma}(w(v_i)),$$

and the *optimal mean-payoff value* is

$$\text{val}^{\text{MP}}(\mathcal{M}) = \sup_{\sigma} \text{MP}(\mathcal{M}, \sigma).$$

Remark 4.1. It is known that (see e.g. [Put94]):

- pure memoryless strategies are sufficient for mean-payoff optimality, that is there exists a pure memoryless strategy σ such that $\text{val}^{\text{MP}}(\mathcal{M}) = \text{MP}(\mathcal{M}, \sigma)$;
- for variants of the definition of mean-payoff expected reward (using \liminf instead of \limsup), or where the \limsup and $\mathbb{E}(\cdot)$ operators are swapped (also known as the expected mean-payoff value), the same pure memoryless strategy is optimal;
- all vertices in an end-component have the same optimal mean-payoff value.

Intuitively, a strategy σ that plays according to an optimal mean-payoff strategy after some time t^* has an asymptotic behaviour that is at least as good as any strategy, in particular any ε -optimal strategy; up to time t^* (thus for finitely many steps), if the strategy σ plays like an ε -optimal strategy, then the sequence of expected reward (defined above as u_t) is also good enough; the only question is whether switching to an optimal mean-payoff strategy may induce a transient loss of reward after t^* that could impede ε -optimality. In fact, we show that (1) the loss is bounded, and (2) the impact of a bounded loss on the expected value is negligible if t^* is large enough. That the loss is bounded, namely:

$$\sup_{\sigma} \sup_t \sum_{i=0}^t \mathbb{E}_{\mu}^{\sigma}(w(v_i)) \text{ is bounded if } \text{val}^{\text{MP}}(\mathcal{M}) \leq 0,$$

may appear intuitively true, but is not simple to prove even in the special case where the mean-payoff value is 0. The proof has several steps, summarized in Figure 7, leading to Theorem 4.7. We start by proving that the loss is bounded in the simple case of Markov

chains with mean-payoff value 0, then for larger classes of MDPs, using reductions that transform an MDP \mathcal{M} of a larger class into an MDP \mathcal{M}' of a smaller class for which a bound on the loss is already established. The transformations may increase the total expected reward (as then, an upper bound for \mathcal{M}' gives an upper bound for \mathcal{M}).

Lemma 4.2. *In aperiodic Markov chains $\langle M, \mu, w \rangle$ with smallest positive transition probability α , if the mean-payoff value, defined as $\limsup_{t \rightarrow \infty} \frac{1}{t} \sum_{i=0}^{t-1} \mathbb{E}(w(v_i))$, is 0, then*

$$\sup_t \left| \sum_{i=0}^t \mathbb{E}(w(v_i)) \right| \leq 4nW \cdot t_0$$

where $t_0 = 3 \cdot n^5 \cdot \left(\frac{1}{\alpha}\right)^{n^2}$, and W is the largest absolute weight according to w .

Proof. Consider the convergence rate of aperiodic Markov chains (Appendix B):

$$|(\mu \cdot M^t)_j - \pi_j| \leq 3 \cdot K^t \quad \text{for all } j \in V$$

where K is a constant with $K < 1$, namely $K = (1 - \alpha^{n^2})^{1/3n^2}$.

Consider time $t_0 = 3 \cdot n^5 \cdot \left(\frac{1}{\alpha}\right)^{n^2}$, which is such that $K^{t_0} \leq \frac{1}{2n^3}$ for all $n \geq 1$ (using Lemma A.1 in Appendix).

Now consider the expected total reward at time t , given by

$$\left| \sum_{i=0}^t \mathbb{E}(w(v_i)) \right| = \left| \sum_{i=0}^t \mu \cdot M^i \cdot w \right|$$

and show that it is bounded by $4nW \cdot t_0$, for all $t \geq 0$. The proof goes by showing a bound on the total reward that can be accumulated within a time unit. At time i , the reward per time unit is $\mu \cdot M^i \cdot w \leq 3nW \cdot K^i$, since the mean-payoff value of the Markov chain is 0, which implies $\pi \cdot w = 0$. For times $t < t_0$, we bound the total reward per time unit trivially by W . For times $k \cdot t_0 \leq t < (k+1) \cdot t_0$ (where $k \geq 1$), we bound the total reward per time unit by $3nW \cdot \left(\frac{1}{2n^3}\right)^k$ since $3 \cdot K^t \leq 3 \cdot K^{k \cdot t_0} \leq 3 \cdot \left(\frac{1}{2n^3}\right)^k$.

It is now sufficient to show that the sum of the bounds on the total reward per time unit is bounded by $t_0 \cdot (W + 1)$ for arbitrarily large t , which we establish as follows:

$$t_0 \cdot W + \sum_{k=1}^{\infty} 3nW \cdot \left(\frac{1}{2n^3}\right)^k \cdot t_0 = t_0 \cdot W + 3nW \cdot \frac{t_0}{2n^3 - 1} \leq 4nW \cdot t_0. \quad \square$$

To prove a similar result for MDPs (Theorem 4.7), we first consider the case of MDPs that consist of a single end-component, and show by contradiction that if it has mean-payoff value 0 and a large expected total reward could be accumulated from a vertex v_0 using some strategy σ_0 , then by reaching v_0 again (which is possible since the MDP is strongly connected) and repeating the same strategy σ_0 , we could get a strictly positive mean-payoff value. A technical difficulty in this proof is that v_0 may be reached by paths of different lengths, but the large expected total reward that can be accumulated from v_0 is obtained in a fixed number of steps.

Lemma 4.3. *In an MDP \mathcal{M} that is an end-component (i.e., V is an end-component), if $val^{\text{MP}}(\mathcal{M}) \leq 0$ and $|V| = n$, then*

$$\sup_{\sigma} \sup_t \sum_{i=0}^t \mathbb{E}_{\mu}^{\sigma}(w(v_i)) \leq 12 \cdot n^6 \cdot W \cdot \left(\frac{1}{\alpha}\right)^{n^3}$$

where α is the smallest positive transition probability in \mathcal{M} , and W its largest absolute weight.

Proof. Assume towards contradiction that there exists an initial distribution μ such that the inequality of the lemma does not hold in \mathcal{M} . It follows that in some initial vertex v_0 (such that $\mu(v_0) > 0$) the inequality does not hold, i.e. there exists a strategy σ_0 and time t_0 such that the expected total reward from v_0 under σ_0 at time t_0 is at least $12 \cdot n^6 \cdot W \cdot \left(\frac{1}{\alpha}\right)^{n^3}$.

