

Algorithms for Singularly Perturbed Limiting Average Markov Control Problems¹

Mohammed Abbad and Jerzy A. Filar
Department of Mathematics and Statistics
University of Maryland at Baltimore County
Baltimore, Maryland

Tomasz R. Bielecki
Department of Mathematics
University of Kansas
Lawrence, Kansas

1. Introduction

In this paper we consider a singularly perturbed Markov Decision Process with the limiting average cost criterion. We assume that the underlying process is composed of n separate irreducible processes, and that the small perturbation is such that it "unites" these processes into a single irreducible process. This structure corresponds to the Markov chains admitting "strong and weak interactions" that arises in many applications, and was studied by a number of authors (e.g., see Delebecque and Quadrat [5], Phillips and Kokotovic [13], Coderch et al. [3], Kokotovic [10], Schweitzer [15] and [16], Rohlicek and Willsky [14], and Aldaheri and Khalil [1]). Our results can also be viewed as a continuation of a line of research initiated by Schweitzer [15] in 1968.

In Section 2 we introduce the formulation and some results given by Bielecki and Filar [2] of the underlying control problem for the singularly perturbed MDP; the so-called "limit Markov Control Problem" (limit MCP). In particular these authors proved that an optimal solution to the perturbed MDP can be approximated by an optimal solution of the limit MCP for sufficiently small perturbation.

In Section 3 we demonstrate that the above limit Markov Control Problem can be solved by a suitably constructed linear program.

In Section 4 we construct an algorithm for solving the limit Markov Control Problem based on the policy improvement method. Recently we learned that this algorithm is similar to one given by Pervozvanskii and Gaitgori [12]. However, these authors did not explicitly consider the limit Markov Control Problem, and worked only in the smaller class of deterministic strategies.

2. Definitions and Preliminaries

A discrete Markovian Decision Process (MDP, for short) is observed at time points $t = 0, 1, 2, \dots$. The state space is denoted by $S = \{1, 2, \dots, N\}$. With each state $s \in S$ we associate a finite action set $A(s) = \{1, 2, \dots, m_s\}$.¹ At any time point t the system is in one of the states s and the controller chooses an action $a \in A(s)$; as a result the following occur: (i) an immediate reward $r(s, a)$ is accrued, and (ii) the process moves to a state $s' \in S$ with transition probability $p(s' | s, a)$, where $p(s' | s, a) \geq 0$ and $\sum_{s' \in S} p(s' | s, a) = 1$. Henceforth, such an MDP will be synonymous with the four-tuple:

$$\Gamma = \langle S, \{A(s); s \in S\}, \{r(s, a); (s, a) \in S \times A(s)\}, \{p(s' | s, a); (s, a, s') \in S \times A(s) \times S\} \rangle.$$

While a general control strategy in Γ may depend on the complete state-action histories of the process, in this paper we shall concern ourselves only with the class Π of all stationary strategies. A stationary strategy $\pi \in \Pi$ is the vector:

$$\pi = \{\pi(s, a) | (s, a) \in S \times A(s)\},$$

¹Note that action $i \in A(s)$ may not be the same as action $i \in A(s')$ if $s \neq s'$. This simplification of notation should not cause ambiguity.

²We are indebted to Marc Teboulle for pointing out reference [10] to us.

where $\pi(s, a)$ is the probability that the controller chooses action $a \in A(s)$ in state s whenever that state is visited, of course, $\sum_{a \in A(s)} \pi(s, a) = 1$ for all s . A strategy $\pi \in \Pi$ will be called deterministic if $\pi(s, a) \in \{0, 1\}$ for all $(s, a) \in S \times A(s)$.

With every $\pi \in \Pi$ we shall associate the following quantities:

$r(\pi) = (r_1(\pi), \dots, r_N(\pi))^T$, the vector of single stage expected rewards in which $r_s(\pi) := \sum_{a \in A(s)} r(s, a)\pi(s, a)$ for each $s \in S$; a Markov matrix $P(\pi) = (p_{ss'}(\pi))_{s, s'=1}^N$, where $p_{ss'}(\pi) := \sum_{a \in A(s)} p(s' | s, a)\pi(s, a)$ for all $s, s' \in S$; the generator of the corresponding Markov Chain, namely, the matrix $G(\pi) := P(\pi) - I$; the corresponding Cesaro-limit matrix (sometimes called the ergodic projection at infinity) defined by:

$$P^*(\pi) = (p_{ss'}^*(\pi))_{s, s'=1}^N := \lim_{t \rightarrow \infty} \frac{1}{t+1} \sum_{k=0}^t P^k(\pi),$$

where $P^0(\pi) := I_N$, an $N \times N$ identity matrix; the overall performance index resulting from the use of π , namely, $J(s, \pi) := [P^*(\pi)r(\pi)]_s$ for each initial state $s \in S$.²

The "classical" limiting average Markov Decision problem is the optimization problem: Find $\pi^0 \in \Pi$ such that

$$J(s, \pi^0) = \max_{\pi} J(s, \pi) \quad \text{for all } s \in S.$$

A strategy π^0 satisfying the above will be called optimal. It is well known that there always exists an optimal deterministic policy and there is a number of finite algorithms for its computation (e.g., Denardo [6], Derman [7], Kallenberg [8]).

