



HAL
open science

Singularly perturbed linear programs and Markov decision processes

Konstantin Avrachenkov, Jerzy A Filar, Vladimir G Gaitsgory, Andrew Stillman

► **To cite this version:**

Konstantin Avrachenkov, Jerzy A Filar, Vladimir G Gaitsgory, Andrew Stillman. Singularly perturbed linear programs and Markov decision processes. *Operations Research Letters*, 2016, 44 (3), pp.297 - 301. 10.1016/j.orl.2016.02.005 . hal-01399403

HAL Id: hal-01399403

<https://inria.hal.science/hal-01399403>

Submitted on 21 Nov 2016

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

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

Singularly perturbed linear programs and Markov decision processes

Konstantin Avrachenkov^a, Jerzy A. Filar^b, Vladimir Gaitsgory^c, Andrew Stillman^b

^a*Inria Sophia Antipolis, 2004 Route des Lucioles, 06902 Sophia Antipolis, France;*
K.Avrachenkov@inria.fr

^b*School of Computer Science, Engineering & Mathematics, Flinders University, Australia*

^c*Department of Mathematics, Macquarie University, Australia*

Abstract

Linear programming formulations for the discounted and long-run average MDPs have evolved along separate trajectories. In 2006, E. Altman conjectured that the two linear programming formulations of discounted and long-run average MDPs are, most likely, a manifestation of general properties of singularly perturbed linear programs. In this note we demonstrate that this is, indeed, the case.

Keywords: Markov Decision Processes (MDPs), Discounted MDPs, Long-run average MDPs, Singularly Perturbed Linear Programs, Limiting Linear Program

1. Introduction

The connection between linear programming and Markov Decision Processes (MDPs) was launched in the 1960's, with the papers by D'Epenoux [8], De Ghellinck [9] and Manne [18]. While the linear programming formulation for the discounted MDP was relatively straightforward, extension to the long-run average, multi-chain, MDP proved challenging and required nearly two decades to arrive at a single linear program supplied, by Hordijk and Kallenberg [11, 12], that completely solves such a multi-chain MDP. We refer the reader to Kallenberg [15, 16], Puterman [21] and Altman [1] for excellent, comprehensive, treatments of linear programming methods for discrete time Markov decision processes. Even though the approaches to discounted and long-run average MDPs evolved along separate trajectories, Tauberian theorems provided a theoretical connection between the two cases with the discount parameter approaching unity from below; e.g. see Blackwell [6] and Veinott [25, 26].

Parametric linear programming has a long history that is well documented in many excellent textbooks (e.g., see Murty [19]). However, majority of the so-called sensitivity analyses presented in operations research books focus on perturbations of the objective function coefficients or of the right hand side vector; sometimes extending also to changes in non-basic columns. To the best of our knowledge, Jeroslow [14] was, perhaps, the first to consider perturbations of the entire coefficient matrix of a linear program. In the context of MDP, the results of [14] have been applied to Blackwell optimality [13] and to perturbed MDPs [2]. In Pervozvanskii and Gaitsgory [20] the authors focus on the singularly perturbed case where a discontinuity can arise as the perturbation parameter

approaches a critical value. In the latter and in the more recent book by Avrachenkov et al [4] the main cause of that discontinuity has been the change in the rank of the coefficient matrix at the critical value of the perturbation parameter. Hence, it was perhaps surprising that such discontinuities can also arise when the rank does not change, as shown very recently in [3].

In 2006, Eitan Altman conjectured that the two linear programming formulations of discounted and long-run average MDPs must be a manifestation of some general properties of singularly perturbed linear programs. In this note we demonstrate that this is, indeed, the case by first extending the results in [3] and then formally applying new singular perturbation results to the MDP problem.

2. General perturbed linear programming problem

Consider the family of linear programming problems parameterized by $\varepsilon > 0$:

$$\begin{aligned} & \max \langle c^{(0)} + \varepsilon c^{(1)}, x \rangle \\ \text{s. t. } & (A^{(0)} + \varepsilon A^{(1)})x = b^{(0)} + \varepsilon b^{(1)}, \\ & x \geq 0, \end{aligned} \tag{1}$$

where $c^{(0)}, c^{(1)} \in \mathbb{R}^n$, $b^{(0)}, b^{(1)} \in \mathbb{R}^m$ and $A^{(0)}, A^{(1)} \in \mathbb{R}^{m \times n}$. The optimal value, the solution set and the feasible set of Problem (1) are denoted as $F^*(\varepsilon), \theta^*(\varepsilon)$ and $\theta(\varepsilon)$, respectively.

The goal of the perturbed linear programming approach is to construct, if possible, a linear programming problem that does not depend on ε and such that its optimal solutions are *feasible limiting optimal* for (1) in the sense prescribed below by Definition 1. The linear program with this property will be called a *limiting LP*.

