

Authors are encouraged to submit new papers to INFORMS journals by means of a style file template, which includes the journal title. However, use of a template does not certify that the paper has been accepted for publication in the named journal. INFORMS journal templates are for the exclusive purpose of submitting to an INFORMS journal and should not be used to distribute the papers in print or online or to submit the papers to another publication.

Developing Effective Service Policies for Multiclass Queues with Abandonment: Asymptotic Optimality and Approximate Policy Improvement

Terry James

STOR-i Doctoral Training Centre, Lancaster University, t.james@lancaster.ac.uk

Kevin Glazebrook

Department of Management Science, Lancaster University, k.glazebrook@lancaster.ac.uk

Kyle Lin

Operations Research Department, Naval Postgraduate School, kyllin@nps.edu

We study a single server queuing model with multiple classes and impatient customers. The goal is to determine a service policy to maximise the long-run reward rate earned from serving customers net of holding costs and penalties respectively due to customers waiting for and leaving before receiving service. We first show that it is without loss of generality to study a pure-reward model. Since standard methods can usually only compute the optimal policy for problems with up to three customer classes, our focus is to develop a suite of heuristic approaches, with a preference for operationally simple policies with good reward characteristics. One such heuristic is the $R\mu\theta$ rule—a priority policy that ranks all customer classes based on the product of reward R , service rate μ , and abandonment rate θ . We show that the $R\mu\theta$ rule is asymptotically optimal as customer abandonment rates approach zero, and often performs well in cases where the simpler $R\mu$ rule performs poorly. The paper also develops an approximate policy improvement method that uses simulation and interpolation to estimate the bias function for use in a dynamic programming recursion. For systems with two or three customer classes, our numerical study indicates that the best of our simple priority policies is near optimal in most cases, and when it is not, the approximate policy improvement method invariably tightens up the gap substantially. For systems with five customer classes, our heuristics typically achieve within 4% of an upper bound for the optimal value, which is computed via a linear program that relies on a relaxation of the original system. The computational requirement of the approximate policy improvement method grows rapidly when the number of customer classes or the traffic intensity increases.

Key words: multiclass queue; customer abandonment; Markov decision process; index policy; approximate policy improvement

History: received April 2014; revised December 2014, May 2015, June 2015, July 2015; accepted August 2015.

1. Introduction

This paper considers a setting in which a single server must preemptively serve impatient customers across k customer classes. Different classes of customers arrive according to independent Poisson processes, with the arrival rate being λ_i for class i customers, $1 \leq i \leq k$. The service time for a class i customer follows an exponential distribution with rate μ_i . However, each class i customer will only remain available for service for a random time that follows an exponential distribution with rate θ_i , after which the customer will abandon the system, whether the customer is still waiting in the queue or is already in service. If a class i customer is served to completion, then a reward R_i is earned, but if he abandons the system before service completion, then a penalty D_i is incurred. In addition, each class i customer in the system incurs a linear holding cost at rate c_i per time unit. We seek to determine a service policy that maximizes the long-run reward rate earned net of penalties and holding costs incurred.

Our model has real-world applications. For security surveillance, a customer in a class corresponds to a suspect in a particular area, with service being the screening of suspects by a security resource. The penalties D_i represent the expected damage incurred when suspects leave the area and evade screening, while it is most natural to set $R_i = c_i = 0$. For call centres, the class of a customer indicates a particular service need, while the server corresponds to an agent. The rewards R_i represent the revenue received by serving a customer, while c_i and D_i respectively model the loss of goodwill incurred due to customers waiting and hang up. Customer abandonments are a key feature in both applications.

The first step of our analysis is to show that the three parameters in the reward structure—namely R_i , D_i , and c_i —can be consolidated into a single parameter through proper transformation. As it turns out, it is without loss of generality to consider a reward structure with only one of these three parameters while setting the other two equal to 0. In this paper, we choose the pure reward model (with $D_i = c_i = 0$ for all i). Although it is possible to formulate our model as a Markov Decision Process (MDP), and use standard methods of dynamic programming (DP) to compute the optimal policy by truncating the state space, the computation is usually only practical for problems with up to three

customer classes. Hence, the paper focuses on developing strongly performing heuristic policies, with a preference for operationally simple policies with strong reward characteristics.

Our first approach is to develop a suite of simple *priority policies*, which are effective across much of the problem’s parameter space. Such policies serve customers according to a strict priority ordering among the customer classes. In the case where the system is overloaded, it has been shown in the literature that the $R\mu$ rule—a priority policy that ranks all customer classes based on the product of reward R and service rate μ —performs well, since it maximizes the instantaneous reward rate (Atar et al. 2010, Ayesta et al. 2011, Verloop 2014, Larrañaga et al. 2014). To complement the $R\mu$ rule in the light-traffic case, we study the $R\mu\theta$ rule, which ranks all customer classes based on the product of R , μ , and the abandonment rate θ . This ranking was proposed in Section 2 of Glazebrook et al. (2004) for batch problems, and our paper extends its application to systems with customer arrivals. We prove that the $R\mu\theta$ rule is asymptotically optimal as customer abandonment rates approach zero in light traffic systems. Finally, we apply a pairwise-swapping (PaS) mechanism to both the $R\mu$ rule and the $R\mu\theta$ rule to search for an improved priority policy.

