

Multi-weighted Markov Decision Processes with Reachability Objectives

Patricia Bouyer Mauricio González

LSV – CNRS & ENS Paris-Saclay – Université Paris-Saclay – France
{patricia.bouyer,mauricio.gonzalez}@lsv.fr

Nicolas Markey

IRISA – CNRS & Univ. Rennes & INRIA – France
nicolas.markey@irisa.fr

Mickael Randour

FNRS & UMONS – Belgium
mickael.randour@umonts.ac.be

In this paper, we are interested in the synthesis of schedulers in double-weighted Markov decision processes, which satisfy both a percentile constraint over a weighted reachability condition, and a quantitative constraint on the expected value of a random variable defined using a weighted reachability condition. This problem is inspired by the modelization of an electric-vehicle charging problem. We study the cartography of the problem, when one parameter varies, and show how a partial cartography can be obtained via two sequences of optimization problems. We discuss completeness and feasibility of the method.

1 Introduction

Importing formal methods in connected fields. Formal methods can help providing algorithmic solutions for control design. The electric-vehicle (EV) charging problem is an example of such an application area. This problem, usually presented as a control problem (see eg. [4]), can actually be modelled using Markov decision processes (MDP) [11, 15, 13]. Probabilities provide a way to model the non-flexible part of the energy network (consumption outside EV, for which large databases exist—and from which precise statistics can be extracted); we can then express an upper bound on the peak load as a safety condition (encoded as a reachability condition in our finite-horizon model), the constraint on the charging of all vehicles as a quantitative reachability objective, and various optimization criteria (e.g. minimizing the ageing of distribution transformers, or the energy price) as optimization of random cost variables.

Due to the specific form of the constructed model (basically acyclicity of the model), an ad-hoc method could be implemented using the tool PRISM, yielding interesting practical results as reported in [13] and in a forthcoming paper. However, the computability of an optimal strategy in a general MDP, as well as the corresponding decision problem, was unexplored.

Markov decision processes. MDPs have been studied for long [17, 12]. An MDP is a finite-state machine, on which a kind of game is played as follows. In each state, several decisions (a.k.a. actions) are available, each yielding a distribution over possible successor states. Once an action is selected, the next state is chosen probabilistically, following the distribution corresponding to the selected action. The game proceeds that way *ad infinitum*, generating an infinite play. The way actions are chosen is according to a *strategy* (also called *policy* in the context of MDPs). Rewards and/or costs can be associated to each action or edge, and various rules for aggregating individual rewards and costs encountered along a play can be applied to obtain various payoff functions. Examples of payoff functions include:

- sum up all the encountered costs (or rewards) along a play, until reaching some target (finite if the target is reached, infinite otherwise)—this is the so-called *truncated-sum* payoff;
- sum up all the encountered costs (or rewards) along a play with a discount factor at each step—this is the so-called *discounted-sum* payoff;
- average over the encountered costs (or rewards) along a play—this is the so-called *mean-payoff*.

Those payoff functions have been extensively studied in the literature; discounted-sum and mean-payoff have been shown to admit optimal memoryless and deterministic strategies, which can be computed using linear programming, yielding a polynomial-time algorithm. Alternative methods, such as value iteration or policy improvement, can be used in practice. On the other hand, the *shortest-path* problem, which aims at minimizing the truncated-sum, has been fully understood only recently [1]: one can decide in polynomial time as well whether the shortest-path is finite, or whether it is equal to $+\infty$ (if one cannot ensure reaching the target almost-surely), or whether it can be smaller than any arbitrary small number (if a negative loop can be enforced—note that in the context of stochastic systems, such a statement may be misleading, but it corresponds to a rough intuition), and corresponding strategies can be computed.

Multi-constrained problems in MDPs. The paradigm of multi-constrained objectives in stochastic systems in general, and in MDPs in particular, has recently arisen. It allows to express various (quantitative or qualitative) properties over the model, and to synthesize strategies accordingly. This new field of research is very rich and ambitious, with various types of objective combinations (see for instance [2, 18] for recent overviews). For recent developments on MDPs, one can cite:

- Pareto curves, or percentile queries, of multiple quantitative objectives: given several payoff functions, evaluate which tradeoff can be made between the probabilities, or the expectations, of the various payoff functions. In [8, 6, 9], solutions based on linear programming are provided for mean-payoff objectives. The percentile-query problem for various quantitative payoff functions is studied in [19].
- probability of conjunctions of objectives: given several payoff functions, evaluate the probability that all constraints are satisfied. This problem is studied in [14] for reachability (that is, for the truncated-sum payoff function), a PSPACE lower bound is proved for that problem, already with a single payoff function.
- the “beyond worst-case” paradigm: satisfy both a safety constraint on all outcomes, and various performance criteria. Variations of this problem for various payoff functions have been studied in [10, 7, 5].
- conditional expectations [3] or conditional values-at-risk [16], which measures likelihoods of properties under some assumptions on the system, have recently been investigated.

Our contributions. The general multi-constrained problem, arising from the EV-charging problem as modelled in [13], takes as an input an MDP with two weights, w_1 and w_2 , and requires the existence (and synthesis) of a strategy ensuring that some (absorbing) target state be reached, with a percentile constraint on the truncated sum of w_1 (lower bound parameterized by ε), and an expectation constraint on the truncated sum of w_2 . The initial EV-charging problem corresponds to the instance of that problem when $\varepsilon = 0$, where w_1 represents the energy that is used for charging and w_2 represents the ageing of the transformer.

As defined above, our problem integrates both the “beyond worst-case” paradigm of [7], and percentile queries as in [9] (mixing probabilities and expectations). While in [9] linear programs are used for solving percentile queries (heavily relying on the fact that mean-payoff objectives are tail objectives), we need different techniques since the truncated-sum payoff is very much prefix-dependent; actually, the PSPACE-lower bound of [14] immediately applies here as well (even without a constraint on the expectation of w_2). We develop here a methodology to describe the cartography of the problem, that is, the set of values of the parameter ε for which the problem has a solution. Our approach is based on two sequences of optimization problems which, in some cases we characterize, allow to have the (almost) full picture. We then discuss computability issues.

2 Preliminary definitions

Let S be a finite set. We write $\text{Dist}(S)$ for the set of distributions over S , that is, the set of functions $\delta: S \rightarrow [0, 1]$ such that $\sum_{s \in S} \delta(s) = 1$. A distribution over S is *Dirac* if $\delta(s) = 1$ for some $s \in S$.

2.1 Definition of the model

In this paper, we mainly focus on doubly-weighted Markov decision processes, but the technical developments mainly rely on simply-weighted Markov decision processes. We therefore define the setting with an arbitrary number of weights.

Definition 1. Let $k \in \mathbb{N}$. A k -weighted Markov decision process (kw -MDP) is a tuple $\mathcal{M} = (S, s_{\text{init}}, \odot, E, (w_i)_{1 \leq i \leq k})$, where:

- S is a finite set of states;
- $s_{\text{init}} \in S$ is the initial state;
- $\odot \in S$ is the target state;
- $E \subseteq S \times \text{Dist}(S)$ is a finite set of stochastic edges;
- for each $1 \leq i \leq k$, the function $w_i: S \times S \rightarrow \mathbb{Q}$ assigns a rational weight to each transition of the complete graph with state space S .

A (finite, infinite) path in \mathcal{M} from s is a (finite, infinite) sequence of states $s_0 s_1 s_2 \dots$ such that $s_0 = s$ and for every i , there is $\delta_i \in \text{Dist}(S)$ such that $(s_i, \delta_i) \in E$ and $\delta_i(s_{i+1}) > 0$. Finite paths are equivalently called histories. We write $\text{Paths}^{\mathcal{M}}(s)$ (resp. $\text{Paths}_{\infty}^{\mathcal{M}}(s)$) for the set of paths (resp. infinite paths), in \mathcal{M} from state s . Given a history $h = s_0 s_1 s_2 \dots s_N$ and $\ell \leq N$, the w_i -accumulated weight of h after ℓ steps is defined as $\text{Acc}_{w_i}^{\ell}(h) = \sum_{j=1}^{\ell} w_i(s_{j-1}, s_j)$. This notion extends straightforwardly to infinite paths.