Definition 1. A vector $x \in \mathbb{R}^n$ is called *feasible limiting optimal* for the perturbed linear program (1) if $x \in \liminf_{\varepsilon \downarrow 0} \theta(\varepsilon)$ and $\lim_{\varepsilon \downarrow 0} F^*(\varepsilon) = \langle c^{(0)}, x \rangle$.

Let us introduce and discuss a set of assumptions:

Assumption (H_0): There exists a positive γ_0 and a bounded set $B \subset \mathbb{R}^n$ such that $\theta(\varepsilon) \subset B$ for every $\varepsilon \in (0, \gamma_0]$.

Assumption (H_0^*): There exists a positive γ_0 and a bounded set $B \subset \mathbb{R}^n$ such that $\theta^*(\varepsilon) \subset B$ for every $\varepsilon \in (0, \gamma_0]$.

Assumption (H_1): The matrix $A^{(0)}$ has rank m .

Assumption (H_2): For all ε sufficiently small and positive, the rank of $A^{(0)} + \varepsilon A^{(1)}$ is equal to m .

Note that Assumption (H_1) implies Assumption (H_2). Also, Assumption (H_0) implies Assumption (H_0^*).

The unperturbed problem is said to satisfy Slater condition if

$$\theta(0) \cap \mathbb{R}_{++}^n \neq \emptyset, \text{ where } \mathbb{R}_{++}^n \stackrel{\text{def}}{=} \{x \in \mathbb{R}^n : x > 0\}. \tag{2}$$

In [20], it has been shown that if Assumptions (H_0) and (H_1) are valid and if the Slater condition (2) is satisfied, then the unperturbed LP is the limiting problem for the perturbed program (1). That is, every optimal solution of the former is limiting optimal for the latter. In [20] it has also been shown that if Assumption (H_1) is not satisfied, the discontinuity of $\theta(\varepsilon)$ at $\varepsilon = 0$ may occur. This is a case of so-called *singular perturbation*. The authors of [20] proposed a limiting LP to deal with the case of singular perturbation. Then, in [3] it has been demonstrated that if the Slater condition is not satisfied for the unperturbed LP, the discontinuity of $\theta(\varepsilon)$ at $\varepsilon = 0$ may occur with Assumptions (H_0) and (H_1) being satisfied. The authors of [3] have constructed a limiting LP for the case when the Slater condition is not satisfied for the unperturbed problem. Below we show that a result similar to that obtained in [3] can be established with the replacement of (H_0) by (H_0^*) .

Assume that (H_1) is satisfied and define the set

$$J_0 := \{i \in \{1, \dots, n\} : \exists x \in \theta(0) \text{ such that } x_i > 0\}. \quad (3)$$

According to this definition, if $j \notin J_0$, then $x_j = 0$ for every $x \in \theta(0)$. Moreover, if $J_0 \neq \emptyset$, convexity of $\theta(0)$ implies that there exists $\hat{x} \in \theta(0)$ such that $\hat{x}_j > 0$ for every $j \in J_0$. Note that J_0 can be determined by solving n independent linear programming problems $\max_{x \in \theta(0)} x_j$, with $j = 1, \dots, n$.

Consider the following linear program

$$\max\{\langle c^{(0)}, x^0 \rangle : x^0 \in \theta_1\} \stackrel{\text{def}}{=} F_1^*, \quad (4)$$

where

$$\theta_1 \stackrel{\text{def}}{=} \{x^0 : \exists (x^0, x^1) \in \Theta_1\}, \quad (5)$$

and

$$\Theta_1 = \{(x^0, x^1) \in \mathbb{R}^n \times \mathbb{R}^n : x^0 \in \theta(0), \quad A^{(0)}x^1 + A^{(1)}x^0 = b^{(1)}, \quad x_j^1 \geq 0 \forall j \notin J_0\}. \quad (6)$$

Note that,

$$\theta_1 \subset \theta(0) \quad \text{and therefore} \quad F_1^* \leq F^*(0).$$

Slater condition (2) is equivalent to having $J_0 = \{1, 2, \dots, n\}$. If this is the case, then $\theta_1 = \theta(0)$ (provided that Assumption (H_1) is satisfied), and the problem (4) is equivalent to the unperturbed problem. If the Slater condition is not satisfied, these two problems are not equivalent.

Following [3], let us introduce the following extended version of the Slater condition.

Definition 2. *We shall say that the extended Slater condition of order 1 (or, for brevity, ES-1) is satisfied if there exists $(\hat{x}^0, \hat{x}^1) \in \Theta_1$ such that $\hat{x}_j^1 > 0$ for every $j \notin J_0$ and $\hat{x}_j^0 > 0$ for every $j \in J_0$.*

Theorem 1. *Let Assumptions (H_0^*) and (H_2) be satisfied. Then*

$$\limsup_{\varepsilon \downarrow 0} \theta^*(\varepsilon) \subset \theta_1 \quad (7)$$

and

$$\limsup_{\varepsilon \downarrow 0} F^*(\varepsilon) \leq F_1^*. \quad (8)$$

If, in addition, Assumption (H_1) and the ES-1 condition are satisfied, then

$$\limsup_{\varepsilon \downarrow 0} \theta^*(\varepsilon) \subset \theta_1^*, \quad (9)$$

where θ_1^* is the set of optimal solutions of problem (4), and

$$\lim_{\varepsilon \downarrow 0} F^*(\varepsilon) = F_1^* . \quad (10)$$

Also, any optimal solution x^0 of the problem (4) is limiting optimal for the perturbed problem (1).

