
STRATEGY COMPLEXITY OF POINT PAYOFF, MEAN PAYOFF AND TOTAL PAYOFF OBJECTIVES IN COUNTABLE MDPs

RICHARD MAYR AND ERIC MUNDAY

University of Edinburgh, School of Informatics, LFCS, 10 Crichton Street, Edinburgh EH89AB, UK

ABSTRACT. We study countably infinite Markov decision processes (MDPs) with real-valued transition rewards. Every infinite run induces the following sequences of payoffs: 1. Point payoff (the sequence of directly seen transition rewards), 2. Mean payoff (the sequence of the sums of all rewards so far, divided by the number of steps), and 3. Total payoff (the sequence of the sums of all rewards so far). For each payoff type, the objective is to maximize the probability that the \liminf is non-negative.

We establish the complete picture of the strategy complexity of these objectives, i.e., how much memory is necessary and sufficient for ε -optimal (resp. optimal) strategies. Some cases can be won with memoryless deterministic strategies, while others require a step counter, a reward counter, or both.

1. INTRODUCTION

Background. Countably infinite Markov decision processes (MDPs) are a standard model for dynamic systems that exhibit both stochastic and controlled behavior; see, e.g., standard textbooks [27, 24, 15, 28] and references therein. Some fundamental results and proof techniques for countable MDPs were established in the framework of Gambling Theory [15, 24]. See also Ornstein’s seminal paper on stationary strategies [26]. Further applications include control theory [7, 1], operations research and finance [25, 3, 5, 30], artificial intelligence and machine learning [33, 31], and formal verification [21, 2, 9, 16, 8, 17, 14, 4]. The latter works often use countable MDPs to describe unbounded structures in computational models such as stacks/recursion, counters, queues, etc.

An MDP is a directed graph where states are either random or controlled. In a random state the next state is chosen according to a fixed probability distribution. In a controlled state the controller can choose a distribution over all possible successor states. By fixing a strategy for the controller (and an initial state), one obtains a probability space of runs of the MDP. The goal of the controller is to optimize the expected value of some objective function on the runs. The type of strategy necessary to achieve an ε -optimal (resp. optimal) value for a given objective is called its *strategy complexity*.

Extended version of results presented at CONCUR 2021.

Key words and phrases: Markov decision processes, Strategy complexity, Mean payoff.

Transition rewards and liminf objectives. MDPs are given a reward structure by assigning a real-valued (resp. integer or rational) reward to each transition. Every run then induces an infinite sequence of seen transition rewards $r_0 r_1 r_2 \dots$. We consider the lim inf of this sequence, as well as two other important derived sequences.

- (1) The point payoff considers the lim inf of the sequence $r_0 r_1 r_2 \dots$ directly.
- (2) The mean payoff considers the lim inf of the sequence $\left\{ \frac{1}{n} \sum_{i=0}^{n-1} r_i \right\}_{n \in \mathbb{N}}$, i.e., the mean of all rewards seen so far in an expanding prefix of the run.
- (3) The total payoff considers the lim inf of the sequence $\left\{ \sum_{i=0}^{n-1} r_i \right\}_{n \in \mathbb{N}}$, i.e., the sum of all rewards seen so far.

For each of the three cases above, the lim inf threshold objective is to maximize the probability that the lim inf of the respective type of sequence is ≥ 0 .

Our contribution. We establish the strategy complexity of all the lim inf threshold objectives above for *countably infinite* MDPs. (For the simpler case of finite MDPs, see the paragraph on related work below.) We show the amount and type of memory that is necessary and sufficient for ε -optimal strategies (and optimal strategies, where they exist).

Classes of strategies are defined via the amount and type of memory used, and whether they are randomized or deterministic. Some canonical types of memory for strategies are the following: No memory (also called memoryless or positional), finite memory, a step counter (i.e., a discrete clock), a reward counter (i.e., a variable that records the sum of all transition rewards seen so far) and general infinite memory. Strategies using only a step counter are also called *Markov strategies* [27]. The reward counter has the same type as the transition rewards in the MDP, i.e., integers, rationals or reals. Moreover, there can be combinations of these, e.g., a step counter plus some finite general purpose memory. Other types of memory are possible, e.g., an unbounded stack or a queue, but they are less common in the literature.

To establish an upper bound X on the strategy complexity of an objective in countable MDPs, it suffices to prove that there always exist good (ε -optimal, resp. optimal) strategies in some class of strategies X . Lower bounds on the strategy complexity of an objective can only be established in the sense of proving that good strategies for the objective do not exist in some classes Y , Z , etc. See Figure 1 for an example.

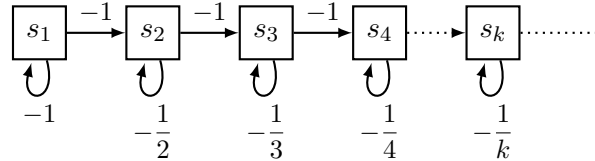


Figure 1: Adapted from [27, Example 8.10.2]. While there is no optimal MD (memoryless deterministic) strategy, the following strategy is optimal for lim inf/lim sup mean payoff: Loop $\exp(\exp(k))$ many times in state s_k for all k . In this particular example, this can be implemented with either just a step counter or just a reward counter, but in general both are needed; cf. Table 1.

Note that different classes of strategies are not always comparable, for several reasons. First, different types of memory may be incomparable. E.g., a step counter uses infinite memory, but it is updated in a very particular way, and thus it does not subsume a finite

general purpose memory. Second, randomized strategies are more general than deterministic ones if they use the same memory, but not if the memory is different. E.g., randomized positional strategies are incomparable to deterministic strategies with finite memory (or a step counter). Since strategy classes are not always comparable, there can be cases with several incomparable upper/lower bounds.

Moreover, there is no weakest type of infinite memory with restricted use. Hence, upper and lower bounds on the strategy complexity of an objective can be only be tight *relative* to the considered alternative strategy classes, e.g., the canonical classes mentioned above.

Our results are summarized in Table 1. By $\text{Rand}(X)$ (resp. $\text{Det}(X)$) we denote the classes of randomized (resp. deterministic) strategies that use memory of size/type X . SC denotes a step counter, RC denotes a reward counter and F denotes arbitrary finite memory. Positional/memoryless means that no memory is used. The simplest type are memoryless deterministic (MD) strategies. The results depend on the type of objective (point, total, or mean payoff) and on whether the MDP is finitely or infinitely branching. For our objectives, the strategy complexities of ε -optimal and optimal strategies (where they exist) coincide, but the proofs are different.

For clarity of presentation, our counterexamples use large transition rewards and high degrees of branching. However, the results can be strengthened such that the lower bounds hold even for just binary branching MDPs with rational transition probabilities and transition rewards in $\{-1, 0, 1\}$; cf. Section 6.

	Point payoff	Mean payoff	Total payoff
ε -opt., fin. branch.	Det(Positional) 3.5	Det(SC+RC) 4.1, 4.12, 4.15	Det(RC) 5.1, 5.7
opt., fin. branch.	Det(Positional) 3.8	Det(SC+RC) 4.2, 4.20, 4.23	Det(RC) 5.4, 5.8
ε -opt., inf. branch.	Det(SC) 3.10, 3.13	Det(SC+RC) 4.1, 4.12, 4.15	Det(SC+RC) 5.2, 5.7, 5.6
opt., inf. branch.	Det(SC) 3.12, 3.13	Det(SC+RC) 4.2, 4.20, 4.23	Det(SC+RC) 5.5, 5.8, 5.6

Table 1: Combined upper and lower bounds on the strategy complexity of ε -optimal (resp. optimal) strategies for point, mean and total payoff objectives in finitely branching and infinitely branching MDPs. All strategies are deterministic and randomization does not help. For each result, we list the numbers of the theorems that show the upper and lower bounds on the strategy complexity. The lower bounds hold wrt. the canonical strategies mentioned above. Explicit lower bounds are listed in the tables in the following sections. The lower bounds hold even for integer transition rewards. The upper bounds hold even for real-valued transition rewards.

Some complex new proof techniques are developed to show these results. E.g., the examples showing the lower bound in cases where both a step counter and a reward counter are required use a finely tuned tradeoff between different risks that can be managed with both counters, but not with just one counter plus arbitrary finite memory. The strategies showing the upper bounds need to take into account convergence effects, e.g., the sequence of point rewards $-1/2, -1/3, -1/4, \dots$ *does* satisfy $\liminf \geq 0$, i.e., one cannot assume that rewards are integers.

Related work. Mean payoff objectives for *finite* MDPs have been widely studied; cf. survey in [11]. There exist optimal MD strategies for \liminf mean payoff (which are also optimal for \limsup mean payoff since the transition rewards are bounded), and the associated computational problems can be solved in polynomial time [11, 27]. Similarly, see [12]

for a survey on lim sup and lim inf point payoff objectives in finite stochastic games and MDPs, where there also exist optimal MD strategies, and the more recent paper by Flesch, Predtetchinski and Sudderth [18] on simplifying optimal strategies.

All this does *not* carry over to countably infinite MDPs. Optimal strategies need not exist (not even for much simpler objectives), (ε) -optimal strategies can require infinite memory, and computational problems are not defined in general, since a countable MDP need not be finitely presented [27, 24, 15, 28, 23]. Moreover, attainment for lim inf mean payoff need not coincide with attainment for lim sup mean payoff, even for very simple examples. E.g., consider the acyclic infinite graph with transitions $s_n \rightarrow s_{n+1}$ for all $n \in \mathbb{N}$ with reward $(-1)^n 2^n$ in the n -th step, which yields a lim inf mean payoff of $-\infty$ and a lim sup mean payoff of $+\infty$.

Mean payoff objectives for countably infinite MDPs have been considered in [27, Section 8.10], e.g., [27, Example 8.10.2] (adapted in Figure 1) shows that there are no optimal MD (memoryless deterministic) strategies for lim inf/lim sup mean payoff. [29, Counterexample 1.3] shows that there are not even ε -optimal memoryless randomized strategies for lim inf/lim sup mean payoff. (We show much stronger lower/upper bounds; cf. Table 1.)

Sudderth [32] considered an objective on countable MDPs that is related to our point payoff threshold objective. However, instead of maximizing the probability that the lim inf/lim sup is non-negative, it asks to maximize the *expectation* of the lim inf/lim sup point payoffs, which is a different problem (e.g., it can tolerate a high probability of a negative lim inf/lim sup if the remaining cases have a huge positive lim inf/lim sup). Hill & Pestien [19] showed the existence of good randomized Markov strategies for the lim sup of the *expected* average reward up-to step n for growing n , and for the *expected* lim inf of the point payoffs.

2. PRELIMINARIES

Markov decision processes. A *probability distribution* over a countable set S is a function $f : S \rightarrow [0, 1]$ with $\sum_{s \in S} f(s) = 1$. We write $\mathcal{D}(S)$ for the set of all probability distributions over S . A *Markov decision process* (MDP) $\mathcal{M} = (S, S_\square, S_\circ, \longrightarrow, P, r)$ consists of a countable set S of *states*, which is partitioned into a set S_\square of *controlled states* and a set S_\circ of *random states*, a *transition relation* $\longrightarrow \subseteq S \times S$, and a *probability function* $P : S_\circ \rightarrow \mathcal{D}(S)$. We write $s \longrightarrow s'$ if $(s, s') \in \longrightarrow$, and refer to s' as a *successor* of s . We assume that every state has at least one successor. The probability function P assigns to each random state $s \in S_\circ$ a probability distribution $P(s)$ over its (non-empty) set of successor states. A *sink* in \mathcal{M} is a subset $T \subseteq S$ closed under the \longrightarrow relation, that is, $s \in T$ and $s \longrightarrow s'$ implies that $s' \in T$.

An MDP is *acyclic* if the underlying directed graph (S, \longrightarrow) is acyclic, i.e., there is no directed cycle. It is *finitely branching* if every state has finitely many successors and *infinitely branching* otherwise. An MDP without controlled states ($S_\square = \emptyset$) is called a *Markov chain*.

In order to specify our point/mean/total payoff objectives (see below), we define a function $r : S \times S \rightarrow \mathbb{R}$ that assigns numeric rewards to transitions.

Strategies and Probability Measures. A *run* ρ is an infinite sequence of states and transitions $s_0 e_0 s_1 e_1 \dots$ such that $e_i = (s_i, s_{i+1}) \in \longrightarrow$ for all $i \in \mathbb{N}$. Let $Runs_{\mathcal{M}}^{s_0}$ be the set of all runs from s_0 in the MDP \mathcal{M} . A *partial run* is a finite prefix of a run, $pRuns_{\mathcal{M}}^{s_0}$ is the set of all partial runs from s_0 and $pRuns_{\mathcal{M}}$ the set of partial runs from any state.

We write $\rho_s(i) \stackrel{\text{def}}{=} s_i$ for the i -th state along ρ and $\rho_e(i) \stackrel{\text{def}}{=} e_i$ for the i -th transition along ρ . We sometimes write runs as $s_0 s_1 \dots$, leaving the transitions implicit. We say that a (partial) run ρ *visits* s if $s = \rho_s(i)$ for some i , and that ρ starts in s if $s = \rho_s(0)$.

A *strategy* is a function $\sigma : pRuns_{\mathcal{M}} \cdot S_{\square} \rightarrow \mathcal{D}(S)$ that assigns to partial runs ρs , where $s \in S_{\square}$, a distribution over the successors $\{s' \in S \mid s \rightarrow s'\}$. The set of all strategies in \mathcal{M} is denoted by $\Sigma_{\mathcal{M}}$ (we omit the subscript and write Σ if \mathcal{M} is clear from the context). A (partial) run $s_0 e_0 s_1 e_1 \dots$ is consistent with a strategy σ if for all i either $s_i \in S_{\square}$ and $\sigma(s_0 e_0 s_1 e_1 \dots s_i)(s_{i+1}) > 0$, or $s_i \in S_{\circ}$ and $P(s_i)(s_{i+1}) > 0$.

An MDP $\mathcal{M} = (S, S_{\square}, S_{\circ}, \rightarrow, P, r)$, an initial state $s_0 \in S$, and a strategy σ induce a probability space in which the outcomes are runs starting in s_0 and with measure $\mathcal{P}_{\mathcal{M}, s_0, \sigma}$ defined as follows. It is first defined on *cylinders* $s_0 e_0 s_1 e_1 \dots s_n Runs_{\mathcal{M}}^{s_n}$: if $s_0 e_0 s_1 e_1 \dots s_n$ is not a partial run consistent with σ then $\mathcal{P}_{\mathcal{M}, s_0, \sigma}(s_0 e_0 s_1 e_1 \dots s_n Runs_{\mathcal{M}}^{s_n}) \stackrel{\text{def}}{=} 0$. Otherwise, $\mathcal{P}_{\mathcal{M}, s_0, \sigma}(s_0 e_0 s_1 e_1 \dots s_n Runs_{\mathcal{M}}^{s_n}) \stackrel{\text{def}}{=} \prod_{i=0}^{n-1} \bar{\sigma}(s_0 e_0 s_1 \dots s_i)(s_{i+1})$, where $\bar{\sigma}$ is the map that extends σ by $\bar{\sigma}(ws) = P(s)$ for all partial runs $ws \in pRuns_{\mathcal{M}} \cdot S_{\circ}$. By Carathéodory's theorem [6], this extends uniquely to a probability measure $\mathcal{P}_{\mathcal{M}, s_0, \sigma}$ on the Borel σ -algebra \mathcal{F} of subsets of $Runs_{\mathcal{M}}^{s_0}$. Elements of \mathcal{F} , i.e., measurable sets of runs, are called *events* or *objectives* here. For $X \in \mathcal{F}$ we will write $\bar{X} \stackrel{\text{def}}{=} Runs_{\mathcal{M}}^{s_0} \setminus X \in \mathcal{F}$ for its complement and $\mathcal{E}_{\mathcal{M}, s_0, \sigma}$ for the expectation wrt. $\mathcal{P}_{\mathcal{M}, s_0, \sigma}$. We drop the indices if possible without ambiguity.

Objectives. We consider objectives that are determined by a predicate on infinite runs. We assume familiarity with the syntax and semantics of the temporal logic LTL [13]. Formulas are interpreted on the structure (S, \rightarrow) . We use $\llbracket \varphi \rrbracket^s$ to denote the set of runs starting from s that satisfy the LTL formula φ , which is a measurable set [34]. We also write $\llbracket \varphi \rrbracket$ for $\bigcup_{s \in S} \llbracket \varphi \rrbracket^s$. Where it does not cause confusion we will identify φ and $\llbracket \varphi \rrbracket$ and just write $\mathcal{P}_{\mathcal{M}, s, \sigma}(\varphi)$ instead of $\mathcal{P}_{\mathcal{M}, s, \sigma}(\llbracket \varphi \rrbracket^s)$. The reachability objective of eventually visiting a set of states X can be expressed by $\llbracket FX \rrbracket \stackrel{\text{def}}{=} \{\rho \mid \exists i. \rho_s(i) \in X\}$. Reaching X within at most k steps is expressed by $\llbracket F^{\leq k} X \rrbracket \stackrel{\text{def}}{=} \{\rho \mid \exists i \leq k. \rho_s(i) \in X\}$. The definitions for eventually visiting certain transitions are analogous. The operator G (always) is defined as $\neg F \neg$. So the safety objective of avoiding X is expressed by $G \neg X$.

We consider the following objectives.

- The $PP_{\liminf \geq 0}$ objective is to maximize the probability that the *lim inf* of the *point* payoffs (the immediate transition rewards) is ≥ 0 , i.e., $PP_{\liminf \geq 0} \stackrel{\text{def}}{=} \{\rho \mid \liminf_{n \in \mathbb{N}} r(\rho_e(n)) \geq 0\}$.
- The $MP_{\liminf \geq 0}$ objective is to maximize the probability that the *lim inf* of the *mean* payoff is ≥ 0 , i.e., $MP_{\liminf \geq 0} \stackrel{\text{def}}{=} \{\rho \mid \liminf_{n \in \mathbb{N}} \frac{1}{n} \sum_{j=0}^{n-1} r(\rho_e(j)) \geq 0\}$.
- The $TP_{\liminf \geq 0}$ objective is to maximize the probability that the *lim inf* of the *total* payoff (the sum of the transition rewards seen so far) is ≥ 0 , i.e., $TP_{\liminf \geq 0} \stackrel{\text{def}}{=} \{\rho \mid \liminf_{n \in \mathbb{N}} \sum_{j=0}^{n-1} r(\rho_e(j)) \geq 0\}$.

An objective φ is called *shift invariant* in \mathcal{M} if for every run $\rho' \rho$ in \mathcal{M} with some finite prefix ρ' we have $\rho' \rho \in \llbracket \varphi \rrbracket \Leftrightarrow \rho \in \llbracket \varphi \rrbracket$. An objective is called a *shift invariant objective* if it is shift invariant in every MDP. $PP_{\liminf \geq 0}$ and $MP_{\liminf \geq 0}$ are shift invariant objectives, but $TP_{\liminf \geq 0}$ is not. Also $PP_{\liminf \geq 0}$ is more general than co-Büchi. (The special case of integer transition rewards coincides with co-Büchi, since rewards ≤ -1 and accepting states can be encoded into each other.)

Strategy Classes. Strategies are in general *randomized* (R) in the sense that they take values in $\mathcal{D}(S)$. A strategy σ is *deterministic* (D) if $\sigma(\rho)$ is a Dirac distribution for all ρ . General strategies can be *history dependent* (H), while others are restricted by the size or type of memory they use, see below. We consider certain classes of strategies:

- A strategy σ is *memoryless* (M) (also called *positional*) if it can be implemented with a memory of size 1. We may view M-strategies as functions $\sigma : S_{\square} \rightarrow \mathcal{D}(S)$.
- A strategy σ is *finite memory* (F) if there exists a finite memory \mathbf{M} implementing σ . Hence FR stands for finite memory randomized.
- A *step counter (SC) strategy* bases decisions only on the current state and the number of steps taken so far, i.e., it uses an unbounded integer counter that gets incremented by 1 in every step. Such strategies are also called *Markov strategies* [27].
- *k-bit Markov strategies* use k extra bits of general purpose memory in addition to a step counter [21].
- A *reward counter (RC) strategy* uses infinite memory, but only in the form of a counter that always contains the sum of all transition rewards seen so far.
- A *step counter + reward counter strategy* uses both a step counter and a reward counter.

Step counters and reward counters are very restricted forms of memory, since the memory update is not directly under the control of the player. These counters merely record an aspect of the partial run.

Memory and strategies. Here we give a formal definition how strategies use memory. Let \mathbf{M} be a countable set of memory modes, and let $\tau : \mathbf{M} \times S \rightarrow \mathcal{D}(\mathbf{M} \times S)$ be a function that meets the following two conditions: for all modes $\mathbf{m} \in \mathbf{M}$,

- for all controlled states $s \in S_{\square}$, the distribution $\tau(\mathbf{m}, s)$ is over $\mathbf{M} \times \{s' \mid s \rightarrow s'\}$.
- for all random states $s \in S_{\circ}$, and $s' \in S$, we have $\sum_{\mathbf{m}' \in \mathbf{M}} \tau(\mathbf{m}, s)(\mathbf{m}', s') = P(s)(s')$.

