# Eventually-stationary policies for Markov decision models with non-constant discounting

## (Invited paper)

Yair Carmon
Department of Electrical Engineering
Technion — IIT
Haifa 32000, Israel

Adam Shwartz[*]
Department of Electrical Engineering
Technion — IIT
Haifa 32000, Israel
adam@ee.technion.ac.il

## ABSTRACT

We investigate the existance of simple policies in finite discounted cost Markov Decision Processes, when the discount factor is not constant. We introduce a class called "exponentially representable" discount functions. Within this class we prove existence of optimal policies which are eventually stationary—from some time $N$ onward, and provide an algorithm for their computation. Outside this class, optimal policies with this structure in general do not exist.

## Keywords

Markov Decision Processes, Discounted Cost, Mixed Discounting, Hyperbolic Discounting, General Discounting Function

## 1. INTRODUCTION

In Markov decision models (MDPs), discounting is used to model the fact that the further in the future something happens, the less important it is. Simple discounting, where the reward is multiplied by a constant discount factor at each epoch, arises naturally from economic considerations, when constant rates of interest (bonds) or constant inflation are assumed. Such models are relatively easy to analyze. In particular, in this case there exists an optimal policy which is stationary, namely deterministic and independent of the time and of past states [5]. Because it can be intuitively understood and handily analyzed, simple discounting has been thoroughly researched and applied to countless models—from machine learning, computer networks to game theory and psychology.

However, in general the decrease in value of the future need not be geometric. In models of "learning curves", the cost of "getting to know" the system is added to the original criterion. The discounting in this criterion typically

---

[*]Corresponding Author

has a power-law form, though some models have geometric learning curves. Additional geometric decreases arise from models such as "Moore's law," in which a cost (in this case of a unit of computation power) decreases geometrically fast. The addition of an exponentially-decreasing component (with a rate different from the discount factor) to a discounted Markov decision model results in a weighted discounted criterion—that is, a criterion that is the sum of several standard discounted criteria. A theory for finite models with weighted discounted criteria was developed by Feinberg and Shwartz [2]. The main results are that for such criteria there are optimal policies that are stationary from some finite time onwards, called $N$-stationary policies, and an algorithm for the computation of these policies is given.

In models of human preferences, it makes sense to use discounting with a decreasing rate. Intuitively, while tomorrow may be less important than today, a year and a day from now is just about the same as a year from now. "Hyperbolic discount functions", of the form $(1 + \alpha n)^{-\gamma/\alpha}$ with $\alpha, \gamma > 0$, feature a decreasing discounting rate, and are reported to effectively model psychological preferences (see [3] for presentation, and [6] for critique). Even Moore's Law seems to be breaking down, leading to non-constant discounting.

Because of the difficulty of analyzing decision processes with general discount functions, most theoretical results are obtained with "toy functions", for example $f = [1, \delta\beta, \delta\beta^2, \delta\beta^3, \ldots]$, which in a sense has a decreasing discounting rate for $0 < \delta < 1$, and often serves as a replacement of the hyperbolic function mentioned above (see, e.g. [4]). Following the lines of the weighted discounted theory, we define the class of "exponentially representable" functions, prove that when they are used as discount functions there exist $N$-stationary optimal policies, and describe a computation algorithm. These functions display the decreasing discount rate of the hyperbolic discount functions. However, the hyperbolic discount functions are not exponentially representable, and moreover—exponentially representable discount functions cannot be used to model power-law learning curves.

Below we define the model and the exponentially representable functions, and state our main result. In section 2 we develop the algorithm for the computation of the optimal policy and through it prove our result. In section 3 we further discuss the meaning of exponential representability and which functions belong in that class. Finally, in section 4 we show that the $N$-stationary property is not always assured.

## 1.1 Markov Decision Processes

Consider a discrete time process with a finite state space $\mathbf{X}$, where $x_n \in \mathbf{X}$ denotes the state at time $n = 0, 1, 2, \ldots$. At each time, an action is chosen from a finite action set $\mathbf{A}(x)$. Let $\mathbf{A} = \bigcup_{x \in \mathbf{X}} \mathbf{A}(x)$ be the (finite) action space—the set of possible actions and $a_n \in \mathbf{A}(x_n)$ the chosen action at time $n$. The state and chosen action at time $n$ determine the probability distribution of the next state through the transition probability $p(x_{n+1}|x_n, a_n)$, assumed to be independent of the time and of further information on the past. Additionally, the state and action determine the immediate reward $r(x_n, a_n)$ given at time $n$.

