# Infinite Markov decision processes with decision independent subset of states[*]

András Mészáros

MTA-BME Information Systems Research Group

1117, Magyar Tudosok krt. 2, Budapest, Hungary

meszarosa@hit.bme.hu

Miklós Telek

Dept. of Networked Systems and Services

Budapest University of Technology and Economics

1117, Magyar Tudosok krt. 2, Budapest, Hungary

telek@hit.bme.hu

October 3, 2018

**Abstract**

In this work, we examine Markov Decision Processes (MDPs) that are composed of a finite subset of states with decisions and a (potentially infinite) subset of states without decisions. We show that, if some parameters of these MDPs can be computed efficiently, then a transformation to a condensed representation is possible, in which only the states with decisions are kept, and the rewards and transition rates are modified such that the optimal policy is the same in the original and the modified MDP.

When the subset of states without a decision is infinite, the analysis is based on some structural regularity of the process. Two practically important structures are considered the M/G/1-type and the G/M/1-type.

Keywords: Markov decision process, equivalent representation, state space reduction, Markov chain with regular structure.

---

# 1  Introduction

Markov decision processes (MDPs) give a powerful yet simple tool to formalize and solve decision problems and are commonly used in a wide variety of fields from machine learning [11] through telecommunication [2] to finance [3]. One of the biggest issues when using MDPs is known as state space explosion, that is, the number of states of the MDP usually grows exponentially with the number of variables of the analysed system. Because of this, the number of states in the MDP can easily increase to a point where the classical MDP solving methods cannot be used. To overcome this problem, some techniques have been developed. One possibility is to prove that the optimal policy in the MDP has threshold form, which means that the optimal decision in any state can be determined based on whether a certain parameter is above a fixed value (threshold). In this case, finding the threshold is enough to solve the optimization problem. For instance, accepting requests to a queue may be optimal until the queue length reaches a certain value. See e.g. [12] or [8] for more examples. Another noteworthy approach is presented in [6], where more efficient solution methods for MDPs with factored representations are considered. These cannot be applied, however for MDPs that are infinite or cannot be factorized.

Apart from exact optimal solutions, one can get a quasi-optimal solution by using certain approximation techniques. One possible approach is the truncation of the state space. This may happen based on the physical model (e.g. the size of the buffer is constrained), as in [17] and [10] for example. Alternatively, one can truncate based on only mathematical considerations as discussed by [1]. Another interesting method is presented in [16], where a so-called deterministic simulative model is introduced. The essence of this model is that the original MDP is transformed in such a way, that transitions of the new model all become deterministic.

In this work, we propose an efficient reduction method that can be used for MDPs which are composed of a finite subset of states with decisions and a finite or infinite subset of states without decisions. More specifically we give a method that compresses the MDP to the size of the subset with decisions. The presented method requires that some important parameters of the MDP can be calculated efficiently. We also show how to use the presented reduction method for infinite state MDPs with QBD, M/G/1-type and G/M/1-type structures, which are the most prevalent classes of infinite MDPs in queueing problems.

The problem of reducing the size of the state space of such systems has been considered in [14], but the solution proposed there was suboptimal in the sense that the reduced state space was much larger than the set of states with a decision. The solution proposed here is optimal in this sense.

The rest of the paper is organized as follows. The paper starts with an example in Section 2 to motivate the forthcoming analysis. Section 3 provides a

summary of MDPs. The parameters of the MDP compression method are provided in Section 4, while the compression method is stated and proved in Section 5. The special cases of infinite MDPs with M/G/1 type and G/M/1 type structures are provided in Section 6, whose reward parameters a computed in Section 7.

## 2 A motivating example

The problem considered in this paper has a strong practical motivation which is detailed in [5, 14]. Here we summarize the problem and the related model for completeness.

In case of multi-server systems with identical but state dependent servers, it is an interesting optimization problem to properly assign new incoming jobs with one of the free servers, if more than one server is idle at customer arrival. Consider a simple queueing systems with multiple MAP servers, where the incoming customers can be freely assigned with service unit in case of more than one available free servers.

In particular, [14] considers the optimal control of MAP/MAP/n queues. The simplest version of such models is the M/MAP/2 queue where customers arrive according to a Poisson process with rate $\lambda$, the service process of each server is a MAP with two states characterized by the matrix pair $(\boldsymbol{S_0}, \boldsymbol{S_1})$. The associated state transition structure is

$$\boldsymbol{Q}(a) = \begin{pmatrix} \boldsymbol{L}_0 & \boldsymbol{F}_0(a) & \boldsymbol{0} & \cdots & \\ \boldsymbol{B}_1 & \boldsymbol{L}_1 & \boldsymbol{F}_1 & 0 & \cdots \\ 0 & \boldsymbol{B}_2 & \boldsymbol{L} & \boldsymbol{F} & 0 \\ \vdots & 0 & \boldsymbol{B} & \boldsymbol{L} & \boldsymbol{F} & \ddots \\ & & \ddots & \ddots & \ddots & \ddots \end{pmatrix},$$

where

$$\boldsymbol{L}_0 = -\lambda \boldsymbol{I} \otimes \boldsymbol{I}, \boldsymbol{F}_0(a) = \lambda \left( (\boldsymbol{I} \otimes \boldsymbol{I}) \boldsymbol{P}(a), \quad (\boldsymbol{I} \otimes \boldsymbol{I}) (\boldsymbol{I} - \boldsymbol{P}(a)) \right),$$

$$\boldsymbol{B}_1 = \begin{pmatrix} \boldsymbol{I} \otimes \boldsymbol{S}_1 \\ \boldsymbol{S}_1 \otimes \boldsymbol{I} \end{pmatrix}, \boldsymbol{L}_1 = \begin{pmatrix} -\lambda \boldsymbol{I} \otimes \boldsymbol{I} + \boldsymbol{I} \otimes \boldsymbol{S}_0 & 0 \\ 0 & -\lambda \boldsymbol{I} \otimes \boldsymbol{I} + \boldsymbol{S}_0 \otimes \boldsymbol{I} \end{pmatrix},$$

$$\boldsymbol{F}_1 = \lambda \begin{pmatrix} \boldsymbol{I} \otimes \boldsymbol{I} \\ \boldsymbol{I} \otimes \boldsymbol{I} \end{pmatrix}, \boldsymbol{B}_2 = \begin{pmatrix} \boldsymbol{I} \otimes \boldsymbol{S}_1 \\ \boldsymbol{S}_1 \otimes \boldsymbol{I} \end{pmatrix}, \boldsymbol{L} = -\lambda \boldsymbol{S}_0 \oplus \boldsymbol{S}_0,$$

$$\boldsymbol{F} = \boldsymbol{A}_1 \otimes \boldsymbol{I} \otimes \boldsymbol{I}, \boldsymbol{B} = \boldsymbol{I} \otimes \boldsymbol{S}_1 + \boldsymbol{S}_1 \otimes \boldsymbol{I}.$$

and the matrix which is responsible for the decision upon customer arrival to the idle system is

$$\boldsymbol{P}(a_1) = \operatorname{diag}(1/2, 1, 0, 1/2) \text{ and } \boldsymbol{P}(a_2) = \operatorname{diag}(1/2, 0, 1, 1/2).$$

3

According to matrix $\boldsymbol{P}(a)$, at a customer arrival to the empty system the customer is directed to the server in phase 1 by choosing action $a_1$ and to the server in phase 2 by choosing action $a_1$, if the servers are in different phases. If idle servers are in the same phase, the service units are chosen evenly.

The associated reward matrix

$$
\boldsymbol{C}(a) = \begin{pmatrix} \boldsymbol{I} & & & & \\ & \frac{1}{2}\boldsymbol{I} & & & \\ & & \boldsymbol{0} & & \\ & & & \boldsymbol{0} & \\ & & & & \ddots \end{pmatrix},
$$

which is decision independent, intends to maximize the server idle time.

For curiosity, we note that the counter intuitive conclusion gained by the analysis of this model in [5] is that it worth to choose the slower server, because it results in a better system state for higher levels of system saturation.

# 3 Theoretical background

## 3.1 Markov Decision Processes

In the paper, we consider continuous time, time-homogeneous, non-discounted, MDPs with the following definition.

**Definition 3.1.** *Let $X(t)$ be a continuous time Markov chain with state space $S$, $A$ a set of decisions, $\alpha_0$ an initial probability vector, $\boldsymbol{Q}(a)$ a decision dependent generator matrix satisfying $\boldsymbol{Q}(a)\underline{\boldsymbol{1}} = 0$ for $\forall a \in A$ (where $\underline{\boldsymbol{1}}$ is the column vector of $1$s with appropriate size), $\boldsymbol{C}(a)$ is a decision dependent diagonal reward rate matrix. We say that the tuple $(S, A, \alpha_0, \boldsymbol{Q}(a), \boldsymbol{C}(a))$ is a continuous time Markov decision process. The generator of the MDP, $\boldsymbol{Q}(a)$,*

For such MDPs the usual optimization problem is to find a policy (state-decision mapping) $\pi^*(s) \in \{\pi(s) : S \to A\}$ such that

$$
\pi^* = \arg\max_{\pi} \mathrm{E}_{\pi} \left[ \lim_{T \to \infty} \frac{1}{T} \int_{t=0}^{T} \boldsymbol{C}_{X(t),X(t)}(\pi(X(t))dt \right].
$$

Throughout the paper we assume that for every policy the Markov model is composed of exactly one communicating block and potentially one transient block. In this case, the optimal policy can also be expressed as

$$
\pi^* = \arg\max_{\pi} \alpha(\pi)\boldsymbol{C}(\pi)\underline{\boldsymbol{1}}, \tag{1}
$$

4

where $\alpha(\pi)$ is the steady state probability vector for policy $\pi$, that satisfies,

$$\alpha(\pi)\boldsymbol{Q}(\pi) = \underline{\boldsymbol{0}}, \quad \alpha(\pi)\underline{\boldsymbol{1}} = 1,$$

where $\underline{\boldsymbol{0}}$ is the column vector of 0s of appropriate size [9].

The definition of MDPs does not constrain the sign of the elements of the $\boldsymbol{C}(a)$ reward matrix, however offsetting $\boldsymbol{C}(a)$ with a constant value $c$ does not change the optimal policy since for any policy $\pi$

$$\alpha(\pi)\left(\boldsymbol{C}(\pi) + c\boldsymbol{I}\right)\underline{\boldsymbol{1}} = \alpha(\pi)\boldsymbol{C}(\pi)\underline{\boldsymbol{1}} + c,$$

that is, the optimal policy is the same for an MDP with reward matrix $\boldsymbol{C}(a)$ and $\boldsymbol{C}'(a) = \boldsymbol{C}(a) + c\boldsymbol{I}, \forall c \in \mathbb{R}$, therefore in the following we assume that $\min_{a,i} \boldsymbol{C}(a)_{ii} > 0$, that is, the reward rate in every state is positive for every decision.

## 3.2 Basic transformation of MDPs

Our main goal is to examine MDPs for which the state space $S$ can be partitioned into two disjoint subsets $S_U$ and $S_D$ ($S_U \cup S_D = S$, $S_U \cap S_D = \emptyset$) where $S_U$ is finite and contains all the states in which decisions can be made and $S_D$ is potentially infinite and contains only states where decisions are not made, or decisions have the same effect (i.e., $\boldsymbol{Q}_{ij}(a_k) = \boldsymbol{Q}_{ij}(a_\ell), \forall i, j \in S_D, k, \ell \in A$). Without loss of generality we assume that the states in $S_U$ have lower indexes than the states in $S_D$ (i.e., $i < j, \forall i \in S_U, j \in S_D$), thus the $\boldsymbol{Q}(\pi)$ generator matrix and the $\boldsymbol{C}(\pi)$ reward-rate matrix have the following block structure

$$\boldsymbol{Q}(\pi) = \begin{pmatrix} \boldsymbol{Q_U}(\pi) & \boldsymbol{Q_{UD}}(\pi) \\ \boldsymbol{Q_{DU}} & \boldsymbol{Q_D} \end{pmatrix}, \quad \boldsymbol{C}(\pi) = \begin{pmatrix} \boldsymbol{C_U}(\pi) & 0 \\ 0 & \boldsymbol{C_D} \end{pmatrix}, \qquad (2)$$

where $\pi$ in the argument indicates that the respective part of the matrix depends on the actual policy. In the rest of the paper we assume that $S_D$ is transient with finite sojourn time, consequently $\boldsymbol{Q_D}$ is non-singular and the $(i, j)$ element of $(-\boldsymbol{Q_D})^{-1}$ is the mean time spent in state $j \in S_D$ before leaving $S_D$ starting from $i \in S_D$.

### 3.2.1 Transformation of MDPs with no decisions in $S_D$

Let $\alpha(\pi)$ be the stationary probability vector of the Markov chain with generator $\boldsymbol{Q}(\pi)$. Then $\alpha(\pi)$ is the solution of the linear system $\alpha(\pi)\boldsymbol{Q}(\pi) = 0$ with normalizing equation $\alpha(\pi)\underline{\boldsymbol{1}} = 1$. Let $\alpha_U(\pi)$ and $\alpha_D(\pi)$ be the parts of vector $\alpha(\pi)$

associated with subsets $S_U$ and $S_D$, respectively. Using (2), the partitioned form of the linear system is

$$\alpha_U(\pi)\boldsymbol{Q_U}(\pi) + \alpha_D(\pi)\boldsymbol{Q_{DU}} = 0$$
$$\alpha_U(\pi)\boldsymbol{Q_{UD}}(\pi) + \alpha_D(\pi)\boldsymbol{Q_D} = 0, \qquad (3)$$

from which we obtain a linear system for $\alpha_U$

$$\alpha_U(\pi)(\boldsymbol{Q_U}(\pi) - \boldsymbol{Q_{UD}}(\pi)\boldsymbol{Q_D}^{-1}\boldsymbol{Q_{DU}}) = \alpha_U(\pi)\boldsymbol{Q_c}(\pi) = 0, \qquad (4)$$

where

$$\boldsymbol{Q_c}(\pi) = \boldsymbol{Q_U}(\pi) + \boldsymbol{Q_{UD}}(\pi)(-\boldsymbol{Q_D})^{-1}\boldsymbol{Q_{DU}}. \qquad (5)$$

The Markov chain with state space $S_U$ and generator $\boldsymbol{Q_c}(\pi)$ is referred to as censored Markov chain. It is obtained from the original Markov chain by "switching off the clock when the Markov chain visits $S_D$ and switching on the clock when the Markov chain visits $S_U$" [13].

