

**stichting
mathematisch
centrum**



AFDELING MATHEMATISCHE BESLISKUNDE
(DEPARTMENT OF OPERATIONS RESEARCH)

BW 44/76

AUGUSTUS

P.J. SCHWEITZER & A. FEDERGRUEN

THE ASYMPTOTIC BEHAVIOUR OF UNDISCOUNTED VALUE
ITERATION IN MARKOV DECISION PROBLEMS

Prepublication

2e boerhaavestraat 49 amsterdam

BIBLIOTHEEK MATHEMATISCH CENTRUM
—AMSTERDAM—

Printed at the Mathematical Centre, 49, 2e Boerhaavestraat, Amsterdam.

The Mathematical Centre, founded the 11-th of February 1946, is a non-profit institution aiming at the promotion of pure mathematics and its applications. It is sponsored by the Netherlands Government through the Netherlands Organization for the Advancement of Pure Research (Z.W.O), by the Municipality of Amsterdam, by the University of Amsterdam, by the Free University at Amsterdam, and by industries.

The asymptotic behaviour of undiscounted value iteration in Markov Decision Problems^{*)}

by

P.J. Schweitzer

and

A. Federgruen

ABSTRACT

This paper considers undiscounted Markov Decision Problems. For the general multichain case, we obtain necessary and sufficient conditions which guarantee that the maximal total expected reward for a planning horizon of n epochs minus n times the long run average expected reward has a finite limit as $n \rightarrow \infty$ for each initial state and each final reward vector. In addition, we obtain a characterization of the chain- and periodicity structure of the set of one-step and J -step maximal gain policies. Finally, we discuss the asymptotic properties of the undiscounted value-iteration method.

KEY WORDS & PHRASES: *Markov Decision Problems; average cost criterion, chain and periodicity structure, asymptotic behaviour: value-iteration method*

^{*)} This paper is not for review; it is meant for publication elsewhere.

§1. INTRODUCTION

The value-iteration equations for undiscounted Markov Decision Processes (MDPs) with finite state- and action space, were first studied by BELLMAN [2] and HOWARD [5]:

$$(1.1) \quad v(n+1)_i = Qv(n)_i, \quad i = 1, \dots, N$$

where the Q operator is defined by:

$$(1.2) \quad Qx_i = \max_{k \in K(i)} \left\{ q_i^k + \sum_{j=1}^N P_{ij}^k x_j \right\}, \quad i = 1, \dots, N$$

and $v(0)$ is a given N -vector.

$\Omega = \{1, \dots, N\}$ denotes the state space, $K(i)$ the finite set of alternatives in state i , q_i^k the one-step expected reward and $P_{ij}^k \geq 0$ the transition probability to state j , when alternative $k \in K(i)$ is chosen in state i . ($i=1, \dots, N$).

For all $n = 1, 2, \dots$ and $i \in \Omega$, $v(n)_i$ may be interpreted as the maximal total expected reward for a planning horizon of n epochs, when starting at state i and given an amount $v(0)_j$ is obtained when ending up at state j . BELLMAN [2] showed that if every P_{ij}^k is strictly positive, then $v(n)_i \sim ng^*$, $n \rightarrow \infty$, the scalar g^* being the maximal gain rate and HOWARD [5] conjectured that there generally exist two N -vectors g^* and v^* , such that

$$(1.3) \quad \lim_{n \rightarrow \infty} v(n) - ng^* - v^* = 0.$$

Although BROWN ([3],[4.3]) showed that $v(n) - ng^*$ is bounded, provided g^* is taken as the maximal gain rate vector, the limit in (1.3) may not exist for arbitrary $v(0)$ if some of the transition probability matrices (*tpm's*) are periodic. The identification of sufficient conditions for the existence of the limit in (1.3) is of particular importance:

- (a) when considering the infinite horizon-model with the average return per unit time criterion, as an approximation to the model where the planning horizon is finite though large.

(b) for the case $N \gg 1$, where the value-iteration method is the only practical way of locating maximal-gain policies. If the limit in (1.3) exists, then a generalization of ODoni [9] shows that any policy achieving the maxima in (1.2) for large n is maximal gain. However, if the limit in (1.3) fails to exist, then example 4 in LANERY [6] shows that policies achieving the maxima for large n in (1.2) need not be maximal gain.

Sufficiency conditions for the existence of the limit in (1.3) have been established by WHITE [16] and SCHWEITZER ([11],[12]) in the unichain case, where $g_i^* = g^*$ (say) for all $i \in \Omega$.

Related convergence results for MDPs with compact action spaces, the denumerable and general state space case and for continuous time Markov Decision Processes were obtained in respectively BATHER [1], HORDIJK, SCHWEITZER and TIJMS [4], TIJMS [15] and LEMBERSKY [7].

In this paper we establish the weakest sufficient condition. It holds for the general *multichain* case, and states that the limit in (1.3) exists for every $v(0) \in E^N$, *if and only if* there exists a randomized maximal gain policy whose tpm is aperiodic (but not necessarily unichained) and has $R^* = \{i \in \Omega \mid i \text{ is recurrent for some pure maximal gain policy}\}$ as its set of recurrent states.

In addition, we show that in general the sequence $\{v(n) - ng^*\}_{n=1}^\infty$ is asymptotically periodic, i.e. there exists an integer d^* (which merely depends upon the chain- and periodicity structure of the maximal gain policies), such that

$$(1.4) \quad \lim_{n \rightarrow \infty} v(nJ+r) - (nJ+r)g^* \quad \text{exists for all } v(0) \in E^N$$

if and only if J is a multiple of d^* .

The sufficiency parts of the above mentioned results were treated in LANERY [6]. However, it appears that the proof of proposition 19 in [6] from which the main result is derived, is either incomplete or incorrect (Note 1).

Moreover, our methods use the set of all randomized policies, and involve the analysis of the chain- and periodicity structure of the one- and J -step (randomized) maximal gain policies ($J \geq 1$). This enables a full characterization of the asymptotic period.

In section 2, we give some notation and preliminaries. In section 3, we analyze the periodicity structure of the maximal gain policies, while in section 4 the chain- and periodicity-structure of the multi-step maximal gain policies is characterized. In section 5, we obtain inter alia the above mentioned results with respect to the asymptotic periodicity, and the necessary and sufficient condition for the existence of the limit in (1.3) for all $v(0) \in E^N$.

Finally, we show how the convergences of the various sequences $\{v(nJ+r)_i - (nJ+r)g_i^*\}_{n=1}^\infty$ ($r=1, \dots, J$; $i=1, \dots, N$) interdepend.

In section 6, we give some properties of the policies that attain the maxima in (1.1) for large n .

§2. NOTATION AND PRELIMINARIES

A (stationary) randomized policy f is a tableau $[f_{ik}]$ satisfying $f_{ik} \geq 0$ and $\sum_{k \in K(i)} f_{ik} = 1$, where f_{ik} is the probability that the k -th alternative is chosen when entering state i .

We let S_R denote the set of all randomized policies, and S_P the set of all pure (non-randomized) policies (i.e. each $f_{ik}=0$ or 1). Associated with each $f \in S_R$, are a N -component reward $q(f)$ and $N \times N$ -matrix $P(f)$:

$$(2.1) \quad q(f)_i = \sum_{k \in K(i)} f_{ik} q_i^k; \quad P(f)_{ij} = \sum_{k \in K(i)} f_{ik} P_{ij}^k, \quad 1 \leq i, j \leq N.$$

Note that $P(f)$ is a stochastic matrix ($P(f)_{ij} \geq 0$, $\sum_{j=1}^N P(f)_{ij} = 1$; $1 \leq i, j \leq N$). For any $f \in S_R$, we define the stochastic matrix $\Pi(f)$ as the Cesaro limit of the sequence $\{P^n(f)\}_{n=1}^\infty$, which always exists and has the following properties:

$$(2.2) \quad P(f)\Pi(f) = \Pi(f) = \Pi(f)P(f).$$

Denote by $n(f)$ the number of subchains (closed, irreducible sets of states) for $P(f)$. Then:

$$(2.3) \quad \Pi(f)_{ij} = \sum_{m=1}^{n(f)} \phi_i^m(f) \pi^m(f)_j,$$

where $\pi^m(f)$ is the unique equilibrium distribution of $P(f)$ on the m^{th} sub-chain $C^m(f)$, and $\phi_i^m(f)$ is the probability of absorption in $C^m(f)$, starting from state i . Let $R(f) = \{j | \Pi(f)_{jj} > 0\}$, i.e. $R(f)$ is the set of recurrent states for $P(f)$.

Let $d^m(f) \geq 1$ denote the period of $C^m(f)$, and let $\{C^{m,\beta}(f) \mid \beta = 1, \dots, d^m(f)\}$ indicate the set of cyclically moving subsets (c.m.s.) of $C^m(f)$ numbered such that for any $m = 1, \dots, n(f)$ and $\beta = 1, \dots, d^m(f)$ (cf. [10]):

$$(2.4) \quad i \in C^{m,\beta}(f) \Rightarrow P(f)_{ij} > 0 \text{ only if } j \in C^{m,\beta+1}(f)$$

with the convention that hereafter β in $C^{m,\beta}(f)$ is taken modulo $d^m(f)$ e.g. $C^{m,\beta+1}(f) = C^{m,1}(f)$ if $\beta = d^m(f)$.

$$(2.5) \quad \begin{aligned} &\text{For all } i \in C^m(f): \\ &d^m(f) = \text{greatest common divisor (g.c.d.) of } \{n \mid P(f)_{ii}^n > 0\} \\ &= \text{g.c.d. } \{n \mid \text{there exists a cycle } (s_0=i, s_1, \dots, s_n=i) \\ &\quad \text{for } P(f)\} \end{aligned}$$

where $(s_0=i, s_1, \dots, s_n=i)$ is called a *cycle* for $P(f)$ if $P(f)_{s_1 s_1+1} > 0$ and if all the s_l are distinct ($l=0, \dots, n-1$).

$$(2.6) \quad \lim_{n \rightarrow \infty} P^{nd^m(f)+r}(f)_{ij} > 0, \text{ for all } i \in C^{m,\beta}(f) \text{ and } j \in C^{m,\beta+r}(f) \\ (r = 1, 2, \dots).$$

For each $f \in S_R$, we define the gain rate vector $g(f) = \Pi(f)q(f)$, such that $g(f)_i$ represents the long run average expected return per unit time, when the initial state is i , and policy f is used. We thus have

$$(2.7) \quad g(f)_i = \sum_{m=1}^{n(f)} \phi_i^m(f) g^m(f), \quad i \in \Omega$$

with $g^m(f) = \langle \pi^m(f), q(f) \rangle, \quad m = 1, \dots, n(f).$

Next define:

$$(2.8) \quad g_i^* = \sup_{f \in S_R} g(f)_i; \quad i = 1, \dots, N.$$

Since HOWARD [5] proved the existence of pure policies f which attain the N suprema in (2.8) simultaneously, we can define:

$$(2.9) \quad S_{PMG} = \{f \in S_P \mid g(f) = g^*\}; \quad S_{RMG} = \{f \in S_R \mid g(f) = g^*\}$$

as the set of all pure and the set of randomized maximal gain policies.

Finally define R^* as the set of states that are recurrent under some maximal gain policy:

$$R^* = \{i \mid i \in R(f) \text{ for some } f \in S_{RMG}\}.$$

The following lemma which was proved in SCHWEITZER & FEDERGRUEN [13], (th. 3.2) provides a basic characterization of this set:

LEMMA 2.1

- (a) $R^* = \{i \mid i \in R(f) \text{ for some } f \in S_{PMG}\}.$
- (b) *The set $\{f \in S_{RMG} \mid R(f) = R^*\}$ is not empty.*
- (c) *Define $n^* = \min\{n(f) \mid f \in S_{RMG} \text{ with } R(f) = R^*\}$ and*

$$S_{RMG}^* = \{f \in S_{RMG} \mid R(f) = R^* \text{ and } n(f) = n^*\}.$$

Fix $f^ \in S_{RMG}^*$.*

Any subchain of any $f \in S_{RMG}$ is contained within a subchain of $P(f^)$.*

- (d) *All $f^* \in S_{RMG}^*$ have the same collection of subchains $\{R^{*\alpha}, \alpha = 1, \dots, n^*\}$.*
 - (e) *For any $\alpha \in \{1, \dots, n^*\}$, $g_i^* = g^{*\alpha}$ (say) for all $i \in R^{*\alpha}$.*
 - (f) *Let $R^{(1)}, \dots, R^{(m)}$ be disjoint sets of states such that*
 - (1) *if C is a subchain of some $f \in S_{RMG}$, then $C \subseteq R^{(k)}$, for some k , $1 \leq k \leq m$*
 - (2) *there exists a $f \in S_{RMG}$, with $\{R^{(k)} \mid k = 1, \dots, m\}$ as its set of subchains.*
- Then $m = n^*$ and after renumbering $R^{(\alpha)} = R^{*\alpha}$.*

Define the operator

$$(2.10) \quad Tx = \max_{k \in L(i)} \{q_i^k + \sum_j p_{ij}^k x_j\},$$

where

$$L(i) = \{k \in K(i) \mid g_i^* = \sum_j p_{ij}^k g_j^*\}, \quad \text{for all } i \in \Omega.$$

Let Q^n (and T^n) denote the n -fold application of the operator $Q(T)$:

$$Q^n x = Q(Q^{n-1} x); \quad T^n x = T(T^{n-1} x); \quad n = 2, \dots \text{ and } x \in E^N.$$

The basic properties of both operators were studied in SCHWEITZER & FEDERGRUEN [14]. In particular, it was shown that the Q operator reduces to T in the following two ways:

$$(2.11) \quad \text{for each } x \in E^N, \text{ there exists a scalar } t_0(x), \text{ such that}$$

$$Q^n(x + tg^*) = T^n(x + tg^*) \text{ for } n = 1, 2, \dots \text{ and } t \geq t_0(x) \text{ (cf. [14], lemma 2.2 part (c))}$$

$$(2.12) \quad \text{for each } x \in E^N \text{ there exists an integer } n_0(x) \text{ such that}$$

$$Q^{n+1} x = T(Q^n x) = T^{n+1-n_0(x)} Q^{n_0(x)} x,$$

for all $n \geq n_0(x)$ (cf. [3] and [14], lemma 2.2 part (c))

We next consider the functional equation:

$$(2.13) \quad v + g^* = Tv.$$

Let $V = \{v \in E^N \mid v \text{ satisfies (2.13)}\}$ and define for any $v \in V$:

$$(2.14) \quad b(v)_i^k = q_i^k - g_i^* + \sum_{j=1}^N p_{ij}^k v_j - v_i, \quad i \in \Omega, k \in K(i)$$

$$b(v, f)_i = \sum_{k \in K(i)} f_{ik} b(v)_i^k = [q(f) - g^* + P(f)v - v]_i, \quad i \in \Omega, \quad f \in S_R$$

Observe, that for all $v \in V$, $\max_{k \in L(i)} b(v)_i^k = 0$, for all $i \in \Omega$.