The rule used to choose the action at each time is called a policy. We call $h_n = x_0 a_0 \cdots x_{n-1} a_{n-1} x_n$ the history at time $n$. Since the choice of action is required to be causal, $h_n$ contains all the information available for making it. The most general policy is therefore a mapping of every history $h_n$ to a probability measure $\pi(\cdot|h_n)$ on $\mathbf{A}(x_n)$. Policies in which only one action is possible for each given history are called deterministic, so that $a_n = \pi(h_n)$. Policies which depend only on the current state and the time, i.e. $\pi(\cdot|h_n) = \pi(\cdot|x_n, n)$ are called Markov policies, and Markov policies which do not depend on the time are called stationary.

The discounted criterion assigns to each policy and initial state numerical value

$$V^\beta(x; \pi) = \mathbb{E}_x^\pi \sum_{n=0}^\infty \beta^n r(x_n, a_n)$$

where $\mathbb{E}_x^\pi$ is the expectation induced by policy $\pi$, given that $h_0 = x$, and $0 < \beta < 1$ is the discount factor. Since $r(x, a)$ is bounded, the value is always finite.

We discuss a more general discounted criterion, in which $\beta^n$ is replaced with $f(n)$:

$$V^{\text{gen}}(x; \pi) = \mathbb{E}_x^\pi \sum_{n=0}^\infty f(n) r(x_n, a_n) . \qquad (1)$$

A sufficient condition for the above summation to be well defined is $|f(n)| \le K\beta^n$ for some $0 < \beta < 1$ or, equivalently, $f(n) = \beta^n g(n)$ for $0 < \beta < 1$ and some bounded function $g(n)$. We call a function that satisfies this condition *exponentially bounded*.

The weighted discounted criterion is a sum of a finite number of standard discounted criteria, each with a possibly different immediate reward function:

$$V^{\text{wd}}(x; \pi)$$
$$= \mathbb{E}_x^\pi \sum_{n=0}^\infty \sum_{k=1}^K \beta_k^n r_k(x_n, a_n) \text{ with } \beta_1 > \beta_2 > \cdots > \beta_K . \quad (2)$$

Let us define the maximal and minimal values of an MDP, respectively:

$$V(x) \equiv \sup_\pi V(x; \pi) \qquad V^-(x) \equiv \inf_\pi V(x; \pi) . \qquad (3)$$

An optimal policy is a policy for which $V(x; \pi) = V(x)$, for all $x \in \mathbf{X}$.

DEFINITION 1.1. *A Markov policy $\pi$ is called* N-stationary *if*

$$\pi(x, n) = \pi(x, N) \quad \forall x \in \mathbf{X}, n \ge N .$$

DEFINITION 1.2. *A function $f : \{0, 1, \ldots\} \to \mathbb{R}$ is called* exponentially representable *(ER) if there exist sequences $\{c_k\}_{k=1}^\infty$ and $\{\beta_k\}_{k=1}^\infty$ such that:*

- $\{\beta_k\}_{k=1}^\infty$ *is positive, strictly decreasing and $\beta_1 < 1$.*

- $f(n) = \sum_{k=1}^\infty c_k \beta_k^n$, *the sum converges absolutely after some $N < \infty$.*

EXAMPLE 1.3. *The function*

$$f(n) = \frac{1}{e-1} \sum_{k=1}^\infty \frac{1}{k!} \left(\beta_0^k\right)^n = \frac{e^{\beta_0^n} - 1}{e - 1}$$

*is ER for $0 < \beta_0 < 1$ and logarithmically convex $((\log f(x))'' > 0)$—hence with a decreasing rate of discounting, since the rate is inversely proportional to $f(n+1)/f(n)$. This is the required property in aforementioned the human preferences models.*

## 1.2 The main result

Our starting point will be the following result on the structure of optimal policies under criteria (2) and (1).

THEOREM 1.4. *In a weighted discounted MDP (2), there exists an optimal policy which is Markov and deterministic. This holds also in MDPs with general discounting (1) when the discount function is exponentially bounded.*

For a weighted discounted MDP, a proof (under more general conditions) is given in [2], Thm. 2.2. The idea is to embed the process in an ordinary discounted MDP with a countable state space $\tilde{\mathbf{X}} = (\mathbf{X} \times \mathbb{N}_+)$, where time is added so that the new state is $\tilde{x}_n = (x_n, n)$, and use the standard result that discounted criteria have deterministic and stationary optimal policies. The same embedding can be carried out in the case of a single general discount function which is exponentially bounded. This theorem also extends straightforwardly to a criterion that is a sum of several criteria with general discount functions, as long as those functions are exponentially bounded. Since this is a straightforward extension, we omit the details.

