

An explicit solution for the value function of a priority queue

Ger Koole

Vrije Universiteit
De Boelelaan 1081a
1081 HV Amsterdam
The Netherlands
koole@few.vu.nl

Philippe Nain

INRIA
2004 Route des Lucioles
06902 Sophia Antipolis
France
nain@sophia.inria.fr

Appeared as *Queueing Systems* **47**: 251–282, 2004

Abstract

We consider a multiclass preemptive-resume priority queue with Poisson arrivals and general service times. We derive explicit expressions for the discounted expected and long-run average weighted queue lengths and switching costs, the latter one only in the case of exponential service times. We illustrate our results with numerical calculations.

1 Introduction

We consider an N -class preemptive-resume priority queue with Poisson arrivals, general class-dependent service time distributions, and holding and switching costs. This queue is equipped with a single server and an infinite buffer. We derive closed-form expressions for the total expected discounted holding and switching costs for any initial state. The expression for the holding costs holds for general service time distributions, the one for the switching costs only for exponential service times. Using limiting arguments we also calculate the long-run average holding and switching costs as well as the bias function. Thus we find the solution of the so-called Poisson equation. A careful description of the system at hand is given in Section 2.

This paper is a follow-up of both [9] and [5]. In these papers the case of two customer classes and exponential service times is treated, for the discounted and for the average cost criterion, respectively. Another related study is the work by Harrison [6] who found a closed-form expression for the total expected discounted holding cost in a multiclass non-preemptive resume priority queue with Poisson inputs and general services (but no switching costs). To the best of our knowledge these works are the only ones where expected total discounted or average costs are computed for multi-dimensional systems with stochastically dependent components.

The derivation of closed-form expressions for the costs associated with a given queueing system is interesting in its own, especially when the system is as standard as the one addressed in this paper. Another interest resides in the use of this work, and of its possible extensions, in the context of the optimal control of queues. For example, the preemptive priority rule is one of the possible policies in the system where a controller can decide at each instant which queue to serve. If the switching costs are positive, then this policy is certainly not optimal. If $N \geq 3$ or 4 the direct methods such as value iteration cannot be executed anymore, due to the curse of dimensionality. Thus we have to rely on approximation methods, for example one-step optimization (see, e.g., [12, 15]) or reinforcement learning (see, e.g., [1, 3]). In the former method it is necessary to start with the solution of the optimality equation for a fixed policy; in the latter method one is interested in the form of the solution. For the preemptive priority rule we give the solution to the optimality equation in this paper. In Section 7 we come back to these applications.

This paper is organized as follows: Section 2 introduces the model and various cost criterions; preliminary results, mostly related to transforms of busy periods, are collected in Section 3. Explicit forms for the discounted expected holding cost and switching are derived in Section 4 and in Section 5, respectively. The average cost criteria are considered in Section 6. The applications in Section 7 conclude the paper.

2 Model description

Consider a multiclass single-server queue with infinite waiting room under the preemptive resume priority service discipline [7, p. 53]. Under this service discipline the preempted customer resumes service from the point where it was interrupted. There are N classes of customers and we assume that customers of class m have priority over customers of class n if $m < n$.

Customers of class n ($n = 1, 2, \dots, N$) arrive at the system according to a Poisson process with constant intensity $\lambda_n > 0$. The consecutive service times required by customers of class n form an i.i.d. sequence of random variables. Denote by $S_n(s)$, $s \geq 0$, the Laplace-Stieltjes Transform (LST) of the service times for customers of class n . We further assume that all (Poisson) arrival processes and service times processes are mutually independent.

There are deterministic holding costs ($c_n \geq 0$ for customers of class n) and switching costs. We denote by $s_{m,n} \geq 0$, $m \neq n$, the instantaneous cost incurred when the server switches from class- m to class- n . By convention we set $s_{n,n} = 0$ for $n = 1, 2, \dots, N$.

Let $X_n(t)$ be the number of customers of class n in the system at time $t \geq 0$ and let $T_k^{m,n}$ be the time when the server switches from class- m to class n for the k th time.

Given that there are x_n customers of class n ($n = 1, 2, \dots, N$) at time $t = 0$ which have received no service yet, and that the server is attending customers of class z ($z = 1, 2, \dots, N$) at $t = 0$, the overall discounted cost incurred in $[0, \infty)$ is given by (with $\beta > 0$ a constant

held fixed throughout)

$$V_{\mathbf{x},z}(\beta) = V_{\mathbf{x},z}^h(\beta) + V_{\mathbf{x},z}^s(\beta) \quad (1)$$

with

$$V_{\mathbf{x},z}^h(\beta) := \mathbf{E}_{\mathbf{x},z} \left[\int_0^\infty e^{-\beta t} \sum_{n=1}^N c_n X_n(t) dt \right] \quad (2)$$

$$V_{\mathbf{x},z}^s(\beta) := \mathbf{E}_{\mathbf{x},z} \left[\sum_{k=1}^\infty \sum_{\substack{1 \leq m,n \leq N \\ m \neq n}} e^{-\beta T_k^{m,n}} s_{m,n} \right] \quad (3)$$

and $\mathbf{x} = (x_1, \dots, x_N)$, where $\mathbf{E}_{\mathbf{x},z}$ is the conditional expectation operator given that the system is in state (\mathbf{x}, z) at time 0.

Observe that the r.h.s. of (2) does not depend on the position of the server since we are only considering holding costs. In the following, the subscript z will be dropped in $V_{\mathbf{x},z}^h(\beta)$.

Our main objective in this paper is to determine $V_{\mathbf{x},z}(\beta)$ in explicit form. We will also be interested in the explicit computation of the average holding cost $g_{\mathbf{x}}^h$ and switching cost $g_{\mathbf{x},z}^s$ defined as

$$g_{\mathbf{x}}^h := \lim_{\beta \rightarrow 0} \beta V_{\mathbf{x}}^h(\beta) \quad (4)$$

$$g_{\mathbf{x},z}^s := \lim_{\beta \rightarrow 0} \beta V_{\mathbf{x},z}^s(\beta) \quad (5)$$

as well as on the explicit derivation of the bias function $B_{\mathbf{x},z;\mathbf{x}_0,z_0}$, given by

$$B_{\mathbf{x},z;\mathbf{x}_0,z_0} := \lim_{\beta \rightarrow 0} (V_{\mathbf{x},z}(\beta) - V_{\mathbf{x}_0,z_0}(\beta)) \quad (6)$$

where (\mathbf{x}_0, z_0) is an arbitrary reference state.

3 Preliminaries

In this section we establish some basic results that we will use later on and that are related to LSTs of busy periods.

A word on the notation in use: $\mathbb{N} = \{0, 1, \dots\}$ will denote the set of all nonnegative integers; \mathbf{e}_i^n will denote the unit vector of dimension n with all components equal to 0 except the i -th that is equal to 1; when $n = N$ we will omit the superscript N in \mathbf{e}_i^N .

For any vector $\mathbf{x} = (x_1, \dots, x_N) \in \mathbb{N}^N$, $\mathbf{x}^n = (x_1, \dots, x_n)$ for $n = 1, 2, \dots, N$. The vector $\mathbf{x} = (x_1, \dots, x_N) \in \mathbb{N}^N$ such that $x_1 = \dots = x_{i-1} = 0$ and $x_i > 0$ will be denoted as $\mathbf{x}_{[i]}$. In particular, the last two definitions imply that $\mathbf{x}_{[i]}^n = (0, \dots, 0, x_i, \dots, x_n) \in \mathbb{N}^n$ if $i \leq n \leq N$ and $\mathbf{x}_{[i]}^n = (0, \dots, 0) \in \mathbb{N}^n$ if $1 \leq n < i$.

Central to the analysis of priority queues are the concepts of busy periods and service completion times ([7]). For the purpose of our analysis let us introduce:

- (1) The *busy period of an arbitrary customer of class* $1, 2, \dots, n$ (with LST $\tilde{\gamma}_n(s)$) is the time needed to clear the system of all customers of class $1, 2, \dots, n$ given that there are no customers of class $1, 2, \dots, n$ in the system just before time $t = 0$ and that a customer of class less than $n + 1$ enters the queue at time $t = 0$ (this customer is of class i with probability proportional to the arrival rate);
- (2) The *service completion time of a customer of class* n (with LST $C_n(s)$) is defined as the time that elapses between the first entry of a customer of class n in the server and its departure from the system. Note that under the enforced service policy when a customer of class n enters the server the queue does not contain any customers of class $1, 2, \dots, n - 1$;
- (3) The *busy period of customers of class* n (with LST $\gamma_n(s)$) is the time needed to empty the system of customers of class n given that there are no customers of class $1, 2, \dots, n$ in the system just before time $t = 0$ and that a customer of class n enters the queue at time $t = 0$;
- (4) The *busy period of customers of class* $1, 2, \dots, n$ with initial workload $\mathbf{x}^n = (x_1, \dots, x_n)$ (with LST $\tau_{\mathbf{x}^n}(s)$) is the time needed to empty the system of all customers of class $1, 2, \dots, n$ given that there are x_i customers of class i , $i = 1, 2, \dots, n$, in the system at time 0; by convention $\tau_{\mathbf{x}^n}(\beta) = 1$ if $n = 0$.

The following lemma reports two basic results on the M/G/1 queue that will be used to compute the four LSTs introduced above.

Lemma 3.1 *Consider an M/G/1 queue with arrival intensity $\lambda > 0$ and with LST of the service times $S(s)$. Let $\Delta(s)$ be the LST of the length of a busy period (defined as the time needed to empty the queue given that a customer arrives in an empty system at time $t = 0$). Assume that at time $t = 0$ there is an initial waiting time with LST $V(s)$.*

(i) *The LST $\Omega(s)$ of the time needed to empty the system is given by*

$$\Omega(s) = V(s + \lambda(1 - \Delta(s))) \in (0, 1), \quad s > 0. \quad (7)$$

(ii) *Assume that $V = S$, i.e., the initial workload corresponds to the service time of one customer. Then $\Omega(s) = \Delta(s)$, and $\Delta(s)$ is the single root in $(0, 1)$ of the equation*

$$z = S(s + \lambda(1 - z)), \quad s > 0. \quad (8)$$

Proof. For the proof of both equations, in case $\Re(s) \geq 0$ (with $\Re(s)$ the real part of $s \in \mathbb{C}$), see [16, pp. 58–59, Theorem 3]. Δ is the LST of some proper or improper distribution,

therefore $\Delta(s) \in (0, 1)$ if $s > 0$, and thus also $\Omega(s) \in (0, 1)$. The proof that (8) has a single root in $|z| < 1$ can be found in [16, pp. 47–48, Lemma 1]. This combined with the facts that $z < S(s + \lambda(1 - z))$ if $z = 0$, $z > S(s + \lambda(1 - z))$ if $z = 1$, and the continuity of $z - S(s + \lambda(1 - z))$ gives that $z \in (0, 1)$. ■

We are now in a position to address the computation of $\tilde{\gamma}_n(s)$, $C_n(s)$, $\gamma_n(s)$ and $\tau_{\mathbf{x}^n}(s)$. This is done in Lemma 3.2. Throughout $f \circ g(s)$ stands for $f(g(s))$ for any mappings f and g .

Define $\Lambda_n := \sum_{i=1}^n \lambda_i$, $n = 1, 2, \dots, N$.

Lemma 3.2 *Let $s > 0$ and $n = 1, 2, \dots, N$. Then,*

- $\tilde{\gamma}_n(s)$ is the single root in $(0, 1)$ of the equation

$$z = \tilde{S}_n(s + \Lambda_n(1 - z)) \quad (9)$$

with $\tilde{S}_n(s) := (1/\Lambda_n) \sum_{i=1}^n \lambda_i S_i(s)$;

- $C_n(s)$ is given by

$$C_n(s) = S_n(s + \Lambda_{n-1}(1 - \tilde{\gamma}_{n-1}(s))); \quad (10)$$

- $\gamma_n(s)$ is the single root in $(0, 1)$ of the equation

$$z = C_n(s + \lambda_n(1 - z)); \quad (11)$$

- for all $\mathbf{x} = (x_1, \dots, x_n) \in \mathbb{N}^N$,

$$\tau_{\mathbf{x}^n}(s) = \tau_{\mathbf{x}^{n-1}}(s + \lambda_n(1 - \gamma_n(s))) \gamma_n(s)^{x_n} \quad (12)$$

$$= \prod_{i=1}^n \gamma_i(g_{i+1} \circ \dots \circ g_n(s))^{x_i}, \quad (13)$$

with $g_n(s) := s + \lambda_n(1 - \gamma_n(s))$ for $n \geq 1$, and, by convention, $\gamma_n(g_{n+1} \circ \dots \circ g_n(s)) = \gamma_n(s)$.

Proof.

Proof of (9). Under the enforced statistical assumptions (Poisson arrivals, independent service times) and service policy, the length of a busy period for customers of class $1, 2, \dots, n$ in the original system is the same as the length of a busy period in a single-class M/G/1 queue with arrival intensity Λ_n and with LST of the service times $\tilde{S}_n(s)$. Hence, by the second statement in Lemma 3.1, $z = \tilde{\gamma}_n(s)$ is the only root in $(0, 1)$ of equation (9).

Proof of (10). The service completion time for a customer of class n can be seen as the time needed to clear the system of all customers of class $1, 2, \dots, n - 1$ given that at time $t = 0$

there was an initial workload (corresponding to the service time of a customer of class n) with LST $S_n(s)$. By applying the first statement in Lemma 3.1 we see that $C_n(s)$ is given by the r.h.s. of (10). (This result is not new and can be found, for instance, in [7, p. 109, formula (7.10)].)