The censored Markov chain defines the stationary probability of the states in $S_U$ through (4) apart from a normalizing constant, because $\sum_{i \in S_u} \alpha_i(\pi) = \alpha_U(\pi)\mathbf{1}_U$ is not known based on (4). We can also express $\alpha_D(\pi)$ from (3) as

$$\alpha_D(\pi) = \alpha_U(\pi)\boldsymbol{Q_{UD}}(\pi)(-\boldsymbol{Q_D})^{-1}, \qquad (6)$$

from which

$$1 = \alpha(\pi)\mathbf{1} = \alpha_D(\pi)\mathbf{1} + \alpha_U(\pi)\mathbf{1} = \alpha_U(\pi)(\mathbf{1} + \boldsymbol{Q_{UD}}(\pi)(-\boldsymbol{Q_D})^{-1}\mathbf{1}). \qquad (7)$$

Using (6), we can rewrite (1) as

$$\begin{aligned}
\pi^* &= \arg\max_{\pi} \alpha(\pi)\boldsymbol{C}(\pi)\mathbf{1} \\
&= \arg\max_{\pi} \alpha_U(\pi)\boldsymbol{C_U}(\pi)\mathbf{1} + \alpha_D(\pi)\boldsymbol{C_D}(\pi)\mathbf{1} \\
&= \arg\max_{\pi} \alpha_U(\pi)\big(\boldsymbol{C_U}(\pi)\mathbf{1} + \boldsymbol{Q_{UD}}(\pi)(-\boldsymbol{Q_D})^{-1}\boldsymbol{C_D}\mathbf{1}\big)
\end{aligned} \qquad (8)$$

# 4  Parameters for the MDP compression method

In this section, we compute some measures of interest that are needed for the MDP compression method.

Let us consider a continuous time MDP with partition $S_U$ (states with decision) and $S_D$ (states with no decision). For $i \in S_U$, let $\rho_{S_U \setminus i}$ be the time to visit a state in $S_U$ different from $i$, that is

$$\rho_{S_U \setminus i} = min(t | X(t) \in S_U \setminus i).$$

Based on the decision dependent state partitioning we define the $\boldsymbol{P}(\pi)$ matrix and the $\tau(\pi)$ and $c(\pi)$ vectors by their elements as follows:

$$\boldsymbol{P}_{ij}(\pi) = Pr(X(\rho_{S_U \setminus i}) = j \mid X(0) = i), \tag{9}$$

$$\tau_i(\pi) = E[\rho_{S_U \setminus i} \mid X(0) = i], \tag{10}$$

$$c_i(\pi) = E[\int_{t=0}^{\rho_{S_U \setminus i}} C_{X(t)X(t)}dt \mid X(0) = i], \tag{11}$$

where $i, j \in S_U$ and $i \neq j$. That is, assuming policy $\pi$, $\boldsymbol{P}_{ij}(\pi)$ is the probability that the process starting from $i \in S_U$ first enters to $S_U \setminus i$ in state $j$, $\tau_i(\pi)$ is the expected time to the first visit in $S_U \setminus i$ and $c_i(\pi)$ is the expected reward accumulated until this visit. $\boldsymbol{P}(\pi)_{ii} = 0$ by definition. The following theorem expresses $\boldsymbol{P}(\pi)$, $\tau(\pi)$ and $c(\pi)$ based on the partitioned description of the MDP.

**Theorem 1.** $\boldsymbol{P}(\pi)$, $\tau(\pi)$, and $c(\pi)$ can be obtained as

$$\boldsymbol{P}(\pi) = (-\boldsymbol{diagm}\langle\boldsymbol{Q_c}(\pi)\rangle)^{-1}(\boldsymbol{Q_c}(\pi) - \boldsymbol{diagm}\langle\boldsymbol{Q_c}(\pi)\rangle) \tag{12}$$

$$\tau(\pi) = (-\boldsymbol{diagm}\langle\boldsymbol{Q_c}(\pi)\rangle)^{-1}(\boldsymbol{1} + \boldsymbol{Q_{UD}}(\pi)\boldsymbol{A1}), \tag{13}$$

$$c(\pi) = (-\boldsymbol{diagm}\langle\boldsymbol{Q_c}(\pi)\rangle)^{-1}(\boldsymbol{C_U}(\pi)\boldsymbol{1} + \boldsymbol{Q_{UD}}(\pi)\boldsymbol{M1}), \tag{14}$$

where $\boldsymbol{Q_c}(\pi) = \boldsymbol{Q_U}(\pi) + \boldsymbol{Q_{UD}}(\pi)(-\boldsymbol{Q_D})^{-1}\boldsymbol{Q_{DU}}$, $\boldsymbol{A} = (-\boldsymbol{Q_D})^{-2}\boldsymbol{Q_{DU}}$, $\boldsymbol{M} = (-\boldsymbol{Q_D})^{-1}\boldsymbol{C_D}(-\boldsymbol{Q_D})^{-1}\boldsymbol{Q_{DU}}$ and $\boldsymbol{diagm}\langle\rangle$ is the operator that creates a diagonal matrix from an input matrix by setting all its non-diagonal elements to zero.

*Proof.* Although the formulas for $\boldsymbol{P}(\pi)$ and $\tau(\pi)$ can be derived alternatively by an easier approach, we will use a unified approach, for all three measures. Let $\rho_{S_U}$ be the first time when the process visits $S_U$. We define matrix $\boldsymbol{G}(t)$ such that for $i \in S_D$ and $j \in S_U$, $\boldsymbol{G}_{ij}(t) = Pr(X(\rho_{S_U}) = j, \rho_{S_U} < t | X(0) = i)$ and matrix $\boldsymbol{g}(t)$ as $\boldsymbol{g}(t) = \frac{d}{dt}\boldsymbol{G}(t)$. That is, $\boldsymbol{G}_{ij}(t)$ is the probability that the process starting from state $i \in S_D$ will visit $S_U$ before time $t$ and the first visit will be to state $j \in S_U$. We can express $\boldsymbol{g}_{ij}(t)$ based on the first state transition as

$$\boldsymbol{g}_{ij}(t) = -\boldsymbol{Q_{D_{ii}}}e^{\boldsymbol{Q_{D_{ii}}}t}\frac{\boldsymbol{Q_{DU_{ij}}}}{-\boldsymbol{Q_{D_{ii}}}} + \int_{\tau=0}^t -\boldsymbol{Q_{D_{ii}}}e^{\boldsymbol{Q_{D_{ii}}}\tau}\sum_{k,k\neq i}\frac{\boldsymbol{Q_{D_{ik}}}}{-\boldsymbol{Q_{D_{ii}}}}\boldsymbol{g}_{kj}(t-\tau)d\tau. \tag{15}$$

Here $-\boldsymbol{Q_{D_{ii}}}e^{\boldsymbol{Q_{D_{ii}}}\tau}$ corresponds to the density that the first state transition happens at time $\tau$, $\frac{\boldsymbol{Q_{DU_{ij}}}}{-\boldsymbol{Q_{D_{ii}}}}$ is the probability that the process goes directly to state $j$ at

the first transition and $\sum_{k,k\neq i} \frac{\boldsymbol{Q_{D_{ik}}}}{-\boldsymbol{Q_{D_{ii}}}}\boldsymbol{g}_{kj}(t-\tau)$ is the probability density that the process goes to some other state in $Q_D$ and it enters $S_U$ in state $j$ at time $t - \tau$.

For the Laplace transform of $\boldsymbol{g}_{ij}(t)$, $\boldsymbol{g^*}_{ij}(s) = \int_t e^{-st}\boldsymbol{g_{ij}}(t)dt$ we get

$$\boldsymbol{g^*}_{ij}(s) = \frac{-\boldsymbol{Q_{D_{ii}}}}{s - \boldsymbol{Q_{D_{ii}}}}\left(\frac{\boldsymbol{Q_{DU_{ij}}}}{-\boldsymbol{Q_{D_{ii}}}} + \sum_{k \in S_D, k \neq i} \frac{\boldsymbol{Q_{D_{ik}}}}{-\boldsymbol{Q_{D_{ii}}}}\boldsymbol{g^*}_{kj}(s)\right).$$

By multiplying both sides by $s - \boldsymbol{Q_{D_{ii}}}$ and adding $\boldsymbol{Q_{D_{ii}}}\boldsymbol{g^*}_{ii}(s)$ we obtain

$$s\boldsymbol{g^*}_{ij}(s) = \boldsymbol{Q_{DU_{ij}}} + \sum_{k \in S_D}\boldsymbol{Q_{D_{ik}}}\boldsymbol{g^*}_{kj}(s),$$

which can be written in matrix form as

$$s\boldsymbol{g^*}(s) = \boldsymbol{Q_{DU}} + \boldsymbol{Q_D}\boldsymbol{g^*}(s), \tag{16}$$

from which

$$\boldsymbol{g^*}(s) = (s\boldsymbol{I} - \boldsymbol{Q_D})^{-1}\boldsymbol{Q_{DU}}. \tag{17}$$

Since $\boldsymbol{Q_{DU}}$ and $\boldsymbol{Q_D}$ are policy independent, $\boldsymbol{g^*}(s)$ is policy independent as well. We define $\boldsymbol{G}_{ij} = \int_{t=0}^{\infty}\boldsymbol{g}_{ij}(t)dt = Pr(X(\rho_{S_U}) = j \mid X(0) = i \in S_D)$, which is the probability that the process starting from state $i \in S_D$ enters $S_U$ in state $j$. From $\boldsymbol{G} = \int_{t=0}^{\infty}\boldsymbol{g}(t)dt = \lim_{s\to 0}\boldsymbol{g^*}(s)$ and (17) we have

$$\boldsymbol{G} = (-\boldsymbol{Q_D})^{-1}\boldsymbol{Q_{DU}}. \tag{18}$$

We note here that, since the generator of the MDP is stationary, and $\boldsymbol{Q_D}$ is transient, that is, $Pr(\rho_{S_U} < \infty \mid X(0) = i) = 1, \forall i \in S_D$, thus $\boldsymbol{G}\underline{\boldsymbol{1}} = \underline{\boldsymbol{1}}$.

From the moment generating property of the Laplace transform [7] we also have

$$\boldsymbol{A}_{ij} \overset{def}{=} \mathrm{E}[\rho_{S_U}I_{\{\rho_{S_U}<\infty,X(\rho_{S_U})=j\}}|X(0) = i] = -\frac{d}{ds}\boldsymbol{g^*}_{ij}(s)\big|_{s=0}.$$

By differentiating (16) according to $s$ in $s = 0$ we obtain $\boldsymbol{G} = -\boldsymbol{Q_D}\boldsymbol{A}$, from which

$$\boldsymbol{A} = (-\boldsymbol{Q_D})^{-1}\boldsymbol{G} = (-\boldsymbol{Q_D})^{-2}\boldsymbol{Q_{DU}}. \tag{19}$$

Similar to $\boldsymbol{G}_{ij}(t)$ we define $\boldsymbol{K}_{ij}(r)$ as

$$\boldsymbol{K}_{ij}(r) = Pr\left(X(\rho_{S_U}) = j, \int_{t=0}^{\rho_{S_U}}\boldsymbol{C_{DX(t),X(t)}}dt < r|X(0) = i \in S_D\right), \tag{20}$$

that is, $\boldsymbol{K}_{ij}(r)$ is the probability that the process starting from state $i$ will visit $S_U$ before reward $r$ is accumulated, and the first visit will be to state $j \in S_U$. In

a similar fashion to $g_{ij}(t)$ we express $k_{ij}(r) = \frac{d}{dr}K_{ij}(r)$ based on the first state transition as

$$k_{ij}(r) = \frac{-Q_{D_{ii}}}{C_{D_{ii}}} e^{\frac{Q_{D_{ii}}}{C_{D_{ii}}}r} \frac{Q_{DU_{ij}}}{-Q_{D_{ii}}} \tag{21}$$
$$+ \int_{u=0}^{r} \frac{-Q_{D_{ii}}}{C_{D_{ii}}} e^{\frac{Q_{D_{ii}}}{C_{D_{ii}}}u} \sum_{k \in S_D, k \neq i} \frac{Q_{D_{ik}}}{-Q_{D_{ii}}} k_{kj}(r-u)du.$$

Compared to (15), the difference is that instead of the elapsed time we consider the reward accumulated up to the first state transition, which is exponentially distributed with rate $\frac{-Q_{D_{ii}}}{C_{D_{ii}}}$. Following the same steps as in the case of $g_{ij}(t)$ we get that the Laplace transform of $k(r)$, $k^*(s)$, satisfies

$$sC_D k^*(s) = Q_{DU} + Q_D k^*(s). \tag{22}$$

Just like $g^*(s)$, $k^*(s)$ is also independent of the actual policy, since it depends on the process behaviour during a visit in $S_D$. From (22) we have

$$K \stackrel{def}{=} \lim_{r \to \infty} K(r) = \int_{r=0}^{\infty} k(r)dr = \lim_{s \to 0} k^*(s) = (-Q_D)^{-1}Q_{DU},$$

where
$$K_{ij} = Pr(X(\rho_{S_U}) = j \,|\, X(0) = i \in S_D) = G_{ij}.$$

Similar to $A$, we introduce

$$M_{ij} \stackrel{def}{=} E\left[ I_{\{X(\gamma_U)=j\}} \int_{t=0}^{\rho_U} C_{DX(t),X(t)}dt \,\Big|\, X(0) = i \in S_D \right] = -\frac{d}{ds}k^*_{ij}(s)\Big|_{s=0},$$

which can be obtained from the derivative of (22) in $s = 0$, from which

$$C_D K = -Q_D M,$$

from which
$$M = (-Q_D)^{-1}C_D(-Q_D)^{-1}Q_{DU}. \tag{23}$$

Having the $G$, $A$, $M$ matrices, describing the behaviour of the process in $S_D$, we are ready to express $P(\pi)$, $\tau(\pi)$ and $c(\pi)$. The required derivations follow the same pattern. In each cases the formulas will be broken up into three terms according to the following three cases:

- Case 1: the process moves to state $j \in S_U \setminus i$ during the first state transition.

- Case 2: the process first moves to $k \in S_D$, spends some time in $S_D$, then enters $S_U$ in state $j \neq i$.

- Case 3: the process first moves to $k \in S_D$, spends some time in $S_D$, then enters $S_U$ in state $i$. This case adds a recursive term to the formulas.