Finally, we define for any $i \in R^*$, the set $K^*(i)$ as the set of actions which a pure maximal gain policy that has i among its recurrent states, could prescribe:

$$(2.15) \quad K^*(i) = \{k \in K(i) \mid \text{there exists a } f \in S_{PMG}, \text{ with } i \in R(f) \text{ and } f_{ik} = 1\}.$$

The following lemma gives the necessary and sufficient condition for a policy to be maximal gain, characterizes the sets $K^*(i)$ and shows that any policy that randomizes among all actions in $K^*(i)$, in each of the states in R^* , and among all actions in $L(i)$ for the states in $\Omega - R^*$, belongs to S_{RMG}^* :

LEMMA 2.2

- (a) Fix $v \in V$. A policy $f \in S_R$ is maximal gain (i.e. $f \in S_{RMG}^*$) if and only if
- (1) for all $i \in \Omega$, $f_{ik} > 0 \Rightarrow k \in L(i)$ i.e. $P(f)g^* = g^*$
 - (2) for all $i \in R(f)$, $f_{ik} > 0 \Rightarrow b(v)_i^k = 0$ i.e. $\Pi(f)b(v, f) = 0$.
- (b) $K^*(i) = \{k \in L(i) \mid \text{there exists a } f \in S_{RMG}, \text{ with } i \in R(f), \text{ and } f_{ik} > 0\}$,
 $i \in R^*$.
- (c) For any $v \in V$,
- $$K^*(i) = \{k \in L(i) \mid b(v)_i^k = 0 \text{ and } \sum_{j \in R^{*\alpha}} P_{ij}^k = 1\}, \text{ for all } i \in R^{*\alpha},$$
- $$\alpha = 1, \dots, n^*.$$
- (d) Define $f^* \in S_R$ such that

$$\{k \mid f_{ik}^* > 0\} = \begin{cases} K^*(i), & i \in R^* \\ L(i), & i \in \Omega - R^*. \end{cases}$$

Then $f^* \in S_{RMG}^*$.

PROOF

(a) cf. theorem 3.1 part (a) in [13].

(b) Clearly, $K^*(i)$ is contained within the set on the right hand side.

Next, fix $i \in R^*$, $k \in K(i)$ and $f \in S_{RMG}$, such that $i \in R(f)$ and $f_{ik} > 0$, and use lemma 2.1 in [13] in order to show that there exists a $h \in S_{PMG}$,

with $i \in R(h)$, and $h_{ik} = 1$, as well, which proves the reversed inclusion.

(c) Fix $\alpha \in \{1, \dots, n^*\}$, $i_0 \in R^{\alpha}$.

First, let $k \in K^*(i)$ and $f \in S_{\text{RMG}}$, with $i \in R(f)$ and $f_{ik} > 0$, and apply part (a) of this lemma, and part (c) of lemma 2.1, in order to prove that $K^*(i)$ is contained within the set on the right hand side of the equality.

Next, take $k_0 \in L(i_0)$ such that $b(v)_{i_0}^{k_0} = 0$ and $\sum_{j \in R^{\alpha}} P_{ij}^{k_0} = 1$, and fix $f^* \in S_{\text{RMG}}^*$. Define f^{**} such that

$$f_{i_0 k_0}^{**} = 1, \text{ and } f_{jk}^{**} = f_{jk}^*, \quad \text{for all } j \neq i_0, k \in K(i).$$

Use part (d) of lemma 2.1, in order to show that all states in $R^{\alpha} \setminus \{i_0\}$ can reach state i_0 under $P(f^{**})$ whereas state i_0 can only reach states within R^{α} . We conclude that $i_0 \in R(f^{**})$, while $f^{**} \in S_{\text{RMG}}$, as can be verified using part (a) of this lemma, thus proving the reversed inclusion.

(d) cf. remark 1 in [13]. \square

We finally need the following lemma:

LEMMA 2.3

(a) Fix $f^1, f^2 \in S_R$, and let C^1 and C^2 be two subchains of $P(f^1)$ and $P(f^2)$ with period d^1 and d^2 respectively, such that $C^1 \cap C^2 \neq \emptyset$. Define f^3 such that

$$\{k \mid f_{ik}^3 > 0\} = \begin{cases} \{k \mid f_{ik}^2 > 0\} & \text{for all } i \in C^2 \setminus C^1 \\ \{k \mid f_{ik}^1 > 0\} \cup \{k \mid f_{ik}^2 > 0\} & \text{for all } i \in C^1 \cap C^2 \\ \{k \mid f_{ik}^1 > 0\} & \text{otherwise.} \end{cases}$$

Then

- (1) $C^1 \cup C^2$ is a subchain of $P(f^3)$, the period d^3 of which is a common divisor of d^1 and d^2 .
- (2) if $f^1, f^2 \in S_{\text{RMG}}$, then $f^3 \in S_{\text{RMG}}$.

(b) For any $f \in S_R$, define the set of pure policies $S_P(f) = \bigcap_{i \in \Omega} \{k \mid f_{ik} > 0\}$.

Then for all $m = 1, \dots, n(f)$:

$$(2.16) \quad d^m(f) = \text{g.c.d.}\{d^r(h) \mid h \in S_p(f), 1 \leq r \leq n(h), C^r(h) \subseteq C^m(f)\}.$$

PROOF

(a) (1) Show that $C^1 \cup C^2$ is a closed and communicating set of states for $R(f^3)$. The former is immediate; the latter holds since any state in $C^1 \cap C^2$ communicates with $C^1 \cup C^2$.

Since $\{n \mid P(f^3)_{ii}^n > 0\} \supseteq \{n \mid P(f^1)_{ii}^n > 0\} \cup \{n \mid P(f^2)_{ii}^n > 0\}$, it follows (cf. (2.5)) that $d^3 = \text{g.c.d.}\{n \mid P(f^3)_{ii}^n > 0\}$ is a common divisor of d^1 and d^2 .

(2) Observe that for each $i \in \Omega$, $f_{ik}^3 > 0$ only for $k \in L(i)$ since it follows from lemma 2.2 part (a) that $f_{ik}^1 > 0$ and $f_{ik}^2 > 0$ only for $k \in L(i)$. Using the fact that $R(f^3) \subseteq R(f^1) \cup C^2$, and applying lemma 2.2 part (a2) one verifies that $f^3 \in S_{\text{RMG}}$.

(b) Fix $m \in \{1, \dots, n(f)\}$ and $h \in S_p(f)$. Since $C^m(f)$ is closed under any policy in $S_p(f)$, $P(h)$ has a subchain $C^r(h) \subseteq C^m(h)$ ($1 \leq r \leq n(h)$). Since $P(h)_{ij} > 0$ only if $P(f)_{ij} > 0$, and since $i \in C^m(f)$ implies that

$P(f)_{ii}^t > 0$ only if t is a multiple of $d^m(f)$, it follows that for

$i \in C^r(h)$, $P(h)_{ii}^t > 0$ only if t is a multiple of $d^m(f)$. Thus (2.5)

implies that the left hand side of (2.16) is less than or equal to its right hand side.

To prove the reversed inequality in (2.16) fix $i \in C^m(f)$ and recall from (2.5) that

$$(2.17) \quad d^m(f) = \text{g.c.d.}\{n \mid \text{there exists a cycle } (s_0=i, \dots, s_n=i) \text{ of } P(f)\}.$$

We next show that

$$(2.18) \quad \text{for each cycle } S = \{s_0 = i, s_1, \dots, s_n = i\} \text{ of } P(f), \text{ there exists a pure policy } g \in S_p(f) \text{ which has } i \text{ recurrent and contains the same cycle.}$$

As a consequence, we obtain that each of the elements in the set to the left of (2.17) is a multiple of the period of a subchain of a pure policy

that lies within $C^m(f)$, thus proving the reversed inequality in (2.16) and hence part (b).

In order to show (2.18), construct the policy $h \in S_P(f)$ as follows:

Let $h_{s_1 k} = 1$ for any one k such that $f_{s_1 k} > 0$ and $p_{s_1 s_{1+1}}^k > 0$ ($1=0, \dots, n-1$);

for $j \notin C^m(f)$, let $h_{jk} = 1$ for any one k such that $f_{jk} > 0$.

If $S \neq C^m(f)$, let Δ initially be equal to S , and define $\bar{\Delta} = C^m(f) - \Delta$.

Next, the following step is performed:

Choose a state $j \in \bar{\Delta}$ and an alternative k such that $f_{jk} > 0$ and $p_{jt}^k > 0$ for some $t \in \Delta$, transfer j from $\bar{\Delta}$ to Δ and define $h_{jk} = 1$. Such k and t can always be found since all states in $C^m(f)$ communicate under $P(f)$. Repeat this step for the new Δ and $\bar{\Delta}$, until $\bar{\Delta}$ is empty. This construction shows that S is a cycle for $P(h)$, with $i \in R(h)$ since i can be reached from any state in $C^m(f)$, and $C^m(f)$ is closed. \square

REMARK 1. The period d^3 , defined in part (a) of the previous lemma, does not necessarily have to be the greatest common divisor of d^1 and d^2 . Take

$$P(f^1) = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix} \text{ and } P(f^2) = \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \text{ with } d^1 = d^2 = 3 \text{ and } d^3 = 1.$$

However, it can be shown that $d^3 = \text{g.c.d.}\{d^1, d^2\}$ does hold, when $P(f^1)$ and $P(f^2)$ merely differ in one row, the corresponding state being recurrent for both chains (cf. part (b)).

§3. THE PERIODICITY STRUCTURE OF THE POLICIES IN S_{RMG}

We first define

$$(3.1) \quad d(\alpha) = \min \{d^m(f) \mid f \in S_{\text{RMG}}, 1 \leq m \leq n(f), C^m(f) \subseteq R^{*\alpha}\},$$

$$\alpha = 1, \dots, n^*$$

$$(3.2) \quad d_i = \min \{d^m(f) \mid f \in S_{\text{RMG}}, 1 \leq m \leq n(f), i \in C^m(f)\}, i \in R^*$$

i.e. $d(\alpha) [d_i]$ denotes the minimum of the periods of the subchains of the maximal gain policies that lie within $R^{*\alpha}$ [that contain the state i]. Let $f^* \in S_{RMG}^*$ be defined as in lemma 2.2 part (d), i.e. let

$$\{k \mid f_{ik}^* > 0\} = \begin{cases} K^*(i), & i \in R^* \\ L(i), & i \in \Omega \setminus R^*. \end{cases}$$

For each $\alpha = 1, \dots, n^*$ and $t = 1, \dots, d^\alpha(f^*)$ let $R^{*\alpha, t} = C^{\alpha, t}(f^*)$ with the convention that hereafter t in $R^{*\alpha, t}$ is taken modulo $d^\alpha(f^*)$ (e.g. $R^{*\alpha, t} = R^{*\alpha, 1}$ if $t = d^\alpha(f^*) + 1$).

THEOREM 3.1: (Periodicity structure) (cf. lemma 2.1)

- (a) $d^\alpha(f^*) = d(\alpha)$, $\alpha = 1, \dots, n^*$.
- (b) Fix $\alpha \in \{1, \dots, n^*\}$. Let $h \in S_{RMG}$ and $C^m(h) \subseteq R^{*\alpha}$. Then $d^m(h)$ is a multiple of $d(\alpha)$.
- (c) $d(\alpha) = \text{g.c.d.}\{d^m(f) \mid f \in S_{PMG}, 1 \leq m \leq n(f), C^m(f) \subseteq R^{*\alpha}\}$, $\alpha = 1, \dots, n^*$.
- (d) $d_i = d(\alpha)$ for all $i \in R^{*\alpha}$, $\alpha = 1, \dots, n^*$.
- (e) $d(\alpha) = \min\{d^\alpha(f) \mid f \in S_{RMG}^*\}$, $\alpha = 1, \dots, n^*$.
- (f) The set $S_{RMG}^{**} = \{f \in S_{RMG}^* \mid d^\alpha(f) = d(\alpha), \alpha = 1, \dots, n^*\}$ is non-empty.
- (g) For each $i \in R^*$, say $i \in R^{*\alpha, t}$ ($1 \leq \alpha \leq n^*$; $1 \leq t \leq d(\alpha)$) and $k \in K^*(i)$:
 $P_{ij}^k > 0 \Rightarrow j \in R^{*\alpha, t+1}$.
- (h) For each $h \in S_{RMG}$, and $i \in R(h) \cap R^{*\alpha, t}$ ($1 \leq \alpha \leq n^*$; $1 \leq t \leq d(\alpha)$)
 $P(h)_{ij} > 0$ only for $j \in R^{*\alpha, t+1} \cap R(h)$.
- (i) Fix $h \in S_{RMG}$, with $C^m(h) \subseteq R^{*\alpha}$ ($1 \leq m \leq n(h)$; $1 \leq \alpha \leq n^*$). $C^m(h)$ has $d^m(h)/d(\alpha)$ c.m.s. within each of the sets $R^{*\alpha, t}$ ($1 \leq t \leq d(\alpha)$).
- (j) All $f \in S_{RMG}^{**}$ have the same collection of c.m.s.
 $\{R^{*\alpha, t} \mid (\alpha=1, \dots, n^*; t=1, \dots, d(\alpha))\}$.
- (k) Let $R^{(1)}, \dots, R^{(M)}$ be disjoint sets of states, such that
 - (1) if C is a c.m.s. of some subchain of some $f \in S_{RMG}$, then $C \subseteq R^{(k)}$, for some k , $1 \leq k \leq M$.
 - (2) there exists a $f \in S_{RMG}$, with $\{R^{(k)} \mid k = 1, \dots, M\}$ as its collection of c.m.s.

