

Petr Mandl; Ma.Rosario Romera Ayllón
On adaptive control of Markov processes

Kybernetika, Vol. 23 (1987), No. 2, 89--103

Persistent URL: <http://dml.cz/dmlcz/124869>

Terms of use:

© Institute of Information Theory and Automation AS CR, 1987

Institute of Mathematics of the Academy of Sciences of the Czech Republic provides access to digitized documents strictly for personal use. Each copy of any part of this document must contain these *Terms of use*.



This paper has been digitized, optimized for electronic delivery and stamped with digital signature within the project *DML-CZ: The Czech Digital Mathematics Library*
<http://project.dml.cz>

ON ADAPTIVE CONTROL OF MARKOV PROCESSES

PETR MANDL, M^a ROSARIO ROMERA AYLÓN

We consider a countable controlled Markov process. The control policy is the adaptive one. The differential generator is represented by the transition rate matrix. The case of infinite planning horizon and average reward per unit time as criterion is considered.

Sufficient conditions to ensure the existence of an optimal stationary strategy are given. The dependence of the asymptotic behaviour of the criterion function on the asymptotic behaviour of the control is also investigated. The paper is an extension to the continuous time case of the results in [3].

1. INTRODUCTION

We consider a system S with countable state space I . The process $X = \{X_t, t \geq 0\}$ describes the trajectory of the system. The control process is $Z = \{Z_t, t \geq 0\}$ ranging in a compact metric space \mathcal{Z} . The transition rate matrix Q is conservative, and the transition rates are continuous functions of z ,

$$Q = \|q(i, j; z)\|_{i, j \in I}, \quad z \in \mathcal{Z}, \quad q(i, z) = -q(i, i; z), \quad i \in I, \quad z \in \mathcal{Z}.$$

The criterion function is interpreted as the reward from S ; up to time T it is given by

$$R_T = \int_0^T r(X_t, Z_t) dt,$$

where $r(i, z)$ is a continuous function in z , $i \in I$, $z \in \mathcal{Z}$.

Admissible controls are defined following [4].

We consider the non-decreasing system of σ -algebras $\mathcal{F} = \{\mathcal{F}_t, t \geq 0\}$ generated by the process $X = \{X_t, t \geq 0\}$, which takes values in I , having piecewise constant right-continuous trajectories,

$$\mathcal{F}_t = \sigma\{X_s, s \leq t; \text{events of probability } 0\}.$$

The control process $Z = \{Z_t, t \geq 0\}$ with values in \mathcal{Z} is right-continuous and adapted to \mathcal{F} .

We say that X is the evolution of S under admissible control Z if

$${}^i M_t = \chi\{X_t = i\} - \int_0^t q(X_s, i; Z_s) ds, \quad t \geq 0, \quad i \in I,$$

are martingales with respect to \mathcal{F} .

In the general case the control strategies considered are the non-anticipative ones.

For stationary controls the control parameter at time t equals

$$Z_t = (\sigma)_{X_t}, \quad t \geq 0,$$

where $\sigma \in \mathcal{Z} \times \mathcal{Z} \times \mathcal{Z} \dots = \mathcal{Z}^\infty$.

The transition rate matrix under a stationary strategy has the expression

$$Q_\sigma = \|q(i, k; (\sigma)_i)\|_{i, k \in I} = \|q_\sigma(i, k)\|_{i, k \in I},$$

and the reward up to time T it is given by

$$R_T = \int_0^T r(X_t, (\sigma)_{X_t}) dt = \int_0^T (r_\sigma)_{X_t} dt.$$

To each stationary strategy σ we associate the mean reward defined as

$$\lim_{T \rightarrow \infty} T^{-1} R_T = \mu_\sigma,$$

provided that the limit exists and is constant with probability one.

We take optimality in the following sense; $\hat{\sigma}$ is an optimal stationary policy, if

$$\mu_{\hat{\sigma}} = \hat{\mu} = \sup_{\sigma \in \mathcal{Z}^\infty} \mu_\sigma.$$

Trajectory of S Corresponding to a Stationary Policy

Without loss of generality we select the state 1 to be the exceptional one. Following the idea expressed by Hordijk ([2]) in order to obtain sufficient conditions for the existence of an optimal stationary strategy, we introduce the \tilde{Q}_σ matrix

$$\tilde{q}_\sigma(i, j) = \begin{cases} q_\sigma(i, j), & i, j \in I, \quad i \neq 1, \\ q_\sigma(1, j) \delta_{1j}, & j \in I, \end{cases}$$

which is the restriction of Q_σ , obtained by replacing every element in the first column with exception of the diagonal one by 0. δ is the Kronecker delta function.

The matrix \tilde{Q}_σ then corresponds to the situation when the trajectory of S vanishes at the jump into state 1.

Let $\tilde{P}_\sigma(t)$, $t \geq 0$, be the minimal solution of the backward equation

$$\frac{d}{dt} \tilde{P}_\sigma(t) = \tilde{Q}_\sigma \tilde{P}_\sigma(t)$$

$$\tilde{P}_\sigma(0) = I \text{ (unit matrix).}$$

With a stationary policy σ we associate $X = \{X_t, t \geq 0\}$ describing the evolution

of S under σ as follows. From the initial state $X_0 = i$ the process X follows the transition law \tilde{P}_σ up to the first vanishing of the trajectory. The trajectory is right continuous. After the vanishing, the process starts afresh in state 1, and follows \tilde{P}_σ until it vanishes again and so on.

