

# On Mean Reward Variance in Semi-Markov Processes

Karel Sladký

Institute of Information Theory and Automation  
Academy of Sciences of the Czech Republic  
Pod Vodárenskou věží 4, 18208 Praha 8, Czech Republic

Received: date / Revised version: date

**Abstract** As an extension of the discrete-time case, this note investigates the variance of the total cumulative reward for the embedded Markov chain of semi-Markov processes. Under the assumption that the chain is aperiodic and contains a single class of recurrent states recursive formulae for the variance are obtained which show that the variance growth rate is asymptotically linear in time. Expressions are provided to compute this growth rate.

**Keywords.** Markov and semi-Markov processes with rewards, variance of cumulative reward, asymptotic behaviour

**AMS Subject Classification.** Primary 90C47  
Secondary 60J27

## 1 Introduction

The usual average criteria examined in the literature on Markov reward processes can be insufficient to fully capture the various aspects for a decision maker. It may be preferable to also include more sophisticated criteria to reflect variability-risk features (for details see [17]). Most notably, the variance of cumulative rewards can be indicative and seems of interest.

The variance of cumulative reward structures has been studied extensively for discrete-time Markov reward chains. The majority of these results has been involved with the development of further optimization criterion (see [1], [2], [3], [4], [5], [7], [8], [10], [13], [15], [16]).

For the continuous-time case only structural results were reported in [5] while just recently explicit expressions were developed in [14]. (The present paper is related to the latter reference but focuses on the embedded Markov chain to also include non-exponential densities).

No results, however, seem to be reported for semi-Markov processes. Such an extension is not direct as higher moments and stochastic durations are to be taken into account.

In this paper therefore we aim to establish a first step in this direction. This step can be seen as essential for further optimization results in line with the discrete-time case.

## 2 Formulation and Notation

Consider a semi-Markov reward process  $Y = \{Y(t), t \geq 0\}$  with finite state space  $\mathcal{I} = \{1, 2, \dots, N\}$  along with the embedded Markov chain  $X = \{X_n, n = 0, 1, \dots\}$ . The transition and reward structure is characterized by

- $p_{ij}$  : transition probability from  $i \rightarrow j$  ( $i, j \in \mathcal{I}, j \neq i$ ) of the embedded Markov chain  $X$  with generic stochastic matrix  $\mathbf{P} = [p_{ij}]_{i,j=1}^N$ ,
- $\eta_{ij}$  : random time of the transition from  $i \rightarrow j$ , hence  $\eta_i = \sum_{j \in \mathcal{I}} p_{ij} \eta_{ij}$  is the random time spent by the semi-Markov process  $Y$  in state  $i$ ,
- $F_{ij}(\tau)$  : distribution function representing the conditional probability  $\mathbb{P}(\eta_{ij} \leq \tau)$ ,
- $r_{ij}$  : instantaneous transition reward for a transition from  $i \rightarrow j$ ,
- $r_i$  : reward rate per unit of time incurred in state  $i$ .

We will be interested in the variance of the total reward per unit of time and its behavior for the infinite time horizon.

To this end, let the vector  $\mathcal{R}(t)$  denote the expected total reward of the semi-Markov process  $Y(t)$  up to time  $t$  given its initial state at time  $t = 0$ . More precisely, we are interested in

$$g = \lim_{t \rightarrow \infty} \frac{1}{t} \cdot \mathcal{R}_i(t) \quad \text{with} \quad \mathcal{R}_i(t) = \mathbb{E}[\xi(t) | Y(0) = i] \quad (1)$$

where

$$\xi(t) = \left[ \int_0^t r_{Y(s)} ds + \sum_{k=0}^{N(t)} r_{Y(\tau_k^-), Y(\tau_k^+)} \right]$$

with  $Y(s)$ , denoting the state of the system at time  $s$ ,  $Y(\tau_k^-)$  and  $Y(\tau_k^+)$  the state just prior and after the  $k$ -th jump,  $N(t)$  the number of jumps up to time  $t$  and  $\mathbf{E}$  and  $\sigma^2$  the standard symbols for expectation and variance.

Now consider the embedded Markov chain  $X = \{X_n, n = 0, 1, 2, \dots\}$  and let the vectors  $\mathbf{R}(n)$ ,  $\mathbf{S}(n)$ , and  $\mathbf{V}(n)$  denote the first moment, the second moment and the variance of the (random) total reward  $\xi_n$  received in the first  $n$  transitions of the embedded Markov chain  $X$  given its initial state at time  $t = 0$  respectively.