Then, $M = \sum_{\alpha=1}^{n^*} d(\alpha)$ and after renumbering $R^{(k)} = R^{*\alpha, t}$, $k = 1, \dots, \sum_{\alpha=1}^{n^*} d(\alpha)$.

PROOF

(a), (b) Fix $\alpha \in \{1, \dots, n^*\}$ and let $h \in S_{\text{RMG}}$, with $C^m(h) \subseteq R^{*\alpha}$ (for some m , $1 \leq m \leq n(h)$). Define f^{**} such that

$$\{k \mid f_{ik}^{**} > 0\} = \begin{cases} \{k \mid h_{ik} > 0\} \cup \{k \mid f_{ik}^* > 0\}, & \text{for all } i \in C^m(h) \\ \{k \mid f_{ik}^* > 0\}, & \text{otherwise.} \end{cases}$$

It then follows from the definitions of the policy f^* and the sets $K^*(i)$ (cf. lemma 2.2 part (b)) that

$$\{k \mid f_{ik}^{**} > 0\} = \begin{cases} K^*(i) & \text{for } i \in R^* \\ L(i) & \text{for } i \in \Omega \setminus R^* \end{cases}$$

which implies that f^* and f^{**} have the same chain- and periodicity structure. In particular,

$$d^\alpha(f^{**}) = d^\alpha(f^*).$$

On the other hand, applying lemma (2.3), part (a), it follows that $d^\alpha(f^{**})$ is a divisor of $d^m(h)$, hence

$$(3.3) \quad d^\alpha(f^*) \text{ divides } d^m(h),$$

such that

$$\begin{aligned} d(\alpha) &\leq d^\alpha(f^*) \leq \min\{d^m(h) \mid h \in S_{\text{RMG}}, 1 \leq m \leq n(h), C^m(h) \subseteq R^{*\alpha}\} \\ &= d(\alpha). \end{aligned}$$

This proves part(a), whereas the combination of part (a) and (3.3) proves part (b).

(c) Define f^* as in part (a), use the fact that $d(\alpha) = d^\alpha(f^*)$ and apply lemma (2.3) part (b).

(d) Fix $i \in R^{*\alpha}$. Clearly $d_i \geq d(\alpha)$ (cf. (3.1) and (3.2)) and use part (a) to show $d_i \leq d(\alpha)$ as well.

(e), (f) immediate from part (a).

(g) Observe that $P(f^*)_{ij} > 0 \Rightarrow j \in R^{*\alpha, t+1}$ (cf. (2.4)) and use lemma 2.2 part (d).

(h) Use the fact that $h_{ik} > 0$ only for $k \in K^*(i)$ (cf. lemma 2.2 part (b)) and apply part (g).

(i) Recall from part (b) that $d^m(h)$ is a multiple of $d(\alpha)$. Take $i \in C^{m,1}(h)$, assume $i \in R^{*\alpha, t}$ ($1 \leq t \leq d(\alpha)$) and fix $s \in \{0, \dots, d(\alpha) - 1\}$. In view of part (h), we obtain for $r = 0, \dots, \frac{d^m(h)}{d(\alpha)} - 1$:

$$P(h)_{ij}^{nd^m(h)+rd(\alpha)+s} > 0 \text{ only for } j \in R^{*\alpha, t+s} \quad n = 1, 2, \dots$$

Since $\lim_{n \rightarrow \infty} P(h)_{ij}^{nd^m(h)+rd(\alpha)+s} > 0$ for all $j \in C^{m, rd(\alpha)+s+1}(h)$, (cf. (2.4)) we conclude that $C^{m, rd(\alpha)+s+1}(h) \subseteq R^{*\alpha, t+s}$ for $r = 0, \dots, \frac{d^m(h)}{d(\alpha)} - 1$ which proves part (i).

(j) Let $f \in S_{RMG}^{**}$ and fix $\alpha \in \{1, \dots, n^*\}$. It follows from part (i) that each of the sets $R^{*\alpha, t}$ ($1 \leq t \leq d(\alpha)$) contains exactly *one* c.m.s. $C^{\alpha, s}(f)$ (for some $1 \leq s \leq d(\alpha)$) of $P(f)$.

Since $R^{*\alpha} = \bigcup_{s=1}^{d(\alpha)} C^{\alpha, s}(f) = \bigcup_{t=1}^{d(\alpha)} R^{*\alpha, t}$, we conclude that for any $1 \leq s \leq d(\alpha)$:

$$C^{\alpha, s}(f) = R^{*\alpha, t} \quad \text{for some } t = t(s)$$

which proves that all $f \in S_{RMG}^{**}$ have the same collection of c.m.s.

(k) Apply property (1) to conclude that $R^{*\alpha, t} \subseteq R^{(k(\alpha, t))}$ for $\alpha = 1, \dots, n^*$; $t = 1, \dots, d(\alpha)$, and apply property (2) and part (i) to conclude

$$R^{(k)} \subseteq \text{some } R^{*\alpha, t}, \quad k = 1, \dots, M. \quad \square$$

REMARK 2. In [13], a finite procedure was given for calculating R^*, n^* and each $R^{*\alpha}$ after using the Policy Improvement Algorithm to find g^* and $a \in V$.

Part (a) of the previous theorem shows that this procedure can be extended in order to find the $d(\alpha)$, the sets $R^{*\alpha, \beta}$ and a $f \in S_{RMG}^{**}$ in a finite number of calculations, as well:

- (1) For each $i \in R^*$, determine the sets $K^*(i)$ (use lemma 2.2 part (c))
 (2) Define $f^* \in S_{RMG}^{**}$ by

$$\{k \mid f_{ik}^* > 0\} = \begin{cases} K^*(i), & i \in R^* \\ L(i), & i \in \Omega \setminus R^* \end{cases}$$

Then the cyclically moving subsets of each subchain $R^{*\alpha}$ of $P(f^*)$ form the $\{R^{*\alpha, \beta}\}_{\beta=1}^{n^*}$.

Consider the following example:

EXAMPLE 1:

i	k	P_{i1}^k	P_{i2}^k	P_{i3}^k	P_{i4}^k	P_{i5}^k	q_i^k	S_p	$n(f)$	$C^1(f)$	$C^2(f)$	$d^1(f)$	$d^2(f)$
1	1	1	0	0	0	0	0	f^1 (1,1,1,1,1)	2	{1}	{2,3}	1	2
2	1	0	0	1	0	0	$q_2^1 \leq 0$	f^2 (1,1,1,1,2)	2	{1}	{2,3}	1	2
	2	0	0	0	1	0	0	f^3 (1,1,1,1,3)	2	{1}	{2,3}	1	2
3	1	0	1	0	0	0	0	f^4 (1,2,1,1,1)	2	{1}	{2,3,4,5}	1	4
4	1	0	0	0	0	1	0	f^5 (1,2,1,1,2)	2	{1}	{2,4,5}	1	3
5	1	0	0	1	0	0	0	f^6 (1,2,1,1,3)	1	{1}	—	1	—
	2	0	1	0	0	0	$q_5^2 \leq 0$						
	3	1	0	0	0	0	0						

Table 1:

Table 1 lists the six pure policies, their subchains and periods. Observe that (whatever the specific value of q_2^1, q_5^2):

$g^* = (0,0,0,0,0)$; $K(i) = L(i)$ for all $i \in \Omega$ and $V = \{(x_1, \dots, x_5) \mid x_1 = x_2 = x_3 = x_4 \geq x_5\}$, $n^* = 2$, $R^{*1} = \{1\}$; $R^{*2} = \{2,3,4,5\}$; since $d(1) = 1$, $R^{*1,1} = \{1\}$.

Next, consider the following cases:

case	q_2^1	q_5^2	S_{PMG}	$K^*(2)$	$K^*(5)$
1	0	0	$\{f^1, f^2, f^3, f^4, f^5, f^6\}$	{1,2}	{1,2}
2	<0	0	$\{f^4, f^5, f^6\}$	{2}	{1,2}
3	0	<0	$\{f^1, f^2, f^3, f^4, f^6\}$	{1,2}	{2}
4	<0	<0	$\{f^4, f^6\}$	{2}	{2}

Define $f^* \in S_R$ as in lemma 2.2 part (d):

$$P(f^*) = \begin{array}{cc} \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & x & x & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 \\ 0 & x & x & 0 & 0 \end{bmatrix} & \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 \\ 0 & x & x & 0 & 0 \end{bmatrix} \\ \text{case 1} & \text{case 2} \end{array}$$

$$P(f^*) = \begin{array}{cc} \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & x & x & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 \end{bmatrix} & \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 \end{bmatrix} \\ \text{case 3} & \text{case 4} \end{array}$$

In case 1, $P(f^*)$ is aperiodic and $d(2) = 1 = \text{g.c.d. } \{2, 2, 2, 4, 3\}$ (cf. th. 3.1(a)).

In case 2, $P(f^*)$ is aperiodic and $d(2) = 1 = \text{g.c.d. } \{4, 3\}$ (").

In case 3, $P(f^*)$ has R^{*2} periodic with $d(2) = 2 = \text{g.c.d. } \{2, 2, 2, 4\}$ (").

$$R^{*2,1} = \{2, 5\}; R^{*2,2} = \{3, 4\}$$

In case 4, $P(f^*)$ has R^{*2} periodic with $d(2) = 4 = \text{g.c.d. } \{4\}$ (").

$$R^{*2,1} = \{2\}; R^{*2,2} = \{4\}; R^{*2,3} = \{5\}; R^{*2,4} = \{3\}.$$

Thus randomization is essential for both the recurrency properties and the periodicity structure: it plays the indispensable role of coalescing subchains and of decreasing periods. In general, there may fail to exist a pure maximal gain policy f with $R(f) = R^*$, or which achieves the minimal number n^* of subchains, or which achieves the minimal period in every subchain. For instance, case 1 of example 1 with state 1 and action 3 in state 5 omitted, shows that

- (a) all pure (maximal gain) policies have periodic tpm's, while a randomized (maximal gain) policy is aperiodic.
- (b) none of the pure (max. gain) policies has R^* as its recurrent set, although a randomized (max. gain) policy does.

Observe that , whereas $d(\alpha) = \text{g.c.d. } \{d^m(f) \mid f \in S_{\text{PMG}}, 1 \leq m \leq n(f), C^m(f) \subseteq R^{*\alpha}\}$ for all $\alpha = 1, \dots, n^*$ (cf. part (a) of th.3.1), we may have

$$d_i = \text{g.c.d. } \{d^m(f) \mid f \in S_{\text{RMG}}, 1 \leq m \leq n(f), i \in C^m(f)\} <$$

$$< \text{g.c.d. } \{d^m(f) \mid f \in S_{\text{PMG}}, 1 \leq m \leq n(f), i \in C^m(f)\}$$

(Take case 1 of example 1, and $i = 3$).

§4. THE MULTI-STEP POLICIES

Fix an integer $J \geq 2$, and observe from (1.2) that

$$(4.1) \quad Q^J x_i = \max_{\xi \in \tilde{K}(i)} \{ \tilde{q}_i^\xi + \sum_j \tilde{p}_{ij}^\xi x_j \} \text{ where}$$

$$\tilde{K}(i) = \{(f^1, \dots, f^J) \mid f^1, \dots, f^J \in S_p\}$$

$$\tilde{q}_i^\xi = q(f^1)_i + P(f^1)q(f^2)_i + \dots + P(f^1) \dots P(f^{J-1})q(f^J)_i,$$

$$i \in \Omega, \xi = (f^1, \dots, f^J) \in \tilde{K}(i)$$

$$\tilde{p}_{ij}^\xi = P(f^1) \dots P(f^J)_{ij}; \quad 1 \leq i, j \leq N \text{ and } \xi = (f^1, \dots, f^J) \in \tilde{K}(i).$$

Let $\tilde{Q} = Q^J$, and define a related "J-step"-MDP, denoted by a tilde, with Ω as its state space, $\tilde{K}(i)$ as the (finite) set of alternatives in state $i \in \Omega$, \tilde{q}_i^ξ as the one-step expected reward and \tilde{p}_{ij}^ξ as the transition probability to state j , when alternative $\xi \in \tilde{K}(i)$ is chosen when entering state i .