First we modify the MDP \mathcal{M} to obtain an MDP \mathcal{M}' as follows, in a way that does not decrease the value of $\sup_{\sigma} \sup_t \sum_{i=0}^t \mathbb{E}_{\mu}^{\sigma}(w(v_i))$:

- increase the weight of every vertex by $|val^{\text{MP}}(\mathcal{M})|$ (define $w'(v) = w(v) - val^{\text{MP}}(\mathcal{M})$), and
- add a copy \hat{v}_0 of v_0 with weight $-W$ and a self-loop; formally, define a new action set $A' = A \cup \{\hat{a} \mid a \in A\}$, a new state space $V' = V \cup \{\hat{v}_0\}$ with $w'(\hat{v}_0) = -W$ and $w'(v) = w(v)$ for all $v \in V$, and transitions on \hat{a} that replace v_0 by \hat{v}_0 as follows, for all $v \in V$ and all $a \in A$: $\theta'(v, \hat{a})(\hat{v}_0) = \theta(v, a)(v_0)$, $\theta'(v, \hat{a})(v_0) = 0$, and $\theta'(v, \hat{a})(v') = \theta(v, a)(v')$ for all $v' \in V \setminus \{v_0\}$. The self-loop on \hat{v}_0 is defined on all actions \hat{a} (i.e., $\theta'(\hat{v}_0, \hat{a})(\hat{v}_0) = 1$), and the other actions have the same effect as from v_0 (i.e., $\theta'(\hat{v}_0, a) = \theta(v_0, a)$ and $\theta'(v, a) = \theta(v, a)$ for all $v \in V$).

In the new MDP $\mathcal{M}' = \langle V', A', \theta', \mu, w' \rangle$, we note that:

- the expected total reward from v_0 is not smaller in \mathcal{M}' than in \mathcal{M} , since we increased weights of existing transitions, and we added new transitions, which cannot decrease the expected total reward (strategies of \mathcal{M} can still be played in \mathcal{M}');
- the optimal mean-payoff value of \mathcal{M}' is $val^{\text{MP}}(\mathcal{M}') = 0$ since increasing all weights by $|val^{\text{MP}}(\mathcal{M})|$ has the effect to increase the mean-payoff value by the same amount; moreover, adding the copy \hat{v}_0 with weight $-W$ does not change the optimal mean-payoff value. To see this, fix a memoryless strategy σ , and consider the recurrent classes of the resulting Markov chain. If a recurrent class C contains \hat{v}_0 , then either the self-loop on \hat{v}_0 is used by the strategy σ and then the mean-payoff value of C is $-W$, or the self-loop on \hat{v}_0 is not used and the mean-payoff value of C is less than the mean-payoff value of $C' = C \cup v_0 \setminus \{\hat{v}_0\}$ which is a recurrent class that can be obtained in \mathcal{M} using the strategy that copies σ but plays a whenever σ plays \hat{a} . Since $val^{\text{MP}}(\mathcal{M}) \leq 0$, it follows that the mean-payoff value of C' (and thus of C) is at most 0. On the other hand, if C does not contain \hat{v}_0 , then it can be obtained in \mathcal{M} and thus its mean-payoff value is also at most 0.
- the state space of \mathcal{M}' is of size $|V'| = n + 1 \leq n^2$, and \mathcal{M}' is still an end-component.

Given the strategy σ_0 and time t_0 as above, we show that there exists a strategy σ^* and a time t^* such that from all vertices $v \in V$, we have

$$\sum_{i=0}^{t^*-1} \mathbb{E}_v^{\sigma^*}(w'(v_i)) \geq 5 \cdot n^6 \cdot W \cdot \left(\frac{1}{\alpha}\right)^{n^3} > 0 \quad (4.1)$$

which entails, since the bound t^* is the same for all vertices, that by repeating the strategy σ_0 every t_0 steps the mean-payoff value of \mathcal{M}' is positive, $val^{\text{MP}}(\mathcal{M}') > 0$, in contradiction with the fact that $val^{\text{MP}}(\mathcal{M}') = 0$.

Let $t^* = t_{\text{reach}} + t_0$ where $t_{\text{reach}} = \frac{n}{\alpha^n}$ and we construct σ^* that plays as follows:

- for t_{reach} steps, play a pure memoryless strategy σ_{reach} to reach \hat{v}_0 almost-surely (and play the self-loop with weight $-W$ in \hat{v}_0); such a strategy exists because the MDP is an end-component [dA97];
- after t_{reach} steps: if the current vertex is \hat{v}_0 , play for the next t_0 steps the strategy σ_0 ; if the current vertex is not \hat{v}_0 , play for the next t_0 steps a memoryless optimal strategy for the mean-payoff value (thus using only actions in A and without visiting \hat{v}_0).

The value t_{reach} is such that the probability mass in \hat{v}_0 after t_{reach} steps is at least $\frac{1}{2}$: since the strategy σ_{reach} is pure memoryless and $|V'| = n + 1$, we can use the analysis of Markov chain reachability to claim that the probability to have reached target vertex \hat{v}_0 after $t_{\text{reach}} = k \cdot n$ steps is at least $1 - (1 - \alpha^n)^k > \frac{1}{2}$ since $k = \frac{1}{\alpha^n}$ and $(1 - \frac{1}{x})^x < e^{-1} < \frac{1}{2}$ (by Lemma A.1 in Appendix).

We bound the expected total reward of σ^* as follows:

- after t_{reach} steps, since all weights are bounded by W , the expected total reward is at least $-t_{\text{reach}} \cdot W$;
- in the next t_0 steps, the collected reward from \hat{v}_0 (in which the probability mass is at least $\frac{1}{2}$) is at least $12 \cdot n^6 \cdot W \cdot (\frac{1}{\alpha})^{n^3}$ (by the definition of σ_0 and t_0), and the collected reward from other vertices is at least $-12 \cdot n^6 \cdot W \cdot (\frac{1}{\alpha})^{n^2}$ (by Lemma 4.2, since the optimal strategy for the mean-payoff value is memoryless and plays only actions in A , which gives a Markov chain with n vertices and mean-payoff value equal to 0).

It follows that the expected total reward of σ^* after t^* steps is at least:

$$\begin{aligned} & -\frac{nW}{\alpha^n} + 6 \cdot n^6 \cdot W \cdot \left(\frac{1}{\alpha}\right)^{n^3} - 6 \cdot n^6 \cdot W \cdot \left(\frac{1}{\alpha}\right)^{n^2} \\ & \geq 6 \cdot n^6 \cdot W \cdot \left(\frac{1}{\alpha}\right)^{n^3} - 7 \cdot n^6 \cdot W \cdot \left(\frac{1}{\alpha}\right)^{n^2} \\ & \geq 5 \cdot n^6 \cdot W \cdot \frac{1}{\alpha^{n^3}} \text{ since } n \geq 2 \text{ and } \alpha \leq \frac{1}{2}, \end{aligned}$$

which establishes the bound (4.1) and concludes the proof. \square

We can easily extend the result to MDPs with several end-components, if all of them have mean-payoff value at most 0.

Lemma 4.4. *In an MDP \mathcal{M} with n vertices in which all end-components have an optimal mean-payoff value at most 0, we have*

$$\sup_{\sigma} \sup_t \sum_{i=0}^t \mathbb{E}_{\mu}^{\sigma}(w(v_i)) \leq 12 \cdot n^8 \cdot W \cdot \left(\frac{1}{\alpha}\right)^{n^3+n}$$

where α is the smallest positive transition probability in \mathcal{M} , and W its largest absolute weight.