The function τ together with an initial memory mode \mathbf{m}_0 induce a strategy σ_{τ} as follows. Consider the Markov chain with the set $\mathbf{M} \times S$ of states and the probability function τ . A sequence $\rho = s_0 \cdots s_i$ corresponds to a set $H(\rho) = \{(\mathbf{m}_0, s_0) \cdots (\mathbf{m}_i, s_i) \mid \mathbf{m}_0, \dots, \mathbf{m}_i \in \mathbf{M}\}$ of runs in this Markov chain. Each $\rho s \in s_0 S^* S_{\square}$ induces a probability distribution $\mu_{\rho s} \in \mathcal{D}(\mathbf{M})$, the probability of being in state (\mathbf{m}, s) conditioned on having taken some partial run from $H(\rho s)$. We define σ_{τ} such that $\sigma_{\tau}(\rho s)(s') = \sum_{\mathbf{m}, \mathbf{m}' \in \mathbf{M}} \mu_{\rho s}(\mathbf{m}) \tau(\mathbf{m}, s)(\mathbf{m}', s')$ for all $\rho s \in S^* S_{\square}$ and all $s' \in S$.

We say that a strategy σ can be *implemented* with memory \mathbf{M} if there exist $\mathbf{m}_0 \in \mathbf{M}$ and τ such that $\sigma_{\tau} = \sigma$.

Optimal and ε -optimal Strategies. Given an objective φ , the value of state s in an MDP \mathcal{M} , denoted by $\text{val}_{\mathcal{M}, \varphi}(s)$, is the supremum probability of achieving φ . Formally, $\text{val}_{\mathcal{M}, \varphi}(s) \stackrel{\text{def}}{=} \sup_{\sigma \in \Sigma} \mathcal{P}_{\mathcal{M}, s, \sigma}(\varphi)$ where Σ is the set of all strategies. For $\varepsilon \geq 0$ and state $s \in S$, we say that a strategy is ε -optimal from s if $\mathcal{P}_{\mathcal{M}, s, \sigma}(\varphi) \geq \text{val}_{\mathcal{M}, \varphi}(s) - \varepsilon$. A 0-optimal strategy is called *optimal*. An optimal strategy is *almost-surely winning* if $\text{val}_{\mathcal{M}, \varphi}(s) = 1$. Considering an MD strategy as a function $\sigma : S_{\square} \rightarrow S$ and $\varepsilon \geq 0$, σ is *uniformly ε -optimal* (resp. *uniformly optimal*) if it is ε -optimal (resp. optimal) from *every* $s \in S$.

MDP variants. In order to show our results, we will sometimes use derived MDPs. Given an MDP \mathcal{M} , we define three different MDPs $S(\mathcal{M})$, $R(\mathcal{M})$ and $A(\mathcal{M})$. These new MDPs will be used in order to reduce objectives to $PP_{\liminf \geq 0}$ in a simpler setting with the step and/or reward counter encoded into the states.

Definition 2.1. Let \mathcal{M} be an MDP with a given initial state s_0 . We construct the MDP $S(\mathcal{M}) \stackrel{\text{def}}{=} (S', S'_\square, S'_\circ, \longrightarrow_{S(\mathcal{M})}, P')$ that encodes the step counter into the states as follows:

- The state space of $S(\mathcal{M})$ is $S' \stackrel{\text{def}}{=} \{(s, n) \mid s \in S \text{ and } n \in \mathbb{N}\}$. Note that S' is countable. We write s'_0 for the initial state $(s_0, 0)$.
- $S'_\square \stackrel{\text{def}}{=} \{(s, n) \in S' \mid s \in S_\square\}$ and $S'_\circ \stackrel{\text{def}}{=} S' \setminus S'_\square$.
- The set of transitions in $S(\mathcal{M})$ is

$$\longrightarrow_{S(\mathcal{M})} \stackrel{\text{def}}{=} \{((s, n), (s', n+1)) \mid (s, n), (s', n+1) \in S', s \longrightarrow_{\mathcal{M}} s'\}.$$

- $P' : S'_\circ \rightarrow \mathcal{D}(S')$ is defined such that

$$P'(s, n)(s', m) \stackrel{\text{def}}{=} \begin{cases} P(s)(s') & \text{if } (s, n) \longrightarrow_{S(\mathcal{M})} (s', m) \\ 0 & \text{otherwise} \end{cases}$$

- The reward $r((s, n) \longrightarrow_{S(\mathcal{M})} (s', n+1)) \stackrel{\text{def}}{=} r(s \longrightarrow_{\mathcal{M}} s')$.

By labeling the state with the path length from s_0 , we effectively encode a step counter into the MDP $S(\mathcal{M})$.

Lemma 2.2. *Let \mathcal{M} be an MDP with initial state s_0 . Then given an MD strategy σ' in $S(\mathcal{M})$ attaining $c \in [0, 1]$ for $PP_{\liminf \geq 0}$ from $(s_0, 0)$, there exists a strategy σ attaining c for $PP_{\liminf \geq 0}$ in \mathcal{M} from s_0 which uses the same memory as σ' plus a step counter.*

Proof. Let σ' be an MD strategy in $S(\mathcal{M})$ attaining $c \in [0, 1]$ for $PP_{\liminf \geq 0}$ from $(s_0, 0)$. We define a strategy σ on \mathcal{M} from s_0 that uses the same memory as σ' plus a step counter. Then σ plays on \mathcal{M} exactly like σ' plays on $S(\mathcal{M})$, keeping the step counter in its memory instead of in the state. I.e., at a given state s and step counter value n , σ plays exactly as σ' plays in state (s, n) . By our construction of $S(\mathcal{M})$ and the definition of σ , the sequences of point rewards seen by σ' in runs on $S(\mathcal{M})$ coincide with the sequences of point rewards seen by σ in runs in \mathcal{M} . Hence we obtain $\mathcal{P}_{S(\mathcal{M}), (s_0, 0), \sigma'}(PP_{\liminf \geq 0}) = \mathcal{P}_{\mathcal{M}, s_0, \sigma}(PP_{\liminf \geq 0})$ \square

Definition 2.3. Let \mathcal{M} be an MDP. From a given initial state s_0 , the reward level in each state $s \in S$ can be any of the countably many values r_1, r_2, \dots corresponding to the rewards accumulated along all the possible paths leading to s from s_0 . We construct the MDP $R(\mathcal{M}) \stackrel{\text{def}}{=} (S', S'_\square, S'_\circ, \longrightarrow_{R(\mathcal{M})}, P')$ that encodes the reward counter into the state as follows:

- The state space of $R(\mathcal{M})$ is $S' \stackrel{\text{def}}{=} \{(s, r) \mid s \in S, r \in \mathbb{R} \text{ is a reward level attainable at } s\}$. Note that S' is countable. We write s'_0 for the initial state $(s_0, 0)$.
- $S'_\square \stackrel{\text{def}}{=} \{(s, r) \in S' \mid s \in S_\square\}$ and $S'_\circ \stackrel{\text{def}}{=} S' \setminus S'_\square$.
- The set of transitions in $R(\mathcal{M})$ is

$$\begin{aligned} \longrightarrow_{R(\mathcal{M})} \stackrel{\text{def}}{=} \{ & ((s, r), (s', r')) \mid (s, r), (s', r') \in S', \\ & s \longrightarrow s' \text{ in } \mathcal{M} \text{ and } r' \stackrel{\text{def}}{=} r + r(s \rightarrow s')\}. \end{aligned}$$

- $P' : S'_\circ \rightarrow \mathcal{D}(S')$ is defined such that

$$P'(s, r)(s', r') \stackrel{\text{def}}{=} \begin{cases} P(s)(s') & \text{if } (s, r) \rightarrow_{R(\mathcal{M})} (s', r') \\ 0 & \text{otherwise} \end{cases}$$

- The reward for taking transition $(s, r) \rightarrow (s', r')$ is r' .

By labeling transitions in $R(\mathcal{M})$ with the state encoded total reward of the target state, we ensure that the point rewards in $R(\mathcal{M})$ correspond exactly to the total rewards in \mathcal{M} .

Lemma 2.4. *Let \mathcal{M} be an MDP with initial state s_0 . Then given an MD (resp. Markov) strategy σ' in $R(\mathcal{M})$ attaining $c \in [0, 1]$ for $PP_{\liminf \geq 0}$ from $(s_0, 0)$, there exists a strategy σ attaining c for $TP_{\liminf \geq 0}$ in \mathcal{M} from s_0 which uses the same memory as σ' plus a reward counter.*

Proof. Let σ' be an MD (resp. Markov) strategy in $R(\mathcal{M})$ attaining $c \in [0, 1]$ for $PP_{\liminf \geq 0}$ from $(s_0, 0)$. We define a strategy σ on \mathcal{M} from s_0 that uses the same memory as σ' plus a reward counter. Then σ plays on \mathcal{M} exactly like σ' plays on $R(\mathcal{M})$, keeping the reward counter in its memory instead of in the state. I.e., at a given state s (and step counter value m , in case σ' was a Markov strategy) and reward level r , σ plays exactly as σ' plays in state (s, r) (and step counter value m , in case σ' was a Markov strategy). By our construction of $R(\mathcal{M})$ and the definition of σ , the sequences of point rewards seen by σ' in runs on $R(\mathcal{M})$ coincide with the sequences of total rewards seen by σ in runs in \mathcal{M} . Hence we obtain $\mathcal{P}_{R(\mathcal{M}), (s_0, 0), \sigma'}(PP_{\liminf \geq 0}) = \mathcal{P}_{\mathcal{M}, s_0, \sigma}(TP_{\liminf \geq 0})$ as required. \square

Definition 2.5. Given an MDP \mathcal{M} with initial state s_0 , we define the new MDP $A(\mathcal{M})$ that encodes the mean payoffs of partial runs of \mathcal{M} into the seen transition rewards in $A(\mathcal{M})$. To this end, the states of $A(\mathcal{M})$ encode both the step counter and the total reward of the run so far. However, the transition rewards in $A(\mathcal{M})$ reflect the *mean* payoff, i.e., the total reward divided by the number of steps.

We construct $A(\mathcal{M}) \stackrel{\text{def}}{=} (S', S'_\square, S'_\circ, \rightarrow_{A(\mathcal{M})}, P')$ as follows:

- The state space of $A(\mathcal{M})$ is

$$S' \stackrel{\text{def}}{=} \{(s, n, r) \mid s \in S, n \in \mathbb{N} \text{ and } r \in \mathbb{R} \text{ is a reward level attainable at } s \text{ at step } n\}$$

Note that S' is countable. We write s'_0 for the initial state $(s_0, 0, 0)$ of $A(\mathcal{M})$.

- $S'_\square \stackrel{\text{def}}{=} \{(s, n, r) \in S' \mid s \in S_\square\}$ and $S'_\circ \stackrel{\text{def}}{=} S' \setminus S'_\square$.
- The set of transitions in $A(\mathcal{M})$ is

$$\begin{aligned} \rightarrow_{A(\mathcal{M})} \stackrel{\text{def}}{=} \{ & ((s, n, r), (s', n+1, r')) \mid \\ & (s, n, r), (s', n+1, r') \in S', \\ & s \rightarrow s' \text{ in } \mathcal{M} \text{ and } r' = r + r(s \rightarrow s') \}. \end{aligned}$$

- $P' : S'_\circ \rightarrow \mathcal{D}(S')$ is defined such that

$$P'(s, n, r)(s', n', r') \stackrel{\text{def}}{=} \begin{cases} P(s)(s') & \text{if } (s, n, r) \rightarrow_{A(\mathcal{M})} (s', n', r') \\ 0 & \text{otherwise} \end{cases}$$

- The reward for taking transition $(s, n, r) \rightarrow (s', n', r')$ is r'/n' , i.e., the transition reward is the *mean* payoff of the partial run so far.

Lemma 2.6. *Let \mathcal{M} be an MDP with initial state s_0 . Then given an MD strategy σ' in $A(\mathcal{M})$ attaining $c \in [0, 1]$ for $PP_{\liminf \geq 0}$ from $(s_0, 0, 0)$, there exists a strategy σ attaining c for $MP_{\liminf \geq 0}$ in \mathcal{M} from s_0 which uses just a reward counter and a step counter.*

Proof. The proof is very similar to that of Lemma 2.4. \square

Sums and products. In our proofs we will use the following basic properties of sums and products (see, e.g., [10]).

Proposition 2.7. *Given an infinite sequence of real numbers a_n with $0 \leq a_n < 1$, we have*

$$\prod_{n=1}^{\infty} (1 - a_n) > 0 \quad \Leftrightarrow \quad \sum_{n=1}^{\infty} a_n < \infty.$$

and the “ \Rightarrow ” implication holds even for the weaker assumption $0 \leq a_n \leq 1$.

Proof. If $a_n = 1$ for any n then the “ \Rightarrow ” implication is vacuously true, but the “ \Leftarrow ” implication does not hold in general. In the following we assume $0 \leq a_n < 1$.

In the case where a_n does not converge to zero, the property is trivial. In the case where $a_n \rightarrow 0$, it is shown by taking the logarithm of the product and using the limit comparison test as follows.

Taking the logarithm of the product gives the series

$$\sum_{n=1}^{\infty} \ln(1 - a_n)$$

whose convergence (to a finite number ≤ 0) is equivalent to the positivity of the product. It is also equivalent to the convergence (to a number ≥ 0) of its negation $\sum_{n=1}^{\infty} -\ln(1 - a_n)$. But observe that (by L'Hôpital's rule)

$$\lim_{x \rightarrow 0} \frac{-\ln(1 - x)}{x} = 1.$$

Since $a_n \rightarrow 0$ we have

$$\lim_{n \rightarrow \infty} \frac{-\ln(1 - a_n)}{a_n} = 1.$$

By the limit comparison test, the series $\sum_{n=1}^{\infty} -\ln(1 - a_n)$ converges if and only if the series $\sum_{n=1}^{\infty} a_n$ converges. \square

Proposition 2.8. *Given an infinite sequence of real numbers a_n with $0 \leq a_n \leq 1$,*

$$\prod_{n=1}^{\infty} a_n > 0 \quad \Rightarrow \quad \forall \varepsilon > 0 \exists N. \prod_{n=N}^{\infty} a_n \geq (1 - \varepsilon).$$

Proof. If $a_n = 0$ for any n then the property is vacuously true. In the following we assume $a_n > 0$. Since $\prod_{n=1}^{\infty} a_n > 0$, by taking the logarithm we obtain $\sum_{n=1}^{\infty} \ln(a_n) > -\infty$. Thus for every $\delta > 0$ there exists an N s.t. $\sum_{n=N}^{\infty} \ln(a_n) \geq -\delta$. By exponentiation we obtain $\prod_{n=N}^{\infty} a_n \geq \exp(-\delta)$. By picking $\delta = -\ln(1 - \varepsilon)$ the result follows. \square

Point Payoff		ε -optimal	Optimal
Finitely branching	Upper Bound	Det(Positional) 3.5	Det(Positional) 3.8
	Lower Bound	n/a	n/a
Infinitely branching	Upper Bound	Det(SC) 3.10	Det(SC) 3.12
	Lower Bond	\neg Rand(F) 3.13	\neg Rand(F) 3.13

Table 2: Strategy complexity of ε -optimal/optimal strategies for the point payoff objective in infinitely/finitely branching MDPs.

3. POINT PAYOFF

3.1. Upper Bounds. In this section we show that for *finitely branching* MDPs, there exist ε -optimal MD strategies for $PP_{\liminf \geq 0}$. Whereas for *infinitely branching* MDPs, a step counter suffices in order to achieve $PP_{\liminf \geq 0}$ ε -optimally.

These two very technical results will form the basis of our upper bound analysis of $MP_{\liminf \geq 0}$ and $TP_{\liminf \geq 0}$ in later sections.

Lemma 3.1. ([21, Lemma 23]) *For every acyclic MDP with a safety objective and every $\varepsilon > 0$, there exists an MD strategy that is uniformly ε -optimal.*

Theorem 3.2. ([22, Theorem 7]) *Let $\mathcal{M} = (S, S_\square, S_\circ, \longrightarrow, P, r)$ be a countable MDP, and let φ be an event that is shift invariant in \mathcal{M} . Suppose for every $s \in S$ there exist ε -optimal MD strategies for φ . Then:*

- (1) *There exist uniform ε -optimal MD strategies for φ .*
- (2) *There exists a single MD strategy that is optimal from every state that has an optimal strategy.*

3.1.1. Finitely Branching Case. In order to prove the main result of this section, we use the following result on the **Transience** objective, which is the set of runs that do not visit any state infinitely often. Given an MDP $\mathcal{M} = (S, S_\square, S_\circ, \longrightarrow, P, r)$, $\text{Transience} \stackrel{\text{def}}{=} \bigwedge_{s \in S} \text{FG} \neg s$.

Theorem 3.3. ([22, Theorem 8]) *In every countable MDP there exist uniform ε -optimal MD strategies for Transience.*

In this section we consider finitely branching MDPs. We need the following technical lemma that holds only for finitely branching MDPs.

Lemma 3.4. *Given a finitely branching countable MDP \mathcal{M} , a subset $T \subseteq \longrightarrow$ of the transitions and a state s , we have*

$$\text{val}_{\mathcal{M}, \neg FT}(s) < 1 \Rightarrow \exists k \in \mathbb{N}. \text{val}_{\mathcal{M}, \neg F^{\leq k} T}(s) < 1$$

i.e., if it is impossible to completely avoid T then there is a bounded threshold k and a fixed nonzero chance of seeing T within $\leq k$ steps, regardless of the strategy.

Proof. It suffices to show that $\forall k \in \mathbb{N}. \text{val}_{\mathcal{M}, \neg F^{\leq k} T}(s) = 1$ implies $\text{val}_{\mathcal{M}, \neg FT}(s) = 1$. Since \mathcal{M} is finitely branching, the state s has only finitely many successors $\{s_1, \dots, s_n\}$.

Consider the case where s is a controlled state. If we had the property $\forall 1 \leq i \leq n \exists k_i \in \mathbb{N}. \text{val}_{\mathcal{M}, \neg F^{\leq k_i} T}(s_i) < 1$ then we would have $\text{val}_{\mathcal{M}, \neg F^{\leq k} T}(s) < 1$ for $k = (\max_{1 \leq i \leq n} k_i) + 1$ which contradicts our assumption. Thus there must exist an $i \in \{1, \dots, n\}$ with $\forall k \in \mathbb{N}$

$\mathbb{N} \cdot \mathbf{val}_{\mathcal{M}, \neg F \leq kT}(s_i) = 1$. We define a strategy σ that chooses the successor state s_i when in state s .

Similarly, if s is a random state, we must have $\forall k \in \mathbb{N} \cdot \mathbf{val}_{\mathcal{M}, \neg F \leq kT}(s_i) = 1$ for all its successors s_i .

By using our constructed strategy σ , we obtain $\mathcal{P}_{\mathcal{M}, s, \sigma}(\neg FT) = 1$ and thus $\mathbf{val}_{\mathcal{M}, \neg FT}(s) = 1$ as required. \square

Theorem 3.5. *Consider a finitely branching MDP $\mathcal{M} = (S, S_{\square}, S_{\circ}, \longrightarrow, P, r)$ with initial state s_0 and a $PP_{\liminf \geq 0}$ objective. Then there exist ε -optimal MD strategies.*

Proof. Let $\varepsilon > 0$. We begin by partitioning the state space into two sets, S_{safe} and $S \setminus S_{\text{safe}}$. The set S_{safe} is the subset of states which is surely winning for the safety objective of only using transitions with non-negative rewards (i.e., never using transitions with negative rewards at all). Since \mathcal{M} is finitely branching, there exists a uniformly optimal MD strategy σ_{safe} for this safety objective [27, 23].

We construct a new MDP \mathcal{M}' by modifying \mathcal{M} . We create a gadget G_{safe} composed of a sequence of new controlled states x_0, x_1, x_2, \dots where all transitions $x_i \rightarrow x_{i+1}$ have reward 0. Hence any run entering G_{safe} is winning for $PP_{\liminf \geq 0}$. We insert G_{safe} into \mathcal{M} by replacing all incoming transitions to S_{safe} with transitions that lead to x_0 . The idea behind this construction is that when playing in \mathcal{M} , once you hit a state in S_{safe} , you can win surely by playing an optimal MD strategy for safety. So we replace S_{safe} with the surely winning gadget G_{safe} . Thus

$$\mathbf{val}_{\mathcal{M}, PP_{\liminf \geq 0}}(s_0) = \mathbf{val}_{\mathcal{M}', PP_{\liminf \geq 0}}(s_0) \quad (3.1)$$

and if an ε -optimal MD strategy exists in \mathcal{M} , then there exists a corresponding one in \mathcal{M}' , and vice-versa.

We now consider a general (not necessarily MD) ε -optimal strategy σ for $PP_{\liminf \geq 0}$ from s_0 on \mathcal{M}' , i.e.,

$$\mathcal{P}_{\mathcal{M}', s_0, \sigma}(PP_{\liminf \geq 0}) \geq \mathbf{val}_{\mathcal{M}', PP_{\liminf \geq 0}}(s_0) - \varepsilon. \quad (3.2)$$

Define the safety objective Safety_i which is the objective of never seeing any point rewards $< -2^{-i}$. This then allows us to characterize $PP_{\liminf \geq 0}$ in terms of safety objectives.

$$PP_{\liminf \geq 0} = \bigcap_{i \in \mathbb{N}} F(\text{Safety}_i). \quad (3.3)$$

Now we define the safety objective $\text{Safety}_i^k \stackrel{\text{def}}{=} F^{\leq k}(\text{Safety}_i)$ to attain Safety_i within at most k steps. This allows us to write

$$F(\text{Safety}_i) = \bigcup_{k \in \mathbb{N}} \text{Safety}_i^k. \quad (3.4)$$

By continuity of measures from above we get

$$0 = \mathcal{P}_{\mathcal{M}', s_0, \sigma} \left(F(\text{Safety}_i) \cap \bigcap_{k \in \mathbb{N}} \overline{\text{Safety}_i^k} \right) = \lim_{k \rightarrow \infty} \mathcal{P}_{\mathcal{M}', s_0, \sigma} \left(F(\text{Safety}_i) \cap \overline{\text{Safety}_i^k} \right).$$

Hence for every $i \in \mathbb{N}$ and $\varepsilon_i \stackrel{\text{def}}{=} \varepsilon \cdot 2^{-i}$ there exists n_i such that

$$\mathcal{P}_{\mathcal{M}', s_0, \sigma} \left(F(\text{Safety}_i) \cap \overline{\text{Safety}_i^{n_i}} \right) \leq \varepsilon_i. \quad (3.5)$$

Now we can show the following claim.