Let \tilde{S}_R dentote the set of all (stationary) randomized policies with respect to the above defined MDP, and observe that

$$\tilde{S}_R = \bigcup_{i \in \Omega} \bigcap_{r=1}^J S_R.$$

In complete analogy to the definitions given in section 2, we define the operator \tilde{T} , the sets \tilde{S}_P , \tilde{S}_{PMG} , \tilde{S}_{RMG} , \tilde{S}_{RMG}^* , $\tilde{S}_{\text{RMG}}^{**}$, \tilde{R}^* , $\tilde{R}^{*\alpha}$, $\tilde{R}^{*\alpha, \beta}$, \tilde{V} the integers \tilde{n}^* , $\tilde{d}(\alpha)$, \tilde{d}_i , for each $\phi \in \tilde{S}_R$, the quantities $\tilde{q}(\phi)$, $\tilde{p}(\phi)$, $\tilde{\Pi}(\phi)$, $\tilde{g}(\phi)$, $\tilde{n}(\phi)$, $\tilde{d}^m(\phi)$, and for each $i \in \Omega$, the set $\tilde{L}(i)$.

Observe that a "J-step policy" $\phi \in \tilde{S}_R$ is specified by NJ "one-step" policies $\{\phi^{r,i} \mid r = 1, \dots, J; i = 1, \dots, N\}$ such that policy ϕ uses "action" $(\phi^{1,i}, \phi^{2,i}, \dots, \phi^{J,i}) \in \tilde{K}(i)$ while in state i :

$$\begin{aligned}\tilde{q}(\phi)_i &= q(\phi^{1,i})_i + P(\phi^{1,i})q(\phi^{2,i})_i + \dots P(\phi^{1,i}) \dots P(\phi^{J-1,i})q(\phi^{J,i})_i \\ \tilde{P}(\phi)_{ij} &= P(\phi^{1,i}) \dots P(\phi^{J,i})_{ij}.\end{aligned}$$

The following theorem characterizes the "J-step" maximal gain policies and shows how their chain- and periodicity structure are connected with the corresponding ones in our original MDP.

First, define for any $\phi \in \tilde{S}_R$:

$$(4.2) \quad \begin{aligned}T^{r,i}(\phi) &= \{j \mid P(\phi^{1,i}) \dots P(\phi^{r,i})_{ij} > 0\}, \quad i \in \Omega, r = 1, \dots, J \\ \tilde{T}^{0,i}(\phi) &= \{i\}, \quad i \in \Omega.\end{aligned}$$

THEOREM 4.1. Fix $J \geq 2$. Then

(a) $\tilde{g}^* = Jg^*$ and $\{\phi \mid \text{there exists } f \in S_{\text{RMG}} \text{ such that } \phi^{r,i} = f \text{ for all } r = 1, \dots, J; i = 1, \dots, N\} \subseteq \tilde{S}_{\text{RMG}}$.

(b) Let $\xi = (f^1, \dots, f^J) \in \tilde{K}(i)$. The following statements are equivalent:

(1) $\xi \in \tilde{L}(i)$.

(2) $f_{ik}^1 = 1 \Rightarrow k \in L(i)$.

$f_{jk}^r = 1 \Rightarrow k \in L(j)$ for $2 \leq r \leq J$ and all j , such that

$$P(f^1) \dots P(f^{r-1})_{ij} > 0.$$

(c) V is an n^* -dimensional subset of the \tilde{n}^* -dimensional set \tilde{V} .

(d) Fix $v \in V$. Then $\phi \in \tilde{S}_{\text{RMG}}$ if and only if

$$(4.3) \quad \phi_{jk}^{r+1,i} > 0 \Rightarrow k \in L(i), \text{ for all } j \in T^{r,i}(\phi), i \in \Omega, r = 0, \dots, J-1$$

$$b(v, \phi^{r+1,i})_j = 0 \text{ for all } j \in T^{r,i}(\phi), i \in \tilde{R}(\phi), r = 0, \dots, J-1.$$

(e) Fix $f \in S_{\text{RMG}}^{**}$, and take $\phi \in \tilde{S}_R$ such that $\phi^{i,r} = f$ for all $i \in \Omega$,

$r = 1, \dots, J$. Then

(1) $\tilde{R}(\phi) = R^*$.

(2) The collection of subchains of $\tilde{P}(\phi)$ is given by:

$$(4.4) \quad \left\{ \bigcup_{k=1}^{\infty} R^{*\alpha, r+kJ} \mid \alpha = 1, \dots, n^*; r = 1, \dots, \text{g.c.d.}(J, d(\alpha)) \right\}$$

(3) Each of the $R^{*\alpha, \beta}$ ($\alpha=1, \dots, n^*$; $\beta=1, \dots, d(\alpha)$) is a cyclically moving subset of $\tilde{P}(\phi)$.

(f) $\tilde{R}^* = R^*$.

$$(g) \quad \{\tilde{R}^{*\gamma} \mid \gamma = 1, \dots, \tilde{n}^*\} = \left\{ \bigcup_{k=1}^{\infty} R^{*\alpha, r+kJ} \mid \alpha = 1, \dots, n^*; r = 1, \dots, \text{g.c.d.}(J, d(\alpha)) \right\}$$

i.e. $\tilde{n}^* = \sum_{\alpha=1}^{n^*} \text{g.c.d.}(J, d(\alpha)) \geq n^*$.

(h) $\{\tilde{R}^{*\alpha, \beta}\} = \{R^{*\alpha, \beta}\}$; i.e. fix $\alpha \in \{1, \dots, n^*\}$. Then

$$\tilde{d}(\beta) = d(\alpha) / \text{g.c.d.}(J, d(\alpha)) \text{ for all } \tilde{R}^{*\beta} \subseteq R^{*\alpha}.$$

PROOF

(a) Let $\phi \in \tilde{S}_{\text{RMG}}$. Observe that $v(nJ) = Q^J v((n-1)J) \geq \tilde{q}(\phi) + \tilde{P}(\phi)v((n-1)J) \geq [I + \dots + \tilde{P}^{n-1}(\phi)] \tilde{q}(\phi) + \tilde{P}^n(\phi)v(0)$. Hence,

$$(4.5) \quad g^* = \lim_{n \rightarrow \infty} \frac{v(nJ)}{nJ} \geq \frac{\tilde{\Pi}(\phi)\tilde{q}(\phi)}{J} = \tilde{g}^*/J.$$

Next, let $f \in S_{\text{RMG}}$, and define $\phi \in \tilde{S}_R$, such that $\phi^{r,i} = f$ for all $i \in \Omega$, $r = 1, \dots, J$; Observe that

$$\begin{aligned} \tilde{g}^* &\geq \tilde{g}(\phi) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=0}^{n-1} \tilde{P}^k(\phi) \tilde{q}(\phi) = \\ &= \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=0}^{n-1} P^{kJ}(f) [I + \dots + P^{J-1}(f)] q(f) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=0}^{nJ-1} P(f)^k q(f) = \\ &= J(\Pi(f)q(f)) = Jg^* \end{aligned}$$

which together with (4.5) proves part (a).

(b): Recall that $g^* \geq P(f)g^*$ for any $f \in S_R$. If $\xi \in \tilde{L}(i)$, then for each r ,

$$\begin{aligned} P(f^1) \dots P(f^r) g_i^* &\leq g_i^* = \sum_j \tilde{P}_{ij}^\xi g_j^* = P(f^1) \dots P(f^r) [P(f^{r+1}) \dots P(f^J) g^*]_i \\ &\leq P(f^1) \dots P(f^r) g_i^*. \end{aligned}$$

Hence, $P(f^1) \dots P(f^r) g_i^* = g_i^*$. When $r = 1$, this implies $g_i^* = \sum_j P(f^1)_{ij} g_j^*$ and when $r \geq 2$, this implies that $[P(f^r) g_j^*] = g_j^*$ for all j , such that $P(f^1) \dots P(f^{r-1})_{ij} > 0$.

(c) Fix $v \in V$, and $i \in \Omega$, take $\xi = (f^1, \dots, f^J) \in \tilde{L}(i)$ and observe from part (b) that

$$\begin{aligned} v_i^* &\geq q(f^1)_i - g_i^* + [P(f^1)v^*]_i. \\ v_j^* &\geq q(f^2)_j - g_j^* + [P(f^2)v^*]_j, \text{ for all } j, \text{ such that } P(f^1)_{ij} > 0. \\ &\vdots \\ v_j^* &\geq q(f^J)_j - g_j^* + [P(f^J)v^*]_j, \text{ for all } j, \text{ such that } P(f^1) \dots P(f^{J-1})_{ij} > 0. \end{aligned}$$

Insert the J inequalities successively into each other and conclude that

$$v_i^* \geq \tilde{q}_i^\xi + \sum_j \tilde{P}_{ij}^\xi v_j^* - Jg_i^*, \text{ for all } \xi \in \tilde{L}(i)$$

whereas the equality sign holds for $\xi = (f^1, \dots, f^J)$ iff

$$\begin{aligned} (4.6) \quad b(v^*, f^1)_i &= 0 \\ b(v^*, f^r)_j &= 0 \text{ for all } j \text{ such that } P(f^1) \dots P(f^{r-1})_{ij} > 0; r = 2, \dots \end{aligned}$$

We conclude that $v_i^* + \tilde{g}_i^* = \max_{\xi \in \tilde{L}(i)} \left\{ \tilde{q}_i^\xi + \sum_j \tilde{P}_{ij}^\xi v_j^* \right\} = \tilde{T}v_i$ for all $i \in \Omega$,
or $v \in \tilde{V}$.

Hence $V \subseteq \tilde{V}$. The dimensions of V and \tilde{V} follow from th.5.5 in [13].

(d) Apply lemma 2.2 part (a) to the "J-step" MDP, and use the fact that $v \in \tilde{V}$ (cf. part (c)), in order to show that $\phi \in \tilde{S}_{RMG}$, iff

$$(4.7) \quad \tilde{\phi}_{i\xi} > 0 \Rightarrow \xi \in \tilde{L}(i) \quad \text{for all } i \in \Omega$$

$$\tilde{b}(v, \phi)_i = 0 \quad \text{for all } i \in \tilde{R}(\phi).$$

Use part (b), (4.6) and (4.2) in order to prove that (4.7) is equivalent to 4.3).

(e) Fix $\alpha \in \{1, \dots, n^*\}$ and $r, \beta \in \{1, \dots, d(\alpha)\}$ such that $\beta = r + kJ$ (modulo $d(\alpha)$) for some $k = 1, 2, \dots$

It then follows from th.3.1 part (j) and (2.5) that $P(f)_{ij}^{nd(\alpha)+kJ} > 0$

for all n sufficiently large, $i \in R^{*\alpha, r}$ and $j \in R^{*\alpha, \beta}$.

Since $P(\phi) = P(f)^J$, it follows that

$\tilde{P}(\phi)_{ij}^{nJ+k} > 0$, for all n sufficiently large, $i \in R^{*\alpha, r}$ and $j \in R^{*\alpha, \beta}$

which shows that all the states in each of the sets in (4.4) communicate with each other for $\tilde{P}(\phi)$. In addition, we observe, using th.3.1 part (g) that each of the sets in (4.4) is closed under $\tilde{P}(\phi)$ as well which proves that all of these sets are subchains of $\tilde{P}(\phi)$, and $\tilde{R}(\phi) \supseteq R^*$. We complete the proof of parts (e) (1) and (2), by showing the reversed inclusion $R(\phi) \subseteq R^*$, merely noting that for all $i \in \Omega \setminus R^*$,

$$\tilde{\Pi}(\phi)_{ii} = \lim_{n \rightarrow \infty} \tilde{P}(\phi)_{ii}^n = \lim_{n \rightarrow \infty} P(f)_{ii}^{nJ} = 0.$$

We next fix $\alpha \in \{1, \dots, n^*\}$, $\beta \in \{1, \dots, d(\alpha)\}$ and a state $i \in R^{*\alpha, \beta}$. Observe from th.3.1 part (j) that $R^{*\alpha, \beta}$ is a cyclically moving subset of $P(f)$ and use (2.4) and (2.6) in order to show

$$(4.8) \quad P(f)_{ii}^n > 0 \Rightarrow P(f)_{ij}^n > 0 \text{ for all } n \text{ sufficiently large, and all } j \in R^{*\alpha, \beta}$$

$$P(f)_{ii}^n > 0 \Rightarrow P(f)_{ij}^n = 0 \text{ for all } n = 1, 2, \dots \text{ and } j \notin R^{*\alpha, \beta}.$$

Note, using $\tilde{P}(\phi) = P(f)^J$ that (4.8) holds for $\tilde{P}(\phi)$ as well and conclude that each of the $R^{*\alpha, \beta}$ is a cyclically moving subset of $\tilde{P}(\phi)$, thus proving part (e) (3).

(f), (g) and (h) Fix $\phi \in \tilde{S}_{RMG}$ and let \tilde{C} be a subchain of $\tilde{P}(\phi)$.

Define

$$\bar{T} = \bigcup_{i \in \tilde{C}} \bigcup_{r=1}^J T^{r,i}(\phi)$$

(cf. (4.2)) and observe that

$$\tilde{C} = \bigcup_{i \in \tilde{C}} T^{J,i}(\phi),$$

hence

$$(4.9) \quad \tilde{C} \subseteq \bar{T}.$$

For each $j \in \bar{T}$, let $A_j = \{(r,i) \mid 0 \leq r < J, i \in \tilde{C} \text{ and } j \in T^{r,i}(\phi)\}$.

Next fix $v \in V$, and define $f \in S_R$ such that

$$\{k \mid f_{jk} > 0\} = \begin{cases} \bigcup_{(r,i) \in A_j} \{k \mid \phi_{jk}^{r+1,i} > 0\} & \text{for } j \in \bar{T}, \\ \{k \in L(j) \mid b(v)_j^k = 0\} & \text{for } j \notin \bar{T}. \end{cases}$$

Use part (d) in order to show that for all $i \in \Omega$:

$b(v,f)_i = 0$ and $f_{ik} > 0$ only for $k \in L(i)$, hence $f \in S_{RMG}$ via lemma (2.2) part (a). Since \bar{T} is closed, and the states in \bar{T} communicate with each other for $P(f)$, we conclude, that \bar{T} is a subchain. This implies using lemma 2.1 part (e) that

$$(4.10) \quad \tilde{C} \subseteq \bar{T} \subseteq R^{*\alpha} \text{ (for one } \alpha, 1 \leq \alpha \leq n^*)$$

which proves $\tilde{R}^* \subseteq R^*$ and hence *part (f)*, the reversed inclusion $\tilde{R}^* \supseteq R^*$ following from part (e) (1).