Proof. Consider an MDP \mathcal{M} as in the lemma statement, and assume without loss of generality that the initial distribution μ is a Dirac distribution, namely $\mu(v_0) = 1$ for some vertex v_0 .

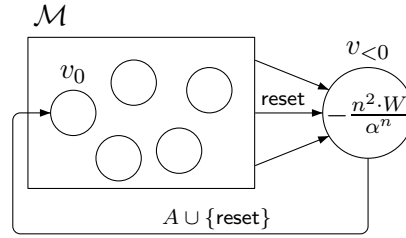


FIGURE 8. Back-edge transformation (Proof of Lemma 4.4).

We present the back-edge transformation from \mathcal{M} to an MDP \mathcal{M}' as follows (Figure 8). We add a special action `reset`, thus the action set of \mathcal{M}' is $A \cup \{\text{reset}\}$. The state space of \mathcal{M}' is $V \cup \{v_{<0}\}$ where $v_{<0}$ is a new vertex with weight $-\frac{n^2 \cdot W}{\alpha^n}$, and the transition functions of \mathcal{M}' and \mathcal{M} are identical for actions in A and vertices in V . On action `reset`, the transition function θ' of \mathcal{M}' has a “back-edge” from every vertex to $v_{<0}$; from $v_{<0}$ there is an edge to v_0 on every action. Thus $\theta'(v, \text{reset}) = v_{<0}$ for all vertices $v \in V$, and $\theta'(v_{<0}, a) = \mu$ for all actions $a \in A \cup \{\text{reset}\}$.

First note that the supremum of expected total reward is not smaller in \mathcal{M}' than in \mathcal{M} , since \mathcal{M}' contains all vertices and transitions of \mathcal{M} .

We now show that the mean-payoff value of \mathcal{M}' is at most 0, which establishes the bound in the lemma statement as follows: since the whole \mathcal{M}' is an end-component, we can apply Lemma 4.3 where the largest absolute weight in \mathcal{M}' is $\frac{n^2 \cdot W}{\alpha^n}$, which gives the announced bound for \mathcal{M}' , thus also for \mathcal{M} .

To show that the mean-payoff value of \mathcal{M}' is at most 0, fix an optimal strategy σ for mean-payoff in \mathcal{M}' , which we can assume to be pure memoryless (Remark 4.1). In the resulting Markov chain, we show that all recurrent classes have mean-payoff value at most 0 as follows: if a recurrent class does not contain $v_{<0}$, then it is an end-component in the original MDP \mathcal{M} , and therefore its mean-payoff value is at most 0; otherwise it contains $v_{<0}$ and since the frequency³ f_0 of a recurrent vertex is at least $\frac{\alpha^n}{n+1} \geq \frac{\alpha^n}{n^2}$, the mean-payoff value of the recurrent class is at most $-f_0 \cdot \frac{n^2 \cdot W}{\alpha^n} + (1 - f_0) \cdot W \leq -W + W = 0$. This shows that all recurrent classes have mean-payoff value at most 0, and thus the optimal mean-payoff value of the MDP \mathcal{M}' is at most 0. \square

In an arbitrary MDP with mean-payoff value at most 0, some end-components may have positive value, and others negative value, as in the example of Figure 9: the three end-components $\{v_0\}$, $\{v_1, v_2\}$, $\{v_3\}$ have respective mean-payoff value -1 , 1 , and -2 . From the initial distribution μ where $\mu(v_0) = \mu(v_1) = \frac{1}{2}$, the mean-payoff value is 0. The case where the MDP has some end-components with positive mean-payoff value requires a slightly more technical proof (see also Figure 7): we first show in Lemma 4.5 that the supremum of expected total reward in MDPs is bounded if all end-components are *uniform* (an end-component is uniform if all its vertices have the same weight); then we present uniformization in Lemma 4.6 to transform arbitrary MDPs into uniform MDPs.

³The frequency of the vertex with largest frequency is at least $f_n = \frac{1}{n+1}$, and the frequency of the $(k+1)$ -th frequent vertex is at least $\alpha^k \cdot f_n$.

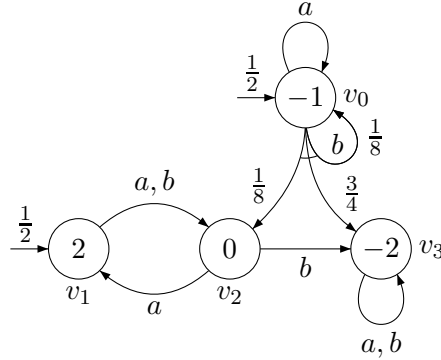


FIGURE 9. An MDP with positive and negative end-components. Its mean-payoff value is 0.

Given an MDP with weight vector w , let \mathcal{E} be the union of all its end-components. Define the vector w_{trans} and w_{ec} as follows:

$$w_{\text{trans}}(v) = \begin{cases} w(v) & \text{if } v \in V \setminus \mathcal{E} \\ 0 & \text{if } v \in \mathcal{E} \end{cases}$$

$$w_{\text{ec}}(v) = \begin{cases} 0 & \text{if } v \in V \setminus \mathcal{E} \\ w(v) & \text{if } v \in \mathcal{E} \end{cases}$$

It follows that $w = w_{\text{trans}} + w_{\text{ec}}$ and by the triangular inequality, we have

$$\sup_k \sum_{i=0}^k \mathbb{E}_\mu^\sigma(w(v_i)) \leq \sup_k \sum_{i=0}^k \mathbb{E}_\mu^\sigma(w_{\text{trans}}(v_i)) + \sup_k \sum_{i=0}^k \mathbb{E}_\mu^\sigma(w_{\text{ec}}(v_i)).$$

Using Lemma 4.4, it is easy to bound the supremum of expected total reward for w_{trans} , and we present a bound on the supremum of expected total reward for w_{ec} in uniform MDPs as follows.

Lemma 4.5. *Given an MDP \mathcal{M} with n vertices, let w_{trans} and w_{ec} be the weight vectors of the transient vertices and of the end-components, respectively. We have*

$$\sup_\sigma \sup_t \sum_{i=0}^t \mathbb{E}_\mu^\sigma(w_{\text{trans}}(v_i)) \leq 12 \cdot n^8 \cdot W \cdot \left(\frac{1}{\alpha}\right)^{n^3+n},$$

and if $\text{val}^{\text{MP}}(\mathcal{M}) \leq 0$ and all end-components of M are uniform, then

$$\sup_\sigma \sup_t \sum_{i=0}^t \mathbb{E}_\mu^\sigma(w_{\text{ec}}(v_i)) \leq 12 \cdot n^8 \cdot W \cdot \left(\frac{1}{\alpha}\right)^{n^3+n},$$

where α is the smallest positive transition probability in \mathcal{M} , and W its largest absolute weight.

Proof. For the first part, the bound for w_{trans} is given by Lemma 4.4, since the MDP with weight vector w_{trans} has all its end-component with mean-payoff value 0.