Proof. Most steps of the proof are similar to the corresponding steps of the proof of Theorem 2.1 in [3], and we will only indicate the steps that differ from those used in the aforementioned proof.

Let us introduce the following notations. Given a finite set S , denote by $|S|$ the number of elements of S . Let $S_m := \{J \subset \{1, 2, \dots, n\} : |J| = m\}$, so $|S_m| = \binom{n}{m}$. Given a matrix $D \in \mathbb{R}^{m \times n}$ and an index set $J \in S_m$, the matrix $D_J \in \mathbb{R}^{m \times m}$ is constructed by extracting from D the set of m columns indexed by the elements of J . In a similar way, given a vector $x \in \mathbb{R}^n$ and $J \in S_m$, we denote by x_J the vector of \mathbb{R}^m constructed by extracting from x the coordinates x_j , $j \in J$ (that is, $x_J \stackrel{\text{def}}{=} \{x_j\}$, $j \in J$).

In Lemmas 3.1 and 3.2 of [3] it was established that

$$S_m = \Omega_1 \cup \Omega_2 \quad \text{with} \quad \Omega_1 \cap \Omega_2 = \emptyset,$$

where Ω_1 and Ω_2 are defined by the equations

$$\Omega_1 := \{J \in S_m : (A^{(0)} + \varepsilon A^{(1)})_J \text{ is nonsingular for } \varepsilon \in (0, \gamma)\} \neq \emptyset,$$

$$\Omega_2 := \{J \in S_m : (A^{(0)} + \varepsilon A^{(1)})_J \text{ is singular for all } \varepsilon \in [0, \gamma)\}$$

(here and in what follows, γ stands for a positive number small enough).

Also, it was established that, if

$$x_J(\varepsilon) := [(A^{(0)} + \varepsilon A^{(1)})_J]^{-1}(b_0 + \varepsilon b_1) \quad (11)$$

($J \in \Omega_1$) and if

$$\limsup_{\varepsilon \downarrow 0} \|x_J(\varepsilon)\| < \infty, \quad (12)$$

then $x_J(\varepsilon)$ allows the power series expansion

$$x_J(\varepsilon) = \sum_{l=0}^{\infty} \varepsilon^l x_J^l, \quad \forall \varepsilon \in (0, \gamma). \quad (13)$$

Let $\Omega_1^* \subset \Omega_1$ be such that $J \in \Omega_1^*$ if and only if there exists a subsequence $\varepsilon' \rightarrow 0$ such that the vector $x(\varepsilon) = \{x_j(\varepsilon)\}$, $j = 1, \dots, n$, the non-zero elements of which are equal to the corresponding non-zero elements of $x_J(\varepsilon) = \{x_j(\varepsilon)\}$, $j \in J$ (with $x_J(\varepsilon)$ being as in (11)) satisfies the inclusion

$$x(\varepsilon') \in \theta^*(\varepsilon').$$

Since (due to Assumption (H_0^*)) (12) is satisfied, $x_J(\varepsilon)$ allows the expansion (13), and hence

$$x(\varepsilon) = \sum_{l=0}^{\infty} \varepsilon^l x^l, \quad \forall \varepsilon \in (0, \gamma), \quad (14)$$

where non-basic components ($j \notin J$) are equal to zero in both left and right hand sides. From (14) it follows, in particular, that

$$\lim_{\varepsilon \rightarrow 0} x(\varepsilon) = x^0 \quad (15)$$

By substituting (14) into the constraints of the perturbed problem (1), one can readily verify that $(x^0, x^1) \in \Theta_1$. Hence $x^0 \in \theta_1$.

The argument above proves that any partial limit (cluster point) of any basic optimal solution of the problem is contained in θ_1 . Since any element of $\theta^*(\varepsilon)$ can be presented as a convex combination of the optimal basic solutions and since θ_1 is convex, this proves the validity of (7), which, in turn, implies (8).