Next, fix $i \in \tilde{C}$. We then have in view of (4.10) that $i \in R^{*\alpha,\beta}$ (for some $\beta, 1 \leq \beta \leq d(\alpha)$). Use the fact that $\bar{T} \subseteq R^{*\alpha}$, and th.3.1 part (g) in order to show successively that

$$T^{r,i}(\phi) \subseteq R^{*\alpha,\beta+r} \quad \text{for } r = 1, \dots, J.$$

In particular, we obtain that

$$(4.11) \quad \{j \mid \tilde{P}(\phi)_{ij} > 0\} = T^{J,i}(\phi) \subseteq R^{*\alpha,\beta+J} \text{ such that}$$

$$\tilde{C} = \{j \mid \tilde{P}(\phi)_{ij}^k > 0 \text{ for some } k = 1, 2, \dots\} \subseteq \bigcup_{k=1}^{\infty} R^{*\alpha,\beta+kJ}$$

which together with part (e) (2) proves *part (g)*, using lemma 2.1 part (f).

Finally, a repeated application of (4.11) shows that

$$\tilde{P}(\phi)_{ii}^n > 0 \Rightarrow \tilde{P}(\phi)_{ij}^n = 0 \quad \text{for all } j \notin R^{*\alpha, \beta}, \text{ and all } n=1,2,\dots$$

which in view of (2.6) shows that each of the cyclically moving subsets of each of the policies in \tilde{S}_{RMG} lies within *one* $R^{*\alpha, \beta}$. This, in combination with part (e) (3), proves *part (h)*, using th.3.1 part (k). \square

REMARK 3. It is well known from Markov Chain Theory that the chain structure of the J -th power of a single stochastic matrix $P(f)$ is related to the chain structure of $P(f)$ in the following way:

- (a) the states that are transient (recurrent) for $P(f)$ are transient (recurrent) for $P^J(f)$.
 - (b) One obtains the subchain of $P^J(f)$ as follows
for each subchain $C^m(f)$ ($m=1,\dots,n(f)$), partition the collection of cyclically moving subsets $\{C^{\alpha, \beta}(f) \mid \beta = 1,\dots,d^m(f)\}$ (where the numbering of the c.m.s. satisfies (2.4)) into g.c.d. $\{J, d^m(f)\}$ subcollections, such that
 - (1) each of the subcollections contains exactly $d^m(f)/\text{g.c.d.}\{J, d^m(f)\}$ c.m.s.
 - (2) the rank numbers of the c.m.s. within the same subcollection differ a multiple of $\text{g.c.d.}\{J, d^m(f)\}$ (modulo $d^m(f)$).
 - (c) the collection of all the c.m.s. of $P(f)$ and the one of $P^J(f)$ coincide.
- Parts (f), (g) and (h) of the previous theorem show that the same correspondence holds with respect to the chain structure of the set of "one-step" maximal gain policies, and the one of the set of "J-step" maximal gain policies.

Consider, for instance, the "2-step" MDP in example 1:

case	J	\tilde{n}^*	\tilde{R}^{*1}	$\tilde{d}(1)$	$\tilde{R}^{*1,1}$	\tilde{R}^{*2}	$\tilde{d}(2)$	$\tilde{R}^{*2,1}$	$\tilde{R}^{*2,2}$
3	2	3	{1}	1	{1}	{2,5}	1	{2,5}	—
4	2	3	{1}	1	{1}	{2,5}	2	{2}	{5}
3	4	3	{1}	1	{1}	{2,5}	1	{2,5}	—
4	4	5	{1}	1	{1}	{2}	1	{2}	—

\tilde{R}^{*3}	$\tilde{d}(3)$	$\tilde{R}^{*3,1}$	$\tilde{R}^{*3,2}$	\tilde{R}^{*4}	$\tilde{d}(4)$	$\tilde{R}^{*4,1}$	\tilde{R}^{*5}	$\tilde{d}(5)$	$\tilde{R}^{*5,1}$
$\{3,4\}$	1	$\{3,4\}$	—	—	—	—	—	—	—
$\{3,4\}$	2	$\{3\}$	$\{4\}$	—	—	—	—	—	—
$\{3,4\}$	1	$\{3,4\}$	—	—	—	—	—	—	—
$\{3\}$	1	$\{3\}$	—	$\{4\}$	1	$\{4\}$	$\{5\}$	1	$\{5\}$

Table 3:

(Verify that $\tilde{n}^* = \sum_{\alpha=1}^2 \text{g.c.d. } \{J, d(\alpha)\}$ and that $\tilde{d}(\alpha) = d(2)/\text{g.c.d. } \{J, d(2)\}$ for $\alpha = 2, \dots, \tilde{n}^*$.)

Define $d^* = \text{least common multiple of } \{d(\alpha) \mid \alpha = 1, \dots, \tilde{n}^*\}$.

The following corollary will be needed for the analysis of the asymptotic behaviour of $v(n)$:

COROLLARY 4.2. Let $J = d^*$. Then

- (a) $\{\tilde{R}^\gamma \mid \gamma = 1, \dots, \tilde{n}^*\} = \{\tilde{R}^{*\alpha, \beta} \mid \alpha = 1, \dots, \tilde{n}^* : \beta = 1, \dots, d(\alpha)\}$.
- (b) $\tilde{n}^* = \sum_{\alpha=1}^{\tilde{n}^*} d(\alpha)$.
- (c) $\tilde{d}(\gamma) = 1$ for all $\gamma = 1, \dots, \tilde{n}^*$.

§5. THE ASYMPTOTIC BEHAVIOUR OF $v(n)$

In this section we study the asymptotic behaviour of $v(n)$. We show that $\{v(nJ+r) - (nJ+r)g^*\}_{n=1}^\infty$ converges for every final reward vector $v(0)$, if and only if J is a multiple of d^* , and as a consequence that $\{v(n) - ng^*\}_{n=1}^\infty$ converges for every vector $v(0)$ if and only if there exists an aperiodic randomized maximal gain policy that has R^* as its set of recurrent states.

THEOREM 5.1.

- (a) $\{v(n) - ng^*\}_{n=1}^\infty$ is bounded.
- (b) (cf. LANERY [6] proposition 7). If $f \in S_{\text{RMG}}$, and C is a subchain of $P(f)$,

with period d , then $\lim_{n \rightarrow \infty} [v(nd+r) - (nd+r)g^*]_i$ exists for all $i \in C$,
 $= 0, \dots, d-1$ and $v(0) \in E^N$.

(c) $\lim_{n \rightarrow \infty} [v(nd(\alpha)+r) - (nd(\alpha)+r)g^*]_i$ exists for all $i \in R^{*\alpha}$, $\alpha = 1, \dots, n^*$,
 $r = 1, \dots, d(\alpha)$ and $v(0) \in E^N$.

(d) $\lim_{n \rightarrow \infty} [v(nd^*+r) - (nd^*+r)g^*]_i$ exists for all $i \in \Omega$, $r = 1, \dots, d^*$ and all
 $v(0) \in E^N$.

PROOF.

(a) cf. BROWN [3] (corr. 4.3) and SCHWEITZER & FEDERGRUEN [14], remark 1.

(b) Note that

$$v(n+1)_i \geq q(f)_i + P(f)v(n)_i, \quad i \in C.$$

$$(n+1)g_i^* = g_i^* + nP(f)g_i^*, \quad i \in C, \text{ since } f \in S_{\text{RMG}} \text{ (cf. lemma 2.2, part (a)).}$$

$$v_i^* = q(f)_i - g_i^* + P(f)v_i^*, \quad i \in C, \text{ for any } v^* \in V.$$

Fix $v^* \in V$, let $e(n) = v(n) - ng^* - v^*$, and subtract the above equalities from the inequality, in order to get $e(n+1)_i \geq P(f)e(n)_i$, $i \in C$ and by induction

$$(5.1) \quad e(md+nd+r)_i \geq P(f)^{md} e(nd+r)_i, \quad i \in C.$$

It follows from part (a) that the sequence $\{v(nd+r)_i - (nd+r)g_i^*\}_{n=1}^\infty$ and hence $\{e(nd+r)_i\}_{n=1}^\infty$ has at least one cluster point.

For any $i \in C$, let x_i and y_i be two cluster points of the sequence $\{e(nd+r)_i\}_{n=1}^\infty$

Consider (sub)sequences $\{n_k\}_{k=1}^\infty$ and $\{m_k\}_{k=1}^\infty$ of the sequence of positive integers, such that $\lim_{k \rightarrow \infty} e(n_k d+r)_i = x_i$, $i \in C$ and $\lim_{k \rightarrow \infty} e(m_k d+n_k d+r)_i = y_i$, $i \in C$. Replace in (5.1) n and m by n_k and m_k , and let k tend to infinity, in order to conclude

$$(5.2) \quad y_i \geq \sum_{j \in C} \bar{\pi}_{ij} x_j, \quad i \in C$$

where $\bar{\pi}_{ij} = \lim_{n \rightarrow \infty} P(f)_{ij}^{nd}$; $i, j \in C$.

Multiply (5.2) by $\bar{\pi} \geq 0$ to get $\bar{\pi}y \geq \bar{\pi}x$. Since x and y are arbitrary cluster points, we have the reversed inequality $\bar{\pi}x \geq \bar{\pi}y$ as well, hence $\bar{\pi}x = \bar{\pi}y$. As a consequence, (5.2) becomes

$$y_i \geq \sum_{j \in C} \bar{\pi}_{ij} y_j, \quad i \in C.$$

Multiply these inequalities by $\bar{\pi} \geq 0$, and note $\bar{\pi}_{ii} > 0$, for all $i \in C$ (cf. (2.6)), to conclude that

$$y_i = [\bar{\pi}y]_i, \quad i \in C.$$

Thus,

$$y_i = \bar{\pi}y_i = \bar{\pi}x_i = x_i \quad \text{for all } i \in C$$

which proves that $\{e(nd+r)_i\}_{n=1}^{\infty}$ has exactly one cluster point, for any $i \in C$.

(c) Take f^* as in lemma (2.2) part (d), and apply part (b), using th.3.1 part (a).

(d) It suffices to prove that $\lim_{n \rightarrow \infty} [Q^{nd^*} v(0) - nd^* g^*]$ exists for all $v(0)$, because then $\lim_{n \rightarrow \infty} [v(nd^* + r) - (nd^* + r)g^*] = \lim_{n \rightarrow \infty} [Q^{nd^*} v(r) - nd^* g^*] - rg^*$ will also exist for all $v(0)$ and all r .

Define $\tilde{Q} = Q^{d^*}$ and consider the d^* -step MDP, as described in section 4.

Note $v(nd^*) - nd^* g^* = \tilde{Q}^n v(0) - ng^*$ (cf. th.4.1 part (a)). Fix $v(0)$ and define

$$x_i = \liminf_{n \rightarrow \infty} [\tilde{Q}^n v(0) - ng^*]_i; \quad X_i = \limsup_{n \rightarrow \infty} [\tilde{Q}^n v(0) - ng^*]_i, \quad i \in \Omega.$$

From part (a), it follows that $-\infty < x_i \leq X_i < \infty$ for all i .

Observe, using (2.11) that for all n sufficiently large

$$\begin{aligned} (5.3) \quad [\tilde{Q}^{n+1} v(0) - (n+1)g^*]_i &= [\tilde{T}\tilde{Q}^n v(0) - (n+1)g^*]_i = [\tilde{T}[\tilde{Q}^n v(0) - ng^*] - g^*]_i \\ &= \max_{\xi \in \tilde{L}(i)} \{ \tilde{q}_i^\xi - g_i^* + \sum_j \tilde{p}_{ij}^\xi [\tilde{Q}^n v(0) - ng^*]_j \}, \\ &\quad i \in \Omega. \end{aligned}$$

Fix $i \in \Omega$, take (sub)sequences $\{n_k\}_{k=1}^\infty$ (with $\lim_{k \rightarrow \infty} n_k = \infty$) such that $\lim_{k \rightarrow \infty} [\tilde{Q}^{n_k} v(0) - n_k \tilde{g}^*]$ exists and $\lim_{k \rightarrow \infty} [\tilde{Q}^{n_k+1} v(0) - (n_k+1) \tilde{g}^*]_i = x_i$ (or X_i resp.). Replace n by n_k in (5.3), and let k tend to infinity in order to conclude

$$(5.4) \quad X_i \leq \max_{\xi \in \tilde{L}(i)} [\tilde{q}_i^\xi - \tilde{g}_i^* + \sum_j \tilde{P}_{ij}^\xi X_j], \quad i \in \Omega.$$

$$(5.5) \quad x_i \geq \max_{\xi \in \tilde{L}(i)} [\tilde{q}_i^\xi - \tilde{g}_i^* + \sum_j \tilde{P}_{ij}^\xi x_j], \quad i \in \Omega.$$

If ϕ achieves the N maxima in (5.4), we have

$$(5.6) \quad \tilde{q}(\phi) - \tilde{g}^* + \tilde{P}(\phi)x \leq x \leq X \leq \tilde{q}(\phi) - \tilde{g}^* + \tilde{P}(\phi)X$$

or
$$0 \leq X - x \leq \tilde{P}(\phi)(X - x),$$

whence we get, by iterating this inequality

$$0 \leq X - x \leq \tilde{\Pi}(\phi)(X - x).$$

We complete the proof of showing $X - x = 0$ by demonstrating that $(X - x)_i = 0$ for all $i \in \tilde{R}(\phi)$.