Proof of (11). From the definition of a busy period of a customer of class n it is easily seen that $\gamma_n(s)$ is equal to the LST of the length of a busy period in an M/G/1 queue with arrival intensity λ_n and with the LST of the service completion time of a customers of class n , $C_n(s)$. Hence, by the second statement of Lemma 3.1, $\gamma_n(s)$ is the only root in $(0, 1)$ of equation (11).

Proof of (12)-(13). The identity (13) is easily derived from (12), so that we only have to prove (12). Assume that there are x_i customers of class $i = 1, 2, \dots, n$ at time $t = 0$. Let T be the first time after time $t = 0$ when there are no more customers of class $1, 2, \dots, n - 1$ in the system and where there are exactly x_n customers of class n . Let $T(s) = \mathbf{E}[\exp(-sT)]$ be the LST of the random variable T . Clearly, $T(s)$ is equal to the LST of the length of a busy period in an M/G/1 queue with arrival rate λ_n , with LST of the service times $C_n(s)$, and with an initial waiting with LST $\tau_{\mathbf{x}^{n-1}}(s)$. Hence, by lemma 3.1, we have that $T(s) = \tau_{\mathbf{x}^{n-1}}(s + \lambda_n(1 - \gamma_n(s)))$. Define now T_2 as the first time after time T when there are no more customers of class $1, 2, \dots, n$ in the system. The LST of $T_2 - T$ is clearly given by $\gamma_n(s)^{x_n}$. Therefore, since T and $T_2 - T$ are independent random variables under the statistical assumptions placed on the model, the LST of T_2 , which is by definition equal to $\tau_{\mathbf{x}^n}(s)$, is given by the product of $\tau_{\mathbf{x}^{n-1}}(s + \lambda_n(1 - \gamma_n(s)))$ and $\gamma_n(s)^{x_n}$, which concludes the proof. ■

The third and last technical lemma establishes a useful relationship between the LSTs $\tilde{\gamma}_n$ and γ_i , $i = 1, 2, \dots, n$.

Lemma 3.3 For $s > 0$, $n = 1, 2, \dots, N$,

$$\tilde{\gamma}_n(s) = \sum_{j=1}^n \frac{\lambda_j}{\Lambda_n} \gamma_j(g_{j+1} \circ \dots \circ g_n(s)). \quad (14)$$

As a consequence, if class- n customers have exponentially distributed service times with mean $1/\mu_n$, then

$$\gamma_n(s) = \frac{\mu_n}{\mu_n + g_n(s) + \Lambda_{n-1} (1 - \tilde{\gamma}_{n-1}(g_n(s)))} \quad (15)$$

$$= \frac{\mu_n}{\mu_n + s + \sum_{j=1}^n \lambda_j (1 - \gamma_j(g_{j+1} \circ \dots \circ g_n(s)))}. \quad (16)$$

Proof. We have

$$\tilde{\gamma}_n(s) = \sum_{j=1}^n \frac{\lambda_j}{\Lambda_n} \tau_{\mathbf{e}_j^n}(s)$$

by definition of $\tau_{\mathbf{x}^n}(s)$. By replacing $\tau_{\mathbf{e}_j^n}(s)$ by its value $\gamma_j(g_{j+1} \circ \dots \circ g_n(s))$ given in Lemma 3.2 we find (14). The identity (15) follows from (10)-(11), while (16) is a direct consequence of (14). \blacksquare

4 Holding costs

Our objective in this section is to compute the expected overall discounted holding cost $V_{\mathbf{x}}^h(\beta)$ defined in (2).

Define $H_{\mathbf{x}^n}^n(\beta)$ as the expected discounted total number of customers of class n in $[0, \infty)$ given that x_i class- i customers ($i = 1, 2, \dots, n$) were in the system at time 0, namely,

$$H_{\mathbf{x}^n}^n(\beta) = \mathbb{E}_{\mathbf{x}^n} \left[\int_0^\infty e^{-\beta t} X_n(t) dt \right]. \quad (17)$$

When $x_1 = \dots = x_n = 0$ then $H_{\mathbf{x}^n}^n(\beta)$ is denoted by $H_0^n(\beta)$.

Since under the preemptive resume priority rule customers of class n are not affected by the presence of customers of class $n + 1, \dots, N$, we deduce from (2) and (17) that

$$V^h(\mathbf{x}) = \sum_{n=1}^N c_n H_{\mathbf{x}^n}^n(\beta) \quad (18)$$

for all $\mathbf{x} \in \mathbb{N}^N$, $z \in \{1, 2, \dots, N\}$.

Below, we first derive a closed-form expression for $H_{\mathbf{x}^n}^n(\beta)$ through an intuitive argument. A rigorous proof follows in Theorem 4.1.

Consider a state \mathbf{x} with $\mathbf{x}^n \neq 0$. If there were always class- n customers in the system then the expected discounted total number of customers of class n in $[0, \infty)$ would be

$$\frac{x_n}{\beta} + \frac{\lambda_n}{\beta^2} - \tau_{\mathbf{x}^{n-1}}(\beta) \frac{C_n(\beta)}{\beta(1 - C_n(\beta))}. \quad (19)$$

The first term corresponds to the holding costs of customers initially present; the second corresponds to the newly arriving class- n customers; the last term corresponds to class- n departures. Before starting service on a class- n customer the first $n - 1$ queues should be empty, explaining the $\tau_{\mathbf{x}^{n-1}}(\beta)$ -term in the expression.

Equation (19) is only valid if there are always class- n customers in the system. However, after a time with LST $\tau_{\mathbf{x}^n}(\beta)$ there are no longer customers of class n and we need to correct (19) for that. The correction term is necessarily of the form

$$\tau_{\mathbf{x}^n}(\beta) \left(H_0^n(\beta) - \left(\frac{\lambda_n}{\beta^2} - \frac{C_n(\beta)}{\beta(1 - C_n(\beta))} \right) \right).$$

Hence, we postulate that

$$H_{\mathbf{x}^n}^n(\beta) = \frac{x_n}{\beta} + \frac{\lambda_n}{\beta^2} - \tau_{\mathbf{x}^{n-1}}(\beta) \frac{C_n(\beta)}{\beta(1-C_n(\beta))} + \tau_{\mathbf{x}^n}(\beta) \left(H_0^n(\beta) - \left(\frac{\lambda_n}{\beta^2} - \frac{C_n(\beta)}{\beta(1-C_n(\beta))} \right) \right) \quad \text{for } n = 1, 2, \dots, N. \quad (20)$$

The unknown function $H_0^n(\beta)$ can be determined by conditioning on the first event to occur after time $t = 0$ in queue $1, 2, \dots, n$. This gives

$$H_0^n(\beta) = \frac{1}{\Lambda_n + \beta} \sum_{i=1}^n \lambda_i H_{\mathbf{e}_i}^n(\beta) \quad \text{for } n = 1, 2, \dots, N. \quad (21)$$

Straightforward algebra using (20) and (13) give

$$H_0^n(\beta) = \left(\beta + \sum_{i=1}^{n-1} \lambda_i (1 - \gamma_i(g_{i+1} \circ \dots \circ g_n(\beta))) + \lambda_n (1 - \gamma_n(\beta)) \right)^{-1} \times \left\{ \frac{\Lambda_n \lambda_n}{\beta^2} + \frac{\lambda_n}{\beta} - \Gamma_n(\beta) \left[\sum_{i=1}^{n-2} \lambda_i \gamma_i(g_{i+1} \circ \dots \circ g_{n-1}(\beta)) + \lambda_{n-1} \gamma_{n-1}(\beta) + \lambda_n \right] - \left(\frac{\lambda_n}{\beta^2} - \Gamma_n(\beta) \right) \left(\sum_{i=1}^{n-1} \lambda_i \gamma_i(g_{i+1} \circ \dots \circ g_n(\beta)) + \lambda_n \gamma_n(\beta) \right) \right\}$$

with

$$\Gamma_n(\beta) = \frac{C_n(\beta)}{\beta(1-C_n(\beta))}, \quad n = 1, 2, \dots, N. \quad (22)$$

With the help of Lemma 3.3 $H_0^n(\beta)$ rewrites as

$$H_0^n(\beta) = (\beta + \Lambda_n (1 - \tilde{\gamma}_n(\beta)))^{-1} \times \left\{ \frac{\Lambda_n \lambda_n}{\beta^2} + \frac{\lambda_n}{\beta} - \Gamma_n(\beta) [\Lambda_{n-1} \tilde{\gamma}_{n-1}(\beta) + \lambda_n] - \left(\frac{\lambda_n}{\beta^2} - \Gamma_n(\beta) \right) \Lambda_n \tilde{\gamma}_n(\beta) \right\} = \frac{\lambda_n}{\beta^2} - \Gamma_n(\beta) + \Gamma_n(\beta) \frac{\beta + \Lambda_{n-1} (1 - \tilde{\gamma}_{n-1}(\beta))}{\beta + \Lambda_n (1 - \tilde{\gamma}_n(\beta))}.$$

It remains to plug this value of $H_0^n(\beta)$ into (20) to finally find

$$H_{\mathbf{x}^n}^n(\beta) = \frac{x_n}{\beta} + \frac{\lambda_n}{\beta^2} - \tau_{\mathbf{x}^{n-1}}(\beta) \Gamma_n(\beta) + \tau_{\mathbf{x}^n}(\beta) \Gamma_n(\beta) \Psi_n(\beta) \quad (23)$$

for $n = 1, 2, \dots, N$ and $\mathbf{x}^n \in \mathbb{N}^n$, with

$$\Psi_n(\beta) = \frac{\beta + \Lambda_{n-1} (1 - \tilde{\gamma}_{n-1}(\beta))}{\beta + \Lambda_n (1 - \tilde{\gamma}_n(\beta))}. \quad (24)$$

We now provide a rigorous proof of (23)-(24). Until Theorem 4.1 n is held fixed in $\{1, \dots, N\}$.

We introduce some further notation: Let $\hat{\tau}_{\mathbf{x}^{n-1}}$ be the time needed to clear the system of all customers of class $1, 2, \dots, n-1$ given that there are x_i customers of class $i = 1, 2, \dots, n-1$ in the system at time 0. Also define \hat{C}_n as the service completion time for customers of class n . Recall (see the beginning of Section 3) that the LSTs of $\hat{\tau}_{\mathbf{x}^{n-1}}$ and \hat{C}_n are given by $\tau_{\mathbf{x}^{n-1}}(\beta)$ and $C_n(\beta)$, respectively.

Define the mapping $w : \mathbb{N}^n \rightarrow \mathbb{R}$ as $w(\mathbf{x}^n) = \max\{x_n, 1\}$. For each mapping $u : \mathbb{N}^n \rightarrow \mathbb{R}$, define the norm $\|u\|$ by

$$\|u\| = \sup_{\mathbf{x}^n \in \mathbb{N}^n} \{u(\mathbf{x}^n)/w(\mathbf{x}^n)\}$$

and let \mathcal{B} be the Banach space of all such u for which $\|u\| < \infty$.

We begin with the following lemma that establishes a recursive scheme to be satisfied by $H_{\mathbf{x}^n}^n(\beta)$.

Lemma 4.1 *For all $\mathbf{x}^n = (x_1, \dots, x_n) \in \mathbb{N}^n$, $H_{\mathbf{x}^n}^n(\beta)$ is the unique solution in \mathcal{B} of*

$$\begin{aligned} F_{\mathbf{x}^n} &= \frac{1 - \tau_{\mathbf{x}^{n-1}}(\beta)}{\beta} \left(x_n + \frac{\lambda_n}{\beta} \right) - \frac{\lambda_n}{\beta} \int_0^\infty e^{-\beta y} y dP(\hat{\tau}_{\mathbf{x}^{n-1}} < y) \\ &+ \int_{y=0}^\infty e^{-(\beta+\lambda_n)y} \sum_{k=0}^\infty \frac{(\lambda_n y)^k}{k!} F_{(x_n+k)\mathbf{e}_n^n} dP(\hat{\tau}_{\mathbf{x}^{n-1}} < y) \end{aligned} \quad (25)$$

if $(x_1, \dots, x_{n-1}) \neq 0$, $x_n \geq 0$,

$$\begin{aligned} F_{\mathbf{x}^n} &= \frac{1 - C_n(\beta)}{\beta} \left(x_n + \frac{\lambda_n}{\beta} \right) - \frac{\lambda_n}{\beta} \int_0^\infty e^{-\beta y} y dP(\hat{C}_n < y) \\ &+ \int_{y=0}^\infty e^{-(\beta+\lambda_n)y} \sum_{k=0}^\infty \frac{(\lambda_n y)^k}{k!} F_{(x_{n-1}+k)\mathbf{e}_n^n} dP(\hat{C}_n < y) \end{aligned} \quad (26)$$

if $(x_1, \dots, x_{n-1}) = 0$, $x_n > 0$,

$$F_0 = \frac{1}{\Lambda_n + \beta} \sum_{i=1}^n \lambda_i F_{\mathbf{e}_i^n}. \quad (27)$$

Proof. First we show that $H_{\mathbf{x}^n}^n(\beta)$, as given in (23), satisfies (25)-(27).

Fix $\mathbf{x}^n = (x_1, \dots, x_n) \in \mathbb{N}^n$ and recall that $\mathbf{e}_n^n = (0, \dots, 0, 1) \in \mathbb{N}^n$.