From now on we focus on ER discount functions. Since they are exponentially bounded, using Theorem 1.4 we restrict our policies to be Markov and deterministic.

THEOREM 1.5. *Consider a finite MDP with an ER discount function. There exists an N-stationary optimal policy for this problem, with $N < \infty$. This policy can be found using Algorithm 2.6.*

## 2. OPTIMAL POLICIES

The generalized discounted criterion in (1), with $f(n)$ ER, is an infinite version of the weighted discounted criterion. To see this, find $\{c_k\}_{k=1}^\infty$ and decreasing $\{\beta_k\}_{k=1}^\infty$ such that $f(n) = \sum_{k=1}^\infty c_k \beta_k^n$, and rewrite the criterion as

$$V^{\text{gen}}(x; \pi) = \mathbb{E}_x^\pi \sum_{n=0}^\infty f(n) r(x_n, a_n)$$

$$= \mathbb{E}_x^\pi \sum_{n=0}^\infty \sum_{k=1}^\infty \beta_k^n c_k r(x_n, a_n) \qquad (4)$$

which is an infinite weighted discounted criterion with $r_k(x_n, a_n) = c_k r(x_n, a_n)$.

We now adapt the algorithm from part 3 of [2] to infinite weighted discounted criteria induced by an ER discount function. We will also prove that this algorithm halts after a finite number of iterations, and provide a bound on that number. Let

$$V_k(x; \pi) = \mathbb{E}_x^\pi \sum_{n=0}^\infty \beta_k^n c_k r(x_n, a_n) \tag{5}$$

denote the value of the $k^{\text{th}}$ summand in (4), and let $V_k(x)$ and $V_k^-(x)$ be the maximal and minimal value for initial state $x$. Define a "conserving set":

$$\Gamma_1(x) \equiv$$
$$\left\{ a \in \mathbf{A}(x) \mid V_1(x) = c_1 r(x, a) + \beta_1 \sum_{y \in \mathbf{X}} p(y|x, a) V_1(y) \right\}. \tag{6}$$

It is easy to see that a policy is optimal for this criterion if and only if it chooses actions from the (conserving) set $\Gamma_1(x)$ when in state $x$: see Lemma 3.1 in [2].

Let $\mathbf{X}_1 = \{x \in \mathbf{X} \mid \Gamma_1(x) \neq \mathbf{A}(x)\}$ be the set of states for which suboptimal actions for criterion $V_1$ exist. If $\mathbf{X}_1 \neq \emptyset$, define:

$$\varepsilon_1 \equiv \min_{x \in \mathbf{X}_1, a \in \mathbf{A}(x) \backslash \Gamma_1(x)}$$
$$\left( V_1(x) - c_1 r(x, a) - \beta_1 \sum_{y \in \mathbf{X}} p(y|x, a) V_1(y) \right). \tag{7}$$

$\varepsilon_1$ is the value of the smallest "mistake" one can make in the choice of a single action, with regard to criterion $V_1$. If $\mathbf{X}_1 = \emptyset$ define $N_1 \equiv 0$. Otherwise define:

$$N_1 = \min \left\{ n \geq 0 \mid \varepsilon_1 > \sum_{k=2}^\infty \left( \frac{\beta_k}{\beta_1} \right)^n \max_{x \in \mathbf{X}} \left( V_k(x) - V_k^-(x) \right) \right\}. \tag{8}$$

LEMMA 2.1. *If $f(n)$ is ER, $N_1$ is well defined and finite.*

PROOF. Define $S(n) = \sum_{k=2}^\infty (\beta_k/\beta_1)^n \max_{x \in \mathbf{X}} (V_k(x) - V_k^-(x))$. Denote the span semi-norm of $r(x, a)$ by $M = \max_{x \in \mathbf{X}, a \in \mathbf{A}(x)} r(x, a) - \min_{x \in \mathbf{X}, a \in \mathbf{A}(x)} r(x, a)$. Then

$$\forall k : \max_{x \in \mathbf{X}} \left( V_k(x) - V_k^-(x) \right) \leq |c_k| \frac{M}{1 - \beta_k} \leq |c_k| \frac{M}{1 - \beta_1} \tag{9}$$

$$S(n) \leq \frac{\beta_1^{-n} M}{1 - \beta_1} \sum_{k=2}^\infty \beta_k^n |c_k|. \tag{10}$$