Let us now establish the validity of the second part of the theorem. Let $x^0 \in \theta_1^*$ and let x^1 be such that $(x^0, x^1) \in \Theta_1$. Define $(x^0(\delta), x^1(\delta))$ by the equation

$$x^0(\delta) \stackrel{\text{def}}{=} (1 - \delta)x^0 + \delta \hat{x}^0, \quad x^1(\delta) \stackrel{\text{def}}{=} (1 - \delta)x^1 + \delta \hat{x}^1, \quad \delta \in (0, 1), \quad (16)$$

where (\hat{x}^0, \hat{x}^1) are as in the ES-1 condition. Note that $(x^0(\delta), x^1(\delta)) \in \Theta_1$ (due to convexity of Θ_1) and also that

$$x_j^0(\delta) \geq \delta \hat{x}_j^0 \geq \delta a \quad \forall j \in J_0, \quad x_j^1(\delta) \geq \delta \hat{x}_j^1 \geq \delta a \quad \forall j \notin J_0, \quad (17)$$

where

$$a \stackrel{\text{def}}{=} \min\left\{ \min_{j' \in J_0} \hat{x}_{j'}^0, \min_{j' \notin J_0} \hat{x}_{j'}^1 \right\} > 0. \quad (18)$$

In the proof of Theorem 2.1 in [3], it has been established that there exists $x^2(\delta, \varepsilon)$ such that

$$x(\delta, \varepsilon) \stackrel{\text{def}}{=} x^0(\delta) + \varepsilon x^1(\delta) + \varepsilon^2 x^2(\delta, \varepsilon) \in \theta(\varepsilon) \quad \forall \delta \in [c_1 \varepsilon, 1], \quad \forall \varepsilon \in (0, \gamma) \quad (19)$$

and such that

$$\|x^2(\delta, \varepsilon)\| \leq c_2 \quad \forall \delta \in (0, 1), \quad \forall \varepsilon \in (0, \gamma), \quad (20)$$

where c_1 and c_2 are sufficiently large constants. Take $\delta(\varepsilon) \stackrel{\text{def}}{=} c_1 \varepsilon$. Then

$$\tilde{x}(\varepsilon) \stackrel{\text{def}}{=} x(\delta(\varepsilon), \varepsilon) \in \theta(\varepsilon) \quad \forall \varepsilon \in (0, \gamma), \quad (21)$$

and

$$\lim_{\varepsilon \rightarrow 0} \tilde{x}(\varepsilon) = x^0. \quad (22)$$

Since

$$\langle c^{(0)} + \varepsilon c^{(1)}, \tilde{x}(\varepsilon) \rangle \leq F^*(\varepsilon),$$

from (22) it follows that

$$F_1^* = \langle c^{(0)}, x^0 \rangle \leq \liminf_{\varepsilon \rightarrow 0} F^*(\varepsilon) \quad (23)$$

(the equality being due to the fact that x^0 was chosen to be an optimal solution of (4)). The validity of (23) and (8) implies the validity of (10). The latter along with (7) imply (9). Finally, the fact that any optimal solution x^0 of the problem (4) is limiting optimal in the perturbed problem (1) follows from (22) and from that

$$\lim_{\varepsilon \rightarrow 0} \langle c^{(0)} + \varepsilon c^{(1)}, \tilde{x}(\varepsilon) \rangle = \langle c^{(0)}, x^0 \rangle = F_1^* = \lim_{\varepsilon \rightarrow 0} F^*(\varepsilon).$$

□

Instead of problem (4), it may be more convenient to deal with the following problem

$$\begin{aligned} & \max \langle c^{(0)}, x^0 \rangle \\ & \text{s. t. } A^{(0)}x^0 = b^{(0)}, \\ & A^{(0)}x^1 + A^{(1)}x^0 = b^{(1)}, \\ & x^0 \geq 0, \\ & x^1 \geq 0, \end{aligned} \tag{24}$$

the statement of which does not involve the set J_0 . Let us give a sufficient condition, under which the problems (4) and (24) are equivalent in the sense of Definition 3 introduced below (the latter makes use of the fact that the objective function in (4) and (24) do not explicitly depend on x^1).

Definition 3. We will say that the problems (4) and (24) are equivalent when the sets θ_1 and $\tilde{\theta}_1$,

$$\tilde{\theta}_1 \stackrel{\text{def}}{=} \{x^0 \in \theta(0) : \text{there exists } x^1 \in \mathbb{R}^n \text{ such that } A^{(0)}x^1 + A^{(1)}x^0 = b^{(1)}, x^1 \geq 0\}, \tag{25}$$

coincide.

A sufficient condition for problems (4) and (24) to be equivalent is provided by the following result.

Proposition 1. Let J_0 be as in (3). If there exists $\alpha := \{\alpha_j\}_{j \in J_0}$ such that

$$A_{J_0}^{(0)}\alpha = 0, \text{ with } \alpha_j > 0 \quad \forall j \in J_0. \tag{26}$$

Then problems (4) and (24) are equivalent.