For the second part, assuming $\text{val}^{\text{MP}}(\mathcal{M}) \leq 0$ and all end-components of M are uniform, the bound for w_{ec} is established as follows. First consider the MDP \mathcal{M} with weight function w' defined by $w'(v) = w(v)$ if $v \in \mathcal{E}$, and $w'(v) = -W$ otherwise (i.e., for transient vertices).

Note that $w' = w_{\text{ec}} - w_0$ where $w_0(v) = 0$ if $v \in \mathcal{E}$, and $w_0(v) = W$ otherwise. In \mathcal{M} with w' , we will show that:

$$\sup_{\sigma} \sup_t \sum_{i=0}^t \mathbb{E}_{\mu}^{\sigma}(w'(v_i)) \leq 0. \quad (4.2)$$

To obtain the bound for $w_{\text{ec}} = w' + w_0$ and conclude the proof, we use the triangular inequality which entails that $\sup_{\sigma} \sup_t \sum_{i=0}^t \mathbb{E}_{\mu}^{\sigma}(w_{\text{ec}}(v_i)) \leq 0 + B$ where B is the bound given by the first part of the lemma for w_0 .

To show (4.2), we consider an arbitrary strategy σ and we show that for all $n \geq 0$, the expected reward at step n is at most 0, that is:

$$\sum_{v \in V} \delta_n^{\sigma}(v) \cdot w'(v) \leq 0$$

where δ_n^{σ} is the vertex distribution of \mathcal{M} after n steps under strategy σ .

Given δ_n^{σ} (as an initial distribution), consider a memoryless strategy $\sigma^{\mathcal{E}}$ such that end-components are never left, defined as follows for all $v \in V$: if $v \in \mathcal{E}$, then $\sigma^{\mathcal{E}}(v)$ is an action to stay in the end-component of v in the next step; otherwise, $\sigma^{\mathcal{E}}(v)$ is an arbitrary action. Such a strategy $\sigma^{\mathcal{E}}$ exists by definition of end-components.

By the assumption that $\text{val}^{\text{MP}}(\mathcal{M}) \leq 0$ (with weight vector w thus also with w'), it follows that the limit distributions δ^* satisfy $\sum_{v \in \mathcal{E}} \delta^*(v) \cdot \eta(v) \leq 0$.

Then, within the distribution δ_n^{σ} , the probability mass $p_{\mathcal{E}}$ in \mathcal{E} never leaves an end-component, and the probability mass $1 - p_{\mathcal{E}}$ in $V \setminus \mathcal{E}$ eventually (in the limit) gets injected in \mathcal{E} (and then never leaves). The (future) contribution of the probability mass $1 - p_{\mathcal{E}}$ to the expected reward of limit distributions is bounded below by $-W$ (since $\eta(v) \geq -W$ for all $v \in \mathcal{E}$, where $\eta(v)$ is the mean-payoff value of v and of the end-component containing v since \mathcal{M} is uniform). It follows that:

$$\begin{aligned} 0 &\geq \sum_{v \in \mathcal{E}} \delta^*(v) \cdot \eta(v) \\ &\geq \sum_{v \in \mathcal{E}} \delta_n^{\sigma}(v) \cdot \eta(v) + \sum_{v \in V \setminus \mathcal{E}} \delta_n^{\sigma}(v) \cdot (-W) \\ &= \sum_{v \in V} \delta_n^{\sigma}(v) \cdot w'(v) \end{aligned}$$

which concludes the proof by entailing (4.2). \square

Uniformization. We present a *uniformization* procedure that, given an MDP \mathcal{M} with mean-payoff value at most 0, constructs an MDP \mathcal{M}' with the same mean-payoff value as \mathcal{M} , with a larger supremum of expected total reward, and in which all end-components are uniform.

The procedure has two steps. First we construct \mathcal{M}' and weight vector w_1 by transforming the structure of \mathcal{M} in such a way that for each end-component, there is a single vertex from which the end-component can be entered (as illustrated on the left of Figure 10). This shape of MDP can be obtained as follows: for each end-component E , create a new vertex v_E (with weight $w_1(v_E) = 0$) with an edge from v_E to every vertex in E , and modify the transition function from every vertex v outside E to redirect all the probability mass that was going from v to E to go to v_E . Analogously we transfer the probability of the initial distribution that was in E to v_E . By doubling the weight of every vertex to define w_1 (and inserting, for every vertex, a new intermediate vertex with weight 0 that is entered

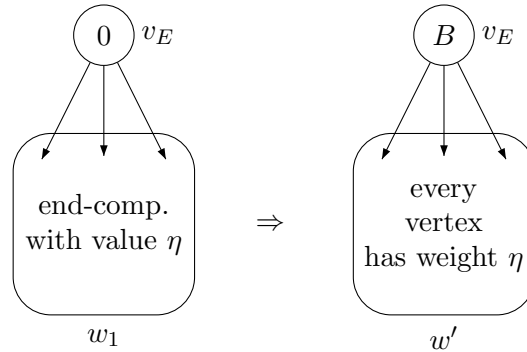


FIGURE 10. Uniformization (including transformation of the weight vector w_1 into w') for Lemma 4.6. The number B is twice the bound given by Lemma 4.4.

before going to the original vertex), it is easy but tedious to show that the mean-payoff value according to w_1 remains the same, and that the supremum of expected total reward is at least twice the original one (note also that the number of vertices has at most tripled).

In the second step, we construct a new vector w' of weights for \mathcal{M}' such that every end-component becomes uniform (illustrated on the right of Figure 10). Given an end-component E and the vertex v_E , let $w'(v_E) = B$ and $w'(v) = \eta$ where η is the mean-payoff value of E (according to w_1 , or equivalently according to w) and $B = 12 \cdot n^8 \cdot (3W) \cdot (\frac{1}{\alpha})^{n^3+n}$ is three times the bound given by Lemma 4.4.

Lemma 4.6. *Given an MDP \mathcal{M} with n vertices, we can construct an MDP \mathcal{M}' with the following properties:*

- (1) *all end-components of \mathcal{M}' are uniform.*
- (2) *\mathcal{M} and \mathcal{M}' have the same mean-payoff value;*
- (3) *the number of vertices in \mathcal{M}' is at most $3n$;*
- (4) *the supremum of expected total reward in \mathcal{M} is less than half the supremum of expected total reward in \mathcal{M}' ;*

Proof. Consider \mathcal{M}' obtained from \mathcal{M} by the uniformization procedure. Item 1. holds by construction, and the proof of item 2. and item 3. has been sketched along with the uniformization procedure (note that the definition of w' does not change the mean-payoff value of the end-components, as compared to weight vector w_1 and w).

We show item 4. for the transformation of one end-component E (Figure 10), and the lemma follows by applying the result successively to each end-component.