This proves $N_1$ is finite—see [1]. $\square$

REMARK 2.2. *Let $\mathbf{X} = \{x_0\}$ and $\mathbf{A} = \{a_1, a_2\}$ with $r(x, a_1) = 1$ and $r(x, a_2) = 0$. Then, $\max_{x \in \mathbf{X}} \left( V_k(x) - V_k^-(x) \right) = |c_k|/1 - \beta_k$ for any $k$, so $S(n) = \beta_1^{-n} \sum_{k=2}^\infty \beta_k^n |c_k|$. Here $S(n) \xrightarrow[n \to \infty]{} 0$ only if $\sum_{k=2}^\infty \beta_k^n c_k$ converges absolutely for some $N < \infty$, i.e. only if $f(n)$ is ER. Thus for a given discount function $f$, the bound $N_1$ is well-defined for any model if and only if the discount function is ER.*

*Using the definitions (6) and (8) of $\Gamma_1(x)$ and $N_1$ respectively,*

LEMMA 2.3. *Consider a finite MDP with an ER discount function. If $\sigma$ is an optimal Markov policy then $n \geq N_1$ and $\mathbb{P}_x^\sigma \{x_n = z\} > 0$ imply $\sigma(z, n) \in \Gamma_1(z)$.*

PROOF. The proof extends that of Lemma 3.3 in [2], using the same basic ideas. It is simple to show (see [2]) that the optimality of $\sigma$ implies that for any stationary policy $\phi$, time $m$ and state $z \in \mathbf{X}$ such that $\mathbb{P}_x^\sigma \{x_m = z\} > 0$:

$$\mathbb{E}_x^\sigma \left\{ \sum_{n=m}^\infty f(n) r(x_n, a_n) | x_m = z \right\}$$
$$\geq \mathbb{E}_x^\phi \left\{ \sum_{n=m}^\infty f(n) r(x_n, a_n) | x_m = z \right\}.$$

By substituting $f(n) = \sum_{k=1}^\infty c_k \beta_k^n$ in the above equation, denoting the time-shifted optimal policy by $\sigma^m(x, n) = \sigma(x, n + m)$, and using $V_k$ as defined in (5), the above inequality can be rewritten ask

$$\sum_{k=1}^\infty \beta_k^m V_k(z; \sigma^m) \geq \sum_{k=1}^\infty \beta_k^m V_k(z; \phi)$$

and therefore

$$\sum_{k=2}^\infty \beta_k^m \left( V_k(z; \sigma^m) - V_k(z; \phi) \right) \geq \beta_1^m \left( V_1(z; \phi) - V_1(z; \sigma^m) \right). \tag{11}$$

We now suppose by contradiction that for some time $m \geq N_1$ and for some state $z \in \mathbf{X}$ such that $\mathbb{P}_x^\sigma \{x_m = z\} > 0$, $\sigma(z, m) \notin \Gamma_1(z)$. Let $\phi$ be the optimal stationary policy for criterion $V_1$, so that $V_1(x, \phi) = V_1(x)$ for all $x \in \mathbf{X}$. Since $m \geq N_1$, and by the definition of $N_1$ (which is meaningful due to the fact that $f$ is ER) we have:

$$\varepsilon_1 > \sum_{k=2}^\infty \left( \frac{\beta_k}{\beta_1} \right)^m [V_k(z) - V_k^-(z)]$$
$$\geq \sum_{k=2}^\infty \left( \frac{\beta_k}{\beta_1} \right)^m [V_k(z; \sigma^m) - V_k^-(z; \phi)]$$
$$\geq V_1(z; \phi) - V_1(z; \sigma^m) \tag{12}$$

where the second inequality comes from the definitions and the last from (11). On the other hand, since $\sigma(z, m) \notin \Gamma_1(z)$ and from the definition of $\varepsilon_1$:

$$V_1(z) - V_1(z; \sigma^m) \geq V_1(z)$$
$$- \left( r(z, \sigma(z, m)) + \beta_1 \sum_{y \in \mathbf{X}} p(y|z, \sigma(z, m)) V_1(y) \right) \geq \varepsilon_1 \tag{13}$$

Inequalities (12) and (13) contradict each other, thus proving this lemma. $\square$