That is,

$$\begin{aligned} R_i(n) &= \mathbb{E}[\xi_n | X(0) = i], \\ S_i(n) &= \mathbb{E}[\xi_n^2 | X(0) = i], \\ V_i(n) &= \sigma^2[\xi_n | X(0) = i] \end{aligned}$$

where

$$\xi_n = \sum_{k=0}^{n-1} [r_{X_k} \cdot \eta_{X_k, X_{k+1}} + r_{X_k, X_{k+1}}].$$

Similarly,  $\mathbf{D}(n)$  is the vector of expected times  $\zeta_n$  spent in the first  $n$  transitions with elements

$$D_i(n) = \mathbb{E}[\zeta_n | X(0) = i] \quad \text{where} \quad \zeta_n = \sum_{k=0}^{n-1} [\eta_{X_k, X_{k+1}}].$$

From the theory of dynamic programming it is well known that under assumption AS1 below

$$g^r := \lim_{n \rightarrow \infty} \frac{1}{n} \cdot R_i(n), \quad g^t := \lim_{n \rightarrow \infty} \frac{1}{n} \cdot D_i(n) \quad \text{exist for all } i \in \mathcal{I}.$$

It can then also be shown (see [12], Chapt. 7.3) that under assumptions AS1–AS2 below then also:

$$g = \lim_{t \rightarrow \infty} \frac{1}{t} \cdot \mathcal{R}_i(t) =: \frac{g^r}{g^t} = \frac{\lim_{n \rightarrow \infty} R_i(n)}{\lim_{n \rightarrow \infty} D_i(n)}. \quad (2)$$

As mentioned in the introduction, such as for second order optimization purposes, and in analogy with (2) as our main interest we aim to study the “mean” variance defined by

$$G := \frac{\lim_{n \rightarrow \infty} V_i(n)}{\lim_{n \rightarrow \infty} D_i(n)}. \quad (3)$$

*Notation.* In what follows  $\mathbf{I}$  denotes an identity matrix, and  $\mathbf{e}$  is reserved for a unit column vector. By  $\varepsilon(t)$  we denote a function in  $t$  such that  $\varepsilon(t) \rightarrow 0$  exponentially fast as  $t \rightarrow \infty$ , i.e. for some  $\alpha$  and  $\beta$ :  $|\varepsilon(t)| \leq \alpha \cdot e^{-\beta t}$ . By  $\boldsymbol{\varepsilon}(t)$  we denote a vector function such that each component  $\varepsilon_i(t) \rightarrow 0$  as  $t \rightarrow \infty$  exponentially fast.

*Assumptions.* We make the following assumptions:

**AS 1.** The transition probability matrix  $\mathbf{P}$  has a single class of recurrent states and is aperiodic.

**AS 2.**  $F_{ij}(\tau)$  is a non-lattice distribution. The first two moments  $d_i^{(1)}, d_i^{(2)}$  of the of (random) time  $\eta_i$  ( $i \in \mathcal{I}$ ) spent by the process  $Y$  in any state during each visit are positive and finite, i.e. we assume that for  $\ell = 1, 2$  and any  $i, j = 1, \dots, N$

$$0 < d_{ij}^{(\ell)} = \int_0^\infty \tau^\ell dF_{ij}(\tau) < \infty \quad \text{hence also} \quad 0 < d_i^{(\ell)} = \sum_{j=1}^N p_{ij} d_{ij}^{(\ell)} < \infty.$$

### 3 Preliminaries

As well-known, under AS 1 the rows of the limiting matrix  $\mathbf{P}^* = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=0}^{n-1} \mathbf{P}^k$  of the embedded Markov chain  $X$  are identical and equal to the (row) vector of steady state probabilities  $\boldsymbol{\pi} = [\pi_1, \dots, \pi_N]$  as determined by  $\boldsymbol{\pi} = \boldsymbol{\pi} \cdot \mathbf{P}$  along with the normalizing condition  $\boldsymbol{\pi} \cdot \mathbf{e} = 1$ . Moreover, since  $\mathbf{P}$  is aperiodic  $\mathbf{P}^* = \lim_{k \rightarrow \infty} \mathbf{P}^k$  and the convergence is geometrically fast. Then

$$g = \frac{\sum_{j \in \mathcal{I}} \pi_j \cdot r_j^{(1)}}{\sum_{j \in \mathcal{I}} \pi_j \cdot d_j^{(1)}} \quad (4)$$

with  $r_j^{(1)}$  defined by (16) and  $d_j^{(1)}$  the first moment of the (random) time  $\eta_j$ . Also the average reward  $g^r$  and the average time  $g^t$  per transition of the embedded Markov chain are well defined by:

$$g^r = \sum_{j \in \mathcal{I}} \pi_j \cdot r_j^{(1)}, \quad g^t = \sum_{j \in \mathcal{I}} \pi_j \cdot d_j^{(1)}.$$