Given an arbitrary strategy σ , let $s_k(v) = \sum_{i=0}^k \mathbb{E}_v^\sigma(w_1(v_i))$ and $s'_k(v) = \sum_{i=0}^k \mathbb{E}_v^\sigma(w'(v_i))$ be the expected total reward from initial vertex v after k steps, according to the weight vector w_1 and w' respectively. We show that for all $v \in V \setminus E$, and for all k , we have $s_k(v) \leq s'_k(v)$, which establishes the claim that \mathcal{M}' (with w') has a larger supremum of expected total reward than \mathcal{M}' (with w_1), since the initial distribution of \mathcal{M}' has support in $V \setminus E$ (and we showed that the supremum of expected total reward in \mathcal{M}' (with w_1) is at least twice the one in \mathcal{M}). We show this by induction on k . The base case $k = 0$ holds

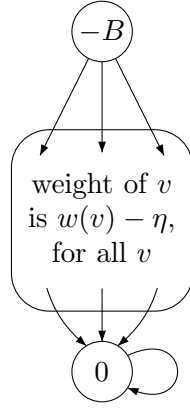


FIGURE 11. Over-approximation of the MDP with weight vector $w_{\text{dif}} = w - w'$ (from Figure 10).

since $w_1(v) \leq w'(v)$ for all $v \in V \setminus E$. For the induction case, assume $s_t(v) \leq s'_t(v)$ for all $v \in V \setminus E$, and for all $t \leq k - 1$. Let $v \in V \setminus E$ and we consider two cases:

- If $v \neq v_E$, the claim $s_k(v) \leq s'_k(v)$ holds since all successors of v are in $V \setminus E$ and we can use the induction hypothesis:

$$\begin{aligned} s_k(v) &= w_1(v) + \max_{a \in A} \sum_{u \in V \setminus E} s_{k-1}(u) \cdot \theta(v, a)(u) \\ &\leq w_1(v) + \max_{a \in A} \sum_{u \in V \setminus E} s'_{k-1}(u) \cdot \theta(v, a)(u) \\ &= s'_k(v) \end{aligned}$$

- If $v = v_E$, then all successors of v belong to E , and we cannot directly use the induction hypothesis. Consider the weight vector $w_{\text{dif}} = w_1 - w'$, and show that for $s_k^{\text{dif}}(v) = \sum_{i=0}^k \mathbb{E}_v^\sigma(w_{\text{dif}}(v_i))$, we have $s_k^{\text{dif}}(v) \leq 0$, which implies $s_k(v) \leq s'_k(v)$. Given the weight vector w_{dif} , we know that starting from $v = v_E$, as soon as a path leaves the end-component E , its contribution to the expected total reward is at most 0 (by induction hypothesis, since at most $k - 1$ steps remain after exiting). Therefore, it is sufficient to show that the expected total reward in k steps is at most 0 in the MDP of Figure 11 where the edges going out of the end-component E are directed to a sink with weight 0. The weights defined by w_{dif} in E are bounded by $\|w_1\| + \|w'\| \leq 2W + W$. The number of vertices in E is less than $2n$ (where n is the number of vertices in the original MDP \mathcal{M}), and only half of them are relevant to define the expected total reward. By Lemma 4.4, since all end-components have mean-payoff value 0 in the MDP of Figure 11, and B is the bound given by Lemma 4.4, we have $\sup_k s_k^{\text{dif}}(v) \leq 0$ and therefore $s_k^{\text{dif}}(v) \leq 0$, which concludes the proof of the induction case. \square

We finally obtain an analogue of Lemma 4.2 for MDPs: the expected total reward is bounded in MDPs with non-positive mean-payoff value.

Theorem 4.7. *Given an MDP \mathcal{M} with n vertices and $\text{val}^{\text{MP}}(\mathcal{M}) \leq 0$, we have:*

$$\sup_{\sigma} \sup_t \sum_{i=0}^t \mathbb{E}_{\mu}^{\sigma}(w(v_i)) \in O\left(n^{16} \cdot W \cdot \left(\frac{1}{\alpha}\right)^{O(n^3)}\right)$$

where α is the smallest positive transition probability in \mathcal{M} , and W its largest absolute weight.

Proof. Given MDP \mathcal{M} , we use the triangular inequality to bound the supremum of expected total reward by the sum the supremum on the transient vertices and on the end-components. For transient vertices, we use directly the bound in Lemma 4.5, and for end-components, we use the construction of \mathcal{M}' in Lemma 4.6, and then apply Lemma 4.5 where the number of vertices is $3n$, and the largest weight is $B = 12 \cdot n^8 \cdot (3W) \cdot \left(\frac{1}{\alpha}\right)^{n^3+n}$. Since the supremum of expected total reward in \mathcal{M}' is twice the supremum of expected total reward in \mathcal{M} , we get the following bound:

$$\begin{aligned} & \sup_{\sigma} \sup_t \sum_{i=0}^t \mathbb{E}_{\mu}^{\sigma}(w(v_i)) \\ & \leq 12 \cdot n^8 \cdot W \cdot \left(\frac{1}{\alpha}\right)^{n^3+n} \\ & \quad + \frac{12}{2} \cdot (3n)^8 \cdot \left(12 \cdot n^8 \cdot (3W) \cdot \left(\frac{1}{\alpha}\right)^{n^3+n}\right) \cdot \left(\frac{1}{\alpha}\right)^{27n^3+3n} \\ & = \underbrace{12 \cdot n^8 \cdot W \cdot \left(\frac{1}{\alpha}\right)^{n^3+n} + 2^3 \cdot 3^{11} \cdot n^{16} \cdot W \cdot \left(\frac{1}{\alpha}\right)^{28n^3+4n}}_{B^*} \\ & \in O\left(n^{16} \cdot W \cdot \left(\frac{1}{\alpha}\right)^{O(n^3)}\right). \quad \square \end{aligned}$$

Using Theorem 4.7, for all $\varepsilon > 0$ we can compute a bound t^* such that there exists an ε -optimal strategy (for expected value) that plays according to an optimal mean-payoff strategy after time t^* .

Lemma 4.8. *Given an MDP \mathcal{M} and $\varepsilon > 0$, there exists an ε -optimal strategy that plays, after time $t^* = \frac{T \cdot (2B^* + \varepsilon)}{\varepsilon}$ (where B^* is the bound given by Theorem 4.7), according to a memoryless optimal strategy σ_{MP} for the mean-payoff value.*

Proof. Consider an arbitrary strategy σ in \mathcal{M} (under expected stopping time T), and given $t^* \geq T$, consider a strategy σ^* that plays like σ up to time t^* , and then switches to a memoryless mean-payoff optimal strategy σ_{MP} , in the MDP \mathcal{M} with initial distribution $\mu^* = \delta_{t^*}^{\sigma}$ (the vertex distribution of \mathcal{M} after t^* steps under strategy σ). Let η^* be the optimal mean-payoff value from μ^* in \mathcal{M} , and let $w' = w - \eta^*$ (where $w'(v) = w(v) - \eta^*$ for all $v \in V$). With weight vector w' , the optimal mean-payoff value of \mathcal{M} is 0 from μ^* .