Claim 3.6.

$$\mathcal{P}_{\mathcal{M}', s_0, \sigma} \left(\bigcap_{i \in \mathbb{N}} \text{Safety}_i^{n_i} \right) \geq \text{val}_{\mathcal{M}', PP_{\liminf \geq 0}}(s_0) - 2\varepsilon.$$

Proof.

$$\begin{aligned}
& \mathcal{P}_{\mathcal{M}', s_0, \sigma} \left(\bigcap_{i \in \mathbb{N}} \text{Safety}_i^{n_i} \right) \\
& \geq \mathcal{P}_{\mathcal{M}', s_0, \sigma} \left(\bigcap_{k \in \mathbb{N}} \text{F}(\text{Safety}_k) \cap \bigcap_{i \in \mathbb{N}} \text{Safety}_i^{n_i} \right) \\
& = \mathcal{P}_{\mathcal{M}', s_0, \sigma} \left(\left(\bigcap_{k \in \mathbb{N}} \text{F}(\text{Safety}_k) \cap \bigcap_{i \in \mathbb{N}} \text{Safety}_i^{n_i} \right) \cup \left(\overline{\bigcap_{k \in \mathbb{N}} \text{F}(\text{Safety}_k)} \cap \bigcap_{k \in \mathbb{N}} \text{F}(\text{Safety}_k) \right) \right) \\
& = \mathcal{P}_{\mathcal{M}', s_0, \sigma} \left(\bigcap_{k \in \mathbb{N}} \text{F}(\text{Safety}_k) \cap \left(\bigcap_{i \in \mathbb{N}} \text{Safety}_i^{n_i} \cup \overline{\bigcap_{k \in \mathbb{N}} \text{F}(\text{Safety}_k)} \right) \right) \\
& = 1 - \mathcal{P}_{\mathcal{M}', s_0, \sigma} \left(\overline{\bigcap_{k \in \mathbb{N}} \text{F}(\text{Safety}_k)} \cup \left(\overline{\bigcap_{i \in \mathbb{N}} \text{Safety}_i^{n_i}} \cap \bigcap_{k \in \mathbb{N}} \text{F}(\text{Safety}_k) \right) \right) \\
& \geq 1 - \mathcal{P}_{\mathcal{M}', s_0, \sigma} \left(\overline{\bigcap_{k \in \mathbb{N}} \text{F}(\text{Safety}_k)} \right) - \mathcal{P}_{\mathcal{M}', s_0, \sigma} \left(\overline{\bigcap_{i \in \mathbb{N}} \text{Safety}_i^{n_i}} \cap \bigcap_{k \in \mathbb{N}} \text{F}(\text{Safety}_k) \right) \\
& = \mathcal{P}_{\mathcal{M}', s_0, \sigma} (PP_{\liminf \geq 0}) - \mathcal{P}_{\mathcal{M}', s_0, \sigma} \left(\overline{\bigcap_{i \in \mathbb{N}} \text{Safety}_i^{n_i}} \cap \bigcap_{k \in \mathbb{N}} \text{F}(\text{Safety}_k) \right) \quad \text{by (3.3)} \\
& \geq \text{val}_{\mathcal{M}', PP_{\liminf \geq 0}}(s_0) - \varepsilon - \mathcal{P}_{\mathcal{M}', s_0, \sigma} \left(\bigcup_{i \in \mathbb{N}} \overline{\text{Safety}_i^{n_i}} \cap \bigcap_{k \in \mathbb{N}} \text{F}(\text{Safety}_k) \right) \quad \text{by (3.2)} \\
& \geq \text{val}_{\mathcal{M}', PP_{\liminf \geq 0}}(s_0) - \varepsilon - \sum_{i \in \mathbb{N}} \mathcal{P}_{\mathcal{M}', s_0, \sigma} \left(\overline{\text{Safety}_i^{n_i}} \cap \bigcap_{k \in \mathbb{N}} \text{F}(\text{Safety}_k) \right) \\
& \geq \text{val}_{\mathcal{M}', PP_{\liminf \geq 0}}(s_0) - \varepsilon - \sum_{i \in \mathbb{N}} \varepsilon_i \quad \text{by (3.5)} \\
& = \text{val}_{\mathcal{M}', PP_{\liminf \geq 0}}(s_0) - 2\varepsilon
\end{aligned}$$

□

Since \mathcal{M}' does not have an implicit step counter, we use the following construction to approximate one. We define the distance $d(s)$ from s_0 to a state s as the length of the shortest path from s_0 to s . Let $\text{Bubble}_n(s_0) \stackrel{\text{def}}{=} \{s \in S \mid d(s) \leq n\}$ be those states that can be reached within n steps from s_0 . Since \mathcal{M}' is finitely branching, $\text{Bubble}_n(s_0)$ is finite for every fixed n . Let

$$\text{Bad}_i \stackrel{\text{def}}{=} \{t \in \rightarrow_{\mathcal{M}'} \mid t = s \rightarrow_{\mathcal{M}'} s', s \notin \text{Bubble}_{n_i}(s_0) \text{ and } r(t) < -2^{-i}\}$$

be the set of transitions originating outside $\text{Bubble}_{n_i}(s_0)$ whose reward is too negative. Thus a run from s_0 that satisfies $\text{Safety}_i^{n_i}$ cannot use any transition in Bad_i , since (by definition of $\text{Bubble}_{n_i}(s_0)$) they would come after the n_i -th step.

Now we create a new state \perp whose only outgoing transition is a self loop with reward -1 . We transform \mathcal{M}' into \mathcal{M}'' by re-directing all transitions in Bad_i to the new target state \perp for every i . Notice that any run visiting \perp must be losing for $PP_{\liminf \geq 0}$ due to the negative reward on the self loop, but it must also be losing for **Transience** because of the self loop.

We now show that the change from \mathcal{M}' to \mathcal{M}'' has decreased the value of s_0 for $PP_{\liminf \geq 0}$ by at most 2ε , i.e.,

$$\text{val}_{\mathcal{M}'', PP_{\liminf \geq 0}}(s_0) \geq \text{val}_{\mathcal{M}', PP_{\liminf \geq 0}}(s_0) - 2\varepsilon. \quad (3.6)$$

Equation (3.6) follows from the following steps.

$$\begin{aligned} \text{val}_{\mathcal{M}'', PP_{\liminf \geq 0}}(s_0) &\geq \mathcal{P}_{\mathcal{M}'', s_0, \sigma} \left(\bigcap_{i \in \mathbb{N}} \text{Safety}_i^{n_i} \right) \\ &= \mathcal{P}_{\mathcal{M}', s_0, \sigma} \left(\bigcap_{i \in \mathbb{N}} \text{Safety}_i^{n_i} \right) && \text{by def. of } \mathcal{M}'' \\ &\geq \text{val}_{\mathcal{M}', PP_{\liminf \geq 0}}(s_0) - 2\varepsilon && \text{by Claim 3.6} \end{aligned}$$

In the next step we argue that under *every* strategy σ'' from s_0 in \mathcal{M}'' the attainment for $PP_{\liminf \geq 0}$ and **Transience** coincide, i.e.,

Claim 3.7.

$$\forall \sigma''. \mathcal{P}_{\mathcal{M}'', s_0, \sigma''}(PP_{\liminf \geq 0}) = \mathcal{P}_{\mathcal{M}'', s_0, \sigma''}(\text{Transience}).$$

Proof. First we show that

$$\text{Transience} \subseteq PP_{\liminf \geq 0} \text{ in } \mathcal{M}''. \quad (3.7)$$

Let $\rho \in \text{Transience}$ be a transient run. Then ρ can never visit the state \perp . Moreover, ρ must eventually leave every finite set forever. In particular ρ must satisfy $\text{FG}(\neg \text{Bubble}_{n_i}(s_0))$ for every i , since $\text{Bubble}_{n_i}(s_0)$ is finite, because \mathcal{M}'' is finitely branching. Thus ρ must either fall into G_{safe} , in which case it satisfies $PP_{\liminf \geq 0}$, or for every i , ρ must eventually leave $\text{Bubble}_{n_i}(s_0)$ forever. By definition of $\text{Bubble}_{n_i}(s_0)$ and \mathcal{M}'' , the run ρ must eventually stop seeing rewards $< -2^{-i}$ for every i . In this case ρ also satisfies $PP_{\liminf \geq 0}$. Thus (3.7).

Secondly, we show that

$$\forall \sigma''. \mathcal{P}_{\mathcal{M}'', s_0, \sigma''}(PP_{\liminf \geq 0} \cap \overline{\text{Transience}}) = 0. \quad (3.8)$$

i.e., except for a null-set, $PP_{\liminf \geq 0}$ implies **Transience** in \mathcal{M}'' .

Let σ'' be an arbitrary strategy from s_0 in \mathcal{M}'' and \mathfrak{R} be the set of all runs induced by it. For every $s \in S$, let $\mathfrak{R}_s \stackrel{\text{def}}{=} \{\rho \in \mathfrak{R} \mid \rho \text{ satisfies } \text{GF}(s)\}$ be the set of runs seeing state s infinitely often. In particular, any run $\rho \in \mathfrak{R}_s$ is not transient. Indeed, $\overline{\text{Transience}} = \bigcup_{s \in S} \mathfrak{R}_s$. We want to show that for every state $s \in S$ and strategy σ''

$$\mathcal{P}_{\mathcal{M}'', s_0, \sigma''}(PP_{\liminf \geq 0} \cap \mathfrak{R}_s) = 0. \quad (3.9)$$

Since all runs seeing a state in G_{safe} are transient, every \mathfrak{R}_s with $s \in G_{\text{safe}}$ must be empty. Similarly, every run seeing \perp is losing for $PP_{\liminf \geq 0}$ by construction. Hence we have (3.9) for any state s where $s = \perp$ or $s \in G_{\text{safe}}$.

Now consider \mathfrak{R}_s where s is neither in G_{safe} nor \perp . Let $T_{\text{neg}} \stackrel{\text{def}}{=} \{t \in \longrightarrow \mid r(t) < 0\}$ be the subset of transitions with negative rewards in \mathcal{M}'' .

We now show that $\text{val}_{\mathcal{M}'', \neg FT_{neg}}(s) < 1$ by assuming the opposite and deriving a contradiction. Assume that $\text{val}_{\mathcal{M}'', \neg FT_{neg}}(s) = 1$. The objective $\neg FT_{neg}$ is a safety objective. Thus, since \mathcal{M}'' is finitely branching, there exists a strategy from s that surely avoids T_{neg} (always pick an optimal move) [27, 23]. (This does not hold in infinitely branching MDPs where optimal moves might not exist.) However, by construction of \mathcal{M}'' , this implies that $s \in G_{\text{safe}}$. Contradiction. Thus $\text{val}_{\mathcal{M}'', \neg FT_{neg}}(s) < 1$.

Since \mathcal{M}'' is finitely branching, we can apply Lemma 3.4 and obtain that there exists a threshold k_s such that $\text{val}_{\mathcal{M}'', \neg F \leq k_s T_{neg}}(s) < 1$. Therefore $\delta_s \stackrel{\text{def}}{=} 1 - \text{val}_{\mathcal{M}'', \neg F \leq k_s T_{neg}}(s) > 0$. Thus, under every strategy, upon visiting s there is a chance $\geq \delta_s$ of seeing a transition in T_{neg} within the next $\leq k_s$ steps. Moreover, the subset $T_{neg}^s \subseteq T_{neg}$ of transitions that can be reached in $\leq k_s$ steps from s is finite, since \mathcal{M}'' is finitely branching. The finiteness of T_{neg}^s implies that the maximum of the rewards in T_{neg}^s exists and is still negative, i.e., $\ell_s \stackrel{\text{def}}{=} \max\{r(t) \mid t \in T_{neg}^s\} < 0$. (This would not be true for an infinite set, since the sup over an infinite set of negative numbers could be zero.) Let $T_{\leq \ell} \stackrel{\text{def}}{=} \{t \in T \mid r(t) \leq \ell_s\}$ be the subset of transitions with rewards $\leq \ell_s$ in \mathcal{M}'' .

Thus, under *every* strategy, upon visiting s there is a chance $\geq \delta_s$ of seeing a transition in $T_{\leq \ell}$ within the next $\leq k_s$ steps.

For every state $s \in S$, let $\mathfrak{R}_s^i \stackrel{\text{def}}{=} \{\rho \in \mathfrak{R} \mid \rho \text{ visits } s \text{ at least } i \text{ times}\}$, so we get $\mathfrak{R}_s = \bigcap_{i \in \mathbb{N}} \mathfrak{R}_s^i$. We obtain

$$\begin{aligned}
& \sup_{\sigma''} \mathcal{P}_{\mathcal{M}'', s_0, \sigma''}(PP_{\liminf \geq 0} \cap \mathfrak{R}_s) \\
& \leq \sup_{\sigma''} \mathcal{P}_{\mathcal{M}'', s_0, \sigma''}(FG \neg T_{\leq \ell} \cap \mathfrak{R}_s) && \text{set inclusion} \\
& = \sup_{\sigma''} \lim_{n \rightarrow \infty} \mathcal{P}_{\mathcal{M}'', s_0, \sigma''}(F^{\leq n} G \neg T_{\leq \ell} \cap \mathfrak{R}_s) && \text{continuity of measures} \\
& \leq \sup_{\sigma'''} \mathcal{P}_{\mathcal{M}'', s, \sigma'''}(G \neg T_{\leq \ell} \cap \mathfrak{R}_s) && s \text{ visited after } > n \text{ steps} \\
& = \sup_{\sigma'''} \mathcal{P}_{\mathcal{M}'', s, \sigma'''}(G \neg T_{\leq \ell} \cap \bigcap_{i \in \mathbb{N}} \mathfrak{R}_s^i) && \text{def. of } \mathfrak{R}_s^i \\
& = \sup_{\sigma'''} \lim_{i \rightarrow \infty} \mathcal{P}_{\mathcal{M}'', s, \sigma'''}(G \neg T_{\leq \ell} \cap \mathfrak{R}_s^i) && \text{continuity of measures} \\
& \leq \lim_{i \rightarrow \infty} (1 - \delta_s)^i = 0 && \text{by def. of } \mathfrak{R}_s^i \text{ and } \delta_s
\end{aligned}$$

and thus (3.9).

From this we obtain $\mathcal{P}_{\mathcal{M}'', s_0, \sigma''}(PP_{\liminf \geq 0} \cap \overline{\text{Transience}}) = \mathcal{P}_{\mathcal{M}'', s_0, \sigma''}(PP_{\liminf \geq 0} \cap \bigcup_{s \in S} \mathfrak{R}_s) \leq \sum_{s \in S} \mathcal{P}_{\mathcal{M}'', s_0, \sigma''}(PP_{\liminf \geq 0} \cap \mathfrak{R}_s) = 0$ and thus (3.8).

From (3.7) and (3.8) we obtain that for every σ'' we have

$$\begin{aligned}
& \mathcal{P}_{\mathcal{M}'', s_0, \sigma''}(PP_{\liminf \geq 0}) \\
& = \mathcal{P}_{\mathcal{M}'', s_0, \sigma''}(PP_{\liminf \geq 0} \cap \text{Transience}) + \mathcal{P}_{\mathcal{M}'', s_0, \sigma''}(PP_{\liminf \geq 0} \cap \overline{\text{Transience}}) \\
& = \mathcal{P}_{\mathcal{M}'', s_0, \sigma''}(\text{Transience}) + 0 \\
& = \mathcal{P}_{\mathcal{M}'', s_0, \sigma''}(\text{Transience})
\end{aligned}$$

and thus Claim 3.7. □

By Theorem 3.3, there exists a uniformly ε -optimal MD strategy $\hat{\sigma}$ from s_0 for **Transience** in \mathcal{M}'' , i.e.,

$$\mathcal{P}_{\mathcal{M}'', s_0, \hat{\sigma}}(\text{Transience}) \geq \text{val}_{\mathcal{M}'', \text{Transience}}(s_0) - \varepsilon. \quad (3.10)$$

We construct an MD strategy σ^* in \mathcal{M} which plays like σ_{safe} in S_{safe} and plays like $\hat{\sigma}$ everywhere else.

$$\begin{aligned} \mathcal{P}_{\mathcal{M}, s_0, \sigma^*}(PP_{\liminf \geq 0}) &= \mathcal{P}_{\mathcal{M}', s_0, \hat{\sigma}}(PP_{\liminf \geq 0}) && \text{def. of } \sigma^* \text{ and } \sigma_{\text{safe}} \\ &\geq \mathcal{P}_{\mathcal{M}'', s_0, \hat{\sigma}}(PP_{\liminf \geq 0}) && \text{new losing sink in } \mathcal{M}'' \\ &= \mathcal{P}_{\mathcal{M}'', s_0, \hat{\sigma}}(\text{Transience}) && \text{by Claim 3.7} \\ &\geq \text{val}_{\mathcal{M}'', \text{Transience}}(s_0) - \varepsilon && \text{by (3.10)} \\ &= \text{val}_{\mathcal{M}'', PP_{\liminf \geq 0}}(s_0) - \varepsilon && \text{by Claim 3.7} \\ &\geq \text{val}_{\mathcal{M}', PP_{\liminf \geq 0}}(s_0) - 2\varepsilon - \varepsilon && \text{by (3.6)} \\ &= \text{val}_{\mathcal{M}, PP_{\liminf \geq 0}}(s_0) - 3\varepsilon && \text{by (3.1)} \end{aligned}$$

Hence σ^* is a 3ε -optimal MD strategy for $PP_{\liminf \geq 0}$ from s_0 in \mathcal{M} . Since ε can be chosen arbitrarily small, the result follows. \square

Corollary 3.8. *Given a finitely branching MDP \mathcal{M} and initial state s_0 , optimal strategies, where they exist, can be chosen MD for $PP_{\liminf \geq 0}$.*

Proof. Since $PP_{\liminf \geq 0}$ is shift invariant, the result follows from Theorem 3.5 and Theorem 3.2. \square

Remark 3.9. The proof of Theorem 3.5 (and thus also Corollary 3.8) uses the result about the **Transience** objective from Theorem 3.3. This is unavoidable, since, at least for finitely branching MDPs, Theorem 3.5 also conversely implies Theorem 3.3 as follows.

Consider a finitely branching MDP $\mathcal{M} = (S, S_{\square}, S_{\circ}, \rightarrow, P, r)$ with initial state s_0 . Then we can define a reward structure on \mathcal{M} such that **Transience** and $PP_{\liminf \geq 0}$ coincide. For each transition t from state s to state s' let $n(t) \stackrel{\text{def}}{=} \min\{n \mid s \in \text{Bubble}_n(s_0)\}$ and $r(t) \stackrel{\text{def}}{=} -1/n(t)$. Since \mathcal{M} is finitely branching, $\text{Bubble}_n(s_0)$ is finite for every n . Transient runs eventually leave every finite set forever. Thus for all runs starting in s_0 we have **Transience** $\subseteq PP_{\liminf \geq 0}$, because $\lim_{n \rightarrow \infty} -1/n = 0$. For the reverse inclusion, consider a non-transient run from s_0 . This run visits some state s infinitely often. Since \mathcal{M} is finitely branching, by the Pigeonhole principle, it also visits some transition t from s infinitely often. So it infinitely often sees some reward $r(t) < 0$ and thus does not satisfy $PP_{\liminf \geq 0}$. I.e., $\overline{\text{Transience}} \subseteq \overline{PP_{\liminf \geq 0}}$. Now, since **Transience** $= PP_{\liminf \geq 0}$, Theorem 3.5 implies Theorem 3.3 (for the finitely branching case).

However, the connection between **Transience** and $PP_{\liminf \geq 0}$ only holds for finitely branching MDPs. In infinitely branching MDPs, the result about **Transience** (Theorem 3.3) still holds, but the result for $PP_{\liminf \geq 0}$ is different, as shown in Theorem 3.10 in the next section.

3.1.2. Infinitely Branching Case. In this section we consider infinitely branching MDPs. In this setting, ε -optimal strategies for $PP_{\liminf \geq 0}$ require more memory than in the finitely branching case. In the following theorem we show how to obtain ε -optimal deterministic

Markov strategies for $PP_{\liminf \geq 0}$. We do this by deriving ε -optimal MD strategies in $S(\mathcal{M})$ via a reduction to a safety objective.

Theorem 3.10. *Consider an MDP \mathcal{M} with initial state s_0 and a $PP_{\liminf \geq 0}$ objective. For every $\varepsilon > 0$ there exist*

- ε -optimal MD strategies in $S(\mathcal{M})$.
- ε -optimal deterministic Markov strategies in \mathcal{M} .

Proof. Let $\varepsilon > 0$. We work in $S(\mathcal{M})$ by encoding the step counter into the states of \mathcal{M} . Thus $S(\mathcal{M})$ is an acyclic MDP with implicit step counter and corresponding initial state $s'_0 = (s_0, 0)$.