Obviously for the constant vectors  $\mathbf{g}^r$  and  $\mathbf{g}^t$  with elements  $g^r$  and  $g^t$  we have  $\mathbf{g}^r = \mathbf{P}^* \cdot \mathbf{r}^{(1)}$ ,  $\mathbf{g}^t = \mathbf{P}^* \cdot \mathbf{d}^{(1)}$ , respectively.

From the theory of dynamic programming (e.g. [11], [12]) it is well-known that the vectors of expected reward and times for  $n$  transitions fulfil the recursive formulas

$$\mathbf{R}(n+1) = \mathbf{r}^{(1)} + \mathbf{P} \cdot \mathbf{R}(n) \quad \text{with} \quad \mathbf{R}(0) = 0 \quad (5)$$

$$\mathbf{D}(n+1) = \mathbf{d}^{(1)} + \mathbf{P} \cdot \mathbf{D}(n) \quad \text{with} \quad \mathbf{D}(0) = 0. \quad (6)$$

Furthermore, under assumption AS 1, there exists vectors  $\mathbf{w}^r, \mathbf{w}^t$  such that

$$\mathbf{R}(n) = \mathbf{g}^r \cdot n + \mathbf{w}^r + \boldsymbol{\varepsilon}(n) \quad (7)$$

$$\mathbf{D}(n) = \mathbf{g}^t \cdot n + \mathbf{w}^t + \boldsymbol{\varepsilon}(n). \quad (8)$$

Hence, both  $\mathbf{R}(n)$  and  $\mathbf{D}(n)$  possess a linear growth rate  $g^r$  and  $g^t$  in  $n$  up to a geometric convergence to the null vector and vectors  $\mathbf{w}^r$  and  $\mathbf{w}^t$ , respectively. The constant vectors  $\mathbf{g}^r$ ,  $\mathbf{g}^t$  along with vectors  $\mathbf{w}^r$ ,  $\mathbf{w}^t$  are uniquely determined by

$$\mathbf{w}^r + \mathbf{g}^r = \mathbf{r}^{(1)} + \mathbf{P} \cdot \mathbf{w}^r, \quad \mathbf{P}^* \cdot \mathbf{w}^r = 0 \quad (9)$$

$$\mathbf{w}^t + \mathbf{g}^t = \mathbf{d}^{(1)} + \mathbf{P} \cdot \mathbf{w}^t, \quad \mathbf{P}^* \cdot \mathbf{w}^t = 0. \quad (10)$$

From (9), (10) we can also immediately conclude the existence of a constant vector  $\mathbf{g}$  along with vector  $\mathbf{w}$  being the unique solution of the equation

$$\mathbf{w} + \mathbf{D}^{(1)} \cdot \mathbf{g} = \mathbf{r}^{(1)} + \mathbf{P} \cdot \mathbf{w}, \quad \mathbf{P}^* \cdot \mathbf{w} = 0, \quad (11)$$

where the  $i$ th element of the diagonal matrix  $\mathbf{D}^{(1)}$  equals  $d_i^{(1)}$ , where the elements of the constant vector  $\mathbf{g}$  are equal to  $g = g^r/g^t$  and where for the  $i$ th element of  $\mathbf{w}$ ,  $w_i = w_i^r - g \cdot w_i^t$ .

*Conditional Expectations of Reward.* Let  $\mathbf{E}_i$  denote the conditional expectation given  $X(0) = i$ . Since the process  $X$  is time homogeneous, for  $n > m$  we can conclude that

$$\mathbf{E}_i[\xi_n] = \mathbf{E}_i[\xi_m] + \mathbf{E}_i\left\{\sum_{j \in \mathcal{I}} \mathbf{P}(X_m = j) \cdot \mathbf{E}_j[\xi_{n-m}]\right\}. \quad (12)$$