We first assume that $(x_1, \dots, x_{n-1}) \neq 0$. We have, cf. (17),

$$H_{\mathbf{x}^n}^n(\beta) = \mathbb{E}_{\mathbf{x}^n} \left[\int_0^{\hat{\tau}_{\mathbf{x}^{n-1}}} e^{-\beta t} X_n(t) dt \right] + \mathbb{E}_{\mathbf{x}^n} \left[\int_{\hat{\tau}_{\mathbf{x}^{n-1}}}^\infty e^{-\beta t} X_n(t) dt \right]. \quad (28)$$

We will consider separately both integrals in the r.h.s. of (28).

Conditioning on $\hat{\tau}_{\mathbf{x}^{n-1}}$ in the first integral yields

$$\begin{aligned}
& \mathbb{E}_{\mathbf{x}^n} \left[\int_0^{\hat{\tau}_{\mathbf{x}^{n-1}}} e^{-\beta t} X_n(t) dt \right] \\
&= \int_{y=0}^{\infty} \int_{t=0}^y e^{-\beta t} \mathbb{E}_{\mathbf{x}^n} [X_n(t) | \hat{\tau}_{\mathbf{x}^{n-1}} = y] dt dP(\hat{\tau}_{\mathbf{x}^{n-1}} < y) \\
&= \int_{y=0}^{\infty} \int_{t=0}^y e^{-\beta t} (x_n + \lambda_n t) dt dP(\hat{\tau}_{\mathbf{x}^{n-1}} < y) \\
&= \left(\frac{1 - \tau_{\mathbf{x}^{n-1}}(\beta)}{\beta} \right) \left(x_n + \frac{\lambda_n}{\beta} \right) - \frac{\lambda_n}{\beta} \int_0^{\infty} e^{-\beta y} y dP(\hat{\tau}_{\mathbf{x}^{n-1}} < y) \tag{29}
\end{aligned}$$

after elementary algebra.

Let us now focus on the second integral in the r.h.s. of (28). Let $A_n(t)$ be the number of customers of class n that have arrived in $[0, t)$. Conditioning on $\hat{\tau}_{\mathbf{x}^{n-1}}$, then on $A_n(\hat{\tau}_{\mathbf{x}^{n-1}})$, yields

$$\begin{aligned}
& \mathbb{E}_{\mathbf{x}^n} \left[\int_{\hat{\tau}_{\mathbf{x}^{n-1}}}^{\infty} e^{-\beta t} X_n(t) dt \right] \\
&= \int_{y=0}^{\infty} e^{-\lambda_n y} \sum_{k \geq 0} \frac{(\lambda_n y)^k}{k!} \int_{t=y}^{\infty} e^{-\beta t} \mathbb{E}_{\mathbf{x}^n} [X_n(t) | \hat{\tau}_{\mathbf{x}^{n-1}} = y, A_n(y) = k] dt dP(\hat{\tau}_{\mathbf{x}^{n-1}} < y) \\
&= \int_{y=0}^{\infty} e^{-(\beta + \lambda_n)y} \sum_{k \geq 0} \frac{(\lambda_n y)^k}{k!} \int_{u=0}^{\infty} e^{-\beta u} \mathbb{E}_{\mathbf{x}^n} [X_n(u + y) | \hat{\tau}_{\mathbf{x}^{n-1}} = y, A_n(y) = k] du \\
&\quad \times dP(\hat{\tau}_{\mathbf{x}^{n-1}} < y) \\
&= \int_{y=0}^{\infty} e^{-(\beta + \lambda_n)y} \sum_{k \geq 0} \frac{(\lambda_n y)^k}{k!} H_{(x_n+k)}^n \mathbf{e}_n^n(\beta) dP(\hat{\tau}_{\mathbf{x}^{n-1}} < y). \tag{30}
\end{aligned}$$

Summing up the r.h.s. of (29) and (30) gives (25).

The proof of (30) is analogous to that of (29) and is left to the reader (Hint: condition on \hat{C}_n). Finally, (27) has already been derived in (21). This concludes the proof that $H_{\mathbf{x}^n}^n$ solves equations (25)-(27).

We now show that $H_{\mathbf{x}^n}^n \in \mathcal{B}$. Let π be any policy that *never* serves customers of class n and let $\hat{H}_{\mathbf{x}^n}^n$ be the expected discounted total number of customers of class n in $[0, \infty)$ given that there are x_i customers of class $i = 1, 2, \dots, n$ at time 0. Clearly,

$$H_{\mathbf{x}^n}^n \leq \hat{H}_{\mathbf{x}^n}^n = \frac{x_n}{\beta} + \frac{\lambda_n}{\beta^2}$$

which shows that $H_{\mathbf{x}^n}^n \in \mathcal{B}$ since $\hat{H}_{\mathbf{x}^n}^n \in \mathcal{B}$.

It remains to show that equations (25)-(27) have a unique solution in \mathcal{B} .

Let B be the set of functions $u : \mathbb{N}^n \rightarrow \mathbb{R}$. Equations (25)-(27) define an operator $T : B \rightarrow B$ such that $F = TF$. This operator is actually the dynamic programming operator associated with the μ -rule and is such that $TF = C + QF$, with C the direct cost and Q the (defective) transition matrix.

It is readily seen that $\|C\| < \infty$. Next, we derive a bound for $\sum_{\mathbf{y}^n \in \mathbb{N}^n} Q(\mathbf{x}^n, \mathbf{y}^n)w(\mathbf{y}^n)$, where $w \in B$ has been defined earlier.

First define

$$\alpha = \max \left\{ \max_{\mathbf{x}^{n-1} \neq 0} \{\tau_{\mathbf{x}^{n-1}}(\beta)\}, C_n(\beta), \frac{\Lambda_n}{\Lambda_n + \beta} \right\}.$$

Note that $\alpha < 1$. Define also

$$b = \max \left\{ \max_{\mathbf{x}^{n-1} \neq 0} \left\{ \frac{\lambda_n}{\alpha} \int_0^\infty e^{-\beta y} y dP(\hat{\tau}_{\mathbf{x}^{n-1}} < y) \right\}, \frac{\lambda_n}{\alpha} \int_0^\infty e^{-\beta y} y dP(\hat{C}_n < y), \frac{\Lambda_n}{\alpha(\Lambda_n + \beta)} \right\}.$$

It is easily seen that $\sum_{\mathbf{y}^n} Q(\mathbf{x}^n, \mathbf{y}^n)w(\mathbf{y}^n) \leq \alpha(w(\mathbf{x}^n) + b)$. From this we can derive, following the arguments in [10, p. 1228] that T is closed under the $\|\cdot\|$ -norm, and that T^J for some well chosen J is a contraction.

We start with the closedness property. For an arbitrary mapping $F \in \mathcal{B}$ and $\mathbf{x}^{n-1} \neq 0$ we have (see (25)-(27))

$$\begin{aligned} |TF_{\mathbf{x}^n}| &\leq \frac{x_n}{\beta} + \frac{\lambda_n}{\beta^2} + \sum_{\mathbf{y}^n \in \mathbb{N}^n} Q(\mathbf{x}^n, \mathbf{y}^n) \frac{|F_{\mathbf{y}^n}|}{w(\mathbf{y}^n)} w(\mathbf{y}^n) \\ &\leq \frac{w(\mathbf{x}^n)}{\beta} + \frac{\lambda_n}{\beta^2} + \|F\| \sum_{\mathbf{y}^n} Q(\mathbf{x}^n, \mathbf{y}^n)w(\mathbf{y}^n) \\ &\leq \frac{w(\mathbf{x}^n)}{\beta} + \frac{\lambda_n}{\beta^2} + \alpha\|F\|(w(\mathbf{x}^n) + b). \end{aligned}$$

A similar analysis holds for \mathbf{x} such that $\mathbf{x}^{n-1} = 0$ and $x_n > 0$. For $\mathbf{x}^n = 0$ we have $|TF_0| \leq (\Lambda_n + \beta)^{-1} \sum_i |F_{\mathbf{e}_i^n}| \leq \Lambda_n(\Lambda_n + \beta)^{-1}\|F\|$. Dividing all three equations by $w(\mathbf{x}^n)$ and taking the supremum shows that TF is $\|\cdot\|$ -bounded.

Finally, let show that T^J is a contraction for some J large enough. We have for $F, \tilde{F} \in B$

$$\begin{aligned} |TF_{\mathbf{x}^n} - T\tilde{F}_{\mathbf{x}^n}| &= |QF_{\mathbf{x}^n} - Q\tilde{F}_{\mathbf{x}^n}| \\ &\leq \sum_{\mathbf{y}^n \in \mathbb{N}^n} Q(\mathbf{x}^n, \mathbf{y}^n) \frac{|F_{\mathbf{y}^n} - \tilde{F}_{\mathbf{y}^n}|}{w(\mathbf{y}^n)} w(\mathbf{y}^n) \\ &\leq \alpha\|F - \tilde{F}\|(w(\mathbf{x}^n) + b). \end{aligned}$$

Assume that

$$|T^k F_{\mathbf{x}^n} - T^k \tilde{F}_{\mathbf{x}^n}| \leq \alpha^k \|F - \tilde{F}\|(w(\mathbf{x}^m) + kb) \quad (31)$$

for $k = 1, 2, \dots, m$ and let us show that this inequality holds for $k = m + 1$.

We have,

$$\begin{aligned}
|T^{m+1}F_{\mathbf{x}^n} - T^{m+1}\tilde{F}_{\mathbf{x}^n}| &\leq \sum_{\mathbf{y}^n \in \mathbb{N}^n} Q(\mathbf{x}^n, \mathbf{y}^n) |T^m F_{\mathbf{y}^n} - T^m \tilde{F}_{\mathbf{y}^n}| \\
&\leq \alpha^m \|F - \tilde{F}\| \sum_{\mathbf{y}^n \in \mathbb{N}^n} Q(\mathbf{x}^n, \mathbf{y}^n) (w(\mathbf{x}^n) + mb) \\
&\leq \alpha^{m+1} \|F - \tilde{F}\| (w(\mathbf{x}^n) + (m+1)b).
\end{aligned}$$

This proves the induction step. Now take J such that $\alpha^J (w(\mathbf{x}^n) + Jb) < 1$ (such a J exists since $\alpha < 1$). Then $\|T^J F - T^J \tilde{F}\| < \|F - \tilde{F}\|$, and therefore T^J is a contraction. From this it follows (see e.g. [10]) that a $\|\cdot\|$ -bounded solution to $F = TF$ exists and that it is unique.

In conclusion, we have shown that $H_{\mathbf{x}^n}^n$ is the unique solution in \mathcal{B} of the equation $H_{\mathbf{x}^n}^n = TH_{\mathbf{x}^n}^n$, which concludes the proof. \blacksquare

We are now ready to prove the main result of this section.

Theorem 4.1 (Expected total discounted holding cost)

For $n = 1, 2, \dots, N$, $\mathbf{x} = (x_1, \dots, x_N)$ and $\mathbf{x}^n = (x_1, \dots, x_n)$, the expected total discounting holding cost is given by

$$V_{\mathbf{x}}^h(\beta) = \sum_{i=1}^N c_n H_{\mathbf{x}^n}^n(\beta) \quad (32)$$

where $H_{\mathbf{x}^n}^n(\beta)$ is given in (23).

Proof. Because of (18) it suffices to prove that (23) holds. The proof will consist in checking that the r.h.s. of (23) satisfies the set of equations (25)-(27) in Lemma 4.1.

We start with (25): substituting $H_{(x_n+\beta)\mathbf{e}_n}^n(\beta)$ by the value given in (23) yields

$$\begin{aligned}
&\int_{y=0}^{\infty} e^{-(\beta+\lambda_n)y} \sum_{k \geq 0} \frac{(\lambda_n y)^k}{k!} H_{(x_n+k)\mathbf{e}_n}^n(\beta) dP(\hat{\tau}_{\mathbf{x}^{n-1}} < y) \\
&= \int_{y=0}^{\infty} e^{-(\beta+\lambda_n)y} \sum_{k \geq 0} \frac{(\lambda_n y)^k}{k!} \left[\frac{x_n + k}{\beta} + \frac{\lambda_n}{\beta^2} - \Gamma_n(\beta) + \gamma(\beta)^{x_n+\beta} \Gamma_n \Psi_n(\beta) \right] \\
&\quad \times dP(\hat{\tau}_{\mathbf{x}^{n-1}} < y) \\
&= \left(\frac{x_n}{\beta} + \frac{\lambda_n}{\beta^2} \right) \tau_{\mathbf{x}^{n-1}}(\beta) - \Gamma_n(\beta) \tau_{\mathbf{x}^{n-1}}(\beta) + \Gamma_n(\beta) \Psi_n(\beta) \gamma_n(\beta)^{x_n} \tau_{\mathbf{x}^{n-1}}(\beta + \lambda_n(1 - \gamma_n(\beta))) \\
&\quad + \frac{\lambda_n}{\beta} \int_0^{\infty} e^{-\beta y} y dP(\hat{\tau}_{\mathbf{x}^{n-1}} < y)
\end{aligned}$$

$$\begin{aligned}
&= \left(\frac{x_n}{\beta} + \frac{\lambda_n}{\beta^2} \right) \tau_{\mathbf{x}^{n-1}}(\beta) - \Gamma_n(\beta) \tau_{\mathbf{x}^{n-1}}(\beta) + \Gamma_n(\beta) \Psi_n(\beta) \tau_{\mathbf{x}^n}(\beta) \\
&\quad + \frac{\lambda_n}{\beta} \int_0^\infty e^{-\beta y} y dP(\hat{\tau}_{\mathbf{x}^{n-1}} < y)
\end{aligned} \tag{33}$$

where we have used (12) to derive the last equation.