If $\Gamma_1(x)$ is a singleton for all $x \in \mathbf{X}$, then the lemma requires any optimal policy to be $N_1$-stationary, and determines the stationary part of the policy. If it is not a singleton, we know that after time $N_1$ our action sets reduce to $\Gamma_1(x)$ and for every admissible policy, $V_1$ will attain its maximum value and therefore be irrelevant. Our task therefore becomes finding the optimal policy for the weighted sum starting from the second discount factor, with the action sets restricted to $\Gamma_1$. To iterate the above process, define recursively for $k > 1$, the restricted action sets in iteration $k$ — $\mathbf{A}_k(x) = \Gamma_{k-1}(x)$,

the $m^{\text{th}}$ value function restricted to the $k^{\text{th}}$ action set —
$V_m^{\mathbf{A}_k}(x)$, and similarly the minimal value function $V_m^{-,\mathbf{A}_k}(x)$.
Additionally:

$$\Gamma_k(x) \equiv \Big\{ a \in \mathbf{A}_k(x) \mid V_k^{\mathbf{A}_k}(x) = c_k r(x,a) + \beta_k$$
$$\sum_{y \in \mathbf{X}} p(y|x,a) V_k^{\mathbf{A}_k}(y) \Big\} \quad (14)$$

$$\mathbf{X}_k = \{ x \in \mathbf{X} \mid \Gamma_k(x) \neq \mathbf{A}_k(x) \} \quad (15)$$

$$\varepsilon_k \equiv \min_{x \in \mathbf{X}_k, a \in \mathbf{A}_k(x) \setminus \Gamma_k(x)}$$
$$\left( V_k^{\mathbf{A}_k}(x) - c_k r(x,a) - \beta_k \sum_{y \in \mathbf{X}} p(y|x,a) V_k^{\mathbf{A}_k}(y) \right) \quad (16)$$

$$N_k = \min \left\{ n \geq N_{k-1} \mid \varepsilon_k > \sum_{m=k+1}^{\infty} \left( \frac{\beta_m}{\beta_k} \right)^n \right.$$
$$\left. \max_{x \in \mathbf{X}} \left( V_m^{\mathbf{A}_k}(x) - V_m^{-,\mathbf{A}_k}(x) \right) \right\} \quad (17)$$

where $\varepsilon_k$ is set to $\infty$ in the case that $\mathbf{X}_k = \emptyset$. Again, $N_k$
is well defined when $f(n)$ is ER. With these definitions, the
following is evident:

LEMMA 2.4. *Consider a finite MDP with ER discount func-*
*tion. If $\sigma$ is an optimal Markov policy, then for every $k \geq 1$,*
*$n \geq N_k$ and state $z \in \mathbf{X}$ such that $\mathbb{P}_x^\sigma \{ x_n = z \} > 0$, we have*
*$\sigma(z,n) \in \Gamma_k(z)$.*

PROOF. By induction using Lemma 2.3 and the above
definitions. $\square$

We will now prove that iterating this procedure does indeed
provide us with an $N$-stationary policy after a finite and
bounded number of computations.

LEMMA 2.5. *Consider a finite MDP with an ER discount*
*function. Let $S = |\mathbf{X}|$. Then for all $k \geq 2S - 1$ and every*
*$x \in \mathbf{X}$, $\Gamma_k(x) = \Gamma_{2S-1}(x)$.*

PROOF. If $\Gamma_{2S-1}(x)$ is a singleton for all $x \in \mathbf{X}$, then the
lemma is immediate. Otherwise, let $\Phi = \{\phi_1, \phi_2, ..., \phi_L\}$ be
the set of stationary policies such that $\phi_i(x) \in \Gamma_{2S-1}(x)$ for
all $x \in \mathbf{X}$, $i = 1, 2, ..., L$. For $\phi \in \Phi$, define $f_\phi : [0,1) \to \mathbb{R}^S$
as

$$[f_\phi(\beta)]_s = \mathbb{E}_{x^s}^\phi \sum_{n=0}^{\infty} \beta^n r(x_n, a_n) ,$$

so that $V_k(x^s; \phi) = c_k(f_\phi(\beta_k))_s$. Let $[P_\phi]_{m,n} \equiv p(x^n|x^m, \phi$
$(x^s))$ and $[r_\phi]_s = r(x^s, \phi(x^s))$ be the state transition matrix
and reward vector induced by $\phi_i$. Then

$$f_\phi(\beta) = r_\phi + \beta P_\phi f_\phi(\beta) \Rightarrow f_\phi(\beta) = (I - \beta P_\phi)^{-1} r_\phi . \quad (18)$$