We start with $\boldsymbol{P}_{ij}(\pi)$ for $i \neq j$, which gives the probability of entering set $S_U \setminus i$ in state $j \in S_U$ when the process starts in $i \in S_U$:

$$\boldsymbol{P}_{ij}(\pi) = Pr(\rho_{S_U \setminus i} < \infty, X(\rho_{S_U \setminus i}) = j \mid X(0) = i) =$$
$$\underbrace{\frac{\boldsymbol{Q}_{U\,ij}(\pi)}{-\boldsymbol{Q}_{U\,ii}(\pi)}}_{Case\ 1} + \underbrace{\sum_{k \in S_D} \frac{\boldsymbol{Q}_{UD\,ik}(\pi)}{-\boldsymbol{Q}_{U\,ii}(\pi)} \boldsymbol{G}_{kj}}_{Case\ 2} + \underbrace{\sum_{k \in S_D} \frac{\boldsymbol{Q}_{UD\,ik}(\pi)}{-\boldsymbol{Q}_{U\,ii}(\pi)} \boldsymbol{G}_{ki} \boldsymbol{P}_{ij}(\pi)}_{Case\ 3},$$

from which

$$-\left( \boldsymbol{Q}_{U\,ii}(\pi) + \sum_{k \in S_D} \boldsymbol{Q}_{UD\,ik}(\pi)\boldsymbol{G}_{ki} \right) \boldsymbol{P}_{ij}(\pi) = \boldsymbol{Q}_{U\,ij}(\pi) + \sum_{k \in S_D} \boldsymbol{Q}_{UD\,ik}(\pi)\boldsymbol{G}_{kj}.$$
$$(24)$$

As defined before let **diagm**$\langle\rangle$ be the operator that creates a diagonal matrix from an input matrix such that all non-diagonal elements of the original matrix are set to zero. Using this notation we can write (24) in matrix form as

$$\boldsymbol{P}(\pi) = (-\mathbf{diagm}\langle \boldsymbol{Q}_U(\pi) + \boldsymbol{Q}_{UD}(\pi)\boldsymbol{G}\rangle)^{-1}$$
$$\cdot (\boldsymbol{Q}_U(\pi) + \boldsymbol{Q}_{UD}(\pi)\boldsymbol{G} - \mathbf{diagm}\langle \boldsymbol{Q}_U(\pi) + \boldsymbol{Q}_{UD}(\pi)\boldsymbol{G}\rangle) \quad (25)$$
$$= (-\mathbf{diagm}\langle \boldsymbol{Q}_c(\pi)\rangle)^{-1}(\boldsymbol{Q}_c(\pi) - \mathbf{diagm}\langle \boldsymbol{Q}_c(\pi)\rangle)$$

where $-\mathbf{diagm}\langle \boldsymbol{Q}_c(\pi)\rangle$ in the second term of the right side ensures that the diagonal of $\boldsymbol{P}(\pi)$ is equal to zero according to the definition of $\boldsymbol{P}_{ii}(\pi)$ and we used that $\boldsymbol{Q}_c(\pi) = \boldsymbol{Q}_U(\pi) + \boldsymbol{Q}_{UD}(\pi)(-\boldsymbol{Q}_D)^{-1}\boldsymbol{Q}_{DU} = \boldsymbol{Q}_U(\pi) + \boldsymbol{Q}_{UD}(\pi)\boldsymbol{G}$ (which comes from the definition of $\boldsymbol{Q}_c(\pi)$ in (5) and the definition of $\boldsymbol{G}$ in (19)).

The formula for $\tau_i(\pi)$, which describes the expected time before entering $S_U \setminus i$ when the process starts in $i \in S_U$, has a similar structure:

$$\tau_i(\pi) = E\left[\rho_{U\setminus i}|X(0)=i, \exists \rho_{S_U\setminus i}\right] = \underbrace{\sum_{j\in S_U\setminus i} \frac{\boldsymbol{Q_U}_{ij}(\pi)}{-\boldsymbol{Q_U}_{ii}(\pi)} \frac{1}{-\boldsymbol{Q_U}_{ii}(\pi)}}_{Case\ 1} +$$

$$\underbrace{\sum_{j\in S_U\setminus i}\sum_{k\in S_D} \frac{\boldsymbol{Q_{UD}}_{ik}(\pi)}{-\boldsymbol{Q_U}_{ii}(\pi)}\left(\boldsymbol{G}_{kj}\frac{1}{-\boldsymbol{Q_U}_{ii}(\pi)} + \boldsymbol{A}_{kj}\right)}_{Case\ 2} +$$

$$\underbrace{\sum_{k\in S_D} \frac{\boldsymbol{Q_{UD}}_{ik}(\pi)}{-\boldsymbol{Q_U}_{ii}(\pi)}\left(\boldsymbol{G}_{ki}\frac{1}{-\boldsymbol{Q_U}_{ii}(\pi)} + \boldsymbol{A}_{ki} + \boldsymbol{G}_{ki}\tau_i(\pi)\right)}_{Case\ 3}$$

$$= \frac{1}{-\boldsymbol{Q_U}_{ii}(\pi)} + \sum_{j\in S_U}\sum_{k\in S_D} \frac{\boldsymbol{Q_{UD}}_{ik}(\pi)}{-\boldsymbol{Q_U}_{ii}(\pi)}\boldsymbol{A}_{kj} + \sum_{k\in S_D} \frac{\boldsymbol{Q_{UD}}_{ik}(\pi)}{-\boldsymbol{Q_U}_{ii}(\pi)}\boldsymbol{G}_{ki}\tau_i(\pi), \quad (26)$$

where we used that $\boldsymbol{Q_D}$ is transient, thus $\sum_{j\in S_U}\boldsymbol{G}_{kj}=1, \forall k\in S_D$. The first term (Case 1) is the mean time until the first transition multiplied by the probability of Case 1. In the second term (Case 2) $\sum_{j\in S_U\setminus i}\sum_{k\in S_D} \frac{\boldsymbol{Q_{UD}}_{ik}(\pi)}{-\boldsymbol{Q_U}_{ii}(\pi)}\boldsymbol{G}_{kj}\frac{1}{-\boldsymbol{Q_U}_{ii}(\pi)}$ is the time until the first transition multiplied by the probability of Case 2 and summed for all $j \in S_U \setminus i$ and $\sum_{j\in S_U\setminus i}\sum_{k\in S_D} \frac{\boldsymbol{Q_{UD}}_{ik}(\pi)}{-\boldsymbol{Q_U}_{ii}(\pi)}\boldsymbol{A}_{kj}$ is the remaining time until visiting state $j$ multiplied by the probability of Case 2 and summed for all $j \in S_U \setminus i$. We note that in this part multiplication by $\boldsymbol{G}_{kj}$ is not necessary, because $\boldsymbol{A}_{kj}$ already contains the probability $\boldsymbol{G}_{kj}$. The third term contains analogous components to the second term, but for Case 3, when the process enters $S_U$ in state $i$, we also have the added $\sum_{k\in S_D} \frac{\boldsymbol{Q_{UD}}_{ik}(\pi)}{-\boldsymbol{Q_U}_{ii}(\pi)}\boldsymbol{G}_{ki}\tau_i(\pi)$, because if the process returns to state $i$, an additional $\tau_i(\pi)$ time is needed to visit some state in $S_U \setminus i$. We can express $\tau_i(\pi)$ from (26) as

$$-\left(\boldsymbol{Q_U}_{ii}(\pi) + \sum_{k\in S_D}\boldsymbol{Q_{UD}}_{ik}(\pi)\boldsymbol{G}_{ki}\right)\tau_i(\pi) = 1 + \sum_{j\in S_U}\sum_{k\in S_D}\boldsymbol{Q_{UD}}_{ik}(\pi)\boldsymbol{A}_{kj}.$$
$$(27)$$

Using $\boldsymbol{G} = (-\boldsymbol{Q_D})^{-1}\boldsymbol{Q_{DU}}$ from (18), the matrix form of (27) is

$$\tau(\pi) = (-\mathbf{diagm}\langle\boldsymbol{Q_U}(\pi) + \boldsymbol{Q_{UD}}(\pi)\boldsymbol{G}\rangle)^{-1}(\mathbf{1} + \boldsymbol{Q_{UD}}(\pi)\boldsymbol{A1}), \qquad (28)$$

which is identical with (13), the corresponding equation of Theorem 1.

Finally we have a very similar expression for the elements of the $c(\pi)$ vector:

$$c_i(\pi) = E[\int_{t=0}^{\rho_{S_U \setminus i}} C_{X(t)X(t)} \mid X(0) = i] = \underbrace{\sum_{j \in S_U \setminus i} \frac{Q_{Uij}(\pi)}{-Q_{Uii}(\pi)} \frac{C_{Uii}(\pi)}{-Q_{Uii}(\pi)}}_{Case\ 1} +$$

$$\underbrace{\sum_{j \in S_U \setminus i} \sum_{k \in S_D} \frac{Q_{UDik}(\pi)}{-Q_{Uii}(\pi)} \left( G_{kj} \frac{C_{Uii}(\pi)}{-Q_{Uii}(\pi)} + M_{kj} \right)}_{Case\ 2} +$$

$$\underbrace{\sum_{k \in S_D} \frac{Q_{UDik}(\pi)}{-Q_{Uii}(\pi)} \left( G_{ki} \frac{C_{Uii}(\pi)}{-Q_{Uii}(\pi)} + M_{ki} + G_{ki}c_i(\pi) \right)}_{Case\ 3}$$

$$= \frac{C_{Uii}(\pi)}{-Q_{Uii}(\pi)} + \sum_{j \in S_U} \sum_{k \in S_D} \frac{Q_{UDik}(\pi)}{-Q_{Uii}(\pi)} M_{kj} + \sum_{k \in S_D} \frac{Q_{UDik}(\pi)}{-Q_{Uii}(\pi)} G_{ki}c_i(\pi). \quad (29)$$

This expression follows the same logic as (26), the only difference is that instead of times we have rewards accumulated over those times, thus the $\frac{1}{-Q_{Uii}}(\pi)$ terms in (26) are replaced by $\frac{C_{Uii}(\pi)}{-Q_{Uii}}(\pi)$, $A$ is changed to $M$ and $\tau_i(\pi)$ is changed to $c_i(\pi)$. We can express $c_i(\pi)$ from (29) as

$$-\left( Q_{Uii}(\pi) + \sum_{k \in S_D} Q_{UDik}(\pi)G_{ki} \right) c_i(\pi) = C_{Uii}(\pi) + \sum_{j \in S_U} \sum_{k \in S_D} Q_{UDik}(\pi)M_{kj}.$$

Using the same substitutions as in (12) and (13), the matrix form of this equation is identical to (14), the final equation of Theorem 1. $\qquad \square$

# 5 Compression of partitioned MDPs

In this section, we present the main contribution of the paper, the compressed representation of partitioned MDPs. We discuss the idea behind the compression, and provide analytical proof for the equivalence of the original and the compressed forms.

**Theorem 2.** *Let $(S, A, Q(\pi), C(\pi))$ be an MDP with irreducible $Q(\pi)$, where the generator and reward-rate matrices can be partitioned according to (2). The MDP defined by $(S', A', Q'(\pi), C'(\pi))$ and the one defined by $(S, A, Q(\pi), C(\pi))$ have the same optimal policy, where $S' = S_U$, $A' = A$,*

$$Q'_{ij}(\pi) = \begin{cases} -\frac{1}{\tau_i(\pi)}, & \text{if } i = j, \\ \frac{P_{ij}(\pi)}{\tau_i(\pi)}, & \text{otherwise,} \end{cases} \quad \text{and} \quad C'_{ij}(\pi) = \begin{cases} \frac{c_i(\pi)}{\tau_i(\pi)}, & \text{if } i = j, \\ 0, & \text{otherwise.} \end{cases}$$

**Remark.** *The core idea of the compression is the following. When the process exits a state $i \in S_U$, it can take multiple trajectories. It can either directly transition to another state in $S_U \setminus \{i\}$, or it can go through a number of transitions in $S_D \cup \{i\}$ before transitioning to a state in $S_U \setminus \{i\}$. When the decision is made in state $i$, the distribution of the next state reached in $S_U \setminus \{i\}$, the distribution of the time until reaching this state, and the distribution of the accumulated reward until reaching this state can be calculated (and no further knowledge of decisions in other states is needed). The theorem states that, to optimise the long-term expected reward rate, we can change the distribution of the time needed to reach $S_U \setminus \{i\}$ and the associated accumulated reward to exponentially distributed variables; as long as their expected values match with the expected value of the original distributions, the optimal policy will not change.*

*Proof.* Let us denote by $\mathbf{diagv}\langle v \rangle$ the diagonal matrix created from the elements of vector $v$, let $\boldsymbol{S}(\pi) = \mathbf{diagv}\langle \tau(\pi) \rangle$ and $\boldsymbol{Z}(\pi) = -\mathbf{diagm}\langle Q_c(\pi) \rangle$. Then, using the formula for $\boldsymbol{P}(\pi)$ from (12) we can write

$$
\begin{aligned}
\boldsymbol{Q}'(\pi) &= \mathbf{diagv}\langle \tau(\pi) \rangle^{-1}(-\boldsymbol{I} + \boldsymbol{P}(\pi)) \\
&= \boldsymbol{S}^{-1}(\pi)(-\boldsymbol{I} + (-\mathbf{diagm}\langle \boldsymbol{Q_c}(\pi) \rangle)^{-1}(\boldsymbol{Q_c}(\pi) - \mathbf{diagm}\langle \boldsymbol{Q_c}(\pi) \rangle)) \quad (30) \\
&= \boldsymbol{S}^{-1}(\pi)\boldsymbol{Z}^{-1}(\pi)\boldsymbol{Q_c}(\pi).
\end{aligned}
$$

and using the formula for $c(\pi)$ from (14) we have

$$
\begin{aligned}
\boldsymbol{C}'(\pi) &= \mathbf{diagv}\langle \tau(\pi) \rangle^{-1}\mathbf{diagm}\langle c(\pi) \rangle \\
&= \boldsymbol{S}^{-1}(\pi)(-\mathbf{diagm}\langle \boldsymbol{Q_c}(\pi) \rangle)^{-1}(\boldsymbol{C_U}(\pi)\mathbf{1} + \boldsymbol{Q_{UD}}(\pi)M\mathbf{1}) \\
&= \boldsymbol{S}^{-1}(\pi)\boldsymbol{Z}^{-1}(\pi)(\boldsymbol{C_U}(\pi)\mathbf{1} + \boldsymbol{Q_{UD}}(\pi)(-\boldsymbol{Q_D})^{-1}\boldsymbol{C_D}(-\boldsymbol{Q_D})^{-1}\boldsymbol{Q_{DU}}\mathbf{1}) \\
&= \boldsymbol{S}^{-1}(\pi)\boldsymbol{Z}^{-1}(\pi)(\boldsymbol{C_U}(\pi)\mathbf{1} + \boldsymbol{Q_{UD}}(\pi)(-\boldsymbol{Q_D})^{-1}\boldsymbol{C_D}\mathbf{1}),
\end{aligned}
$$
$$(31)$$