Our second approach—after having produced a set of simple priority policies—is to develop an effective approximate policy improvement (API) method. For a given policy, the API method uses simulation to estimate bias values for a set of carefully chosen states, and then uses these values to interpolate the bias function for all states. This approximate bias function allows us to run policy improvement to obtain a new policy. Our numerical results indicate that, in most cases, the best priority policy is nearly optimal in systems with 2 or 3 customer classes; in the cases where it is not, the API method invariably tightens up the gap substantially. In one instance, the API method improves our best priority policy—which is 4.26% suboptimal—in yielding a policy that is only 0.04% suboptimal. In some applications such as security surveillance, even a small difference in reward rate performance can be of practical importance.

In order to evaluate our heuristic approaches for systems with more than 3 customer classes, where obtaining the optimal policy is computationally too intensive, we also develop a novel upper bound for the maximal long-run reward rate using linear programming methods. The linear program relies on a relaxation of the original system, and its tightness degrades as the number of customer classes increases. In our numerical study

with 5 customer classes, our best heuristic is typically within 4% of the upper bound for the maximal long-run reward rate.

Our model has been studied in Glazebrook et al. (2004), which develops a heuristic policy via a two-stage process. The first stage analyzes a policy that allocates a fixed-service effort to each class at all time, and computes the optimal allocation, while the second stage performs an exact policy improvement (PI). The parametric optimisation involved at the first stage can pose computational challenges, especially when the number of customer classes increases, and was only implemented for cases with 2 customer classes in Glazebrook et al. (2004). Numerical tests show that our API method performs considerably better for a wide range of parameters. Down et al. (2011) studied a special case of our model with 2 customer classes, and $\mu_1 = \mu_2$. They proved that the priority policy that serves a class 1 customer whenever possible is optimal, if $R_1 \geq R_2$ and $\theta_1 \geq \theta_2$.

There are some recent works on control of queueing systems with impatient customers. An approach based on approximating Brownian control problems in heavy traffic has been explored by Harrison and Zeevi (2004), Ata and Tongarlak (2013), and Kim and Ward (2013). The latter article considers general arrival, service, and abandonment processes. Recent studies of controlled stochastic systems with customer abandonments have featured a range of application domains. Garnett et al. (2002), Jouini et al. (2010), and Bassamboo et al. (2005) consider call centre applications. Abandonments have been used in military applications to model targets that move out of range of defensive forces. See, for example, Gaver et al. (2006), Lin et al. (2009), and Glazebrook and Punton (2008). In patient triage applications, abandonments have been used to model medical emergency patients in danger of dying while awaiting treatment. See, for example, Argon et al. (2008) and Li and Glazebrook (2010).

There are also some recent works on approximate approaches to DP seeking to overcome computational intractability. See, for example, Section 1.2 of Powell (2011). Contributions that deploy value function approximations within a PI approach include those of Krishnan (1987), Glazebrook et al. (2004), and Li and Glazebrook (2010), while API methods which utilise simulation are discussed in Section 10.5 of Powell (2011) and Chapters 6 and 7 of Bertsekas (2012). Our API approach can be viewed as a refined approximate dynamic programming (ADP) implementation with two distinctive features: (1) a suite

of strongly-performing priority policies to initialise the API algorithm, and (2) a simulation/interpolation methodology to fit the bias surface by estimating biases both at states that are frequently visited, and also at a carefully chosen set of widely spread states.

The rest of the paper proceeds as follows. Section 2 first shows that it is without loss of generality to consider a pure reward model, and then formulates the problem as an MDP and describes how, in principle, to compute the optimal policy. Section 3 presents the service policies based on class prioritisation, including the $R\mu\theta$ rule, the $R\mu$ rule, and a potentially improved priority policy achieved by pairwise swapping. Section 4 describes the API method. Section 5 presents a linear program to compute an upper bound for the maximal long-run reward rate in order to evaluate our heuristics where the optimal policy is not available. Section 6 offers a numerical study, and Section 7 concludes the paper.