Hence, cf. (29) and (33),

$$H_{\mathbf{x}^n}^n(\beta) = \frac{x_n}{\beta} + \frac{\lambda_n}{\beta^2} - \tau_{\mathbf{x}^{n-1}}(\beta) \Gamma_n(\beta) + \tau_{\mathbf{x}^n}(\beta) \Gamma_n(\beta) \Psi_n(\beta)$$

which is the value found in the r.h.s. of (23).

We now consider (26). Using (23) we find

$$\begin{aligned}
&\int_{y=0}^\infty e^{-(\beta+\lambda_n)y} \sum_{k \geq 0} \frac{(\lambda_n y)^k}{k!} H_{(x_n-1+k)\mathbf{e}_n}^n(\beta) dP(\hat{C}_n < y) \\
&= \int_{y=0}^\infty e^{-(\beta+\lambda_n)y} \sum_{k \geq 0} \frac{(\lambda_n y)^k}{k!} \left[\frac{x_n - 1 + k}{\beta} + \frac{\lambda_n}{\beta^2} - \Gamma_n(\beta) + \gamma_n(\beta)^{x_n-1+k} \Gamma_n(\beta) \Psi_n(\beta) \right] \\
&\quad \times dP(\hat{C}_n < y) \\
&= \left(\frac{x_n - 1}{\beta} + \frac{\lambda_n}{\beta^2} \right) C_n(\beta) - \Gamma_n(\beta) \tilde{C}_n(\beta) - \gamma_n(\beta)^{x_n-1} \Psi(\beta) \Gamma_n(\beta) C_n(\beta + \lambda(1 - \gamma_n(\beta))) \\
&\quad - \frac{\lambda_n}{\beta} \int_0^\infty e^{-\beta y} y dP(\hat{C}_n < y) \\
&= \left(\frac{x_n - 1}{\beta} + \frac{\lambda_n}{\beta^2} \right) C_n(\beta) - \Gamma_n(\beta) C_n(\beta) - \gamma_n(\beta)^{x_n} \Psi(\beta) \Gamma_n(\beta) \\
&\quad - \frac{\lambda_n}{\beta} \int_0^\infty e^{-\beta y} y dP(\hat{C}_n < y)
\end{aligned} \tag{34}$$

$$= \left(\frac{x_n}{\beta} + \frac{\lambda_n}{\beta^2} \right) C_n(\beta) - \Gamma_n(\beta) - \gamma_n(\beta)^{x_n} \Psi(\beta) \Gamma_n(\beta) - \frac{\lambda_n}{\beta} \int_0^\infty e^{-\beta y} y dP(\hat{C}_n < y) \tag{35}$$

where (34) follows from (11) and (35) is a consequence of the definition of $\Gamma_n(\beta)$ in (22). Plugging (35) in (26) we get

$$H_{\mathbf{x}^n}^n(\beta) = \frac{x_n}{\beta} + \frac{\lambda_n}{\beta^2} - \Gamma_n(\beta) - \gamma_n(\beta)^{x_n} \Psi(\beta) \Gamma_n(\beta)$$

which is the value found in the r.h.s. of (23) for $x_1 = \dots = x_{n-1} = 0$.

Finally, the fact that (27) is satisfied by the r.h.s. of (23) follows from the derivation of the latter expression. \blacksquare

5 Switching costs

In this section we derive a closed-form expression for the expected total discounted switching cost $V_{\mathbf{x},z}^s(\beta)$ defined in (3), in the case where the service times are *exponentially* distributed.

Recall the definitions of vectors \mathbf{x}^n , $\mathbf{x}_{[i]}$ and $\mathbf{x}_{[i]}^n$ introduced at the beginning of Section 3.

Theorem 5.1 (Expected total discounted switching cost)

The expected total discounted switching cost is given by

$$\begin{aligned} V_{\mathbf{x},1}^s(\beta) &= \sum_{j=1}^N r(j)\tau_{\mathbf{x}^j}(\beta) \\ &\quad + \sum_{1 < k \leq l \leq N} r(k,l)\tau_{\mathbf{x}^{k-1}}(\beta + \lambda_k + \cdots + \lambda_l) \mathbf{1}(x_k = \cdots = x_l = 0) \end{aligned} \quad (36)$$

$$V_{\mathbf{x}_{[i]},z}^s(\beta) = s_{z,i} - s_{1,i} + V_{\mathbf{x}_{[i]},1}^s(\beta), \quad i, z = 1, 2, \dots, N \quad (37)$$

$$V_{\mathbf{0},z}^s(\beta) = \frac{1}{\Lambda + \beta} \sum_{j=1}^N \lambda_j (s_{z,j} - s_{1,j} + V_{\mathbf{e}_j,1}^s(\beta)), \quad z = 2, 3, \dots, N \quad (38)$$

where

$$r(k,l) := \begin{cases} s_{k,l} - s_{k,l+1} + s_{k-1,l+1} - s_{k-1,l} & \text{if } l = 1, 2, \dots, N-1 \\ s_{k,N} - s_{k-1,N} + \sum_{m=1}^N \frac{\lambda_m}{\Lambda + \beta} (s_{k-1,m} - s_{k,m}) & \text{if } l = N \end{cases} \quad (39)$$

for $k = 2, \dots, N$, $l = k, \dots, N$.

In (36) the coefficients $r(1), \dots, r(N)$ are recursively defined by

$$\begin{aligned} \sum_{j=1}^{i-1} r(j) g_1 \circ \cdots \circ g_j(\beta) &= \beta s_{1,i} + \sum_{k=1}^{i-1} \lambda_k (s_{1,i} + s_{i,k} - s_{1,k}) \\ &\quad - \sum_{k=2}^{i-1} \sum_{l=k}^{i-1} r(k,l) g_1 \circ \cdots \circ g_{k-1}(\beta + \lambda_k + \cdots + \lambda_l), \quad i = 2, 3, \dots, N, \end{aligned} \quad (40)$$

$$\sum_{j=1}^N r(j) g_1 \circ \cdots \circ g_j(\beta) = - \sum_{k=2}^N \sum_{l=k}^N r(k,l) g_1 \circ \cdots \circ g_{k-1}(\beta + \lambda_k + \cdots + \lambda_l). \quad (41)$$

◇

Proof. Since the switching costs $s_{m,n}$ are bounded, $V_{\mathbf{x},z}^s(\beta)$ is the unique bounded solution of the dynamic programming (DP) equations [14]

$$(\Lambda + \mu_i + \beta)V_{\mathbf{x}_{[i]},i}^s(\beta) = \sum_{m=1}^N \lambda_m V_{\mathbf{x}_{[i]}+\mathbf{e}_m,i}^s(\beta) + \mu_i V_{\mathbf{x}_{[i]}-\mathbf{e}_i,i}^s(\beta) \quad (42)$$

$$(\Lambda + \beta)V_{\mathbf{0},z}^s(\beta) = \sum_{m=1}^N \lambda_m V_{\mathbf{e}_m,z}^s(\beta) \quad (43)$$

$$V_{\mathbf{x}_{[i]},z}^s(\beta) = s_{z,i} + V_{\mathbf{x}_{[i]},i}^s(\beta) \quad (44)$$

for $i, z = 1, 2, \dots, N$.

The proof will consist in checking that $V_{\mathbf{x},z}^s(\beta)$, defined in (36)-(38), satisfies the DP equations (42)-(44).

Equations (43) for $z = 2, 3, \dots, N$ and (44) for $z = 1, 2, \dots, N$ are automatically satisfied thanks to (37) and (38). It remains to check that $V_{\mathbf{x},z}^s(\beta)$ satisfies equations (42) and (43) for $z = 1$.

We begin with (42). The following identities hold ($i = 1, 2, \dots, N$):

$$V_{\mathbf{x}_{[i]},i}^s(\beta) = -s_{1,i} + V_{\mathbf{x}_{[i]},1}^s(\beta); \quad (45)$$

$$\begin{aligned} V_{\mathbf{x}_{[i]},1}^s(\beta) &= \sum_{j=1}^{i-1} r(j) + \sum_{j=i}^N r(j) \tau_{\mathbf{x}_{[i]}^j}(\beta) + \sum_{k=2}^{i-1} \sum_{l=k}^{i-1} r(k, l) \\ &+ \sum_{k=i+1}^N \sum_{l=k}^N r(k, l) \tau_{\mathbf{x}_{[i]}^{k-1}}(\beta + \lambda_k + \dots + \lambda_l) \mathbf{1}(x_k = \dots = x_l = 0); \end{aligned} \quad (46)$$

$$\begin{aligned} V_{\mathbf{x}_{[i]}+\mathbf{e}_m,1}^s(\beta) &= \sum_{j=1}^{m-1} r(j) + \sum_{j=m}^{i-1} r(j) \gamma_m(g_{m+1} \circ \dots \circ g_j(\beta)) \\ &+ \sum_{j=i}^N r(j) \tau_{\mathbf{x}_{[i]}^j+\mathbf{e}_m^j}(\beta) + \sum_{k=2}^{m-1} \sum_{l=k}^{m-1} r(k, l) \\ &+ \sum_{k=m+1}^{i-1} \sum_{l=k}^{i-1} r(k, l) \gamma_m(g_{m+1} \circ \dots \circ g_{k-1}(\beta + \lambda_k + \dots + \lambda_l)) \\ &+ \sum_{k=i+1}^N \sum_{l=k}^N r(k, l) \tau_{\mathbf{x}_{[i]}^{k-1}+\mathbf{e}_m^{k-1}}(\beta + \lambda_k + \dots + \lambda_l) \mathbf{1}(x_k = \dots = x_l = 0) \end{aligned} \quad (47)$$

for $m = 1, 2, \dots, i$;

$$V_{\mathbf{x}_{[i]}+\mathbf{e}_m,1}^s(\beta) = \sum_{j=1}^{i-1} r(j) + \sum_{j=i}^N r(j) \tau_{\mathbf{x}_{[i]}^j+\mathbf{e}_m^j}(\beta) + \sum_{k=2}^{i-1} \sum_{l=k}^{i-1} r(k, l) \quad (48)$$

$$\begin{aligned}
& + \sum_{k=i+1}^{m-1} \sum_{l=k}^{m-1} r(k, l) \tau_{\mathbf{x}_{[i]}^{k-1}}(\beta + \lambda_k + \cdots + \lambda_l) \mathbf{1}(x_k = \cdots = x_l = 0) \\
& + \sum_{k=m+1}^N \sum_{l=k}^N r(k, l) \tau_{\mathbf{x}_{[i]}^{k-1} + \mathbf{e}_m^{k-1}}(\beta + \lambda_k + \cdots + \lambda_l) \mathbf{1}(x_k = \cdots = x_l = 0)
\end{aligned}$$

for $m = i + 1, i + 2, \dots, N$;

$$\begin{aligned}
V_{\mathbf{x}_{[i]}^s - \mathbf{e}_i, 1}(\beta) &= \sum_{j=1}^{i-1} r(j) + \sum_{j=i}^N r(j) \tau_{\mathbf{x}_{[i]}^j - \mathbf{e}_i^j}(\beta) + \sum_{k=2}^{i-1} \sum_{l=k}^{i-1} r(k, l) \\
&+ \sum_{k=i+1}^N \sum_{l=k}^N r(k, l) \tau_{\mathbf{x}_{[i]}^{k-1} - \mathbf{e}_i^{k-1}}(\beta + \lambda_k + \cdots + \lambda_l) \mathbf{1}(x_k = \cdots = x_l = 0).
\end{aligned} \tag{49}$$

Identity (45) follows from (37). The other identities directly come from (37) and from the relation

$$\tau_{\mathbf{e}_m^j}(\beta) = \gamma_m(g_{m+1} \circ \cdots \circ g_j(\beta)) \tag{50}$$

that we used to derive (47), where (50) was obtained from (13).