We consider a general (not necessarily MD) ε -optimal strategy σ for $PP_{\liminf \geq 0}$ from s'_0 on $S(\mathcal{M})$, i.e.,

$$\mathcal{P}_{S(\mathcal{M}), s'_0, \sigma}(PP_{\liminf \geq 0}) \geq \text{val}_{S(\mathcal{M}), PP_{\liminf \geq 0}}(s'_0) - \varepsilon. \quad (3.11)$$

Define the safety objective Safety_i which is the objective of never seeing any point reward $< -2^{-i}$. This then allows us to characterize $PP_{\liminf \geq 0}$ in terms of safety objectives.

$$PP_{\liminf \geq 0} = \bigcap_{i \in \mathbb{N}} \text{F}(\text{Safety}_i) \quad (3.12)$$

Now we define the safety objective $\text{Safety}_i^k \stackrel{\text{def}}{=} \text{F}^{\leq k}(\text{Safety}_i)$ to attain Safety_i within at most k steps. This allows us to write

$$\text{F}(\text{Safety}_i) = \bigcup_{k \in \mathbb{N}} \text{Safety}_i^k. \quad (3.13)$$

By continuity of measures from above we get

$$\begin{aligned} 0 &= \mathcal{P}_{S(\mathcal{M}), s'_0, \sigma} \left(\text{F}(\text{Safety}_i) \cap \bigcap_{k \in \mathbb{N}} \overline{\text{Safety}_i^k} \right) \\ &= \lim_{k \rightarrow \infty} \mathcal{P}_{S(\mathcal{M}), s'_0, \sigma} \left(\text{F}(\text{Safety}_i) \cap \overline{\text{Safety}_i^k} \right). \end{aligned}$$

Hence for every $i \in \mathbb{N}$ and $\varepsilon_i \stackrel{\text{def}}{=} \varepsilon \cdot 2^{-i}$ there exists n_i such that

$$\mathcal{P}_{S(\mathcal{M}), s'_0, \sigma} \left(\text{F}(\text{Safety}_i) \cap \overline{\text{Safety}_i^{n_i}} \right) \leq \varepsilon_i. \quad (3.14)$$

Claim 3.11.

$$\mathcal{P}_{S(\mathcal{M}), s'_0, \sigma} \left(\bigcap_{i \in \mathbb{N}} \text{Safety}_i^{n_i} \right) \geq \text{val}_{S(\mathcal{M}), PP_{\liminf \geq 0}}(s'_0) - 2\varepsilon.$$

Proof. The proof is almost identical to the proof of Claim 3.6. Instead of \mathcal{M}' with initial state s_0 we have $S(\mathcal{M})$ with initial state s'_0 , and instead of equations (3.3), (3.2) and (3.5) we use the corresponding equations (3.12), (3.11) and (3.14), respectively. \square

Let $\varphi \stackrel{\text{def}}{=} \bigcap_{i \in \mathbb{N}} \text{Safety}_i^{n_i} \subseteq PP_{\liminf \geq 0}$. It follows from Claim 3.11 that

$$\text{val}_{S(\mathcal{M}), \varphi}(s'_0) \geq \text{val}_{S(\mathcal{M}), PP_{\liminf \geq 0}}(s'_0) - 2\varepsilon. \quad (3.15)$$

The objective φ is a safety objective on $S(\mathcal{M})$. Therefore, since $S(\mathcal{M})$ is acyclic, we can apply Lemma 3.1 to obtain a uniformly ε -optimal MD strategy σ' for φ . Thus

$$\begin{aligned}
& \mathcal{P}_{S(\mathcal{M}), s'_0, \sigma'}(PP_{\liminf \geq 0}) \\
& \geq \mathcal{P}_{S(\mathcal{M}), s'_0, \sigma'}(\varphi) && \text{set inclusion} \\
& \geq \text{val}_{S(\mathcal{M}), \varphi}(s'_0) - \varepsilon && \sigma' \text{ is } \varepsilon\text{-opt.} \\
& \geq \text{val}_{S(\mathcal{M}), PP_{\liminf \geq 0}}(s'_0) - 3\varepsilon. && \text{by (3.15)}
\end{aligned}$$

Thus σ' is a 3ε -optimal MD strategy for $PP_{\liminf \geq 0}$ in $S(\mathcal{M})$.

By Lemma 2.2 this then yields a 3ε -optimal Markov strategy for $PP_{\liminf \geq 0}$ from s_0 in \mathcal{M} , since runs in \mathcal{M} and $S(\mathcal{M})$ coincide wrt. $PP_{\liminf \geq 0}$. \square

The following corollary allows us to re-purpose Theorem 3.10 to obtain an upper bound for optimal strategies.

Corollary 3.12. *Given an MDP \mathcal{M} and initial state s_0 , optimal strategies, where they exist, can be chosen with just a step counter for $PP_{\liminf \geq 0}$.*

Proof. We work in $S(\mathcal{M})$ and we apply Theorem 3.10 to obtain ε -optimal MD strategies from every state of $S(\mathcal{M})$. Since $PP_{\liminf \geq 0}$ is a shift invariant objective, Theorem 3.2 yields an MD strategy that is optimal from every state of $S(\mathcal{M})$ that has an optimal strategy. By Lemma 2.2 we can translate this MD strategy on $S(\mathcal{M})$ back to a Markov strategy in \mathcal{M} , which is optimal for $PP_{\liminf \geq 0}$ from s_0 (provided that s_0 admits any optimal strategy at all). \square

3.2. Lower Bounds. We have just showed that finitely branching ε -optimal $PP_{\liminf \geq 0}$ can be achieved MD, as a result no lower bound exists within the scope of our memory considerations. For the infinitely branching case, we provide a counterexample in which $PP_{\liminf \geq 0}$ and co-Büchi coincide.

Theorem 3.13. *There exists an infinitely branching MDP \mathcal{M} as in Figure 2 with reward implicit in the state and initial state s such that*

- every FR strategy σ is such that $\mathcal{P}_{\mathcal{M}, s, \sigma}(PP_{\liminf \geq 0}) = 0$
- there exists an HD strategy σ such that $\mathcal{P}_{\mathcal{M}, s, \sigma}(PP_{\liminf \geq 0}) = 1$.

Hence, optimal (and even almost-surely winning) strategies and ε -optimal strategies for $PP_{\liminf \geq 0}$ require infinite memory beyond a reward counter.

Proof. This follows directly from [23, Theorem 4] and the observation that in Figure 2, $PP_{\liminf \geq 0}$ and co-Büchi objectives coincide. \square

Consequently, when the MDP \mathcal{M} is infinitely branching and has the reward counter implicit in the state, $PP_{\liminf \geq 0}$ requires at least a step counter.

Note that $TP_{\liminf \geq 0}$ also coincides with co-Büchi in the MDP \mathcal{M} of Figure 2, hence we restate this theorem in terms of $TP_{\liminf \geq 0}$ later in Theorem 5.6.

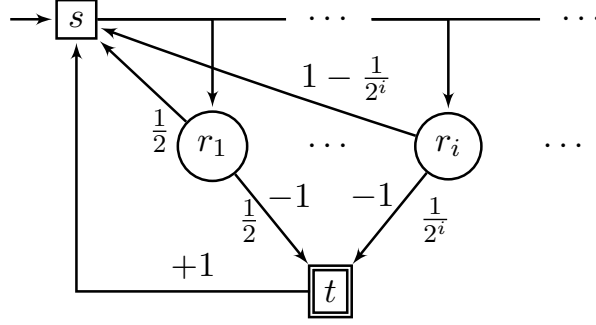


Figure 2: This infinitely branching MDP is adapted from [23, Figure 3] and augmented with a reward structure. (A very similar example has been described in [32, Example 2].) All of the edges carry reward 0 except the edges entering t that carry reward -1 and the edge from t to s carries reward $+1$. As a result, entering t necessarily brings the total reward down to -1 before resetting it to 0. We use a reduction to the co-Büchi objective (i.e., visiting t only finitely often) to show that infinite memory is required for almost-sure as well as ε -optimal strategies for $TP_{\liminf \geq 0}$ as well as $PP_{\liminf \geq 0}$.

Mean Payoff		ε -optimal	Optimal
Finitely branching	Upper Bound	Det(SC+RC) 4.1	Det(SC+RC) 4.2
	Lower Bound	$\neg\text{Rand}(\text{F}+\text{SC})$ 4.12 and $\neg\text{Rand}(\text{F}+\text{RC})$ 4.15	$\neg\text{Rand}(\text{F}+\text{SC})$ 4.20 and $\neg\text{Rand}(\text{F}+\text{RC})$ 4.23
Infinitely branching	Upper Bound	Det(SC+RC) 4.1	Det(SC+RC) 4.2
	Lower Bond	$\neg\text{Rand}(\text{F}+\text{SC})$ 4.12 and $\neg\text{Rand}(\text{F}+\text{RC})$ 4.15	$\neg\text{Rand}(\text{F}+\text{SC})$ 4.20 and $\neg\text{Rand}(\text{F}+\text{RC})$ 4.23

Table 3: Strategy complexity of ε -optimal/optimal strategies for the mean payoff objective in infinitely/finitely branching MDPs.

4. MEAN PAYOFF

4.1. Upper Bounds. In order to tackle the upper bounds for the mean payoff objective $MP_{\liminf \geq 0}$, we work with the acyclic MDP $A(\mathcal{M})$ which encodes both the step counter and the average reward into the state. Once the average reward is encoded into the state, the point payoff coincides with the mean payoff. We use this observation to reduce $MP_{\liminf \geq 0}$ to $PP_{\liminf \geq 0}$ and obtain our upper bounds from the corresponding point payoff results.

Corollary 4.1. *Given an MDP \mathcal{M} and initial state s_0 , there exist ε -optimal strategies σ for $MP_{\liminf \geq 0}$ which use just a step counter and a reward counter.*

Proof. We consider the encoded system $A(\mathcal{M})$ in which both step counter and reward counter are implicit in the state. Recall that the partial mean payoffs in \mathcal{M} correspond exactly to point rewards in $A(\mathcal{M})$. Since $A(\mathcal{M})$ has an encoded step counter, Theorem 3.10 gives us ε -optimal MD strategies for $PP_{\liminf \geq 0}$ in $A(\mathcal{M})$. Lemma 2.6 allows us to translate these strategies back to \mathcal{M} with a memory overhead of just a reward counter and a step counter as required. \square

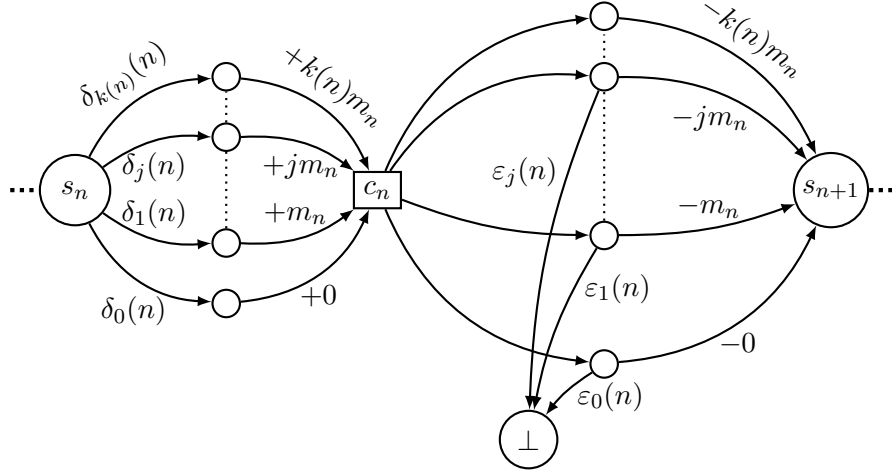


Figure 3: A typical building block with $k(n) + 1$ choices, first random then controlled. The number of choices $k(n) + 1$ grows unboundedly with n . This is the n -th building block of the MDP in Figure 4. The $\delta_i(n)$ and $\varepsilon_i(n)$ are probabilities depending on n and the $\pm im_n$ are transition rewards. We index the successor states of s_n and c_n from 0 to $k(n)$ to match the indexing of the δ 's and ε 's such that the bottom state is indexed with 0 and the top state with $k(n)$.

Corollary 4.2. *Given an MDP \mathcal{M} and initial state s_0 , optimal strategies, where they exist, can be chosen with just a reward counter and a step counter for $MP_{\liminf \geq 0}$*

Proof. We place ourselves in $A(\mathcal{M})$ and apply Theorem 3.10 to obtain ε -optimal MD strategies from every state of $A(\mathcal{M})$. Since $MP_{\liminf \geq 0}$ is a shift invariant objective, Theorem 3.2 yields a single MD strategy that is optimal from every state of $A(\mathcal{M})$ that has an optimal strategy. By Lemma 2.6 we can translate this MD strategy on $A(\mathcal{M})$ back to a strategy on \mathcal{M} with a step counter and a reward counter. Provided that s_0 admits any optimal strategy at all, we obtain an optimal strategy for $MP_{\liminf \geq 0}$ from s_0 that uses only a step counter and a reward counter. \square

4.2. Lower Bounds. In this section we will show that $MP_{\liminf \geq 0}$ always requires at least a step counter and a reward counter, whether it be for ε -optimal or optimal strategies. We introduce in Figure 3 a building block for an MDP which will form the foundation for all of our lower bound results for $MP_{\liminf \geq 0}$. Many of these results also hold for $TP_{\liminf \geq 0}$, so we will restate them in Section 5 in due course.

4.2.1. ε -optimal Strategies. We construct an acyclic MDP \mathcal{M} in which the step counter is implicit in the state as follows.

The system consists of a sequence of gadgets. Figure 3 depicts a typical building block in this system. The system consists of these gadgets chained together as illustrated in Figure 4, starting with n sufficiently high at $n = N^*$. In the controlled choice, there is a small chance in all but the top choice of falling into a \perp state. These \perp states are abbreviations for an infinite chain of states with -1 reward on the transitions and are thus losing. The intuition

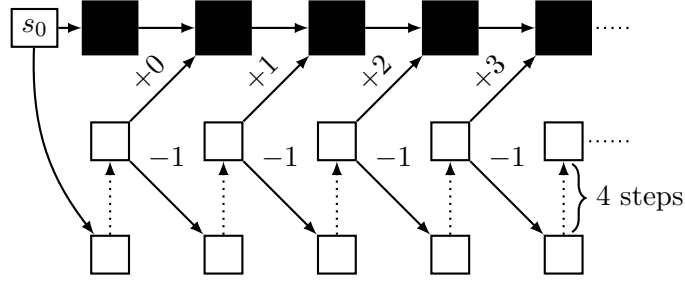


Figure 4: The buildings blocks from Figure 3 represented by black boxes are chained together (n increases as you go to the right). The chain of white boxes allows to skip arbitrarily long prefixes while preserving path length. The positive rewards from the white states to the black boxes reimburse the lost reward accumulated until then. The -1 rewards between white states ensure that skipping gadgets forever is losing.

behind the construction is that there is a random transition with branching degree $k(n) + 1$. Then, the only way to win, in the controlled states, is to play the i -th choice if one arrived from the i -th choice. Thus intuitively, to remember what this choice was, one requires at least $k(n) + 1$ memory modes. That is to say, the one and only way to win is to mimic, and mimicry requires memory. We will present two similar versions of this MDP, initially we present a version in which the step counter is implicit in the state, then later we will present a version in which the reward counter is implicit in the state instead.

Remark 4.3. \mathcal{M} is acyclic, finitely branching and for every state $s \in S$, $\exists n_s \in \mathbb{N}$ such that every path from s_0 to s has length n_s . That is to say the step counter is implicit in the state.

Additionally, the number of transitions in each gadget grows unboundedly with n according to the function $k(n)$. Consequently, we will show that the number of memory modes required to play correctly grows above every finite bound. This will imply that no finite amount of memory suffices for ε -optimal strategies.

Notation: All logarithms are assumed to be in base e .

$$\begin{aligned} \log_1 n &\stackrel{\text{def}}{=} \log n, & \log_{i+1} n &\stackrel{\text{def}}{=} \log(\log_i n) \\ \delta_0(n) &\stackrel{\text{def}}{=} \frac{1}{\log n}, & \delta_i(n) &\stackrel{\text{def}}{=} \frac{1}{\log_{i+1} n}, & \delta_{k(n)}(n) &\stackrel{\text{def}}{=} 1 - \sum_{j=0}^{k(n)-1} \delta_j(n) \\ \varepsilon_0(n) &\stackrel{\text{def}}{=} \frac{1}{n \log n}, & \varepsilon_{i+1}(n) &\stackrel{\text{def}}{=} \frac{\varepsilon_i(n)}{\log_{i+2} n}, & \text{i.e. } \varepsilon_i(n) &= \frac{1}{n \cdot \log n \cdot \log_2 n \cdots \log_{i+1} n}, & \varepsilon_{k(n)}(n) &\stackrel{\text{def}}{=} 0 \\ \text{Tower}(0) &\stackrel{\text{def}}{=} e^0 = 1, & \text{Tower}(i+1) &\stackrel{\text{def}}{=} e^{\text{Tower}(i)}, & N_i &\stackrel{\text{def}}{=} \text{Tower}(i) \end{aligned}$$

Lemma 4.4. *The family of series $\sum_{n > N_j} \delta_j(n) \cdot \varepsilon_i(n)$ is divergent for all $i, j \in \mathbb{N}$, $i < j$. Additionally, the related family of series $\sum_{n > N_i} \delta_i(n) \cdot \varepsilon_i(n)$ is convergent for all $i \in \mathbb{N}$.*

Proof. These are direct consequences of Cauchy's Condensation Test. \square

Definition 4.5. We define $k(n)$, the rate at which the number of transitions grows. We define $k(n)$ in terms of fast growing functions g , Tower and h defined for $i \geq 1$ as follows:

$$g(i) \stackrel{\text{def}}{=} \min \left\{ N : \left(\sum_{n>N} \delta_{i-1}(n) \varepsilon_{i-1}(n) \right) \leq 2^{-i} \right\}, \quad h(1) \stackrel{\text{def}}{=} 2$$

$$h(i+1) \stackrel{\text{def}}{=} \left\lceil \max \left\{ g(i+1), \text{Tower}(i+2), \min \left\{ m+1 \in \mathbb{N} : \sum_{n=h(i)}^m \varepsilon_{i-1}(n) \geq 1 \right\} \right\} \right\rceil.$$

Note that function g is well defined by Lemma 4.4, and $h(i+1)$ is well defined since for all i , $\sum_{n=h(i)}^{\infty} \varepsilon_{i-1}(n)$ diverges to infinity. $k(n)$ is a slow growing unbounded step function defined in terms of h as $k(n) \stackrel{\text{def}}{=} h^{-1}(n)$. The Tower function features in the definition to ensure that the transition probabilities are always well defined. g and h are used to smooth the proofs of Lemma 4.7 and Lemma 4.8 respectively. *Notation:* $N^* \stackrel{\text{def}}{=} \min\{n \in \mathbb{N} : k(n) = 1\}$. This is intuitively the first natural number for which the construction is well defined.

The reward m_n which appears in the n -th gadget is defined such that it outweighs any possible reward accumulated up to that point in previous gadgets. As such we define $m_n \stackrel{\text{def}}{=} 2k(n) \sum_{i=N^*}^{n-1} m_i$, with $m_{N^*} \stackrel{\text{def}}{=} 1$ and where $k(n)$ is the branching degree.

To simplify the notation, the state s_0 in our theorem statements refers to s_{N^*} .

Lemma 4.6. *For $k(n) \geq 1$, the transition probabilities in the gadgets are well defined.*

Proof. Recall that $\text{Tower}(i)$ is i repeated exponentials. Thus, $\log(\text{Tower}(i)) = \text{Tower}(i-1)$.

When checking whether probabilities in a given gadget are well defined, first we choose a gadget. The choice of gadget gives us a branching degree $k(n) + 1$ which in turn lower bounds the value of n in that gadget. So for a branching degree of $k(n) + 1$, we have n lower bounded by $\text{Tower}(k(n) + 1)$ by definition of $k(n)$.

We need to show that $\sum_{i=0}^{k(n)-1} \delta_i(n) \leq 1$. Indeed, we have that:

$$\sum_{i=0}^{k(n)-1} \delta_i(n) \leq \sum_{i=0}^{k(n)-1} \frac{1}{\log_{i+1}(\text{Tower}(k(n) + 1))} = \sum_{i=1}^{k(n)} \frac{1}{\text{Tower}(i)} < \sum_{i=1}^{k(n)} \frac{1}{e^i} < \sum_{i=1}^{k(n)} \frac{1}{2^i} < 1.$$

Hence, for $k(n) \geq 1$, the transition probabilities are well defined, i.e. $\delta_0(n), \delta_1(n), \dots, \delta_{k(n)}(n)$ do indeed sum to 1. \square

Lemma 4.7. *For every $\varepsilon > 0$, there exists a strategy σ_ε with $\mathcal{P}_{\mathcal{M}, s_0, \sigma_\varepsilon}(MP_{\liminf \geq 0}) \geq 1 - \varepsilon$ that cannot fail unless it hits a \perp state. Formally, $\mathcal{P}_{\mathcal{M}, s_0, \sigma_\varepsilon}(MP_{\liminf \geq 0} \wedge \mathbf{G}(\neg \perp)) = \mathcal{P}_{\mathcal{M}, s_0, \sigma_\varepsilon}(\mathbf{G}(\neg \perp)) \geq 1 - \varepsilon$. So in particular, $\text{val}_{\mathcal{M}, MP_{\liminf \geq 0}}(s_0) = 1$.*

Proof. We define a strategy σ which in c_n always mimics the choice in s_n . We first prove that playing this way gives us a positive chance of winning. Then we show that there are strategies σ_ε that attain $1 - \varepsilon$ from s_0 without hitting a \perp state. This implies in particular that $\text{val}_{\mathcal{M}, MP_{\liminf \geq 0}}(s_0) = 1$.