2. Model and MDP Formulation

Recall that our model has three reward parameters. For a class i customer, there is a reward R_i for service completion, a penalty D_i for customer abandonment, and a linear holding cost rate c_i per time unit, for $1 \leq i \leq k$. If we write $N_{i,\theta}^\pi$ for the number of class i customers in the system and α_i^π , β_i^π for, respectively, the rate of class i service completions and abandonments under policy π in steady state, then the optimal long-run system reward rate net of holding costs and abandonment penalties can be written as

$$\max_{\pi} \sum_{i=1}^k (R_i \alpha_i^\pi - D_i \beta_i^\pi - c_i E[N_{i,\theta}^\pi]). \quad (1)$$

However, the guaranteed stability of the system implies that, for all choices of class i and policy π , we have that

$$\lambda_i = \alpha_i^\pi + \beta_i^\pi, \quad (2)$$

and

$$\beta_i^\pi = \theta_i E[N_{i,\theta}^\pi]. \quad (3)$$

Using (2) and (3), we can rewrite (1) in three different ways:

$$\max_{\pi} \sum_{i=1}^k \left(\left(R_i + D_i + \frac{c_i}{\theta_i} \right) \alpha_i^\pi - \left(D_i + \frac{c_i}{\theta_i} \right) \lambda_i \right) \quad (4)$$

$$= \max_{\pi} \sum_{i=1}^k \left(R_i \lambda_i - \left(R_i + D_i + \frac{c_i}{\theta_i} \right) \beta_i^\pi \right) \quad (5)$$

$$= \max_{\pi} \sum_{i=1}^k (R_i \lambda_i - ((R_i + D_i)\theta_i + c_i) E[N_{i,\theta}^\pi]). \quad (6)$$

Equation (4) transforms the original model into an equivalent pure-reward model with $R_i + D_i + c_i/\theta_i$ earned upon every class i service completion. Similarly, equation (5) shows an equivalent model with only penalties upon customer abandonment, while equation (6) shows an equivalent model with only linear holding costs. Without loss of generality, we shall focus on the pure-reward model ($D_i = c_i = 0$ for all i) for the remainder of the paper.

Denote the system state by $\mathbf{n} = (n_1, \dots, n_k)$, with n_i the number of class i customers present in the system. We further write $\mathbf{n}(t)$ for the system state at time t . Further details of the model are as follows:

1. Decision epochs occur at time zero and at all transitions of the system state.
2. At each decision epoch, the server must decide which waiting customer to serve next across all customer classes. The set of admissible actions for state \mathbf{n} is given by

$$A(\mathbf{n}) = \{a : n_a \geq 1, 1 \leq a \leq k\}.$$

We use \mathbf{e}_i for the system state in which only a single customer of class i is present in the system.

3. In state $\mathbf{n} \neq \mathbf{0}$ under admissible action $a \in A(\mathbf{n})$, the effective transition rate is $\Lambda(\mathbf{n}, a) = \mu_a + \sum_{i=1}^k (\lambda_i + n_i \theta_i)$. Transitions to states $\mathbf{n} + \mathbf{e}_i$, $\mathbf{n} - \mathbf{e}_a$, and $\mathbf{n} - \mathbf{e}_j$, $j \neq a$, respectively, occur with probabilities $\lambda_i \{\Lambda(\mathbf{n}, a)\}^{-1}$, $(\mu_a + n_a \theta_a) \{\Lambda(\mathbf{n}, a)\}^{-1}$, and $n_j \theta_j \{\Lambda(\mathbf{n}, a)\}^{-1}$. The effective transition rate in the empty state $\mathbf{0}$ is $\Lambda(\mathbf{0}) = \sum_{i=1}^k \lambda_i$ with a transition from $\mathbf{0}$ to state \mathbf{e}_i occurring with probability $\lambda_i \{\Lambda(\mathbf{0})\}^{-1}$. When a transition from \mathbf{n} to $\mathbf{n} - \mathbf{e}_a$ occurs at a class a service completion, a reward R_a is earned.

4. A service policy is a rule for choosing admissible actions using the history of the process (past states and actions) only. An *admissible, deterministic, stationary, and Markov* policy is determined by a function $\pi : \mathbb{N}^k \rightarrow \{1, \dots, k\}$ satisfying $\pi(\mathbf{n}) \in A(\mathbf{n}), \forall \mathbf{n}$. The theory of MDPs (see, for example, Chapter 8 of Puterman (1994)) implies that, to determine the optimal policy, it is sufficient to consider only policies in this class.

5. The goal of analysis is to determine a policy that maximises the long-run reward rate earned or that will come close to doing so.