Using Lemma 4.2 in the Markov chain obtained by fixing the strategy σ_{MP} in \mathcal{M} with initial distribution μ^* , we obtain:

$$\sup_t \left| \sum_{i=0}^t \mathbb{E}_{\mu^*}^{\sigma_{\text{MP}}}(w'(v_i)) \right| \leq \underbrace{12 \cdot n^6 \cdot W \cdot \left(\frac{1}{\alpha}\right)^{n^2}}_{C^*}. \quad (4.3)$$

Let $u_t = \sum_{i=0}^t \mathbb{E}_{\mu}^{\sigma}(w(v_i))$ and let $u_t^* = \sum_{i=0}^t \mathbb{E}_{\mu^*}^{\sigma^*}(w(v_i))$ be the sequence of expected total reward under strategy σ and σ^* respectively. To show ε -optimality of σ^* , take $t^* \geq \frac{T \cdot (2B^* + \varepsilon)}{\varepsilon}$ and show that:

$$\text{val}(u^*, T) \geq \text{val}(u, T) - \varepsilon$$

The proof is in two steps. First we bound the difference $u_t - u_t^*$ as follows, for all $t \geq 1$:

$$\begin{aligned} u_t - u_t^* &= \sum_{i=0}^t \mathbb{E}_{\mu}^{\sigma}(w(v_i)) - \sum_{i=0}^t \mathbb{E}_{\mu^*}^{\sigma^*}(w(v_i)) \\ &= \sum_{i=0}^t \mathbb{E}_{\mu}^{\sigma}(w'(v_i)) - \sum_{i=0}^t \mathbb{E}_{\mu^*}^{\sigma^*}(w'(v_i)) \\ &\hspace{15em} (\text{since } \mathbb{E}(w(\cdot)) = \mathbb{E}(w'(\cdot)) + \eta^*) \\ &= \sum_{i=t^*}^t \mathbb{E}_{\mu}^{\sigma}(w'(v_i)) - \sum_{i=t^*}^t \mathbb{E}_{\mu^*}^{\sigma^*}(w'(v_i)) \\ &\hspace{15em} (\sigma \text{ and } \sigma^* \text{ agree in the first } t^* \text{ steps}) \\ &\leq B^* + C^* \leq 2B^* \\ &\hspace{15em} (\text{triangular inequality and bounds given by Theorem 4.7 and (4.3)}) \end{aligned}$$

In a second step, consider an arbitrary bi-Dirac distribution δ with support $\{t_1, t_2\}$ and expected stopping-time T , and consider the difference between the value of sequences u_t and u_t^* under δ , if $t_2 \geq t^*$ (the difference is 0 if $t_2 < t^*$):

$$\begin{aligned} &\mathbb{E}_{\delta}(u) - \mathbb{E}_{\delta}(u^*) \\ &= \frac{u_{t_1}(t_2 - T) + u_{t_2}(T - t_1)}{t_2 - t_1} - \frac{u_{t_1}^*(t_2 - T) + u_{t_2}^*(T - t_1)}{t_2 - t_1} \\ &= \frac{T - t_1}{t_2 - t_1} \cdot (u_{t_2} - u_{t_2}^*) \\ &\hspace{15em} (\text{since } \sigma \text{ and } \sigma^* \text{ agree in the first } t^* \text{ steps, and thus } u_{t_1} = u_{t_1}^*) \\ &\leq \frac{T - t_1}{t_2 - t_1} \cdot 2B^* \leq \frac{T}{t^* - T} \cdot 2B^* \leq \varepsilon \\ &\hspace{15em} (\text{since } 0 \leq t_1 \leq T) \end{aligned}$$

It follows that under all bi-Dirac distributions δ with expected stopping-time T , the expected value of the sequence u_t^* is, up to additive error ε , greater than the expected value

of u_t . Therefore, since bi-Dirac distributions are sufficient for optimality (Section 3.2.1), we have $\text{val}(u^*, T) \geq \text{val}(u, T) - \varepsilon$. Hence σ^* is ε -optimal. \square

We can express in the existential theory of the reals that the value of a strategy that eventually plays according to a memoryless strategy (as in Lemma 4.8) is above a given threshold, which entails decidability of computing an approximation of the optimal value up to an additive error ε .

Lemma 4.9. *Given an MDP \mathcal{M} and a time t^* , we can compute to an arbitrary level of precision $\varepsilon > 0$ the optimal value among the strategies that play after time t^* according to a memoryless strategy.*

Proof. We describe the choices of an arbitrary strategy up to time t^* using variables $x_{v,t,a}$ for every $v \in V$, $0 \leq t \leq t^*$, and $a \in A$, where $x_{v,t,a}$ is the probability to play action a at time t in vertex v . Note that we ignore the history of vertices, which is no loss of generality since the utility achieved by a strategy at time t only depends on the probability mass in each vertex at time t , and if a sequence of distribution can be achieved by some strategy, then it can be achieved by a Markov strategy (in which the choice depends only on the time and the current vertex). We can express the probability mass $p_{v,t}$ in v at time t as $p_{v,t} = \sum_{u \in V} \sum_{a \in A} p_{u,t-1} \cdot x_{u,t-1,a} \cdot \theta(u, a)(v)$ where θ is the transition function of \mathcal{M} . It is then easy to express the utility u_t as a function of the variables $p_{v,t}$ and $x_{v,t,a}$.

After time t^* , consider a memoryless strategy and we can express its mean-payoff value η^* as a function of the vertex distribution at time t^* , thus as a function of the variables p_{v,t^*} . Then for $t = t^* + 1, t^* + 2, \dots, \hat{t}$, we express the utility u_t at time t as a function of the variables $x_{v,t,a}$ and $p_{v,t}$, and consider the utility sequence $u_0, \dots, u_{\hat{t}}, u_{\hat{t}} + \eta^*, u_{\hat{t}} + 2\eta^*, \dots$ (corresponding to an ultimately periodic path) using Lemma 3.8 and by an argument similar to the proof of Lemma 3.9 using the bound of Lemma 4.2 for Markov chains, we get a bound on the approximation error as follows: the value after \hat{t} differ by at most $D = n \cdot W \cdot K_1 \cdot K_3^{\hat{t}-t^*}$ from the actual utility, thus the error on the value is at most

$$\frac{D \cdot (T - t_1)}{t_2 - t_1} \leq D \cdot T$$

which is at most ε for $\hat{t} \geq t^* + B$ where $B = \frac{\ln(\frac{\varepsilon}{n \cdot W \cdot T \cdot K_1})}{\ln(K_3)}$ (Lemma 3.9). \square

By Lemma 4.8 and Lemma 4.9, we can compute up to error $\frac{\varepsilon}{2}$ the value of an $\frac{\varepsilon}{2}$ -optimal strategy, and since the error is additive ($\varepsilon = \frac{\varepsilon}{2} + \frac{\varepsilon}{2}$), it follows from the proof of Lemma 4.9 that, by computing (as a symbolic expression in variables $x_{v,t,a}$ and $p_{v,t}$) the sequence of utilities up to time $\hat{t} = \frac{T \cdot (4B^* + \varepsilon)}{\varepsilon} + \frac{\ln(\frac{\varepsilon}{2n \cdot W \cdot T \cdot K_1})}{\ln(K_3)}$ and then considering an increment of η^* at every step, we can compute the value of optimal expected value of the MDP up to error ε in exponential space (since \hat{t} is exponential and the existential theory of the reals can be decided in PSPACE [Can88]). In this way, we obtain the main result of this section: an approximation of the value with expected stopping time can be computed for MDPs up to an arbitrary additive error.