The jump-matrix corresponding to the Markov chain embedded in process X is

$$\tilde{P}_\sigma = \|\tilde{P}_\sigma(i, j)\|_{i, j \in I},$$

where

$$\tilde{p}_\sigma(i, j) = \begin{cases} 0, & i = j, \\ \frac{\tilde{q}_\sigma(i, j)}{\tilde{q}_\sigma(i)}, & i \neq j, \end{cases} = \frac{\tilde{q}_\sigma(i, j)}{\tilde{q}_\sigma(i)} (1 - \delta_{ij}), \quad i, j \in I.$$

The semi-Markovian viewpoint gives the following expressions

$$(1) \quad \int_0^\infty \tilde{P}_\sigma(s) e \, ds = \sum_{n=0}^\infty \tilde{P}_\sigma^n \frac{e}{q_\sigma},$$

$$\int_0^\infty \tilde{P}_\sigma(s) r_\sigma \, ds = \sum_{n=0}^\infty \tilde{P}_\sigma^n \frac{r_\sigma}{q_\sigma},$$

provided that the series converge absolutely. We set $e = (1, 1, \dots, 1)'$.

If v is a vector, then $(|v|)_i = |(v)_i|$, $(v^2)_i = (v)_i^2$. Analogously

$$\begin{pmatrix} v_1 \\ v_2 \end{pmatrix}_i = \begin{pmatrix} (v_1)_i \\ (v_2)_i \end{pmatrix} \quad \text{for vectors } v_1 \text{ and } v_2.$$

2. EXISTENCE OF AN OPTIMAL STATIONARY STRATEGY

We consider optimality in the maximum average reward per unit time sense; $\hat{\sigma} \in \mathcal{Z}^\infty$ is optimal if

$$\mu_{\hat{\sigma}} = \sup_{\sigma \in \mathcal{Z}^\infty} \mu_\sigma = \hat{\mu}.$$

In order to ensure the stability we introduce a Liapunov-type condition, which is similar to that introduced by Hordijk in [2].

Hypothesis H 1. There exists a vector $y_1 \geq 0$ such that

- (i) $|r_\sigma| + e + \tilde{Q}_\sigma y_1 \leq 0, \quad \sigma \in \mathcal{Z}^\infty,$
- (ii) $\lim_{N \rightarrow \infty} \tilde{P}_\sigma^N y_1 = 0, \quad \sigma \in \mathcal{Z}^\infty,$
- (iii) $\lim_{\sigma \rightarrow \sigma_0} \tilde{Q}_\sigma y_1 = \tilde{Q}_{\sigma_0} y_1, \quad \sigma, \sigma_0 \in \mathcal{Z}^\infty.$

The convergence considered in \mathcal{Z}^∞ is the component-wise one given by the metric in the compact space \mathcal{Z}^∞ .

Example. Conditions H 1 do not exclude the runs of the trajectory to infinity.

This is seen on the divergent pure birth process with rates

$$q(j, j+1) = j^2 = -q(j, j), \quad j = 1, 2, \dots$$

The j th inequality in (i) is

$$(2) \quad (|r_\sigma|)_j + 1 - j^2(y_1)_j + j^2(y_1)_{j+1} \leq 0.$$

Let

$$(|r_\sigma|)_j + 1 \leq \text{const.}, \quad j = 1, 2, \dots$$

Hence (2) holds for

$$(y_1)_j = 2 \text{ const.}/j, \quad j = 1, 2, \dots$$

Lemma 1. The Liapunov type condition (i) ensures the existence of μ_σ for all $\sigma \in \mathcal{E}^\infty$.

Proof. Having in consideration that one method to construct the \bar{P}_σ minimal backward Kolmogorov solution is the truncation procedure, we introduce the ${}^n\bar{P}_\sigma$ solution of the equation

$$\frac{d}{dt} {}^n\bar{P}_\sigma(t) = {}^n\bar{Q}_\sigma {}^n\bar{P}_\sigma(t), \quad {}^n\bar{P}_\sigma(0) = I,$$

where ${}^n\bar{Q}_\sigma$ is the $n \times n$ truncation of \bar{Q}_σ .

We also consider truncated vectors ${}^ne, {}^ny_1$. Then we have

$$\begin{aligned} {}^ne + {}^n\bar{Q}_\sigma {}^ny_1 &\leq 0, \\ {}^n\bar{P}_\sigma(s) {}^ne + {}^n\bar{P}_\sigma(s) {}^n\bar{Q}_\sigma {}^ny_1 &\leq 0, \\ {}^n\bar{P}_\sigma(s) {}^ne + \frac{d}{ds} {}^n\bar{P}_\sigma(s) {}^ny_1 &\leq 0, \end{aligned}$$

and hence

$$\int_0^t {}^n\bar{P}_\sigma(s) {}^ne \, ds + {}^n\bar{P}_\sigma(t) {}^ny_1 \leq {}^n\bar{P}_\sigma(0) {}^ny_1 = {}^ny_1.$$