Proof. We must show that $\tilde{\theta}_1 = \theta_1$. The definitions readily imply that $\tilde{\theta}_1 \subset \theta_1$. Let us prove the opposite inclusion. Take $x^0 \in \theta_1$. From the definition of θ_1 it follows that $x^0 \in \theta(0)$. From this definition it also follows that there exists $x^1 = (x_j^1) \in \mathbb{R}^n$ such that

$$x_j^1 \geq 0, \forall j \notin J_0, \quad A^{(0)}x^1 + A^{(1)}x^0 = b^{(1)}. \tag{27}$$

If $x^1 \geq 0$, then by definition $x^0 \in \tilde{\theta}_1$. Otherwise, there exist a component (or components) of x^1 such that $x_j^1 < 0$ for $j \in J_0$. In this case, we can take

$$t > \max_{j \in J_0, x_j^1 < 0} \{-x_j^1/\alpha_j\} > 0,$$

where $\{\alpha_j\}_{j \in J_0}$ are as in (26). Define $\hat{x} \in \mathbb{R}^n$ as

$$\hat{x}_j := \begin{cases} t\alpha_j + x_j^1 & \text{if } j \in J_0, \\ x_j^1 & \text{if } j \notin J_0, \end{cases}$$

The definition of t ensures $\hat{x} \geq 0$. Using (26) and (27), we also have

$$\begin{aligned} A^{(0)}\hat{x} + A^{(1)}x^0 &= A_{J_0}^{(0)}(t\alpha + x_{J_0}^1) + A_{J_0^C}^{(0)}[x^1]_{J_0^C} + A^{(1)}x^0 \\ &= tA_{J_0}^{(0)}\alpha + A^{(0)}x^1 + A^{(1)}x^0 = b^{(1)}, \end{aligned}$$

where we used the notation $J_0^C := \{i : i \notin J_0\}$. The above expression implies that $x^0 \in \tilde{\theta}_1$, because we found a vector $\hat{x} \in \mathbb{R}^n$ such that $\hat{x} \geq 0$ and $A^{(0)}\hat{x} + A^{(1)}x^0 = b^{(1)}$. \square

Corollary 1. *If $b^{(0)} = 0$, then problems (4) and (24) are equivalent.*

Proof. By the very definition of J_0 , there exists $x_{J_0} > 0$ such that

$$A_{J_0}^{(0)}x_{J_0} = b^{(0)} = 0.$$

Thus, the role of α is played by x_{J_0} in the present case. \square

3. Application to Markov Decision Processes

Let us consider a discrete-time Markov Decision process (also called a Controlled Markov Chain) with a finite state space $\mathbb{X} = \{1, \dots, N\}$ and a finite action space $\mathbb{A}(i) = \{1, \dots, m_i\}$ for each state $i \in \mathbb{X}$. At any time point t the system is in one of the states $i \in \mathbb{X}$ and the controller or “decision-maker” chooses an action $a \in \mathbb{A}(i)$; as a result the following occur: (a) the controller gains an immediate reward r_{ia} , and (b) the process moves to a state $j \in \mathbb{X}$ with transition probability p_{iaj} , where $p_{iaj} \geq 0$ and $\sum_{j \in \mathbb{X}} p_{iaj} = 1$.

A decision rule π_t at time t is a function which assigns a probability to the event that any particular action a is taken at time t . In general, π_t may depend on all history $h_t = (i_0, a_0, i_1, a_1, \dots, a_{t-1}, i_t)$ up to time t . The distribution $\pi_t(a_t|h_t)$ defines the probability of selecting the action a_t at time t given the history h_t .

A control (or policy) is a sequence of decision rules $\pi = (\pi_0, \pi_1, \dots, \pi_t, \dots)$. A policy π is called Markov if $\pi_t(\cdot|h_t) = \pi_t(\cdot|i_t)$. If $\pi_t(\cdot|i) = \pi_{t'}(\cdot|i)$ for all $t, t' \in \mathbb{N}$ then the Markov policy π is called stationary. It is defined by a distribution π_{ia} , where π_{ia} is the probability of choosing action a when the system is in state i . Furthermore, a deterministic policy π is a stationary policy whose single decision rule is nonrandomized. It can be defined by the function $f(i) = a, a \in \mathbb{A}(i)$.

Let U, U^S and U^D denote the sets of all policies, all stationary policies and all deterministic policies, respectively. It is known that, in many contexts, there is no loss of generality in restricting consideration to stationary or even deterministic policies (see e.g., [21]).

For any stationary policy $\pi \in U^S$ we can define the corresponding transition matrix $P(\pi) = \{p_{ij}(\pi)\}_{i,j=1}^N$ and the reward vector $r(\pi) = \{r_i(\pi)\}_{i=1}^N$

$$p_{ij}(\pi) := \sum_{a \in \mathbb{A}(i)} p_{iaj} \pi_{ia}, \quad r_i(\pi) := \sum_{a \in \mathbb{A}(i)} r_{ia} \pi_{ia}.$$

The expected average reward $g_i(\pi)$ and the expected discounted reward $v_i^\alpha(\pi)$, associated with policy π , can be expressed as follows:

$$g_i(\pi) := \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T [P^{t-1}(\pi)r(\pi)]_i$$

and

$$v_i^\alpha(\pi) := (1 - \alpha) \sum_{t=1}^{\infty} \alpha^{t-1} [P^{t-1}(\pi)r(\pi)]_i = (1 - \alpha) [(I - \alpha P(\pi))^{-1}r(\pi)]_i,$$

respectively, where $i \in \mathbb{X}$ is an initial state and $\alpha \in (0, 1)$ is a discount factor.