A (randomized) *strategy* in \mathcal{M} is a function σ assigning to every history $h = s_0 s_1 s_2 \dots s_N$ a distribution over $s_N E = \{\delta \in \text{Dist}(S) \mid (s_N, \delta) \in E\}$. A strategy σ is said to be *pure* whenever the distributions it prescribes are Dirac. A path $s_0 s_1 s_2 \dots$ is an outcome of σ whenever for every strict prefix $s_0 s_1 s_2 \dots s_N$, there exists $\delta \in s_N E$ such that $\sigma(h)(\delta) > 0$ and $\delta(s_{N+1}) > 0$. Basically, the outcomes of a strategy are the paths that are activated by the strategy. We write $\text{out}^{\mathcal{M}}(\sigma, s)$ (resp. $\text{out}_{\infty}^{\mathcal{M}}(\sigma, s)$) for the set of finite (resp. infinite) outcomes of σ from state s .

Given a strategy σ and a state s , we denote with $\mathbb{P}_{\sigma, s}^{\mathcal{M}}$ the probability distribution, according to σ , over the infinite paths in $\text{Paths}_{\infty}^{\mathcal{M}}(s)$, defined in the standard way using cylinders based on finite paths from s . If f is a measurable functions from $\mathbb{P}_{\sigma, s}^{\mathcal{M}}$ to \mathbb{R} , we denote by $\mathbb{E}_{\sigma, s}^{\mathcal{M}}(f)$ the expected value of f w.r.t. the probability distribution $\mathbb{P}_{\sigma, s}^{\mathcal{M}}$, that is, $\mathbb{E}_{\sigma, s}^{\mathcal{M}}(f) = \int f \, d\mathbb{P}_{\sigma, s}^{\mathcal{M}}$. In all notations, we may omit to mention

the superscript \mathcal{M} when it is clear in the context, and may omit to mention the starting state s when it is s_{init} , so that \mathbb{P}_σ corresponds to $\mathbb{P}_{\sigma, s_{\text{init}}}^{\mathcal{M}}$.

Example 1. Consider the 2w-MDP \mathcal{M} of Figure 1. It has four states s_0, s_1, s_2 and \odot , and five edges, labelled with their names (here a, b, c, d and e). Weights label pairs of states (but are represented here only for pairs of states that may be activated). Edges a, b, d and e have Dirac distributions, while edge c has a stochastic choice (represented by the small black square). For readability we do not write the exact distributions, but in this example, they are assumed to be uniform.

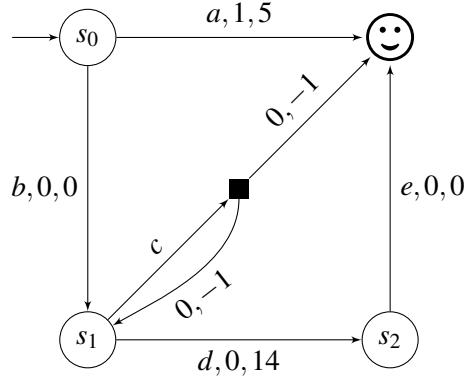


Figure 1: An example of a 2w-MDP

2.2 Payoff functions

We are interested in quantitative reachability properties (also called truncated-sum in the literature), which we formalize as follows. Let $\rho = s_0 s_1 s_2 \dots \in \text{Paths}(s_0)$. We use standard LTL-based notations for properties; for instance, we write $\rho \models \mathbf{F}\odot$ (resp. $\rho \models \mathbf{F}_I\odot$, when I is an interval of \mathbb{N}) when there is j (resp. $j \in I$) such that $s_j = \odot$, and $\rho \models \mathbf{G}\neg\odot$ (resp. $\rho \models \mathbf{G}_I\neg\odot$, when I is an interval of \mathbb{N}) when $s_j \neq \odot$ for every j (resp. for every $j \in I$). We will often use expressions $\sim N$ (in $\{<, \leq, =, \geq, >\} \times \mathbb{N}$ for defining intervals of \mathbb{N}).

If $\rho \models \mathbf{F}\odot$ and $1 \leq i \leq n$, we define the i -th payoff function $\text{TS}_{w_i}^\odot(\rho)$ as $\text{Acc}_{w_i}^N(\rho)$ where N is the least index such that $\rho \models \mathbf{F}_{=N}\odot$. If $\rho \not\models \mathbf{F}\odot$ then $\text{TS}_{w_i}^\odot(\rho) = +\infty$. The function $\rho \mapsto \text{TS}_{w_i}^\odot(\rho)$ is measurable, hence for every $\bowtie v$ in $\{<, \leq, =, \geq, >\} \times \mathbb{Q}$, $\mathbb{P}_{\sigma, s_0}(\{\rho \in \text{Paths}_\infty(s_0) \mid \text{TS}_{w_i}^\odot(\rho) \bowtie v\})$ (simply written as $\mathbb{P}_{\sigma, s_0}(\text{TS}_{w_i}^\odot \bowtie v)$) and $\mathbb{E}_{\sigma, s_0}(\text{TS}_{w_i}^\odot)$ are well-defined. We write $\rho \models (\text{TS}_{w_i}^\odot \bowtie v)$ whenever $\text{TS}_{w_i}^\odot(\rho) \bowtie v$.

In the rest of the paper, we assume that \odot is a sink state, and that there is a single loop on \odot whose weights are all equal to 0. This is w.l.o.g. since we will study payoff functions $\text{TS}_{w_i}^\odot$, which only consider the prefix up to the first visit to \odot .

Example 2. Consider again the example of Figure 1. Consider the strategy σ which selects a or b uniformly at random in s_0 , and always selects c in s_1 . Then,

$$\mathbb{P}_{\sigma, s_0}(\mathbf{F}\odot) = 1 \quad \mathbb{P}_{\sigma, s_0}(\text{TS}_{w_1}^\odot \geq 1) = \frac{1}{2} \quad \mathbb{E}_{\sigma, s_0}(\text{TS}_{w_2}^\odot) = \frac{1}{2} \cdot 5 + \frac{1}{2} \cdot \sum_{i=1}^{\infty} \frac{-i}{2^i} = 1 + \frac{1}{2}$$

3 The problem

The problem that we tackle in this paper arises from a recent study of an EV-charging problem [13]. The general problem we will define is a relaxed version of the original problem, combining several stochastic constraints (a percentile query over some payoff function and a constraint on the expectation of some payoff function) together with a worst-case obligation. While various payoff functions could be relevant, we focus on those payoff functions that were used for the EV-charging problem, that is, quantitative reachability (i.e., the truncated-sum payoff). We will see that the developed techniques are really specific to our choice of payoff functions.

In this paper, we focus on a combination of *sure reachability* of the goal state, of a percentile constraint on the proportion of paths having high value for the first payoff, and of a constraint on the expected value of the second payoff.

Let $\mathcal{M} = (S, s_{\text{init}}, \odot, E, (w_1, w_2))$ be a 2w-MDP. Let $v_1, v_2 \in \mathbb{Q}$. For every $\varepsilon \geq 0$, we define the problem $\text{Problem}_{\mathcal{M}, v_1, v_2}(\varepsilon)$ as follows: there exists a strategy σ_ε such that

1. for all $\rho \in \text{out}_\infty^{\mathcal{M}}(\sigma_\varepsilon, s_{\text{init}})$, it holds $\rho \models \mathbf{F} \odot$;
2. $\mathbb{P}_{\sigma_\varepsilon, s_{\text{init}}}^{\mathcal{M}} \left(\text{TS}_{w_1}^\odot \geq v_1 \right) \geq 1 - \varepsilon$;
3. $\mathbb{E}_{\sigma_\varepsilon, s_{\text{init}}}^{\mathcal{M}} \left(\text{TS}_{w_2}^\odot \right) < v_2$.

We aim at computing the values of ε for which $\text{Problem}_{\mathcal{M}, v_1, v_2}(\varepsilon)$ has a solution. For the rest of this section, we assume that there is a strategy σ such that $\mathbb{E}_{\sigma, s_{\text{init}}}^{\mathcal{M}} \left(\text{TS}_{w_2}^\odot \right) < v_2$. Otherwise $\text{Problem}_{\mathcal{M}, v_1, v_2}(\varepsilon)$ trivially has no solutions, for any ε . This can be decided using the algorithm recently developed in [1].