Fix $i \in \{1, 2, \dots, N\}$. With (45)-(49) the DP equation (42) becomes

$$0 = X_i + Y_i$$

with

$$\begin{aligned}
X_i &:= -(\Lambda + \mu_i + \beta) \left(-s_{1,i} + \sum_{j=1}^{i-1} r(j) + \sum_{k=2}^{i-1} \sum_{l=k}^{i-1} r(k, l) \right) \\
&+ \sum_{m=1}^{i-1} \lambda_m (s_{i,m} - s_{1,m}) - \sum_{m=i}^N \lambda_m s_{1,i} \\
&+ \sum_{m=1}^i \lambda_m \left(\sum_{j=1}^{m-1} r(j) + \sum_{j=m}^{i-1} r(j) \gamma_m(g_{m+1} \circ \cdots \circ g_j(\beta)) + \sum_{k=2}^{m-1} \sum_{l=k}^{m-1} r(k, l) \right. \\
&\left. + \sum_{k=m+1}^{i-1} \sum_{l=k}^{i-1} r(k, l) \gamma_m(g_{m+1} \circ \cdots \circ g_{k-1}(\beta + \lambda_k + \cdots + \lambda_l)) \right) \\
&+ \sum_{m=i+1}^N \lambda_m \left(\sum_{j=1}^{i-1} r(j) + \sum_{k=2}^{i-1} \sum_{l=k}^{i-1} r(k, l) \right) \\
&+ \mu_i \left(-s_{1,i} + \sum_{j=1}^{i-1} r(j) + \sum_{k=2}^{i-1} \sum_{l=k}^{i-1} r(k, l) \right)
\end{aligned}$$

and

$$\begin{aligned}
Y_i := & -(\Lambda + \mu_i + \beta) \left(\sum_{j=i}^N r(j) \tau_{\mathbf{x}_{[i]}}^j(\beta) \right. \\
& + \sum_{k=i+1}^N \sum_{l=k}^N r(k, l) \tau_{\mathbf{x}_{[i]}^{k-1}}(\beta + \lambda_k + \dots + \lambda_l) \mathbf{1}(x_k = \dots = x_l = 0) \Big) \\
& + \sum_{m=1}^i \lambda_m \left(\sum_{j=i}^N r(j) \tau_{\mathbf{x}_{[i]} + \mathbf{e}_m^j}(\beta) \right. \\
& + \sum_{k=i+1}^N \sum_{l=k}^N r(k, l) \tau_{\mathbf{x}_{[i]}^{k-1} + \mathbf{e}_m^{k-1}}(\beta + \lambda_k + \dots + \lambda_l) \mathbf{1}(x_k = \dots = x_l = 0) \Big) \\
& + \sum_{m=i+1}^N \lambda_m \left(\sum_{j=i}^{m-1} r(j) \tau_{\mathbf{x}_{[i]}}^j(\beta) + \sum_{j=m}^N r(j) \tau_{\mathbf{x}_{[i]} + \mathbf{e}_m^j}(\beta) \right. \\
& + \sum_{k=i+1}^{m-1} \sum_{l=k}^{m-1} r(k, l) \tau_{\mathbf{x}_{[i]}^{k-1}}(\beta + \lambda_k + \dots + \lambda_l) \mathbf{1}(x_k = \dots = x_l = 0) \\
& + \sum_{k=m+1}^N \sum_{l=k}^N r(k, l) \tau_{\mathbf{x}_{[i]}^{k-1} + \mathbf{e}_m^{k-1}}(\beta + \lambda_k + \dots + \lambda_l) \mathbf{1}(x_k = \dots = x_l = 0) \Big) \\
& + \mu_i \left(\sum_{j=i}^N r(j) \tau_{\mathbf{x}_{[i]} - \mathbf{e}_i^j}(\beta) \right. \\
& + \sum_{k=i+1}^N \sum_{l=k}^N r(k, l) \tau_{\mathbf{x}_{[i]}^{k-1} - \mathbf{e}_i^{k-1}}(\beta + \lambda_k + \dots + \lambda_l) \mathbf{1}(x_k = \dots = x_l = 0) \\
& + \left[s_{1,i} + V_{\mathbf{x}_{[i]} - \mathbf{e}_i, i}^s(\beta) - \sum_{j=1}^{i-1} r(j) - \sum_{j=i}^N r(j) \tau_{\mathbf{x}_{[i]} - \mathbf{e}_i^j}(\beta) - \sum_{k=2}^{i-1} \sum_{l=k}^{i-1} r(k, l) \right. \\
& \left. - \sum_{k=i+1}^N \sum_{l=k}^N r(k, l) \tau_{\mathbf{x}_{[i]}^{k-1} - \mathbf{e}_i^{k-1}}(\beta + \lambda_k + \dots + \lambda_l) \mathbf{1}(x_k = \dots = x_l = 0) \right] \mathbf{1}(x_i = 1) \Big). \tag{51}
\end{aligned}$$

The terms X_i and Y_i have been obtained as follows: all constant terms in (42) have been collected in X_i , whereas all terms which are functions of x_1, \dots, x_N have been collected in Y_i . Also notice that we have taken into account the fact that $V_{\mathbf{x}_{[i]} - \mathbf{e}_i, i}^s(\beta) \neq -s_{1,i} + V_{\mathbf{x}_{[i]} - \mathbf{e}_i, 1}^s(\beta)$ when $x_i = 1$, which explains the presence of the correction term in (51) when $x_i = 1$.

We now check that $X_i = 0$ and that $Y_i = 0$ for all $i = 1, 2, \dots, N$.

(i) *Checking that $X_i = 0$*

Recall that i is fixed in $\{1, 2, \dots, N\}$. Elementary algebra yields

$$\begin{aligned}
X_i &= - \sum_{j=1}^{i-1} r(j) \left[\beta + \sum_{m=1}^j \lambda_m (1 - \gamma_m(g_{m+1} \circ \dots \circ g_j(\beta))) \right] \\
&\quad + \beta s_{1,i} + \sum_{m=1}^{i-1} \lambda_m (s_{1,i} + s_{i,m} - s_{1,m}) \\
&\quad - \sum_{k=2}^{i-1} \sum_{l=k}^{i-1} r(k, l) \left[\beta + \sum_{m=k}^l \lambda_m \right. \\
&\quad \left. + \sum_{m=1}^{k-1} \lambda_m (1 - \gamma_m(g_{m+1} \circ \dots \circ g_{k-1}(\beta + \lambda_k + \dots + \lambda_l))) \right].
\end{aligned}$$

With the help of the relation

$$g_i \circ \dots \circ g_k(\beta) = \beta + \sum_{m=i}^k \lambda_m (1 - \gamma_m(g_{m+1} \circ \dots \circ g_k(\beta))) \quad (52)$$

that follows from the definition of mappings $\{g_n(\beta)\}_n$, X_i rewrites as

$$\begin{aligned}
X_i &= - \sum_{j=1}^{i-1} r(j) g_1 \circ \dots \circ g_j(\beta) + \beta s_{1,i} + \sum_{m=1}^{i-1} \lambda_m (s_{1,i} + s_{i,m} - s_{1,m}) \\
&\quad - \sum_{k=2}^{i-1} \sum_{l=k}^{i-1} r(k, l) g_1 \circ \dots \circ g_{k-1}(\beta + \lambda_k + \dots + \lambda_l). \quad (53)
\end{aligned}$$

The r.h.s. of (53) vanishes for $i = 1$; it also vanishes for $i = 2, 3, \dots, N$ from the definition of coefficients $r(1), \dots, r(N-1)$ given in (40)-(41).

(ii) *Checking that $Y_i = 0$.*

Write Y_i as $Y_i = Y_{i,1} + Y_{i,2} + \mu_i Y_{i,3}$ with

$$\begin{aligned}
Y_{i,1} &:= -(\Lambda + \mu_i + \beta) \sum_{j=i}^N r(j) \tau_{\mathbf{x}_{[i]}^j}(\beta) + \sum_{m=1}^i \lambda_m \sum_{j=i}^N r(j) \tau_{\mathbf{x}_{[i]}^j + \mathbf{e}_m^j}(\beta) \\
&\quad + \sum_{m=i+1}^N \lambda_m \left(\sum_{j=i}^{m-1} r(j) \tau_{\mathbf{x}_{[i]}^j}(\beta) + \sum_{j=m}^N r(j) \tau_{\mathbf{x}_{[i]}^j + \mathbf{e}_m^j}(\beta) \right) + \mu_i \sum_{j=i}^N r(j) \tau_{\mathbf{x}_{[i]}^j - \mathbf{e}_i^j}(\beta),
\end{aligned}$$

$$Y_{i,2} := -(\Lambda + \mu_i + \beta) \sum_{k=i+1}^N \sum_{l=k}^N r(k, l) \tau_{\mathbf{x}_{[i]}^{k-1}}(\beta + \lambda_k + \dots + \lambda_l) \mathbf{1}(x_k = \dots = x_l = 0)$$

$$\begin{aligned}
& + \sum_{m=1}^i \lambda_m \sum_{k=i+1}^N \sum_{l=k}^N r(k, l) \tau_{\mathbf{x}_{[i]}^{k-1} + \mathbf{e}_m^{k-1}}(\beta + \lambda_k + \dots + \lambda_l) \mathbf{1}(x_k = \dots = x_l = 0) \\
& + \sum_{m=i+1}^N \lambda_m \left(\sum_{k=i+1}^{m-1} \sum_{l=k}^{m-1} r(k, l) \tau_{\mathbf{x}_{[i]}^{k-1}}(\beta + \lambda_k + \dots + \lambda_l) \mathbf{1}(x_k = \dots = x_l = 0) \right. \\
& + \sum_{k=m}^N \sum_{l=k}^N r(k, l) \tau_{\mathbf{x}_{[i]}^{k-1} + \mathbf{e}_m^{k-1}}(\beta + \lambda_k + \dots + \lambda_l) \mathbf{1}(x_k = \dots = x_l = 0) \left. \right) \\
& + \mu_i \sum_{k=i+1}^N \sum_{l=k}^N r(k, l) \tau_{\mathbf{x}_{[i]}^{k-1} - \mathbf{e}_i^{k-1}}(\beta + \lambda_k + \dots + \lambda_l) \mathbf{1}(x_k = \dots = x_l = 0)
\end{aligned}$$

and

$$\begin{aligned}
Y_{i,3} & := \left[s_{1,i} + V_{\mathbf{x}_{[i]}^s - \mathbf{e}_i, i}(\beta) - \sum_{j=1}^{i-1} r(j) - \sum_{j=i}^N r(j) \tau_{\mathbf{x}_{[i]}^j - \mathbf{e}_i^j}(\beta) - \sum_{k=2}^{i-1} \sum_{l=k}^{i-1} r(k, l) \right. \\
& \left. - \sum_{k=i+1}^N \sum_{l=k}^N r(k, l) \tau_{\mathbf{x}_{[i]}^{k-1} - \mathbf{e}_i^{k-1}}(\beta + \lambda_k + \dots + \lambda_l) \mathbf{1}(x_k = \dots = x_l = 0) \right] \mathbf{1}(x_i = 1)
\end{aligned}$$

Let us show that $Y_{i,1} = Y_{i,2} = Y_{i,3} = 0$.

Introduce the mappings ($n = 1, 2, \dots, N$)

$$f_n(s) = \lambda_n \gamma_n(s) - \left(\mu_n + \lambda_n + s + \sum_{j=1}^{n-1} \lambda_j (1 - \gamma_j(g_{j+1} \circ \dots \circ g_n(s))) \right) + \mu_n \gamma_n^{-1}(s). \quad (54)$$

We observe from (16) that

$$f_n(s) = 0 \quad \text{for all } s > 0, n = 1, 2, \dots, N. \quad (55)$$

After elementary algebra we find

$$\begin{aligned}
Y_{i,1} & = \sum_{j=i}^N r(j) \tau_{\mathbf{x}_{[i]}^j}(\beta) \left[f_i(g_{i+1} \circ \dots \circ g_j(\beta)) \right. \\
& \quad \left. + g_{i+1} \circ \dots \circ g_j(\beta) - \beta - \sum_{m=i+1}^j \lambda_m (1 - \gamma_m(g_{m+1} \circ \dots \circ g_j(\beta))) \right] \\
& = 0
\end{aligned} \quad (56)$$

from (52) and (55). To establish (56) we have used the identities

$$\tau_{\mathbf{x}^j + \mathbf{e}_m^j}(\beta) = \gamma_m(g_{m+1} \circ \dots \circ g_j(\beta)) \tau_{\mathbf{x}^j}(\beta), \quad m = 1, 2, \dots, j$$

that we have derived from (13).

Similarly we find

$$\begin{aligned}
Y_{i,2} &= \sum_{k=i+1}^N \sum_{l=k}^N r(k,l) \tau_{\mathbf{x}_{[i]}^{k-1}}(\beta + \lambda_k + \dots + \lambda_l) \mathbf{1}(x_k = \dots = x_l = 0) \\
&\quad \left[f_i(g_{i+1} \circ \dots \circ g_{k-1}(\beta + \lambda_k + \dots + \lambda_l)) + g_{i+1} \circ \dots \circ g_{k-1}(\beta + \lambda_k + \dots + \lambda_l) \right. \\
&\quad \left. - (\beta + \lambda_k + \dots + \lambda_l) - \sum_{m=i+1}^{k-1} \lambda_m (1 - \gamma_m(g_{m+1} \circ \dots \circ g_{k-1}(\beta + \lambda_k + \dots + \lambda_l))) \right] \\
&= 0
\end{aligned}$$

where the latter equality again follows from (55) and (52).

To show that $Y_{i,3} = 0$ for all $i = 1, 2, \dots, N$ we proceed as follows. Assume that $\mathbf{x}_{[i]} = (0, \dots, 0, x_i, x_{i+1}, \dots, x_N)$ with $x_i = 1$, $x_{i+1} = \dots = x_{t-1} = 0$ and $x_t > 0$ for some $t = i + 1, \dots, N$.

Then, $V_{\mathbf{x}_{[i]}^s - \mathbf{e}_i, i}(\beta) = s_{i,t} - s_{1,t} + V_{\mathbf{x}_{[i]}^t, 1}(\beta)$, and $Y_{i,3}$ becomes

$$\begin{aligned}
Y_{i,3} &= s_{1,i} + s_{i,t} - s_{1,t} + V_{\mathbf{x}_{[i]}^t, 1}(\beta) - \sum_{j=1}^{i-1} r(j) - \sum_{j=i}^N r(j) \tau_{\mathbf{x}_{[i]}^j - \mathbf{e}_i^j}(\beta) - \sum_{k=2}^{i-1} \sum_{l=k}^{i-1} r(k,l) \\
&\quad - \sum_{k=i+1}^N \sum_{l=k}^N r(k,l) \tau_{\mathbf{x}_{[i]}^{k-1} - \mathbf{e}_i^{k-1}}(\beta + \lambda_k + \dots + \lambda_l) \mathbf{1}(x_k = \dots = x_l = 0) \\
&= s_{1,i} + s_{i,t} - s_{1,t} + \sum_{k=2}^i \sum_{l=i}^{t-1} r(k,l) \tag{57}
\end{aligned}$$

where we have used (46) with $i = t$ to establish the latter equality. From our convention that $s_{i,i} = 0$ for all $i = 1, 2, \dots, N$ we see that the r.h.s. of (57) vanishes when $i = 1$ so that $Y_{1,3} = 0$. Let us now show that the r.h.s. of (57) also vanishes for $i = 2, 3, \dots, N$ or, equivalently, that

$$\sum_{k=2}^i \sum_{l=i}^{t-1} r(k,l) = -s_{1,i} - s_{i,t} + s_{1,t} \quad \text{for } 2 \leq i < t \leq N. \tag{58}$$

From (58) we find

$$\sum_{k=2}^i \sum_{l=i}^{t-1} r(k,l) - \sum_{k=2}^i \sum_{l=i}^{t-2} r(k,l) = \sum_{k=2}^i r(k, t-1) = s_{i,t-1} - s_{i,t} + s_{1,t} - s_{1,t-1}$$

so that

$$\sum_{k=2}^i r(k, t-1) - \sum_{k=2}^{i-1} r(k, t-1) = r(i, t-1) = s_{i,t-1} - s_{i,t} + s_{i-1,t} - s_{i-1,t-1}$$

for $2 \leq i < t < N$, or equivalently, $r(i, j) = s_{i,j} - s_{i,j+1} + s_{i-1,j+1} - s_{i-1,j}$ for $2 \leq i \leq j < N$, which is the result announced in (39-a).