Theorem 4.10. *The optimal expected value of an MDP with expected stopping time T can be computed to an arbitrary level of precision $\varepsilon > 0$, in exponential space.*

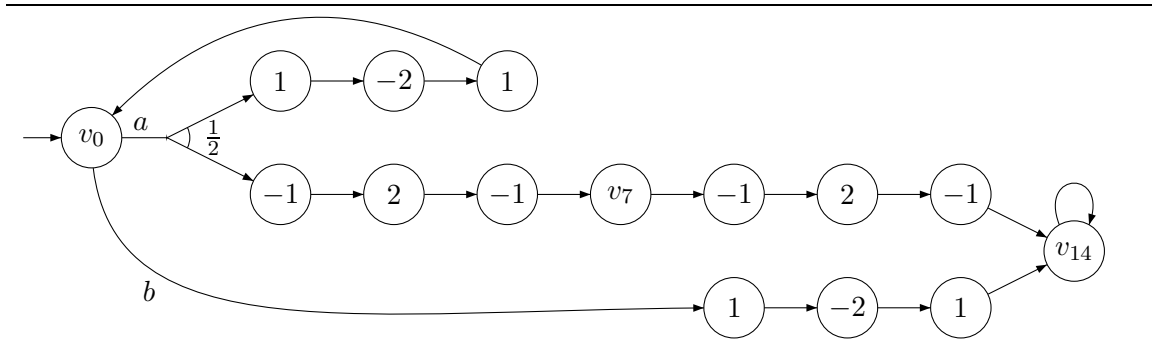


FIGURE 12. An MDP where memory is necessary for optimal expected value in pure strategies.

5. CONCLUSION

We studied Markov chains and MDPs with expected stopping time, and showed the hardness of computing the exact value, as the associated decision problem for Markov chains is inter-reducible with the Positivity problem, thus at least as hard as the Skolem problem. Approximation of the value can be computed in exponential time for Markov chains, and exponential space for MDPs (thus the approximation problem is decidable although optimal strategies require infinite memory).

It is an open question to determine the exact complexity of the approximation problem, and whether approximations can be computed in polynomial time, or if any complexity-theoretic lower bound can be established. We are not aware of any complexity lower bounds for approximation of the Positivity problem. Another direction for future work is to determine the memory requirement for pure strategies in MDPs. Figure 12 shows an MDP where memory is necessary in pure strategies. Consider expected stopping time $T = 8$, and the weight of states v_0, v_7, v_{14} is 0. A pure strategy that plays action a initially in v_0 and action b in the next visit to v_0 ensures expected value of 0 whereas the expected value of the two memoryless strategies (playing either always a , or always b) is negative. This example can be adapted to show that support-based⁴ strategies are not sufficient either.

Acknowledgment. The authors are grateful to the anonymous reviewers of LICS 2021 and of a previous version of this paper for insightful comments that helped improving the presentation. The research presented in this paper was partially supported by the grant ERC CoG 863818 (ForM-SMArt).

REFERENCES

- [AAOW15] S. Akshay, T. Antonopoulos, J. Ouaknine, and J. Worrell. Reachability problems for Markov chains. *Inf. Process. Lett.*, 115(2):155–158, 2015.
- [BCC⁺03] A. Biere, A. Cimatti, E. M. Clarke, O. Strichman, and Y. Zhu. Bounded model checking. *Advances in Computers*, 58:117–148, 2003.
- [BGB12] C. Baier, M. Größer, and N. Bertrand. Probabilistic ω -automata. *J. ACM*, 59(1):1, 2012.

⁴A strategy σ is support-based if it plays according to the current state and the support of the current state distribution. Formally, given a sequence $\rho = v_0 \dots v_k$, define $\text{last}(\rho) = v_k$, and for all $i \geq 1$ define $S_\sigma(i) = \{v \in V \mid \exists \rho \in V^i : \mathbb{P}^\sigma(\rho) > 0 \wedge \text{last}(\rho) = v\}$. Then σ is support-based if for all $\rho, \rho' \in V^+$, $S_\sigma(|\rho|) = S_\sigma(|\rho'|)$ and $\text{last}(\rho) = \text{last}(\rho')$ imply $\sigma(\rho) = \sigma(\rho')$.

- [BK08] C. Baier and J.-P. Katoen. *Principles of Model Checking*. MIT, 2008.
- [BKN⁺19] N. Balaji, S. Kiefer, P. Novotný, G. A. Pérez, and M. Shirmohammadi. On the complexity of value iteration. In *Proc. of ICALP: Automata, Languages, and Programming*, volume 132 of *LIPICs*, pages 102:1–102:15. Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik, 2019.
- [Can88] J. F. Canny. Some algebraic and geometric computations in PSPACE. In *Proc. of STOC: Symposium on Theory of Computing*, pages 460–467. ACM, 1988.
- [CD19] K. Chatterjee and L. Doyen. Graph planning with expected finite horizon. In *Proc. of LICS: Logic in Computer Science*, pages 1–13. IEEE, 2019.
- [CD21] K. Chatterjee and L. Doyen. Stochastic processes with expected stopping time. In *Proc. of LICS: Logic in Computer Science*, pages 1–13. IEEE, 2021.
- [CY95] C. Courcoubetis and M. Yannakakis. The complexity of probabilistic verification. *J. ACM*, 42(4):857–907, 1995.
- [dA97] L. de Alfaro. *Formal Verification of Probabilistic Systems*. PhD thesis, Stanford University, 1997.
- [DB12] I. Dew-Becker. *Essays on Time-Varying Discount Rates*. PhD thesis, Harvard University, 2012.
- [EMSS92] E. A. Emerson, A. K. Mok, A. P. Sistla, and J. Srinivasan. Quantitative temporal reasoning. *Real-Time Systems*, 4(4):331–352, 1992.
- [Erd89] P. Erdős. Ramanujan and I. In *Number Theory*, Lecture Notes in Mathematics 1395, pages 1–20. Springer, 1989.
- [FV97] J. Filar and K. Vrieze. *Competitive Markov Decision Processes*. Springer-Verlag, 1997.
- [Gal13] R. G. Gallager. *Stochastic processes: theory for applications*. Cambridge University Press, 2013.
- [HHH06] V. Halava, T. Harju, and M. Hirvensalo. Positivity of second order linear recurrent sequences. *Discrete Applied Mathematics*, 154(3):447–451, 2006.
- [How60] H. Howard. *Dynamic Programming and Markov Processes*. MIT, 1960.
- [KA04] Y. Kwon and G. Agha. Linear inequality LTL (iLTL): A model checker for discrete time Markov chains. In *Proc. of ICFEM: Formal Methods and Software Engineering*, LNCS 3308, pages 194–208. Springer, 2004.
- [KSK66] J. G. Kemeny, J. L. Snell, and A. W. Knapp. *Denumerable Markov chains*. D. Van Nostrand Company, 1966.
- [MHC03] O. Madani, S. Hanks, and A. Condon. On the undecidability of probabilistic planning and related stochastic optimization problems. *Artif. Intell.*, 147(1-2):5–34, 2003.
- [NR10] P. Norvig and S. J. Russell. *Artificial Intelligence - A Modern Approach (3rd ed.)*. Pearson Education, 2010.
- [OR94] M. J. Osborne and A. Rubinstein. *A Course in Game Theory*. MIT Press, 1994.
- [OW14] J. Ouaknine and J. Worrell. Positivity problems for low-order linear recurrence sequences. In *Proc. of SODA: Symposium on Discrete Algorithms*, pages 366–379. SIAM, 2014.
- [Paz71] A. Paz. *Introduction to probabilistic automata*. Academic Press, 1971.
- [PT87] C. H. Papadimitriou and J. N. Tsitsiklis. The complexity of Markov decision processes. *Math. Oper. Res.*, 12(3):441–450, 1987.
- [Put94] M. L. Puterman. *Markov decision processes*. Wiley and Sons, 1994.
- [Rab63] M. O. Rabin. Probabilistic automata. *Information & Control*, 6:3:230–245, 1963.
- [Rei99] R. Reisz. Decomposition theorems for probabilistic automata over infinite objects. *Informatica, Lithuanian Acad. Sci.*, 10:4:427–440, 1999.
- [SS04] T. Smith and R. G. Simmons. Heuristic search value iteration for POMDPs. In *Proc. of UAI: Uncertainty in Artificial Intelligence*, pages 520–527. AUAI Press, 2004.