A standard approach to determine the ϵ -optimal policies is through the application of DP to a version of the above system with finite state space $\times_{i=1}^k \{0, 1, \dots, N_i\}$. In this truncated version, new class i customers are blocked from entering the system when N_i are already present. The N_i must be chosen large enough to ensure that this system

approximates the original model well enough for the purpose at hand. With a finite state space, it becomes possible to convert the problem to one in discrete time through the process of *uniformisation*. We write $\Delta = \sum_{i=1}^k (\lambda_i + \mu_i + N_i \theta_i)$, a uniform upper bound on the rate of state transitions in the finite state system. By the addition of fictitious transitions from a state to itself, we develop a uniformised system that makes transitions at a uniform rate Δ . We write $V^\pi(\mathbf{n}, t)$ and $V(\mathbf{n}, t)$ for the expected reward earned under the application of policy π and an optimal policy, respectively, over t transitions of the uniformised process, beginning at time zero in system state \mathbf{n} . Standard theory enables us to write $V^\pi(\mathbf{n}, t) = \frac{g^\pi}{\Delta} t + \omega^\pi(\mathbf{n}) + o(1)$ and $V(\mathbf{n}, t) = \frac{g}{\Delta} t + \omega(\mathbf{n}) + o(1)$ as $t \rightarrow \infty$, where g^π and g are the long-run reward rates or gains earned, and ω^π and ω the bias functions under application of π and an optimal policy, respectively. Bias functions yield an estimate of the transient effect on rewards of the starting state \mathbf{n} and will be further discussed in Section 4. Bellman’s equation for the finite state system can now be written

$$\frac{g}{\Delta} + \omega(\mathbf{n}) = \max_a \left\{ \frac{R_a \mu_a}{\Delta} + \sum_{\mathbf{n}' \in \mathcal{S}} p(\mathbf{n}' | \mathbf{n}, a) \omega(\mathbf{n}') \right\}, \quad (7)$$

where the $p(\mathbf{n}' | \mathbf{n}, a)$ are transition probabilities under the uniformisation. It is now possible to compute the optimal gain and associated optimal policy for the finite state approximation by a recursive scheme such as DP value iteration or by linear programming; further details may be found in Chapter 8 of Puterman (1994). However, since the state space grows exponentially in k , in practice, the computations quickly become intractable for $k \geq 4$. Hence, the focus of our paper is to develop near-optimal heuristic policies that require much less computation.

3. Service Policies Based on Class Prioritisation

A policy that is easy to implement is for the server to prioritize all customer classes in an ordered list, and always serve a customer highest on the list among all customers present in the system. In the case of an overloaded system, there are almost always many customers present in the system. It is therefore intuitive that the server should pay little attention to the possibility of idling, and focus on continuously maximizing the instantaneous reward rate. To do so, the server simply needs to always serve a customer having the maximal $R\mu$ value among all customers present in the system—hence the $R\mu$ rule. As seen in equations (4) and (6), the $R\mu$ rule in our pure-reward model is equivalent to the $c\mu/\theta$ rule

in the linear-holding-cost only model. The strong performance of this rule in heavy traffic is affirmed in the work of Atar et al. (2010), Ayesta et al. (2011), Verloop (2014), and Larrañaga et al. (2014).

Away from heavy traffic, lost reward opportunities due to an empty system become a much more important concern. Motivated by this observation, Section 3.1 introduces the $R\mu\theta$ rule, and establishes its asymptotic optimality. Section 3.2 compares the $R\mu\theta$ and the $R\mu$ rules. Section 3.3 presents a mechanism to explore local improvements on a given priority policy.

3.1. The $R\mu\theta$ Rule

If a system is not overloaded with customers, then it becomes important to take into account the lost reward opportunities when the system becomes empty due to customer abandonment. For example, consider a two-class system, with $R_1\mu_1 = R_2\mu_2$ and $\theta_1 < \theta_2$. If there is one customer present from each class, then intuition suggests that the server should first serve the class 2 customer, since there is a better chance that the class 1 customer will still be available later on. Consequently, a class's priority should go up as its abandonment rate θ increases. We call the rule in which the server always serves a customer having the maximal $R\mu\theta$ value among all customers present in the system, the $R\mu\theta$ rule. As seen in equations (4) and (6), the $R\mu\theta$ rule in our pure-reward model is equivalent to the $c\mu$ rule in the linear-holding-cost only model. Whereas the $c\mu$ rule is optimal in queueing systems with no customer abandonment (see, for example, Section 5.2 in Gittins et al. (2011)), it is not optimal in systems with customer abandonment (Down et al. 2011).

The main result of this section is to show that the $R\mu\theta$ rule is asymptotically optimal as $\theta \rightarrow 0$. First, write $R^\pi(\theta)$ for the reward rate achieved by policy π , and $R^{R\mu\theta}(\theta)$ for the reward rate achieved by the $R\mu\theta$ rule. To describe the limiting regime, we suppose that the abandonment rate of each customer class is the multiple of some underlying rate θ , such that $\theta_i = \theta\nu_i$, where $\nu_i > 0$, $1 \leq i \leq k$.