Since $P_\phi$ is stochastic, by the Perron–Frobenius theorem $I -$
$\beta P_\phi$ is invertible for $\beta \in [0,1)$ and singular for $\beta = 1$. Since
$M^{-1} = \text{adj}(M) / \det(M)$, by (18) every entry (coordinate)
of $f_\phi$ is a rational function of $\beta$, with numerator degree $S-1$
and denominator degree $S$, with a pole at $\beta = 1$, which
possibly cancels with a zero in some of the entries. Since
$\phi \in \Phi$ if and only if it is optimal for all criteria $V_k$ for
$k = 1, 2, ..., 2S-1$ (under different action sets for each $k$), all

policies in $\Phi$ must have the same values for $\beta_1, \beta_2, ..., \beta_{2S-1}$.
Consequently, for every $i, j \leq L$:

$$f_{\phi_i}(\beta_k) = f_{\phi_j}(\beta_k) , \ \forall k = 1, 2, ..., 2S - 1 . \quad (19)$$

Fix $i$ and $j$ and consider each entry of $f_{\phi_i}(\beta) - f_{\phi_j}(\beta) = 0$.
It is a polynomial equation of degree $2S - 2$ (since the com-
mon poles at $\beta = 1$ cancel). However, according to (19),
this polynomial has $2S - 1$ distinct roots—and is therefore
identically zero. We conclude that $f_{\phi_i}(\beta) = f_{\phi_j}(\beta)$ for all
$\beta \in [0,1)$ and every two policies $\phi_i, \phi_j \in \Phi$, and accordingly
$V_k(x; \phi)$ is the same over all $\phi \in \Phi$, for each $x \in \mathbf{X}$ and
$k \geq 2S - 1$. This means that for $k \geq 2S - 1$, all possible
policies have identical values, and will therefore all be opti-
mal. Since the set of optimal policies remains constant, so
do the conserving sets. $\square$

The proof of Theorem 1.5 now follows—see [1] for details.

The computation of $\{N_k\}_{k=1}^{2S-1}$ involves evaluations of in-
finite sums, which are usually not feasible. To avoid this, we
can instead find upper bounds $\hat{N}_k \geq N_k$ for each $k$, and com-
pute an $\hat{N}_{2S-1}$-stationary optimal policy with a stationary
part determined by the conserving sets. One way to find $\hat{N}_k$
is to use the semi-norm bounds in (10). In each iteration, the
semi-norm of the reward function should be computed with
respect to the restricted action set, and therefore decrease.

ALGORITHM 2.6.

1. Find $\{\beta_k\}_{k=1}^{\infty}$ and $\{c_k\}_{k=1}^{\infty}$ of Definition 1.2, set $S = |\mathbf{X}|$ and $k = 1$.

2. Compute $\Gamma_k(x)$ for all $x \in \mathbf{X}$, $\varepsilon_k$ and $N_k$ or an appro-
priate upper bound.

3. If $\Gamma_k(x)$ is a singleton for every $x \in \mathbf{X}$, or $k = 2S+1$,
set $N = N_k$ and continue. Else set $\mathbf{A}_{k+1}(\cdot) = \Gamma_k(\cdot)$,
increment $k$ by 1 and go back to step 2.

4. Fix a stationary policy $\psi$, such that $\psi(x) \in \Gamma_k(x)$ for
all $x \in \mathbf{X}$.

5. Compute an optimal Markov policy $\sigma$ for the $N$-step
MDP with immediate rewards $r_n(x,a) = f(n) r(x,a)$,
for $n = 0, 1, ... N - 1$ and terminal reward

$$\mathbb{E}_{x_N}^\psi \sum_{n=0}^{\infty} f(n+N) r(x_{n+N}, a_{n+N}) = \sum_{k=1}^{\infty} \beta_k^N V_k(x_N; \psi) .$$

6. The $N$-stationary optimal policy is defined by $\pi(x,n) = \sigma(x,n)$ for $n < N$ and $\pi(x,n) = \psi(x)$ for $n \geq N$.

Step 5 is standard—see [5].

REMARK 2.7. *Our results can be extended to criteria of*
*the form:*

$$V(x; \pi) = \mathbb{E}_x^\pi \sum_{n=0}^{\infty} \sum_{k=1}^{K} f_k(n) r_k(x_n, a_n) , \ f_k(n) = \sum_{i=1}^{\infty} c_{i,k} \beta_{i,k}^n$$

*where for each $k$, $f_k(n)$ is ER and the additional condition*

$$\beta_{i,k} > \beta_{1,k+1} , \ \forall i, k . \quad (20)$$

*Lemmas 2.3 and 2.4 can be extended by changing the defi-*
*nitions of the $N_k$'s to include the rest of the discount func-*
*tions, with condition (20) making sure they remain well de-*
*fined. The $N$-stationary optimal policy can then be obtained*

by finding $\Gamma_{2S-1,1}(x)$ for the first discount function. In the case it is not a singleton, the action space will be restricted appropriately, and the procedure will be applied to $f_2$. This may continue until $\Gamma_{2S-1,K}(x)$ is computed, from which we may choose the stationary part of the optimal policy arbitrarily. Finally, we remark that if $r_k(\cdot) = b_k r(\cdot)$ for some function $r(x,a)$, the procedure will end in the computation of $\Gamma_{2S-1,1}(x)$, since afterwards all permissible policies for the stationary part will have the same value.