Example 3. *To illustrate the problem, we consider again the example given in Figure 1. Consider $\varepsilon = 0.5$, $v_1 = 1$ and $v_2 = 4.3$. The only way to satisfy the threshold constraint on w_1 is that at least half of the paths use a , impacting 2.5 over the expectation of w_2 . The other paths have to go to s_1 , and then take c for some time (provided the play goes back to s_1) in order to decrease the expectation of w_2 , before it becomes possible to take d and then e (so that the strategy is surely winning). This strategy uses both randomization (at s_0) and memory (counting the number of times c is taken before d can be taken).*

We call the *cartography* of our problem the function which associates to every $\varepsilon \in [0; 1]$, either true if $\text{Problem}_{\mathcal{M}, v_1, v_2}(\varepsilon)$ has a solution, or false otherwise. It is easily seen that the cartography is a threshold function, and can be characterized by an interval $I = \langle \gamma; 1 \rangle$ (which may be left-open or left-closed): In what follows, we describe an algorithmic technique to approximate this interval, and under additional conditions, to compute the bound γ . Whether the bound belongs to the interval remains open in general.

Link with the electric-vehicle (EV) charging problem. The (centralized) EV-charging problem consists in scheduling power loads within a time interval $[0; T]$ (T being a fixed time bound) with uncertain exogenous loads, so as to minimize the impact of loading on the electric distribution network (measured through the ageing of the transformer, which depends on the temperature of its winding). Following standard models, time is discretized, and the instantaneous energy consumption at time t can be written as the sum of the non-flexible load ℓ_t^{nf} (consumption outside EV) and the flexible load ℓ_t^{f} , corresponding to the EV charging. The flexible loads at each time are controllable actions, while the non-flexible part is known, or statistically estimated using past databases.

A first constraint on the transformer is given by its capacity: $\ell_t^{\text{f}} + \ell_t^{\text{nf}} \leq L^{\text{max}}$ (where L^{max} is a constant) for every $0 \leq t \leq T$. A second constraint represents the charge required for charging all vehicles on

schedule: $\sum_{t=0}^T \ell_t^f \geq \text{LoC}^{\max}$, where LoC^{\max} is a constant. The flexible load ℓ_t^f can thus be seen as a weight function w_1 .

While greedy solutions can be used to solve the above constraints, the ageing of transformer has not been taken into account so far. Using a standard model for the ageing of a transformer (see [4, 13] for details), it can be expressed as a weight function w_2 based on a discrete model in which states aggregate information on the system at the the two last timepoints. Globally, a 2w-MDP \mathcal{M} can be built, such that a controller for the EV-charging problem coincides with a solution to $\text{Problem}_{\mathcal{M}, \text{LoC}^{\max}, v_2}(0)$, for some bound v_2 for the expected ageing of the transformer.

4 Approximated cartography

We fix a 2w-MDP $\mathcal{M} = (S, s_0, \ominus, E, w_1, w_2)$ and two thresholds $v_1, v_2 \in \mathbb{Q}$. We introduce two simpler optimization problems related to $\text{Problem}_{\mathcal{M}, v_1, v_2}(\varepsilon)$, from which we derive informations on the good values of ε for which that problem has a solution. As we explain below, our approach is in general not complete. However, we observe that the true part of the cartography of our problem is an interval of the form $\langle \gamma; 1 \rangle$; under some hypotheses, we prove that our approach allows to approximate arbitrarily the bound γ , but may not be able to decide if the interval is left-open or left-closed.

4.1 Optimization problems

Let N be an integer. We write ϕ_N^+ for the property $\mathbf{F}_{\leq N} \ominus \wedge \text{TS}_{w_1}^{\ominus} \geq v_1$ (which specifies that the target is reached in no more than N steps, with a w_1 -weight larger than or equal to v_1), and ϕ_N^- for the property $\mathbf{F}_{\leq N} \ominus \wedge \text{TS}_{w_1}^{\ominus} < v_1$ (which means that the target is reached in no more than N steps, with a w_1 -weight smaller than v_1). We write ψ_N for the property $\mathbf{G}_{\leq N} \neg \ominus$ (the target is not reached during the N first steps). By extension, we write ϕ^+ , ϕ^- and ψ for the properties $\mathbf{F} \ominus \wedge \text{TS}_{w_1}^{\ominus} \geq v_1$, $\mathbf{F} \ominus \wedge \text{TS}_{w_1}^{\ominus} < v_1$ and $\mathbf{G} \neg \ominus$. Finally, we may also (abusively) use such formulas to denote the set of paths that satisfy them.

For every N and every path ρ of \mathcal{M} of length at least N , it holds that: $\rho \models \phi_N^+ \vee \phi_N^- \vee \psi_N$. Moreover, observe that $\phi_N^+ \subseteq \phi^+$ and $\phi^+ \subseteq \phi_N^+ \vee \psi_N$. As a consequence, for every N and every strategy σ , $\mathbb{P}_{\sigma}(\phi_N^+) \leq \mathbb{P}_{\sigma}(\phi^+) \leq \mathbb{P}_{\sigma}(\phi_N^+ \vee \psi_N)$.

4.1.1 First optimization problem

We define

$$\overline{\text{val}}_N = \inf \left\{ \mathbb{P}_{\sigma} \left(\phi_N^- \vee \psi_N \right) \mid \sigma \text{ s.t. } \mathbb{E}_{\sigma} \left(\text{TS}_{w_2}^{\ominus} \right) < v_2 \right\}$$

and for every $\alpha > 0$, we fix a witnessing strategy $\sigma_{N, \alpha}$ for $\overline{\text{val}}_N$ up to α (i.e. $\mathbb{E}_{\sigma_{N, \alpha}} \left(\text{TS}_{w_2}^{\ominus} \right) < v_2$ and $\mathbb{P}_{\sigma_{N, \alpha}} \left(\phi_N^- \vee \psi_N \right) \leq \overline{\text{val}}_N + \alpha$).

Remark. Note that, since we assume that there is a strategy σ such that $\mathbb{E}_{\sigma} \left(\text{TS}_{w_2}^{\ominus} \right) < v_2$, the constraint of this optimization problem is non-empty. Note also that if σ is a strategy such that $\mathbb{E}_{\sigma} \left(\text{TS}_{w_2}^{\ominus} \right) < v_2$, then $\mathbb{P}_{\sigma}(\mathbf{F} \ominus) = 1$, since for every path ρ , $\text{TS}_{w_2}^{\ominus}(\rho) = +\infty$ whenever $\rho \not\models \mathbf{F} \ominus$.

It is not hard to see that the sequence $(\overline{\text{val}}_N)_{N \in \mathbb{N}}$ is non-increasing (see Appendix). We let $\bar{\gamma} = \lim_{N \rightarrow +\infty} \overline{\text{val}}_N$. We then have:

Lemma 2. For every $\varepsilon < \bar{\gamma}$, $\text{Problem}_{\mathcal{M}, v_1, v_2}(\varepsilon)$ has no solution.

Proof. Fix $\varepsilon < \bar{\gamma}$, and assume towards a contradiction that $\text{Problem}_{\mathcal{M}, v_1, v_2}(\varepsilon)$ has a solution.

Fix a winning strategy σ_ε for $\text{Problem}_{\mathcal{M}, v_1, v_2}(\varepsilon)$. By the first winning constraint, there exists N_ε such that any outcome $\rho \in \text{out}_\infty(\sigma_\varepsilon)$ satisfies $\mathbf{F}_{\leq N_\varepsilon} \ominus$ (thanks to König's lemma). Furthermore, since $\mathbb{E}_{\sigma_\varepsilon}(\text{TS}_{w_2}^\ominus) < v_2$, then σ_ε belongs to the domain of the optimization problem defining $\overline{\text{val}}_{N_\varepsilon}$. Hence, we have

$$\overline{\text{val}}_{N_\varepsilon} \leq \mathbb{P}_{\sigma_\varepsilon}(\phi_{N_\varepsilon}^- \vee \psi_{N_\varepsilon}) = 1 - \mathbb{P}_{\sigma_\varepsilon}(\phi_{N_\varepsilon}^+) = 1 - \mathbb{P}_{\sigma_\varepsilon}(\phi^+) \leq \varepsilon.$$

This is then a contradiction with $\varepsilon < \bar{\gamma} \leq \overline{\text{val}}_{N_\varepsilon}$. Hence, we deduce that for every $\varepsilon < \bar{\gamma}$, $\text{Problem}_{\mathcal{M}, v_1, v_2}(\varepsilon)$ has no solution. \square

Lemma 3. *For every N , for every $\varepsilon > \overline{\text{val}}_N \geq \bar{\gamma}$, $\text{Problem}_{\mathcal{M}, v_1, v_2}(\varepsilon)$ has a solution.*

Proof. Let N be an integer, and $\varepsilon > \overline{\text{val}}_N$. Let σ_N be a strategy such that

$$\mathbb{P}_{\sigma_N}(\phi_N^- \vee \psi_N) < \varepsilon \text{ and } \mathbb{E}_{\sigma_N}(\text{TS}_{w_2}^\ominus) < v_2$$

Let σ_{Att} be an attractor (memoryless) strategy on \mathcal{M} , that is, a strategy which enforces reaching \ominus ; write M for a positive upper-bound on the accumulated weight w_2 when playing that strategy (from any state). For $k \geq N$, define σ_N^k as: play σ_N for the first k steps, and if \ominus is not reached, then play σ_{Att} . We show that we can find k large enough such that this strategy is a solution to $\text{Problem}_{\mathcal{M}, v_1, v_2}(\varepsilon)$.