Multiply the right inequality in (5.6) by $\tilde{\Pi}(\phi) \geq 0$, noting that ϕ has support on $X_{i \in \Omega} \tilde{L}(i)$, in order to get

$$0 \leq \tilde{\Pi}(\phi) [\tilde{q}(\phi) - \tilde{g}^*] = \tilde{g}(\phi) - \tilde{g}^* \leq 0,$$

where the last inequality follows from (2.8). Hence $\phi \in \tilde{S}_{\text{RMG}}$ and $\tilde{R}(\phi) \subseteq \tilde{R}^* = R^*$ (cf. th.4.1 part (f)) which proves $(X - x)_i = 0$, $i \in \tilde{R}(\phi)$, since part (c) shows that $(X - x)_i = 0$ for all $i \in R^*$. \square

We next show that the sequences $\{v(nJ+r) - (nJ+r)g^*\}_{n=1}^\infty$ do not converge for all final reward vectors $v(0)$, unless J is a multiple of d^* .

However, we first need the following lemma.

LEMMA 5.2. Define $\tilde{Q} = Q^{d^*}$, and consider the corresponding " d^* -step" MDP. Let \tilde{T}, \tilde{V} be defined as in section 4, and fix $v \in V$.

(a) For all $\tilde{v} \in \tilde{V}$, we have

$\tilde{v} = v + x$, where there are \tilde{n}^* constants $\{y^{\alpha,\beta} \mid \alpha = 1, \dots, n^*; \beta = 1, \dots, d(\alpha)\}$ with the convention that the superscript β in $y^{\alpha,\beta}$ is taken modulo $d(\alpha)$, such that for all $\alpha \in \{1, \dots, n^*\}$, and $\beta \in \{1, \dots, d(\alpha)\}$:

$$(5.7) \quad x_i = y^{\alpha,\beta} \quad \text{for all } i \in R^{*\alpha,\beta},$$

$$(5.8) \quad (T^{\tilde{m}} \tilde{v})_i = v_i + m g_i^* + y^{\alpha,\beta+m} \quad \text{for all } i \in R^{*\alpha,\beta}; m = 0, 1, 2, \dots$$

(b) $\tilde{v} \in \tilde{V}$ can be chosen such that all the $y^{\alpha,\beta}$ are distinct.

PROOF.

(a) Observe, using th.4.1 part (c) that $v \in \tilde{V}$, and use th.5.1 of SCHWEITZER & FEDERGRUEN [13] in order to show (5.7).

Next, take $f \in S_{RMG}^*$ and observe, using lemma (2.2) part (a) that

$$(5.9) \quad T^{\tilde{m}} \tilde{v} \geq q(f) + P(f) T^{\tilde{m}-1} \tilde{v}, \quad m = 1, \dots, d^*.$$

Using the fact that $\tilde{v} \in \tilde{V}$ and inserting the d^* inequalities in (5.9) successively into each other, we obtain

$$(5.10) \quad \tilde{v} + d^* g^* = \tilde{T} \tilde{v} \geq T^{d^*} \tilde{v} \geq [I + \dots + P(f)^{d^*-1}] q(f) + P(f)^{d^*} \tilde{v}.$$

By multiplying (5.10) with $\Pi(f) \geq 0$, we conclude strict equality for all components $i \in R^*$. It next follows from (5.9) that

$$T^{d^*} \tilde{v}_i = [q(f) + P(f) T^{d^*-1} \tilde{v}]_i \quad \text{for all } i \in R^*,$$

and more generally that

$$(5.11) \quad [T^k \tilde{v}]_i = [q(f) + P(f) T^{k-1} \tilde{v}]_i \quad \text{for all } k = 1, \dots, d^* \text{ and}$$

$$i \in \{i \mid P(f)_{ji}^{d^*-k} > 0 \text{ for some } j \in R^*\} = R^*,$$

where the last equality follows from $R(f) = R^*$.

We next prove (5.8) for $m = 0, \dots, d^*$. It then follows that (5.8) holds for any value of m , since for all $m = 1, 2, \dots$ and $m = 1, \dots, d^*$

$$\begin{aligned} T^{nd^*+m} \tilde{v}_i &= T^m(T^{nd^*} \tilde{v})_i = T^m(\tilde{v} + nd^* g^*)_i = nd^* g_i^* + T^m \tilde{v}_i = \\ &= nd^* g_i^* + v_i + mg_i^* + y^{\alpha, \beta+m} = v_i + (nd^*+m)g_i^* + y^{\alpha, \beta+nd^*+m} \end{aligned}$$

for all $i \in R^{*\alpha, \beta}$.

First observe that (5.8) holds for $m = 0$. Next assume it holds for $m = k$, with $0 \leq k < d^*$. It then follows that (5.8) holds for $m = k + 1$, as well since, using (5.11), and th.3.1 part (g)

$$\begin{aligned} (T^{k+1} \tilde{v})_i &= [q(f) + P(f) T^k \tilde{v}]_i = \\ &= q(f)_i + \sum_{j \in R^{*\alpha, \beta+1}} P(f)_{ij} \{v_j + kg_j^* + y^{\alpha, \beta+k+1}\} \\ &= 0 + v_i + (k+1)g_i^* + y^{\alpha, \beta+k+1}. \end{aligned}$$

(b) It follows from th.5.5 in SCHWEITZER & FEDERGRUEN [13] that the \tilde{n}^* parameters $\{y^{\alpha, \beta} \mid \alpha = 1, \dots, n^*; \beta = 1, \dots, d(\alpha)\}$ may be chosen independently over some (finite) region in $E^{\tilde{n}^*}$.

THEOREM 5.3.

- (a) Fix $\alpha \in \{1, \dots, n^*\}$, $i \in R^{*\alpha}$, $J \geq 1$, and $r \in \{0, \dots, J-1\}$.
 $\lim_{n \rightarrow \infty} v(nJ+r)_i - (nJ+r)g_i^*$ exists for all $v(0)$,
 only if J is a multiple of $d_i = d(\alpha)$.
- (b) Fix $J \geq 1$ and $r \in \{0, \dots, J-1\}$.
 $\lim_{n \rightarrow \infty} v(nJ+r) - (nJ+r)g^*$ exists for all $v(0) \in E^N$ only if J is a multiple of d^* .

PROOF.

- (a) Fix $v \in V$, and choose $\tilde{v} \in \tilde{V}$ as in part (b) of the previous lemma. Pick t large enough that $Q^n(\tilde{v} + tg^*) = T^n(\tilde{v} + tg^*)$, for $n = 1, 2, \dots$ (cf. (2.10)). Finally, let $i \in R^{*\alpha, \beta}$ ($1 \leq \beta \leq d(\alpha)$). Observe that $\tilde{v} + tg^* \in \tilde{V}$, and apply lemma 5.2 part (a) in order to show

$$Q^{nJ+r}(\tilde{v} + tg^*)_i = T^{nJ+r}(\tilde{v} + tg^*)_i = tg_i^* + v_i + (nJ+r)g_i^* + y^{\alpha, \beta+nJ+r}.$$

Hence,

$$Q^{nJ+r}(\tilde{v}+tg^*)_i - (nJ+r)g_i^* = v_i + tg_i^* + y^{\alpha, \beta+nJ+r}.$$

Since $\lim_{n \rightarrow \infty} Q^{nJ+r}(\tilde{v}+tg^*)_i - (nJ+r)g_i^*$ exists and since the $y^{\alpha, \beta}$ ($\alpha=1, \dots, n^*$; $\beta=1, \dots, d(\alpha)$) are chosen to be distinct, we must have $\beta + nJ + r$ (modulo $d(\alpha)$) = γ (say) for all n large enough, which implies that J is a multiple of $d(\alpha)$.

- (b) Since $\lim_{n \rightarrow \infty} [v(nJ+r) - (nJ+r)g^*]_i$ exists for all $i \in R^*$ and $v(0) \in E^N$, it follows from part (a) that J must be a multiple of the $d(\alpha)$ ($\alpha=1, \dots, n^*$) hence J is a multiple of d^* .

Combining th.5.1 parts (c) and (d), with th.5.3, we obtain our main result.

THEOREM 5.4.

- (a) Fix $\alpha \in \{1, \dots, n^*\}$, $i \in R^{*\alpha}$, and two integers J and r . Then $\lim_{n \rightarrow \infty} v(nJ+r)_i - (nJ+r)g_i^*$ exists for all $v(0) \in E^N$, if and only if J is a multiple of $d(\alpha) = d_i^*$.
- (b) $\lim_{n \rightarrow \infty} v(nJ+r) - (nJ+r)g^*$ exists for all $v(0) \in E^N$, if and only if J is a multiple of d^* .

REMARK 4. The following conditions are equivalent statements of the necessary and sufficient condition for the convergence of $\{v(n) - ng^*\}_{n=1}^\infty$, for all $v(0) \in E^N$.

- (I) $d^* = 1$.
- (II) There exists an aperiodic *randomized* maximal gain policy f , with $R(f) = R^*$.
- (III) Each state $i \in R^*$ lies within an aperiodic subchain of some *randomized* maximal gain policy.
- (IV) For each $\alpha \in \{1, \dots, n^*\}$ there exists a *randomized* maximal gain policy which has an aperiodic subchain within $R^{*\alpha}$.

(Observe that (I) \Rightarrow (II) as a result of th.3.1 part (a), (II) \Rightarrow (III), and (III) \Rightarrow (IV) are immediate, whereas (IV) \Rightarrow (I) is immediate from (3.1)).

We note that in (II), (III) and (IV) the adjective "randomized" cannot be replaced by "pure"; in fact, the modification of example 1, case 1 where $K(5) = \{1, 2\}$ shows that $d^* = 1$ can occur, with *all* of the *pure* policies being periodic.

Moreover, example 1, case 1, case 3 and 4 show that the addition "with $R(f) = R^*$ " in (II) is indispensable: f^6 is an aperiodic maximal gain policy, however with $R(f^6) \subset R^*$.

Finally example 1, case 3, with $d^* = 2$, shows that $\lim_{n \rightarrow \infty} v(n) - ng^*$ fails to exist for some $v(0) \in E^5$ (take $v(0) = [2q_5^2 \ q_5^2 \ 0 \ 0 \ q_5^2]$, observe that $v(2n+1) = [2q_5^2 \ 0 \ q_5^2 \ q_5^2 \ 0]$ and $v(2n) = [2q_5^2 \ q_5^2 \ 0 \ 0 \ q_5^2]$. Note that $v(0) \in \tilde{V} - V$ and cf. th.5.3).

THEOREM 5.5. *The following conditions are sufficient for the existence of $\lim_{n \rightarrow \infty} [v(n) - ng^*]$ for all $v(0) \in E^N$.*

(I) *All of the transition probabilities are strictly positive.*

$$P_{ij}^k > 0, \quad \text{for all } i, j \in \Omega, \quad \text{and } k \in K(i)$$

(cf. BELLMAN [2], BROWN [3]).

(II) *For all $v(0) \in E^N$, there exists an aperiodic $f \in S_P$, and an integer n_0 , such that*

$$v(n+1) = q(f) + P(f)v(n), \quad \text{for all } n \geq n_0 \text{ (cf. MORTON [8]).}$$

(III) *There exists a state s and an integer $v \geq 1$, such that*

$$P(f^1) \dots P(f^v)_{is} > 0 \quad \text{for all } f^1, f^2, \dots, f^v \in S_P; i \in \Omega$$

(cf. WHITE [16]).

(IV) *Every $f \in S_P$ is aperiodic (cf. SCHWEITZER [11] and [12]).*

(V) *Every $f \in S_{PMG}$ is aperiodic (cf. SCHWEITZER [11] and [12]).*

(VI) *For each $i \in R^*$, there exists a pure maximal gain policy f , such that state i is recurrent and aperiodic for $P(f)$.*

(VII) *Every pure maximal gain policy has a unichained tpm, and at least one of them is aperiodic.*

PROOF. (I) \Rightarrow (III) \Rightarrow (IV) \Rightarrow (V) \Rightarrow (VI) where the last implication follows from lemma (2.1) part (a).

(VI) $\Rightarrow d_i^* = 1$ for all $i \in R^* \Rightarrow d^* = d(\alpha) = 1$ for all $\alpha = 1, \dots, n^*$ (cf. th.3.1. part (c)). The sufficiency of (II) follows from the fact that after n_0 iterations the policy space may be reduced to $S_P^{\text{new}} = \{f\}$ which satisfies (IV).

(VII) $\Rightarrow n^* = 1$, since the subchains of any two tpm's must intersect, and in addition $d^* = d(1) = 1$ as a consequence of th.3.2.

We have seen that for arbitrary $J \geq 1$, and some fixed $v(0)$ the sequences $\{v(nJ+r)_i - (nJ+r)g_i^*\}_{n=1}^\infty$ may fail to converge for some (or all) $i \in \Omega$ and for some (or all) $r = \{0, 1, \dots, J-1\}$.

However, the various sequences interdepend as far as their asymptotic behaviour is concerned.

We conclude this section by exhibiting the various ways in which this interdependence occurs. However we first need the following lemma.

LEMMA 5.6. Fix $f \in S_{\text{RMG}}$.

$$\lim_{n \rightarrow \infty} [v(n+1)_i - q(f)_i - P(f)v(n)_i] = 0, \quad \text{for all } i \in R(f).$$

PROOF. Use the fact that for all $i \in \Omega$, $f_{ik} > 0$ only for $k \in L(i)$ (cf. lemma 2.2 part (a)) in order to show that

$$(5.12) \quad v(n+1) - (n+1)g^* \geq q(f) - g^* + P(f)[v(n) - ng^*].$$

By multiplying (5.12) with $\Pi(f)$, we obtain

$$\Pi(f)(v(n+1) - (n+1)g^*) \geq \Pi(f)(v(n) - ng^*).$$

Observing from th. 5.1 part (a) that $\Pi(f)(v(n) - ng^*)$ is bounded in n , we conclude the existence of $L = \lim_{n \rightarrow \infty} \Pi(f)(v(n) - ng^*)$. Define

$$\delta(n) = v(n+1) - q(f) - P(f)v(n)$$

and note that $\delta(n) \geq 0$ for all n (cf. (1.1)).