Letting $n \rightarrow \infty$ and then $t \rightarrow \infty$, we obtain

$$(3) \quad \int_0^\infty \bar{P}_\sigma(s) e \, ds \leq y_1.$$

Similarly

$$(4) \quad \int_0^\infty \bar{P}_\sigma(s) |r_\sigma| \, ds \leq y_1.$$

From (3) it follows that under σ state 1 is positively recurrent. From (4) it follows that the expected reward earned until the jump into state 1 is finite. Hence, the strong law of large numbers yields

$$\lim_{T \rightarrow \infty} \frac{1}{T} R_T = \frac{\left(\int_0^\infty \bar{P}_\sigma(s) r_\sigma \, ds \right)_1}{\left(\int_0^\infty \bar{P}_\sigma(s) e \, ds \right)_1} = \mu_\sigma \quad \text{a.s.} \quad \square$$

If we make the additional assumption

$$q(i, 1; z) \leq \text{const.}, \quad z \in \mathcal{Z}, \quad i = 2, 3, \dots,$$

it follows with regard to (1) and to (4) that $\hat{\mu}$ is finite.

We also denote $\bar{q} = \sup_{\substack{i \in I \\ z \in \mathcal{Z}}} q(i, 1; z)$.

Theorem 1. Let H 1 hold. Then there exists an optimal stationary policy.

In virtue of (1) the proof can proceed as in [2].

Next define the potential

$$w = \sup_{\sigma_0 \sigma_1 \dots} \sum_{n=0}^{\infty} \bar{P}_{\sigma_0} \bar{P}_{\sigma_1} \dots \bar{P}_{\sigma_{n-1}} \left(\frac{r_{\sigma_n}}{q_{\sigma_n}} - \hat{\mu} \frac{e}{q_{\sigma_n}} \right),$$

which represents the optimal expected reward if $r(i, z) - \hat{\mu}$ is taken as reward function.

We have the optimality equations

$$(5) \quad w = \sup_{\sigma} \left\{ \frac{r_{\sigma}}{q_{\sigma}} - \frac{\hat{\mu}e}{q_{\sigma}} + \bar{P}_{\sigma} w \right\},$$

$$(6) \quad w = \frac{r_{\hat{\sigma}}}{q_{\hat{\sigma}}} - \frac{\hat{\mu}e}{q_{\hat{\sigma}}} + \bar{P}_{\hat{\sigma}} w,$$

which are similar to the discrete case optimality equations obtained in [3].

(5) and (6) are equivalent to

$$(7) \quad 0 = \sup_{\sigma} \{ r_{\sigma} - \hat{\mu}e + Q_{\sigma} w \},$$

$$(8) \quad 0 = r_{\hat{\sigma}} - \hat{\mu}e + Q_{\hat{\sigma}} w.$$

We can conclude that $(\hat{\mu}, \hat{\sigma})$ is the solution of the optimality equations. We also have

$$0 = r_{\bar{\sigma}} - \mu_{\bar{\sigma}}e + Q_{\bar{\sigma}} w_{\bar{\sigma}} \quad \text{for } \bar{\sigma} \in \mathcal{Z}^{\infty}.$$

3. ASYMPTOTIC BEHAVIOUR OF THE CRITERION FUNCTIONAL

In this section we consider the initial state X_0 fixed.

Lemma 2. Let $0 \leq c(i, z)$, $i \in I$, $z \in \mathcal{Z}$, be continuous function of z . Let

$$(9) \quad c_{\sigma} + \bar{Q}_{\sigma} y \leq 0, \quad \sigma \in \mathcal{Z}^{\infty},$$

with $y \geq 0$.

Then, under admissible control Z for $T \geq 0$

$$(10) \quad \mathbb{E} \int_0^T c(X_t, Z_t) dt + \mathbb{E}(y)_{X_T} \leq (y)_{X_0} + (y)_1 \bar{q} T.$$

Proof. Let $n \geq X_0$. Denote $\tau_n = \inf \{t, X_t > n\}$. Then

$$\{ {}^i M_{t \wedge \tau_n}, t \geq 0 \}, \quad i \in I, \quad \text{where } t \wedge \tau_n = \min(t, \tau_n),$$

are martingales, and

$$\sum_{i=1}^m (y)_i {}^i M_{t \wedge \tau_n}, \quad t \geq 0,$$

are also martingales for each m .

Let $T \geq 0$. For $m > X_0$ we have

$$\begin{aligned} (y)_{X_0} &= \mathbb{E} \sum_{i=1}^m (y)_i {}^i M_{T \wedge \tau_n} = \\ &= \mathbb{E} \chi \{ X_{T \wedge \tau_n} \leq m \} (y)_{X_{T \wedge \tau_n}} - \mathbb{E} \int_0^{T \wedge \tau_n} \sum_{i=1}^m q(X_t, i; Z_t) (y)_i dt. \end{aligned}$$

Letting $m \rightarrow \infty$ we obtain

$$(11) \quad (y)_{X_0} = \mathbb{E}(y)_{X_{T \wedge \tau_n}} - \mathbb{E} \int_0^{T \wedge \tau_n} \sum_i q(X_t, i; Z_t) (y)_i dt.$$