The first condition is satisfied, since either the target state is reached during the k first steps (i.e. while playing σ_N), or it will be surely reached by playing σ_{Att} . Since $\mathbb{P}_{\sigma_N^k}(\phi_N^+) = 1 - \mathbb{P}_{\sigma_N^k}(\phi_N^- \vee \psi_N)$ and $\mathbb{P}_{\sigma_N^k}(\phi_N^+) \leq \mathbb{P}_{\sigma_N^k}(\phi^+)$, it is the case that $\mathbb{P}_{\sigma_N^k}(\phi^+) \geq 1 - \varepsilon$, which is the second condition for being a solution to $\text{Problem}_{\mathcal{M}, v_1, v_2}(\varepsilon)$. Finally, thanks to the law of total expectation, we can write:

$$\begin{aligned} \mathbb{E}_{\sigma_N^k}(\text{TS}_{w_2}^\ominus) &= \mathbb{E}_{\sigma_N^k}(\text{TS}_{w_2}^\ominus \mid \mathbf{F}_{\leq k} \ominus) \cdot \mathbb{P}_{\sigma_N^k}(\mathbf{F}_{\leq k} \ominus) + \mathbb{E}_{\sigma_N^k}(\text{TS}_{w_2}^\ominus \mid \mathbf{G}_{\leq k} \neg \ominus) \cdot \mathbb{P}_{\sigma_N^k}(\mathbf{G}_{\leq k} \neg \ominus) \\ &\leq \mathbb{E}_{\sigma_N^k}(\text{Acc}_{w_2}^k \mid \mathbf{F}_{\leq k} \ominus) \cdot \mathbb{P}_{\sigma_N^k}(\mathbf{F}_{\leq k} \ominus) + \mathbb{E}_{\sigma_N^k}(\text{Acc}_{w_2}^k + M \mid \mathbf{G}_{\leq k} \neg \ominus) \cdot \mathbb{P}_{\sigma_N^k}(\mathbf{G}_{\leq k} \neg \ominus) \\ &\quad \text{(since the global impact of playing the strategy } \sigma_{\text{Att}} \text{ is bounded by } M) \\ &= \mathbb{E}_{\sigma_N^k}(\text{Acc}_{w_2}^k) + M \cdot \mathbb{P}_{\sigma_N^k}(\mathbf{G}_{\leq k} \neg \ominus) \\ &\quad \text{(by linearity of expectation and the law of total expectation again)} \\ &= \mathbb{E}_{\sigma_N}(\text{Acc}_{w_2}^k) + M \cdot \mathbb{P}_{\sigma_N}(\mathbf{G}_{\leq k} \neg \ominus) \\ &\quad \text{(since } \sigma_N^k \text{ coincides with } \sigma_N \text{ on the } k \text{ first steps)} \end{aligned}$$

Now, since $\mathbb{E}_{\sigma_N}(\text{TS}_{w_2}^\ominus)$ is finite, it is the case that $\mathbb{P}_{\sigma_N}(\mathbf{F} \ominus) = 1$. Hence:

- $\lim_{k \rightarrow +\infty} \mathbb{E}_{\sigma_N}(\text{Acc}_{w_2}^k) = \mathbb{E}_{\sigma_N}(\text{TS}_{w_2}^\ominus)$, and
- $\lim_{k \rightarrow +\infty} \mathbb{P}_{\sigma_N}(\mathbf{G}_{\leq k} \neg \ominus) = 0$.

Let $\eta = v_2 - \mathbb{E}_{\sigma_N}(\text{TS}_{w_2}^\ominus) > 0$. One can choose k large enough such that

$$\left| \mathbb{E}_{\sigma_N}(\text{Acc}_{w_2}^k) - \mathbb{E}_{\sigma_N}(\text{TS}_{w_2}^\ominus) \right| < \eta/2 \quad \text{and} \quad \mathbb{P}_{\sigma_N}(\mathbf{G}_{\leq k} \neg \ominus) < \eta/2M.$$

We conclude that:

$$\mathbb{E}_{\sigma_N^k}(\text{TS}_{w_2}^\ominus) < \mathbb{E}_{\sigma_N}(\text{TS}_{w_2}^\ominus) + \eta < v_2.$$

The strategy σ_N^k therefore witnesses the fact that $\text{Problem}_{\mathcal{M}, v_1, v_2}(\varepsilon)$ has a solution. \square

4.1.2 Second optimization problem

We now define

$$\underline{\text{val}}_N = \inf \left\{ \mathbb{P}_\sigma(\phi_N^-) \mid \sigma \text{ s.t. } \mathbb{E}_\sigma(\text{TS}_{w_2}^\ominus) < v_2 \right\}.$$

Notice that for any N , $\underline{\text{val}}_N \leq \overline{\text{val}}_N$. For every $\alpha > 0$, we fix a witness strategy $\tilde{\sigma}_{N,\alpha}$ for $\underline{\text{val}}_N$ up to α (so that $\mathbb{E}_{\tilde{\sigma}_{N,\alpha}}(\text{TS}_{w_2}^\ominus) < v_2$ and $\mathbb{P}_{\tilde{\sigma}_{N,\alpha}}(\phi_N^-) \leq \underline{\text{val}}_N + \alpha$).

This time, it can be observed that the sequence $(\underline{\text{val}}_N)_{N \in \mathbb{N}}$ is non-decreasing. We let $\underline{\gamma} = \lim_{N \rightarrow +\infty} \underline{\text{val}}_N$. From the results and remarks above, we have $\underline{\text{val}}_N \leq \underline{\gamma} \leq \bar{\gamma}$ for any N . From Lemma 2, we get:

Lemma 4. *For any N and any $\varepsilon < \underline{\text{val}}_N \leq \underline{\gamma}$, $\text{Problem}_{\mathcal{M},v_1,v_2}(\varepsilon)$ has no solution.*

While the status of $\text{problem}_{\mathcal{M},v_1,v_2}(\bar{\gamma})$ is in general unknown, we still have the following properties:

- Proposition 1.**
- *If $\text{Problem}_{\mathcal{M},v_1,v_2}(\bar{\gamma})$ has a solution, then the sequence $(\overline{\text{val}}_N)_N$ is stationary and ultimately takes value $\bar{\gamma}$. The converse need not hold in general;*
 - *$\underline{\gamma} = \bar{\gamma}$ does neither imply that $\text{Problem}_{\mathcal{M},v_1,v_2}(\bar{\gamma})$ has a solution, nor that $\text{Problem}_{\mathcal{M},v_1,v_2}(\bar{\gamma})$ has no solution.*

4.1.3 Summary

Figure 2 summarizes the previous analysis. The picture seems rather complete, since only the status

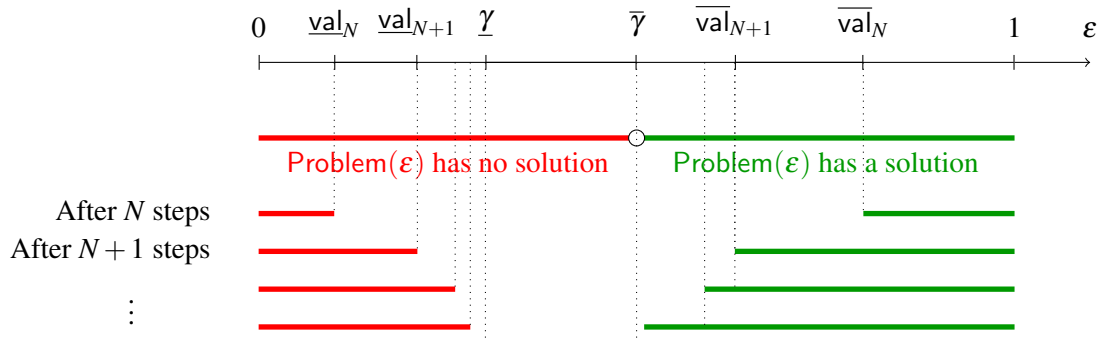


Figure 2: A partial cartography of our problem

of $\text{Problem}_{\mathcal{M},v_1,v_2}(\bar{\gamma})$ remains uncertain. However, it remains to discuss two things: first, the limits $\underline{\gamma}$ and $\bar{\gamma}$ are a priori unknown, hence the cartography is not effective so far. The idea is then to use the sequences $(\underline{\text{val}}_N)_N$ and $(\overline{\text{val}}_N)_N$ to approximate the limits. We will therefore discuss cases where the two limits coincide (we then say that the approach is *almost-complete*), allowing for a converging scheme and hence an algorithm to almost cover the interval $[0, 1]$ with either red (*there are no solutions*) or green (*there is a solution*), that is, to almost compute the full cartography of the problem. Second, we should discuss the effectiveness of the approach.