THEOREM 1. *If $\sum_{i=1}^k \lambda_i/\mu_i < 1$, then*

$$\max_{\pi} R^\pi(\theta) - R^{R\mu\theta}(\theta) \leq O(\theta^2).$$

The proof of Theorem 1 is given in full in Appendix A in the Online Supplement, but we summarise the key elements here to facilitate the subsequent discussion. The main idea of the proof is to bound $R^\pi(\theta)$ below for priority policies and above for general nonidling policies. If we write W_i^π for the waiting time (time to achieve completed service) of a class i job in steady state under π for the no abandonment case ($\theta = 0$), then for a priority policy ϖ , we show that under the conditions of the result,

$$R^\varpi(\theta) \geq \sum_{i=1}^k \lambda_i R_i - \theta \sum_{i=1}^k \lambda_i R_i \nu_i E[W_i^\varpi] + O(\theta^2), \quad (8)$$

while for all nonidling policies π , we have that

$$R^\pi(\theta) \leq \sum_{i=1}^k \lambda_i R_i - \theta \sum_{i=1}^k \lambda_i R_i \nu_i E[W_i^\pi] + O(\theta^2). \quad (9)$$

It must then follow that

$$\max_{\pi} R^\pi(\theta) - R^{R\mu\theta}(\theta) \leq \theta \left\{ \sum_{i=1}^k \lambda_i R_i \nu_i E[W_i^{R\mu\theta}] - \min_{\pi} \sum_{i=1}^k \lambda_i R_i \nu_i E[W_i^\pi] \right\} + O(\theta^2), \quad (10)$$

where the maximum and minimum in (10) are over all policies π . By Little's law, the minimisation in (10) is of a holding cost rate objective for the system without abandonments. A classical queueing control result (the $c\mu$ rule) implies that this minimum is achieved by the $R\mu\theta$ rule. Theorem 1 now follows easily.

In heavy traffic, the $R\mu$ rule appropriately greedily chooses processing actions to maximise the instantaneous reward rate achieved. In the regime of Theorem 1, the focus is on choices of policy π that minimise reward rate loss from the system through abandonments. From the preceding proof, this loss rate is given by $\theta h^\pi + O(\theta^2)$, where

$$h^\pi = \sum_{i=1}^k \lambda_i R_i \nu_i E[W_i^\pi].$$

The strong performance of the $R\mu\theta$ rule resides in its minimisation of the dominant $O(\theta)$ component of this loss rate—a consequence of the optimality of the $c\mu$ rule for linear holding costs in the absence of abandonments.

Close inspection of the proof of Theorem 1 in the Online Supplement will reveal that we make little use of the stochastic structure of the system's service mechanism. The $R\mu\theta$ rule emerges as a priority policy which minimises a holding cost-type objective for the no

abandonment system ($\theta = 0$). It is therefore unsurprising that the result can be generalised to more complex service situations, provided that a priority policy continues to optimise an appropriate holding cost. An important model class to which a natural extension of Theorem 1 applies are *Klimov Networks* (see Klimov (1974) and Klimov (1978)), in which each customer service has a sequence of phases, with movement between phases and toward service completion determined by customer class-specific Markovian routing matrices. This class *inter alia* provides an extension of Theorem 1 to a model in which service requirements are independent, have finite second moment, and are identically distributed within each class. Further details can be found in Appendix B in the Online Supplement.

All of the extensions to Theorem 1 mentioned above concern single-server systems. If we move to a multiserver version of our system with abandonments, with m servers working in parallel, then the required stability condition becomes $\rho = \sum_{i=1}^k \lambda_i / \mu_i < m$ and the $R\mu\theta$ rule now allocates preemptive service to the m customers present in the system whose associated $R\mu\theta$ are maximal. The proof of a suitable version of Theorem 1 for this system goes through up to (10). However, it is no longer true that the $R\mu\theta$ rule achieves the minimum in (10), though it does come close to doing so. To give a theoretical result for this system we need the quantity

$$B(m) = \rho(R\mu\nu)_{\max} \left(\frac{1}{\mu} \right)_{\max} I(m > 1),$$

where I is an indicator and the maxima in the expression are taken over the customer classes. The following result makes use of Theorem 3 in Glazebrook and Niño-Mora (2001), which shows that $B(m)$ bounds above the quantity multiplying θ on the right-hand side of (10) when there are m servers. It generalises Theorem 1 to multiserver systems.

PROPOSITION 1. *When there are m servers and $\sum_{i=1}^k \lambda_i / \mu_i < m$, we have that*

$$\max_{\pi} R^{\pi}(\theta) - R^{R\mu\theta}(\theta) \leq \theta B(m) + O(\theta^2).$$

3.2. Comparing the $R\mu\theta$ and $R\mu$ Rules

It follows from calculations in the proof of Theorem 1 that when $\sum_{i=1}^k \lambda_i / \mu_i < 1$, we have $R^{\pi}(\theta) \rightarrow \sum_{i=1}^k \lambda_i R_i$, as $\theta \rightarrow 0$, for all priority policies (and hence both the $R\mu\theta$ and $R\mu$ rules). Unsurprisingly, all priority policies achieve the maximal reward rate $\sum_{i=1}^k \lambda_i R_i$ in

the no abandonment limit, since in the limit all jobs are served. Think of a surveillance problem in which abandonments of the system are rare, but very damaging, and attention focuses on making the $O(\theta)$ loss rate from abandonments as small as possible. Now consider a situation in which the class orderings determined by the $R\mu\theta$ and $R\mu$ rules are distinct. It follows from (8) and (9) that

$$R^{R\mu\theta}(\theta) - R^{R\mu}(\theta) = \theta(h^{R\mu} - h^{R\mu\theta}) + O(\theta^2).$$