In this paper we shall assume that:

$$(A1) \quad S = \cup_{i=1}^n S_i \text{ where } S_i \cap S_j = \emptyset \text{ if } i \neq j, n > 1, \text{ card } S_i = n_i, n_1 + \dots + n_n = N,$$

and

$$(A2) \quad p(s' | s, a) = 0 \text{ whenever } s \in S_i \text{ and } s' \in S_j, i \neq j.$$

Consequently we can think of Γ as being the "union" of n smaller MDP's Γ_i , defined on the state space S_i , for each $i = 1, 2, \dots, n$, respectively. Note that if Π_i is the space of stationary strategies in Γ_i , then a strategy $\pi \in \Pi$ in Γ can be written in the natural way as $\pi = (\pi^1, \pi^2, \dots, \pi^n)$, where $\pi^i \in \Pi_i$. The probability transition matrix in Γ_i corresponding to π^i is, of course, defined by: $P_i(\pi^i) := (p_{s, s'}(\pi^i))_{s, s' \in S_i}$, and the generator $G_i(\pi^i)$ and the Cesaro-limit $P_i^*(\pi^i)$ matrices can be defined in a manner analogous to that in the original process Γ . In addition, we assume that:

(A3) For every $i = 1, 2, \dots, n$ and for all $\pi^i \in \Pi_i$ the matrix $P_i(\pi^i)$ is an irreducible matrix.

In view of (A3), $P_i(\pi^i)$ is a matrix with identical rows. We shall denote any row of $P_i(\pi^i)$ by $p_i(\pi^i)$.

Remark 2.1 Note that for all $\pi \in \Pi$ we have the following representation of $P(\pi)$:

$$P(\pi) = EM(\pi)$$

where E is an $N \times n$ matrix with entries:

$$e_{sj} = \begin{cases} 1 & \text{if } \sum_{k=1}^{j-1} n_k < s \leq \sum_{k=1}^j n_k \\ 0 & \text{otherwise} \end{cases}$$

for $s = 1, 2, \dots, N$ and $j = 1, 2, \dots, n$, and $M(\pi)$ is an $n \times N$ matrix with entries:

$$m_{js}(\pi) = \begin{cases} [p_j(\pi^j)]_s & \text{if } \sum_{k=1}^{j-1} n_k < s \leq \sum_{k=1}^j n_k \\ 0 & \text{otherwise} \end{cases}$$

for $j = 1, 2, \dots, n$ and $s = 1, 2, \dots, N$. Of course we set $\sum_{k=1}^0 n_k := 0$. Note also that from the above definitions we conclude that:

$$M(\pi)E = I_{n \times n}.$$

We shall now consider the situation where the transition probabilities of Γ are perturbed slightly. Towards this goal we shall define the *disturbance law* as the set:

$$D = \{d(s' | s, a) \mid (s, a, s') \in S \times A(s) \times S\}$$

where the elements of D satisfy $\sum_{s' \in S} d(s' | s, a) = 0$ for all $(s, a) \in S \times A(s)$, $-1 \leq d(s | s, a) \leq 0$, $d(s' | s, a) \geq 0$, and $s \neq s'$.

Now, with every $\pi \in \Pi$ we can associate a *perturbation generator matrix* $D(\pi) = (d_{ss'}(\pi))_{s, s'=1}^N$ where $d_{ss'}(\pi) = \sum_{a \in A(s)} d(s' | s, a) \pi(s, a)$. We shall also require that there exists $\epsilon_0 > 0$ such that for every $\pi \in \Pi$,

$$G_\epsilon(\pi) := G(\pi) + \epsilon D(\pi) \quad (2.1)$$

is a generator of a Markov chain for any $\epsilon \leq \epsilon_0$.

We shall consider a family of *perturbed processes* Γ_ϵ for $\epsilon \in [0, \epsilon_0]$ that differ from the original MDP Γ only in the transition law, namely, in Γ_ϵ for every $(s, a, s') \in S \times A(s) \times S$ we have:

$$p_\epsilon(s' | s, a) := p(s' | s, a) + \epsilon d(s' | s, a). \quad (2.2)$$

By (2.1) and (2.2) we have that every $\pi \in \Pi$ induces in the perturbed process Γ_ϵ the Markov Chain with the probability transition matrix:

$$P_\epsilon(\pi) = G_\epsilon(\pi) + I_N. \quad (2.3)$$

The most important structural assumption of this paper is stated below.

(SPA) *Singular Perturbation Assumption:*

For every $\pi \in \Pi$ and every $\epsilon \in (0, \epsilon_0]$, $P_\epsilon(\pi)$ is an irreducible matrix.

Remark 2.2 Note that as a result of (SPA) and (A1)–(A3) the rank of $P_\epsilon(\pi)$ is 1, which is strictly less than n , the rank of $P^*(\pi)$. Consequently our perturbation is indeed a *singular perturbation* in the sense of Delebecque [4].

Under assumption SPA the MDP Γ_ϵ , $\epsilon \in [0, \epsilon_0]$, defined as: $\Gamma_\epsilon = (S, \{A(s); s \in S\}, \{r(s, a); (s, a) \in S \times A(s)\}, \{p_\epsilon(s' | s, a); (s, a, s') \in S \times A(s) \times S\})$ is called the *singularly perturbed MDP*.

Denoting by $P_\epsilon^*(\pi)$ the ergodic projection at infinity corresponding to $P_\epsilon(\pi)$ we define the overall performance index resulting from the use of $\pi \in \Pi$ by:

$$\mathcal{J}_\epsilon(s, \pi) = [P_\epsilon^*(\pi)r(\pi)]_s, \quad s \in S.$$

The optimal value function \mathcal{J}_ϵ corresponding to Γ_ϵ is given by:

$$\mathcal{J}_\epsilon(s) = \max_{\pi \in \Pi} \mathcal{J}_\epsilon(s, \pi), \quad s \in S.$$

Note that by SPA, $P_\epsilon(\pi)$ is irreducible and hence there is no loss of generality in omitting the argument s in the above, and writing simply $\mathcal{J}_\epsilon(\pi)$, and \mathcal{J}_ϵ respectively.

We now recall from [2] the analysis of the limiting behavior of \mathcal{J}_ϵ as ϵ goes to zero and the validity of the so called *limit control principle* in the present framework, that is, an optimal strategy in a certain limit optimization problem, to be defined in sequel, is δ -optimal in Γ_ϵ for any $\delta > 0$ and for ϵ sufficiently small.