Often an interest rate $\rho = (1 - \alpha)/\alpha$ is used instead of the discount factor. We note that the interest rate is close to zero when the discount factor is close to 1.

The following power series expansion, so-called Blackwell series expansion [6, 21], helps to establish a relation between discount optimality and average optimality

$$v_i^\alpha(\pi) = (1 - \alpha) \left[\frac{g_i(\pi)}{1 - \alpha} + h_i(\pi) + \dots \right] = g_i(\pi) + (1 - \alpha)h_i(\pi) + \dots, \quad (28)$$

where $h(\pi) = (I - P(\pi) + P^*(\pi))^{-1}(I - P^*(\pi))r(\pi)$ is a so-called bias vector. We note that often the expected discount reward vector is introduced without the factor $(1 - \alpha)$. In that case the power series (28) becomes a Laurent power series. However, the factor $(1 - \alpha)$ makes exposition of the results easier in the context of our singular perturbation approach.

We now introduce the discount optimality and the average optimality criteria in MDP optimization problem.

Definition 4. *The stationary policy π_* is called the discount optimal for fixed $\alpha \in (0, 1)$ if*

$$v_i^\alpha(\pi_*) \geq v_i^\alpha(\pi)$$

for each $i \in \mathbb{X}$ and all $\pi \in U^S$.

Definition 5. *The stationary policy π_* is called the average optimal if*

$$g_i(\pi_*) \geq g_i(\pi)$$

for each $i \in \mathbb{X}$ and all $\pi \in U^S$.

The power series (28) suggests another equivalent definition of the average optimality.

Definition 6. *The stationary policy π_* is called the average optimal if*

$$\lim_{\alpha \uparrow 1} [v_i^\alpha(\pi_*) - v_i^\alpha(\pi)] \geq 0$$

for each $i \in \mathbb{X}$ and all $\pi \in U^S$.

We note that this definition corresponds to the concept of limiting optimality in the context of perturbed linear programming.

In the case of discount optimality, the optimal value vector can be found as a solution of the following LP (see e.g., [1, 21]).

$$\min_{\gamma} \sum_j \gamma_j \tilde{v}_j \tag{29}$$

$$\text{subject to } \sum_j [\delta_{ij} - \alpha p_{iaj}] \tilde{v}_j \geq r_{ia}, \forall (i, a) \in \mathbb{X} \times \mathbb{A},$$

where $\gamma_j > 0$ and can be chosen as elements of some probability distribution. Without loss of generality, we may assume that the additional non-negativity constraints

$$\tilde{v}_j \geq 0, \quad \forall j,$$

are satisfied. The latter can be induced by adding a sufficiently large value $r_0 > 0$ to all immediate rewards r_{ia} . This transformation does not change the structure of optimal policies.

In the case of long-run average optimality, the optimal value vector can be found as a solution of another LP (see e.g., [1, 21]).

$$\min_{\gamma} \sum_j \gamma_j \tilde{v}_j \tag{30}$$

$$\text{subject to } \sum_j [\delta_{ij} - p_{iaj}] \tilde{v}_j \geq 0 \quad \forall (i, a) \in \mathbb{X} \times \mathbb{A},$$

$$\tilde{v}_i + \sum_j [\delta_{ij} - p_{iaj}] \tilde{u}_j \geq r_{ia} \quad \forall (i, a) \in \mathbb{X} \times \mathbb{A}.$$

Again, by adding a sufficiently large value $r_0 > 0$ to all immediate rewards and noticing that

$$\sum_j \delta_{ij} = \sum_j p_{iaj} = 1 \quad \forall (i, a), \tag{31}$$

one may assume, without loss of generality, that the non-negativity constraints

$$\tilde{v}_j \geq 0, \quad \tilde{u}_j \geq 0 \quad \forall j,$$

are satisfied.

Our aim is to demonstrate that LP (30) for the long-run average MDP can be derived from LP (29) for the discounted MDP by the formal singular perturbation methods [3, 20].

Take $\varepsilon \stackrel{\text{def}}{=} (1 - \alpha)/\alpha$. By making a change of variables $v_j = \varepsilon/(1 + \varepsilon)\tilde{v}_j$, one can rewrite the LP problem (29) in the form

$$\begin{aligned} & \min_{\gamma} \sum_j \gamma_j v_j & (32) \\ \text{subject to} & \sum_j [(1 + \varepsilon)\delta_{ij} - p_{iaj}]v_j \geq \varepsilon r_{ia}, \forall (i, a) \in \mathbb{X} \times \mathbb{A}, \\ & v_j \geq 0, \quad j = 1, \dots, n. \end{aligned}$$