where we used that $(-\boldsymbol{Q_D})^{-1}\boldsymbol{Q_{DU}}\mathbf{1} = \mathbf{1}$ for stationary MDPs. The $\alpha'(\pi)$ stationary probability vector of the compressed MDP has to satisfy the

$$
\begin{aligned}
\alpha'(\pi)\boldsymbol{Q}'(\pi) &= 0 \\
\alpha'(\pi)\mathbf{1} &= 1
\end{aligned}
$$

system of equations. The first equation can be transformed as

$$
0 = \alpha'(\pi)\boldsymbol{Q}'(\pi) = \alpha'(\pi)\boldsymbol{S}^{-1}(\pi)\boldsymbol{Z}^{-1}(\pi)\boldsymbol{Q_c}(\pi) = \alpha''(\pi)\boldsymbol{Q_c}(\pi), \quad (32)
$$

where $\alpha''(\pi) = \alpha'(\pi)\boldsymbol{S}^{-1}(\pi)\boldsymbol{Z}^{-1}(\pi)$. The second equation can be transformed as

$$
\begin{aligned}
1 &= \alpha'(\pi)\underline{\mathbf{1}} \\
&= \alpha'(\pi)\boldsymbol{S}^{-1}(\pi)\boldsymbol{Z}^{-1}(\pi)\boldsymbol{S}(\pi)\boldsymbol{Z}(\pi)\underline{\mathbf{1}} \\
&= \alpha''(\pi)(-\mathbf{diagm}\langle Q_c(\pi)\rangle)\mathbf{diagv}\langle\tau(\pi)\rangle\underline{\mathbf{1}} \\
&= \alpha''(\pi)(-\mathbf{diagm}\langle Q_c(\pi)\rangle)(-\mathbf{diagm}\langle Q_c\rangle(\pi))^{-1}(\underline{\mathbf{1}} + \boldsymbol{Q_{UD}}(\pi)\boldsymbol{A}\underline{\mathbf{1}}) \\
&= \alpha''(\pi)(\underline{\mathbf{1}} + \boldsymbol{Q_{UD}}(\pi)\boldsymbol{A}\underline{\mathbf{1}}).
\end{aligned}
$$

Using the formula for $\boldsymbol{A}$ from (18) and using that $(-\boldsymbol{Q_D})^{-1}\boldsymbol{Q_{DU}}\underline{\mathbf{1}} = \underline{\mathbf{1}}$ we can further transform the expression as

$$
1 = \alpha''(\pi)(\underline{\mathbf{1}} + \boldsymbol{Q_{UD}}(\pi)(-\boldsymbol{Q_D})^{-1}\underline{\mathbf{1}}). \tag{33}
$$

Thus, from (32) and (33) we obtain the

$$
0 = \alpha''(\pi)\boldsymbol{Q_c}(\pi) \,, \ 1 = \alpha''(\pi)(\underline{\mathbf{1}} + \boldsymbol{Q_{UD}}(\pi)(-\boldsymbol{Q_D})^{-1}\underline{\mathbf{1}})
$$

system of linear equations for $\alpha''(\pi)$ that completely determine $\alpha''(\pi)$. However these are the same as the linear equations for $\alpha_U(\pi)$, (4) and (7), thus $\alpha''(\pi) = \alpha_U(\pi)$.

The mean reward rate of the compressed MDP can be given using (31) as

$$
\begin{aligned}
\alpha'(\pi)\boldsymbol{C}'(\pi)\underline{\mathbf{1}} &= \alpha'(\pi)\boldsymbol{S}^{-1}(\pi)\boldsymbol{Z}^{-1}(\pi)(\boldsymbol{C_U}(\pi)\underline{\mathbf{1}} + \boldsymbol{Q_{UD}}(\pi)(-\boldsymbol{Q_D})^{-1}\boldsymbol{C_D}\underline{\mathbf{1}}) \\
&= \alpha_U(\pi)(\boldsymbol{C_U}(\pi)\underline{\mathbf{1}} + \boldsymbol{Q_{UD}}(\pi)(-\boldsymbol{Q_D})^{-1}\boldsymbol{C_D}\underline{\mathbf{1}})
\end{aligned}
$$

where we used that $\alpha''(\pi) = \alpha_U(\pi) = \alpha'(\pi)\boldsymbol{S}^{-1}(\pi)\boldsymbol{Z}^{-1}(\pi)$. This, however, is the same as the mean reward rate for the original MDP, as can be seen from (8), thus the average reward rate of the original and the compressed MDP is the same for any given policy, thus their optimal policies are also the same. $\qquad\square$

# 6 Compression of partitioned MDPs with special structures

In the previous section, we presented the general formulas of the proposed MDP compression method. In this section, we discuss the application of the method in some special, practically important cases. The general formulas for the compression method rely on the calculation of the $\boldsymbol{G}$, $\boldsymbol{A}$, and $\boldsymbol{M}$ matrices. For MDPs with finite $S_D$ the calculation of these matrices can be done based on (18), (19) and (23), which are based on the computation of the inverse of $\boldsymbol{Q_D}$. If the $S_D$ subset is infinite, however, the calculation of these matrices is not trivial and it has

to rely on the structural regularity of the MDP. We discuss two important cases where the calculation of $G$, $A$, and $M$ matrices are possible, which are the cases when the structures of the MDP in $S_D$ are spatially homogeneous M/G/1-type and G/M/1-type.

During the analysis of matrix $G$ and $A$ we are going to utilize some known results of M/G/1-type and G/M/1-type processes, while the analysis presented for matrix $M$ was not discussed in the literature to the best of the authors' knowledge.

## 6.1 The M/G/1 type process

An MDP is of M/G/1 type (with the considered $\{S_U, S_D\}$ partitioning) if its generator matrix has the following block structure

$$
Q(\pi) = \left( \begin{array}{c|cccc}
\bar{L}(\pi) & \bar{F}_1(\pi) & \bar{F}_2(\pi) & \ldots & \\
\hline
B & L & F_1 & F_2 & \ldots \\
 & B & L & F_1 & F_2 & \ldots \\
 & & \ddots & \ddots & \ddots & \ddots
\end{array} \right). \tag{34}
$$

The MDP can be partitioned to levels according to the blocks so that each row in (34) corresponds to a separate level. In the following we assume that the reward rate matrix also has some level based regularity, thus

$$
C(\pi) = \left( \begin{array}{c|cccc}
C_1(\pi) & & & \\
\hline
 & C_2 & & \\
 & & C_3 & \\
 & & & \ddots
\end{array} \right), \tag{35}
$$

where $C_n = f(C_0, n)$ for $n = 2, 3, \ldots,$. We discuss some specific cases in Section 7.2.

In the following we denote the $i$th state of the $n$th level by $v_{n,i}$, where $i \in \Phi = \{1, 2, \ldots, \phi\}$ and $n \in \{1, 2, \ldots\}$. Using the level based partitioning the process cannot descend more than one level during a single state transition (i.e., a transition cannot happen from level $n$ to level $n - k$, $k \geq 2$), and any downward transition between neighbouring levels is described by the $B$ matrix. We define the levels such that level 1 is identical with $S_U$ and the rest of the levels are in $S_D$,

thus the $S_U$, $S_D$ based decomposition of the $\boldsymbol{Q}(\pi)$ and $\boldsymbol{C}(\pi)$ is

$$\boldsymbol{Q_U}(\pi) = \bar{\boldsymbol{L}}(\pi), \quad \boldsymbol{Q_{UD}}(\pi) = \begin{pmatrix} \bar{\boldsymbol{F}}_1(\pi) & \bar{\boldsymbol{F}}_2(\pi) & \dots \end{pmatrix},$$

$$\boldsymbol{Q_{DU}} = \begin{pmatrix} \boldsymbol{B} \\ 0 \\ \vdots \end{pmatrix}, \quad \boldsymbol{Q_D} = \begin{pmatrix} \boldsymbol{L} & \boldsymbol{F_1} & \boldsymbol{F_2} & \dots \\ \boldsymbol{B} & \boldsymbol{L} & \boldsymbol{F_1} & \boldsymbol{F_2} & \dots \\ & \ddots & \ddots & \ddots & \ddots \end{pmatrix},$$

$$\boldsymbol{C_U}(\pi) = \boldsymbol{C_1}(\pi), \quad \boldsymbol{C_D} = \begin{pmatrix} \boldsymbol{C_2} & & \\ & \boldsymbol{C_3} & \\ & & \ddots \end{pmatrix}.$$

To compute the $\boldsymbol{G}$, $\boldsymbol{A}$, and $\boldsymbol{M}$ matrices for the infinite $S_D$ we define the $\hat{\boldsymbol{G}}$, $\hat{\boldsymbol{A}}$, and $\hat{\boldsymbol{M}}_n$ matrices that are similar, however, instead of trajectories from $S_D$ to $S_U$ they describe trajectories from level $n$ to level $n-1$ inside $S_D$. Matrix $\hat{\boldsymbol{G}}$ is known to be the characteristic matrix of the M/G/1-type process [15] and is well discussed in the literature, unlike the analysis of matrix $\hat{\boldsymbol{M}}_n$.

**Theorem 3.** *For an MDP with M/G/1-type structure in $S_D$ according to (34) and (35) the parameters of the reduced MDP representation are*

$$\boldsymbol{P}(\pi) = (-\boldsymbol{diagm}\langle\boldsymbol{Q_c}(\pi)\rangle)^{-1}(\boldsymbol{Q_c}(\pi) - \boldsymbol{diagm}\langle\boldsymbol{Q_c}(\pi)\rangle), \tag{36}$$

$$\tau(\pi) = (-\boldsymbol{diagm}\langle\boldsymbol{Q_c}(\pi)\rangle)^{-1}\left(\underline{\boldsymbol{1}} + \sum_{i=1}^{\infty}\bar{\boldsymbol{F}}_i(\pi)\sum_{\ell=0}^{i}\hat{\boldsymbol{G}}^{\ell}\hat{\boldsymbol{A}}\hat{\boldsymbol{G}}^{i-\ell}\underline{\boldsymbol{1}}\right), \tag{37}$$

$$c(\pi) = (-\boldsymbol{diagm}\langle\boldsymbol{Q_c}(\pi)\rangle)^{-1}\left(\boldsymbol{C_1}(\pi)\underline{\boldsymbol{1}} + \sum_{i=1}^{\infty}\bar{\boldsymbol{F}}_i(\pi)\sum_{\ell=0}^{i}\hat{\boldsymbol{G}}^{\ell}\hat{\boldsymbol{M}}_{i+1-\ell}\hat{\boldsymbol{G}}^{i-\ell}\underline{\boldsymbol{1}}\right), \tag{38}$$

*where $\boldsymbol{Q_c}(\pi) = \bar{\boldsymbol{L}}(\pi) + \sum_{i=1}^{\infty}\bar{\boldsymbol{F}}_i(\pi)\hat{\boldsymbol{G}}^i$, $\hat{\boldsymbol{G}}$ is the minimal non-negative solution of*

$$0 = \boldsymbol{B} + \boldsymbol{L}\hat{\boldsymbol{G}} + \sum_{m=1}^{\infty}\boldsymbol{F_m}\hat{\boldsymbol{G}}^{m+1}, \tag{39}$$

*$\hat{\boldsymbol{A}}$ is the solution of the linear matrix equation*

$$\hat{\boldsymbol{G}} = \boldsymbol{L}\hat{\boldsymbol{A}} + \sum_{m=1}^{\infty}\boldsymbol{F_m}\sum_{\ell=0}^{m}\hat{\boldsymbol{G}}^{\ell}\hat{\boldsymbol{A}}\hat{\boldsymbol{G}}^{m-\ell}, \tag{40}$$

*and $\hat{\boldsymbol{M}}_n$ is the solution of*

$$\boldsymbol{C_n}\hat{\boldsymbol{G}} = \boldsymbol{L}\hat{\boldsymbol{M}}_n + \sum_{m=1}^{\infty}\boldsymbol{F_m}\sum_{\ell=0}^{m}\hat{\boldsymbol{G}}^{\ell}\hat{\boldsymbol{M}}_{n+m-\ell}\hat{\boldsymbol{G}}^{m-\ell}. \tag{41}$$

**Remark.** *Efficient numerical methods are available for the solution of* (39) *e.g., in [4]. The solution of* (40) *can be achieved e.g., with the use of the column stacking vec operator, for which* $vec(\boldsymbol{ABC}) = (\boldsymbol{C}^T \otimes \boldsymbol{A})vec(\boldsymbol{B})$. *Applying vec for* (40) *gives*

$$vec(\hat{\boldsymbol{G}}) = (\boldsymbol{I} \otimes \boldsymbol{L}) \, vec(\hat{\boldsymbol{A}}) + \sum_{m=1}^{\infty} \sum_{\ell=0}^{m} \left( \hat{\boldsymbol{G}}^{m-\ell T} \otimes \boldsymbol{F_m} \hat{\boldsymbol{G}}^{\ell} \right) vec(\hat{\boldsymbol{A}}),$$

*from which*

$$vec(\hat{\boldsymbol{A}}) = \left( (\boldsymbol{I} \otimes \boldsymbol{L}) + \sum_{m=1}^{\infty} \sum_{\ell=0}^{m} \left( \hat{\boldsymbol{G}}^{m-\ell T} \otimes \boldsymbol{F_m} \hat{\boldsymbol{G}}^{\ell} \right) \right)^{-1} vec(\hat{\boldsymbol{G}}).$$

*The solution of* (41) *is more difficult in general. It is discussed in Section 7.2 for some special* $\boldsymbol{C_n}$ *(and* $\boldsymbol{F_n}$*) series.*

*We also note that the expressions in the theorem can be further simplified based on the fact that the characteristic matrix of a positive recurrent M/G/1 type process is a stochastic matrix, that is* $\hat{\boldsymbol{G}}\underline{\boldsymbol{1}} = \underline{\boldsymbol{1}}$. *When applying this simplification it is enough to compute vectors* $\hat{\boldsymbol{A}}\underline{\boldsymbol{1}}$ *and* $\hat{\boldsymbol{M_i}}\underline{\boldsymbol{1}}$ *instead of matrices* $\hat{\boldsymbol{A}}$ *and* $\hat{\boldsymbol{M_i}}$.