Note that the last expectation is finite, since from (9) follows

$$\left| \sum_i q(X_t, i; Z_t) (y)_i \right| \leq q(X_t, Z_t) (y)_{X_t} + |q(X_t, 1; Z_t)| (y)_1,$$

and in the integrand is $X_t \in \{1, \dots, n\}$ for $t \in [0, T \wedge \tau_n]$. From (11) and (9)

$$\begin{aligned} (y)_{X_0} - \mathbb{E} \int_0^{T \wedge \tau_n} c(X_t, Z_t) dt &= \mathbb{E}(y)_{X_{T \wedge \tau_n}} - \mathbb{E} \int_0^{T \wedge \tau_n} \chi \{ X_t \neq 1 \} q(X_t, 1; Z_t) (y)_1 dt - \\ &- \mathbb{E} \int_0^{T \wedge \tau_n} (c(X_t, Z_t) + \sum_i \hat{q}(X_t, i; Z_t) (y)_i) dt \geq \mathbb{E} \chi \{ \tau_n > T \} (y)_{X_T} - \bar{q}(y)_1 \mathbb{E}(T \wedge \tau_n). \end{aligned}$$

Letting $n \rightarrow \infty$ (10) is obtained. \square

Law of Large Numbers

In the following we fix $\bar{\sigma} \in \mathcal{L}^\infty$. To $\bar{\sigma}$ we can associate the scalar-vector couple $(\mu_{\bar{\sigma}}, w_{\bar{\sigma}})$ satisfying

$$0 = r_{\bar{\sigma}} - \mu_{\bar{\sigma}} e + Q_{\bar{\sigma}} w_{\bar{\sigma}}.$$

We shall write briefly μ, w .

We define

$$\varphi(i, z) = r(i, z) - \mu + \sum_j q(i, j; z) (w)_j, \quad i \in I, \quad z \in \mathcal{L},$$

which is a continuous function in z . $\varphi(X_t, Z_t)$ is a measure of difference between the actual parameter value Z_t and the value $(\bar{\sigma})_{X_t}$ corresponding to the stationary strategy $\bar{\sigma}$.

We introduce the process

$$(12) \quad M_T = R_T - T\mu + (w)_{X_T} - (w)_{X_0} - \int_0^T \varphi(X_t, Z_t) dt, \quad T \geq 0,$$

and the functions for $m = 2, 3, \dots$

$$r_m(i, z) = \sum_{j \neq i} q(i, j; z) ((w)_j - (w)_i)^m,$$

$$h_m(i, z) = \sum_{j \neq i} q(i, j; z) ((y_1)_j^m + (y_1)_i^m), \quad i \in I, \quad z \in \mathcal{Z}.$$

In order to prove that the Law of Large Numbers holds for R_T , we need to make a stronger hypothesis.

Hypothesis H 2. There exists a $y_2 \geq y_1^2$ such that

$$h_{2\sigma} + \bar{Q}_\sigma y_2 \leq 0, \quad \sigma \in \mathcal{Z}^\infty.$$

Since $|w| \leq \text{const. } y_1$ we have,

$$|r_m(i, z)| \leq \text{const. } h_m(i, z).$$

Lemma 3. Let H 1 and H 2 hold. Then for any admissible Z ,

- (i) $M = \{M_T, T \geq 0\}$ is a martingale with respect to \mathcal{F} .
- (ii) $\mathbb{E}(M_T - M_S)^2 = \mathbb{E} \int_S^T r_2(X_t, Z_t) dt$, $0 \leq S \leq T$.
- (iii) $\mathbb{E}M_T^2 + \mathbb{E}(y_2)_{X_T} \leq (y_2)_{X_0} + (y_2)_1 \bar{q}T$, $T \geq 0$.

Proof. Detailed proof is given in [5], Lemma III.2.1. From Lemma 2 one obtains

$$\mathbb{E} \int_0^{T \wedge \tau_n} r_2(X_t, Z_t) dt + \mathbb{E}(y_2)_{X_{T \wedge \tau_n}} \leq (y_2)_{X_0} + (y_2)_1 \bar{q}T, \quad T \geq 0,$$

where $\tau_n = \inf \{t, X_t > n\}$.

First it is shown that $\{M_{T \wedge \tau_n}, T \geq 0\}$, $n = 1, 2, \dots$, are martingales and that

$$\mathbb{E}(M_{T \wedge \tau_n})^2 = \mathbb{E} \int_0^{T \wedge \tau_n} r_2 dt.$$

Letting $n \rightarrow \infty$ (i), (ii), (iii) follow.

From this lemma we can get the following results.

Theorem 1. Let H 1, H 2 hold. Then

$$\lim_{T \rightarrow \infty} T^{-1} R_T = \mu$$

in probability, in first order mean, in second order mean if and only if

$$\lim_{T \rightarrow \infty} T^{-1} \int_0^T \varphi(X_t, Z_t) dt = 0,$$

respectively in probability, in first order mean and in second order mean.

Theorem 2. Let H 1, H 2 hold. Then

$$(13) \quad \lim_{T \rightarrow \infty} T^{-1} M_T = 0 \quad \text{a.s.}$$

Theorem 3. Let H 1, H 2 hold. Then under any admissible control

$$\lim_{T \rightarrow \infty} P(T^{-1} R_T \geq \hat{\mu} + \varepsilon) = 0 \quad \text{for } \varepsilon > 0,$$

and

$$\overline{\lim}_{T \rightarrow \infty} T^{-1} E R_T \leq \hat{\mu}.$$

The proof of these theorems can proceed as in [5].