Playing according to σ , the only way to lose is by dropping into the \perp state. This is because by mimicking, the player finishes each gadget with a reward of 0. In the n -th gadget, the chance of reaching the \perp state is $\sum_{j=0}^{k(n)-1} \delta_j(n) \cdot \varepsilon_j(n)$. Thus, the probability of

surviving while playing in all the gadgets is

$$\prod_{n \geq N^*} \left(1 - \sum_{j=0}^{k(n)-1} \delta_j(n) \cdot \varepsilon_j(n) \right).$$

However, by Proposition 2.7, this product is strictly greater than 0 if and only if the sum

$$\sum_{n \geq N^*} \left(\sum_{i=0}^{k(n)-1} \delta_i(n) \varepsilon_i(n) \right)$$

is finite. With some rearranging exploiting the definition of $k(n)$ we see that this is indeed the case:

$$\begin{aligned} & \sum_{n \geq N^*} \left(\sum_{i=0}^{k(n)-1} \delta_i(n) \varepsilon_i(n) \right) \\ & \leq \sum_{i \geq 1} \left(\sum_{n=g(i)}^{\infty} \delta_{i-1}(n) \varepsilon_{i-1}(n) \right) && \text{by definition of } k(n) \\ & \leq \sum_{i \geq 1} 2^{-i} && \text{by definition of } g(n) \\ & \leq 1 \end{aligned}$$

Hence the player has a non zero chance of winning.

When playing with the ability to skip gadgets, as illustrated in Figure 4, all runs not visiting a \perp state are winning since the total reward never dips below 0. Hence $\mathcal{P}_{\mathcal{M}, s_0, \sigma_\varepsilon}(MP_{\liminf \geq 0} \wedge \neg \perp) = \mathcal{P}_{\mathcal{M}, s_0, \sigma_\varepsilon}(\neg \perp)$. Thus the idea is to skip an arbitrarily long prefix of gadgets to push the chance of winning ε close to 1 by pushing the chance of visiting a \perp state ε close to 0. From the N -th state, for $N \geq N^*$, the chance of winning is

$$\prod_{n \geq N} \left(1 - \sum_{j=0}^{k(n)-1} \delta_j(n) \cdot \varepsilon_j(n) \right) > 0$$

By Proposition 2.8 this can be made arbitrarily close to 1 by choosing N sufficiently large.

Let $N_\varepsilon \stackrel{\text{def}}{=} \min \left\{ N \in \mathbb{N} \mid \prod_{n \geq N} \left(1 - \sum_{j=0}^{k(n)-1} \delta_j(n) \cdot \varepsilon_j(n) \right) \geq 1 - \varepsilon \right\}$. Now define the strategy σ_ε to be the strategy that plays like σ after skipping forwards by N_ε gadgets. Thus, by definition σ_ε attains $1 - \varepsilon$ for all $\varepsilon > 0$.

Thus, by playing σ_ε for an arbitrarily small ε the chance of winning must be arbitrarily close to 1. Hence, $\text{val}_{\mathcal{M}, MP_{\liminf \geq 0}}(s_0) = 1$. \square

Lemma 4.8. $\sum_{n=k^{-1}(2)}^{\infty} \frac{1}{2} \delta_{j(n)}(n) \varepsilon_{i(n)}(n)$ diverges for all $i(n), j(n) \in \{0, 1, \dots, k(n) - 1\}$ with $i(n) < j(n)$.

Proof. This result is not immediate, because the indexing functions $i(n)$ and $j(n)$ may grow with $k(n)$ as n increases.

Under the assumption that $i(n) < j(n)$ we have that

$$\delta_{j(n)}(n) \varepsilon_{i(n)}(n) \geq \delta_{j(n)}(n) \varepsilon_{j(n)-1}(n) \geq \delta_{k(n)-1}(n) \varepsilon_{k(n)-2}(n) = \varepsilon_{k(n)-1}(n).$$

Thus it suffices to show that $\sum_{n=k^{-1}(2)}^{\infty} \varepsilon_{k(n)-1}(n)$ diverges:

$$\begin{aligned}
\sum_{n=k^{-1}(2)}^{\infty} \varepsilon_{k(n)-1}(n) &= \sum_{a=2}^{\infty} \sum_{n=k^{-1}(a)}^{k^{-1}(a+1)-1} \varepsilon_{a-1}(n) && \text{splitting the sum up} \\
&= \sum_{a=2}^{\infty} \sum_{n=h(a)}^{h(a+1)-1} \varepsilon_{a-1}(n) && k(n) = h^{-1}(n) \\
&\geq \sum_{a=2}^{\infty} 1 && \text{definition of } h(n)
\end{aligned}$$

Note that the definition of $h(i)$ says exactly that a block of the form $\sum_{n=h(a)}^{h(a+1)-1} \varepsilon_{a-1}(n)$ is at least 1. Hence $\sum_{n=k^{-1}(2)}^{\infty} \frac{1}{2} \delta_{j(n)}(n) \varepsilon_{i(n)}(n)$ diverges as required. \square

Lemma 4.9. *For any sequence $\{\alpha_n\}$, where $\alpha_n \in [0, 1]$ for all n , and any functions $i(n), j(n) : \mathbb{N} \rightarrow \mathbb{N}$ with $i(n), j(n) \in \{0, 1, \dots, k(n) - 1\}$, $i(n) < j(n)$ for all n , the following sum diverges:*

$$\sum_{n=k^{-1}(2)}^{\infty} \left(\delta_{j(n)}(n) (\alpha_n \varepsilon_{j(n)}(n) + (1 - \alpha_n) \varepsilon_{i(n)}(n)) + \delta_{i(n)}(n) (\alpha_n + (1 - \alpha_n) \varepsilon_{i(n)}(n)) \right). \quad (4.1)$$

Proof. We can narrow our focus by noticing that

$$\begin{aligned}
&\sum_{n=k^{-1}(2)}^{\infty} \left(\delta_{j(n)}(n) (\alpha_n \varepsilon_{j(n)}(n) + (1 - \alpha_n) \varepsilon_{i(n)}(n)) + \delta_{i(n)}(n) (\alpha_n + (1 - \alpha_n) \varepsilon_{i(n)}(n)) \right) \\
&= \sum_{n=k^{-1}(2)}^{\infty} \alpha_n \delta_{j(n)}(n) \varepsilon_{j(n)}(n) + (1 - \alpha_n) \delta_{i(n)}(n) \varepsilon_{i(n)}(n) && \text{Convergent by def. of } \delta_i(n), \varepsilon_i(n) \\
&+ \sum_{n=k^{-1}(2)}^{\infty} (1 - \alpha_n) \delta_{j(n)}(n) \varepsilon_{i(n)}(n) + \alpha_n \delta_{i(n)}(n)
\end{aligned}$$

Hence the divergence of (4.1) depends only on the divergence of

$$\sum_{n=k^{-1}(2)}^{\infty} (1 - \alpha_n) \delta_{j(n)}(n) \varepsilon_{i(n)}(n) + \alpha_n \delta_{i(n)}(n).$$

No matter how the sequence $\{\alpha_n\}$ behaves, for every n we have that either $\alpha_n \geq 1/2$ or $1 - \alpha_n \geq 1/2$. Hence for every n it is the case that

$$\begin{aligned}
(1 - \alpha_n) \delta_{j(n)}(n) \varepsilon_{i(n)}(n) + \alpha_n \delta_{i(n)}(n) &\geq \frac{1}{2} \delta_{j(n)}(n) \varepsilon_{i(n)}(n) \\
&\text{or} \\
&\geq \frac{1}{2} \delta_{i(n)}(n)
\end{aligned}$$

Define the function f as follows:

$$f(n) = \begin{cases} \frac{1}{2}\delta_{i(n)}(n) & \text{if } \alpha_n \geq 1/2 \\ \frac{1}{2}\delta_{j(n)}(n)\varepsilon_{i(n)}(n) & \text{otherwise} \end{cases}$$

Hence no matter how $\{\alpha_n\}$ behaves, we have that

$$\begin{aligned} \sum_{n=k^{-1}(2)}^{\infty} \left(\delta_{j(n)}(n)(\alpha_n\varepsilon_{j(n)}(n) + (1 - \alpha_n)\varepsilon_{i(n)}(n)) + \delta_{i(n)}(n)(\alpha_n + (1 - \alpha_n)\varepsilon_{i(n)}(n)) \right) \\ \geq \sum_{n=k^{-1}(2)}^{\infty} f(n). \end{aligned}$$

We know that both $\sum_{n=k^{-1}(2)}^{\infty} \frac{1}{2}\delta_{j(n)}(n)\varepsilon_{i(n)}(n)$ and $\sum_{n=k^{-1}(2)}^{\infty} \frac{1}{2}\delta_{i(n)}(n)$ diverge for all $i(n), j(n) \in \{0, 1, \dots, k(n) - 1\}$, $i(n) < j(n)$, as shown in Lemma 4.8.

Thus $\sum_{n=k^{-1}(2)}^{\infty} \frac{1}{2}\delta_{j(n)}(n)\varepsilon_{i(n)}(n)$ and $\sum_{n=k^{-1}(2)}^{\infty} \frac{1}{2}\delta_{i(n)}(n)$ must also diverge no matter how $i(n)$ and $j(n)$ behave. As a result it must be the case that $\sum_{n=k^{-1}(2)}^{\infty} f(n)$ diverges. Hence (4.1) must be divergent as desired as $i(n)$ and $j(n)$ vary for $n \geq k^{-1}(2)$. \square

Lemma 4.10. *For any FR strategy σ , almost surely either the mean payoff dips below -1 infinitely often, or the run hits a \perp state, i.e. $\mathcal{P}_{\mathcal{M}, \sigma, s_0}(MP_{\liminf \geq 0}) = 0$.*

Outline of the proof. Let σ be some FR strategy with k memory modes. We prove a *lower bound* e_n on the probability of a local error (reaching a \perp state, or seeing a mean payoff ≤ -1) in the current n -th gadget. This lower bound e_n holds regardless of events in past gadgets, regardless of the memory mode of σ upon entering the n -th gadget, and cannot be improved by σ randomizing its memory updates.

The main idea is that, once $k(n) > k + 1$ (which holds for $n \geq N'$ sufficiently large) by the Pigeonhole Principle there will always be a memory mode confusing at least two different branches $i(n), j(n) \neq k(n)$ of the previous random choice at state s_n . This confusion yields a probability $\geq e_n$ of reaching a \perp state or seeing a mean payoff ≤ -1 , regardless of events in past gadgets and regardless of the memory upon entering the n -th gadget. We show that $\sum_{n \geq N'} e_n$ is a *divergent* series. Thus, by Proposition 2.7, $\prod_{n \geq N'} (1 - e_n) = 0$. Hence, $\mathcal{P}_{\mathcal{M}, \sigma, s_0}(MP_{\liminf \geq 0}) \leq \prod_{n \geq N'} (1 - e_n) = 0$. \square

Full proof. Let σ be some FR strategy with k memory modes. Our MDP consists of a linear sequence of gadgets (Figure 3) and is in particular acyclic. The n -th gadget is entered at state s_n and takes 4 steps. Locally in the n -th gadget there are 3 possible scenarios:

- (1) The random transition picks some branch i at s_n and the strategy then picks a branch $j > i$ at c_n .

By the definition of the payoffs (multiples of m_n ; cf. Definition 4.5), this means that we see a mean payoff ≤ -1 , regardless of events in past gadgets. This is because the numbers m_n grow so quickly with n that even the combined maximal possible rewards of all past gadgets are so small in comparison that they do not matter for the outcome in the n -th gadget, i.e., rewards from past gadgets cannot help to avoid seeing a mean payoff ≤ -1 in the above scenario.

- (2) We reach the losing sink \perp (and thus will keep seeing a mean payoff ≤ -1 forever). This happens with probability $\varepsilon_j(n)$ if the strategy picks some branch j at c_n , regardless of past events.
- (3) All other cases.

As explained above, due to the definition of the rewards (Definition 4.5), events in past gadgets do not make the difference between (1),(2),(3) in the current gadget. It just depends on the choices of the strategy σ in the current gadget.

Let Bad_n be the event of seeing either of the two unfavorable outcomes (1) or (2) in the n -th gadget. Let p_n be the probability of Bad_n under strategy σ . Since σ has memory, the probabilities p_n are not necessarily independent. However, we show *lower bounds* $e_n \leq p_n$ that hold universally for every FR strategy σ with $\leq k$ memory modes and every n such that $k(n) > k + 1$. The lower bound e_n will hold regardless of the memory mode of σ upon entering the n -th gadget.

Memory updates. First we show that σ randomizing its memory update after observing the random transition from state s_n does *not* help to reduce the probability of event Bad_n . I.e., we show that without restriction σ can update its memory deterministically after observing the transition from state s_n .

Once in the controlled state c_n , the strategy σ can base its choice only on the current state (always c_n in the n -th gadget) and on the current memory mode. Thus, in state c_n , in each memory mode m , the strategy has to pick a distribution $\mathcal{D}_m^{c_n}$ over the available transitions from c_n . By the finiteness of the number of memory modes of σ (just $\leq k$ by our assumption above), for each possible reward level x (obtained in the step from the preceding random transition from s_n) there is a best memory mode $m(x)$ such that $\mathcal{D}_{m(x)}^{c_n}$ is optimal (in the sense of minimizing the probability of event Bad_n) for that particular reward level x . (In case of a tie, just use an arbitrary tie break, e.g., some pre-defined linear order on the memory modes.)

Therefore, upon witnessing a reward level x in the random transition from state s_n , the strategy σ can minimize the probability of event Bad_n by *deterministically* setting its memory to $m(x)$. Thus randomizing its memory update does not help to reduce the probability of Bad_n , and we may assume without restriction that σ updates its memory deterministically.

(Note that the above argument only works because it is local to the current gadget where we have a finite number of decisions (here just one), we have a finite number of memory modes, and a one-dimensional criterion for local optimality (minimizing the probability of event Bad_n). We do *not* claim that randomized memory updates are useless for every strategy in every MDP and every objective.)

Claim 4.11. Assume that the transitions $i(n)$ and $j(n)$ (with $i(n) < j(n)$) leading to state c_n are confused in the memory of the strategy. Then we can assume without restriction that the strategy only plays transitions $i(n)$ and $j(n)$ with nonzero probability from state c_n , since every other behavior yields a higher probability of the event Bad_n (cf. Figure 5).

Proof. When confusing transitions $i(n)$ and $j(n)$ with $i(n) < j(n)$, the player's choice of transition from c_n can be broken down into 5 distinct cases. The player can choose transition $x(n)$ as follows.

- (1) $x(n) = i(n)$
- (2) $x(n) = j(n)$
- (3) $x(n) > j(n)$

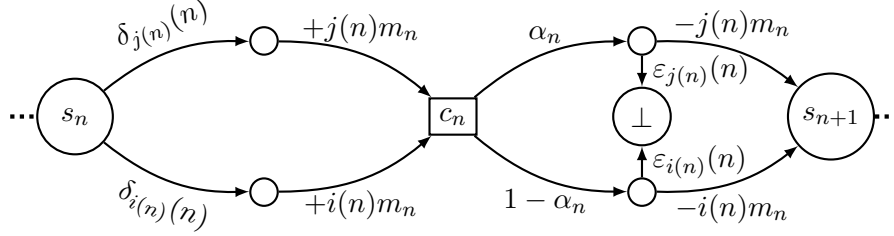


Figure 5: When transitions $i(n)$ and $j(n)$ are confused in the player's memory, the player's choice is at least as bad as the reduced payoff in this simplified gadget.

(4) $x(n) < i(n)$

(5) $i(n) < x(n) < j(n)$

Case 1 leads to a probability of Bad_n of $\delta_{j(n)}(n)\varepsilon_{i(n)}(n) + \delta_{i(n)}(n)\varepsilon_{i(n)}(n)$.

Case 2 leads to a probability of Bad_n of $\delta_{j(n)}(n)\varepsilon_{j(n)}(n) + \delta_{i(n)}(n)$.

Case 3 leads to a mean payoff ≤ -1 (and thus Bad_n) with probability 1. This is the worst possible case.

Case 4 leads to a probability of Bad_n of $\delta_{j(n)}(n)\varepsilon_{x(n)}(n) + \delta_{i(n)}(n)\varepsilon_{x(n)}(n) > \delta_{j(n)}(n)\varepsilon_{i(n)}(n) + \delta_{i(n)}(n)\varepsilon_{i(n)}(n)$, i.e., this is worse than Case 1.

Case 5 leads to a probability of Bad_n of $\delta_{j(n)}(n)\varepsilon_{x(n)}(n) + \delta_{i(n)}(n) > \delta_{j(n)}(n)\varepsilon_{j(n)}(n) + \delta_{i(n)}(n)$, i.e., this is worse than Case 2.

Hence, without restriction we can assume that only cases 1 and 2 will get played with positive probability, that is to say that in state c_n the strategy will only randomize over the outgoing transitions $i(n)$ and $j(n)$. \square

The lower bounds e_n . Now we consider an FR strategy σ that without restriction updates its memory *deterministically* after each random choice (from state s_n) in the n -th gadget. It can still randomize its actions, however.

Let N' be the minimal number such that for all $n \geq N'$ we have $k(n) > k + 1$. In particular, this implies $N' \geq k^{-1}(2)$, and thus we can apply Lemma 4.9 later.

Once $n \geq N'$, then by the Pigeonhole Principle there will always be a memory mode confusing at least two different transitions $i(n), j(n) \neq k(n)$ from state s_n to c_n . Note that this holds regardless of the memory mode of σ upon entering the n -th gadget. (The strategy might confuse many other scenarios, but just one confused pair $i(n), j(n) \neq k(n)$ is enough for our lower bound.) Without loss of generality, let $j(n)$ be larger of the two confused transitions, i.e., $i(n) < j(n)$. Let $i(n)$ and $j(n)$ be two functions taking values in $\{0, 1, \dots, k(n) - 1\}$ where $i(n) < j(n)$ for all n .

Confusing two transitions $i(n)$ and $j(n)$ from s_n to c_n (where without restriction $i(n) < j(n)$), the strategy is in the same memory mode afterwards. However, it can still randomize its choices in state c_n . To prove our lower bound on the probability of Bad_n , it suffices to consider the case where the strategy only randomizes over the outgoing transitions $i(n)$ and $j(n)$ from state c_n . This is because, by Claim 4.11, every other behavior would perform even worse, in the sense of yielding a higher probability of Bad_n .

That is to say that the strategy picks the higher $j(n)$ -th branch with some probability α_n and the lower $i(n)$ -th branch with probability $1 - \alpha_n$. (We leave the probabilities α_n

unspecified here. Using Lemma 4.9, we'll show that our result holds regardless of their values.)

The local chance of the event Bad_n is then lower bounded by

$$e_n \stackrel{\text{def}}{=} \delta_{j(n)}(n)(\alpha_n \varepsilon_{j(n)}(n) + (1 - \alpha_n) \varepsilon_{i(n)}(n)) + \delta_{i(n)}(n)(\alpha_n + (1 - \alpha_n) \varepsilon_{i(n)}(n)).$$

The term above just expresses a case distinction. In the first scenario, the random transition chooses the $j(n)$ -th branch (with probability $\delta_{j(n)}(n)$) and then the strategy chooses the $j(n)$ -th branch with probability α_n and the lower $i(n)$ -th branch with probability $1 - \alpha_n$, and you obtain the respective chances of reaching the sink \perp . In the second scenario, the random transition chooses the $i(n)$ -th branch (with probability $\delta_{i(n)}(n)$). If the strategy then chooses the higher $j(n)$ -th branch (with probability α_n) then we have outcome (1), yielding a mean payoff ≤ -1 . If the strategy chooses the $i(n)$ -th branch (with probability $1 - \alpha_n$) then we still have a chance of $\varepsilon_{i(n)}(n)$ of reaching the sink.

Since, as shown above, randomized memory updates do not help to reduce the probability of Bad_n , the lower bound e_n for deterministic updates carries over to the general case. Thus, even for general randomized FR strategies σ with k memory modes, the probability of event Bad_n in the n -th gadget (for $n \geq N'$) is lower bounded by e_n , regardless of the memory mode \mathbf{m} upon entering the gadget and regardless of events in past gadgets. We write $\sigma[\mathbf{m}]$ for the strategy σ in memory mode \mathbf{m} and obtain

$$\forall n \geq N'. \forall \mathbf{m}. \mathcal{P}_{\mathcal{M}, \sigma[\mathbf{m}], s_n}(Bad_n) \geq e_n \quad (4.2)$$

The final step. Let $Bad \stackrel{\text{def}}{=} \cup_n Bad_n$.

Since $i(n), j(n) \neq k(n)$ and $N' \geq k^{-1}(2)$, we apply Lemma 4.9 to conclude that the series $\sum_{n=N'}^{\infty} e_n = \sum_{n=N'}^{\infty} \delta_{j(n)}(n)(\alpha_n \varepsilon_{j(n)}(n) + (1 - \alpha_n) \varepsilon_{i(n)}(n)) + \delta_{i(n)}(n)(\alpha_n + (1 - \alpha_n) \varepsilon_{i(n)}(n))$ is divergent, regardless of the behavior of $i(n), j(n)$ or the sequence $\{\alpha_n\}$.