5 Almost-completeness of the approach

In this section, we discuss the almost-completeness of our approach, and describe situations where one can show that $\underline{\gamma} = \bar{\gamma} \stackrel{\text{def}}{=} \gamma$, which allows to reduce the unknown part of the cartography to the singleton $\{\gamma\}$.

The situations for completeness we describe below are conditions over cycles, either on weight w_1 or on weight w_2 . When we assume that cycles have positive w_i -weights, we mean it for every cycle, except for cycles containing \odot , which we assumed are self-loops with weight 0.

5.1 When all cycles have a positive w_2 -weight

In this subsection, we assume that the w_2 -weight of each cycle of \mathcal{M} is positive (this is the case for instance when the w_2 -weight of each edge is 1, i.e., when w_2 counts the number of steps). We let n be the number of states of \mathcal{M} .

Lemma 5. *There exists a constant $\kappa \geq 0$ such that, for any strategy σ satisfying $\mathbb{E}_\sigma(\text{TS}_{w_2}^\odot) < v_2$, and any $N > n$, it holds:*

$$\mathbb{P}_\sigma(\phi_N^- \vee \phi_N^+) \geq 1 - \frac{n}{N-n} \cdot \kappa.$$

Proof. Assuming otherwise, the impact of all runs that do not belong to $\phi_N^- \vee \phi_N^+$ would be too large for the constraint on $\text{TS}_{w_2}^\odot$. Indeed, applying the law of total expectation, we can write for every $N > n$:

$$\mathbb{E}_\sigma^\mathcal{M}(\text{TS}_{w_2}^\odot) = \mathbb{E}_\sigma^\mathcal{M}(\text{TS}_{w_2}^\odot \mid \mathbf{F}_{\leq N}^\odot) \cdot \mathbb{P}_\sigma^\mathcal{M}(\mathbf{F}_{\leq N}^\odot) + \mathbb{E}_\sigma^\mathcal{M}(\text{TS}_{w_2}^\odot \mid \mathbf{G}_{\leq N} \neg \odot) \cdot \mathbb{P}_\sigma^\mathcal{M}(\mathbf{G}_{\leq N} \neg \odot)$$

Write W_2 for the minimal (possibly negative) w_2 -weight appearing in \mathcal{M} , and c_2 for the minimal (positive by hypothesis) w_2 -weight of cycles in \mathcal{M} . Noticing that, along any path, at most n edges may be outside any cycle, we get

$$\mathbb{E}_\sigma^\mathcal{M}(\text{TS}_{w_2}^\odot \mid \mathbf{F}_{\leq N}^\odot) \geq n \cdot W_2$$

and

$$\mathbb{E}_\sigma^\mathcal{M}(\text{TS}_{w_2}^\odot \mid \mathbf{G}_{\leq N} \neg \odot) \geq n \cdot W_2 + \frac{N-n}{n} \cdot c_2.$$

We get:

$$\mathbb{E}_\sigma^\mathcal{M}(\text{TS}_{w_2}^\odot) \geq n \cdot W_2 \cdot (\mathbb{P}_\sigma^\mathcal{M}(\mathbf{F}_{\leq N}^\odot) + \mathbb{P}_\sigma^\mathcal{M}(\mathbf{G}_{\leq N} \neg \odot)) + \frac{N-n}{n} \cdot c_2 \cdot \mathbb{P}_\sigma^\mathcal{M}(\mathbf{G}_{\leq N} \neg \odot)$$

Since the left-hand side is strictly smaller than v_2 , we get

$$\mathbb{P}_\sigma^\mathcal{M}(\mathbf{G}_{\leq N} \neg \odot) = \mathbb{P}_\sigma^\mathcal{M}(\psi_N) = 1 - \mathbb{P}_\sigma^\mathcal{M}(\phi_N^+ \vee \phi_N^-) \leq \frac{n}{N-n} \cdot \left[\frac{v_2 - W_2 \cdot n}{c_2} \right].$$

□

Lemma 6. *For any constant κ satisfying Lemma 5, and any $N > n$, we have*

$$0 \leq \overline{\text{val}}_N - \underline{\text{val}}_N \leq \frac{n}{N-n} \cdot \kappa.$$

Proof. We already remarked that $\underline{\text{val}}_N \leq \overline{\text{val}}_N$. Now, from Lemma 5, for every strategy σ such that $\mathbb{E}_\sigma(\text{TS}_{w_2}^\odot) < v_2$, it holds for any $N > n$ that

$$\mathbb{P}_\sigma(\psi_N) < \frac{n}{N-n} \cdot \kappa.$$

Hence, for any strategy σ and any $N > n$,

$$\mathbb{P}_\sigma(\phi_N^- \vee \psi_N) = \mathbb{P}_\sigma(\phi_N^-) + \mathbb{P}_\sigma(\psi_N) \leq \mathbb{P}_\sigma(\phi_N^-) + \frac{n}{N-n} \cdot \kappa.$$

Taking the infimum over σ , first in the left-hand side, and then in the right-hand side, we get the expected bound. \square

Corollary 7. $\underline{\gamma} = \bar{\gamma}$.

Remark. Notice that the result does not hold without the assumption. Indeed, consider the 2w-MDP defined by the two deterministic edges $s \xrightarrow{a,0,0} s$ and $s \xrightarrow{b,-1,0} \ominus$, with $v_1 = 0$. Then, for every N , $\underline{\text{val}}_N = 0$, while $\overline{\text{val}}_N = 1$.

5.2 When all cycles have a positive w_1 -weight

We assume that each cycle of \mathcal{M} has a positive w_1 -weight. We first notice that:

Lemma 8. There exists an integer N_0 such that for every path ρ from s_{init} of length $N \geq N_0$ not visiting the goal state, it holds $\text{Acc}_{w_1}^N(\rho) \geq v_1$. In particular, if ρ satisfies $\mathbf{F}_{\geq N_0}(\neg \ominus \wedge \mathbf{F} \ominus)$, then $\text{TS}_{w_1}^\ominus(\rho) \geq v_1$.

Using this remark, we can prove:

Lemma 9. $\underline{\gamma} = \bar{\gamma}$

Proof. We fix the index N_0 as in Lemma 8. For any $N \geq N_0$ and any path ρ of length larger than N , we have

$$\rho \models \phi_N^- \iff \rho \models \phi_{N_0}^- \quad \text{and} \quad \rho \models \phi_N^+ \iff \rho \models \phi_{N_0}^+ \vee (\mathbf{G}_{\leq N_0} \neg \ominus \wedge \mathbf{F}_{(N_0;N]} \ominus).$$

From the first equivalence, we infer that for every $N \geq N_0$, $\underline{\text{val}}_N = \underline{\text{val}}_{N_0}$.