It remains to handle the situation where $x_i = 1$ and $x_{i+1} = \dots = x_N = 0$ for $i = 1, 2, \dots, N$. In this case $V_{\mathbf{x}_{[i]}^s - \mathbf{e}_{i,i}}(\beta) = V_{\mathbf{0},i}^s(\beta)$, where $V_{\mathbf{0},i}^s(\beta)$ is given in (43), and

$$\begin{aligned} Y_{i,3} &= s_{1,i} + \frac{1}{\Lambda + \beta} \sum_{j=1}^N \lambda_j V_{\mathbf{e}_j,j}^s(\beta) - \sum_{j=1}^N r(j) - \sum_{k=1}^{i-1} \sum_{l=k}^{i-1} r(k, l) - \sum_{k=i+1}^N \sum_{l=k}^N r(k, l) \\ &= s_{1,i} + \frac{1}{\Lambda + \beta} \sum_{j=1}^N \lambda_j (s_{i,j} - s_{1,j}) + \frac{1}{\Lambda + \beta} \sum_{j=1}^N \lambda_j V_{\mathbf{e}_j,1}^s(\beta) - \sum_{j=1}^N r(j) \end{aligned} \quad (59)$$

$$\begin{aligned} &- \sum_{k=2}^N \sum_{l=k}^N r(k, l) + \sum_{k=2}^i \sum_{l=i}^N r(k, l) \\ &= s_{1,i} + \frac{1}{\Lambda + \beta} \sum_{j=1}^N \lambda_j (s_{i,j} - s_{1,j}) + \sum_{k=2}^i \sum_{l=i}^N r(k, l) \end{aligned} \quad (60)$$

for $i = 1, 2, \dots, N$. The relation (59) follows from the identity $V_{\mathbf{e}_j,i}^s(\beta) = s_{i,j} - s_{1,j} + V_{\mathbf{e}_j,1}^s(\beta)$ and (60) has been obtained by using (41), (46) and (52).

Note that $Y_{1,3} = 0$ since $s_{1,1} = 0$ by convention. On the other hand, $Y_{i,3}$ will vanish for $i = 2, 3, \dots, N$ if and only if

$$\sum_{k=2}^i \sum_{l=i}^N r(k, l) = -s_{1,i} - \frac{1}{\Lambda + \beta} \sum_{j=1}^N \lambda_j (s_{i,j} - s_{1,j})$$

for $i = 2, 3, \dots, N$. From this relation and (58) with $t = N$ we deduce that

$$\sum_{k=2}^i \sum_{l=i}^N r(k, l) - \sum_{k=2}^i \sum_{l=i}^{N-1} r(k, l) = \sum_{k=2}^i r(k, N) = s_{i,N} - s_{1,N} - \frac{1}{\Lambda + \beta} \sum_{j=1}^N \lambda_j (s_{i,j} - s_{1,j})$$

and finally,

$$\sum_{k=2}^i r(k, N) - \sum_{k=2}^{i-1} r(k, N) = r(i, N) = s_{i,N} - s_{i-1,N} + \frac{1}{\Lambda + \beta} \sum_{j=1}^N \lambda_j (s_{i-1,j} - s_{i,j})$$

for $i = 2, 3, \dots, N$, which is the result announced in (39-b).

It remains to check that the DP equation (43) with $z = 1$ is satisfied. We have (hint: use (46))

$$-(\Lambda + \beta) V_{\mathbf{0},1}^s(\beta) + \sum_{m=1}^N \lambda_m V_{\mathbf{e}_m,1}^s(\beta)$$

$$\begin{aligned}
&= - \sum_{j=1}^N \left[\beta + \sum_{m=1}^j \lambda_m (1 - \gamma_m (g_{m+1} \circ \dots \circ g_j(\beta))) \right] \\
&\quad - \sum_{k=2}^N \sum_{l=k}^N r(k, l) \left[\beta + \sum_{m=k}^l \lambda_m + \sum_{m=1}^{k-1} \lambda_m (1 - \gamma_m (g_{m+1} \circ \dots \circ g_j(\beta + \lambda_k + \dots + \lambda_l))) \right] \\
&= - \sum_{j=1}^N r(j) g_1 \circ \dots \circ g_j(\beta) - \sum_{k=2}^N \sum_{l=k}^N r(k, l) g_1 \circ \dots \circ g_{k-1}(\beta + \lambda_k + \dots + \lambda_l) \quad \text{from (56)} \\
&= 0
\end{aligned}$$

where the latter equality follows from the definition of $r(N)$ given in (41). This concludes the proof. \blacksquare

The fact that (36)-(38) only depend on the LST of the service time distributions suggests that Theorem 5.1 might hold for arbitrary service time distributions, up to a modification of the state of the system (to incorporate remaining service times). However, we were unable to show this.

Unlike for holding costs, it is not easy to find an intuitive explanation for the form of the switching cost function $V_{\mathbf{x},z}^s(\beta)$ for an arbitrary number of classes. Such an intuitive argument was developed in [9] for two classes of customers.

6 Average costs

In this section we derive, in explicit form, the average holding cost $g_{\mathbf{x}}^h$ and the switching cost $g_{\mathbf{x},z}^s$ defined in (4) and (5), respectively, along with the bias function $B_{\mathbf{x},z;\mathbf{x}_0,z_0}$ introduced in (6).

We assume in this section that $\rho = \sum_{i=1}^n \rho_i = -\sum_{i=1}^n \lambda_i S_i'(0) < 1$, i.e., the system is stable. Under this condition all LSTs from Section 3 represent proper probability distributions and therefore converge to 1 as their arguments tend to 0 (see [16, pp. 47–48, Lemma 1]). We also assume that the 2nd moments of the service time distributions are finite.

6.1 Average holding cost

Thanks to Theorem (4.1) we can write $g_{\mathbf{x}}^h$ as

$$g_{\mathbf{x}}^h = \sum_{n=1}^N c_n g_{\mathbf{x}^n}^{h,n} \quad (61)$$

with

$$g_{\mathbf{x}^n}^{h,n} := \lim_{\beta \rightarrow 0} \beta H_{\mathbf{x}^n}^n(\beta). \quad (62)$$

It is therefore sufficient to compute $g_{\mathbf{x}^n}^{h,n}$.

In Section 4 we have found that

$$\beta H_{\mathbf{x}^n}^n(\beta) = x_n + \frac{\lambda_n \frac{1 - C_n(\beta)}{\beta} \phi_n(\beta) + C_n(\beta) [\tau_{\mathbf{x}^n}(\beta) \phi_{n-1}(\beta) - \tau_{\mathbf{x}^{n-1}}(\beta) \phi_n(\beta)]}{(1 - C_n(\beta)) \phi_n(\beta)} \quad (63)$$

where we have set

$$\phi_n(\beta) := \beta + \Lambda_n (1 - \tilde{\gamma}_n(\beta)).$$

Since the numerator and the denominator in the r.h.s. of (63) vanish when $\beta \rightarrow 0$, we need to expand them in Taylor series in order to compute $\lim_{\beta \rightarrow 0} \beta H_{\mathbf{x}^n}^n(\beta)$. We find

$$\begin{aligned} \lambda_n \frac{1 - C_n(\beta)}{\beta} \phi_n(\beta) &= -\lambda_n C_n^{(1)}(0) \phi_n^{(1)}(0) \beta \\ &\quad - \frac{\lambda_n}{2} \left[C_n^{(1)}(0) \phi_n^{(2)}(0) + C_n^{(2)}(0) \phi_n^{(1)}(0) \right] \beta^2 + o(\beta^2), \end{aligned} \quad (64)$$

$$\begin{aligned} C_n(\beta) [\tau_{\mathbf{x}^n}(\beta) \phi_{n-1}(\beta) - \tau_{\mathbf{x}^{n-1}}(\beta) \phi_n(\beta)] &= \left[\phi_{n-1}^{(1)}(0) - \phi_n^{(1)}(0) \right] \beta \\ &\quad + \left[C_n^{(1)}(0) \left(\phi_{n-1}^{(1)}(0) - \phi_n^{(1)}(0) \right) + \tau_{\mathbf{x}^n}^{(1)}(0) \phi_{n-1}^{(1)}(0) - \tau_{\mathbf{x}^{n-1}}^{(1)}(0) \phi_n^{(1)}(0) \right. \\ &\quad \left. + \frac{1}{2} \phi_{n-1}^{(2)}(0) - \frac{1}{2} \phi_n^{(2)}(0) \right] \beta^2 + o(\beta^2) \end{aligned} \quad (65)$$

and

$$(1 - C_n(\beta)) \phi_n(\beta) = -C_n^{(1)}(0) \phi_n^{(1)}(0) \beta^2 + o(\beta^2) \quad (66)$$

where $f^{(k)}(\beta)$ denotes the k -th derivative of the $f(\beta)$ in the variable β .

It remains to evaluate the constants $C_n^{(k)}(0)$, $\phi_n^{(k)}(0)$ and $\tau_{\mathbf{x}^n}^{(1)}(0)$ for $k = 1, 2$. We easily find from Lemma 3.2

$$\phi_n^{(1)}(0) = \frac{1}{1 - \sum_{i=1}^n \rho_i} \quad (67)$$

$$\phi_n^{(2)}(0) = -\frac{\sum_{i=1}^n \lambda_i \mathbf{E}[S_i^2]}{(1 - \sum_{i=1}^n \rho_i)^3} \quad (68)$$

$$C_n^{(1)}(0) = -\frac{1}{\mu_n (1 - \sum_{i=1}^{n-1} \rho_i)} \quad (69)$$

$$C_n^{(2)}(0) = \frac{\sum_{i=1}^{n-1} \lambda_i \mathbf{E}[S_i^2]}{\mu_n \left(1 - \sum_{i=1}^{n-1} \rho_i\right)^3} + \frac{\mathbf{E}[S_n^2]}{\left(1 - \sum_{i=1}^{n-1} \rho_i\right)^2} \quad (70)$$

$$\begin{aligned}
\tau_{\mathbf{x}^n}^{(1)}(0) &= \left(\frac{1 - \sum_{i=1}^{n-1} \rho_i}{1 - \sum_{i=1}^n \rho_i} \right) \tau_{\mathbf{x}^{n-1}}^{(1)}(0) - \frac{x_n}{\mu_n (1 - \sum_{i=1}^n \rho_i)} \\
&= -\frac{\sum_{i=1}^n x_i / \mu_i}{1 - \sum_{i=1}^n \rho_i}
\end{aligned} \tag{71}$$

for $n = 1, 2, \dots, N$, where $\mathbf{E}[S_n^2]$ is the second-order moment of service times of class- n customers.

With (67) and (69) we find that $-\lambda_n C_n^{(1)}(0) \phi_n^{(1)}(0) + \phi_{n-1}^{(1)}(0) - \phi_n^{(1)}(0)$, the coefficient of β in the Taylor series expansion of the numerator in the r.h.s. of (63), vanishes. Therefore,

$$\begin{aligned}
\lim_{\beta \rightarrow 0} \beta H_{\mathbf{x}^n}^n(\beta) &= \left[x_n - \tau_{\mathbf{x}^n}^{(1)}(0) \phi_{n-1}^{(1)}(0) + \tau_{\mathbf{x}^{n-1}}^{(1)}(0) \phi_n^{(1)}(0) \right] / C_n^{(1)}(0) \phi_n^{(1)}(0) \\
&+ \left[\frac{\lambda_n}{2} \left(C_n^{(1)}(0) \phi_n^{(2)}(0) + C_n^{(2)}(0) \phi_n^{(1)}(0) \right) - C_n^{(1)}(0) \left(\phi_{n-1}^{(1)}(0) - \phi_n^{(1)}(0) \right) \right. \\
&\left. - \frac{1}{2} \left(\phi_{n-1}^{(2)}(0) - \phi_n^{(2)}(0) \right) \right] / C_n^{(1)}(0) \phi_n^{(1)}(0).
\end{aligned} \tag{72}$$

By using (67), (69) and (71) we see that that the first term between square brackets in the r.h.s. of (72) vanishes; the second term between square brackets in the r.h.s. of (72) can also be easily evaluated with the help of identities (67)-(70), to yield

$$g_{\mathbf{x}^n}^{h,n} = \lim_{\beta \rightarrow 0} \beta H_{\mathbf{x}^n}^n(\beta) = \frac{\rho_n}{1 - \sum_{i=1}^{n-1} \rho_i} + \frac{\lambda_n \sum_{i=1}^n \lambda_i \mathbf{E}[S_i^2]}{2 \left(1 - \sum_{i=1}^{n-1} \rho_i \right) \left(1 - \sum_{i=1}^n \rho_i \right)} \tag{73}$$

Note that the r.h.s. of (73) is independent of the initial state \mathbf{x}^n .