*Proof.* Similar to the proof of Theorem 1 we apply a unified approach for the analysis of all required measures. For the analysis of the level process in $S_D$, we define $\hat{\boldsymbol{G}}_{ij}(t) = Pr(X(\rho_{n-1}) = v_{n-1,j}, \rho_{n-1} < t | X(0) = v_{n,i})$, where $\rho_{n-1}$ is the time of the first visit to level $n-1$, i.e., $\hat{\boldsymbol{G}}_{ij}(t)$ is the probability that the process, starting from state $i$ of level $n$ reaches level $n-1$ before time $t$ and the first visit is to state $j$ on this level. We also define the multi level version of this measure, $\hat{\boldsymbol{G}}_{\boldsymbol{m}ij}$, that describe trajectories from level $n$ to level $n-m$,

$$\hat{\boldsymbol{G}}_{\boldsymbol{m}ij}(t) = Pr(X(\rho_{n-m}) = v_{n-1,j}, \rho_{n-m} < t | X(0) = v_{n,i}).$$

Furthermore, $\hat{\boldsymbol{g}}(t) = \frac{d}{dt}\hat{\boldsymbol{G}}(t)$ and $\hat{\boldsymbol{g}}_{\boldsymbol{m}}(t) = \frac{d}{dt}\hat{\boldsymbol{G}}_{\boldsymbol{m}}(t)$. By definition, we have $\hat{\boldsymbol{g}}_{\boldsymbol{1}}(t) = \hat{\boldsymbol{g}}(t)$ and by the fact that the first visit from level $n$ to level $n-m$ can be decomposed into the first visit from level $n$ to level $n-1$ and then the first visit from level $n-1$ to level $n-m$ we also have $\hat{\boldsymbol{g}}_{\boldsymbol{m}}(t) = (\hat{\boldsymbol{g}} * \hat{\boldsymbol{g}}_{\boldsymbol{m-1}})(t), \forall m \geq 2$, where $*$ is the convolution operator, i.e., $(\boldsymbol{a} * \boldsymbol{b})(t) = \int_{\tau=0}^{t} \boldsymbol{a}(\tau)\boldsymbol{b}(t-\tau)d\tau$. Similar to $\boldsymbol{g}_{ij}(t)$, we can express $\hat{\boldsymbol{g}}_{ij}(t)$ based on the first transition from state $v_{n,i}$ as

$$\hat{\boldsymbol{g}}_{ij}(t) = -\boldsymbol{L}_{ii}\mathrm{e}^{\boldsymbol{L}_{ii}t}\frac{\boldsymbol{B}_{ij}}{-\boldsymbol{L}_{ii}} +$$

$$\int_{\tau=0}^{t} -\boldsymbol{L}_{ii}\mathrm{e}^{\boldsymbol{L}_{ii}\tau} \left( \sum_{\substack{k \in \Phi, \\ k \neq i}} \frac{\boldsymbol{L}_{ik}}{-\boldsymbol{L}_{ii}} \hat{\boldsymbol{g}}_{kj}(t-\tau) + \sum_{m=1}^{\infty} \sum_{\substack{k \in \Phi, \\ k \neq i}} \frac{\boldsymbol{F_m}_{ik}}{-\boldsymbol{L}_{ii}} \hat{\boldsymbol{g}}_{\boldsymbol{m+1}kj}(t-\tau) \right) d\tau.$$

$$(42)$$

The Laplace transform of (42) gives

$$\hat{\boldsymbol{g}}^*{}_{ij}(s) = \frac{-\boldsymbol{L}_{ii}}{s - \boldsymbol{L}_{ii}}\left(\frac{\boldsymbol{B}_{ij}}{-\boldsymbol{L}_{ii}} + \sum_{\substack{k\in\Phi,\\k\neq i}}\frac{\boldsymbol{L}_{ik}}{-\boldsymbol{L}_{ii}}\hat{\boldsymbol{g}}^*{}_{kj}(s) + \sum_{m=1}^{\infty}\sum_{k\in\Phi}\frac{\boldsymbol{F}_{mik}}{-\boldsymbol{L}_{ii}}\hat{\boldsymbol{g}}^*_{m+1\,kj}(s)\right).$$

Multiplying both sides by $s - \boldsymbol{L}_{ii}$ and adding $\boldsymbol{L}_{ii}\hat{\boldsymbol{g}}^*{}_{ii}(s)$ we obtain

$$s\hat{\boldsymbol{g}}^*{}_{ij}(s) = \boldsymbol{B}_{ij} + \sum_{k\in\Phi}\frac{\boldsymbol{L}_{ik}}{-\boldsymbol{L}_{ii}}\hat{\boldsymbol{g}}^*{}_{kj}(s) + \sum_{m=1}^{\infty}\sum_{k\in\Phi}\frac{\boldsymbol{F}_{mik}}{-\boldsymbol{L}_{ii}}\hat{\boldsymbol{g}}^*_{m+1\,kj}(s),$$

which can be written in matrix form, using $\hat{\boldsymbol{g}}^*_{\boldsymbol{m}}(s) = \hat{\boldsymbol{g}}^*(s)^m$, as

$$s\hat{\boldsymbol{g}}^*(s) = \boldsymbol{B} + \boldsymbol{L}\hat{\boldsymbol{g}}^*(s) + \sum_{m=1}^{\infty}\boldsymbol{F_m}\hat{\boldsymbol{g}}^*(s)^{m+1}. \qquad (43)$$

Similar to the case of $\boldsymbol{g}^*(s)$, using the final value theorem we have that

$$\hat{\boldsymbol{G}}_{ij} \overset{def}{=} Pr(X(\rho_{n-1}) = v_{n-1,j}|X(0) = v_{n,i})$$

$$= \lim_{t\to\infty}G(t) = \int_{t=0}^{\infty}\hat{\boldsymbol{g}}_{ij}(t)dt = \lim_{s\to 0}\hat{\boldsymbol{g}}^*{}_{ij}(s),$$

i.e., $\hat{\boldsymbol{G}}_{ij}$ is the probability that the process, starting from state $i$ of level $n$ reaches level $n-1$ ($n > 2$) in state $j$. Substituting $s = 0$ into (43) we get

$$0 = \boldsymbol{B} + \boldsymbol{L}\hat{\boldsymbol{G}} + \sum_{m=1}^{\infty}\boldsymbol{F_m}\hat{\boldsymbol{G}}^{m+1}, \qquad (44)$$

Which is the well-known matrix equation for computing the characteristic matrix of an M/G/1 type process [15]. Similar to matrix $\boldsymbol{A}$ we define

$$\hat{\boldsymbol{A}}_{ij} \overset{def}{=} \mathrm{E}[\rho_{n-1}I_{\{X(\rho_{n-1})=j\}}|X(0) = i] = -\frac{d}{ds}\hat{\boldsymbol{g}}^*{}_{ij}(s)\big|_{s=0},$$

which we obtain from the moment generating property of the Laplace transform.

Taking the derivative of (43) according to $s$ in $s = 0$ we get

$$\hat{\boldsymbol{G}} = \boldsymbol{L}\hat{\boldsymbol{A}} + \sum_{m=1}^{\infty}\boldsymbol{F_m}\sum_{\ell=0}^{m}\hat{\boldsymbol{G}}^{\ell}\hat{\boldsymbol{A}}\hat{\boldsymbol{G}}^{m-\ell}, \qquad (45)$$

which is the linear equation for computing $\hat{\boldsymbol{A}}$ based on $\hat{\boldsymbol{G}}$.

To obtain the matrix of mean accumulated rewards till the first visit to level $n-m$, starting from level $n$, similar to $\boldsymbol{K}(r)$ we define $\hat{\boldsymbol{K}}_{n,m}(r)$ by its $i,j$ element as

$$\hat{\boldsymbol{K}}_{n,m_{ij}}(r) = Pr\left( X(\rho_{n-m}) = v_{n-m,j}, \int_{t=0}^{\rho_{n-m}} \boldsymbol{C}_{X(t),X(t)} dt < r | X(0) = v_{n,i} \right),$$
(46)

that is, $\hat{\boldsymbol{K}}_{n,m_{ij}}(r)$ is the probability that the process starting from state $i$ of level $n$ will visit level $n-m$ before reward $r$ is accumulated and the first visit to level $n-m$ will be in state $j$. We also define $\hat{\boldsymbol{k}}_{n,m_{ij}}(r) = \frac{d}{dr}\hat{\boldsymbol{K}}_{n,m_{ij}}(r)$. Element $\hat{\boldsymbol{k}}_{n,1_{ij}}(r)$ can be expressed very similar to $\boldsymbol{k}_{ij}(r)$ in (21). Using notation $\hat{\boldsymbol{k}}_n(r) = \hat{\boldsymbol{k}}_{n,1}(r)$ and $\hat{\boldsymbol{k}}_{n,m}(r) = (\hat{\boldsymbol{k}}_{n,n-1} * \hat{\boldsymbol{k}}_{n-1,m-1})(r)$, $\forall n \geq m+1, m \geq 1$ we can write

$$\hat{\boldsymbol{k}}_{n,1_{ij}}(r) = \frac{-\boldsymbol{L}_{ii}}{\boldsymbol{C}_{n_{ii}}} e^{\frac{\boldsymbol{L}_{ii}}{\boldsymbol{C}_{n_{ii}}}r} \frac{\boldsymbol{B}_{ij}}{-\boldsymbol{L}_{ii}} + \int_{u=0}^{r} \frac{-\boldsymbol{L}_{ii}}{\boldsymbol{C}_{n_{ii}}} e^{\frac{\boldsymbol{L}_{ii}}{\boldsymbol{C}_{n_{ii}}}u} \left( \sum_{\substack{k\in\Phi,\\ k\neq i}} \frac{\boldsymbol{L}_{ik}}{-\boldsymbol{L}_{ii}} \hat{\boldsymbol{k}}_{n,1_{kj}}(r-u) \right.$$

$$\left. + \sum_{k\in\Phi}\sum_{m=1}^{\infty} \frac{\boldsymbol{F}_{m_{ik}}}{-\boldsymbol{L}_{ii}} \hat{\boldsymbol{k}}_{n+m,m+1_{kj}}(r-u) \right) du.$$

Taking the Laplace transform of the above equation and rearranging the result in a similar manner as before, we obtain

$$s\boldsymbol{C}_n\hat{\boldsymbol{k}}_n^*(s) = \boldsymbol{B} + \boldsymbol{L}\hat{\boldsymbol{k}}_n^*(s) + \sum_{m=1}^{\infty} \boldsymbol{F}_m\hat{\boldsymbol{k}}_{n+m,m+1}^*(s),$$
(47)

where

$$\hat{\boldsymbol{k}}_{n+m,m+1}^*(s) = \prod_{k=n+m}^{n} \hat{\boldsymbol{k}}_k^*(s) = \hat{\boldsymbol{k}}_{n+m}^*(s)\hat{\boldsymbol{k}}_{n+m-1}^*(s) \ldots \hat{\boldsymbol{k}}_{n+1}^*(s)\hat{\boldsymbol{k}}_n^*(s).$$

Similar to $\lim_{s\to 0}\boldsymbol{k}^*(s) = \boldsymbol{G}$, we have $\lim_{s\to 0}\hat{\boldsymbol{k}}_n^*(s) = \int_{r=0}^{\infty}\hat{\boldsymbol{k}}_n(r)dr = Pr(X(\rho_{n-1}) = v_{n-1,j}|X(0) = v_{n,i}) = \hat{\boldsymbol{G}}$. As a consequence, $\lim_{s\to 0}\hat{\boldsymbol{k}}_n^*(s)$ is level independent even though $\hat{\boldsymbol{k}}_n^*(s)$ is level dependent.

Using the moment generating property of the Laplace transform we define

$$\hat{\boldsymbol{M}}_{n_{ij}} \stackrel{def}{=} -\frac{d}{ds}\hat{\boldsymbol{k}}_{n_{ij}}^*(s)\Big|_{s=0} = \mathrm{E}\left[ I_{\{X(\rho_{n-1})=v_{n-1,j}\}} \int_{t=0}^{\rho_{n-1}} \boldsymbol{C}_{X(t),X(t)}dt \Big| X(0) = v_{n,i} \right].$$

The derivative of (47) at $s = 0$ gives

$$\boldsymbol{C}_n\hat{\boldsymbol{G}} = \boldsymbol{L}\hat{\boldsymbol{M}}_n + \sum_{m=1}^{\infty}\boldsymbol{F}_m\sum_{\ell=0}^{m}\hat{\boldsymbol{G}}^{\ell}\hat{\boldsymbol{M}}_{n+m-\ell}\hat{\boldsymbol{G}}^{m-\ell},$$
(48)

19

where we used that $\hat{G} = \lim_{s \to 0} \hat{k}_n^*(s)$ and $\hat{M}_n = -\frac{d}{ds}\hat{k}_n^*(s)\big|_{s=0}$. This equation has to be solved to obtain the required $\hat{M}_n$ matrices for $n = 1, 2 \ldots$.

The level dependence of $C_n$ makes $\hat{M}_n$ level dependent as well, while $\hat{A}$ is level independent. To separate the level dependent and level independent elements of $\hat{M}_n$, for $i, j, \ell \in \Phi$, $n \geq 1$, $m \geq 0$ we introduce

$$\hat{T}_{mij\ell} \overset{def}{=} \mathrm{E}\left[I_{\{X(\rho_{n-1})=v_{n-1,j}\}}\int_{t=0}^{\rho_{n-1}} I_{\{X(t)=v_{n+m,\ell}\}}dt\,\Big|\,X(0)=v_{n,i}\right]. \tag{49}$$

Due to the spatial homogeneity of the M/G/1 type process $\hat{T}_{mij\ell}$ is level independent, that is, it does not depend on $n$. Based on $\hat{T}_{mij\ell}$, $\hat{A}$ and $\hat{M}_n$ can be obtained as

$$\hat{A}_{ij} = \sum_{m=0}^{\infty}\sum_{\ell\in\Phi}\hat{T}_{mij\ell} \quad\text{and}\quad \hat{M}_{nij} = \sum_{m=0}^{\infty}\sum_{\ell\in\Phi}\hat{T}_{mij\ell}C_{n+m\ell\ell}. \tag{50}$$

The next task is to compute the global $S_D$ related measures $G$, $A$, and $M$ from the level related measures $\hat{G}$, $\hat{A}$, and $\hat{M}_n$. For states $i' = v_{m+1,i}$ and $j' = v_{1,j}$, we have that

$$g^*(s)_{i'j'} = \hat{g}_m^*(s)_{ij} = [\hat{g}^*(s)^m]_{ij},$$

and

$$k^*(s)_{i'j'} = \hat{k}_{m+1,m}^*(s)_{ij} = \left[\prod_{k=m+1}^{2}\hat{k}_k^*(s)\right]_{ij}.$$

Using these relations of the transform we obtain

$$G_{ij} = \lim_{s\to 0}g^*(s)_{i'j'} = \lim_{s\to 0}[\hat{g}^*(s)^m]_{ij} = [\hat{G}^m]_{ij},$$

$$A_{ij} = \frac{d}{ds}g^*(s)_{i'j'}\Big|_{s=0} = \frac{d}{ds}\hat{g}_m^*(s)_{ij}\Big|_{s=0} = \left[\frac{d}{ds}\hat{g}^*(s)^m\Big|_{s=0}\right]_{ij} =$$

$$\left[\sum_{\ell=0}^{m-1}\hat{g}^*(s)^\ell\frac{d}{ds}\hat{g}^*(s)\hat{g}^*(s)^{m-1-\ell}\Big|_{s=0}\right]_{ij} = \left[\sum_{\ell=0}^{m-1}\hat{G}^\ell\hat{A}\hat{G}^{m-1-\ell}\right]_{ij},$$

and

$$M_{i'j'} = \frac{d}{ds}k^*(s)_{i'j'}\Big|_{s=0} = \left[\sum_{\ell=0}^{m-1}\hat{G}^\ell\hat{M}_{m+1-\ell}\hat{G}^{m-1-\ell}\right]_{ij},$$

where the derivation of $M_{ij}$ follows the same patterns as the one of $A_{ij}$ and we used that $\lim_{s\to 0}\hat{k}_k^*(s) = \hat{G}$ is independent of level $k$.