Thus,

$$\begin{aligned} \lim_{n \rightarrow \infty} \Pi(f) \delta(n) &= \lim_{n \rightarrow \infty} \{ \Pi(f) [v(n+1) - (n+1)g^*] - \Pi(f)(q(f)-g^*) - \\ &- \Pi(f)(v(n)-ng^*) \} = L - L = 0, \end{aligned}$$

which proves the lemma using $\delta(n) \geq 0$ and the fact that $\Pi(f) \geq 0$ with $\Pi(f)_{jj} > 0$ for all $j \in R(f)$.

THEOREM 5.7.

(a) Fix $\alpha \in \{1, \dots, n^*\}$; $\lim_{n \rightarrow \infty} v(n)_i - ng_i^*$ exists either for all $i \in R^{*\alpha}$ or for none of them.

(b) Fix $J \geq 1$ and $t \in R^{*\alpha, \beta}$ ($1 \leq \alpha \leq n^*$; $1 \leq \beta \leq d(\alpha)$). Assume $\lim_{n \rightarrow \infty} v(nJ+r)_t - (nJ+r)g_t^*$ exists for some integer r . Then

$$(5.13) \quad \lim_{n \rightarrow \infty} v(nJ+r+s)_i - (nJ+r+s)g_i^* \text{ exists for all } i \in \bigcup_{k=1}^{R^{*\alpha, \beta+kJ-s}} (s=1, 2, \dots).$$

(c) Fix $J \geq 1$, and $\alpha \in \{1, \dots, n^*\}$.

$\lim_{n \rightarrow \infty} v(nJ+r)_i - (nJ+r)g_i^*$ exists for all $r = 1, \dots, J$, either for all $i \in R^{*\alpha}$ or for none of them.

(d) Fix $J \geq 1$, $r_0 \in \{1, \dots, J\}$ and $\alpha \in \{1, \dots, n^*\}$.

If $\lim_{n \rightarrow \infty} v(nJ+r_0)_i - (nJ+r_0)g_i^*$ exists for all $i \in R^{*\alpha}$,

then $\lim_{n \rightarrow \infty} v(nJ+r)_i - (nJ+r)g_i^*$ exists for all $i \in R^{*\alpha}$ and all $r = 1, 2, \dots$

(e) Fix $i \in \Omega$. Assume $\lim_{n \rightarrow \infty} v(nJ^1+r)_i - (nJ^1+r)g_i^*$, and

$\lim_{n \rightarrow \infty} v(nJ^2+s)_i - (nJ^2+s)g_i^*$ exist for all $r \in \{1, \dots, J^1\}$ and $s \in \{1, \dots, J^2\}$.

Let $J^3 = \text{g.c.d}\{J^1, J^2\}$. Then

$\lim_{n \rightarrow \infty} v(nJ^3+t)_i - (nJ^3+t)g_i^*$ exists for all $t = 1, \dots, J^3$,

and hence, if in addition $i \in R^{*\alpha}$ (for some $1 \leq \alpha \leq n^*$) then

$\lim_{n \rightarrow \infty} v(nJ^3+t)_j - (nJ^3+t)g_j^*$ exists for all $t = 1, \dots, J^3$ and all $j \in R^{*\alpha}$.

(f) Fix $i \in R^{*\alpha}$ ($1 \leq \alpha \leq n^*$) and $J \geq 1$. Assume

$\lim_{n \rightarrow \infty} v(nJ+r)_i - (nJ+r)g_i^*$ exists for all $r = 1, \dots, J$.

Let $\hat{J} = \text{g.c.d.}\{J, d(\alpha)\}$. Then

$\lim_{n \rightarrow \infty} v(n\hat{J}+s)_j - (n\hat{J}+s)g_j^*$ exists for all $s = 1, \dots, \hat{J}$ and all $j \in R^{*\alpha}$.

(g) Fix $J \geq 1$. If $\lim_{n \rightarrow \infty} v(nJ+r)_i - (nJ+r)g_i^*$ exists for all $i \in R^*$ and some $r \in \{1, \dots, J\}$, then $\lim_{n \rightarrow \infty} v(nJ+r)_i - (nJ+r)g_i^*$ exists for all $i \in \Omega$ and all $r = 1, 2, \dots$.

PROOF.

(a) Assume $\lim_{n \rightarrow \infty} v(n)_t - ng_t^*$ exists for some $t \in R^{*\alpha}$.

Define $x_i = \liminf_{n \rightarrow \infty} [v(n) - ng^*]_i$; $X_i = \limsup_{n \rightarrow \infty} [v(n) - ng^*]_i$, and observe that $-\infty < x_i \leq X_i < \infty$ as a result of th.5.1 part (a).

Fix $i \in R^{*\alpha}$, pick $\varepsilon > 0$ and apply lemma 5.6 with $f^* \in S_{\text{RMG}}^*$ in order to show that there exists an integer $n(\varepsilon)$, such that for all $n > n(\varepsilon)$

$$(5.14) \quad q(f^*)_i - g_i^* + P(f^*)[v(n) - ng^*]_i \leq v(n+1)_i - (n+1)g_i^* \leq q(f^*)_i - g_i^* + P(f^*)[v(n) - ng^*]_i + \varepsilon.$$

Take (sub)sequences $\{n_k\}_{k=1}^\infty$ (with $\lim_{k \rightarrow \infty} n_k = \infty$) such that $\lim_{k \rightarrow \infty} [v(n_k) - n_k g^*]$ exists and

$$\lim_{k \rightarrow \infty} v(n_k+1)_i - (n_k+1)g_i^* = x_i \text{ (or } X_i \text{)}.$$

Replace n by n_k in (5.14) and let k tend to infinity in order to conclude

$$q(f^*)_i - g_i^* + P(f^*) x_i \leq x_i \leq X_i \leq q(f^*)_i - g_i^* + P(f^*) X_i + \varepsilon,$$

or

$$0 \leq X_i - x_i \leq P(f^*)(X - x)_i, \quad \text{for all } i \in R^{*\alpha},$$

whence we get by iterating this inequality

$$(5.15) \quad 0 \leq X_i - x_i \leq \Pi(f^*)^\alpha, X - x, \quad i \in R^{*\alpha}.$$

Multiply this inequality by $\Pi(f^*) \geq 0$ in order to conclude strict equality on the right of (5.15), thus proving $X_i - x_i = X_t - x_t = 0$, for all i .

- (b) Without loss of generality, we take $r = 0$. Define $\tilde{Q} = Q^J$ and consider the J -step MDP, as defined in section 4. Let f^* be defined as in lemma 2.2 part (d), and let $\bar{R}(s) = \bigcup_{k=1}^{\infty} R^{*\alpha, \beta+kJ-s}$ for $s = 1, 2, \dots$.

Observe that $v(nJ)_i - nJg_i^* = [\tilde{Q}^n v(0)]_i - \tilde{n}g_i^*$ (cf. th.4.1 part (a)).

Apply part (a) of this theorem to the J -step MDP, and use th.4.1

part (g) in order to obtain that $v(nJ)_i - nJg_i^*$ exists for all $i \in \bar{R}(0)$, thus proving (5.13) for $s = 0$. Assume (5.19) holds for $s = S$.

Note, using th.3.1 part (f) that for all $i \in \bar{R}(S+1)$, $P(f^*)_{ij} > 0$ only for $j \in \bar{R}(S)$. It then follows from lemma 5.6 that for all $i \in \bar{R}(S+1)$

$$\lim_{n \rightarrow \infty} v(nJ+S+1)_i - (nJ+S+1)g_i^* = q(f^*)_i - g_i^* + \sum_{j \in \bar{R}(S)} P(f^*)_{ij} x_j$$

where

$$x_i = \lim_{n \rightarrow \infty} v(nJ+S)_i - (nJ+S)g_i^*, \quad \text{for all } i \in \bar{R}(S),$$

which proves part (b) by complete induction.

- (c) Assume $\lim_{n \rightarrow \infty} v(nJ+r)_t - (nJ+r)g_t^*$ exists for all $r = 1, \dots, J$ and $t \in R^{*\alpha, \beta}$

($1 \leq \beta \leq d(\alpha)$). Take $i \in R^{*\alpha, \gamma}$ ($1 \leq \gamma \leq d(\alpha)$) and $s \in \{1, \dots, J\}$. Then

$\lim_{n \rightarrow \infty} v(nJ+s)_i - (nJ+s)g_i^*$ exists as a result of part (b).

- (d) Take $i \in R^{*\alpha, \beta}$ ($1 \leq \beta < d(\alpha)$) and $r \in \{1, \dots, J\}$; $\lim_{n \rightarrow \infty} v(nJ+r)_i - (nJ+r)g_i^*$ exists as a result of part (b).

- (e) Let $p^1 = J^1/J^3$ and $p^2 = J^2/J^3$.

Fix $t \in \{1, \dots, J^3\}$, and define $a(n) = v(nJ^3+t)_i - (nJ^3+t)g_i^*$.

Observe that $A(m) = \lim_{n \rightarrow \infty} a(np^2+m)$ exists for all $m = 1, \dots, p^2$, just as

$A = \lim_{n \rightarrow \infty} a(np^1)$ exists. Observe that there exist two integers $\alpha, \beta \geq 1$ such

that $\alpha p^1 - \beta p^2 = 1$, as a consequence of p^1 and p^2 being relatively prime.

Since $A(m) = \lim_{k \rightarrow \infty} a[(kp^1 + \beta m)p^2 + m] = \lim_{k \rightarrow \infty} a[(kp^2 + \alpha)^m p^1] = A$, for all

$m = 1, \dots, p^2$, it follows that $\lim_{n \rightarrow \infty} a(n)$ exists, thus proving the first assertion, whereas the second one follows immediately from part (c).

- (f) Use part (e) with $J^1 = J$ and $J^2 = d(\alpha)$ (cf. th.5.1 part (c)).

- (g) It follows from part (c) that

$$\lim_{n \rightarrow \infty} v(nJ+r)_i - (nJ+r)g_i^*$$

exists for all $i \in R^*$ and all $r \in \{1, \dots, J\}$ whereas convergence on $\Omega - R^*$ is deduced, using the proof of th.5.1 part (d).

REMARK 5. The following statements illustrate the degree of interdependence with respect to the asymptotic behaviour of the N sequences $\{v(n)_i - ng_i^*\}$ ($i \in \Omega$), and may be proved using the above theorem, merely verifying all possible combinations.

- (a) $\lim_{n \rightarrow \infty} v(n)_i - ng_i^*$ cannot exist for all values of i , but one (cf. SCHWEITZER [12], th.1 part (3)).
- (b) If $\lim_{n \rightarrow \infty} v(n)_i - ng_i^*$ exists for all values of i except two, then these two special states comprise one $R^{*\alpha}$, with $d(\alpha) = 2$.
Moreover, for every randomized maximal gain policy these two states either form a periodic subchain, or are both transient.
- (c) If $\lim_{n \rightarrow \infty} v(n)_i - ng_i^*$ exists for all values of i except three, then either the three states comprise one $R^{*\alpha}$ with $d(\alpha) = 2$ or 3, or else two of them comprise one $R^{*\alpha}$ with $d(\alpha) = 2$, and the third one lies in $\Omega - R^*$, having positive probability to reach $R^{*\alpha}$.

The generalization of th.5.4 for the case of one fixed $v(0)$ is

THEOREM 5-8:

- (a) Fix $v(0)$, and $\alpha \in \{1, \dots, n^*\}$. There exists an integer $J^{0\alpha} \geq 1$, dependent upon $v(0)$, such that $\lim_{n \rightarrow \infty} [v(nJ+r) - (nJ+r)g_i^*]_i$ exists for all $i \in R^{*\alpha}$ and some r if and only if the integer $J \geq 1$ is a multiple of $J^{0\alpha}$. If this condition is met, the limit exists for all r . The integer $d(\alpha)$ is a multiple of $J^{0\alpha}$. If $d(\alpha) \geq 2$, then there exist choices of $v(0)$ such that $J^{0\alpha} < d(\alpha)$ can occur.
- (b) Fix $v(0)$ and define the integer

$$J^0 = \text{l.c.m. } \{J^{0\alpha} \mid 1 \leq \alpha \leq n^*\}$$

which depends upon $v(0)$. Then $\lim_{n \rightarrow \infty} v(nJ+r) - (nJ+r)g_i^*$ exists for some r if and only if the integer $J \geq 1$ is a multiple of J^0 . If this condition

is met, the limit exists for all r . The integer d^* is a multiple of J^0 . If $d^* \geq 2$, then there exist choices of $v(0)$ such that $J^0 < d^*$ can occur.

PROOF:

(a) Let

$$(5.16) \quad J^{0\alpha} = \text{g.c.d.} \{ J \geq 1 \mid \lim_{n \rightarrow \infty} [v(nJ+r)_i - (nJ+r)g_i^*] \text{ exists for all } i \in R^{*\alpha} \}$$

Observe that $J^{0\alpha}$ can be obtained as the g.c.d. of a finite number of integers and apply th.5-7 part (e) to conclude that $J^{0\alpha}$ belongs to the set to the right of (5.16), thus proving the first assertion. The second and third assertion follow from th.5-7 part (d) and th.5-1 part (c), whereas the last one may be verified by choosing

$$(5.17) \quad v(0) = v + tg^* \text{ with } v \in V \text{ and } t \text{ sufficiently large that}$$

$$Q^n(v(0)) = T^n(v(0)) \text{ for } n = 1, 2, \dots$$

(b) Observe from part (a) that $\lim_{n \rightarrow \infty} [v(nJ+r)_i - (nJ+r)g_i^*]$ exists for all $i \in R^*$, and some r if and only if J is a multiple of J^0 , and apply part (g) of th.5-7 to verify the first two assertions. The third assertion follows from th.5-1 part (d) whereas the existence of $v(0)$ with $J^0 = 1$ may be verified by choosing $v(0)$ as in (5-17).