Finally, we obtain

$$\begin{aligned} & \mathcal{P}_{\mathcal{M}, \sigma, s_0}(MP_{\liminf \geq 0}) \\ & \leq \mathcal{P}_{\mathcal{M}, \sigma, s_0}(\text{FG} \neg Bad) && \text{set inclusion} \\ & = \mathcal{P}_{\mathcal{M}, \sigma, s_0} \left(\bigcup_l \text{F}^{\leq l} \text{G} \neg Bad \right) && \text{def. of F} \\ & = \lim_{l \rightarrow \infty} \mathcal{P}_{\mathcal{M}, \sigma, s_0}(\text{F}^{\leq l} \text{G} \neg Bad) && \text{continuity of measures} \\ & \leq \lim_{l \rightarrow \infty} \mathcal{P}_{\mathcal{M}, \sigma, s_0} \left(\bigcap_{n \geq l/4} \neg Bad_n \right) && \text{4 steps per gadget} \\ & \leq \lim_{4N' \leq l \rightarrow \infty} \prod_{n \geq l/4 \geq N'} (\max_{\mathbf{m}} \mathcal{P}_{\mathcal{M}, \sigma[\mathbf{m}], s_n}(\neg Bad_n)) && \begin{array}{l} \text{linear sequence of gadgets,} \\ \text{finite memory and past events} \\ \text{do not help to avoid } Bad_n \end{array} \\ & \leq \lim_{4N' \leq l \rightarrow \infty} \prod_{n \geq l/4 \geq N'} (1 - e_n) && \text{by (4.2)} \\ & = \lim_{4N' \leq l \rightarrow \infty} 0 && \text{divergence of } \sum_{n=N'}^{\infty} e_n \text{ and Proposition 2.7} \\ & = 0 \end{aligned}$$

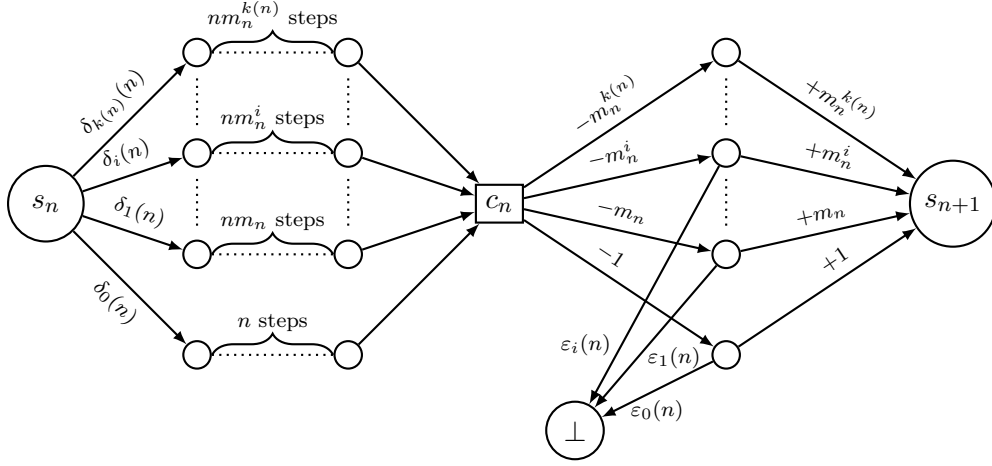


Figure 6: All transition rewards are 0 unless specified. Recall that $\sum \delta_i(n) \cdot \varepsilon_i(n)$ is convergent and $\sum \delta_j(n) \cdot \varepsilon_i(n)$ is divergent for all i, j with $j > i$. The negative reward incurred before falling into the \perp state is reimbursed. We do not show it in the figure for readability. In the state before s_{n+1} , if the correct transition was chosen, the mean payoff is $-1/n$. If the incorrect transition was chosen, then either the mean payoff is $< -m_n/n$, or the risk of falling into \perp is too high.

□

Lemma 4.7 and Lemma 4.10 yield the following theorem.

Theorem 4.12. *There exists a countable, finitely branching and acyclic MDP \mathcal{M} whose step counter is implicit in the state for which $\text{val}_{\mathcal{M}, MP_{\liminf \geq 0}}(s_0) = 1$ and any FR strategy σ is such that $\mathcal{P}_{\mathcal{M}, s_0, \sigma}(MP_{\liminf \geq 0}) = 0$. In particular, there are no ε -optimal k -bit Markov strategies for any $k \in \mathbb{N}$ and any $\varepsilon < 1$ for $MP_{\liminf \geq 0}$ in countable MDPs.*

Proof. Proved by Lemma 4.7 and Lemma 4.10. □

Theorem 4.12 shows that even if the step counter is implicit in the state, infinite memory is still required. We now adapt our construction from Figure 3 such that instead the current total reward is implicit in the state, in order to show that a reward counter plus arbitrary finite memory does not suffice for (ε) -optimal strategies for $MP_{\liminf \geq 0}$ either.

We use the example from Figure 6. It is very similar to Figure 3, but differs in the following ways.

- The current total reward level is implicit in each state.
- The step counter is no longer implicit in the state.
- In the random choice, instead of changing the reward levels in each choice, it is the path length that differs.
- The definition of m_n is different, it is now $m_n \stackrel{\text{def}}{=} \sum_{i=N^*}^{n-1} m_i^{k(n)}$ with $m_{N^*} \stackrel{\text{def}}{=} 1$.

We construct a finitely branching acyclic MDP \mathcal{M}_{RI} (Reward Implicit) which has the total reward implicit in the state. We do so by chaining together the gadgets from Figure 6 as is shown in Figure 4.

In order to convince ourselves that the history of the play in past gadgets does not affect the outcome of the current gadget, we do a brief analysis of the path length and total reward involved in a run going through the n th gadget. Consider the scenario where the play took the i -th random choice. In this case, the path length is upper bounded by $\left(\sum_{i=N^*}^{n-1} 4 + m_i^{k(i)}\right) + 4 + m_n^i \leq 2m_n^i$. In the case where the player chooses the j th controlled choice with $j \geq i$, this gives us an average reward of $-\frac{m_n^j}{2nm_n^i}$. This is < -1 when $j > i$ and converges to 0 with $-\frac{1}{n}$ when $j = i$. The choices $j < i$ are losing due to the risk of falling into the losing sink as described previously.

Hence, the analysis within each gadget still reduces to mimicking the random choice in the controlled choice. This allows us to simply reuse the results from the step counter encoded case in order to obtain symmetrical results for the reward counter encoded case.

Lemma 4.13. $\text{val}_{\mathcal{M}_{\text{RI}}, MP_{\liminf \geq 0}}((s_0, 0)) = 1$.

Proof. We define a strategy σ which, in c_n always mimics the random choice in s_n . Playing according to σ , the only way to lose is by dropping into the bottom state. This is because by mimicking, the mean payoff in each gadget is lower bounded by $-1/n$. The rest of the proof is identical to Lemma 4.7. \square

Lemma 4.14. Any FR strategy σ in \mathcal{M}_{RI} is such that $\mathcal{P}_{\mathcal{M}_{\text{RI}}, s_0, \sigma}(MP_{\liminf \geq 0}) = 0$.

Proof. When playing with finitely many memory modes, there are two ways for a run in \mathcal{M}_{RI} to lose. Either it falls into a losing sink, or it never falls into a sink but its mean payoff is < -1 . The proof that either of these occurs with probability 1 is the same as in Lemma 4.10. \square

Theorem 4.15. There exists a countable, finitely branching, acyclic MDP \mathcal{M}_{RI} with initial state $(s_0, 0)$ with the total reward implicit in the state such that

- $\text{val}_{\mathcal{M}_{\text{RI}}, MP_{\liminf \geq 0}}((s_0, 0)) = 1$,
- for all FR strategies σ , we have $\mathcal{P}_{\mathcal{M}_{\text{RI}}, (s_0, 0), \sigma}(MP_{\liminf \geq 0}) = 0$.

Proof. This follows from Lemma 4.13 and Lemma 4.14. \square

Remark 4.16. The MDPs from Figure 3 and Figure 6 show that good strategies for $MP_{\liminf \geq 0}$ require “at least” a reward counter and a step counter, respectively. There does, of course, exist a *single* MDP where good strategies for $MP_{\liminf \geq 0}$ require at least both a step counter and a reward counter. We construct such an MDP by ‘gluing’ the two different MDPs together via an initial random state which points to each with probability $1/2$.

4.2.2. Optimal Strategies. Even for acyclic MDPs with the step counter implicit in the state, optimal (and even almost sure winning) strategies for $MP_{\liminf \geq 0}$ require infinite memory. To prove this, we consider a variant of the MDP from the previous section which has been augmented to include restarts from the \perp states. For the rest of the section, \mathcal{M} is the MDP constructed in Figure 7. Initially the gadgets used are like Figure 3, then we present similar results using Figure 6 as the gadgets instead.

Remark 4.17. \mathcal{M} is acyclic, finitely branching and the step counter is implicit in the state. We now refer to the rows of Figure 7 as gadgets, i.e., a gadget is a single instance of Figure 4 where the \perp states lead to the next row.

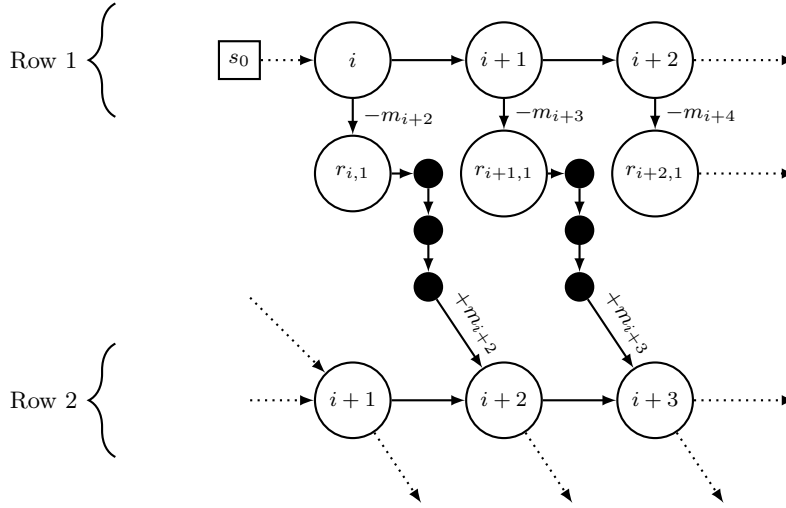


Figure 7: Each row represents a copy of the MDP depicted in Figure 4. Each white circle labeled with a number i represents the correspondingly numbered gadget (like in Figure 3) from that MDP. Now, instead of the bottom states in each gadget leading to an infinite losing chain, they lead to a restart state $r_{i,j}$ which leads to a fresh copy of the MDP (in the next row). Each restart incurs a penalty guaranteeing that the mean payoff dips below -1 before refunding it and continuing on in the next copy of the MDP. The states $r_{i,j}$ are labeled such that the j indicates that if a run sees this state, then it is the j th restart. The i indicates that the run entered the restart state from the i th gadget of the current copy of the MDP. The black states are dummy states inserted in order to preserve path length throughout.

Lemma 4.18. *There exists a strategy σ such that $\mathcal{P}_{\mathcal{M},\sigma,s_0}(MP_{\liminf \geq 0}) = 1$.*

Outline of the proof. Recall the strategy $\sigma_{1/2}$ defined in Lemma 4.7 which achieves at least $1/2$ in each gadget that it is played in. We then construct the almost surely winning strategy σ by concatenating $\sigma_{1/2}$ strategies in the sense that σ plays just like $\sigma_{1/2}$ in each gadget from each gadget's start state.

Since σ achieves at least $1/2$ in every gadget that it sees, with probability 1, runs generated by σ restart only finitely many times. The intuition is then that a run restarting finitely many times must spend an infinite tail in some final gadget. Since σ mimics in every controlled state, not restarting anymore directly implies that the total payoff is eventually always ≥ 0 . Hence all runs generated by σ and restarting only finitely many times satisfy $MP_{\liminf \geq 0}$. Therefore all but a nullset of runs generated by σ are winning, i.e. $\mathcal{P}_{\mathcal{M},s_0,\sigma}(MP_{\liminf \geq 0}) = 1$. \square

Full Proof. We will show that there exists a strategy σ that satisfies the mean payoff objective with probability 1 from s_0 . Towards this objective we recall the strategy $\sigma_{1/2}$ defined in Lemma 4.7. In a given gadget of this MDP with restarts, playing $\sigma_{1/2}$ in said gadget, there is a probability of at most $1/2$ of restarting in that gadget. We then construct strategy σ by concatenating $\sigma_{1/2}$ strategies in the sense that σ plays just like $\sigma_{1/2}$ in each gadget from each gadget's start state.

Let \mathfrak{R} be the set of runs induced by σ from s_0 . We partition \mathfrak{R} into the sets \mathfrak{R}_i and \mathfrak{R}_∞ of runs such that $\mathfrak{R} = (\bigcup_{i=0}^\infty \mathfrak{R}_i) \cup \mathfrak{R}_\infty$. We define for $i = 0$

$$\mathfrak{R}_0 \stackrel{\text{def}}{=} \{\rho \in \mathfrak{R} \mid \forall \ell \in \mathbb{N}. \neg F(r_{\ell,1})\},$$

for $i \geq 1$

$$\mathfrak{R}_i \stackrel{\text{def}}{=} \{\rho \in \mathfrak{R} \mid \exists j \in \mathbb{N}. F(r_{j,i}) \wedge \forall \ell \in \mathbb{N}. \neg F(r_{\ell,i+1})\}$$

and

$$\mathfrak{R}_\infty \stackrel{\text{def}}{=} \{\rho \in \mathfrak{R} \mid \forall i \in \mathbb{N} \exists j \in \mathbb{N}. F(r_{j,i})\}.$$

That is to say for all $i \in \mathbb{N}$, \mathfrak{R}_i is the set of runs in \mathfrak{R} that restart exactly i times and \mathfrak{R}_∞ is the set of runs in \mathfrak{R} that restart infinitely many times.

We go on to define the sets of runs $\mathfrak{R}_{\geq i} \stackrel{\text{def}}{=} \bigcup_{j=i}^\infty \mathfrak{R}_j$ which are those runs which restart at least i times. In particular note that $\mathfrak{R}_\infty = \bigcap_{i=0}^\infty \mathfrak{R}_{\geq i}$ and $\mathfrak{R}_{\geq i+1} \subseteq \mathfrak{R}_{\geq i}$.

By construction, any run $\rho \in \mathfrak{R}_\infty$ is losing since the negative reward that is collected upon restarting instantly brings the mean payoff below -1 by definition of m_n . Thus restarting infinitely many times translates directly into the mean payoff dropping below -1 infinitely many times and thus a strictly negative lim inf mean payoff. As a result it must be the case that $\mathfrak{R}_\infty \subseteq \neg MP_{\liminf \geq 0}$.

After every restart, the negative reward is reimbursed. Intuitively, going through finitely many restarts does not damage the chances of winning. We now show that, except for a nullset, the runs restarting only finitely many times satisfy the objective. Indeed, every run with only finitely many restarts must spend an infinite tail in some final gadget in which it does not restart. In this final gadget, the strategy plays just like $\sigma_{1/2}$, which means that it mimics the random choice in every controlled state. Since, by assumption, there are no more restarts, we obtain $\mathcal{P}_{\mathcal{M},s_0,\sigma}(\mathfrak{R}_i) = \mathcal{P}_{\mathcal{M},s_0,\sigma}(\mathfrak{R}_i \wedge \forall j \in \mathbb{N}. G(\neg r_{j,i+1}))$. We then apply Lemma 4.7 to obtain that

$$\mathcal{P}_{\mathcal{M},s_0,\sigma}(\mathfrak{R}_i) = \mathcal{P}_{\mathcal{M},s_0,\sigma}(\mathfrak{R}_i \wedge \forall j \in \mathbb{N}. G(\neg r_{j,i+1})) = \mathcal{P}_{\mathcal{M},s_0,\sigma}(\mathfrak{R}_i \wedge MP_{\liminf \geq 0}). \quad (4.3)$$

In other words, except for a nullset, the runs restarting finitely often (here i times) satisfy $MP_{\liminf \geq 0}$. Furthermore, notice that from this observation, the sets \mathfrak{R}_i partition the set of winning runs.

We show now that $\mathcal{P}_{\mathcal{M},s_0,\sigma}(\mathfrak{R}_\infty) = 0$. We do so firstly by showing by induction that $\mathcal{P}_{\mathcal{M},s_0,\sigma}(\mathfrak{R}_{\geq i}) \leq 2^{-i}$ for $i \geq 1$, then applying the continuity of measures from above to obtain that $\mathcal{P}_{\mathcal{M},s_0,\sigma}(\mathfrak{R}_\infty) = 0$.

Our base case is $i = 1$. \mathfrak{R} , by definition of σ , is the set of runs induced by playing $\sigma_{1/2}$ in every gadget. By Lemma 4.7, σ attains $\geq 1/2$ in every gadget. Therefore in particular the probability of a run leaving the first gadget is no more than $1/2$, i.e. $\mathcal{P}_{\mathcal{M},s_0,\sigma}(\mathfrak{R}_{\geq 1}) \leq 1/2$.

Now suppose that $\mathcal{P}_{\mathcal{M},s_0,\sigma}(\mathfrak{R}_{\geq i}) \leq 2^{-i}$. After restarting at least i times, the probability of a run restarting at least once more is still $\leq 1/2$ since the strategy being played in every gadget is $\sigma_{1/2}$. Hence

$$\mathcal{P}_{\mathcal{M},s_0,\sigma}(\mathfrak{R}_{\geq i+1}) \leq \mathcal{P}_{\mathcal{M},s_0,\sigma}(\mathfrak{R}_{\geq i}) \cdot \frac{1}{2} \leq 2^{-(i+1)}$$

which is what we wanted.

Now we use the fact that $\mathfrak{R}_\infty = \bigcap_{i=0}^\infty \mathfrak{R}_{\geq i}$ and $\mathfrak{R}_{\geq i+1} \subseteq \mathfrak{R}_{\geq i}$ to apply continuity of measures from above and obtain:

$$\mathcal{P}_{\mathcal{M},s_0,\sigma}(\mathfrak{R}_\infty) = \mathcal{P}_{\mathcal{M},s_0,\sigma}\left(\bigcap_{i=0}^\infty \mathfrak{R}_{\geq i}\right) = \lim_{i \rightarrow \infty} \mathcal{P}_{\mathcal{M},s_0,\sigma}(\mathfrak{R}_{\geq i}) \leq \lim_{i \rightarrow \infty} 2^{-i} = 0.$$

Hence \mathfrak{R}_∞ is a null set.

We can now write down the following:

$$\begin{aligned} 1 &= \mathcal{P}_{\mathcal{M},s_0,\sigma}(\mathfrak{R}) \\ &= \left(\sum_{i=0}^\infty \mathcal{P}_{\mathcal{M},s_0,\sigma}(\mathfrak{R}_i)\right) + \mathcal{P}_{\mathcal{M},s_0,\sigma}(\mathfrak{R}_\infty) && \text{by partition of } \mathfrak{R} \\ &= \left(\sum_{i=0}^\infty \mathcal{P}_{\mathcal{M},s_0,\sigma}(\mathfrak{R}_i \wedge MP_{\liminf \geq 0})\right) + \mathcal{P}_{\mathcal{M},s_0,\sigma}(\mathfrak{R}_\infty) && \text{by Equation (4.3)} \\ &= \left(\sum_{i=0}^\infty \mathcal{P}_{\mathcal{M},s_0,\sigma}(\mathfrak{R}_i \wedge MP_{\liminf \geq 0})\right) \\ &\quad + \mathcal{P}_{\mathcal{M},s_0,\sigma}(\mathfrak{R}_\infty \wedge MP_{\liminf \geq 0}) && \text{by } \mathcal{P}_{\mathcal{M},s_0,\sigma}(\mathfrak{R}_\infty) = 0 \\ &= \mathcal{P}_{\mathcal{M},s_0,\sigma}(MP_{\liminf \geq 0}) && \text{by partition of } MP_{\liminf \geq 0} \end{aligned}$$

Thus $\mathcal{P}_{\mathcal{M},s_0,\sigma}(\mathfrak{R}) = \mathcal{P}_{\mathcal{M},s_0,\sigma}(MP_{\liminf \geq 0}) = 1$, i.e. σ wins almost surely. \square

Lemma 4.19. *For any FR strategy σ , $\mathcal{P}_{\mathcal{M},\sigma,s_0}(MP_{\liminf \geq 0}) = 0$.*

Outline of the proof. Let σ be any FR strategy. We partition the runs generated by σ into runs restarting infinitely often, and those restarting only finitely many times. Any runs restarting infinitely often are losing by construction. The runs restarting only finitely many times spend an infinite tail in a given gadget, letting the mean payoff dip below -1 infinitely many times with probability 1 by Lemma 4.10. Hence we have that $\mathcal{P}_{\mathcal{M},\sigma,s_0}(MP_{\liminf \geq 0}) = 0$. \square

Full proof. There are two ways to lose when playing in this MDP: either the mean payoff dips below -1 infinitely often because the run takes infinitely many restarts, or the run only takes finitely many restarts, but the mean payoff drops below -1 infinitely many times in the last copy of the gadget that the run stays in. Recall that in Lemma 4.10 we showed that any FR strategy with probability 1 either restarts or lets the mean payoff dip below -1 infinitely often.

Let σ be any FR strategy and let \mathfrak{R} to be the set of runs induced by σ from s_0 . We partition \mathfrak{R} into the sets \mathfrak{R}_i and \mathfrak{R}_∞ of runs such that $\mathfrak{R} = (\bigcup_{i=0}^\infty \mathfrak{R}_i) \cup \mathfrak{R}_\infty$. Where we define for $i = 0$

$$\mathfrak{R}_0 \stackrel{\text{def}}{=} \{\rho \in \mathfrak{R} \mid \forall \ell \in \mathbb{N}, \neg F(r_{\ell,1})\},$$

for $i \geq 1$

$$\mathfrak{R}_i \stackrel{\text{def}}{=} \{\rho \in \mathfrak{R} \mid \exists j \in \mathbb{N}, F(r_{j,i}) \wedge \forall \ell \in \mathbb{N}, \neg F(r_{\ell,i+1})\}$$

and

$$\mathfrak{R}_\infty \stackrel{\text{def}}{=} \{\rho \in \mathfrak{R} \mid \forall i, \exists j F(r_{j,i})\}.$$

That is to say for all $i \in \mathbb{N}$, \mathfrak{R}_i is the set of runs in \mathfrak{R} that restart exactly i times and \mathfrak{R}_∞ is the set of runs in \mathfrak{R} that restart infinitely many times.