When $R\mu\theta$ and $R\mu$ are distinct, the quantity that multiplies θ in the above expression is strictly positive. Consequently, there exists θ^* , such that for $\theta < \theta^*$, we have that $R^{R\mu\theta}(\theta) > R^{R\mu}(\theta)$. Therefore,

$$\frac{\max_{\pi} R^{\pi}(\theta) - R^{R\mu\theta}(\theta)}{\max_{\pi} R^{\pi}(\theta) - R^{R\mu}(\theta)} \rightarrow 0, \quad \text{as } \theta \rightarrow 0.$$

It follows that the percentage loss of reward rate due to abandonment from the use of $R\mu\theta$ relative to that experienced from the use of $R\mu$ becomes negligible in the limit $\theta \rightarrow 0$. Numerical support for this conclusion can be found in Appendix C in the Online Supplement.

A similar conclusion can be drawn for a multiserver system, if the condition in Proposition 1 is met. It follows from the analysis in Glazebrook and Niño-Mora (2001) that, for any case in which $R\mu$ and $R\mu\theta$ rules differ in their choice of lowest priority customer class, the difference $h^{R\mu} - \min_{\pi} h^{\pi}$ diverges in the limit. It then follows from (8), (9), and Proposition 1 that, when θ is small for problems close to the $\rho \rightarrow m$ limit, the suboptimality gap $\max_{\pi} R^{\pi}(\theta) - R^{R\mu\theta}(\theta)$ will be negligible compared to $\max_{\pi} R^{\pi}(\theta) - R^{R\mu}(\theta)$.

In the special case with $k = 2$, $\mu_1 = \mu_2$, $R_1 \geq R_2$, and $\theta_1 \geq \theta_2$, both the $R\mu\theta$ and $R\mu$ rules prescribe the priority policy $1 \rightarrow 2$, which is optimal according to Down et al. (2011).

3.3. The PaS Class of Priority Policies

We conclude this section by describing a simple *pairwise-swapping* mechanism to explore local improvements in any given priority policy. Given any class ordering $(\pi_1, \pi_2, \dots, \pi_k)$, we take the classes in order from π_2 to π_k and explore, in turn, whether each class should be promoted up the order. This is achieved for each class by a sequence of pairwise comparisons with the next highest class in the list to determine how high up the list the class can be promoted. In comparing classes i and j , we consider the two-class subsystem comprising them alone (with their class parameters inherited from the full problem) and use value

iteration to compute the respective performance of the two priority policies $i \rightarrow j$ and $j \rightarrow i$. If the better policy contradicts the current class ordering, then a pairwise swap is performed, and the procedure is repeated until a comparison does not result in a swapping. We then examine the potential promotion for the next class on the original list π_2 to π_k . We label this priority policy PaS.

4. Approximate Policy Improvement Algorithm

This section introduces an approximate policy improvement (API) algorithm to improve our suite of simple priority policies. Section 4.1 overviews the methodology, and Section 4.2 discusses the algorithm in detail.

4.1. Heuristic Based on Policy Improvement

Policy improvement (PI) develops optimal policies for MDPs by using the DP recursion to produce a sequence of successively improving policies (Howard 1960). In our problem, we truncate the state space and uniformise, as in Section 2, to develop an ergodic system with optimality equation in (7). To develop a PI step from policy π , let $\omega^\pi(\mathbf{n})$ be the bias associated with system state \mathbf{n} under policy π . A new policy PI_π , say, is obtained as follows:

$$\text{PI}_\pi(\mathbf{n}) = \arg \max_a \left\{ \frac{R_a \mu_a}{\Delta} + \sum_{\mathbf{n}' \in \mathcal{S}} p(\mathbf{n}' | \mathbf{n}, a) \omega^\pi(\mathbf{n}') \right\}. \quad (11)$$

Accordingly, policy PI_π always takes the current decision optimally, given that all future decisions are made according to π . Tijms (1994) noted that the first few PI iterations usually yield the greatest improvement.

The challenge to implementation of PI in large systems lies in the intractability of the computation of the bias ω^π . Hence, approximations are required, and the PI step in (11) can be replaced by

$$\text{API}_\pi(\mathbf{n}) = \arg \max_a \left\{ \frac{R_a \mu_a}{\Delta} + \sum_{\mathbf{n}' \in \mathcal{S}} p(\mathbf{n}' | \mathbf{n}, a) \tilde{\omega}^\pi(\mathbf{n}') \right\}, \quad (12)$$

where $\tilde{\omega}^\pi$ approximates ω^π .