§6. THE ASYMPTOTIC PROPERTIES OF THE POLICIES GENERATED BY THE Q-OPERATOR

In this final section, we indicate some of the properties of the policies that attain the N maxima in (1.2), and the resulting consequences with respect to the use of the value-iteration method.

First, define for any $\varepsilon \geq 0$, and $i \in \Omega$

$$K(i, n, \epsilon) = \{k \in K(i) \mid v(n)_i - \epsilon \leq q_i^k + \sum_j P_{ij}^k v(n-1)_j \leq v(n)_i\},$$

$$n = 1, 2, \dots$$

i.e. the set of actions, which at the n -th step of the value-iteration are ϵ -optimal, in the sense that they achieve the maximum in (1.2) within ϵ .

For any $v \in V$, define $L(i, v) = \{k \in L(i) \mid b(v)_i^k = 0\}$, $i \in \Omega$. Note from lemma 2-2 part (a) that $K^*(i) \subseteq L(i, v)$ for all $i \in R^*$.

THEOREM 6.1. Fix $v(0) \in E^N$.

(a) If $\lim_{n \rightarrow \infty} v(n) - ng^* = v^*$ exists, then

(1) $v^* \in V$.

(2) There exists an integer n_0 , such that if a policy $f \in S_R$ satisfies $v(n) = q(f) + P(f)v(n-1)$ for some $n \geq n_0$ then $f \in S_{RMG}$, with $b(v^*, f) = 0$.

(b) For any $\epsilon > 0$, and all $i \in R^*$,

$K^*(i) \subseteq K(i, n, \epsilon) \subseteq L(i)$ for all n sufficiently large,

(c) Let $\{\epsilon_n\}_{n=1}^{\infty}$ be a sequence of non-negative numbers, such that

(1) $\lim_{n \rightarrow \infty} \epsilon_n = 0$.

(2) $\lim_{n \rightarrow \infty} \lambda^n / \epsilon_n = 0$ for all $\lambda \in (0, 1)$.

If $\lim_{n \rightarrow \infty} v(n) - ng^* = v^*$ exists, then

$K(i, n, \epsilon_n) = L(i, v^*)$ for all n sufficiently large, and all $i \in \Omega$.

(d) If there exists a $f \in S_R$, and an integer n_0 , such that

$v(n+1) = q(f) + P(f)v(n)$, for all $n \geq n_0$, then $f \in S_{RMG}$.

PROOF.

(a) (1) cf. lemma 2.2 part (f) in [14].

(2) The fact that for large n , and any $i \in \Omega$, only alternatives in $L(i, v^*)$ can attain the maximum in (1.2) was shown in [14] (3.3); it then follows from lemma 2.2 part (a) that $f \in S_{RMG}$.

(b) Fix $i \in \Omega$.

Observe that $k \in K(i, n, \varepsilon)$ implies

$$v(n)_i - ng_i^* - \varepsilon \leq q_i^k - g_i^* + \sum_j p_{ij}^k \left[v(n-1)_j - (n-1)g_j^* \right] + (n-1) \left\{ \sum_j p_{ij}^k g_j^* - g_i^* \right\}$$

which cannot hold for large n , if $\sum_j p_{ij}^k g_j^* - g_i^* < 0$, since $[v(n) - ng^*]$ is bounded (cf. th.5.1 part (a)). Hence $K(i, n, \varepsilon) \subseteq L(i)$.

Next, fix $k \in K^*(i)$ and let $f \in S_{\text{PMG}}$ with $f_{ik} = 1$ and $i \in R(f)$ (cf. (2.12)). Apply lemma 5.6 in order to show that

$$\lim_{n \rightarrow \infty} v(n+1)_i - q_i^k - \sum_j p_{ij}^k v(n)_j = 0,$$

which proves $K^*(i) \subseteq K(i, n, \varepsilon)$ for n sufficiently large.

(c) Define $y(n) = v(n) - ng^* - v^*$, and note that $k \in K(i, n, \varepsilon_n)$, if and only if

$$(6.1) \quad (n-1) \left(\sum_j p_{ij}^k g_j^* - g_i^* \right) + b(v^*)_i^k \geq - \sum_j p_{ij}^k y(n-1)_j + y(n)_i - \varepsilon_n.$$

First let $k \in K(i, n, \varepsilon_n)$. Since the right hand side of (6.1) tends to zero as n tends to infinity it follows that for n large, $\sum_j p_{ij}^k g_j^* - g_i^* = 0$ (i.e. $k \in L(i)$) and $b(v^*)_i^k = 0$, which proves

$$(6.2) \quad K(i, n, \varepsilon_n) \subseteq L(i, v^*) \quad \text{for all } i \in \Omega, \text{ and } n \text{ sufficiently large.}$$

In order to prove the reversed inclusion, we recall from [14] that $y(n)$ converges geometrically to zero, i.e. there exist two numbers K and λ with $0 \leq \lambda < 1$ such that $|y(n)_i| \leq K\lambda^n$ for all $n = 1, 2, \dots$ and $i \in \Omega$. Let $k \in L(i, v^*)$. It then follows that k satisfies (6.1) since the left hand side of (6.1) equals zero, whereas its right hand side is strictly negative, for large n , as a consequence of $\lim_{n \rightarrow \infty} \lambda^n / \varepsilon_n = 0$.

This proves the reversed inclusion in (6.2) and hence part (c).

(d) Observe that $v(n_0+m) = [I + P(f) + \dots + P(f)^{m-1}]q(f) + P(f)^m v(n_0)$, divide this equality by m , let m tend to infinity and conclude that

$$g^* = \lim_{m \rightarrow \infty} \frac{v(n_0+m)}{m} = \Pi(f)q(f) = g(f). \quad \square$$

Part (a) of the above theorem shows that if $[v(n) - ng^*]$ converges, then the value-iteration method will generate only *maximal gain* policies, after a finite number of steps.

On the other hand, example 1 shows that a *non-maximal gain* policy may be generated infinitely often, although not at each iteration step (cf. part (d)), if $\lim_{n \rightarrow \infty} v(n) - ng^*$ fails to exist.

(Take case 4. Let $v(0) = [2q_5^2 \ q_5^2 \ 0 \ 0 \ q_5^2]$ as in Remark 4, and observe that the value-iteration method may generate the sequence of policies $(f^4, f^5, f^4, f^5, \dots)$, where $f^5 \notin S_{PMG}$.)

This phenomenon was first noticed by LANERY [6], example 4, where in the *only possible* sequence of policies to be generated by the value-iteration method, a *maximal* and a *non-maximal gain policy* alternate.

Moreover, the value-iteration method may even generate exclusively *non-maximal gain* policies, after a finite number of steps, if $v(n) - ng^*$ fails to converge.

(Take example 4 in LANERY [6] or example 1, case 4 with

$v(0) = [2q_5^2 \ q_5^2 \ 0 \ 0 \ q_5^2]$ and define a new MDP as follows.

Let $\bar{\Omega} = \{(i, r) \mid i \in \bar{\Omega}, r = 0, 1\}$ with $\bar{K}(i, r) = K(i)$; $\bar{q}_{(i, r)}^k = q_i^k$ and

$$\bar{P}_{(i, r)}^k(j, s) = \begin{cases} p_{ij}^k & \text{if } r = s \quad \text{for all } i, j \in \bar{\Omega}; r = 0, 1, \text{ and } k \in \bar{K}(i, r) \\ 0 & \text{if } r \neq s \end{cases}$$

and choose $\bar{v}(0)_{(i, 0)} = v(0)_i$ and $\bar{v}(0)_{(i, 1)} = v(1)_i$.

Lemma 5.6 shows that *every* maximal gain policy comes closer and closer to attaining the maxima in (1.2) even if $v(n) - ng^*$ fails to exist. Part (b) of the above theorem shows that this may be used in order to "localize" the sets $K^*(i)$ ($i \in R^*$) by determining at each iteration step the sets $K(n, i, \epsilon)$ (for some $\epsilon > 0$) rather than the sets $K(n, i, 0)$.

In case $\lim_{n \rightarrow \infty} v(n) - ng^* = v^*$ does exist part (c) of th.6.1 provides a method in order to determine the sets $L(i, v^*)(i \in \Omega)$ in the course of the value-iterations. Observe that any sequence $\{\epsilon_n\}_{n=1}^{\infty}$, which has ϵ_n^{-1} polynomially bounded in n , may be used, e.g. $\epsilon_n = 1/n$.

If $v^* = \lim_{n \rightarrow \infty} [v(n) - ng^*]$ exists, then the sets $K(i, n, \epsilon_n)$ may be better behaved for large n than the optimizing sets $K(i, n, 0)$.

The former set has a limit $L(i, v^*)$ while the latter may oscillate (cf. BROWN [3], example) or have a limit which is a strict subset of $L(i, v^*)$ (cf. example 2 below)

EXAMPLE 2:

i	k	q_i^k	p_{i1}^k	p_{i2}^k	p_{i3}^k
1	1	0	1	0	0
2	1	0	0	1	0
3	1	0	1	0	0
3	-1	-1	0	.5	.5

Note that $g^* = (0, 0, 0)$, $d^* = 1$ and that both pure policies are maximal gain. Starting with $v(0) = (0, 2, 1)$ we find that $v(n) = (0, 2, .5^n)$ for all $n \geq 0$, hence $v^* = (0, 2, 0)$. For this special choice of $v(0)$, $K(3, n, 0) = \{2\}$ for all n which is a strict subset of $L(3, v^*) = \{1, 2\}$. Consequently, value iteration will fail to generate all functional optimal policies.

Finally, we have seen that regardless whether $v(n) - ng^*$ converges or not, none of the sequences of policies that may be generated by value-iteration, needs to converge, i.e. $\liminf_{n \rightarrow \infty} \bigcup_i K(i, n, 0)$ may be empty. Part (d) of theorem 6.1 shows that in case of foolproof policy convergence, i.e. existence of $\lim_{n \rightarrow \infty} \bigcup_i K(i, n, 0)$, only maximal gain policies ultimately occur. In addition, BATHER [1] part 1 example 1, which has every policy unchained and aperiodic, *falsified* the conjecture (cf. BROWN [3]) that the optimal policy sequence is asymptotically periodic, i.e. the existence of a J -tuple of policies (f^1, \dots, f^J) such that $v(n_0 + kJ + r) = q(f^r) + P(f^r)v(n_0 + kJ + r - 1)$ for all $k = 1, 2, \dots$ and some $n_0 \geq 1$.

NOTE 1. More specifically, the following argument was used in (VII-64) and (VII-75).

$$\max_{k=1, \dots, K} \max_{i, j \in A_k} \{f_i - f_j\} = \max_{1 \leq i, j \leq n} \{f_i - f_j\}$$

where f is an n -vector and the $\{A_k, k = 1, \dots, K\}$ constitute a partition of $\{1, \dots, n\}$.

The assertions (VII-64) and (VII-75) are repeatedly used in the remainder of the proof.

REFERENCES

- [1] BATHER, J., *Optimal decision procedures for finite Markov Chains, Part I, II*, Adv. Appl. Prob.5 (1973), 328-339, 521-540.
- [2] BELLMAN, R., *A Markovian Decision Process*, J. Math. Mech.6 (1957), 679-684.
- [3] BROWN, B., *On the iterative method of dynamic programming on a finite state space discrete time Markov Process*, Ann. Math. Statistics 36 (1965), 1279-1285.
- [4] HORDIJK, A., P.J. SCHWEITZER & H. TIJMS, *The asymptotic behaviour of the minimal total expected cost for the denumerable state Markov Decision Model*, J. Appl. Prob.12 (1975), 298-305.
- [5] HOWARD, R., *Dynamic Programming and Markov Processes*, John Wiley, New York (1960).
- [6] LANERY, E., *Etude asymptotique des systèmes Markoviens à commande*, R.I.R.O.3, (1967) 3-56.
- [7] LEMBERSKY, M.R., *On maximal rewards and ϵ -optimal policies in continuous time Markov decision chains*, Ann. Statist 2, (1974) 159-169.
- [8] MORTON, T., *On the asymptotic convergence rate of cost differences for Markovian Decision Processes*, O.R.19 (1971), 244-248.
- [9] ODONI, A., *On finding the maximal gain for Markov Decision Processes*, O.R.17 (1969), 857-860.
- [10] ROMANOVSKII, V., *Discrete Markov Chains*, Wolters-Noordhoff, Groningen (1970)
- [11] SCHWEITZER, P.J., *Perturbation Theory and Markovian Decision Processes*, Ph.D. dissertation, MIT (1965) (MITORC report H15).

- [12] SCHWEITZER, P.J., *A turnpike theorem for undiscounted Markovian Decision Processes*, presented at ORSA/TIMS, national meeting, May 1968.
- [13] SCHWEITZER, P.J., & A. FEDERGRUEN, *Functional Equations of Undiscounted Markov Renewal Programming*, Math. Center report BW 60/76.
- [14] SCHWEITZER, P.J., & A. FEDERGRUEN, *Geometric convergence of value-iteration in multichain Markov Decision Problems*, (to appear).
- [15] TIJMS, H., *On Dynamic Programming with arbitrary state space, compact action space and the average return criterion*, (1975), Math. Center report BW 55/75.
- [16] WHITE, D., *Dynamic Programming, Markov Chains, and the method of successive approximations*, J. of Math. Anal. and Appl. 6 (1963), 373-376.