We go on to define the sets of runs $\mathfrak{R}_{\geq i} \stackrel{\text{def}}{=} \bigcup_{j=i}^{\infty} \mathfrak{R}_j$ which are those runs which restart at least i times. In particular note that $\mathfrak{R}_{\infty} = \bigcap_{i=0}^{\infty} \mathfrak{R}_{\geq i}$ and $\mathfrak{R}_{\geq i+1} \subseteq \mathfrak{R}_{\geq i}$.

Note that any run in \mathfrak{R}_{∞} is losing by construction. The negative reward that is collected upon restarting instantly brings the mean payoff below -1 by definition of m_n . Thus restarting infinitely many times translates directly into the mean payoff dropping below -1 infinitely many times. Thus $\mathfrak{R}_{\infty} \subseteq \neg MP_{\liminf \geq 0}$ and so it follows that $\mathcal{P}_{\mathcal{M}, s_0, \sigma}(\mathfrak{R}_{\infty}) = \mathcal{P}_{\mathcal{M}, s_0, \sigma}(\mathfrak{R}_{\infty} \wedge \neg MP_{\liminf \geq 0})$. Since the sets \mathfrak{R}_i and \mathfrak{R}_{∞} partition \mathfrak{R} we have that:

$$\mathcal{P}_{\mathcal{M}, s_0, \sigma}(\mathfrak{R}) = \left(\sum_{i=0}^{\infty} \mathcal{P}_{\mathcal{M}, s_0, \sigma}(\mathfrak{R}_i) \right) + \mathcal{P}_{\mathcal{M}, s_0, \sigma}(\mathfrak{R}_{\infty}).$$

It remains to show that every set \mathfrak{R}_i is almost surely losing, i.e. $\mathcal{P}_{\mathcal{M}, s_0, \sigma}(\mathfrak{R}_i) = \mathcal{P}_{\mathcal{M}, s_0, \sigma}(\mathfrak{R}_i \wedge \neg MP_{\liminf \geq 0})$. Consider a run $\rho \in \mathfrak{R}_i$. By definition it restarts exactly i times. As a result, it spends infinitely long in the $i+1$ st gadget. Because σ is an FR strategy, it must be the case that any substrategy σ^* induced by σ that is played in a given gadget is also an FR strategy. This allows us to apply Lemma 4.10 to obtain that

$$\mathcal{P}_{\mathcal{M}, s_0, \sigma}(\mathfrak{R}_i) = \mathcal{P}_{\mathcal{M}, s_0, \sigma}(\mathfrak{R}_i \wedge (\neg MP_{\liminf \geq 0} \vee \exists j \in \mathbb{N}, F(r_{j, i+1}))). \quad (4.4)$$

However, any run $\rho \in \mathfrak{R}_i$ never sees any state $r_{j, i+1}$ for any j by definition. Therefore it follows that

$$\mathcal{P}_{\mathcal{M}, s_0, \sigma}(\mathfrak{R}_i \wedge (\neg MP_{\liminf \geq 0} \vee \exists j \in \mathbb{N}, F(r_{j, i+1}))) = \mathcal{P}_{\mathcal{M}, s_0, \sigma}(\mathfrak{R}_i \wedge (\neg MP_{\liminf \geq 0}))$$

Hence $\mathcal{P}_{\mathcal{M}, s_0, \sigma}(\mathfrak{R}_i) = \mathcal{P}_{\mathcal{M}, s_0, \sigma}(\mathfrak{R}_i \wedge \neg MP_{\liminf \geq 0})$ as required.

As a result we have that

$$\begin{aligned} 1 &= \mathcal{P}_{\mathcal{M}, s_0, \sigma}(\mathfrak{R}) \\ &= \left(\sum_{i=0}^{\infty} \mathcal{P}_{\mathcal{M}, s_0, \sigma}(\mathfrak{R}_i) \right) + \mathcal{P}_{\mathcal{M}, s_0, \sigma}(\mathfrak{R}_{\infty}) && \text{by partition of } \mathfrak{R} \\ &= \left(\sum_{i=0}^{\infty} \mathcal{P}_{\mathcal{M}, s_0, \sigma}(\mathfrak{R}_i \wedge \neg MP_{\liminf \geq 0}) \right) + \mathcal{P}_{\mathcal{M}, s_0, \sigma}(\mathfrak{R}_{\infty} \wedge \neg MP_{\liminf \geq 0}) && \text{by Equation (4.4)} \\ &= \mathcal{P}_{\mathcal{M}, s_0, \sigma}(\neg MP_{\liminf \geq 0}) && \text{by partition of } \mathfrak{R} \end{aligned}$$

That is to say that for any FR strategy σ , $\mathcal{P}_{\mathcal{M}, s_0, \sigma}(MP_{\liminf \geq 0}) = 0$. \square

From Lemma 4.18 and Lemma 4.19 we obtain the following theorem.

Theorem 4.20. *There exists a countable, finitely branching and acyclic MDP \mathcal{M} whose step counter is implicit in the state for which s_0 is almost surely winning $MP_{\liminf \geq 0}$, i.e., $\exists \delta \mathcal{P}_{\mathcal{M}, s_0, \delta}(MP_{\liminf \geq 0}) = 1$, but every FR strategy σ is such that $\mathcal{P}_{\mathcal{M}, s_0, \sigma}(MP_{\liminf \geq 0}) = 0$. In particular, almost sure winning strategies, when they exist, cannot be chosen k -bit Markov for any $k \in \mathbb{N}$ for countable MDPs.*

Proof. Proved by Lemma 4.18 and Lemma 4.19. \square

Now we construct the MDP $\mathcal{M}_{\text{Restart}}$ by using Figure 7, but we substitute the instances of Figure 3 gadgets with instances of Figure 6 gadgets. This allows us to obtain the following results which state that optimal strategies for $MP_{\liminf \geq 0}$ requires infinite memory, even when the reward counter is implicit in the state.

Lemma 4.21. *There exists an HD strategy σ such that $\mathcal{P}_{\mathcal{M}_{\text{Restart}}, s_0, \sigma}(MP_{\liminf \geq 0}) = 1$.*

Proof. The proof is identical to that of Lemma 4.18. \square

Lemma 4.22. *For any FR strategy σ , $\mathcal{P}_{\mathcal{M}_{\text{Restart}}, s_0, \sigma}(MP_{\liminf \geq 0}) = 0$.*

Proof. The proof is identical to that of Lemma 4.19. \square

Theorem 4.23. *There exists a countable, finitely branching and acyclic MDP $\mathcal{M}_{\text{Restart}}$ whose total reward is implicit in the state where, for the initial state s_0 ,*

- *there exists an HD strategy σ s.t. $\mathcal{P}_{\mathcal{M}_{\text{Restart}}, s_0, \sigma}(MP_{\liminf \geq 0}) = 1$.*
- *for every FR strategy σ , $\mathcal{P}_{\mathcal{M}_{\text{Restart}}, s_0, \sigma}(MP_{\liminf \geq 0}) = 0$.*

Proof. This follows from Lemma 4.21 and Lemma 4.22. \square

5. TOTAL PAYOFF

Total Payoff		ε -optimal	Optimal
Finitely branching	Upper Bound	Det(RC) 5.1	Det(RC) 5.4
	Lower Bound	\neg Rand(F+SC) 5.7	\neg Rand(F+SC) 5.8
Infinitely branching	Upper Bound	Det(SC+RC) 5.2	Det(SC+RC) 5.5
	Lower Bond	\neg Rand(F+SC) 5.7 and \neg Rand(F+RC) 5.6	\neg Rand(F+SC) 5.8 and \neg Rand(F+RC) 5.6

Table 4: Strategy complexity of ε -optimal/optimal strategies for the total payoff objective in infinitely/finitely branching MDPs.

5.1. Upper Bounds. In order to tackle the upper bounds for the total payoff objective $TP_{\liminf \geq 0}$, we work with the derived MDPs $R(\mathcal{M})$ and $S(\mathcal{M})$ which encode the total reward and the step counter into the state respectively. Once the total reward is encoded into the state, the point payoff coincides with the total payoff. We use this observation to reduce $TP_{\liminf \geq 0}$ to $PP_{\liminf \geq 0}$ and obtain our upper bounds from the corresponding point payoff results.

Corollary 5.1. *Given a finitely branching MDP \mathcal{M} , there exist ε -optimal strategies for $TP_{\liminf \geq 0}$ which use just a reward counter.*

Proof. We place ourselves in $R(\mathcal{M})$ where $TP_{\liminf \geq 0}$ and $PP_{\liminf \geq 0}$ coincide. Thus we can apply Theorem 3.5 to obtain ε -optimal MD strategies for $TP_{\liminf \geq 0}$ from every state of $R(\mathcal{M})$. By Lemma 2.4 we can translate these MD strategies on $R(\mathcal{M})$ back to strategies on \mathcal{M} with just a reward counter. \square

Corollary 5.2. *Given an MDP \mathcal{M} with initial state s_0 ,*

- *there exist ε -optimal MD strategies for $TP_{\liminf \geq 0}$ in $S(R(\mathcal{M}))$,*
- *there exist ε -optimal strategies for $TP_{\liminf \geq 0}$ which use a step counter and a reward counter.*

Proof. We consider the encoded system $R(\mathcal{M})$ in which the reward counter is implicit in the state. Recall that total rewards in \mathcal{M} correspond exactly to point rewards in $R(\mathcal{M})$. We then apply Theorem 3.10 to $R(\mathcal{M})$ to obtain ε -optimal MD strategies for $PP_{\liminf \geq 0}$ in $S(R(\mathcal{M}))$. Lemma 2.2 allows us to translate these MD strategies back to $R(\mathcal{M})$ with a memory overhead of just a step counter. Then we apply Lemma 2.4 to translate these Markov strategies back to \mathcal{M} with a memory overhead of just a reward counter. Hence ε -optimal strategies for $TP_{\liminf \geq 0}$ in \mathcal{M} just use a step counter and a reward counter as required. \square

Remark 5.3. While ε -optimal strategies for mean payoff and total payoff (in infinitely branching MDPs) have the same memory requirements, the step counter and the reward counter do not arise in the same way. Both the step counter and reward counter used in ε -optimal strategies for mean payoff arise from the construction of $A(\mathcal{M})$. However, in the case for total payoff, only the reward counter arises from the construction of $R(\mathcal{M})$. The step counter on the other hand arises from the Markov strategy needed for point payoff in $R(\mathcal{M})$.

Corollary 5.4. *Given a finitely branching MDP \mathcal{M} and initial state s_0 , optimal strategies, where they exist, can be chosen with just a reward counter for $TP_{\liminf \geq 0}$.*

Proof. We place ourselves in $R(\mathcal{M})$ where $TP_{\liminf \geq 0}$ is shift invariant. Moreover, in $R(\mathcal{M})$ the objectives $TP_{\liminf \geq 0}$ and $PP_{\liminf \geq 0}$ coincide. Thus we can apply Theorem 3.5 to obtain ε -optimal MD strategies for $TP_{\liminf \geq 0}$ from every state of $R(\mathcal{M})$. From Theorem 3.2 we obtain a single MD strategy that is optimal from every state of $R(\mathcal{M})$ that has an optimal strategy. By Lemma 2.4 we can translate this MD strategy on $R(\mathcal{M})$ back to a strategy on \mathcal{M} with just a reward counter. \square

Corollary 5.5. *Given an MDP \mathcal{M} and initial state s_0 , optimal strategies, where they exist, can be chosen with just a reward counter and a step counter for $TP_{\liminf \geq 0}$.*

Proof. We place ourselves in $S(R(\mathcal{M}))$ and apply Corollary 5.2 to obtain ε -optimal MD strategies for $TP_{\liminf \geq 0}$ from every state of $S(R(\mathcal{M}))$. While $TP_{\liminf \geq 0}$ is not shift invariant in \mathcal{M} , it is shift invariant in $S(R(\mathcal{M}))$, and thus we can apply Theorem 3.2 to obtain a single MD strategy that is optimal from every state of $S(R(\mathcal{M}))$ that has an optimal strategy. The result then follows from Lemma 2.2 and Lemma 2.4. \square

5.2. Lower Bounds.

5.2.1. ε -optimal strategies.

Theorem 5.6. *There exists an infinitely branching MDP \mathcal{M} as in Figure 2 with reward implicit in the state and initial state s such that*

- every FR strategy σ is such that $\mathcal{P}_{\mathcal{M},s,\sigma}(TP_{\liminf \geq 0}) = 0$
- there exists an HD strategy σ such that $\mathcal{P}_{\mathcal{M},s,\sigma}(TP_{\liminf \geq 0}) = 1$.

Hence, optimal (and even almost-surely winning) strategies and ε -optimal strategies for $TP_{\liminf \geq 0}$ require infinite memory beyond a reward counter.

Proof. This follows directly from [23, Theorem 4] and the observation that in Figure 2, $TP_{\liminf \geq 0}$, and co-Büchi objectives coincide. \square

The statements and proofs of Lemma 4.7 and Lemma 4.10 also hold for $TP_{\liminf \geq 0}$, giving us the following theorem.

Theorem 5.7. *There exists a countable, finitely branching and acyclic MDP \mathcal{M} whose step counter is implicit in the state for which $\text{val}_{\mathcal{M}, TP_{\liminf \geq 0}}(s_0) = 1$ and any FR strategy σ is such that $\mathcal{P}_{\mathcal{M}, s_0, \sigma}(TP_{\liminf \geq 0}) = 0$. In particular, there are no ε -optimal k -bit Markov strategies for any $k \in \mathbb{N}$ and any $\varepsilon < 1$ for $TP_{\liminf \geq 0}$ in countable MDPs.*

5.2.2. *Optimal strategies.* The statements and proofs of Lemma 4.18 and Lemma 4.19 also hold for $TP_{\liminf \geq 0}$, giving us the following theorem.

Theorem 5.8. *There exists a countable, finitely branching and acyclic MDP \mathcal{M} whose step counter is implicit in the state for which s_0 is almost surely winning $TP_{\liminf \geq 0}$, i.e., $\exists \hat{\sigma} \mathcal{P}_{\mathcal{M}, s_0, \hat{\sigma}}(TP_{\liminf \geq 0}) = 1$, but every FR strategy σ is such that $\mathcal{P}_{\mathcal{M}, s_0, \sigma}(TP_{\liminf \geq 0}) = 0$. In particular, almost sure winning strategies, when they exist, cannot be chosen k -bit Markov for any $k \in \mathbb{N}$ for countable MDPs.*

Proof. Proved by Lemma 4.18 and Lemma 4.19. □

6. STRENGTHENING RESULTS

The counterexamples we present in Section 4 feature finite but unbounded branching degree, unbounded rewards and irrational transition probabilities. In this section we show that the hardness does not depend on these aspects by strengthening the counterexamples to have binary branching, bounded rewards and rational transition probabilities.

Consider a new MDP \mathcal{M} based on the MDP constructed in Figure 3 which now undergoes the following changes. First we bound the branching degree by 2. We do so by replacing the outgoing transitions in states s_n and c_n of each gadget by binary trees with accordingly adjusted probabilities such that there is still a probability of $\delta_i(n)$ of receiving reward $i \cdot m_n$ in each gadget for $i \in \{0, 1, \dots, k(n)\}$.

To adjust for the increased path lengths incurred by the modifications to each gadget, the construction in Figure 4 is accordingly modified by padding each vertical column of white states with extra transitions based on the number of transitions present in the matching gadget. As a result, path length is preserved even when skipping gadgets. The construction in Figure 7 is similarly modified.

Second, we restrict the transition probabilities to rationals. The transition probabilities $\delta_i(n)$ and $\varepsilon_i(n)$ are replaced by close rationals in $(\delta_i(n), \delta_i(n) + 2^{-n})$ and in $(\varepsilon_i(n), \varepsilon_i(n) + 2^{-n})$, respectively. These rationals are constructible, for example by approximation of $\delta_i(n)$ and $\varepsilon_i(n)$ themselves. Since these new rational probabilities are so close to the original ones, all of the relevant convergence and divergence of series is preserved.

Definition 6.1 (Binary branching). We formally define how to modify the MDPs in Section 4.2 such that they have a branching degree of no more than 2.

In each gadget in Figure 3 and Figure 6, we do as is shown in Figure 8. I.e. the outgoing transitions from s_n and c_n are replaced by a binary tree of depth at most $\lceil \lg(k(n) + 1) \rceil$. Because c_n is player controlled, we do not need to define any new transition probabilities. For the outgoing transitions from s_n , new probabilities must be defined such that the probability

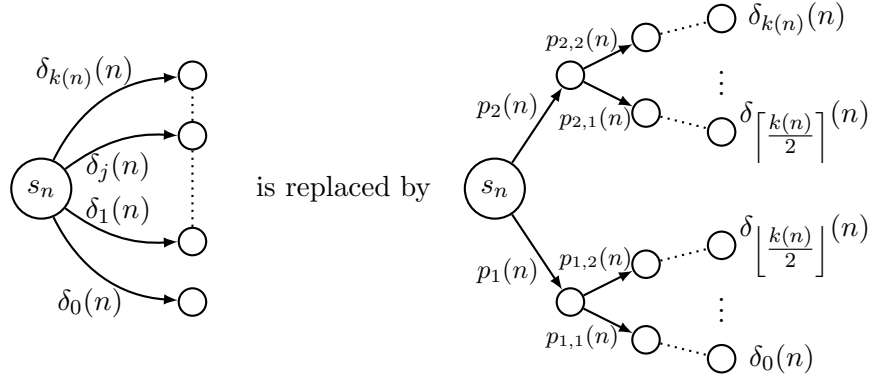


Figure 8: Schema for replacing arbitrary finite branching with binary branching in Figure 3.

of receiving reward $im(n)$ is still $\delta_i(n)$. For illustrative purposes, the transition probabilities for the initial branching are as follows.

$$p_1(n) \stackrel{\text{def}}{=} \sum_{i=0}^{\lfloor \frac{k(n)}{2} \rfloor} \delta_i(n) \text{ and } p_2(n) \stackrel{\text{def}}{=} \sum_{i=\lceil \frac{k(n)}{2} \rceil}^{k(n)} \delta_i(n)$$

All other transition probabilities are obtained inductively.

The extended path lengths are mirrored in a modified Figure 4 by padding the path lengths by $4 + 2\lceil \lg(k(n)) \rceil$ steps instead of 4 steps. Similarly, in Figure 7 we pad the length of the chains of black states by an extra $2\lceil \lg(k(n)) \rceil$ steps. So the step counter is still implicit in the state.

Definition 6.2 (Rational probabilities). We define the probabilities $\gamma_i(n)$ and $\theta_i(n)$ as follows. We set $\gamma_i(n) \in (\delta_i(n), \delta_i(n) + 2^{-n}) \cap \mathbb{Q}$ and similarly set $\theta_i(n) \in (\varepsilon_i(n), \varepsilon_i(n) + 2^{-n}) \cap \mathbb{Q}$. Furthermore, we define new functions g^* , h^* and k^* as follows. We set

$$g^*(i) \stackrel{\text{def}}{=} \min \left\{ N : \left(\sum_{n>N} \gamma_{i-1}(n) \theta_{i-1}(n) \right) \leq 2^{-i} \right\}, \quad h^*(1) \stackrel{\text{def}}{=} 2$$

$$h^*(i+1) \stackrel{\text{def}}{=} \left\lceil \max \left\{ g^*(i+1), \text{Tower}(i+2), \min \left\{ m+1 \in \mathbb{N} : \sum_{n=h^*(i)}^m \theta_{i-1}(n) \geq 1 \right\} \right\} \right\rceil.$$

which yield $k^*(i) \stackrel{\text{def}}{=} h^{*-1}(n)$.

Note that $g^*(i)$ is well defined since

$$\begin{aligned} \sum_{n>N} \gamma_{i-1}(n) \theta_{i-1}(n) &< \sum_{n>N} \delta_{i-1}(n) \varepsilon_{i-1}(n) + 2^{-n} (\delta_{i-1}(n) + \varepsilon_{i-1}(n) + 2^{-n}) \\ &< \sum_{n>N} \delta_{i-1}(n) \varepsilon_{i-1}(n) + 2^{-n} \cdot 3 \end{aligned}$$

is convergent for all i . Similarly, $h^*(i)$ is also well defined since

$$\sum_{n=h^*(i)}^{\infty} \theta_{i-1}(n) > \sum_{n=h^*(i)}^{\infty} \varepsilon_{i-1}(n)$$

which diverges for all i .

Remark 6.3. We now make sure that the results from Section 4 still hold when we change Figure 3 by replacing the transition probabilities $\delta_i(n)$ and $\varepsilon_i(n)$ with $\gamma_i(n)$ and $\theta_i(n)$, respectively. To this end we must check that some crucial results still hold. Namely, that Lemma 4.6 and Lemma 4.9 still hold given the modified transition probabilities.