Computation of the bias ω^π involves specification of a reference state \mathbf{m} , which we take to be one frequently visited under π . We introduce the following quantities:

- $r^\pi(\mathbf{n})$ is the expected reward received starting from state \mathbf{n} until the system enters the reference state \mathbf{m} for the first time, if policy π is used.

- $t^\pi(\mathbf{n})$ is the expected time starting from state \mathbf{n} until the system enters the reference state \mathbf{m} for the first time, if policy π is used.

The system evolving under policy π is ergodic, so $r^\pi(\mathbf{n})$ and $t^\pi(\mathbf{n})$ are guaranteed to be finite for all states. Using the fact that the system regenerates upon entry to the reference state, the theory of regenerative processes (Tijms 1994) indicates that

$$\omega^\pi(\mathbf{n}) = r^\pi(\mathbf{n}) - g^\pi t^\pi(\mathbf{n}), \quad (13)$$

where g^π is the gain of policy π .

From (13), the approximations $\tilde{\omega}^\pi(\mathbf{n})$ can then be obtained by approximating the quantities $r^\pi(\mathbf{n})$, $t^\pi(\mathbf{n})$, and g^π . The heuristic policy can then be defined from (12).

4.2. The Algorithm

The implementation of an API step depends crucially on the approximation scheme used for the bias function. As the bias function does not have an analytical form, we use simulation to estimate it. However, since simulation carries a computational cost, our constrained computational resource needs to be effectively managed through a carefully designed algorithm. The algorithm consists of five sequential, complementary stages, taking an initial policy π as an input to produce a new policy API_π . The five steps are summarized below, with more details to follow.

1. Pilot: Simulate the steady state of initial policy π to estimate its gain, and the frequency each state is visited.
2. Selection: Based on the pilot run, select a set of states at which we estimate the bias function via simulation.
3. Sampling: For each state \mathbf{n} selected, simulate the system under π from that state until some chosen reference state \mathbf{m} is entered and estimate $\omega^\pi(\mathbf{n})$ using (13).
4. Interpolation: Use the simulation results for selected states to interpolate the bias function for all other unselected states.
5. Improvement: Use (12) to produce a new policy API_π .

In step 1, we run a pilot steady state simulation to estimate the gain g^π required to estimate ω^π from (13). To facilitate steps 2 and 3, we also collect data on how often each state is visited in steady state under π .

Step 2 consists of the selection of a small number of states, denoted S_{sel} , at which the bias ω^π will be estimated by simulation in step 3. Interpolation of ω^π at other states will follow

in step 4. The set S_{sel} consists of the *anchor set* together with a *support set*. The anchor set consists of the states most frequently visited in the pilot and hence influential to policy performance. However, anchor states are likely to be tightly grouped together, so alone they will not create an adequate basis for the construction of an effective interpolation scheme. The support set will complement the anchor set to ensure adequate coverage and wider exploration of the state space.

To select M support states, we adopt lattice points of the following form:

$$P_M = \{((\mathbf{z}j \bmod M)/M) = ((z_1j \bmod M)/M, \dots, (z_kj \bmod M)/M) \mid 0 \leq j \leq M - 1\},$$

where \mathbf{z} is an integer vector modulo M . The components of \mathbf{z} are chosen to be relatively prime to each other and to M . In what follows, policies will be constructed for numerous problems with $k = 2, 3$, and 5 making use of the choices $\mathbf{z} = (2, 3)$, $(2, 3, 5)$, and $(2, 3, 5, 7, 11)$, respectively. These lattice points are then appropriately scaled and rounded from the unit hypercube to the state space to obtain the support states. Such well-spread points were proposed in the field of Quasi-Monte Carlo methods for numerical integration and shown to enable good approximations of integrals (Niederreiter 1978).

In step 3, we choose reference state \mathbf{m} to be the one most visited in the pilot and use Monte Carlo simulation to estimate $r^\pi(\mathbf{n})$ and $t^\pi(\mathbf{n})$ for each $\mathbf{n} \in S_{\text{sel}}$. In what follows, we use n for the size of S_{sel} and m for the number of simulated realisations of the system from each $\mathbf{n} \in S_{\text{sel}}$ until entry into reference state \mathbf{m} . If we write $R^\pi(\mathbf{n})$ and $T^\pi(\mathbf{n})$ for the simulation-based estimators of $r^\pi(\mathbf{n})$ and $t^\pi(\mathbf{n})$, respectively, and G^π for the estimator of g^π available from the pilot, then from (13) our estimator of $\omega^\pi(\mathbf{n})$ for $\mathbf{n} \in S_{\text{sel}}$ is $\Omega^\pi(\mathbf{n}) = R^\pi(\mathbf{n}) - G^\pi T^\pi(\mathbf{n})$. Since all estimators are unbiased, and G^π is independent of $R^\pi(\mathbf{n})$ and $T^\pi(\mathbf{n})$, we conclude by conditioning on G^π , that

$$\text{Var}\{\Omega^\pi(\mathbf{n})\} = \text{Var}\{R^\pi(\mathbf{n}) - g^\pi T^\pi(\mathbf{n})\} + \text{Var}(G^\pi) E[(T^\pi(\mathbf{n}))^2]. \quad (14)$$