The r.h.s. of (74) is actually the mean number of class n customers in steady-state (see e.g. [7, Formula (7.41), p. 116], [8, Chapter 3]). This result should not come as a surprise since

$$g_{\mathbf{x}^n}^{h,n} = \lim_{\beta \rightarrow 0} \beta \mathbf{E}_{\mathbf{x}^n} \left[\int_0^\infty e^{-\beta t} X_n(t) dt \right]$$

from (62) and (17). Therefore, if we knew that $\lim_{t \rightarrow \infty} \mathbf{E}_{\mathbf{x}^n}[X_n(t)] = \mathbf{E}[X_n]$, independent of the initial state \mathbf{x}^n , where X_n is the stationary number of class- n customer then, by a Tauberian theorem, we would have that $g_{\mathbf{x}^n}^{h,n} = \mathbf{E}[X_n]$, which is exactly (74). Unfortunately, showing that $\lim_{t \rightarrow \infty} \mathbf{E}_{\mathbf{x}^n}[X_n(t)] = \mathbf{E}[X_n]$ is not easy, despite the fact that it is known that the r.v. $X_n(t)$ converges weakly to X_n as $t \rightarrow \infty$, independent of the initial state, under the stability condition [18]. This explains why we had to go through a lengthy procedure for computing $g_{\mathbf{x}^n}^{h,n}$. Note in passing that the above explanation justifies the fact that $g_{\mathbf{x}^n}^{h,n}$ was introduced as an ‘‘average cost’’.

In conclusion, cf. (61) and (73),

$$g_{\mathbf{x}}^h = \sum_{n=1}^N c_n \left(\frac{\rho_n}{1 - \sum_{i=1}^{n-1} \rho_i} + \frac{\lambda_n \sum_{i=1}^n \lambda_i \mathbf{E}[S_i^2]}{2 \left(1 - \sum_{i=1}^{n-1} \rho_i\right) \left(1 - \sum_{i=1}^n \rho_i\right)} \right). \quad (74)$$

6.2 Average switching cost

We need to compute $g_{\mathbf{x},z}^s = \lim_{\beta \rightarrow 0} \beta V_{\mathbf{x},z}^s$. A glance at Theorem 5.1 shows that

$$g_{\mathbf{x},z}^s = \sum_{j=1}^N \hat{r}(j) \quad (75)$$

for all $\mathbf{x} \in \mathbb{N}^N$ and $z = 1, 2, \dots, N$, where

$$\hat{r}(j) := \lim_{\beta \rightarrow 0} \beta r(j)$$

Let us compute $\hat{r}(j)$ for $j = 1, 2, \dots, N$. To this end, introduce the constants

$$A(j) := \lim_{\beta \rightarrow 0} \frac{g_1 \circ \dots \circ g_j(\beta)}{\beta}, \quad j = 1, 2, \dots, N.$$

Assume that the constants $A(j)$ are well-defined. Then, we see from (40) and (41) that the constants $\hat{r}(j)$, $j = 1, 2, \dots, N$, are also well-defined and can be computed with the following recursive scheme [Hint: write $r(j) g_1 \circ \dots \circ g_j(\beta)$ in the l.h.s. of (40) and (41) as $\beta r(j) \beta^{-1} g_1 \circ \dots \circ g_j(\beta)$ and let $\beta \rightarrow 0$]:

$$\begin{aligned} \sum_{j=1}^{i-1} \hat{r}(j) A(j) &= \sum_{k=1}^{i-1} \lambda_k (s_{1,i} + s_{i,k} - s_{1,k}) \\ &\quad - \sum_{k=2}^{i-1} \sum_{l=k}^{i-1} r(k,l) g_1 \circ \dots \circ g_{k-1}(\lambda_k + \dots + \lambda_l), \quad i = 2, 3, \dots, N, \end{aligned} \quad (76)$$

$$\sum_{j=1}^N \hat{r}(j) A(j) = - \sum_{k=2}^N \sum_{l=k}^N r(k,l) g_1 \circ \dots \circ g_{k-1}(\lambda_k + \dots + \lambda_l). \quad (77)$$

It remains to compute $A(j)$ for $j = 1, 2, \dots, N$. This is done in the following lemma.

Lemma 6.1 For $j = 1, 2, \dots, N$,

$$A(j) = \prod_{k=1}^j \left(1 - \lambda_k \gamma_k^{(1)}(0)\right) \quad (78)$$

$$= \prod_{k=1}^j \frac{1 - \sum_{l=1}^{k-1} \rho_l}{1 - \sum_{l=1}^k \rho_l}, \quad j = 1, 2, \dots, N. \quad (79)$$

Proof. Recall that $g_n(\beta) = \beta + \lambda_n(1 - \gamma_n(\beta))$. We have

$$\lim_{\beta \rightarrow 0} \frac{g_n(\beta)}{\beta} = 1 - \lambda_n \gamma_n^{(1)}(0)$$

which shows that (78) holds for $j = 1$. Assume that (78) holds for $j = 1, 2, \dots, i$ and let us show that it still holds for $j = i + 1$. Since $\lim_{\beta \rightarrow 0} g_n(\beta) = 0$ we have

$$\begin{aligned} A(j+1) &= \lim_{\beta \rightarrow 0} \frac{g_1 \circ \dots \circ g_j(g_{j+1}(\beta))}{g_{j+1}(\beta)} \frac{g_{j+1}(\beta)}{\beta} \\ &= A(j) \left(1 - \lambda_{j+1} \gamma_{j+1}^{(1)}(0)\right) \\ &= \prod_{k=1}^{j+1} \left(1 - \lambda_k \gamma_k^{(1)}(0)\right) \end{aligned} \tag{80}$$

by using the induction hypothesis, which proves (78).

On the other hand, we easily find from (11) that

$$\begin{aligned} \gamma_n^{(1)}(0) &= \frac{C_n^{(1)}(0)}{1 + \lambda_n C_n^{(1)}(0)} \\ &= -\frac{1}{\mu_n (1 - \sum_{i=1}^n \rho_i)} \end{aligned} \tag{81}$$

where we have used (69) to establish (81). Plugging (81) into (80) gives (79), which concludes the proof. \blacksquare

In conclusion, the average switching cost $g_{\mathbf{x},z}^s$ is given by (75), where the constants $\hat{r}(j)$, $j = 1, 2, \dots, N$, can be computed from equations (76) and (77), with constants $A(1), \dots, A(N)$ given in Lemma 6.1. Note, as expected, that $g_{\mathbf{x},z}^s$ does not depend on the initial state (\mathbf{x}, z) .

6.3 Bias function

In this section we will determine, in explicit form, the bias function associated with the reference state $(\mathbf{x}_0, z_0) = (\mathbf{0}, 1)$, namely,

$$B_{\mathbf{x},z} := B_{\mathbf{x},z;\mathbf{0},1} = \lim_{\beta \rightarrow 0} (V_{\mathbf{x},z}(\beta) - V_{\mathbf{0},1}(\beta)).$$

There is of course no loss of generality in choosing that particular reference state since $B_{\mathbf{x},z;\mathbf{x}_0,z_0}$, as defined in (6), can easily be computed for any reference state (\mathbf{x}_0, z_0) if one knows the bias function $B_{\mathbf{x},z}$, thanks to the relation $B_{\mathbf{x},z;\mathbf{x}_0,z_0} = B_{\mathbf{x},z} - B_{\mathbf{x}_0,z_0}$.

The bias function $B_{\mathbf{x},z}$ can be decomposed as the sum of bias functions for holding costs and of a bias function for switching costs. More precisely (see (1) and Theorem 4.1)

$$\begin{aligned} B_{\mathbf{x},z} &= \lim_{\beta \rightarrow 0} \left(V_{\mathbf{x},z}^h(\beta) - V_{\mathbf{0},1}^h(\beta) \right) + \lim_{\beta \rightarrow 0} \left(V_{\mathbf{x},z}^s(\beta) - V_{\mathbf{0},1}^s(\beta) \right) \\ &= \sum_{n=1}^N c_n B_{\mathbf{x}^n}^{n,h} + B_{\mathbf{x},z}^s \end{aligned} \quad (82)$$

where

$$B_{\mathbf{x}^n}^{n,h} := \lim_{\beta \rightarrow 0} \left(H_{\mathbf{x}^n}^n(\beta) - H_{\mathbf{0}}^n(\beta) \right) \quad (83)$$

$$B_{\mathbf{x},z}^s := \lim_{\beta \rightarrow 0} \left(V_{\mathbf{x},z}^s(\beta) - V_{\mathbf{0},1}^s(\beta) \right). \quad (84)$$

Let us first determine $B_{\mathbf{x}^n}^{n,h}$.

We have from (20)

$$\begin{aligned} B_{\mathbf{x}^n}^{n,h} &= \lim_{\beta \rightarrow 0} \frac{(1 - C_n(\beta))(x_n + \lambda_n(1 - \tau_{\mathbf{x}^n}(\beta))/\beta) + C_n(\beta)(\tau_{\mathbf{x}^n}(\beta) - \tau_{\mathbf{x}^{n-1}}(\beta))}{\beta(1 - C_n(\beta))} \\ &\quad - \lim_{\beta \rightarrow 0} \beta H_{\mathbf{0}}^n(\beta) \left(\frac{1 - \tau_{\mathbf{x}^n}(\beta)}{\beta} \right) \\ &= \frac{x_n C_n^{(2)}(0)}{2C_n^{(1)}(0)} + \tau_{\mathbf{x}^n}^{(1)}(0) \left(1 - \frac{\lambda_n C_n^{(2)}(0)}{2C_n^{(1)}(0)} + g^{h,n} \right) + \tau_{\mathbf{x}^{n-1}}^{(1)}(0) \\ &\quad - \frac{\tau_{\mathbf{x}^n}^{(2)}(0)}{2} \left(\lambda_n + \frac{1}{C_n^{(1)}(0)} \right) + \frac{\tau_{\mathbf{x}^{n-1}}^{(2)}(0)}{2C_n^{(1)}(0)} \end{aligned} \quad (85)$$

where $g^{h,n} := \lim_{\beta \rightarrow 0} \beta H_{\mathbf{0}}^n(\beta)$ is given in (73).

With (69) and (71) the equation (85) rewrites as

$$\begin{aligned} B_{\mathbf{x}^n}^{n,h} &= \frac{x_n C_n^{(2)}(0)}{2C_n^{(1)}(0)} - \frac{\sum_{i=1}^n x_i / \mu_i}{1 - \sum_{i=1}^n \rho_i} \left(1 - \frac{\lambda_n C_n^{(2)}(0)}{2C_n^{(1)}(0)} + g^{h,n} \right) - \frac{\sum_{i=1}^{n-1} x_i / \mu_i}{1 - \sum_{i=1}^{n-1} \rho_i} \\ &\quad + \frac{\mu_n}{2} \left(1 - \sum_{i=1}^n \rho_i \right) \tau_{\mathbf{x}^n}^{(2)}(0) - \frac{\mu_n}{2} \left(1 - \sum_{i=1}^{n-1} \rho_i \right) \tau_{\mathbf{x}^{n-1}}^{(2)}(0) \end{aligned} \quad (86)$$

We are left with the evaluation of $\tau_{\mathbf{x}^n}^{(2)}(0)$.

First we compute $\gamma_n^{(1)}(0)$ and $\gamma_n^{(2)}(0)$ that will be needed in the evaluation of $\tau_{\mathbf{x}^n}^{(2)}(0)$. From (11) we find, with the use of (69) and (70),

$$\gamma_n^{(1)}(0) = -\frac{1}{\mu_n (1 - \sum_{i=1}^n \rho_i)} \quad (87)$$

$$\gamma_n^{(2)}(0) = \frac{\sum_{i=1}^{n-1} \lambda_i \mathbf{E}[S_i^2]}{\mu_n (1 - \sum_{i=1}^n \rho_i)^3} + \frac{\mathbf{E}[S_n^2] (1 - \sum_{i=1}^{n-1} \rho_i)}{(1 - \sum_{i=1}^n \rho_i)^3}. \quad (88)$$

Differentiating now twice (12) w.r.t. s and then letting $s = 0$ gives

$$\begin{aligned} \tau_{\mathbf{x}^n}^{(2)}(0) &= \left(1 - \lambda_n \gamma_n^{(1)}(0)\right)^2 \tau_{\mathbf{x}^{n-1}}^{(2)}(0) + \tau_{\mathbf{x}^{n-1}}^{(1)}(0) \left(2x_n \gamma_n^{(1)}(0) \left(1 - \lambda_n \gamma_n^{(1)}(0)\right) - \lambda_n \gamma_n^{(2)}(0)\right) \\ &\quad + x_n(x_n - 1) \gamma_n^{(1)}(0)^2 + x_n \gamma_n^{(2)}(0) \\ &= \left(\frac{1 - \sum_{i=1}^{n-1} \rho_i}{1 - \sum_{i=1}^n \rho_i}\right)^2 \tau_{\mathbf{x}^{n-1}}^{(2)}(0) + \tau_{\mathbf{x}^{n-1}}^{(1)}(0) \left(2x_n \gamma_n^{(1)}(0) \left(1 - \lambda_n \gamma_n^{(1)}(0)\right) - \lambda_n \gamma_n^{(2)}(0)\right) \\ &\quad + x_n(x_n - 1) \gamma_n^{(1)}(0)^2 + x_n \gamma_n^{(2)}(0) \\ &= \frac{1}{(1 - \sum_{i=1}^n \rho_i)^2} \sum_{i=1}^n \left[\tau_{\mathbf{x}^{i-1}}^{(1)}(0) \left(2x_i \gamma_i^{(1)}(0) \left(1 - \lambda_i \gamma_i^{(1)}(0)\right) - \lambda_i \gamma_i^{(2)}(0)\right) \right. \\ &\quad \left. + x_i(x_i - 1) \gamma_i^{(1)}(0)^2 + x_i \gamma_i^{(2)}(0) \right] \\ &= \frac{1}{(1 - \sum_{i=1}^n \rho_i)^2} \sum_{i=1}^n \left[-\frac{\sum_{k=1}^{i-1} x_k / \mu_k}{1 - \sum_{k=1}^{i-1} \rho_k} \left(2x_i \gamma_i^{(1)}(0) \left(1 - \lambda_i \gamma_i^{(1)}(0)\right) - \lambda_i \gamma_i^{(2)}(0)\right) \right. \\ &\quad \left. + x_i(x_i - 1) \gamma_i^{(1)}(0)^2 + x_i \gamma_i^{(2)}(0) \right] \end{aligned} \quad (89)$$

where by convention $\tau_{\mathbf{x}^{i-1}}^{(1)}(0) = 0$ if $i = 1$.