The problems described above pertain to the class of MDP's with singularly perturbed transition kernels or more generally to the class of control problems for Markov Processes with singularly perturbed generators. Related problems for discounted cost criteria have been addressed in Delebecque & Quadrat [5] and Delebecque [4] in the case of discrete time Markov chains and in Philips & Kokotović [13] for both continuous and discrete time Markov chains.

For each $\pi \in \Pi$ let us define the $n \times n$ matrix $\bar{B}(\pi)$ by:

$$\bar{B}(\pi) \stackrel{\text{def}}{=} M(\pi)D(\pi)E.$$

We note that by Remark 2.1 $\bar{B}(\pi)$ is a generator of an "aggregated" Markov chain on a state space $\bar{S} \stackrel{\text{def}}{=} \{1, 2, \dots, n\}$.

Remark 2.3 By A1–A3 and SPA the operator $\bar{B}(\pi)$ defines an irreducible Markov chain on \bar{S} . This can be verified by direct inspection, (see also remark 2 p. 338 in Delebecque [4]).

Now let $\bar{P}^*(\pi)$ denote the ergodic projection at infinity corresponding to $\bar{B}(\pi)$, for each $\pi \in \Pi$. Define, for $\pi \in \Pi$, an $N \times N$ matrix:

$$\hat{P}^*(\pi) = E\bar{P}^*(\pi)M(\pi). \quad (2.4)$$

From Delebecque [4] Theorem 3, and the results of Kato [9], chapter II, Bielecki and Filar [2] obtained an asymptotic result stated in Lemma 2.1.

Lemma 2.1 Under the assumptions A1–A3 and SPA

$$\lim_{\epsilon \rightarrow 0} \max_{\pi \in \Pi} \|P_\epsilon^*(\pi) - \hat{P}^*(\pi)\| = 0, \quad (2.5)$$

where $\|\cdot\|$ denotes any matrix norm.

Using this Lemma, the following Corollary was derived in [2].

Corollary 2.1 *Assume A1-A3 and SPA. Then for all $s \in S_j$, $j = 1, 2, \dots, n$.*

$$\lim_{\epsilon \rightarrow 0} |\max_{\pi \in \Pi} [\bar{P}^*(\pi)\bar{r}(\pi)]_j - \max_{\pi \in \Pi} [P_\epsilon^*(\pi)r(\pi)]_s| = 0 \quad (2.6)$$

where $\bar{r}(\pi) := M(\pi)r(\pi)$.

Remark 2.4 *The optimization problem:*

$$\max_{\pi \in \Pi} [\bar{P}^*(\pi)r(\pi)]_s, \quad s \in S \quad (L)$$

is called the the limit Markov Control Problem.

The optimization problem:

$$\max_{\pi \in \Pi} [\bar{P}^*(\pi)\bar{r}(\pi)]_j, \quad j = 1, 2, \dots, n \quad (AL)$$

is called the aggregated limit Markov Control Problem.

Remark 2.5 *Note that in view of (2.4) we have:*

$$[\bar{P}^*(\pi)r(\pi)]_s = [\bar{P}^*(\pi)\bar{r}(\pi)]_j \text{ for } s \in S_j, \quad j = 1, 2, \dots, n, \quad \pi \in \Pi. \quad (2.9)$$

It follows that any maximizing strategy π^0 for (AL) is also a maximizing strategy for (L) and vice-versa. The existence of a maximizing strategy π^0 is clear. In view of (2.7) and the irreducibility properties it is also clear that π^0 does not depend on the initial state $s \in S$.

Using the results above, the following theorem was proved in [2]:

Theorem 2.1 (Limit Control Principle) *Assume A1-A3 and SPA. Let $\pi^0 \in \Pi$ be any maximizer in (AL). Then for all $\delta > 0$, there exists $\epsilon_\delta > 0$ such that for all $\epsilon < \epsilon_\delta$.*

$$\|P_\epsilon^*(\pi^0)r(\pi^0) - \mathcal{J}_\epsilon\| \leq \delta. \quad (2.8)$$

3. Linear Programming and the Limit Markov Control Problem

We have seen in Section 2 that the optimization problem:

$$\max_{\pi \in \Pi} [\bar{P}^*(\pi)r(\pi)]_s, \quad s \in S \quad (L)$$

is the natural problem to attempt to solve in the case of a singularly perturbed Markov Decision Process. We shall demonstrate that this problem can be converted to an equivalent problem in the space of the long-run state-action frequencies. Towards this goal we shall now consider one of the irreducible MDP's:

$\Gamma_i = \langle S_i, \{A(s); s \in S_i\}, \{r(s, a); (s, a) \in S_i \times A(s)\}, \{p(s' | s, a); (s, a, s') \in S_i \times A(s) \times S_i\} \rangle$ that were already introduced in Section 2, with $i = 1, 2, \dots, n$. The following results are well-known (e.g., see Kallenberg [8]). With each Γ_i we can associate a polyhedral set:

$$X_i = \{x^i = \{x_{sa}^i | (s, a) \in S_i \times A(s)\} \mid \sum_{s \in S_i} \sum_{a \in A(s)} x_{sa}^i = 1; \sum_{s \in S_i} \sum_{a \in A(s)} (\delta_{s,s'} - p(s' | s, a)) x_{sa}^i = 0, \quad s' \in S_i; x_{sa}^i \geq 0, \quad (s, a) \in S_i \times A(s)\};$$

and a bijective map $T : X_i \rightarrow \Pi_i$ such that $T(x^i) := \pi^i$ where

$$\pi^i(s, a) := x_{sa}^i / \sum_{a \in A(s)} x_{sa}^i; \quad (s, a) \in S_i \times A(s).$$