Let $N > N_0$, and write:

$$\begin{aligned} \overline{\text{val}}_N &= \inf \left\{ \mathbb{P}_\sigma(\phi_N^- \vee \psi_N) \mid \sigma \text{ s.t. } \mathbb{E}_\sigma(\text{TS}_{w_2}^\ominus) < v_2 \right\} \\ &= 1 - \sup \left\{ \mathbb{P}_\sigma(\phi_N^+) \mid \sigma \text{ s.t. } \mathbb{E}_\sigma(\text{TS}_{w_2}^\ominus) < v_2 \right\} \\ &= 1 - \sup \left\{ \mathbb{P}_\sigma(\phi_{N_0}^+ \vee (\mathbf{G}_{\leq N_0} \neg \ominus \wedge \mathbf{F}_{(N_0;N]} \ominus)) \mid \sigma \text{ s.t. } \mathbb{E}_\sigma(\text{TS}_{w_2}^\ominus) < v_2 \right\} \end{aligned}$$

We claim that:

Lemma 10.

$$\begin{aligned} \lim_{N \rightarrow +\infty} \sup \left\{ \mathbb{P}_\sigma(\phi_{N_0}^+ \vee (\mathbf{G}_{\leq N_0} \neg \ominus \wedge \mathbf{F}_{(N_0;N]} \ominus)) \mid \sigma \text{ s.t. } \mathbb{E}_\sigma(\text{TS}_{w_2}^\ominus) < v_2 \right\} \\ = \sup \left\{ \mathbb{P}_\sigma(\phi_{N_0}^+ \vee (\mathbf{G}_{\leq N_0} \neg \ominus \wedge \mathbf{F}_{>N_0} \ominus)) \mid \sigma \text{ s.t. } \mathbb{E}_\sigma(\text{TS}_{w_2}^\ominus) < v_2 \right\} \end{aligned}$$

From this lemma, we get that:

$$\begin{aligned}
\lim_{N \rightarrow +\infty} \overline{\text{val}}_N &= 1 - \sup \left\{ \mathbb{P}_\sigma \left(\phi_{N_0}^+ \vee \left(\mathbf{G}_{\leq N_0} \neg \odot \wedge \mathbf{F}_{> N_0} \odot \right) \right) \mid \sigma \text{ s.t. } \mathbb{E}_\sigma \left(\text{TS}_{w_2}^\odot \right) < v_2 \right\} \\
&= \inf \left\{ \mathbb{P}_\sigma \left(\phi_{N_0}^- \vee \mathbf{G} \neg \odot \right) \mid \sigma \text{ s.t. } \mathbb{E}_\sigma \left(\text{TS}_{w_2}^\odot \right) < v_2 \right\} \\
&\quad \text{(because a path not satisfying } \phi_{N_0}^+ \vee \left(\mathbf{G}_{\leq N_0} \neg \odot \wedge \mathbf{F}_{> N_0} \odot \right) \text{ satisfies } \phi_{N_0}^- \vee \mathbf{G} \neg \odot) \\
&= \inf \left\{ \mathbb{P}_\sigma \left(\phi_{N_0}^- \right) \mid \sigma \text{ s.t. } \mathbb{E}_\sigma \left(\text{TS}_{w_2}^\odot \right) < v_2 \right\} \\
&\quad \text{(since } \mathbb{P}_\sigma \left(\mathbf{G} \neg \odot \right) = 0 \text{ when } \mathbb{E}_\sigma \left(\text{TS}_{w_2}^\odot \right) < v_2) \\
&= \underline{\text{val}}_{N_0}
\end{aligned}$$

Hence we conclude that $\underline{\gamma} = \overline{\gamma}$ (and both limits are reached after finitely many steps). \square

Remark. This result requires the w_1 -positivity of cycles, as witnessed by the remark at the end of the previous section.

Also, one could think that assuming w_1 -negativity of cycles would be very similar, but this is not the case, as witnessed by the 2w-MDP defined by $s \xrightarrow{a,-1,0} s$ and $s \xrightarrow{b,1,0} \odot$. Then, for every $N \geq 2$, $\underline{\text{val}}_N = 0$ while $\overline{\text{val}}_N = 1$.

6 Effectiveness of the approach

We now explain how the two optimization problems can be solved. We first unfold our original 2w-MDP \mathcal{M} up to depth N as a tree, keeping a copy of \mathcal{M} below each leaf; write \mathcal{T}_N for this new 2w-MDP. There is a natural one-to-one mapping from paths in \mathcal{M} and paths in \mathcal{T}_N , from which we derive another one-to-one mapping ι_N between strategies in \mathcal{M} and strategies in \mathcal{T}_N . Furthermore two corresponding strategies assign the same probabilities and the same accumulated weights to the paths. As a consequence, for any N , any $\kappa_N \in \{\phi_N^+, \phi_N^-, \psi_N\}$, and any strategy σ in \mathcal{M} , we have

$$\mathbb{P}_\sigma^\mathcal{M}(\kappa_N) = \mathbb{P}_{\iota_N(\sigma)}^{\mathcal{T}_N}(\kappa_N).$$

We do not formalize the relation between \mathcal{M} and \mathcal{T}_N further, as it is rather straightforward.

Our two optimization problems can then be rephrased in \mathcal{T}_N as follows:

$$\overline{\text{val}}_N = \inf \left\{ \mathbb{P}_{\iota_N(\sigma)}^{\mathcal{T}_N} \left(\phi_N^- \vee \psi_N \right) \mid \sigma \text{ s.t. } \mathbb{E}_{\iota_N(\sigma)}^{\mathcal{T}_N} \left(\text{TS}_{w_2}^\odot \right) < v_2 \right\}$$

and

$$\underline{\text{val}}_N = \inf \left\{ \mathbb{P}_{\iota_N(\sigma)}^{\mathcal{T}_N} \left(\phi_N^- \right) \mid \sigma \text{ s.t. } \mathbb{E}_{\iota_N(\sigma)}^{\mathcal{T}_N} \left(\text{TS}_{w_2}^\odot \right) < v_2 \right\}.$$

From \mathcal{T}_N , we build the finite tree $\widehat{\mathcal{T}}_N$ as follows: we keep the first N levels of \mathcal{T}_N , add a fresh state \odot (at level $N+1$), and from each leaf at level N , corresponding to some state s of \mathcal{M} , we add an edge to \odot labelled by the w_2 -stochastic shortest path value from s in \mathcal{M} , that is, $\inf_\sigma \{ \mathbb{E}_{\sigma,s}^\mathcal{M}(\text{TS}_{w_2}^\odot) \}$. Those can be computed [1] (note that each can either be $-\infty$ or a finite value, or $+\infty$ if \odot cannot be almost-surely reached).

Every strategy σ_N in \mathcal{T}_N can then be partly mimicked in $\widehat{\mathcal{T}}_N$ (up to the N -th level of the tree); at level N , there is a single transition, which directly reaches \odot while increasing weight w_2 by the

shortest-path value mentioned above. We write $\widehat{\sigma}_N$ for this strategy in $\widehat{\mathcal{T}}_N$ derived from σ_N . Then we have $\mathbb{E}_{\widehat{\sigma}_N}^{\widehat{\mathcal{T}}_N}(\text{TS}_{w_2}^{\ominus}) \leq \mathbb{E}_{\sigma_N}^{\mathcal{T}_N}(\text{TS}_{w_2}^{\ominus})$, since the transitions from the nodes at level N to the \ominus state in $\widehat{\mathcal{T}}_N$ somehow acts as an “optimal” strategy after the first N levels.

Conversely, for every strategy $\widehat{\sigma}_N$ in $\widehat{\mathcal{T}}_N$ such that $\mathbb{E}_{\widehat{\sigma}_N}^{\widehat{\mathcal{T}}_N}(\text{TS}_{w_2}^{\ominus}) < v_2$:

- if $\mathbb{E}_{\widehat{\sigma}_N}^{\widehat{\mathcal{T}}_N}(\text{TS}_{w_2}^{\ominus}) = -\infty$, then for every $r \in \mathbb{R}$, one can extend $\widehat{\sigma}_N$ into a strategy $\sigma_{N,r}$ in \mathcal{T}_N such that $\mathbb{E}_{\sigma_{N,r}}^{\mathcal{T}_N}(\text{TS}_{w_2}^{\ominus}) < r$;
- otherwise, if $\mathbb{E}_{\widehat{\sigma}_N}^{\widehat{\mathcal{T}}_N}(\text{TS}_{w_2}^{\ominus}) \neq -\infty$, then one can extend $\widehat{\sigma}_N$ into a strategy σ_N in \mathcal{T}_N such that $\mathbb{E}_{\widehat{\sigma}_N}^{\widehat{\mathcal{T}}_N}(\text{TS}_{w_2}^{\ominus}) = \mathbb{E}_{\sigma_N}^{\mathcal{T}_N}(\text{TS}_{w_2}^{\ominus})$.