In order to obtain the main result of this subsection, we introduce the vectors e_k

$$e_k: (e_k)_i = 0, \quad i < k, \quad (e_k)_i = 1, \quad i \geq k,$$

and set

$$\varepsilon_\sigma(k) = \left(\int_0^\infty \bar{P}_\sigma f(s) e_k ds \right) / \left(\int_0^\infty \bar{P}_\sigma f(s) e ds \right)_1, \quad \sigma \in \mathcal{Z}^\infty,$$

$$\hat{\delta}(k) = \sup_{\sigma \in \mathcal{Z}^\infty} \varepsilon_\sigma(k).$$

Following [3], Lemma 10, we can derive

$$\lim_{k \rightarrow \infty} \hat{\delta}(k) = 0.$$

Since H 1 implies

$$e_k + e + 2\bar{Q}_\sigma y_1 \leq 0, \quad \sigma \in \mathcal{Z}^\infty,$$

we can use Theorem 3 to conclude that under H 1, H 2

$$\overline{\lim}_{T \rightarrow \infty} T^{-1} E \int_0^T (e_k)_{X_t} dt = \overline{\lim}_{T \rightarrow \infty} T^{-1} \int_0^T P(X_t \geq k) dt \leq \hat{\delta}(k),$$

for Z admissible.

The analogy between $\varepsilon_\sigma(k)$, $\hat{\delta}(k)$ and the quantities μ_σ , $\hat{\mu}$ is obvious.

Hypothesis H 3. There exists a $y_3 \geq 0$ such that

$$(q_\sigma y_1)^2 + \bar{Q}_\sigma y_3 \leq 0, \quad \sigma \in \mathcal{Z}^\infty.$$

Theorem 4. Let H 1, H 2, H 3 hold. If the admissible control Z is such that

$$(14) \quad \lim_{t \rightarrow \infty} \varrho(Z_t, (\bar{\sigma})_{X_t}) = 0 \quad \text{in probability,}$$

(ϱ is the distance), then

$$\lim_{T \rightarrow \infty} E |T^{-1} R_T - \mu| = 0.$$

Proof. Let (14) hold. According to Theorem 1 we have to establish

$$\lim_{T \rightarrow \infty} T^{-1} E \left| \int_0^T \varphi(X_t, Z_t) dt \right| = 0.$$

H 1 implies

$$\begin{aligned} |r(i, z)| &\leq q(i, z)(y_1)_i, \\ \left| \sum_j q(i, j, z)(w)_j \right| &\leq \text{const.} \left| \sum_{j \neq i} q(i, j, z)(y_1)_j + q(i, z)(y_1)_i \right| \leq \\ &\leq \text{const.} (q(i, z)(y_1)_i + \bar{q}(y_1)_i). \end{aligned}$$

Hence we conclude that

$$(15) \quad |\varphi(i, z)|^2 \leq C[1 + (q(i, z)(y_1)_i)^2],$$

where C is a constant. φ is a continuous function in z , in virtue of H 1, (iii).

Take $L > 0$, $k > 0$, k integer. It holds

$$(16) \quad \begin{aligned} T^{-1} \mathbb{E} \left| \int_0^T \varphi(X_t, Z_t) dt \right| &\leq T^{-1} \int_0^T \mathbb{E} |\chi\{X_t < k\} \varphi(X_t, Z_t)| dt + \\ &+ T^{-1} L \int_0^T \mathbb{P}(X_t \geq k) dt + T^{-1} \int_0^T \mathbb{E} |\chi\{\varphi(X_t, Z_t) > L\} \varphi(X_t, Z_t)| dt. \end{aligned}$$

We have

$$\lim_{t \rightarrow \infty} \mathbb{E} |\chi\{X_t < k\} \varphi(X_t, Z_t)| = 0,$$

since $\chi\{X_t < k\} \varphi(X_t, Z_t)$ is bounded and tends to 0 in probability as $t \rightarrow \infty$, by (14) and by $\varphi(X_t, (\bar{\sigma})_{X_t}) = 0$.

The second term on the right in (15) is estimated from the expression obtained before,

$$\overline{\lim}_{T \rightarrow \infty} T^{-1} \int_0^T \mathbb{P}(X_t \geq k) dt \leq \delta(k).$$

As far as the third term, we can use H 3, (15) and Lemma 2 to derive

$$\int_0^T \mathbb{E} \varphi^2(X_t, Z_t) dt \leq C \left(T + \int_0^T (q(X_t, Z_t)(y_1)_{X_t})^2 dt \right) \leq C(T + (y_3)_{X_0} + (y_3)_1 \bar{q}T).$$

Since $\mathbb{E} |\chi\{\varphi(X_t, Z_t) > L\} \varphi(X_t, Z_t)| \leq L^{-1} \mathbb{E} \varphi^2(X_t, Z_t)$, we have that

$$\overline{\lim}_{T \rightarrow \infty} T^{-1} \int_0^T \mathbb{E} |\chi\{\varphi(X_t, Z_t) > L\} \varphi(X_t, Z_t)| dt \leq L^{-1} C(1 + (y_3)_1 \bar{q}).$$

Consequently,

$$\overline{\lim}_{T \rightarrow \infty} T^{-1} \left| \int_0^T \varphi(X_t, Z_t) dt \right| \leq L\delta(k) + L^{-1} C(1 + (y_3)_1 \bar{q}).$$

The right-hand side can be made arbitrarily small by taking L and then k sufficiently large. \square

The results aforementioned allow to say that the validity of the Law of Large Numbers for R_T holds if the difference (measured in distance) between the actual parameter control value Z_t and that corresponding to the stationary strategy $(\bar{\sigma})_{X_t}$ tends to zero if t goes to infinity.