Now, if $\pi^i \in \Pi_i$ and $p_i^*(\pi^i)$ is the corresponding stationary distribution in Γ_i it can be shown that the inverse map is $T^{-1} : \Pi_i \rightarrow X_i$ such that $T^{-1}(\pi^i) := x^i$, where

$$x_{sa}^i := [p_i^*(\pi^i)]_s \pi^i(s, a); \quad (s, a) \in S_i \times A(s).$$

We shall now rewrite the objective function of (L) (with the help of (2.4) and structural assumptions) as:

$$\begin{aligned} \mathcal{J}(\pi) &:= [\bar{P}^*(\pi)r(\pi)]_s = [E\bar{P}^*(\pi)M(\pi)r(\pi)]_s \\ &= \sum_{i=1}^n [\bar{P}^*(\pi)]_i [p_i^*(\pi^i) \cdot r_i(\pi^i)], \end{aligned} \quad (3.1)$$

where $\bar{p}^*(\pi)$ is the stationary distribution vector of $\bar{P}(\pi) = \bar{B}(\pi) + I_{n \times n}$, $s \in S$, and $\pi = (\pi^1, \pi^2, \dots, \pi^n)$ as in Section 2.

It can be easily checked that with T as above we have:

$$[p_i^*(\pi^i)]_s = [p_i^*(T(x^i))]_s = \sum_{a \in A(s)} x_{sa}^i \quad (3.2)$$

for each $i = 1, 2, \dots, n$ and $s \in S_i$. Hence (3.1) - (3.2) yield:

$$\mathcal{J}(\pi) = \sum_{i=1}^n [\bar{P}^*(\pi)]_i \left(\sum_{s \in S_i} \sum_{a \in A(s)} r(s, a) x_{sa}^i \right) \quad (3.3)$$

for every $s \in S$, where $x^i = T^{-1}(\pi^i)$.

Recall that $\bar{p}^*(\pi)$ is the unique solution of the system of equations:

$$\begin{aligned} \alpha^T \bar{P}(\pi) &= \alpha^T \\ \sum_{i=1}^n \alpha_i &= 1, \end{aligned}$$

since $\bar{P}(\pi)$ is irreducible for all $\pi \in \Pi$. The first of these equations can, by definition of $\bar{B}(\pi)$, be expressed as:

$$\alpha^T M(\pi)(D(\pi) + I_N)E = \alpha^T \quad (3.4)$$

Equivalently, using the fact that $M(\pi)I_N E = I_n$, we can transform (3.4) to:

$$\alpha^T M(\pi)D(\pi)E = \mathbf{0}^T. \quad (3.5)$$

Note that $U(\pi) := M(\pi)D(\pi)E$ is an $n \times n$ matrix whose (i, j) -th entry is:

$$u_{ij}(\pi) := \sum_{s' \in S_j} \sum_{s \in S_i} [p_i^*(\pi^i)]_s d_{ss'}(\pi^i). \quad (3.6)$$

Now using (3.2) we see that with $x^i := T^{-1}(\pi^i)$ for each $i = 1, 2, \dots, n$, we have:

$$v_{ij}(x) := u_{ij}(\pi) = \sum_{s' \in S_j} \sum_{s \in S_i} \sum_{a \in A(s)} d(s' | s, a) x_{sa}^i, \quad (3.7)$$

where $x = (x^1, x^2, \dots, x^n)$.

Setting $V(x) = (v_{ij}(x))_{i,j=1}^n$, we are now led to consider the following nonlinear programming problem (NL):

$$\text{maximize } \sum_{i=1}^n \sum_{s \in S_i} \sum_{a \in A(s)} r(s, a) x_{sa}^i \alpha_i$$

Subject to:

- (i) $x^i = \{x_{sa}^i | (s, a) \in S_i \times A(s)\} \in X_i$, $i = 1, 2, \dots, n$,
- (ii) $\alpha_i \geq 0$, $i = 1, 2, \dots, n$, and $\sum_{i=1}^n \alpha_i = 1$,
- (iii) $\alpha^T V(x) = \mathbf{0}^T$.

Note that $\mathbf{x} = (x^1, \dots, x^n)$ and $\alpha = (\alpha_1, \dots, \alpha_n)$ are the variables in the problem (NL).

We can now state the following result [2]:

Theorem 3.1 Let $(\bar{x}, \bar{\alpha})$ be an optimal solution of the nonlinear program (NL). Define $\bar{\pi}^i := T(\bar{x}^i)$ for $i = 1, 2, \dots, n$ and $\bar{\pi} := (\bar{\pi}^1, \dots, \bar{\pi}^n)$, then $\bar{\pi}$ is optimal in the limit problem (L).

Remark 3.6 The importance of the above theorem stems from the fact that it converts the limit problem (L) into the problem (NL). In what follows, we shall demonstrate that an optimal solution in the program (NL) can be obtained from an optimal solution of an appropriate linear program that can be solved by efficient linear programming techniques.

Consider the following linear programming problem (P):

$$\text{maximize } \sum_{i=1}^n \sum_{s \in S_i} \sum_{a \in A(s)} r(s, a) z_{sa}^i$$

Subject to:

$$\sum_{s \in S_i} \sum_{a \in A(s)} (\delta_{ss'} - p(s' | s, a)) z_{sa}^i = 0; \quad s' \in S_i; \quad i = 1, \dots, n \quad (3.8)$$

$$\sum_{i=1}^n \sum_{s' \in S_i} \sum_{s \in S_i} \sum_{a \in A(s)} d(s' | s, a) z_{sa}^i = 0; \quad j = 1, 2, \dots, n \quad (3.9)$$

$$\sum_{i=1}^n \sum_{s \in S_i} \sum_{a \in A(s)} z_{sa}^i = 1 \quad (3.10)$$

$$z_{sa}^i \geq 0; \quad i = 1, 2, \dots, n; \quad s \in S_i; \quad a \in A(s) \quad (3.11)$$

Lemma 3.1 For any feasible solution \mathbf{z} of (P),

$$\sum_{s \in S_i} \sum_{a \in A(s)} z_{sa}^i > 0 \quad \text{for all } i \in \{1, 2, \dots, n\}$$

Proof: Define $F(\mathbf{z}) := \{i \in \{1, 2, \dots, n\} : \sum_{s \in S_i} \sum_{a \in A(s)} z_{sa}^i > 0\}$ and $\bar{F}(\mathbf{z}) := \{i \in \{1, 2, \dots, n\} : \sum_{s \in S_i} \sum_{a \in A(s)} z_{sa}^i = 0\}$.