Similarly we get

$$\begin{aligned} \mathbf{E}_i[\xi_n]^2 &= \mathbf{E}_i[\xi_m]^2 + \mathbf{E}_i\left\{\sum_{j \in \mathcal{I}} \mathbf{P}(X_m = j) \cdot \mathbf{E}_j[\xi_{n-m}]^2\right\} \\ &\quad + 2 \cdot \mathbf{E}_i[\xi_m] \sum_{j \in \mathcal{I}} \mathbf{P}(X_m = j) \cdot \mathbf{E}_j[\xi_{n-m}]. \end{aligned} \quad (13)$$

#### 4 Reward Variance of Embedded Markov Chains

In this section we focus our attention on the total reward variance generated by the embedded Markov chain  $X = \{X_n, n = 0, 1, \dots\}$  characterized by the transition probability matrix  $\mathbf{P}$  and one-stage (random) rewards  $\eta_{ij} \cdot r_i + r_{ij}$  ( $i, j \in \mathcal{I}, j \neq i$ ) accrued to the transition from state  $i$  into state  $j$ .

From (12) and (13) we directly conclude

$$\mathbf{E}_i[\xi_n] = r_i^{(1)} + \sum_{j \in \mathcal{I}} p_{ij} \cdot \mathbf{E}_j[\xi_{n-1}] \quad (14)$$

$$\begin{aligned} \mathbf{E}_i[\xi_n]^2 &= r_i^{(2)} + \sum_{j \in \mathcal{I}} p_{ij} \cdot \mathbf{E}_j[\xi_{n-1}]^2 \\ &\quad + 2 \cdot \sum_{j \in \mathcal{I}} p_{ij} \cdot [r_i \cdot d_{ij}^{(1)} + r_{ij}] \cdot \mathbf{E}_j[\xi_{n-1}] \end{aligned} \quad (15)$$

where

$$r_i^{(1)} = \sum_{j \in \mathcal{I}} p_{ij} \cdot [r_i \cdot d_{ij}^{(1)} + r_{ij}] \quad (16)$$

$$r_i^{(2)} = \sum_{j \in \mathcal{I}} p_{ij} \cdot \{[r_i]^2 \cdot d_{ij}^{(2)} + [r_{ij}]^2 + 2 \cdot r_i \cdot d_{ij}^{(1)} \cdot r_{ij}\}. \quad (17)$$

By using the more appealing notation  $R_i(n) = \mathbb{E}_i[\xi_n]$ ,  $S_i(n) = \mathbb{E}_i[\xi_n]^2$ , (14), (15) take on the forms:

$$R_i(n+1) = r_i^{(1)} + \sum_{j \in \mathcal{I}} p_{ij} \cdot R_j(n), \quad (18)$$

$$\begin{aligned} S_i(n+1) &= r_i^{(2)} + 2 \cdot \sum_{j \in \mathcal{I}} p_{ij} \cdot [r_i \cdot d_{ij}^{(1)} + r_{ij}] \cdot R_j(n) \\ &\quad + \sum_{j \in \mathcal{I}} p_{ij} \cdot S_j(n). \end{aligned} \quad (19)$$

For the variance  $V_i(\cdot) = S_i(\cdot) - [R_i(\cdot)]^2$  we now arrive after some algebra by using (18), (19)

$$\begin{aligned} V_i(n+1) &= r_i^{(2)} + 2 \cdot \sum_{j \in \mathcal{I}} p_{ij} \cdot [r_i \cdot d_{ij}^{(1)} + r_{ij}] \cdot R_j(n) \\ &\quad - \sum_{j \in \mathcal{I}} p_{ij} \cdot [R_i(n+1) + R_j(n)] \cdot [R_i(n+1) - R_j(n)] \\ &\quad + \sum_{j \in \mathcal{I}} p_{ij} \cdot V_j(n). \end{aligned} \quad (20)$$

*Remark 1 (Transient case.)* For the so-called transient case, i.e. when

$$\lim_{n \rightarrow \infty} R_i(n) = R_i, \quad \lim_{n \rightarrow \infty} S_i(n) = S_i, \quad \lim_{n \rightarrow \infty} V_i(n) = V_i,$$