# 3. STRUCTURE OF ER FUNCTIONS

ER functions behave asymptotically as exponential functions:

LEMMA 3.1. *Let $f(n)$ be an ER function. Then there exist $0 < \beta < 1$ such that*

$$\lim_{n\to\infty} \beta^{-n} f(n) = c \neq 0 \text{ and } c < \infty .$$

PROOF. Write $f(n) = \sum_{k=1}^{\infty} c_k \beta_k^n$. Without loss of generality, we may assume that $c_1 \neq 0$. Since $f$ is ER, we have absolute convergence from some time $N < \infty$. Therefore, for $n > N$ and some $C < \infty$:

$$\beta_1^{-n} \left| \sum_{k=2}^{\infty} c_k \beta_k^n \right| \leq \beta_1^{-n} \sum_{k=2}^{\infty} |c_k| \beta_k^n$$
$$< \frac{\beta_2^{n-N}}{\beta_1^n} \sum_{k=2}^{\infty} |c_k| \beta_k^N = C \left( \frac{\beta_2}{\beta_1} \right)^n \underset{n\to\infty}{\to} 0 .$$

Consequently, $\lim_{n\to\infty} \beta_1^{-n} \sum_{k=2}^{\infty} c_k \beta_k^n = 0$ and choosing $\beta \equiv \beta_1$ we have,

$$\lim_{n\to\infty} \beta^{-n} f(n) = \lim_{n\to\infty} c_1 + \beta_1^{-n} \sum_{k=2}^{\infty} c_k \beta_k^n = c_1 \neq 0$$

and $c_1 < \infty$. $\square$

Functions with power-law form, like $(1+n^2)^{-1}$ or the hyperbolic discount function do not satisfy the conclusion of Lemma 3.1, and are therefore not ER. The same holds for sub-exponential functions, like $1/n!$ and $e^{-n^2}$. Moreover, functions of the form $g(n)\beta^n$, where $g(n) \to 0$ or $g(n) \to \infty$ non-exponentially, are also not ER for the same reason. Examples are $n\beta^n$ and $\beta^n/(1+n)$ for some $0 < \beta < 1$.

# 4. A DISCOUNT FUNCTION WITH NO $N$-STATIONARY OPTIMAL POLICY

When a discount function decreases monotonically it seems natural that it should produce a behavior that is monotonic, or stationary. However, this is not true: we provide an example of a discount function and a model for which there is no $N$-stationary optimal policy. By our results, the discount function is not ER.

Consider the function $f(n) = \beta^n h(n)$, with some $0 < \beta < 1$ and

$$h(n) = \begin{cases} 2 & n \bmod 6 = 0 \\ 1 & \text{otherwise} \end{cases} = [2,1,1,1,1,1,2,1,1,1,\ldots]$$

which is periodic with period 6. The condition of Lemma 3.1 does not hold for $f(n)$, and it is therefore not ER.

Now consider the following (deterministic) model:

$$\mathbf{X} = \{1,2,3,4,5\} \ , \ \mathbf{A}(1) = \{a_1, a_2\} \ ,$$
$$\mathbf{A}(2) = \mathbf{A}(3) = \mathbf{A}(4) = \mathbf{A}(5) = \{a\} \tag{21}$$
$$p(2|1,a_1) = p(3|1,a_2) = p(4|3,a)$$
$$= p(5|4,a) = p(1|5,a) = p(1|2,a) = 1$$

with the immediate reward function

$$r(1,a_1) = 3 \ , \ r(1,a_2) = 4 \ ,$$
$$r(2,a) = r(3,a) = r(4,a) = r(5,a) = 0 \ .$$

Let $\sigma$ be a hypothetical $N$-stationary, optimal and possibly randomized policy for this process. Since state 1 is recurrent, there exists a time $M_0 \geq N$ such that $\mathbb{P}_1^\sigma(x_{M_0} = 1) > 0$. Consequently, $\mathbb{P}_1^\sigma(x_{M_0+4} = 1) > 0$ and $\mathbb{P}_1^\sigma(x_{M_0+8} = 1) > 0$, since at those times there must be a positive probability that a single action is used repeatedly. We know $M_0$ is even because every return to state 1 takes either 2 or 4 steps, and therefore, either $M_0$, $M_0 + 4$ or $M_0 + 8$ divides by 6. We may thus choose $M \geq N$ such that $M$ is a multiple of 6, and $\mathbb{P}_1^\sigma(x_M = 1) > 0$.