In conclusion, the mapping $\mathbf{x}^n \rightarrow B_{\mathbf{x}^n}^{n,h}$ is given in (86) with $C_n^{(1)}(0)$, $C_n^{(2)}(0)$, $g^{h,n}$ and $\tau_{\mathbf{x}^n}^{(2)}(0)$ given in (69), (70), (73) and (89), respectively. The function $\mathbf{x}^n \rightarrow B_{\mathbf{x}^n}^{n,h}$ appears to be a quadratic function in the variables x_1, x_2, \dots, x_n .

We now address the computation of the bias function $B_{\mathbf{x},z}^s$ associated with the switching costs.

Fix $\mathbf{x} = (x_1, \dots, x_N) \in \mathbb{N}^N$ with $\mathbf{x} \neq \mathbf{0}$. Then, there exists $i^* \in \{1, 2, \dots, N\}$ such that $\mathbf{x} = \mathbf{x}_{[i^*]}$ (see definition of the vector $x_{[i^*]}$ in Section 3).

From the definition of the bias function $B_{\mathbf{x},z}^s$ (see (84)) and Theorem 5.1, we see that

$$\begin{aligned} B_{\mathbf{x},z}^s &= \lim_{\beta \rightarrow 0} \left(V_{\mathbf{x}_{[i^*]},z}^s(\beta) - V_{\mathbf{0},1}^s(\beta) \right) \\ &= s_{z,i^*} - s_{1,i^*} + \lim_{\beta \rightarrow 0} \left(V_{\mathbf{x}_{[i^*]},1}^s(\beta) - V_{\mathbf{0},1}^s(\beta) \right) \\ &= s_{z,i^*} - s_{1,i^*} + \sum_{j=1}^N \hat{r}(j) \lim_{\beta \rightarrow 0} \left(\frac{\tau_{\mathbf{x}_{[i^*]}^j}(\beta) - 1}{\beta} \right) \end{aligned}$$

$$\begin{aligned}
& - \sum_{1 < k \leq l \leq N} \hat{r}(k, l) \left(\tau_{\mathbf{x}_{[i^*]}^{k-1}}(\lambda_k + \cdots + \lambda_l) - 1 \right) \mathbf{1}(x_k = \cdots = x_l = 0) \\
& = s_{z, i^*} - s_{1, i^*} - \sum_{j=1}^N \hat{r}(j) \left(\frac{\sum_{i=i^*}^j x_i / \mu_i}{1 - \sum_{i=1}^j \rho_i} \right) \\
& - \sum_{1 < k \leq l \leq N} \hat{r}(k, l) \left(\tau_{\mathbf{x}_{[i^*]}^{k-1}}(\lambda_k + \cdots + \lambda_l) - 1 \right) \mathbf{1}(x_k = \cdots = x_l = 0) \quad (90)
\end{aligned}$$

where

$$\begin{aligned}
\hat{r}(k, l) & := \lim_{\beta \rightarrow 0} r(k, l) \\
& = \begin{cases} s_{k, l} - s_{k, l+1} + s_{k-1, l+1} - s_{k-1, l} & \text{if } l = 1, 2, \dots, N-1 \\ s_{k, N} - s_{k-1, N} + \sum_{m=1}^N \frac{\lambda_m}{\Lambda} (s_{k-1, m} - s_{k, m}) & \text{if } l = N. \end{cases} \quad (91)
\end{aligned}$$

To derive (90) we have first used the identity $\lim_{\beta \rightarrow 0} \left(\tau_{\mathbf{x}_{[i^*]}^j}(\beta) - 1 \right) / \beta = \tau_{\mathbf{x}_{[i^*]}^j}^{(1)}(0)$, then (71).

It remains to handle the case when $\mathbf{x} = \mathbf{0}$. In this case it is easily seen from Theorem 5.1 that

$$B_{\mathbf{0}, z}^s = \frac{1}{\Lambda} \sum_{j=1}^N \lambda_j \left(s_{z, j} - s_{1, j} + B_{\mathbf{x}_{[j], 1}}^s \right) \quad (92)$$

where $B_{\mathbf{x}_{[j], 1}}^s$ can be computed from (90).

7 Applications and numerical results

We have computed both the total discounted holding and switching costs, the average expected holding and switching costs, and the bias vector, for a preemptive priority queue. These results represent an extension to existing results on priority queues and can also be used as the basis for an optimization step for a system where, for instance, switches can occur at any point in time.

Consider first general service time distributions. In a numerical experiment there are many parameters to be chosen, and the influences of changes in all these parameters can be studied. We chose to give only a limited number of numerical results. For this reason we calculated the total expected discounted holding costs for a three-queue system with all service times equally distributed. We did this experiment for three different types of service time distributions: deterministic, exponential, and a distribution that can be 0 and exponentially distributed, both with probability 1/2. The latter distribution can be seen as a hyperexponential with one of the parameters equal to ∞ . Note that the squared coefficients of variation, defined as $\sigma^2(A)/(\mathbb{E}A)^2$ for a r.v. A , is, for our choices of service time distributions, 0, 1, and 2,

respectively. We did our calculations with $(\lambda_1, \lambda_2, \lambda_3) = (0.9, 0, 0)$, all $c_i = 1$, and the expected service times equal to 1. The initial state and β are varied. The results can be found in Table 1. In this table we can see the impact of different service time distributions (by comparing the rows), the impact of different initial states (by comparing lines 1 and 2 or 3 and 4), and the impact of different discount factors (by comparing lines 1 and 3 or 2 and 4). Note that for initial state $(5,0,0)$ the system behaves as a standard $M/G/1$ queue, because $\lambda_2 = \lambda_3 = 0$. For exponential service time distributions a different initial state changes only the service order without having consequences for the total queue length. This explains the equal numbers in the one but last column. For small β the long-run behavior is dominating, and we see that a smaller variance lead to lower holding costs, which is to be expected. For $\beta = 0.1$ we see that in the short run a high variance can give lower queue lengths. This phenomenon is amplified by starting in a different state: the customers in queues 2 and 3 can be preempted, and a high variance makes it more likely that they have finished service before preemption. Also, for low values of β this effect diminishes the consequences of high variances, as can be seen by comparing lines 3 and 4. Readers who are interested in further numerical experimentation can obtain the program used to generate Table 1 (and Table 2) from the authors.

initial state	β	deterministic	exponential	hyperexponential
$(5, 0, 0)$	0.1	48.34	48.00	47.68
$(0, 2, 3)$	0.1	51.28	48.00	44.09
$(5, 0, 0)$	0.01	458.6	630.6	874.2
$(0, 2, 3)$	0.01	472.0	630.6	856.4

Table 1: Values for different types of distributions

Next, we address the application of our results in the context of the optimization of queueing systems. Before making this point more precise, consider first regular optimization procedures such as value iteration or policy iteration (e.g., [13] or [14]). Quite often these optimization procedures cannot be applied. The reason for this is the so-called *curse of dimensionality* (Bellman [2]). The size of the state space grows exponentially in the dimension of the problem. This leads, already for moderate dimensions such as 5 or 6, to state spaces for which vectors of that size cannot be contained anymore in computer memory. This single fact is the main reason why Markov decision processes have been used so little over the last few decades as operational tools. In the last couple of years we have seen a renewed interest in approximation methods. We discuss two important methods, one-step improvement and approximative dynamic programming, and show how the current results fit in.

It has long been observed ([15], [17, Sec. 3.2]) that the application of one step of the policy improvement algorithm generates a policy that is very “close” to the optimal policy. Performing one step of the policy improvement algorithm however requires the identification of a policy that will be used to initialize the procedure along with the computation of the value (cost) function associated with that policy. For high-dimensional problems this is unfeasible, again

because of the curse of dimensionality, unless one knows in explicit form the value function for some fixed policy. This is exactly the case for the preemptive priority queue studied in this paper.

To make the use of the current results more concrete, consider the problem of finding a scheduling policy that minimizes the discounted expected overall cost function (respectively the average overall cost) for a system as the priority queue, but with the possibility for the server to switch at any time from one class of customers to another class of customers. In this setting, the objective is to determine what class of customers should be attended by the server in order to minimize the costs for a given cost criterion. Because of the non-zero switching costs, the structure of the optimal policy turns out to be fairly involved as opposed to the situation where the switching costs are zero and where a simple priority policy—the celebrated μc -rule—is optimal [4, 11]. Starting from the μc -rule, which is a preemptive priority rule, it was shown in [9] and [5] that in the case of two queues one-step improvement gives excellent results. The results in this paper allow us to extend these experimentations to an arbitrary number of queues. In Table 2 we report results for systems with 3 queues, $\lambda_i = 1$ for $i = 1, \dots, 3$, $(\mu_1, \dots, \mu_3) = (6, 6, 8)$, $\beta = 0.1$, initial state (1,2,3,1) (the final 1 indicating the location of the server), and all switching costs equal to s . Each time two entries are displayed: the holding costs and the switching costs. By the optimality of the μc -rule we know that the holding costs are minimized by the priority rule, at the expense of very high switching costs. It is also remarkable how well the 1-step improved policy performs.

c	s	priority rule	1-step improved	optimal
(2, 1, 1)	2	(16.09,47.82)	(25.04,25.59)	(23.59,24.93)
(1, 1, 1)	2	(13.95,47.82)	(21.93,25.05)	(22.41,20.77)
(2, 1, 1)	0.1	(16.09,2.39)	(15.46,2.13)	(15.32,2.18)

Table 2: Total costs for different policies and direct costs

We conclude with a few words on *approximative dynamic programming* ([1, 3]). Central in this technique is the idea of approximating the value function by a function involving a few parameters and then estimating these parameters such that the difference between the value function and its approximation is minimized. The structure of the approximation (e.g. quadratic in the queue lengths) is of course crucial to the success of this method. In this paper we have identified the complete structure for a particular policy. Our findings support the usual choice of a quadratic approximation for the holding costs in the average cost case; for the average switching costs and the discounted cost criterion the structure is more complicated indicating that a more complex structure of the approximation should be considered.

References

- [1] A.G. Barto, S.J. Bradtke, and S.P. Singh, Learning to act using real-time dynamic programming. *Artificial Intelligence*, 72:81–138, 1995.

- [2] R. Bellman, *Adaptive Control Processes: A Guided Tour*. Princeton University Press, 1961.
- [3] D.P. Bertsekas and J.N. Tsitsiklis, *Neuro-Dynamic Programming*. Athena Scientific, 1996.
- [4] C. Buyukkoc, P. Varaiya and J. Walrand, The μ -rule revisited. *Advances in Applied Probability*, 17:237–238, 1985.
- [5] R. Groenevelt, G. M. Koole, and P. Nain, On the bias vector of a two-class preemptive priority queue. *Mathematical Methods of Operations Research*, 55:107–120, 2002.
- [6] J. M. Harrison, A priority queue with discounted linear costs. *Operations Research*, 23:260–269, 1975.
- [7] N. K. Jaiswal, *Priority Queues*. Academic Press, 1968.
- [8] L. Kleinrock, *Queueing Systems*. Vol. 2, John Wiley & Sons, New York, 1975.
- [9] G. M. Koole and P. Nain, On the value function of a priority queue with an application to a controlled polling model. *Queueing Systems*, 34:199–214, 2000.
- [10] S. A. Lippman, On dynamic programming with unbounded rewards. *Management Science*, 21:1225–1233, 1975.
- [11] P. Nain and D. Towsley, Optimal scheduling in a machine with stochastic varying processing rates. *IEEE Trans. on Aut. Control*, Vol. AC-39, No. 9, pp. 1853–1855, Sep. 1994.
- [12] T.J. Ott and K.R. Krishnan. Separable routing: A scheme for state-dependent routing of circuit switched telephone traffic. *Annals of Operations Research*, 35:43–68, 1992.
- [13] M. L. Puterman, *Markov Decision Processes*. Wiley, New York, 1994.
- [14] S. M. Ross, *Introduction to Stochastic Dynamic Programming*. Academic Press, 1983.
- [15] S. A. E. Sassen, H. C. Tijms and R. D. Nobel, A heuristic rule for routing customers to parallel servers. *Statistica Neerlandica*, 5:107–121, 1997.
- [16] L. Takács, *Introduction to the Theory of Queues*. Oxford University Press, 1962.
- [17] H. C. Tijms, *Stochastic Modelling and Analysis. A Computational Approach*. John Wiley & Sons, Chichester, 1986
- [18] P. D. Welch, On preemptive resume priority queues. *Ann. Math. Statist.*, 35:600–611, 1964.