Hence the set of strategies $\widehat{\sigma}_N$ in $\widehat{\mathcal{T}}_N$ such that $\mathbb{E}_{\widehat{\sigma}_N}^{\widehat{\mathcal{T}}_N}(\text{TS}_{w_2}^{\ominus}) < v_2$ coincides with the set of strategies obtained as a pruning of a strategy σ_N in \mathcal{T}_N such that $\mathbb{E}_{\sigma_N}^{\mathcal{T}_N}(\text{TS}_{w_2}^{\ominus}) < v_2$. Our two optimization problems can then be rephrased as:

$$\overline{\text{val}}_N = \inf \left\{ \mathbb{P}_{\widehat{\sigma}_N}^{\widehat{\mathcal{T}}_N} \left(\phi_N^- \vee \psi_N \right) \mid \widehat{\sigma}_N \text{ s.t. } \mathbb{E}_{\widehat{\sigma}_N}^{\widehat{\mathcal{T}}_N} \left(\text{TS}_{w_2}^{\ominus} \right) < v_2 \right\}$$

and

$$\underline{\text{val}}_N = \inf \left\{ \mathbb{P}_{\widehat{\sigma}_N}^{\widehat{\mathcal{T}}_N} \left(\phi_N^- \right) \mid \widehat{\sigma}_N \text{ s.t. } \mathbb{E}_{\widehat{\sigma}_N}^{\widehat{\mathcal{T}}_N} \left(\text{TS}_{w_2}^{\ominus} \right) < v_2 \right\}$$

Since $\widehat{\mathcal{T}}_N$ is a (finite) tree, each strategy $\widehat{\sigma}_N$ in that MDP is memoryless, and can be represented as a probability value given to each edge appearing in the tree.

For each node n of $\widehat{\mathcal{T}}_N$, corresponding to some state s of \mathcal{M} , and for each edge $e = (s, \delta)$ from s , we consider a variable $p_{n,e}$, intended to represent the probability of taking edge e at node n . In particular, we will have the constraints $0 \leq p_{n,e} \leq 1$ and $\sum_{e=(s,\delta)} p_{n,e} = 1$. We write $\mathfrak{P} = (p_{n,e})_{n,e}$ for the tuple of all variables.

The two optimization problems above can then be written as:

$$\inf_{\mathfrak{P}} \left\{ P(\mathfrak{P}) \mid Q(\mathfrak{P}) < v_2 \wedge \bigwedge_{n,e} 0 \leq p_{n,e} \leq 1 \wedge \bigwedge_n \sum_{e=(s,\delta)} p_{n,e} = 1 \right\}$$

where $P(\mathfrak{P})$ and $Q(\mathfrak{P})$ are polynomials (of degree at most N).

Such polynomial optimization problems are in general hard to solve, and we have not been able to exploit the particular shape of our optimization problem to get efficient specialized algorithms. For each N , arbitrary under-approximations of $\underline{\text{val}}_N$ and over-approximations of $\overline{\text{val}}_N$ can be obtained by binary search, using the existential theory of the reals: the latter problem can be solved in polynomial space, but the number of variables of our problem is exponential in N . Using Lemmas 3 and 4, we get informations about the cartography of our problem. We can get approximations of $\underline{\gamma}$ and $\overline{\gamma}$ by iterating this procedure for larger values of N .

Remark. *In case the constraint on the expectation of w_2 can be relaxed (for instance if it is trivially satisfied), then the problem (over $\widehat{\mathcal{T}}_N$) becomes a simple optimal reachability problem in an MDP, for which pure strategies are sufficient (we have seen that this cannot be the case in our setting). The above optimization problem then simplifies into a linear-programming problem, with a much better complexity.*

7 The special case of Problem $\mathcal{M}, v_1, v_2(0)$

While the previous developments cannot give a solution to Problem $\mathcal{M}, v_1, v_2(0)$ since it requires not only to show that $\underline{\gamma} = \bar{\gamma} = 0$, but also that there is a solution to the limit point $\bar{\gamma}$ (which we do not have in general). We dedicate special developments to that problem. In this section, we assume (w.l.o.g.) that weights take integer values.

Problem $\mathcal{M}, v_1, v_2(0)$ can be rephrased as follows: *there exists a strategy σ_0 such that:*

1. for all $\rho \in \text{out}_{\infty}^{\mathcal{M}}(\sigma_0, s_{\text{init}})$, $TS_{w_1}^{\ominus}(\rho) \geq v_1$;
2. $\mathbb{E}_{\sigma_0, s_{\text{init}}}^{\mathcal{M}}(TS_{w_2}^{\ominus}) < v_2$.

Note that this problem is somehow a “beyond worst-case problem”, as defined in [19], with a strong constraint on all outcomes, and a stochastic constraint (here defined using expected value).

We describe a solution in the case all cycles of \mathcal{M} have non-negative w_1 -weights, which is inspired from [19, Theorem 13]. As we explain below, our solution extends to multiple weights (with non-negative cycles) with strong constraints (like the one for w_1). However, it is not correct when w_1 may have negative cycles as well. In that case, the status of the problem remains open.

We “unfold” the 2w-MDP $\mathcal{M} = (S, s_{\text{init}}, \ominus, E, w_1, w_2)$ into a 1w-MDP $\mathcal{N} = (Q, q_{\text{init}}, Q_{\ominus}, T, w)$, explicitly keeping track of w_1 in the states of \mathcal{N} :

- $Q = S \times \{-M, -M+1, \dots, 0, 1, \dots, M + \lfloor v_1 \rfloor + 1, \infty\}$, where $M = W \cdot (|S| + 1)$, W is the maximal absolute value of all weights $w_1(s, s')$ in \mathcal{M} , and $\lfloor v_1 \rfloor$ is the integral part of v_1 ;
- $q_{\text{init}} = (s_{\text{init}}, 0)$;
- $Q_{\ominus} = \{(\ominus, k) \mid k = \infty \text{ or } k \geq v_1\}$;
- $T = \left\{ ((s, c), (s', c')) \mid (s, s') \in E \text{ and } \begin{cases} c' = \infty & \text{if } c + w_1(s, s') > M + \lfloor v_1 \rfloor + 1 \text{ or } c = \infty \\ c' = c + w_1(s, s') & \text{otherwise} \end{cases} \right\}$;
- $w((s, c), (s', c')) = w_2(s, s')$.

There is a natural one-to-one correspondence λ between paths in \mathcal{M} and those in \mathcal{N} : $\lambda(s_0 s_1 \dots s_k \dots) = (s_0, 0)(s_1, c_1) \dots (s_k, c_k) \dots$ where, for every k , $c_k = c_{k-1} + w_1(s_{k-1}, s_k)$ if this value is less than or equal to $M + \lfloor v_1 \rfloor + 1$, and $c_k = \infty$ otherwise. Notice that thanks to our hypothesis on cycles, c_k may never be less than $-M$.

Also, by construction, if $(s_0, 0)(s_1, c_1) \dots (s_k, c_k) \dots$ is a path in \mathcal{N} such that $c_k = \infty$, then for every $j \geq k$, $c_j = \infty$; in that case, in the corresponding path $s_0 s_1 \dots s_k \dots$ in \mathcal{M} , it is the case for every $j \geq k$ that $\text{Acc}_{w_1}^j(s_0 s_1 \dots s_k \dots) \geq v_1$ (indeed, once the accumulated weight has become larger than $M + v_1$, it can never be smaller than v_1 again, thanks to the hypothesis on cycles). Conversely, if $c_k < \infty$, then $\text{Acc}_{w_1}^k(s_0 s_1 \dots s_k \dots) = c_k$.

From that correspondence over paths, strategies in \mathcal{M} can equivalently be seen as strategies in \mathcal{N} via a mapping ι . Using this correspondence:

Lemma 11. *There is a solution to Problem $\mathcal{M}, v_1, v_2(0)$ if, and only if, in \mathcal{N} there is a strategy τ such that:*

1. for all $\rho \in \text{out}_{\infty}^{\mathcal{N}}(\tau, q_{\text{init}})$, $\rho \models \mathbf{F}Q_{\ominus}$;
2. $\mathbb{E}_{\tau, q_{\text{init}}}^{\mathcal{N}}(TS_w^{\ominus}) < v_2$.

Furthermore, if σ is a solution to Problem $\mathcal{M}, v_1, v_2(0)$, then $\tau = \iota(\sigma)$ is a solution to the above problem in \mathcal{N} ; and if τ is a solution to the above problem in \mathcal{N} , then $\sigma = \iota^{-1}(\tau)$ is a solution to Problem $\mathcal{M}, v_1, v_2(0)$.