Since Theorem 1 is stated for the linear programs with equality constraints, let us introduce additional variables σ_{ia} to transform (32) to

$$\begin{aligned} & \min_{\gamma} \sum_j \gamma_j v_j & (33) \\ \text{subject to} & \sum_j [(1 + \varepsilon)\delta_{ij} - p_{iaj}]v_j - \sigma_{ia} = \varepsilon r_{ia} \quad \forall (i, a) \in \mathbb{X} \times \mathbb{A}, \\ & v_j \geq 0, \quad \sigma_{ia} \geq 0. \end{aligned}$$

Note that the linear program above is just a particular case of (1) with $A^{(0)} = \{\delta_{ij} - P_{iaj} \mid -I\}$, $A^{(1)} = \{\delta_{ij} \mid 0\}$, $b^{(0)} = \{0\}$, $b^{(1)} = \{r_{ia}\}$, and $c^{(0)} = \{\gamma_j\}$, $c^{(1)} = 0$. The problem (24) (which is equivalent to (4) due to the fact that $b^{(0)} = \{0\}$; see Corollary 1) can in this case be written as follows

$$\min_{\gamma} \sum_j \gamma_j v_j^0 \quad (34)$$

$$\text{subject to} \quad \sum_j [\delta_{ij} - p_{iaj}]v_j^0 - \sigma_{ia}^0 = 0 \quad \forall (i, a) \in \mathbb{X} \times \mathbb{A}, \quad (35)$$

$$v_i^0 + \sum_j [\delta_{ij} - p_{iaj}]v_j^1 - \sigma_{ia}^1 = r_{ia} \quad \forall (i, a) \in \mathbb{X} \times \mathbb{A}, \quad (36)$$

$$v_j^0 \geq 0, \quad \sigma_{ia}^0 \geq 0, \quad v_j^1 \geq 0, \quad \sigma_{ia}^1 \geq 0.$$

Note that this problem is equivalent to (30) (with v_j^0 and v_j^1 playing the roles of \tilde{v}_j and \tilde{u}_j respectively).

Theorem 2. *The Assumptions (H_0^*) , (H_1) and the ES-1 condition are satisfied and, hence, the problem (34) is limiting LP for the problem (33) in the sense that (9) and (10) are satisfied.*

Proof. Since we consider the discounted reward vector normalized by $1 - \alpha = \varepsilon(1 + \varepsilon)$ (see (28)), the optimal value of the problem (33) remains bounded as $\varepsilon \rightarrow 0$. Hence, since also γ_j are assumed to be positive, Assumption (H_0^*) is satisfied. Assumption (H_1) is obviously satisfied as well (as the matrices of constraints contain the identity matrix). Let us now prove that the ES-1 condition is satisfied. Denote by $J_{0,v}$ and $J_{0,\sigma}$ the sets of indices such that from the fact $v^0 = (v_j^0) \geq 0$ and $\sigma^0 = (\sigma_j^0) \geq 0$ satisfy (35) it follows that

$$v_j^0 = 0 \quad \forall j \notin J_{0,v}, \quad \sigma_j^0 = 0 \quad \forall j \notin J_{0,\sigma}.$$

To verify the ES-1 condition, one needs to show that there exist

$$\hat{v}^0 = (\hat{v}_j^0) \geq 0, \quad \hat{\sigma}^0 = (\hat{\sigma}_j^0) \geq 0, \quad \hat{v}^1 = (\hat{v}_j^1) \geq 0, \quad \hat{\sigma}^1 = (\hat{\sigma}_j^1) \geq 0 \quad (37)$$

that satisfy (35), (36) as well as the property that

$$\hat{v}_j^0 > 0 \quad \forall j \in J_{0,v}, \quad \hat{\sigma}_j^0 > 0 \quad \forall j \in J_{0,\sigma}, \quad \hat{v}_j^1 > 0 \quad \forall j \notin J_{0,v}, \quad \hat{\sigma}_j^1 > 0 \quad \forall j \notin J_{0,\sigma}. \quad (38)$$

Note that $J_{0,v}^c = \emptyset$, due to the fact that $\sum_j [\delta_{ij} - p_{iaj}]M = 0$ for any $M > 0$ and any pair (i, a) . Also, if $v^0 = (v_j^0) \geq 0$, $\sigma^0 = (\sigma_j^0) \geq 0$, $v^1 = (v_j^1) \geq 0$, $\sigma^1 = (\sigma_j^1) \geq 0$ satisfy (35), (36) and the inequalities $v_j^0 > 0 \quad \forall j$, $\sigma_j^0 > 0 \quad \forall j \in J_{0,\sigma}$, are valid, then \hat{v}^0 , $\hat{\sigma}^0$, \hat{v}^1 , $\hat{\sigma}^1$, with the components defined as follows

$$\hat{v}_j^0 \stackrel{\text{def}}{=} v_j^0 + M \quad \forall j, \quad \hat{\sigma}_j^0 \stackrel{\text{def}}{=} \sigma_j^0 \quad \forall j, \quad \hat{v}_j^1 \stackrel{\text{def}}{=} v_j^1 \quad \forall j, \quad \hat{\sigma}_j^1 \stackrel{\text{def}}{=} \sigma_j^1 + M \quad \forall j, \quad (39)$$

will satisfy (35), (36) and (38), provided that M is chosen large enough. This completes the proof. \square

4. Acknowledgements

The work of K. Avrachenkov was partially supported by ARC Discovery Grant DP120100532 and EU Project Congas FP7-ICT-2011-8-317672; the work of J. Filar was partially supported by the ARC grant DP150100618; and the work of V. Gaitsgory was partially supported by the ARC Discovery Grants DP130104432, DP120100532 and DP150100618.