We shall show that $\bar{F}(\mathbf{z}) = \emptyset$. Assume $\bar{F}(\mathbf{z}) \neq \emptyset$. For $i \in \bar{F}(\mathbf{z})$ set $\alpha_i := 0$ and take any strategy $\pi^i \in \Pi_i$ and define $\mathbf{x}^i := T^{-1}(\pi^i)$. For $i \in F(\mathbf{z})$ define:

$$\alpha_i := \sum_{s \in S_i} \sum_{a \in A(s)} z_{sa}^i \quad (3.12)$$

and

$$x_{sa}^i := \frac{z_{sa}^i}{\sum_{s \in S_i} \sum_{a \in A(s)} z_{sa}^i} \quad \text{for all } s \in S_i, a \in A(s). \quad (3.13)$$

Note that we have:

$$z_{sa}^i = x_{sa}^i \alpha_i \quad \text{for all } i \in \{1, 2, \dots, n\}, \quad s \in S_i, \quad a \in A(s). \quad (3.14)$$

From (3.8) and (3.13) we have: for all $i \in F(\mathbf{z})$

$$\sum_{s \in S_i} \sum_{a \in A(s)} (\delta_{ss'} - p(s' | s, a)) x_{sa}^i = 0 \quad \text{for all } s' \in S_i.$$

From (3.13) it follows that:

$$x_{sa}^i \geq 0 \quad \text{for all } s \in S_i, \quad a \in A(s), \quad \text{and} \quad \sum_{s \in S_i} \sum_{a \in A(s)} x_{sa}^i = 1$$

It follows that for all $i \in F(\mathbf{z})$, $\mathbf{x}^i \in X_i$.

Since T is bijective, there exists a strategy π^i in Γ_i such that $\mathbf{x}^i = T^{-1}(\pi^i)$.

Now (3.9) and (3.14) imply that for all $j = 1, 2, \dots, n$:

$$\begin{aligned} 0 &= \sum_{i=1}^n \sum_{s' \in S_j} \sum_{s \in S_i} \sum_{a \in A(s)} d(s' | s, a) x_{sa}^i \alpha_i \\ &= \sum_{i=1}^n \alpha_i \left[\sum_{s' \in S_j} \sum_{s \in S_i} \sum_{a \in A(s)} d(s' | s, a) [p_i^*(\pi_i)]_s \pi^i(s, a) \right] \\ &= \sum_{i=1}^n \alpha_i \left[\sum_{s' \in S_j} \sum_{s \in S_i} [p_i^*(\pi^i)]_s d_{ss'}(\pi^i) \right] = [\alpha^T M(\pi) D(\pi) E]_j, \end{aligned}$$

$$\text{where } \alpha^T = (\alpha_1, \alpha_2, \dots, \alpha_n) \quad \text{and } \pi = (\pi^1, \pi^2, \dots, \pi^n).$$

Thus $\alpha^T M(\pi) (D(\pi) + I_N) E = \alpha^T$ since $M(\pi) I_N E = I_n$, which is the same as $\alpha^T \bar{P}(\pi) = \alpha^T$.

From the definitions of α_i we have by (3.10):

$$\sum_{i=1}^n \alpha_i = \sum_{i=1}^n \sum_{s \in S_i} \sum_{a \in A(s)} z_{sa}^i = 1,$$

and by (3.11) $\alpha_i \geq 0$ for $i = 1, \dots, n$.

Now we have:

$$\begin{aligned} \alpha^T \bar{P}(\pi) &= \alpha^T \\ \sum_{i=1}^n \alpha_i &= 1; \quad \alpha_i \geq 0; \quad i = 1, \dots, n. \end{aligned}$$

But $\bar{P}(\pi)$ is irreducible, hence all α_i must be positive, and thus we have a contradiction. Therefore $\bar{F}(\mathbf{z})$ must be empty. \square

Theorem 3.2 If \bar{z} is an optimal solution of (P), then $(\bar{x}, \bar{\alpha})$ is an optimal solution of (NL). Where

$$\bar{x}_{sa}^i := \frac{\bar{z}_{sa}^i}{\sum_{s \in S_i} \sum_{a \in A(s)} \bar{z}_{sa}^i} \quad \text{and} \quad \bar{\alpha}_i := \sum_{s \in S_i} \sum_{a \in A(s)} \bar{z}_{sa}^i$$

for $i = 1, 2, \dots, n; s \in S_i; a \in A(s)$.

Proof: By Lemma 3.1, \bar{x} is well defined. Also, from the proof of Lemma 3.1, it follows that $(\bar{x}, \bar{\alpha})$ is a feasible solution in (NL). To prove that $(\bar{x}, \bar{\alpha})$ is optimal in (NL), let (\mathbf{x}, α) be any feasible solution in (NL) and define \mathbf{z} by $z_{sa}^i := x_{sa}^i \alpha_i; i = 1, \dots, n; s \in S_i; a \in A(s)$. Note that \mathbf{z} is feasible for (P), and we have:

$$\sum_{i=1}^n \sum_{s \in S_i} \sum_{a \in A(s)} r(s, a) x_{sa}^i \alpha_i = \sum_{i=1}^n \sum_{s \in S_i} \sum_{a \in A(s)} r(s, a) z_{sa}^i \leq$$

$$\sum_{i=1}^n \sum_{s \in S_i} \sum_{a \in A(s)} r(s, a) \bar{z}_{sa}^i = \sum_{i=1}^n \sum_{s \in S_i} \sum_{a \in A(s)} \bar{x}_{sa}^i \bar{\alpha}_i.$$

This proves the optimality of $(\bar{x}, \bar{\alpha})$. \square

Next we shall show that an optimal deterministic strategy for the limit Markov Control Problem (L) can be constructed from an extreme optimal solution of the linear programming problem (P).