APPENDIX A. BASIC INEQUALITIES

We recall basic inequalities that follow from the properties of the exponential and logarithm functions.

Lemma A.1. *For all $x \in \mathbb{R}$:*

- (1) *if $x \geq 1$, then $(1 - \frac{1}{x})^x < e^{-1} < \frac{1}{2}$, and*
- (2) *if $x < 1$, then $\ln(1 - x) \leq -x$.*

APPENDIX B. CONVERGENCE RATE IN MARKOV CHAINS

Let $\langle M, \mu, w \rangle$ be an aperiodic Markov chain (i.e., all recurrent classes are aperiodic). We show that there exist a vector π and numbers K_1, K_2 with $K_2 < 1$ such that for all $t \geq 0$ we have:

$$\|\mu \cdot M^t - \pi\|_\infty \leq K_1 \cdot K_2^t.$$

We recall the following results of [Gal13, Chapter 4]. In every recurrent class (or bottom scc) C of a Markov chain M , there is a steady-state vector π such that for all μ with $\text{Supp}(\mu) \subseteq C$, the vector $\mu \cdot M^t$ converges to π as $t \rightarrow \infty$. Moreover by [Gal13, Eq. (4.22)], for all vertices $i, j \in C$ and for all $t \geq 0$, we have

$$|M_{ij}^t - \pi_j| \leq \left(1 - 2\alpha^{n^2}\right)^{\lfloor t/n^2 \rfloor}, \quad (\text{B.1})$$

where $\alpha = \min\{M_{ij} \mid M_{ij} > 0\}$ is the smallest non-zero probability in M , and by [Gal13, Lemma (4.3.6)] for all $i \in \mathcal{T}$ where \mathcal{T} is the set of all transient vertices,

$$\sum_{j \in \mathcal{T}} M_{ij}^t \leq (1 - \alpha^n)^{\lfloor t/n \rfloor}, \quad (\text{B.2})$$

which gives a bound on the probability to remain in the set of transient vertices after t steps. For $i \in V$ and recurrent class $C \subseteq V$, let $\mathbb{P}_i(\diamond C)$ be the probability to eventually reach a vertex in C (and stay there forever since C is a recurrent class) from i . It directly follows from (B.2) that

$$\mathbb{P}_i(\diamond C) \geq \sum_{j \in C} M_{ij}^m \geq \mathbb{P}_i(\diamond C) - (1 - \alpha^n)^{\lfloor m/n \rfloor}. \quad (\text{B.3})$$

Now consider a Markov chain M with only aperiodic recurrent classes, and let \mathcal{T} be the set of transient vertices, and \mathcal{R} be the set of recurrent classes. The sequence $\mu \cdot M^t$ converges to a steady-state vector $\pi = \sum_{i \in V} \mu_i \cdot \sum_{C \in \mathcal{R}} \mathbb{P}_i(\diamond C) \cdot \pi^C$ where π^C is the steady-state vector of the class C , that is $\pi_j^C = \lim_{t \rightarrow \infty} M_{ij}^t$ (for arbitrary $i \in C$, and remember that the limit is independent of i). Let $u = \lfloor t/2 \rfloor$ and $B = \left(1 - 2\alpha^{n^2}\right)^{\lfloor (t-u)/n^2 \rfloor}$. Then for all $j \in V$,

$$\begin{aligned} & \left| \sum_{i \in V} \mu_i \cdot M_{ij}^t - \pi_j \right| \\ &= \left| \sum_{i \in V} \mu_i \cdot \left(\sum_{k \in \mathcal{T}} M_{ik}^u M_{kj}^{t-u} + \sum_{k \in C} M_{ik}^u M_{kj}^{t-u} - \mathbb{P}_i(\diamond C) \cdot \pi_j^C \right) \right| \end{aligned}$$

(where C is a recurrent class that contains j if j is recurrent, and C is an arbitrary recurrent class if j is transient, since then $M_{kj}^{t-u} = \pi_j^C = 0$ for all $k \in C$)

$$\begin{aligned} & \leq \left| \sum_{i \in V} \mu_i \cdot \left(\sum_{k \in \mathcal{T}} M_{ik}^u + \sum_{k \in C} M_{ik}^u (\pi_j + B) - \mathbb{P}_i(\diamond C) \cdot \pi_j^C \right) \right| \\ & \leq \left| \sum_{i \in V} \mu_i \cdot \left((1 - \alpha^n)^{\lfloor u/n \rfloor} + B + \sum_{k \in C} (M_{ik}^u - \mathbb{P}_i(\diamond C)) \cdot \pi_j^C \right) \right| \\ & \leq (1 - \alpha^n)^{\lfloor u/n \rfloor} + B + \left| \sum_{i \in V} \mu_i \cdot (1 - \alpha^n)^{\lfloor u/n \rfloor} \right| \\ & \leq 2(1 - \alpha^n)^{\lfloor u/n \rfloor} + (1 - 2\alpha^{n^2})^{\lfloor (t-u)/n^2 \rfloor} \\ & \leq 3(1 - \alpha^{n^2})^{t/3n^2} \end{aligned}$$

which proves the result (take $K_1 = 3$ and $K_2 = (1 - \alpha^{n^2})^{1/3n^2}$).