We can further show that $\text{Problem}_{\mathcal{M}, v_1, v_2}(0)$ has a solution if, and only if, the stochastic w_2 -shortest-path for reaching Q_{\ominus} in \mathcal{N} from s_{init} is smaller than v_2 . This latter problem can be decided in PTIME [1]. However, the size of \mathcal{N} is exponential (more precisely, it is pseudo-polynomial in the size of \mathcal{M}). In the end:

Theorem 12. *One can decide $\text{Problem}_{\mathcal{M}, v_1, v_2}(0)$ in pseudo-polynomial time, when cycles of \mathcal{M} have a non-negative w_1 -weight.*

Remark. *Notice that we could also have assumed that all cycles have non-positive w_1 -weight: the construction of \mathcal{N} would be similar, but with state space $S \times \{-\infty, \lfloor v_1 \rfloor - 1 - M', \dots, 0, 1, \dots, M' + 1\}$ where $M' = W' \cdot (|S| + 1)$ and W' is the largest absolute value w_1 -weight in \mathcal{M} . The rest of the argumentation follows the same ideas as above.*

Notice also that our algorithm is readily adapted to the case where we may have several constraints similar to those on w_1 (for each extra variable, one should assume that either each cycle is non-negative, or each cycle is non-positive). It suffices to keep track of the extra weights in the states, and take as target all states where the constraints are fulfilled.

8 Conclusion

In this paper, we investigated a multi-constrained reachability problem over MDPs, which originated in the context of electric-vehicle charging [13]. This problem consists in finding a strategy that surely reaches a quantitative goal (e.g., *all vehicles are fully charged and the load of the network remains below a given bound at any time*) while satisfying a condition on the expected value of some variable (*the life expectancy of the transformer is high or the expected cost of charging all vehicles is minimized*). We developed partial solutions to the problem by providing a cartography of the solutions to (a relaxed version of) the problem. We identified realistic conditions under which the cartography is (almost) complete. However, even under these conditions, the general decision problem (given ε , does $\text{Problem}(\varepsilon)$ have a solution?) remains open so far. Also, the case of MDPs not satisfying these conditions remains also open, but we believe that our approximation techniques may give interesting informations which suffice for practical applications such as electric-vehicle charging.

Our approach for $\text{Problem}(0)$, which amounts to explicitly keep track of the worst-case constraint on w_1 , immediately extends to multiple weights with worst-case constraints (with the same assumptions on cycles—note that the more general setting could not be solved, which has to be put in parallel with the undecidability result of [19, Theorem 12]) for the multi-dimensional percentile problem for truncated sum payoffs. The cartography for the relaxed problem $\text{Problem}(\varepsilon)$ requires solving sequences of intermediary optimization problems, which can be expressed as polynomial optimization problems with polynomial constraints. It could be extended to several such weights as well (either by putting an assumption on the w_1 -weights of every cycle, or just on the w_2 -weight of every cycle). A nice continuation of our work would consist in computing (approximations of) Pareto-optimal solutions in such a setting. Improving the complexity and practicality of our approach is also on our agenda for future work.

Acknowledgement. We thank the anonymous reviewers for their careful reading of our submission.

Patricia Bouyer, Mauricio González and Nicolas Markey are supported by ERC project EQUALLIS. Mickael Randour is an F.R.S.-FNRS Research Associate, and he is supported by the F.R.S.-FNRS Incentive Grant ManySynth.

References

- [1] Christel Baier, Nathalie Bertrand, Clemens Dubslaff, Daniel Gburek & Ocan Sankur (2018): *Stochastic Shortest Paths and Weight-Bounded Properties in Markov Decision Processes*. In: *LICS'18*, IEEE, doi:10.1145/3209108.3209184.
- [2] Christel Baier, Clemens Dubslaff & Sascha Klüppelholz (2014): *Trade-off analysis meets probabilistic model checking*. In: *CSL-LICS'14*, ACM, pp. 1:1–1:10, doi:10.1145/2603088.2603089.
- [3] Christel Baier, Joachim Klein, Sascha Klüppelholz & Sascha Wunderlich (2017): *Maximizing the Conditional Expected Reward for Reaching the Goal*. In: *TACAS'17, LNCS 10206*, Springer, pp. 269–285, doi:10.1007/978-3-662-54580-5_16.
- [4] Olivier Beaude, Samson Lasaulce, Martin Hennebel & Ibrahim Mohand-Kaci (2016): *Reducing the Impact of EV Charging Operations on the Distribution Network*. *IEEE Trans. Smart Grid* 7(6), pp. 2666–2679, doi:10.1109/TSG.2015.2489564.
- [5] Raphaël Berthon, Mickael Randour & Jean-François Raskin (2017): *Threshold Constraints with Guarantees for Parity Objectives in Markov Decision Processes*. In: *ICALP'17, LIPIcs 80, LZI*, pp. 121:1–121:15, doi:10.4230/LIPIcs.ICALP.2017.121.
- [6] Tomáš Brázdil, Václav Brožek, Krishnendu Chatterjee, Vojtěch Forejt & Antonín Kučera (2014): *Markov Decision Processes with Multiple Long-Run Average Objectives*. *LMCS* 10(1:13), pp. 1–29, doi:10.2168/LMCS-10(1:13)2014.
- [7] Véronique Bruyère, Emmanuel Filiot, Mickael Randour & Jean-François Raskin (2017): *Meet your expectations with guarantees: Beyond worst-case synthesis in quantitative games*. *Inf. & Comp.* 254, pp. 259–295, doi:10.1016/j.ic.2016.10.011.
- [8] Krishnendu Chatterjee (2007): *Markov Decision Processes with Multiple Long-Run Average Objectives*. In: *FSTTCS'07, LNCS 4855*, Springer, pp. 473–484, doi:10.1007/978-3-540-77050-3_39.
- [9] Krishnendu Chatterjee, Zuzana Křetínská & Jan Křetínský (2017): *Unifying two views on multiple mean-payoff objectives in Markov decision processes*. *LMCS* 13(2:15), pp. 1–50, doi:10.23638/LMCS-13(2:15)2017.
- [10] Lorenzo Clemente & Jean-François Raskin (2015): *Multidimensional beyond worst-case and almost-sure problems for mean-payoff objectives*. In: *LICS'15*, IEEE, pp. 257–268, doi:10.1109/LICS.2015.33.
- [11] Jonathan Donadee & Marija D. Ilic (2014): *Stochastic Optimization of Grid to Vehicle Frequency Regulation Capacity Bids*. *IEEE Trans. on Smart Grid* 5(2), pp. 1061–1069, doi:10.1109/TSG.2013.2290971.
- [12] Jerzy Filar & Koos Vrieze (1997): *Competitive Markov Decision Processes*. Springer, doi:10.1007/978-1-4612-4054-9.
- [13] Mauricio González, Olivier Beaude, Patricia Bouyer, Samson Lasaulce & Nicolas Markey (2017): *Stratégies d'ordonnancement de consommation d'énergie en présence d'information imparfaite de prévision*. In: *GRETSI'17*. Available at <http://www.lsv.fr/Publis/PAPERS/PDF/GBBLM-gretsi17.pdf>.
- [14] Christoph Haase & Stefan Kiefer (2015): *The Odds of Staying on Budget*. In: *ICALP'15, LNCS 9135*, Springer, pp. 234–246, doi:10.1007/978-3-662-47666-6_19.
- [15] Daniel R. Jiang & Warren B. Powell (2016): *Practicality of Nested Risk Measures for Dynamic Electric Vehicle Charging*. Research Report 1605.02848, arXiv.
- [16] Jan Křetínský & Tobias Meggendorfer (2018): *Conditional Value-at-Risk for Reachability and Mean Payoff in Markov Decision Processes*. Research Report 1805.02946, arXiv.
- [17] Martin L. Puterman (1994): *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley and Sons, doi:10.1002/9780470316887.
- [18] Mickael Randour, Jean-François Raskin & Ocan Sankur (2015): *Variations on the Stochastic Shortest Path Problem*. In: *VMCAI'15, LNCS 8931*, Springer, pp. 1–18, doi:10.1007/978-3-662-46081-8_1.
- [19] Mickael Randour, Jean-François Raskin & Ocan Sankur (2017): *Percentile queries in multi-dimensional Markov decision processes*. *FMSD* 50(2-3), pp. 207–248, doi:10.1007/978-3-319-21690-4_8.