Central Limit Theorem

Our next purpose is to prove the Central Limit Theorem for R_T .

We introduce hypotheses H* 1, H* 2 and H* 3. They are analogous to H 1, H 2 and H 3 with r_σ replaced by $h_{2\sigma}$ (or $r(i, z)$ replaced by $h_2(i, z)$).

Set $Y_n = M_{n+1} - M_n$, $n = 0, 1, \dots$, defined from the martingale (12).

We consider a fixed $\bar{\sigma} \in \mathcal{X}^{\otimes 2}$ and denote $\mu_2 = \lim_{N \rightarrow \infty} N^{-1} \int_0^N r_2(X_t, Z_t) dt$.

From martingale theory ([1]) the following result is known.

Proposition 1. Let $\{Y_n, n = 0, 1, \dots\}$ be the martingale differences. Further let

$$\lim_{N \rightarrow \infty} N^{-1} \sum_{n=0}^{N-1} E\{Y_n^2 | \mathcal{F}_n\} = v \quad \text{in probability}$$

where v is a constant, and for each $\varepsilon > 0$ let the Lindeberg condition hold

$$\lim_{N \rightarrow \infty} N^{-1} \sum_{n=0}^{N-1} E Y_n^2 \chi\{|Y_n| \geq \varepsilon \sqrt{N}\} = 0.$$

Then

M_N/\sqrt{N} as $N \rightarrow \infty$ has asymptotically normal distribution $N(0, v)$.

Extending the argument of Lemma 3 we can prove

$$E\{Y_n^2 | \mathcal{F}_n\} = E\left\{\int_n^{n+1} r_2^2(X_t, Z_t) dt | \mathcal{F}_n\right\}.$$

The aim is to use the result from Proposition 1 so we need to verify its two conditions.

Concerning the first one we have the following. If H 1 and H 2 hold then

$$\lim_{N \rightarrow \infty} \left(N^{-1} \int_0^N r_2(X_t, Z_t) dt - N^{-1} \sum_{n=0}^{N-1} E\{Y_n^2 | \mathcal{F}_n\} \right) = 0 \quad \text{a.s.}$$

To prove this, set

$$L_N = \int_0^N r_2(X_t, Z_t) dt - \sum_{n=0}^{N-1} E\left\{\int_n^{n+1} r_2(X_t, Z_t) dt | \mathcal{F}_n\right\}, \quad N = 1, 2, \dots$$

Then $L = \{L_N, N = 1, 2, \dots\}$ is a martingale, and we have that

$$\begin{aligned} E(L_{N+1} - L_N)^2 &= E\left(\int_N^{N+1} r_2(X_t, Z_t) dt\right)^2 - E\left(E\left\{\int_N^{N+1} r_2(X_t, Z_t) dt | \mathcal{F}_N\right\}\right)^2 \leq \\ &\leq E\left(\int_N^{N+1} r_2(X_t, Z_t) dt\right)^2 \leq E \int_N^{N+1} r_2^2(X_t, Z_t) dt, \quad N = 1, 2, \dots \end{aligned}$$

This, as in the proof of Theorem 2, implies

$$\sum_{n=1}^{\infty} \frac{E(L_{n+1} - L_n)^2}{n^2} < \infty.$$

Consequently, L fulfills the Law of Large Numbers. Hence,

$$(17) \quad \lim_{N \rightarrow \infty} \left(N^{-1} \int_0^N r_2(X_t, Z_t) dt - N^{-1} \sum_{n=0}^{N-1} \mathbb{E}\{Y_n^2 | \mathcal{F}_n\} \right) = 0 \quad \text{a.s.}$$

If $H^* 1$, $H^* 2$, $H^* 3$ hold, we conclude from Theorem 4 that

$$\lim_{N \rightarrow \infty} N^{-1} \int_0^N r_2(X_t, Z_t) dt = \mu_2 \quad \text{in probability.}$$

Hence, (17) implies

$$\lim_{N \rightarrow \infty} N^{-1} \sum_{n=0}^{N-1} \mathbb{E}\{Y_n^2 | \mathcal{F}_n\} = \mu_2 \quad \text{in probability.}$$

Concerning the second hypothesis of Proposition 1, namely the Lindeberg condition, we have to state an auxiliary result.

Lemma 4. Let $H 1$, $H 2$, $H^* 1$, $H^* 2$ and $H^* 3$ hold. Then

$$(18) \quad \overline{\lim}_{N \rightarrow \infty} \frac{1}{N} \sum_{n=0}^{N-1} \mathbb{E}Y_n^4 \leq \text{const.}$$

under arbitrary admissible control Z .