(18), (19) and (20) obtain the closed form:

$$R_i = r_i^{(1)} + \sum_{j \in \mathcal{I}} p_{ij} \cdot R_j, \quad (21)$$

$$S_i = r_i^{(2)} + 2 \cdot \sum_{j \in \mathcal{I}} p_{ij} \cdot (r_i \cdot d_{ij}^{(1)} + r_{ij}) \cdot R_j + \sum_{j \in \mathcal{I}} p_{ij} \cdot S_j, \quad (22)$$

$$\begin{aligned} V_i &= r_i^{(2)} - [R_i]^2 + \sum_{j \in \mathcal{I}} p_{ij} \cdot \{2 \cdot (r_i \cdot d_{ij}^{(1)} + r_{ij}) \cdot R_j + [R_j]^2\} \\ &\quad + \sum_{j \in \mathcal{I}} p_{ij} \cdot V_j. \end{aligned} \quad (23)$$

## 5 Mean Reward Variance in Markov Chains: Infinite Horizon

To investigate the asymptotic behaviour of the variance  $V_i(n)$  we focus our attention on the recursive formula (20) and employ well-known facts on the asymptotic behaviour of  $R_i(n)$  as based upon (7).

By (7) in the third term of (20), we can substitute:

$$R_i(n+1) + R_j(n) = 2 \cdot n \cdot g^r + g^r + w_i^r + w_j^r + \varepsilon(n), \quad (24)$$

$$R_i(n+1) - R_j(n) = g^r + w_i^r - w_j^r + \varepsilon(n). \quad (25)$$

Hence

$$\begin{aligned} & \sum_{j \in \mathcal{I}} p_{ij} \cdot [R_i(n+1) + R_j(n)] \cdot [R_i(n+1) - R_j(n)] \\ &= 2 \cdot n \cdot g^r \cdot (g^r + w_i^r - \sum_{j \in \mathcal{I}} p_{ij} \cdot w_j^r) + \sum_{j \in \mathcal{I}} p_{ij} \cdot \{[g^r + w_i^r]^2 - [w_j^r]^2\} + \varepsilon(n) \\ &= 2 \cdot n \cdot g^r \cdot r_i^{(1)} + \sum_{j \in \mathcal{I}} p_{ij} \cdot \{[g^r + w_i^r]^2 - [w_j^r]^2\} + \varepsilon(n). \end{aligned} \quad (26)$$

Similarly for the second term of (20) we obtain by (7)

$$\begin{aligned} & \sum_{j \in \mathcal{I}} p_{ij} \cdot [r_i \cdot d_{ij}^{(1)} + r_{ij}] \cdot R_j(n) \\ &= \sum_{j \in \mathcal{I}} p_{ij} \cdot [r_i \cdot d_{ij}^{(1)} + r_{ij}] \cdot [n \cdot g^r + w_j^r + \varepsilon(n)] \\ &= n \cdot g^r \cdot r_i^{(1)} + \sum_{j \in \mathcal{I}} p_{ij} \cdot [r_i \cdot d_{ij}^{(1)} + r_{ij}] \cdot w_j^r + \varepsilon(n). \end{aligned} \quad (27)$$

Substituting (26) and (27) in (20) now yields

$$\begin{aligned} V_i(n+1) &= \sum_{j \in \mathcal{I}} p_{ij} \cdot V_j(n) + r_i^{(2)} + 2 \cdot \sum_{j \in \mathcal{I}} p_{ij} \cdot [r_i \cdot d_{ij}^{(1)} + r_{ij}] \cdot w_j^r \\ &\quad - \sum_{j \in \mathcal{I}} p_{ij} \cdot \{[g^r + w_i^r]^2 - [w_j^r]^2\} + \varepsilon(n) \\ &= \sum_{j \in \mathcal{I}} p_{ij} \cdot V_j(n) + s_i + \varepsilon(n), \end{aligned} \quad (28)$$

where for the elements  $s_i$  of the vector  $\mathbf{s}$  we obtain after some algebra:

$$\begin{aligned} s_i &= r_i^{(2)} + \sum_{j \in \mathcal{I}} p_{ij} \cdot (2 \cdot [r_i \cdot d_{ij}^{(1)} + r_{ij}] \cdot w_j^r + [w_j^r]^2) - [g^r + w_i^r]^2 \\ &= \sum_{j \in \mathcal{I}} p_{ij} \cdot \{[r_i]^2 \cdot d_{ij}^{(2)} + 2 \cdot r_i \cdot d_{ij}^{(1)} \cdot [r_{ij} + w_j^r] + [r_{ij} + w_j^r]^2\} \\ &\quad - [g^r + w_i^r]^2. \end{aligned} \quad (29)$$