Substituting $G$, $A$ and $M$ into (12), (13) and (14), for the compressed process we obtain (36), (37) and (38). $\qquad\square$

## 6.2 The G/M/1 type process

An MDP is of G/M/1 type with the considered $\{S_U, S_D\}$ partitioning if its generator matrix has the following block structure

$$\boldsymbol{Q}(\pi) = \left( \begin{array}{c|cccc} \bar{\boldsymbol{L}}(\pi) & \bar{\boldsymbol{F}}(\pi) & & & \\ \hline \bar{\boldsymbol{B}}_1 & \boldsymbol{L} & \boldsymbol{F} & & \\ \bar{\boldsymbol{B}}_2 & \boldsymbol{B}_1 & \boldsymbol{L} & \boldsymbol{F} & \\ \vdots & & \ddots & \ddots & \ddots & \ddots \end{array} \right). \tag{51}$$

The MDP can be partitioned to levels according to blocks so that each matrix block row in (51) corresponds to a level. In the following we assume that the reward rate matrix also has some level based regularity, thus

$$\boldsymbol{C}(\pi) = \left( \begin{array}{c|ccc} \boldsymbol{C}_1(\pi) & & & \\ \hline & \boldsymbol{C}_2 & & \\ & & \boldsymbol{C}_3 & \\ & & & \ddots \end{array} \right), \tag{52}$$

where $\boldsymbol{C}_n = f(\boldsymbol{C}_0, n), \forall n \geq 2$. We discuss some specific cases later in Section 7.1.

**Theorem 4.** *For an MDP with G/M/1 type structure in $S_D$, according to (51) and (52) the parameters of the reduced MDP representation are*

$$\boldsymbol{P}(\pi) = (-\boldsymbol{diagm}\langle \boldsymbol{Q_c}(\pi) \rangle)^{-1} (\boldsymbol{Q_c}(\pi) - \boldsymbol{diagm}\langle \boldsymbol{Q_c}(\pi) \rangle), \tag{53}$$

$$\tau(\pi) = (-\boldsymbol{diagm}\langle \boldsymbol{Q_c}(\pi) \rangle)^{-1} \left( \boldsymbol{I} + \sum_{m=1}^{\infty} \hat{\boldsymbol{R}}_1(\pi) \boldsymbol{R}^{m-1} \right) \underline{\boldsymbol{1}}, \tag{54}$$

$$c(\pi) = (-\boldsymbol{diagm}\langle \boldsymbol{Q_c}(\pi) \rangle)^{-1} \left( \boldsymbol{C}_1(\pi) + \sum_{m=1}^{\infty} \hat{\boldsymbol{R}}_1(\pi) \boldsymbol{R}^{m-1} \boldsymbol{C}_{m+1} \right) \underline{\boldsymbol{1}}, \tag{55}$$

*where $\boldsymbol{Q_c}(\pi) = \bar{\boldsymbol{L}}(\pi) + \sum_{m=1}^{\infty} \hat{\boldsymbol{R}}_1(\pi) \boldsymbol{R}^{m-1} \bar{\boldsymbol{B}}_m$,*

$$\hat{\boldsymbol{R}}_1(\pi) = \bar{\boldsymbol{F}}(\pi) \left( -\boldsymbol{L} - \sum_{m=1}^{\infty} \boldsymbol{R}^m \boldsymbol{B}_m \right)^{-1}$$

*and $\boldsymbol{R}$ is the solution of*

$$0 = \boldsymbol{F} + \boldsymbol{R}\boldsymbol{L} + \sum_{m=1}^{\infty} \boldsymbol{R}^{m+1} \boldsymbol{B}_m.$$

21

**Remark.** *Matrix $\boldsymbol{R}$ is the well studied characteristic matrix of the G/M/1 type process [13]. Efficient numerical methods are available for its computation e.g., in [4]. Unlike (38), the $\boldsymbol{C_n}$ matrices appear in (55) directly, which makes the computation of the reward term much simpler for the G/M/1 type process, because in this case there is no need to compute $\hat{\boldsymbol{M}}_{\boldsymbol{n}}$ from (41).*

*Proof.* Let $\eta_n$ be the time of the first visit to a level equal to or lower than $n$, that is, $\eta_n = \min_{k \in \{0,1,\dots,n\}}(\rho_k)$, furthermore we define the $\boldsymbol{V}(t)$ matrix function whose $ij$ element is

$$\boldsymbol{V}_{ij}(t) = Pr(X(t) = v_{n,j}, \eta_{n-1} > t | X(0) = v_{n,i}), \forall n > 1.$$

That is, $\boldsymbol{V}_{ij}(t)$ is the probability that the process, assuming that it starts in state $i$ of level $n$, visits state $j$ of level $n$ at time $t$ such that it does not visit any lower level before $t$. Furthermore, we define $\boldsymbol{R}(t) = \boldsymbol{F}\boldsymbol{V}(t)$. A stochastic interpretation of its $i, j$ element is

$$\boldsymbol{R}_{ij}(t) = \lim_{\Delta \to 0} \frac{1}{\Delta} Pr(X(t) = v_{n+1,j}, \eta_n > t, X(\Delta) \neq v_{n,i} | X(0) = v_{n,i}).$$

We define the multi level version of $\boldsymbol{R}(t)$ as

$$\boldsymbol{R_{m}}_{ij}(t) = \lim_{\Delta \to 0} \frac{1}{\Delta} Pr(X(t) = v_{n+m,j}, \eta_n > t, X(\Delta) \neq v_{n,i} | X(0) = v_{n,i}).$$

Starting from level $n$ and being at level $n + m$ at time $t$, let $\tau$ be the *last* instance when the process is at level $n + 1$ before time $t$. Then

$$\boldsymbol{R_{m}}_{ij}(t) = \sum_k \int_{\tau=0}^{t} \boldsymbol{R}_{ik}(\tau)\boldsymbol{R_{m-1}}_{kj}(t - \tau)d\tau$$

That is, for $m > 1$, $\boldsymbol{R_m}(t) = \boldsymbol{R}(t) * \boldsymbol{R_{m-1}}(t)$, with $\boldsymbol{R_1}(t) = \boldsymbol{R}(t)$ and $*$ denoting the convolution operator.

To evaluate $\boldsymbol{R_m}(t)$ we first compute $\boldsymbol{V}(t)$. We express $\boldsymbol{V}_{ij}(t)$ using the law of total probability as the event $X(t) = v_{n,j} | X(0) = v_{n,i}$ can partitioned the following way:

a) The first transition happens after time $t$ and $i = j$. The probability of this is $\delta_{ij} \int_{h=t}^{\infty} -\boldsymbol{L}_{ii}e^{\boldsymbol{L}_{ii}h}dh = \delta_{ij}e^{\boldsymbol{L}_{ii}t}$.

b) The first transition happens before time $t$ and this transition is inside level $n$. The probability of this is $\int_{h=0}^{t} -\boldsymbol{L}_{ii}e^{\boldsymbol{L}_{ii}h} \sum_{\substack{k \in \Phi \\ k \neq i}} \frac{\boldsymbol{L}_{ik}}{-\boldsymbol{L}_{ii}}\boldsymbol{V}_{kj}(t - h)dh$.

c) The first transition happens before time $t$ and this transition is to level $n+1$.

22

We can break down c) further according to the last level visited by the process before returning to level $n$ for the first time. The probability that this level is $n + m$ (where $m > 0$, and, to avoid special treatment of $m = 1$, assuming $(\boldsymbol{V} * \boldsymbol{R_0})(t) = \boldsymbol{V}(t)$) is

$$\int_{h=0}^{t} -\boldsymbol{L}_{ii}\mathrm{e}^{\boldsymbol{L}_{ii}h}\sum_{k\in\Phi}\frac{\boldsymbol{F}_{ik}}{-\boldsymbol{L}_{ii}}\big[(\boldsymbol{V} * \boldsymbol{R_{m-1}}\boldsymbol{B_m} * \boldsymbol{V})(t - h)\big]_{kj}dh$$

$$= \int_{h=0}^{t} \mathrm{e}^{\boldsymbol{L}_{ii}h}\big[(\boldsymbol{R_m}\boldsymbol{B_m} * \boldsymbol{V})(t - h)\big]_{ij}dh.$$

Combining a), b), and c) we have

$$\boldsymbol{V}_{ij}(t) = \underbrace{\delta_{ij}\mathrm{e}^{\boldsymbol{L}_{ii}t}}_{a)} + \underbrace{\int_{h=0}^{t}\mathrm{e}^{\boldsymbol{L}_{ii}h}\sum_{\substack{k\in\Phi\\k\neq i}}\boldsymbol{L}_{ik}\boldsymbol{V}_{kj}(t - h)dh}_{b)}$$

$$+ \underbrace{\int_{h=0}^{t}\mathrm{e}^{\boldsymbol{L}_{ii}h}\sum_{m=1}^{\infty}\big[(\boldsymbol{R_m}\boldsymbol{B_m} * \boldsymbol{V})(t - h)\big]_{ij}dh}_{c)}$$

Laplace transforming the above equation and multiplying by $s - \boldsymbol{L}_{ii}$ we get

$$(s - \boldsymbol{L}_{ii})\boldsymbol{V}_{ij}^{*}(s) = \delta_{ij} + \sum_{\substack{k\in\Phi\\k\neq i}}\boldsymbol{L}_{ik}\boldsymbol{V}_{kj}^{*}(s) + \big[\boldsymbol{R_m^*}(s)\boldsymbol{B_m}\boldsymbol{V}^{*}(s)\big]_{ij}.$$

After adding $\boldsymbol{L}_{ii}\boldsymbol{V}_{ii}^{*}(s)$ to both sides we can write the equation in matrix form as

$$s\boldsymbol{V}^{*}(s) = \boldsymbol{I} + \boldsymbol{L}\boldsymbol{V}^{*}(s) + \sum_{m=1}^{\infty}\boldsymbol{R_m^*}(s)\boldsymbol{B_m}\boldsymbol{V}^{*}(s),$$

from which

$$\boldsymbol{V}^{*}(s) = \left(s\boldsymbol{I} - \boldsymbol{L} - \sum_{m=1}^{\infty}\boldsymbol{R}^{*}(s)^{m}\boldsymbol{B_m}\right)^{-1}, \tag{56}$$

using $\boldsymbol{R_m^*}(s) = \boldsymbol{R}^{*}(s)^{m}$. Multiplying it with $\boldsymbol{F}$ from the left side we get

$$\boldsymbol{F}\boldsymbol{V}^{*}(s) = \boldsymbol{R}^{*}(s) = \boldsymbol{F}\left(s\boldsymbol{I} - \boldsymbol{L} - \sum_{m=1}^{\infty}\boldsymbol{R}^{*}(s)^{m}\boldsymbol{B_m}\right)^{-1}.$$

By multiplying with the term in the parentheses from the left and rearranging the equation we obtain

$$sR^*(s) = F + R^*(s)L + \sum_{m=1}^{\infty} R^*(s)^{m+1}B_m.$$

At $s \to 0$ this becomes

$$0 = F + RL + \sum_{m=1}^{\infty} R^{m+1}B_m, \tag{57}$$

where $R = \int_{\tau=0}^{\infty} R(\tau)d\tau = \lim_{s \to 0} R(s)$. Equation (57) is the well-known matrix equation to obtain matrix $R$, the characteristic matrix of the G/M/1 type process [13]. If the process starts in the irregular first level, then the same methodology can be applied, but

$$\hat{V}_{ij}(\pi, t) = \delta_{ij} \int_{h=t}^{\infty} -\bar{L}_{ii}(\pi)e^{\bar{L}_{ii}(\pi)h}dh$$

$$+ \int_{h=0}^{t} -\bar{L}_{ii}(\pi)e^{\bar{L}_{ii}(\pi)h} \sum_{\substack{k \in \Phi \\ k \neq i}} \frac{\bar{L}_{ik}(\pi)}{-\bar{L}_{ii}(\pi)} \hat{V}_{kj}\pi, t - h)dh$$

$$+ \int_{h=0}^{t} e^{\bar{L}_{ii}(\pi)h} \sum_{m=1}^{\infty} \left[ (\hat{R}_m(\pi)\bar{B}_m * \hat{V}(\pi))(t - h) \right]_{ij},$$

where $\hat{R}(\pi, t) = \bar{F}(\pi)V(t)$, $\hat{R}_m(\pi, t) = (\hat{R}(\pi) * R_{m-1})(t)$, $\forall m > 1$. Using the same steps as before we obtain

$$s\hat{R}^*(\pi, s) = \bar{F}(\pi) + \hat{R}^*(\pi, s)\bar{L} + \sum_{m=1}^{\infty} \hat{R}_{m+1}^*(\pi, s)B_m. \tag{58}$$

where $\hat{R}_{m+1}^*(\pi, s) = \hat{R}^*(\pi, s)R_m^*(s) = \hat{R}^*(\pi, s)R^*(s)^m$, from which at $s \to 0$ (58) becomes

$$0 = \bar{F}(\pi) + \hat{R}(\pi)L + \sum_{m=1}^{\infty} \hat{R}(\pi)R^mB_m, \tag{59}$$

where $\hat{R}(\pi) = \lim_{s \to 0} \hat{R}^*(\pi, s)$ and $\hat{R}_m(\pi) = \lim_{s \to 0} \hat{R}_m^*(\pi, s)$. By rearranging (59) we get that

$$\hat{R}(\pi) = \bar{F}(\pi) \left( -L - \sum_{m=1}^{\infty} R^mB_m \right)^{-1}, \quad \hat{R}_{m+1}(\pi) = \hat{R}(\pi)R^m. \tag{60}$$

For G/M/1 type processes we cannot adopt the approach used for the M/G/1 type process. In this case, we directly express the $\boldsymbol{P}(\pi), \tau(\pi)$, and $c(\pi)$ parameters.