In order to show that Lemma 4.6 still holds, we must show that $\sum_{i=0}^{k(n)-1} \gamma_i(n) < 1$. I.e. we get

$$\begin{aligned} \sum_{i=0}^{k(n)-1} \gamma_i(n) &\leq \sum_{i=0}^{k(n)-1} \delta_i(n) + \sum_{i=0}^{k(n)-1} 2^{-n} \\ &= \sum_{i=0}^{k(n)-1} \delta_i(n) + (k(n) - 1)2^{-n} \\ &< \sum_{i=1}^{k(n)} 2^{-i} + 2^{k(n)-1} \cdot 2^{-n} \\ &\leq \sum_{i=1}^{k(n)} 2^{-i} + 2^{-\frac{n}{2}} < 1 \end{aligned} \quad \text{since } k(n) < \frac{n}{2}$$

That is to say that the transition probabilities are indeed well defined using rational probabilities $\gamma_i(n)$ in lieu of $\delta_i(n)$.

Similarly, we must now show that the following sum diverges.

$$\sum_{n=k^{-1}(2)}^{\infty} \left(\gamma_{j(n)}(n)(\alpha_n \theta_{j(n)}(n) + (1 - \alpha_n) \theta_{i(n)}(n)) + \gamma_{i(n)}(n)(\alpha_n + (1 - \alpha_n) \theta_{i(n)}(n)) \right)$$

We do so by noticing that

$$\begin{aligned} &\sum_{n=k^{-1}(2)}^{\infty} \left(\gamma_{j(n)}(n)(\alpha_n \theta_{j(n)}(n) + (1 - \alpha_n) \theta_{i(n)}(n)) + \gamma_{i(n)}(n)(\alpha_n + (1 - \alpha_n) \theta_{i(n)}(n)) \right) \\ &\geq \sum_{n=k^{-1}(2)}^{\infty} \left(\delta_{j(n)}(n)(\alpha_n \varepsilon_{j(n)}(n) + (1 - \alpha_n) \varepsilon_{i(n)}(n)) + \delta_{i(n)}(n)(\alpha_n + (1 - \alpha_n) \varepsilon_{i(n)}(n)) \right) \end{aligned}$$

since $\gamma_i(n) \geq \delta_i(n)$ and $\theta_i(n) \geq \varepsilon_i(n)$ for all i .

Hence Lemma 4.9 yields the desired divergence result.

Putting both of the above results together, we can obtain rational probability versions of Lemma 4.7 and Lemma 4.10.

Combining these constructions allows us to obtain the following properties.

Theorem 6.4. *There exists a countable, acyclic MDP \mathcal{M} , whose step counter is implicit in the state, whose transition probabilities are rational and whose branching degree is bounded by 2*

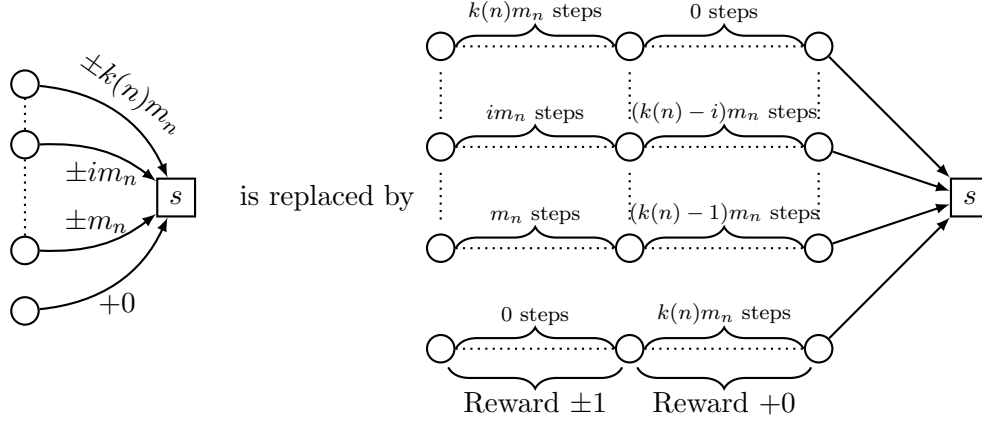


Figure 9: Schema for replacing large rewards with bounded rewards in Figure 3.

for which $\text{val}_{\mathcal{M}, MP_{\liminf \geq 0}}(s_0) = 1$ and any FR strategy σ is such that $\mathcal{P}_{\mathcal{M}, s_0, \sigma}(MP_{\liminf \geq 0}) = 0$. In particular, there are no ε -optimal step counter plus finite memory strategies for any $\varepsilon < 1$ for the $MP_{\liminf \geq 0}$ objective for countable MDPs.

Proof. This follows from Lemma 4.7, Lemma 4.10 by modifying the constructions in Figure 3 and Figure 4 as detailed in Definition 6.1 and Definition 6.2. \square

Theorem 6.5. *There exists a countable, acyclic MDP \mathcal{M} , whose step counter is implicit in the state, whose transition probabilities are rational and whose branching degree is bounded by 2 for which s_0 is almost surely winning for $MP_{\liminf \geq 0}$ and any FR strategy σ is such that $\mathcal{P}_{\mathcal{M}, s_0, \sigma}(MP_{\liminf \geq 0}) = 0$. In particular, almost sure winning strategies, when they exist, cannot be chosen with a step counter plus finite memory for countable MDPs.*

Proof. This follows from Lemma 4.18, Lemma 4.19 by modifying the constructions in Figure 3 and Figure 7 as detailed in Definition 6.1 and Definition 6.2. \square

We now further alter Figure 3 by bounding the rewards. The rewards on transitions are now limited to -1 , 0 or 1 . To compensate for the smaller rewards, in the n -th gadget, each transition bearing a reward is replaced by $k(n) \cdot m_n$ transitions as follows. If the original transition had reward $j \cdot m_n$ then that transition is replaced with $j \cdot m_n$ transitions with reward 1 , and $(k(n) - j) \cdot m_n$ transitions with reward 0 . Symmetrically all negatively weighted transitions are similarly replaced by transitions with rewards -1 and 0 . The extra padding with the transitions with reward 0 is done in order to preserve the path length, i.e., such that the step counter is still implicit in the state.

Definition 6.6 (Bounded rewards). We formally define how to modify the MDPs in Section 4.2 such that their rewards are either -1 , $+1$ or 0 . In each Figure 3 gadget, we do as illustrated in Figure 8. I.e. the incoming transitions to c_n and s_n carry rewards $\pm im_n$ for some i with $0 \leq i \leq k(n)$. We then split this transition carrying reward $\pm im_n$ into a chain of $k(n)m_n$ transitions. The first im_n of which carry reward ± 1 , and the last $(k(n) - i)m_n$ of which carry reward 0 .

The extended path lengths are mirrored in a modified Figure 4 by padding the path lengths by an extra $2k(n)m_n$ steps. Because skipping ahead in Figure 4 reimburses reward $+i$ upon entering state s_{N^*+i+1} , we replace these transitions with i transitions bearing reward

+1 and reflect this increased path length by padding the incoming transitions to s_{N^*+i+1} with an extra i transitions bearing reward 0.

The extended path lengths must also be reflected in Figure 7. This is done by replacing transitions carrying reward $\pm m_i$ by m_i transitions carrying reward ± 1 . We also increase the number of black states from 3 to $2k(n)m_n + 3$ to match the number of steps taken inside the n th gadget.

Theorem 6.7. *There exists a countable, acyclic MDP \mathcal{M} , whose step counter is implicit in the state, whose transition probabilities are rational, whose rewards on transitions are in $\{-1, 0, 1\}$ and whose branching degree is bounded by 2 for which $\text{val}_{\mathcal{M}, TP_{\liminf \geq 0}}(s_0) = 1$ and any FR strategy σ is such that $\mathcal{P}_{\mathcal{M}, s_0, \sigma}(TP_{\liminf \geq 0}) = 0$. In particular, there are no ε -optimal step counter plus finite memory strategies for any $\varepsilon < 1$ for the $TP_{\liminf \geq 0}$ objective for countable MDPs.*

Proof. This follows from Lemma 4.7, Lemma 4.10 by modifying the constructions in Figure 3 and Figure 4 as detailed in Definition 6.1, Definition 6.6 and Definition 6.2. \square

Theorem 6.8. *There exists a countable, acyclic MDP \mathcal{M} , whose step counter is implicit in the state, whose transition probabilities are rational, whose rewards on transitions are in $\{-1, 0, 1\}$ and whose branching degree is bounded by 2 for which s_0 is almost surely winning for $TP_{\liminf \geq 0}$ and any FR strategy σ is such that $\mathcal{P}_{\mathcal{M}, s_0, \sigma}(TP_{\liminf \geq 0}) = 0$. In particular, almost sure winning strategies, when they exist, cannot be chosen with a step counter plus finite memory for countable MDPs.*

Proof. This follows from Lemma 4.18, Lemma 4.19 by modifying the constructions in Figure 3 and Figure 7 as detailed in Definition 6.1, Definition 6.6 and Definition 6.2. \square

Remark 6.9. We draw attention to the fact that we could state the total payoff theorems using bounded rewards, but we did not do so for the equivalent mean payoff results. In the case of mean payoff, with the step counter implicit in the state, having bounded transition rewards, e.g. bounded by $\pm b$, means that the average reward in any given state will always be bounded by $\pm b$. In the context of our example in Figure 3, this means that if we used the construction in Definition 6.6, the absolute worst the average reward can be is ~ -1 , this can only happen going from the 0th random transition to the $k(n)$ th choice. But even worse, using only one bit of memory to remember whether the random transition was $> \frac{k(n)}{2}$ or not, the mean payoff is suddenly at worst $\sim -\frac{1}{k(n)}$ which converges to 0.

In general it would be interesting to consider the mean payoff objective with step counter encoded and bounded rewards since our results do not obviously carry over to this case.

Some extra care is needed to convince ourselves that Theorem 4.15 and Theorem 4.23 can also be strengthened. Consider the construction in Figure 6. In the random choice, the transition rewards are already all 0, so only the branching degree needs to be adjusted by padding the choice with a binary tree as above. In the controlled choice, the transitions carrying reward $\pm m_n^i$ are replaced by m_n^i transitions each bearing reward ± 1 respectively. Therefore, the path lengths increase in the following way in the n -th gadget. In s_n and c_n , the binary trees increase path length by up to $\lceil \lg(k(n) + 1) \rceil$ (where \lg is the logarithm to base 2) and after c_n the path length increases by up to $m_n^{k(n)}$ twice.

Consider the scenario where the play took the i -th random choice and the player makes the ‘best’ mistake where they choose transition $i + 1$. We show that, even in this best error

case (and thus in all other error cases), the newly added path lengths do still not help to prevent seeing a mean payoff $\leq -1/2$ in the n -th gadget. In this case, in the state between c_n and s_{n+1} , the total payoff is $-m_n^{i+1}$ and the total number of steps taken by the play so far is upper bounded by

$$\beta_n \stackrel{\text{def}}{=} \left(\sum_{i=N^*}^{n-1} 2\lceil \lg(k(i) + 1) \rceil + 2m_i^{k(i)} \right) + 2\lceil \lg(k(n) + 1) \rceil + m_n^i + m_n^{i+1}.$$

Recall that $m_n \stackrel{\text{def}}{=} \sum_{i=N^*}^{n-1} m_i^{k(n)}$ with $m_{N^*} \stackrel{\text{def}}{=} 1$, and this is the definition of m_n from Figure 6 which is different from the definition of m_n in Figure 3. Note that $k(n)$ is very slowly growing, so it follows that

$$\beta_n \leq 3m_n + m_n^i + m_n^{i+1} \leq 2m_n^{i+1}.$$

That is to say that the mean payoff is $\leq \frac{-m_n^{i+1}}{2m_n^{i+1}} = -1/2$. As a result, in the case of a bad aggressive decision, the mean payoff will still drop below $-1/2$ in this modified MDP (instead of dropping below -1 in the original MDP). This is just as good to falsify $MP_{\liminf \geq 0}$.

Thus we obtain the following two results.

Theorem 6.10. *There exists a countable, acyclic MDP \mathcal{M} , whose reward counter is implicit in the state, whose transition probabilities are rational, whose rewards on transitions are in $\{-1, 0, 1\}$ and whose branching degree is bounded by 2 for which $\text{val}_{\mathcal{M}, MP_{\liminf \geq 0}}(s_0) = 1$ and any FR strategy σ is such that $\mathcal{P}_{\mathcal{M}, s_0, \sigma}(MP_{\liminf \geq 0}) = 0$. In particular, there are no ε -optimal step counter plus finite memory strategies for any $\varepsilon < 1$ for the $MP_{\liminf \geq 0}$ objective for countable MDPs.*

Proof. This follows from Lemma 4.13, Lemma 4.14, Definition 6.1, Definition 6.6 and Definition 6.2. \square

Theorem 6.11. *There exists a countable, acyclic MDP \mathcal{M} , whose reward counter is implicit in the state, whose transition probabilities are rational, whose rewards on transitions are in $\{-1, 0, 1\}$ and whose branching degree is bounded by 2 for which s_0 is almost surely winning for $MP_{\liminf \geq 0}$ and any FR strategy σ is such that $\mathcal{P}_{\mathcal{M}, s_0, \sigma}(MP_{\liminf \geq 0}) = 0$. In particular, almost sure winning strategies, when they exist, cannot be chosen with a step counter plus finite memory for countable MDPs.*

Proof. This follows from Lemma 4.21, Lemma 4.22, Definition 6.1, Definition 6.6 and Definition 6.2. \square

Remark 6.12. The result from Lemma 4.10 holds even for strategies σ whose memory grows unboundedly, but slower than $k(n) - 1$. That is to say that there exists a countable, acyclic MDP \mathcal{M} , whose step counter is implicit in the state such that $\text{val}_{\mathcal{M}, MP_{\liminf \geq 0}}(s_0) = 1$ and any strategy σ with number of memory modes $< k(n) - 1$ in the n th gadget is such that $\mathcal{P}_{\mathcal{M}, \sigma, s_0}(MP_{\liminf \geq 0}) = 0$. This follows from a slightly modified version of Lemma 4.9 which considers the situation where states $i(n)$ and $k(n)$ are confused in the player's memory. Then the argument used in Lemma 4.10 can be modified to include $i(n), j(n) : \mathbb{N} \rightarrow \{0, 1, \dots, k(n)\}$. The result then follows since in every gadget at least one memory mode will confuse at least two states $i(n), j(n) : \mathbb{N} \rightarrow \{0, 1, \dots, k(n) - 1\}$, which as we have shown is enough to falsify $MP_{\liminf \geq 0}$.

This is in contrast to examples such as Figure 10 where the only requirement on the memory is that it grow unboundedly.

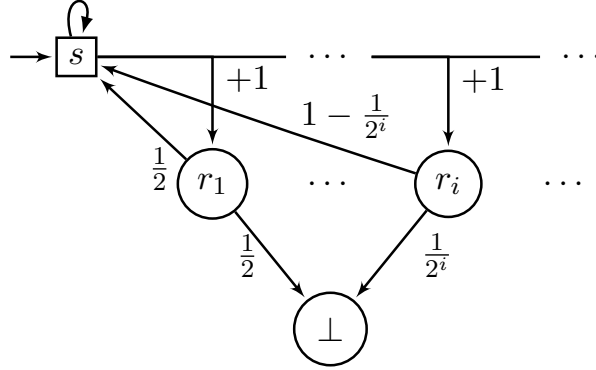


Figure 10: An MDP where ε -optimal strategies for $MP_{\liminf \geq 0}$ require only memory that grows unboundedly with the number of steps taken so far.

7. CONCLUSION AND OUTLOOK

We have established matching lower and upper bounds on the strategy complexity of \liminf threshold objectives for point, total and mean payoff on countably infinite MDPs; cf. Table 1.

The upper bounds hold not only for integer transition rewards, but also for rationals or reals, provided that the reward counter (in those cases where one is required) is of the same type. The lower bounds hold even for integer transition rewards, since all our counterexamples are of this form.

Directions for future work include the corresponding questions for \limsup threshold objectives. While the \liminf point payoff objective generalizes co-Büchi (see Section 2), the \limsup point payoff objective generalizes Büchi. Thus the lower bounds for \limsup point payoff are at least as high as the lower bounds for Büchi objectives [20, 21].

REFERENCES

- [1] P. Abbeel and A. Y. Ng. Learning first-order Markov models for control. In *Advances in Neural Information Processing Systems 17*, pages 1–8. MIT Press, 2004.
- [2] P. A. Abdulla, R. Ciobanu, R. Mayr, A. Sangnier, and J. Sproston. Qualitative analysis of VASS-induced MDPs. In *International Conference on Foundations of Software Science and Computational Structures (FoSSaCS)*, volume 9634, 2016.
- [3] G. Ashkenazi-Golan, J. Flesch, A. Predtetchinski, and E. Solan. Reachability and safety objectives in Markov decision processes on long but finite horizons. *Journal of Optimization Theory and Applications*, 185:945–965, 2020.
- [4] C. Baier and J.-P. Katoen. *Principles of Model Checking*. MIT Press, 2008.
- [5] N. Bäuerle and U. Rieder. *Markov Decision Processes with Applications to Finance*. Springer-Verlag Berlin Heidelberg, 2011.
- [6] P. Billingsley. *Probability and Measure*. Wiley, New York, NY, 1995. Third Edition.
- [7] V. D. Blondel and J. N. Tsitsiklis. A survey of computational complexity results in systems and control. *Automatica*, 36(9):1249–1274, 2000.
- [8] T. Brázdil, V. Brožek, K. Etessami, and A. Kučera. Approximating the Termination Value of One-Counter MDPs and Stochastic Games. *Information and Computation*, 222:121–138, 2013.
- [9] T. Brázdil, V. Brožek, K. Etessami, A. Kučera, and D. Wojtczak. One-counter Markov decision processes. In *ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 863–874. Society for Industrial and Applied Mathematics, 2010.
- [10] T. Bromwich. *An Introduction to the Theory of Infinite Series*. McMillan and Company, London, 1955.
- [11] K. Chatterjee and L. Doyen. Games and Markov decision processes with mean-payoff parity and energy parity objectives. In *Proc. of MEMICS*, volume 7119 of *LNCS*, pages 37–46. Springer, 2011.
- [12] K. Chatterjee, L. Doyen, and T. Henzinger. A survey of stochastic games with limsup and liminf objectives. In *Proc. of ICALP*, volume 5556 of *LNCS*. Springer, 2009.
- [13] E. Clarke, O. Grumberg, and D. Peled. *Model Checking*. MIT Press, Dec. 1999.
- [14] E. M. Clarke, T. A. Henzinger, H. Veith, and R. Bloem, editors. *Handbook of Model Checking*. Springer, 2018.
- [15] L. E. Dubins and L. J. Savage. *How to Gamble If You Must: Inequalities for Stochastic Processes*. Dover Publications Inc., 2014.
- [16] K. Etessami, D. Wojtczak, and M. Yannakakis. Quasi-birth–death processes, tree-like QBDs, probabilistic 1-counter automata, and pushdown systems. *Performance Evaluation*, 67(9):837–857, 2010.
- [17] K. Etessami and M. Yannakakis. Recursive Markov decision processes and recursive stochastic games. *Journal of the ACM*, 62:1–69, 2015.
- [18] J. Flesch, A. Predtetchinski, and W. Sudderth. Simplifying optimal strategies in limsup and liminf stochastic games. *Discrete Applied Mathematics*, 251:40–56, 2018.
- [19] T. Hill and V. Pestien. The existence of good Markov strategies for decision processes with general payoffs. *Stoch. Processes and Appl.*, 24:61–76, 1987.
- [20] S. Kiefer, R. Mayr, M. Shirmohammadi, and P. Totzke. Büchi objectives in countable MDPs. In *ICALP*, volume 132 of *LIPIcs*, pages 119:1–119:14, 2019. Full version at <https://arxiv.org/abs/1904.11573>.
- [21] S. Kiefer, R. Mayr, M. Shirmohammadi, and P. Totzke. Strategy Complexity of Parity Objectives in Countable MDPs. In *International Conference on Concurrency Theory (CONCUR)*, volume 171 of *LIPIcs*, 2020.
- [22] S. Kiefer, R. Mayr, M. Shirmohammadi, and P. Totzke. Transience in countable MDPs. In *Proc. of CONCUR*, volume 203 of *LIPIcs*, 2021. Full version at <https://arxiv.org/abs/2012.13739>.
- [23] S. Kiefer, R. Mayr, M. Shirmohammadi, and D. Wojtczak. Parity Objectives in Countable MDPs. In *LICS*. IEEE, 2017.
- [24] A. P. Maitra and W. D. Sudderth. *Discrete Gambling and Stochastic Games*. Springer-Verlag, 1996.
- [25] A. Nowak. *Advances in dynamic games : applications to economics, finance, optimization, and stochastic control*. Birkhaeuser, Boston, 2005.
- [26] D. Ornstein. On the existence of stationary optimal strategies. *Proceedings of the American Mathematical Society*, 20:563–569, 1969.
- [27] M. L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, Inc., New York, NY, USA, 1st edition, 1994.

- [28] T. Raghaven, T. Ferguson, T. Parthasarathy, and O. Vrieze. *Stochastic games and related topics: in honor of Professor LS Shapley*, volume 7. Springer Science & Business Media, 2012.
- [29] S. M. Ross. *Introduction to Stochastic Dynamic Programming*. Academic Press, New York, 1983.
- [30] M. Schäl. Markov decision processes in finance and dynamic options. In *Handbook of Markov Decision Processes*, pages 461–487. Springer, 2002.
- [31] O. Sigaud and O. Buffet. *Markov Decision Processes in Artificial Intelligence*. John Wiley & Sons, 2013.
- [32] W. D. Sudderth. Optimal Markov strategies. *Decisions in Economics and Finance*, 43:43–54, 2020.
- [33] R. Sutton and A. Barto. *Reinforcement Learning: An Introduction*. Adaptive Computation and Machine Learning. MIT Press, 2018.
- [34] M. Vardi. Automatic verification of probabilistic concurrent finite-state programs. In *Proc. of FOCS’85*, pages 327–338, 1985.