Furthermore, by (9) for the last term of (29) we get

$$\begin{aligned} -[g^r + w_i^r]^2 &= [g^r]^2 - [w_i^r]^2 - 2 \cdot g^r \cdot [g^r + w_i^r] \\ &= [g^r]^2 - [w_i^r]^2 - 2 \cdot g^r \cdot \sum_{j \in \mathcal{I}} p_{ij} \cdot [r_i \cdot d_{ij}^{(1)} + r_{ij} + w_j^r] \end{aligned}$$

which by substitution in (29) and some algebra yields

$$\begin{aligned} s_i &= \sum_{j \in \mathcal{I}} p_{ij} \cdot \{[r_i]^2 \cdot d_{ij}^{(2)} + [r_{ij}]^2 + 2 \cdot r_i \cdot d_{ij}^{(1)} \cdot r_{ij} \\ &\quad + 2 \cdot [r_i \cdot d_{ij}^{(1)} + r_{ij}] \cdot w_j^r + [w_j^r]^2\} + [g^r]^2 - [w_i^r]^2 \\ &\quad - 2 \cdot g^r \cdot \sum_{j \in \mathcal{I}} p_{ij} \cdot [r_i \cdot d_{ij}^{(1)} + r_{ij} + w_j^r] \\ &= \sum_{j \in \mathcal{I}} p_{ij} \cdot \{[r_i]^2 \cdot d_{ij}^{(2)} + 2 \cdot r_i \cdot d_{ij}^{(1)} \cdot [r_{ij} - g^r + w_j^r] \\ &\quad + [r_{ij} - g^r + w_j^r]^2\} - [w_i^r]^2. \end{aligned} \quad (30)$$

Hence, in matrix form we have:

$$\mathbf{V}(n+1) = \mathbf{s} + \mathbf{P} \cdot \mathbf{V}(n) + \boldsymbol{\varepsilon}^{(1)}(n) \quad (31)$$

where the elements of the vector  $\boldsymbol{\varepsilon}^{(1)}(n)$  converge to zero geometrically, that is, for some numbers  $c > 0$  and  $\delta \in (0, 1)$ :  $\|\boldsymbol{\varepsilon}^{(1)}(n)\| \leq c \cdot \delta^n$ .

*Growth rate.* In order to investigate the behaviour of  $\mathbf{V}(n)$  for  $n$  large, let

$$\mathbf{W}(n+1) = \mathbf{s} + \mathbf{P} \cdot \mathbf{W}(n) \quad (32)$$

and

$$\mathbf{X}(n) = \mathbf{V}(n) - \mathbf{W}(n). \quad (33)$$

Hence, in line with (5) for  $\mathbf{R}(n)$ , the vector  $\mathbf{W}(n)$  can be regarded as the total reward vector over  $n$  steps for a Markov reward chain with one step reward vector  $\mathbf{s}$ . As a consequence, similarly to (7), we can conclude the existence of a vector  $\mathbf{w}^{(2)}$  such that

$$\mathbf{W}(n) = \mathbf{g}^{(2)} \cdot n + \mathbf{w}^{(2)} + \boldsymbol{\varepsilon}(n). \quad (34)$$

In words that is,  $\mathbf{W}(n)$  possesses a linear growth rate  $\mathbf{g}^{(2)}$  in  $n$  up to a geometric convergence to the null vector and vector  $\mathbf{w}^{(2)}$ . The constant vector  $\mathbf{g}^{(2)}$  along with vector  $\mathbf{w}^{(2)}$  are uniquely determined by

$$\mathbf{w}^{(2)} + \mathbf{g}^{(2)} = \mathbf{s} + \mathbf{P} \cdot \mathbf{w}^{(2)}, \quad \mathbf{P}^* \cdot \mathbf{w}^{(2)} = 0 \quad (35)$$

where elements of  $\mathbf{s}$  are calculated by (29) or (30).

Iterating (33) we immediately conclude that

$$\|\mathbf{X}(n)\| \leq \left\| \sum_{k=1}^n \mathbf{P}^k \cdot \boldsymbol{\varepsilon}^{(1)}(k) \right\| \leq \sum_{k=1}^n \|\boldsymbol{\varepsilon}^{(1)}(k)\| < c \cdot \frac{1}{1-\delta} =: C. \quad (36)$$