Lemma 3.2 Let \bar{z} be an extreme feasible solution for (P), then for any $i \in \{1, 2, \dots, n\}$ and any $s \in S_i$ there is a unique $a \in A(s)$ such that $\bar{z}_{sa}^i > 0$.

Proof: Since in particular \bar{z} is a feasible solution for (P), it follows from the proof of Lemma 3.1 that for any $i \in \{1, 2, \dots, n\}$ there exists a strategy π^i such that:

$$\sum_{s \in S_i} \frac{z_{sa}^i}{\sum_{a \in A(s)} z_{sa}^i} = [T^{-1}(\pi^i)]_{sa} := [P_i^*(\pi^i)]_s \pi^i(s, a), \quad s \in S_i, a \in A(s).$$

Hence for all $i \in \{1, 2, \dots, n\}$, $s \in S_i$, there exists $a \in A(s)$: $\bar{z}_{sa}^i > 0$. Thus the number of positive elements in \bar{z} is at least $n_1 + n_2 + \dots + n_n := N$. On the other hand, since \bar{z} is an extreme feasible solution for (P) , the number of positive elements in \bar{z} is not greater than the rank k of the matrix corresponding to the linear program (P) . From that matrix we have: for any $i \in \{1, 2, \dots, n\}$, $\sum_{s' \in S_i} [\sum_{s \in S_i} \sum_{a \in A(s)} (\delta_{ss'} - p(s' | s, a)) z_{sa}^i] = \sum_{s \in S_i} \sum_{a \in A(s)} (1 - 1) z_{sa}^i = 0$, and $\sum_{j=1}^n [\sum_{s' \in S_j} \sum_{s \in S_i} \sum_{a \in A(s)} d(s' | s, a) z_{sa}^i] = \sum_{i=1}^n \sum_{s \in S_i} \sum_{a \in A(s)} [\sum_{j=1}^n \sum_{s' \in S_j} d(s' | s, a)] z_{sa}^i = 0$. Hence $k \leq (n_1 - 1) + (n_2 - 1) + \dots + (n_n - 1) + (n - 1) + 1 = n_1 + n_2 + \dots + n_n := N$. Therefore the number of positive elements in \bar{z} is exactly N , and we conclude that $\bar{z}_{sa}^i > 0$ for exactly one $a \in A(s)$ for any $i \in \{1, 2, \dots, n\}$ and any $s \in S_i$. \square

Theorem 3.3 *The limit Markov Control Problem (L) has an optimal deterministic strategy.*

Proof: It is clear that the problem (NL) is feasible. This implies that the problem (P) is feasible (if (x, α) is feasible for (NL) , then $z_{sa}^i = x_{sa}^i \alpha_i$ is feasible for (P)). Since the constraints in (P) define a bounded polyhedron, then (P) has an optimal extreme solution, say \bar{z} . From Theorem 3.2 the corresponding $(\bar{x}, \bar{\alpha})$ is an optimal solution for (NL) . From Theorem 3.1, $\bar{\pi} := (\bar{\pi}^1, \dots, \bar{\pi}^n)$ where $\bar{\pi}^i := T(\bar{x}^i)$ for $i = 1, 2, \dots, n$ is optimal in (L) . From Lemma 3.2, $\bar{\pi}$ must be deterministic. \square

Remark 3.7 *Note that from Lemma 3.2 and from the proof of Theorem 3.3, it follows that if $z := \{z_{sa}^i | s \in S_i; i = 1, \dots, n; a \in A(s)\}$ is an optimal extreme solution of the linear program (P) , then the policy defined by:*

$$f(s) = a; \quad s \in S_i; \quad i = 1, \dots, n \iff z_{sa}^i > 0$$

is optimal in the limit Markov control problem (L).

Remark 3.8 *In this section we proved that the limit problem (L) can be converted into the linear program (P). In view of the special structure of (P), an approach of the Wolfe-Dantzig decomposition method (e.g., see Murty [11]) is suggested. Hence the linear program (P) may not suffer from the curse of dimensionality. In the next section we shall construct another algorithm which does not suffer from the dimensionality of the problem.*

4. Aggregation-Disaggregation Algorithm.

Consider the Aggregated Markov Decision Process $\bar{\Gamma}$ defined as follows:

The State Space of $\bar{\Gamma}$: $\bar{S} := \{1, 2, \dots, n\}$

The Action Spaces of $\bar{\Gamma}$: $\bar{A}(i) := X_{s \in S_i} A(s)$ for each $i \in \bar{S}$

The Transition Law of $\bar{\Gamma}$: for all $(i, \mathbf{a}, j) \in \bar{S} \times \bar{A}(i) \times \bar{S}$,

$$q_{ij}(\mathbf{a}) := \begin{cases} 1 + \sum_{s' \in S_i} \sum_{s \in S_i} [P_i^*(\mathbf{a})]_s d(s' | s, a_s) & i = j \\ \sum_{s' \in S_i} \sum_{s \in S_i} [P_i^*(\mathbf{a})]_s d(s' | s, a_s) & i \neq j \end{cases} \quad (4.1)$$

The Rewards of $\bar{\Gamma}$: for all $i \in \bar{S}$, $\mathbf{a} \in \bar{A}(i)$

$$c(i, \mathbf{a}) := \sum_{s \in S_i} [P_i^*(\mathbf{a})]_s r(s, a_s), \quad (4.2)$$

where $\mathbf{a} = \{a_s | s \in S_i\}$.