From the definition of $\boldsymbol{R_m}$, $\frac{1}{-\boldsymbol{L}_{ii}}\boldsymbol{R_{mij}}$ is the mean time spent in $v_{n+m,j}$ until the process visits some level below $n$, starting from $v_{n,i}$. For $\boldsymbol{P}(\pi)$ we use the same decomposition as for the general case in Section 4, where we had the following cases:

- Case 1: the process moves to state $j$ ($j \in S_U \setminus i$) directly.

- Case 2: the process first moves to $k \in S_D$, spends some time in $S_D$, then enters $S_U$ in state $j$.

- Case 3: the process first moves to $k \in S_D$, spends some time in $S_D$, then enters $S_U$ in state $i$. This case adds a recursive term to the formulas.

Using this decomposition we can write, for $i \neq j$:

$$\boldsymbol{P}_{ij}(\pi) = \underbrace{\frac{\bar{\boldsymbol{L}}_{ij}(\pi)}{-\bar{\boldsymbol{L}}_{ii}(\pi)}}_{Case\ 1} + \underbrace{\sum_{m=1}^{\infty}\sum_{k\in\Phi}\frac{1}{-\bar{\boldsymbol{L}}_{ii}(\pi)}\hat{\boldsymbol{R}}_{\boldsymbol{m}ik}(\pi)\bar{\boldsymbol{B}}_{\boldsymbol{m}kj}}_{Case\ 2}$$

$$+ \underbrace{\sum_{m=1}^{\infty}\sum_{k\in\Phi}\frac{1}{-\bar{\boldsymbol{L}}_{ii}(\pi)}\hat{\boldsymbol{R}}_{\boldsymbol{m}ik}(\pi)\bar{\boldsymbol{B}}_{\boldsymbol{m}ki}\boldsymbol{P}_{ij}(\pi)}_{Case\ 3},$$

and $\boldsymbol{P}_{ii}(\pi) = 0$. That is, for $i \neq j$

$$-\left(\bar{\boldsymbol{L}}_{ii}(\pi) + \sum_{m=1}^{\infty}[\hat{\boldsymbol{R}}_{\boldsymbol{m}}(\pi)\bar{\boldsymbol{B}}_{\boldsymbol{m}}]_{ii}\right)\boldsymbol{P}_{ij}(\pi) = \bar{\boldsymbol{L}}_{ij}(\pi) + \sum_{m=1}^{\infty}[\hat{\boldsymbol{R}}_{\boldsymbol{m}}(\pi)\bar{\boldsymbol{B}}_{\boldsymbol{m}}]_{ij},$$

which can be written in matrix from as

$$\boldsymbol{P}(\pi) = (-\mathbf{diagm}\langle\boldsymbol{Q_c}(\pi)\rangle)^{-1}\left(\boldsymbol{Q_c}(\pi) - \mathbf{diagm}\langle\boldsymbol{Q_c}(\pi)\rangle\right),$$

where $\boldsymbol{Q_c}(\pi) = \bar{\boldsymbol{L}}(\pi) + \sum_{m=1}^{\infty}\hat{\boldsymbol{R}}_{\boldsymbol{m}}(\pi)\bar{\boldsymbol{B}}_{\boldsymbol{m}}$. The subtraction of the diagonal matrix in the second term ensures that the diagonal of $\boldsymbol{P}(\pi)$ is zero. Using $\hat{\boldsymbol{R}}_{\boldsymbol{m}}(\pi) = \hat{\boldsymbol{R}}_{\boldsymbol{1}}(\pi)\boldsymbol{R}^{m-1}$, we can also write $\boldsymbol{Q_c}(\pi) = \bar{\boldsymbol{L}}(\pi) + \sum_{m=1}^{\infty}\hat{\boldsymbol{R}}_{\boldsymbol{1}}(\pi)\boldsymbol{R}^{m-1}\bar{\boldsymbol{B}}_{\boldsymbol{m}}$.

To compute $\tau(\pi)$, let $\gamma_{S_U\setminus i}$ be the first time to reach $S_U \setminus i$, furthermore let $\tau_i(\pi,t) = Pr(\gamma_{S_U\setminus i} > t|X(0) = i)$ and $\tau_i^*(\pi,s) = \int_t \tau_i(\pi,t)e^{-st}dt$. For $\tau_i^*(\pi,s)$

we have

$$\tau_i^*(\pi, s) = \underbrace{\frac{1}{s - \bar{\boldsymbol{L}}_{ii}(\pi)}}_{\text{no transition till } t}$$

$$+ \sum_{m=1}^{\infty} \sum_{k \in \Phi} \sum_{\ell \in \Phi} \underbrace{\frac{-\bar{\boldsymbol{L}}_{ii}(\pi)}{s - \bar{\boldsymbol{L}}_{ii}(\pi)}}_{\text{transition at } x(< t)} \frac{\bar{\boldsymbol{F}}_{ik}(\pi)}{-\bar{\boldsymbol{L}}_{ii}(\pi)} \underbrace{\left(\boldsymbol{V}^*(s)\boldsymbol{R}_{m-1}^*(s)\right)_{k\ell}}_{\text{time in } S_D > t - x}$$

$$+ \sum_{m=1}^{\infty} \sum_{k \in \Phi} \underbrace{\frac{-\bar{\boldsymbol{L}}_{ii}(\pi)}{s - \bar{\boldsymbol{L}}_{ii}(\pi)}}_{\text{tr. at } x} \frac{\bar{\boldsymbol{F}}_{ik}(\pi)}{-\bar{\boldsymbol{L}}_{ii}(\pi)} \underbrace{\left(\boldsymbol{V}^*(s)\boldsymbol{R}_{m-1}^*(s)\bar{\boldsymbol{B}}_m\right)_{ki}}_{\text{return to } i \text{ at } y \ (x < y < t)} \underbrace{\tau_i^*(\pi, s)}_{\text{time to } S_U \setminus i > t - y} .$$

Multiplying with $s - \bar{\boldsymbol{L}}_{ii}(\pi)$ gives

$$(s - \bar{\boldsymbol{L}}_{ii}(\pi))\tau_i^*(\pi, s) = 1 + \sum_{m=1}^{\infty} \sum_{k \in \Phi} \sum_{\ell \in \Phi} \bar{\boldsymbol{F}}_{ik}(\pi) \left(\boldsymbol{V}^*(s)\boldsymbol{R}_{m-1}^*(s)\right)_{k\ell}$$

$$+ \sum_{m=1}^{\infty} \sum_{k \in \Phi} \bar{\boldsymbol{F}}_{ik}(\pi) \left(\boldsymbol{V}^*(s)\boldsymbol{R}_{m-1}^*(s)\bar{\boldsymbol{B}}_m\right)_{ki} \tau_i^*(\pi, s) =$$

$$1 + \sum_{m=1}^{\infty} \sum_{\ell \in \Phi} \hat{\boldsymbol{R}}_m^*(\pi, s)_{i\ell} + \sum_{m=1}^{\infty} \left(\hat{\boldsymbol{R}}_m^*(\pi, s)\bar{\boldsymbol{B}}_m\right)_{ii} \tau_i^*(\pi, s).$$

We are interested in the mean time to get to $S_U \setminus i$ which is $\int_t \tau_i(\pi, t)dt = \tau_i^*(\pi, s)|_{s=0} \triangleq \tau_i(\pi)$. Substituting $s = 0$ we have

$$- \bar{\boldsymbol{L}}_{ii}(\pi)\tau_i(\pi) = 1 + \sum_{m=1}^{\infty} \sum_{\ell \in \Phi} \hat{\boldsymbol{R}}_m(\pi)_{i\ell} + \sum_{m=1}^{\infty} \left(\hat{\boldsymbol{R}}_m(\pi)\bar{\boldsymbol{B}}_m\right)_{ii} \tau_i(\pi),$$

where we used that $\hat{\boldsymbol{R}}_m(\pi) = \hat{\boldsymbol{R}}_m^*(\pi, s)|_{s=0}$. From this

$$- \left(\bar{\boldsymbol{L}}_{ii}(\pi) + \sum_{m=1}^{\infty}[\hat{\boldsymbol{R}}_m(\pi)\bar{\boldsymbol{B}}_m]_{ii}\right) \tau_i(\pi) = 1 + \sum_{m=1}^{\infty} \sum_{\ell \in \Phi} \hat{\boldsymbol{R}}_{mi\ell}(\pi).$$

Using $\hat{\boldsymbol{R}}_m(\pi) = \hat{\boldsymbol{R}}_1(\pi)\boldsymbol{R}^{m-1}$ its matrix form is (54).

Based on a similar argument, for $c_i(\pi)$ we have

$$c_i(\pi) = \frac{\boldsymbol{C}_{1ii}}{-\bar{\boldsymbol{L}}_{ii}(\pi)} + \sum_{m=1}^{\infty} \sum_{k \in \Phi} \frac{1}{-\bar{\boldsymbol{L}}_{ii}(\pi)} \hat{\boldsymbol{R}}_{mik}(\pi)\boldsymbol{C}_{m+1kk}$$

$$+ \sum_{m=1}^{\infty} \sum_{k \in \Phi} \frac{1}{-\bar{\boldsymbol{L}}_{ii}(\pi)} \hat{\boldsymbol{R}}_{mik}(\pi)\bar{\boldsymbol{B}}_{mki}c_i(\pi),$$

26

from which

$$-\left(\bar{L}_{ii}(\pi) + \sum_{m=1}^{\infty}[\hat{R}_m(\pi)\bar{B}_m]_{ii}\right)c_i(\pi) = C_{1ii}(\pi) + \sum_{m=1}^{\infty}\sum_{k\in\Phi}\hat{R}_{mik}(\pi)C_{m+1kk},$$

whose matrix form is (55)                                                                    □

# 7 Calculation of $c(\pi)$ with different reward structures

In this section, we provide methods to calculate $c(\pi)$ when the infinite reward matrix $C(\pi)$ follows different regular structures. Specifically, we consider the following cases:

- geometric series: $C_n = \kappa^n C_0, \forall n > 1$,

- matrix geometric series: $C_n = C_0{}^n, \forall n > 1$,

- polynomial series: $C_n = p(n)C_0, \forall n > 1$,

where $C_0$ is arbitrary diagonal matrix, $\kappa$ is a positive real number and $p(n)$ is a finite order arbitrary non-negative polynomial of $n$.

## 7.1 Calculating the $c(\pi)$ reward function for G/M/1 type processes

For G/M/1 type processes, (55) defines the relation of $c(\pi)$ with $C_1(\pi)$ and $C_m$ ($m \geq 2$). The only non-trivial part of this formula is the computation of the infinite sum $\sum_{m=1}^{\infty} R^{m-1}C_{m+1}$.

### 7.1.1 Reward function for $C_n = \kappa^n C_0$

If $C_n = \kappa^n C_0$, then the infinite sum converges when $\kappa\lambda_R < 1$, where $\lambda_R$ is the spectral radius of $R$. In this case,

$$\sum_{m=1}^{\infty} R^{m-1}C_{m+1} = \sum_{m=1}^{\infty} R^{m-1}\kappa^{m+1}C_0 = \kappa^2(I - \kappa R)^{-1}C_0.$$

### 7.1.2 Reward function for $C_n = C_0{}^n$

If $C_n = C_0{}^n$, then $\sum_{m=1}^{\infty} R^{m-1}C_0{}^{m-1}C_0{}^2 = XC_0{}^2$, where $X = \sum_{m=1}^{\infty} R^{m-1}C_0{}^{m-1}$ is the solution of the Sylvester equation $X = I + RXC_0$.

### 7.1.3 Reward function for $C_n = p(n)C_0$

If $C_n = p(n)C_0$, without loss of generality, we assume $p(n) = \sum_{i=0}^{k} a_i(n-2)^i$. In this case,

$$
\begin{aligned}
\sum_{m=1}^{\infty} \boldsymbol{R}^{m-1}\boldsymbol{C_{m+1}} &= \sum_{m=1}^{\infty} \boldsymbol{R}^{m-1} \sum_{i=0}^{k} a_i(m-1)^i \boldsymbol{C_0} \\
&= a_0(\boldsymbol{I}-\boldsymbol{R})^{-1}\boldsymbol{C_0} + \sum_{i=1}^{k} a_i \sum_{m=2}^{\infty} \frac{\boldsymbol{R}^{m-1}}{(m-1)^{-i}}\boldsymbol{C_0} \\
&= a_0(\boldsymbol{I}-\boldsymbol{R})^{-1}\boldsymbol{C_0} + \sum_{i=1}^{k} a_i Li_{-i}(\boldsymbol{R})\boldsymbol{C_0},
\end{aligned}
$$

where $Li_\ell(\boldsymbol{Y})$ is the polylogarithm function generalised for matrices, i.e.,

$$
Li_\ell(\boldsymbol{R}) = \sum_{m=1}^{\infty} \frac{\boldsymbol{R}^m}{m^\ell}. \tag{61}
$$

If $\ell \in \mathbb{Z}^+$ and the spectral radius of $\boldsymbol{R}$ is less than one (which holds if the generator of the MDP is positive recurrent) $Li_{-\ell}(\boldsymbol{R})$ is finite and can be computed as

$$
Li_{-k}(\boldsymbol{R}) = (-1)^{k+1} \sum_{i=0}^{k-1} i! S(k+1, i+1)(\boldsymbol{R}-\boldsymbol{I})^{i+1},
$$

(see e.g. [18]), where $S(k,i) = \frac{1}{k!}\sum_{j=0}^{k}(-1)^{k-j}\binom{k}{j}j^n$ denotes the Stirling number of second kind.