**Theorem 1** For the (constant) growth rate  $g^{(2)}$  (vector  $\mathbf{g}^{(2)}$ ) and vector  $\mathbf{w}^{(2)}$  as determined by (35), some geometrically converging function  $\varepsilon(n)$ , some constant  $C$  and some bounded vector  $\mathbf{c}(n)$  with  $\|\mathbf{c}(n)\| < C$

$$\mathbf{V}(n) = n \cdot \mathbf{g}^{(2)} + \mathbf{w}^{(2)} + \varepsilon(n) + \mathbf{c}(n) \text{ for all } n \quad (37)$$

and hence,

$$\mathbf{g}^{(2)} = \lim_{n \rightarrow \infty} \frac{\mathbf{V}(n)}{n} = \lim_{n \rightarrow \infty} \frac{\mathbf{W}(n)}{n} = \mathbf{P}^* \cdot \mathbf{s}. \quad (38)$$

*Proof* The relations follow directly combining (31), (32), (33) and (36) with  $\mathbf{c}(n)$  equal to  $\mathbf{X}(n)$ .

*Remark 2* By (9) the coefficients  $w_i^r$ 's are normalized by  $\sum_{i \in \mathcal{I}} \pi_i \cdot w_i^r = 0$ . Since by assumption AS 1 the Markov chain has a single class of recurrent states, the first equation in (9) still holds if we replace  $w_i^r$  by  $\tilde{w}_i^r = w_i^r + c$  ( $i = 1, 2, \dots, N$ ), where  $c$  is an arbitrary constant.

As a consequence, by examining the formulas (29), (30) we can conclude that (38) still holds if in (29), (30) we replace all  $w_i^r$ 's by  $\tilde{w}_i^r$ 's.

Since  $\sum_{i \in \mathcal{I}} \pi_i \cdot p_{ij} = \pi_j$  for arbitrary real  $c_i$ 's it holds

$$\sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{I}} \pi_i \cdot p_{ij} \cdot c_j = \sum_{i \in \mathcal{I}} \pi_i \cdot c_i.$$

Then, since by (9):  $\sum_{i \in \mathcal{I}} \pi_i \cdot w_i^r = 0$  and defining

$$\tilde{s}_i = r_i^{(2)} - [g^r]^2 + 2 \cdot \sum_{j \in \mathcal{I}} p_{ij} \cdot [r_i \cdot d_{ij}^{(1)} + r_{ij}] \cdot w_j^r \quad (39)$$

we can conclude that

$$g^{(2)} = \sum_{i \in \mathcal{I}} \pi_i \cdot s_i = \sum_{i \in \mathcal{I}} \pi_i \cdot \tilde{s}_i. \quad (40)$$

*Remark 3* Assume that  $r_i \equiv 0$  for all  $i \in \mathcal{I}$ . Then (cf. (29), (30))  $s_i$  takes on the form:

$$\begin{aligned} s_i &= \sum_{j \in \mathcal{I}} p_{ij} \cdot \{[r_{ij} + w_j^r]^2\} - [g^r + w_i^r]^2 \\ &= \sum_{j \in \mathcal{I}} p_{ij} \cdot \{[r_{ij} - g^r + w_j^r]^2\} - [w_i^r]^2. \end{aligned} \quad (41)$$

Moreover, by (39) we have

$$\begin{aligned} \tilde{s}_i &= \sum_{j \in \mathcal{I}} p_{ij} \cdot \{[r_{ij}]^2 + 2 \cdot r_{ij} \cdot w_j^r\} - [g^r]^2 \\ &= \sum_{j \in \mathcal{I}} p_{ij} \cdot \{[r_{ij} - g^r]^2 + 2 \cdot r_{ij} \cdot w_j^r\}. \end{aligned} \quad (42)$$

## 6 Mean Reward Variance in Semi-Markov Processes

In analogy with (2) and (4) we are thus interested in the “mean” variance defined by

$$G := \frac{g^{(2)}}{g^t} = \frac{\sum_{j \in \mathcal{I}} \pi_j \cdot s_j}{\sum_{j \in \mathcal{I}} \pi_j \cdot d_j^{(1)}}. \quad (43)$$

Below we present two approaches how to calculate this mean variance  $G$ .