The laborious proof can proceed as in [4] and is based on the relation

$$(19) \quad \mathbb{E}(M_T - M_S)^4 = \mathbb{E} \left(\int_S^T r_4(X_t, Z_t) dt + 4 \int_S^T (M_t - M_S) r_3(X_t, Z_t) dt + \right. \\ \left. + 6 \int_S^T (M_t - M_S)^2 r_2(X_t, Z_t) dt \right).$$

Next step deals with the expression

$$\sum_{n=0}^{N-1} \mathbb{E}(M_{n+1} - M_n)^4 = \sum_{n=0}^{N-1} \mathbb{E} \int_0^1 r_4(X_{t+n}, Z_{t+n}) dt + \\ + 4 \sum_{n=0}^{N-1} \mathbb{E} \int_0^1 (M_{t+n} - M_n) r_3(X_{t+n}, Z_{t+n}) dt + \\ + 6 \sum_{n=0}^{N-1} \mathbb{E} \int_0^1 (M_{t+n} - M_n)^2 r_2(X_{t+n}, Z_{t+n}) dt \leq \mathbb{E} \int_0^N r_4(X_t, Z_t) dt + \\ + 4 \left(\sum_{n=0}^{N-1} \mathbb{E}(M_{n+1} - M_n)^4 \right)^{1/4} \left(\mathbb{E} \int_0^N |r_3(X_t, Z_t)|^{4/3} dt \right)^{3/4} + \\ + 6 \left(\sum_{n=0}^{N-1} \mathbb{E}(M_{n+1} - M_n)^4 \right)^{1/2} \left(\mathbb{E} \int_0^N r_2^2(X_t, Z_t) dt \right)^{1/2}.$$

This expression is obtained from (19), using the Hölder inequality and the fact that

$$\mathbb{E}(M_{t+n} - M_n)^4 \quad \text{is nondecreasing in } t.$$

It is not difficult to see that from here (18) follows, provided that

$$\mathbb{E} \int_0^N r_4(X_t, Z_t) dt + \mathbb{E} \int_0^N |r_3(X_t, Z_t)|^{4/3} dt + \mathbb{E} \int_0^N r_2^2(X_t, Z_t) dt \leq \text{const. } N.$$

From Lemma 4, the Lindeberg condition follows; for $\varepsilon > 0$

$$\overline{\lim}_{N \rightarrow \infty} N^{-1} \sum_{n=0}^{N-1} \mathbb{E} Y_n^2 \chi\{|Y_n| \geq \varepsilon \sqrt{N}\} \leq \overline{\lim}_{N \rightarrow \infty} \varepsilon^{-2} N^{-2} \sum_{n=0}^{N-1} \mathbb{E} Y_n^4 = 0.$$

We can state now the main result in this subsection.

Theorem 5. Let H 1, H 2, H 3, H* 1, H* 2, H* 3 hold, and let Z be an admissible control such that

$$(20) \quad \lim_{t \rightarrow \infty} \varrho(Z_t, (\bar{\sigma})_{x_t}) = 0 \quad \text{in probability,}$$

$$\lim_{T \rightarrow \infty} \frac{1}{\sqrt{T}} \int_0^T \varphi(X_t, Z_t) dt = 0 \quad \text{in probability.}$$

Then $(R_T - T\mu)/\sqrt{T}$ has asymptotically normal distribution $N(0, \mu_2)$ as $T \rightarrow \infty$.

Proof. Set

$$(21) \quad \frac{M_T}{\sqrt{T}} = \frac{R_T - T\mu}{\sqrt{T}} - \frac{(w)_{x_T} - (w)_{x_0}}{\sqrt{T}} + \frac{1}{\sqrt{T}} \int_0^T \varphi(X_t, Z_t) dt.$$

Theorem 4 gives the following relation

$$\lim_{N \rightarrow \infty} N^{-1} \int_0^N r_2(X_t, Z_t) dt = \mu_2 \quad \text{in probability.}$$

Hence from (17)

$$\lim_{N \rightarrow \infty} N^{-1} \sum_{n=0}^{N-1} \mathbb{E}\{Y_n^2 | \mathcal{F}_n\} = \mu_2 \quad \text{in probability.}$$

We see that M_N/\sqrt{N} is asymptotically $N(0, \mu_2)$ as $N \rightarrow \infty$. From Lemma 4 follows

$$(22) \quad \sum_{N=1}^{\infty} \frac{\mathbb{E} Y_N^4}{N^2} < \infty.$$

Further

$$\mathbb{E} \left(\frac{M_T}{\sqrt{[T]}} - \frac{M_{[T]}}{\sqrt{[T]}} \right)^4 \leq \frac{\mathbb{E} Y_{[T]}^4}{[T]^2}.$$

In virtue of (22), the right term in this inequality converges to zero as T tends to ∞ . We conclude that $M_T/\sqrt{[T]}$ and hence M_T/\sqrt{T} is asymptotically $N(0, \mu_2)$ as T tends to infinity. The proof is accomplished by demonstrating the negligibility of

$$[(w)_{x_T} - (w)_{x_0}]/\sqrt{T} \quad \text{as } T \text{ tends to } \infty. \quad \square$$

4. EXAMPLE

As an illustrative example we propose the consideration of a system described as a $M/M/1$ queue, in order to apply the results obtained before.

Let us suppose a known arrival rate given by q , and a known service rate given by z , which is going to be control parameter ranging in the interval $[z', z'']$, where $z' > q$.