Define the shifted value criterion:

$$V^M(1;\pi) \equiv \mathbb{E}^\pi \left\{ \sum_{n=M}^{\infty} \beta^n h(n) r(x_n, a_n) \,|\, x_M = 1 \right\} .$$

Criterion $V^M(1;\pi)$ comprises the part of $V(1;\pi)$ that involves times $M$ and onwards. If an optimal Markov policy for criterion $V(1;\pi)$ has any chance of reaching state 1 in time $M$, it must also optimize criterion $V^M(1;\pi)$, when taken from times $M$ onwards — this is a form of the principle of optimality. Therefore, and considering that $\sigma$ is an optimal Markov policy for $V(1;\pi)$, and $\mathbb{P}_1^\sigma(x_M = 1) > 0$ by our choice of $M$, the stationary policy $\sigma^M(\cdot) = \sigma(\cdot, n+M)$ must be of optimal for criterion $V^M$. However,

$$V^M(1;\pi) = \mathbb{E}_1^\pi \sum_{n=0}^{\infty} \beta^{n+M} h(n+M) r(x_n, a_n) = \beta^M V(1;\pi)$$

where the last equality follows from the periodicity of $h(n)$ and our choice of $M$. Since $V^M(1;\cdot)$ is proportional to $V(1;\cdot)$, $\sigma^M$ is optimal for the original criterion as well, when starting from state 1. Moreover, if $\sigma^M$ is randomized, a stationary and deterministic policy $\hat{\sigma}^M$ with $V(1;\sigma^M) = V(1;\hat{\sigma}^M) = V(1)$ can be obtained. This is done by using, for every state, an action that has positive probability under $\sigma^M$. In conclusion, if this problem has an $N$-stationary optimal policy, a stationary and deterministic policy must maximize $V(1;\cdot)$.

Let $\sigma_1(1) = a_1, \sigma_2(1) = a_2$ be the two stationary policies in this model, and consider the periodic Markov policy:

$$\pi(1,n) = \begin{cases} a_2 & n \bmod 6 = 0 \\ a_1 & n \bmod 6 = 4 \end{cases}$$

Let $\beta = 0.45$. The values of the 3 policies will then be:

$$V(1;\pi) = \frac{8 + 3\beta^4}{1 - \beta^6} \approx 8.19 \tag{22}$$

$$V(1;\sigma_1) = \frac{6 + 3\beta^2 + 3\beta^4}{1 - \beta^6} \approx 6.79 \tag{23}$$

$$V(1;\sigma_2) = \frac{8 + 4\beta^4 + 4\beta^8}{1 - \beta^{12}} \approx 8.17 \ . \tag{24}$$

Evidently, both stationary policies are suboptimal when $x_0 = 1$. Therefore, there cannot be an $N$-stationary optimal policy, since our considerations have shown that it will result in an optimal, stationary and deterministic policy when starting from state 1. Since we chose $\beta < 1/2$, $h(n) > \beta h(n+1)$ for every $n$, making the discount function monotonically decreasing, as required.

## 5. REFERENCES

[1] Y. Carmon and A. Shwartz. "Markov decision processes with exponentially representable discounting," CCIT Report 675, *I*nternal Report, Technion 2008.

[2] E.A. Feinberg and A. Shwartz. Markov decision models with weighted discounted criteria. *M*athematics of Operations Research, 19(1):152–168, 1994.

[3] G. Loewenstein and D. Prelec. Anomalies in intertemporal choice: Evidence and an interpretation. *T*he Quarterly J. Economics, 107(2):573–597, 1992.

[4] D. Laibson. Golden eggs and hyperbolic discounting. *T*he Quarterly J. Economics, 112(2):443–477, 1997.

[5] M.L. Puterman. *M*arkov Decision Processes: Discrete Stochastic Dynamic Programming. John Wiley & Sons, N.Y., 1994.

[6] A. Rubinstein. Economics and psychology? the case of hyperbolic discounting. *I*nternational Economic Review, 44(4):1207–1216, 2003.