Note that for any $i \in \bar{S}$, each action $\mathbf{a} \in \bar{A}(i)$ defines a deterministic strategy which maps $s \in S_i$ onto a_s . Thus in (4.1) and (4.2), $P_i^*(\mathbf{a})$ is well defined.

The validity of the Transition Law, namely:

$$\sum_{j=1}^n q_{ij}(\mathbf{a}) = 1, \quad i \in \bar{S}, \quad \mathbf{a} \in \bar{A}(i) \quad \text{and} \quad q_{ij}(\mathbf{a}) \geq 0, \quad (i, \mathbf{a}, j) \in \bar{S} \times \bar{A}(i) \times \bar{S}$$

can be checked by inspection using the assumptions made on the disturbance law D .

The classical limiting average Markov Decision problem for $\bar{\Gamma}$ is the optimization problem

$$\max_{\bar{\pi} \in \bar{\Pi}} [Q^*(\bar{\pi})c(\bar{\pi})]_i, \quad i \in \bar{S} \quad (AP)$$

where $\bar{\Pi}$ is the class of all stationary strategies in $\bar{\Gamma}$, $Q^*(\bar{\pi})$ is the Cesaro-limit matrix corresponding to the Markov matrix $Q(\bar{\pi})$ which is defined by $Q(\bar{\pi}) := (q_{ij}(\bar{\pi}))_{i,j=1}^n$ where $q_{ij}(\bar{\pi}) := \sum_{\mathbf{a} \in \bar{A}(i)} q_{ij}(\mathbf{a}) \bar{\pi}(i, \mathbf{a})$ for all $i, j \in \bar{S}$; and $c(\bar{\pi}) := (c_1(\bar{\pi}), \dots, c_n(\bar{\pi}))$ is the vector of single stage expected rewards in which $c_i(\bar{\pi}) := \sum_{\mathbf{a} \in \bar{A}(i)} c(i, \mathbf{a}) \bar{\pi}(i, \mathbf{a})$ for each $i \in \bar{S}$.

Note that there is a one to one correspondence between the deterministic strategies in $\bar{\Gamma}$ and Γ .

Let \bar{f} be a deterministic strategy in $\bar{\Gamma}$, the corresponding deterministic strategy f in Γ is defined by:

$$\text{for each } i \in \bar{S} \text{ and } s \in S_i, \quad f(s) := [\bar{f}(i)]_s.$$

Proposition 4.1 *Let \bar{f} be an optimal deterministic strategy for the problem (AP), then the corresponding f is an optimal deterministic strategy for the limit Markov Control Problem (L)*

Proof: Note that from the theory of Markov Decision Processes (e.g., see [6]-[8]) \bar{f} always exists. It can be verified by inspection that $Q(\bar{f}) = \bar{P}(f)$ since $\bar{P}(f) = I_n + M(f)D(f)E$. Thus $Q^*(\bar{f}) = \bar{P}^*(f)$ and $Q^*(\bar{f})c(\bar{f}) = \bar{P}^*(f)\bar{r}(f)$. Therefore f must be an optimal solution for the problem (AL), since the problem (AL) has an optimal deterministic strategy (Remark 2.5 and Theorem 3.3). Now from Remark 2.5 we conclude that f is an optimal strategy for the problem (L). \square

Remark 4.9 *In view of the fact that for any deterministic strategy \bar{f} in the process $\bar{\Gamma}$ we have $Q(\bar{f}) = \bar{P}(f)$ and $\bar{P}(f)$ is irreducible (Remark 2.3), thus $Q(\bar{f})$ is irreducible. This shows that $\bar{\Gamma}$ is an irreducible MDP and therefore it can be solved by using the simplified policy improvement algorithm (e.g., see Denardo [6], Derman [7], Kallenberg [8]).*

The "standard" policy improvement algorithm for the aggregated MDP $\bar{\Gamma}$ can be stated as follows:

Step 1: Select an arbitrary deterministic strategy \bar{f} .

Step 2: Solve the following linear system for the unknowns $\lambda, y_1, \dots, y_{n-1}$:

$$\lambda + y_i = c(i, \bar{f}(i)) + \sum_{j=1}^n q_{ij}(\bar{f}(i)) y_j; \quad i = 1, 2, \dots, n \quad (4.3)$$

where $y_n := 0$.

Step 3: Find a deterministic strategy \bar{g} that satisfies:

$$c(i, \bar{g}(i)) + \sum_{j=1}^n q_{ij}(\bar{g}(i)) y_j > c(i, \bar{f}(i)) + \sum_{j=1}^n q_{ij}(\bar{f}(i)) y_j \quad i = 1, 2, \dots, n. \quad (4.4)$$

If this is not possible for some $i \in \{1, 2, \dots, n\}$, then set $\bar{g}(i) = \bar{f}(i)$.

If $\bar{g} = \bar{f}$, STOP.

Otherwise $\bar{f} \leftarrow \bar{g}$ and go to Step 2.

Now, we shall show that for each $i = 1, 2, \dots, n$ the problem defined in (4.4) can be converted to one iteration of the policy improvement algorithm for an MDP defined by Γ_i except for the rewards which will be defined appropriately.