The process X_t denotes the number of customers in the system at time $t > 0$. The cost functional is defined as follows

$$C_T = \int_0^T (X_t + KZ_t) dt,$$

where K is a positive constant.

The Bellman equation from Section 2 for the minimum mean cost takes the expression

$$0 = \min_{z \in [z', z'']} [Kz - qw_0 + qw_1 - \hat{\mu}],$$

$$0 = \min_{z \in [z', z'']} [i + Kz + zw_{i-1} - (z + q)w_i + qw_{i+1} - \hat{\mu}], \quad i = 1, 2, \dots$$

Then, minimizing the expression in brackets we obtain for z

$$z = z'' \quad \text{if} \quad K + w_{i-1} - w_i < 0,$$

$$z = z' \quad \text{if} \quad K + w_{i-1} - w_i > 0.$$

The optimal strategy of control takes the form

$$\sigma(i) = z', \quad i = 1, 2, \dots, n,$$

$$\sigma(i) = z'', \quad i = n + 1, n + 2, \dots,$$

which is a bang-bang control.

The average cost $\mu(n)$ corresponding to this strategy can be obtained in a direct way. The stationary distribution for the σ strategy is

$$p_i = g \left(\frac{q}{z'} \right)^i, \quad i = 0, 1, \dots, n - 1,$$

$$p_i = g \left(\frac{q}{z'} \right)^n \left(\frac{q}{z''} \right)^{i-n}, \quad i = n, n + 1, \dots,$$

where

$$g = \left(\frac{1 - u^m}{1 - u} + u^n \frac{1}{1 - v} \right)^{-1}, \quad u = \frac{q}{z'}, \quad v = \frac{q}{z''}.$$

According to the definition of the mean cost

$$\mu'(n) = g \left\{ u \frac{1-u^n}{(1-u)^2} - nu^n \left(\frac{1}{1-u} - \frac{1}{1-v} \right) + \right. \\ \left. + u^n \frac{v}{(1-v)^2} + K(z'' - z') u^n \frac{v}{1-v} \right\} + Kz'.$$

The optimal value n verifies

$$\mu(\hat{n}) = \min_n \mu(n) = \hat{\mu}.$$

Then the optimal strategy of control $\hat{\sigma}$ fulfils

$$(\hat{\sigma})_i = z', \quad i = 0, 1, \dots, \hat{n},$$

$$(\hat{\sigma})_i = z'', \quad i = \hat{n} + 1, \dots$$

The optimality equations obtained before can be expressed in the following way

$$0 = Kz' + q(w_1 - w_0) - \hat{\mu},$$

$$0 = i + Kz' - z'(w_i - w_{i-1}) + q(w_{i+1} - w_i) - \hat{\mu}, \quad i = 1, \dots, \hat{n},$$

$$0 = i + Kz'' - z''(w_i - w_{i-1}) + q(w_{i+1} - w_i) - \hat{\mu}, \quad i = \hat{n} + 1, \hat{n} + 2 \dots$$

w_i is a quadratic function for $i = \hat{n} + 1, \hat{n} + 2, \dots$. The fact that $\hat{\sigma}$ is an optimal stationary strategy is expressed by the inequalities

$$w_{\hat{n}} - w_{\hat{n}-1} \leq K \leq w_{\hat{n}+1} - w_{\hat{n}}.$$

Assume now that the rate q is unknown to the controller, and write $\hat{\sigma} = \hat{\sigma}(q)$. Let N_T denote the number of arrivals into the system up to time T . Then the ratio

$$q_T^* = N_T/T$$

is the maximum likelihood estimate of q . A natural way of defining an adaptive control is to set

$$Z_t = (\hat{\sigma}(q_t^*))_{X_t}, \quad t \geq 0.$$

Since

$$\lim_{T \rightarrow \infty} q_T^* = q \quad \text{a.s.},$$

(20) holds with $\hat{\sigma} = \hat{\sigma}(q)$. It can be shown that the other hypotheses of Theorem 5 are fulfilled as well.

(Received March 25, 1986.)

REFERENCES

- [1] B. M. Brown: Martingale Central Limit Theorems. *Ann. Math. Statist.* 42 (1971), 59-66.
 [2] A. Hordijk: Dynamic Programming and Markov Potential Theory. *Math. Centrum, Amsterdam* 1974.

- [3] P. Mandl: On the adaptive control of countable Markov chains. Prob. Theory, Banach Center Publications, Vol. 5, 159—173, Warsaw 1979.
- [4] P. Mandl: Martingale Methods in Discrete State Random Processes. Supplement to the Journal Kybernetika 18 (1982).
- [5] M^a R. Romera: Adaptive Control of Markov Processes with Countable State Space. Doctoral Thesis (in Spanish). Universidad Complutense, Madrid 1985.

RNDr. Petr Mandl, Dr.Sc., matematicko-fyzikální fakulta Univerzity Karlovy (Faculty of Mathematics and Physics, Charles University), Sokolovská 83, 186 00 Praha 8, Czechoslovakia.
Prof. Dr. M^a Rosario Romera-Ayllón, Faculty of Informatics, Polytechnical University of Madrid, Km. 7 Carretera de Valencia, 28031 Madrid, Spain.