- [1] Altman, E., 1999, *Constrained Markov decision processes*, CRC Press.
- [2] Altman, E., Avrachenkov, K.E., and Filar, J.A., 1999, "Asymptotic linear programming and policy improvement for singularly perturbed Markov decision processes", *Mathematical Methods of Operations Research*, v.49(1), pp.97-109.
- [3] Avrachenkov, K., Burachik, R.S., Filar J.A., and Gaitsgory, V., 2012, "Constraint augmentation in pseudo-singularly perturbed linear programs", *Mathematical Programming, Ser. A*, v.132(1-2), pp.179-208.
- [4] Avrachenkov, K.E., Filar, J.A., and Howlett, P.G., 2013, *Analytic perturbation theory and its applications*, SIAM.
- [5] Avrachenkov, K.E., Haviv, M., Howlett, P.G., 2001, "Inversion of analytic matrix functions that are singular at the origin", *SIAM Journal on Matrix Analysis and Applications*, v.22, no.4, pp.1175-1189.
- [6] Blackwell, D., 1962, "Discrete dynamic programming", *The Annals of Mathematical Statistics*, pp.719-726.
- [7] Conway, J. B., 1973, *Functions of one Complex Variable*, Springer-Verlag, Berlin.
- [8] D'Epenoux, F., 1960, "Sur un probleme de production et de stockage dans l'aléatoire.", *Revue Française de Recherche Opérationnelle*, v.14, pp.3-16.
- [9] De Ghellinck, G., 1960, "Les problemes de decisions sequentielles", *Cahiers du Centre d'Etudes de Recherche Opérationnelle*, v.2(2), pp.161-179.
- [10] Filar, J.A., Altman, E., Avrachenkov, K.E., 2002, "An asymptotic simplex method for singularly perturbed linear programs", *Operations Research Letters*, v.30, no.5, pp.295-307.
- [11] Hordijk, A., and Kallenberg, L.C.M., 1979, "Linear programming and Markov decision chains", *Management Science*, v.25(4), pp.352-362.
- [12] Hordijk, A., and Kallenberg, L.C.M., 1984, "Constrained undiscounted stochastic dynamic programming", *Mathematics of Operations Research*, v.9(2), pp.276-289.
- [13] Hordijk, A., Dekker, R., and Kallenberg, L.C.M., 1985, "Sensitivity-analysis in discounted Markovian decision problems", *OR Spektrum*, v.7(3), pp.143-151.

- [14] Jeroslow, R.G., 1973, "Asymptotic linear programming". *Operations Research*, v.21(5), pp.1128-1141.
- [15] Kallenberg, L.C.M., 1983, *Linear programming and finite Markovian control problems*. MC Tracts, v.148, pp.1-245.
- [16] Kallenberg, L., 2002, "Finite state and action MDPs", in *Handbook of Markov Decision Processes: Methods and applications*, (eds.) E.A. Feinberg and A. Shwartz, Kluwer.
- [17] Mangasarian, O., 1994, *Nonlinear Programming*. Classics in Applied Mathematics, Society for Industrial and Applied Mathematics-SIAM Publishers, Philadelphia.
- [18] Manne, A.S., 1960, "Linear programming and sequential decisions", *Management Science*, v.6(3), pp.259-267.
- [19] Murty, K.G., 1983, *Linear programming*, John Wiley & Sons.
- [20] Pervozvanskii, A.A., Gaitsgori, V.G., 1988, *Theory of suboptimal decisions: Decomposition and aggregation*, Kluwer, Dordrecht.
- [21] Puterman, M. L., 1994, *Markov Decision Processes: Discrete Stochastic Dynamic Programming*, John Wiley and Sons, New York.
- [22] Rockafellar, R.T., Wets, R.J.-B. 1998, *Variational Analysis*. Springer, Berlin.
- [23] Ross K.W., Varadarajan, R., "Multichain Markov decision processes with a sample path constraint: A decomposition approach", *Math. Oper. Res.*, v.16, no.1, pp.195-207, 1991.
- [24] Schrijver, A., 1998, *Theory of Linear and Integer Programming*, Wiley.
- [25] Veinott, A.F., 1966, "On finding optimal policies in discrete dynamic programming with no discounting", *The Annals of Mathematical Statistics*, pp.1284-1294.
- [26] Veinott, A.F., 1969, "Discrete dynamic programming with sensitive discount optimality criteria", *The Annals of Mathematical Statistics*, pp.1635-1660.