**Theorem 2** *The mean variance  $G$  given by (43) can be calculated using the limiting probabilities of the semi-Markov processes as*

$$G = \sum_{j \in \mathcal{I}} \bar{\pi}_j \cdot \bar{s}_j = \sum_{j \in \mathcal{I}} \bar{\pi}_j \cdot \hat{s}_j \quad \text{independently of } i \in \mathcal{I} \quad (44)$$

where

$$\bar{s}_i = s_i / d_i^{(1)}, \quad \hat{s}_i = \tilde{s}_i / d_i^{(1)} \quad (45)$$

and

$$\bar{\pi}_i = \frac{\pi_i \cdot d_i^{(1)}}{\sum_{j=1}^N \pi_j \cdot d_j^{(1)}}. \quad (46)$$

*Proof* Since under assumptions AS 1 and AS 2, the rows of the limiting matrix  $\bar{\mathbf{P}}^*$  of the considered semi-Markov process  $Y(t)$  are identical and equal to the row vector  $\bar{\boldsymbol{\pi}} = [\bar{\pi}_1, \dots, \bar{\pi}_N]$  where elements  $\bar{\pi}_i = \lim_{t \rightarrow \infty} \{Y(t) | Y(0) = i\}$  are given by (46) (see e.g. [11], [12]), the proof follows immediately by inserting (45) and (46) in (40).

In what follows we show how to calculate the mean variance using formulas similar to those for calculating the mean reward. To this end, let  $\mathcal{G}$  be a diagonal matrix with all diagonal elements equal to  $G = g^{(2)}/g^t$ . Premultiplying (10) by  $\mathcal{G}$  and subtracting from (35) we immediately get

$$\mathbf{u} = \mathbf{s} - \mathcal{G} \cdot \mathbf{d}^{(1)} + \mathbf{P} \cdot \mathbf{u} \quad (47)$$

where

$$\mathbf{u} = \mathbf{w}^{(2)} - \mathcal{G} \cdot \mathbf{w}^t. \quad (48)$$

So we have arrived at

**Theorem 3** *The mean variance  $G$  can be calculated as a solution of (47). The solution  $(G, \mathbf{u})$  of (47) is unique up to an additive constant to vectors  $\mathbf{u}$ , and unique under the additional normalizing condition  $\mathbf{P}^* \cdot \mathbf{u} = 0$ .*

### Acknowledgement

The author is grateful to Arie Hordijk for his stimulating work in the field of Markov Decision Theory.

The research was supported by the Grant Agency of the Czech Republic under Grants 402/05/0115 and 402/04/1294.

## References

1. Benito F (1982) Calculating the variance in Markov processes with random reward. *Trabajos de Estadística y de Investigación Operativa* 33: 73–85
2. Filar J, Kallenberg LCM and Lee H-M (1989) Variance penalized Markov decision processes. *Mathem. Oper. Research* 14: 147–161
3. Huang Y and Kallenberg LCM (1994) On finding optimal policies for Markov decision chains: a unifying framework for mean-variance-tradeoffs. *Mathem. Oper. Research* 19: 434–448
4. Jaquette SC (1972) Markov decision processes with a new optimality criterion: Small interest rates. *Ann. Math. Statist.* 43: 1894–1901
5. Jaquette SC (1973) Markov decision processes with a new optimality criterion: Discrete time. *Ann. Statist.* 1: 496–505
6. Jaquette SC (1975) Markov decision processes with a new optimality criterion: Continuous time. *Ann. Statist.* 3: 547–553
7. Kadota Y (1997) A minimum average-variance in Markov decision processes. *Bulletin of Informatics and Cybernetics* 29: 83–89
8. Kawai H (1987) A variance minimization problem for a Markov decision process. *European J. Oper. Research* 31: 140–145
9. Kurano M (1987) Markov decision processes with a minimum-variance criterion. *J. Mathem. Anal. Appl.* 123: 572–583
10. Mandl P (1971) On the variance in controlled Markov chains. *Kybernetika* 7: 1–12
11. Puterman ML (1994) *Markov decision processes – discrete stochastic dynamic programming*. Wiley, New York
12. Ross SM (1970) *Applied probability models with optimization applications*. Holden-Day, San Francisco, Calif.
13. Sladký K and Sitař M (2004) Optimal solutions for undiscounted variance penalized Markov decision chains. In: Marti K, Ermoliev Y and Pflug G (eds.) *Dynamic Stochastic Optimization*. Springer-Verlag, Berlin, pp. 43–66
14. Sladký K and van Dijk NM (2004) On total reward variance for continuous-time Markov reward chains. Submitted for publication
15. Sobel MJ (1982) The variance of discounted Markov decision processes. *J. Appl. Probab.* 19: 794–802
16. Sobel MJ (1985) Maximal mean/standard deviation ratio in an undiscounted MDP. *Oper. Research Lett.* 4: 157–159
17. White DJ (1988) Mean variance and probability criteria in finite Markov decision processes: a review. *J. Optim. Theory Appl.* 56: 1–29