For every $i = 1, \dots, n$ and $\mathbf{a} \in \bar{A}(i)$ we have:

$$\begin{aligned} c(i, \mathbf{a}) + \sum_{j=1}^n q_{ij}(\mathbf{a}) y_j &= \sum_{s \in S_i} [P_i^*(\mathbf{a})]_s r(s, a_s) + \sum_{j=1}^n \left\{ \sum_{s' \in S_j} \right. \\ &\left. \sum_{s \in S_i} [P_i^*(\mathbf{a})]_s d(s' | s, a_s) \right\} y_j + \{1 + \sum_{s' \in S_i} \sum_{s \in S_i} [P_i^*(\mathbf{a})]_s d(s' | s, a_s)\} y_i \\ &= y_i + \sum_{s \in S_i} [P_i^*(\mathbf{a})]_s r(s, a_s) + \sum_{j=1}^n \sum_{s' \in S_j} \sum_{s \in S_i} [P_i^*(\mathbf{a})]_s d(s' | s, a_s) y_j \\ &= y_i + \sum_{s \in S_i} [P_i^*(\mathbf{a})]_s \{r(s, a_s) + \sum_{j=1}^n \sum_{s' \in S_j} d(s' | s, a_s) y_j\} = y_i + \\ &\sum_{s \in S_i} [P_i^*(\mathbf{a})]_s \bar{c}_i(s, a_s) = y_i + P_i^*(\mathbf{a})^T \bar{c}_i(\mathbf{a}) \end{aligned}$$

where $\bar{c}_i(s, a_s) := r(s, a_s) + \sum_{j=1}^n \sum_{s' \in S_j} d(s' | s, a_s) y_j$, $s \in S_i$ and $\bar{c}_i(\mathbf{a}) := (\bar{c}_i(s, a_s))_{s \in S_i}$.

It follows that for each $i = 1, \dots, n$, (4.4) is the same as:

$$P_i^*(\bar{g}(i))^T \bar{c}_i(\bar{g}(i)) > P_i^*(\bar{f}(i))^T \bar{c}_i(\bar{f}(i)). \quad (4.5)$$

Since $P_i^*(\mathbf{a})^T \bar{c}_i(\mathbf{a})$ is the value of the "strategy" \mathbf{a} ($s \rightarrow a_s$, $s \in S_i$) in the irreducible MDP Γ_i in which the rewards are defined by \bar{c}_i , then the strategy \bar{g} in (4.5) can be computed by one iteration of the simplified policy improvement algorithm.

From the previous results, our Aggregation-Disaggregation Algorithm for solving the limit Markov Control Problem (L) is stated as follows:

Step 1: Select an arbitrary deterministic strategy f in Γ , and set:

$$[f(i)]_s := f(s); \quad s \in S_i; \quad i = 1, 2, \dots, n.$$

Step 2: Compute $P_i^*(f(i))$; $q_{ij}(f(i))$; and $c(i, f(i))$; $i = 1, 2, \dots, n$

$j = 1, 2, \dots, n$. For each $i = 1, \dots, n$ the computation of $P_i^*(f(i))$ is done by solving the linear system:

$$\mathbf{x}^i P_i(f(i)) = \mathbf{x}^i, \quad \sum_{s \in S_i} x_s^i = 1.$$

Step 3: Solve, for the unknowns $\lambda, y_1, y_2, \dots, y_{n-1}$, the linear system:

$$\lambda + y_i = c(i, f(i)) + \sum_{j=1}^n q_{ij}(f(i)) y_j \quad i = 1, \dots, n; \quad y_n = 0.$$

Step 4: For each $i = 1, \dots, n$ compute the deterministic strategy $g(i)$ obtained after one iteration of the simplified policy improvement algorithm for the MDP Γ_i with reward \bar{c}_i (the starting strategy is $f(i)$).

Step 5: If $g(i) = f(i)$ for all $i = 1, \dots, n$ STOP.

Otherwise $f(i) \leftarrow g(i)$; $i = 1, 2, \dots, n$ and go to Step 2.

References

- [1] R. Aldhaheri and H. Khalil, *Aggregation and optimal control of nearly completely decomposable markov chains*, in Proceedings of the 28th CDC, IEEE, 1989, pp. 1277-1282.
- [2] T. R. Bielecki and J. A. Filar, *Singularly Perturbed Markov Control Problem: Limiting Average Cost*, Tech. Rep. 89-04, University of Maryland at Baltimore County, 1989.
- [3] M. Cordech, A. Willsky, S. Sastry, and D. Castanon, *Hierarchical aggregation of linear systems with multiple time scales*, IEEE Transactions on Automatic Control, **AC-28** (1983), pp. 1017-1029.
- [4] F. Delebecque, *A reduction process for perturbed markov chains*, SIAM Journal of Applied Mathematics, **48** (1983), pp. 325-350.
- [5] F. Delebecque and J. Quadrat, *Optimal control of markov chains admitting strong and weak interactions*, Automatica, **17** (1981), pp. 281-296.
- [6] E. V. Denardo, *Dynamic Programming*, Prentice-Hall, Englewood Cliffs, New Jersey, 1982.
- [7] C. Derman, *Finite State Markovian Decision Process*, Academic Press, New York, 1970.
- [8] L. C. M. Kallenberg, *Linear Programming and Finite Markovian Control Problems*, Mathematical Center Tracts 148, Amsterdam, 1983.
- [9] T. Kato, *Perturbation Theory for Linear Operators*, Springer-Verlag, Berlin, 1980.
- [10] P. Kokotovic, *Application of singular perturbation techniques to control problems*, SIAM Review, **26** (1984), pp. 501-550.
- [11] K. G. Murty, *Linear Programming*, Wiley, New York, 1983.
- [12] A. A. Pervozvanskii and V. G. Gaitsgori, *Theory of Suboptimal Decisions*, Kluwer, Dordrecht, 1988.
- [13] R. G. Phillips and P. Kokotovic, *A singular perturbation approach to modelling and control of markov chains*, IEEE Transactions on Automatic Control, **AC-26** (1981), pp. 1087-1094.
- [14] J. Rohlicek and A. Willsky, *Multiple time scale decomposition of discrete time markov chains*, Systems and Control Letters, **11** (1988), pp. 309-314.
- [15] P. J. Schweitzer, *Perturbation theory and finite markov chains*, Journal of Applied Probability, **5** (1968), pp. 401-413.
- [16] P. J. Schweitzer, *Perturbation series expansions for nearly completely-decomposable markov chains*, Teletraffic Analysis and Computer Performance Evaluation, 1986, pp. 319-328.