## 7.2 Calculating the $c(\pi)$ reward function for M/G/1 type processes

For M/G/1 type processes, (41) needs to be solved for $\hat{\boldsymbol{M}}_{\boldsymbol{n}}$, which is hard in general. To simplify the discussion, in this section, we utilize the fact that the M/G/1 type processes is positive recurrent and consequently $\hat{\boldsymbol{G}}\underline{\boldsymbol{1}} = \underline{\boldsymbol{1}}$. Multiplying both sides of (41) with $\underline{\boldsymbol{1}}$ and using $\hat{\boldsymbol{G}}\underline{\boldsymbol{1}} = \underline{\boldsymbol{1}}$ we have

$$
\begin{aligned}
\boldsymbol{C_n}\underline{\boldsymbol{1}} &= \boldsymbol{L}\boldsymbol{\mu_n} + \sum_{m=1}^{\infty} \boldsymbol{F_m} \sum_{\ell=0}^{m} \hat{\boldsymbol{G}}^\ell \boldsymbol{\mu_{n+m-\ell}}, \\
&= \boldsymbol{L}\boldsymbol{\mu_n} + \sum_{m=1}^{\infty} \boldsymbol{F_m}\boldsymbol{\mu_{n+m}} + \sum_{\ell=1}^{\infty}\sum_{m=\ell}^{\infty} \boldsymbol{F_m}\hat{\boldsymbol{G}}^\ell \boldsymbol{\mu_{n+m-\ell}}, \tag{62}
\end{aligned}
$$

where $\boldsymbol{\mu_n} = \hat{\boldsymbol{M}}_{\boldsymbol{n}}\underline{\boldsymbol{1}}$, $\forall n > 0$, which is sufficient to compute $c(\pi)$ according to (38).

28

### 7.2.1 Reward function for $C_n = \kappa^n C_0$

If $C_n = \kappa^n C_0$, then from (50) we have the following relation for $n > 0$

$$\hat{M}_{nij} = \sum_{m=0}^{\infty} \sum_{\ell \in \Phi} \hat{T}_{mij\ell} C_{n+m\,\ell\ell} = \sum_{m=0}^{\infty} \sum_{\ell \in \Phi} \hat{T}_{mij\ell} \kappa^n C_{m\ell\ell} = \kappa^n \hat{M}_{0ij},$$

and consequently, $\mu_n = \kappa^n \mu_0$. Substituting this into (62), $\mu_0$ can be computed from

$$C_0 \underline{1} = \left( L + \sum_{m=1}^{\infty} \kappa^m F_m + \sum_{\ell=1}^{\infty} \sum_{m=\ell}^{\infty} \kappa^{m-\ell} F_m \hat{G}^\ell \right) \mu_0, \qquad (63)$$

where the existence and the singularity of the matrix in bracket depends on $L$, $F_m$ and $\hat{G}$.

### 7.2.2 Reward function for $C_n = C_0{}^n$

This case does not provide a simple relation for the $\mu_n$ series, which makes the solution of (41) possible in general. A practically important case, when $F_m = F^m$, allows analytically compact description and is discussed below.

### 7.2.3 Reward function for $C_n = p(n)C_0$

For $C_n = p(n)C_0 = \sum_{u=0}^{k} a_u n^u C_0$ we utilize the linear structure of (62) and separate the solution into the following sub-problems

$$C_{n,u} \underline{1} = L\mu_{n,u} + \sum_{m=1}^{\infty} F_m \mu_{n+m,u} + \sum_{\ell=1}^{\infty} \sum_{m=\ell}^{\infty} F_m \hat{G}^\ell \mu_{n+m-\ell,u}, \qquad (64)$$

where $C_{n,u} = n^u C_0$, for $u = \{0, 1, \ldots, k\}$. From the solutions for the sub-problems, $\mu_{n,u}$, the solution for $C_n = \sum_{i=0}^{k} a_u C_{n,u}$ is obtained as

$$\mu_n = \sum_{u=0}^{k} a_u \mu_{n,u}.$$

For $\hat{M}_{n,0}$ we have

$$\hat{M}_{n,0_{ij}} = \sum_{m=0}^{\infty} \sum_{\ell \in \Phi} \hat{T}_{mij\ell} C_{n+m,0_{\ell\ell}}$$

$$= \sum_{m=0}^{\infty} \sum_{\ell \in \Phi} \hat{T}_{mij\ell} (n+m)^0 C_{0\ell\ell} = \sum_{m=0}^{\infty} \sum_{\ell \in \Phi} \hat{T}_{mij\ell} C_{0\ell\ell} \overset{def}{=} \hat{M},$$

from which $\boldsymbol{\mu_{n,0}} = \hat{\boldsymbol{M}}\underline{\boldsymbol{1}} \overset{def}{=} \boldsymbol{\mu}$ is independent of $n$ and can be computed from

$$C_0\underline{1} = \left( L + \sum_{m=1}^{\infty} F_m + \sum_{\ell=1}^{\infty}\sum_{m=\ell}^{\infty} F_m \hat{G}^{\ell} \right) \boldsymbol{\mu}. \tag{65}$$

For $u > 0$ we have

$$\hat{M}_{n+v,u_{ij}} = \sum_{m=0}^{\infty}\sum_{\ell\in\Phi} \hat{T}_{mij\ell} C_{n+v+m,u_{\ell\ell}} = \sum_{m=0}^{\infty}\sum_{\ell\in\Phi} \hat{T}_{mij\ell}(n+v+m)^u C_{0\ell\ell}$$

$$= \sum_{m=0}^{\infty}\sum_{\ell\in\Phi} \hat{T}_{mij\ell} \sum_{r=0}^{u} \binom{u}{r} n^{u-r}(v+m)^r C_{0\ell\ell}$$

$$= \sum_{r=0}^{u} \binom{u}{r} n^{u-r} \sum_{m=0}^{\infty}\sum_{\ell\in\Phi} \hat{T}_{mij\ell} C_{v+m,r_{\ell\ell}} = \sum_{r=0}^{u} \binom{u}{r} n^{u-r} \hat{M}_{v,r_{ij}},$$

that is $\boldsymbol{\mu_{n+v,u}} = \sum_{r=0}^{u} \binom{u}{r} n^{u-r} \boldsymbol{\mu_{v,r}}$.

For $v = 1$ this gives, $\boldsymbol{\mu_{n+1,u}} = \sum_{r=0}^{u} \binom{u}{r} n^{u-r} \boldsymbol{\mu_{1,r}}$. Substituting it into (64) for $n = 1$ and $u > 0$ gives an equation in which the unknowns are $\boldsymbol{\mu_{1,r}}$ for $r = 0, 1, \ldots, u$. E.g., for $n = 1, u = 1$ we have

$$C_0\underline{1} - \left( \sum_{m=1}^{\infty} m F_m + \sum_{\ell=1}^{\infty}\sum_{m=\ell}^{\infty}(m-\ell) F_m \hat{G}^{\ell} \right) \boldsymbol{\mu_{1,0}} =$$

$$\left( L + \sum_{m=1}^{\infty} F_m + \sum_{\ell=1}^{\infty}\sum_{m=\ell}^{\infty} F_m \hat{G}^{\ell} \right) \boldsymbol{\mu_{1,1}},$$

from which $\boldsymbol{\mu_{1,1}}$ can be computed, since $\boldsymbol{\mu_{1,0}} = \boldsymbol{\mu}$ is known from (65). Recursively, applying the same procedure for $u = 1, 2, \ldots, k$ provides the required $\boldsymbol{\mu_{1,u}}$ matrices, from which all $\boldsymbol{\mu_{n,u}}$ matrices ($n \geq 1, u \geq 0$) can be calculated using $\boldsymbol{\mu_{n,u}} = \sum_{r=0}^{u} \binom{u}{r}(n-1)^{u-r} \boldsymbol{\mu_{1,r}}$.

### 7.2.4 Special case of matrix geometric $F_m$ series

When $\boldsymbol{F_m} = \boldsymbol{F}^m$, the infinite summations of the previous subsections containing $\boldsymbol{F_m}$ simplifies significantly. E.g., in (63)

$$\sum_{\ell=1}^{\infty}\sum_{m=\ell}^{\infty} \kappa^{m-\ell} F_m \hat{G}^{\ell} = \sum_{\ell=1}^{\infty}\sum_{m=\ell}^{\infty} \kappa^{m-\ell} F^{m-\ell} F^{\ell} \hat{G}^{\ell} = (I - \kappa F)^{-1} \underbrace{\sum_{\ell=1}^{\infty} F^{\ell} \hat{G}^{\ell}}_{X-I}$$

$$= (I - \kappa F)^{-1}(X - I),$$

30

where $X$ is the solution of the Sylvester equation $X = I + FX\hat{G}$.

In addition to this analytical simplicity, $F_m = F^m$ makes it possible to compute the solution for the matrix geometric reward function.

**Theorem 5.** *If $C_n = C_0{}^n$ and $F_m = F^m$ then the solution of (62) is $\mu_n = HC_0{}^n\mathbf{1}$, where $H$ is the solution of the Sylvester equation*

$$H = (L - I + X)^{-1}(I - FC_0) + (L - I + X)^{-1}F(L - I)\ H\ C_0 \quad (66)$$

*and $X$ is the solution of the Sylvester equation $X = I + FX\hat{G}$.*

*Proof.* If $C_n = C_0{}^n$ and $F_m = F^m$ then (62) takes the form

$$C_0{}^n\mathbf{1} = L\mu_n + \sum_{m=1}^{\infty}\sum_{\ell=0}^{m}F^m\hat{G}^{\ell}\mu_{n+m-\ell}$$

$$= (L - I)\mu_n + \sum_{m=0}^{\infty}\sum_{\ell=0}^{m}F^m\hat{G}^{\ell}\mu_{n+m-\ell}$$

$$= (L - I)\mu_n + \sum_{k=0}^{\infty}\sum_{m=k}^{\infty}F^m\hat{G}^{m-k}\mu_{n+k}$$

$$= (L - I)\mu_n + \sum_{k=0}^{\infty}F^kX\mu_{n+k},$$

which suggests a matrix geometric solution, $\mu_n = W^n\mu$. Substituting this solution, for $n > 0$ we get

$$C_0{}^n\mathbf{1} = (L - I)\mu_n + \sum_{k=0}^{\infty}F^kXW^{n+k}\mu = \left(L - I + \sum_{k=0}^{\infty}F^kXW^k\right)W^n\mu,$$

Since $C_0$ is a diagonal matrix the spectral decomposition of $W$ should be $W = HC_0H^{-1}$, where the unknowns (matrix $H$ and vector $\mu$) are defined by $H^{-1}\mu = \mathbf{1}$ and

$$I = \left(L - I + \sum_{k=0}^{\infty}F^kXHC_0{}^kH^{-1}\right)H, \quad (67)$$

since

$$C_0{}^n\mathbf{1} = \underbrace{\left(L - I + \sum_{k=0}^{\infty}F^kXHC_0{}^kH^{-1}\right)}_{I}H\,C_0{}^n\,\underbrace{H^{-1}\mu}_{\mathbf{1}}.$$

Let $Z = \sum_{k=0}^{\infty} F^k X H C_0^{\,k}$. On the one hand $Z$ is the solution of the Sylvester equation $Z = XH + FZC_0$. On the other hand, from (67), we have

$$I - (L - I) H = Z,$$

which gives the following linear equation for $H$

$$I - (L - I) H = XH + F (I - (L - I) H) C_0,$$

whose standard Sylvester form is (66), and from $\mu = H\underline{1}$ and $\mu_n = W^n \mu$ we have $\mu_n = H_0^{\,n} \underline{1}$ which was to be proven. $\qquad\square$

# 8   Conclusions

We presented a methodology for computing a reduced representation of MDPs when there is a finite subset of states with decisions and a potentially infinite subset of decision independent states. This methodology requires the computation of some state dependent reward measures, which we preformed for two practically important cases when the infinite subset of decision independent states has M/G/1-type and G/M/1-type structure. The special case, when the decision independent part has a QBD structure, can be computed by either of the two general cases. Some required measures are already provided in the literature of M/G/1-type and G/M/1-type processes, but the reward related measures, e.g., the ones discussed in the previous section has not been considered before.

# References

[1] Eitan Altman. *Constrained Markov decision processes*, volume 7. CRC Press, 1999.

[2] Eitan Altman. *Applications of Markov decision processes in communication networks: A survey*. PhD thesis, INRIA, 2000.

[3] Nicole Bäuerle and Ulrich Rieder. *Markov decision processes with applications to finance*. Springer Science & Business Media, 2011.

[4] Dario A. Bini, Guy Latouche, and Beatrice Meini. *Numerical Methods for Structured Markov Chains (Numerical Mathematics and Scientific Computation)*. Oxford University Press, Inc., New York, NY, USA, 2005.

[5] L. Bodrog, M. Gribaudo, G. Horváth, A. Mészáros, and M. Telek. Control of queues with MAP servers: experimental results. In *The Eighth International Conference on Matrix-Analytic Methods in Stochastic Models (MAM8)*, pages 9–11, 2014.

[6] Craig Boutilier, Richard Dearden, and Moisés Goldszmidt. Stochastic dynamic programming with factored representations. *Artificial Intelligence*, 121(1):49–107, 2000.

[7] Michael George Bulmer. *Principles of statistics*. Courier Corporation, 1979.

[8] Dmitry Efrosinin. *Controlled Queueing Systems with Heterogeneous Servers*. PhD thesis, University of Trier, 2004.

[9] Ronald A Howard. *Dynamic programming and Markov processes*. John Wiley, 1960.

[10] Krishna Jagannathan, Shie Mannor, Ishai Menache, and Eytan Modiano. A state action frequency approach to throughput maximization over uncertain wireless channels. *Internet Mathematics*, 9(2-3):136–160, 2013.

[11] Michael Kearns and Satinder Singh. Near-optimal reinforcement learning in polynomial time. *Machine Learning*, 49(2-3):209–232, 2002.

[12] Yasar Levent Kocaga and Amy R. Ward. Admission control for a multi-server queue with abandonment. *Queueing Systems*, 65(3):275–323, 2010.

[13] G. Latouche and V. Ramaswami. *Introduction to Matrix-Analytic Methods in Stochastic Modeling*. Series on statistics and applied probability. ASA-SIAM, 1999.

[14] A. Mészáros and M. Telek. Markov decision process and linear programming based control of MAP/MAP/N queues. In *Computer Performance Engineering, EPEW*, volume 8721 of *LNCS*, pages 179–193, Sept. 2014.

[15] Marcel Neuts. *Structured Stochastic Matrices of M/G/1 Type and Their Applications*. Probability: Pure and Applied. Taylor & Francis, 1989.

[16] Andrew Y Ng and Michael Jordan. Pegasus: A policy search method for large mdps and pomdps. In *Proceedings of the Sixteenth conference on Uncertainty in artificial intelligence*, pages 406–415. Morgan Kaufmann Publishers Inc., 2000.

[17] Joris Slegers, Isi Mitrani, and Nigel Thomas. Optimal dynamic server allocation in systems with on/off sources. In *Formal Methods and Stochastic Models for Performance Evaluation*, pages 186–199. Springer, 2007.

[18] David Wood. The computation of polylogarithms. Technical Report 15-92*, University of Kent, Computing Laboratory, University of Kent, Canterbury, UK, June 1992.