Equation (14) decomposes the variance of the bias estimators into two terms. The first term is controlled by the number of replicates m used in the simulation relating to state \mathbf{n} in step 3, and the second term is controlled by the size of the pilot study in step 1. The computational challenge is dominated by the need to control the first term in (14), as designing a pilot study large enough to control the second term has not proved to be an issue. One feature that helps reduce the first term is that $R^\pi(\mathbf{n})$ and $T^\pi(\mathbf{n})$ are positively

associated. In addition, our choice of reference state \mathbf{m} means that the biases at anchor states (with smaller $R^\pi(\mathbf{n})$ and $T^\pi(\mathbf{n})$) tend to be estimated with greater precision than those at support states, which is a feature shared with other approaches to ADP (see Section 10.10 of Powell (2011)). The central trade-off for the quality of the method for given computational effort is that between large n supporting the quality of the interpolation, and large m supporting precision at the selected states.

In step 4, we use the bias estimates in S_{sel} to interpolate a bias function approximation for the entire state space S . While there are many interpolation algorithms, we use the *radial basis function* method (see Powell (1987)) for its simplicity. Assume that some function $f : S \rightarrow \mathbb{R}$ has known values at each $\mathbf{x}_i \in S_{\text{sel}}$. An augmented radial basis function $h : S \rightarrow \mathbb{R}$ which takes the form

$$h(\mathbf{x}) = \sum_{i=1}^n \alpha_i \phi(\|\mathbf{x} - \mathbf{x}_i\|) + \sum_{j=1}^d \beta_j p_j(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^k, \quad (15)$$

will be designed as a smooth interpolator of f , taking the values $f(\mathbf{x}_i)$ for $\mathbf{x}_i \in S_{\text{sel}}$. From (15), $h(\mathbf{x})$ is a weighted sum of $n = |S_{\text{sel}}|$ radial basis functions $\phi(\cdot)$, one centred on each $\mathbf{x}_i \in S_{\text{sel}}$, together with d low order polynomials $p_j(\cdot)$. Note that $\|\cdot\|$ denotes the Euclidean norm. For $\phi(\cdot)$, we take the *thin plate spline* $\phi(r) = r^2 \log(r)$; for low order polynomials, we set $d = k + 1$ and use $p_1(\mathbf{x}) = 1$, $p_j(\mathbf{x}) = x_{j-1}$, $2 \leq j \leq k + 1$. These choices produce a surface which minimises a measure of smoothness (Powell 1999).

We write A for the $n \times n$ matrix with elements $A_{ij} = \phi(\|\mathbf{x}_i - \mathbf{x}_j\|)$, $1 \leq i, j \leq n$ and P for the $n \times (k + 1)$ matrix with elements $P_{ij} = p_j(\mathbf{x}_i)$, $1 \leq i \leq n$, $1 \leq j \leq k + 1$. We write \mathbf{f} for the n -vector with $f_i = f(\mathbf{x}_i)$, $1 \leq i \leq n$. Let $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$ be corresponding vectors of coefficients. The matrix form of the interpolation problem is

$$\begin{pmatrix} A & P \\ P^T & 0 \end{pmatrix} \begin{pmatrix} \boldsymbol{\alpha} \\ \boldsymbol{\beta} \end{pmatrix} = \begin{pmatrix} \mathbf{f} \\ \mathbf{0} \end{pmatrix}.$$

The equations $A\boldsymbol{\alpha} + P\boldsymbol{\beta} = \mathbf{f}$ ensure that $h(\mathbf{x}_i) = f(\mathbf{x}_i)$, $\mathbf{x}_i \in S_{\text{sel}}$, while the $k + 1$ equations $P^T\boldsymbol{\alpha} = \mathbf{0}$ take up the extra degrees of freedom in the problem, which ensures the radial basis function $h(\cdot)$ is conditionally positive definite and the interpolation problem solvable. Consequently, the interpolation matrix delivers a unique solution in the coefficients and hence in h . If we take $f(\mathbf{x}_i)$, $\mathbf{x}_i \in S_{\text{sel}}$, in the above to be the estimates of bias from step 3, we can then use the resulting $h(\mathbf{x})$, $\mathbf{x} \in S$, as bias estimates for all states.

In step 5, we design a new policy API_π by using the function h from step 4 in place of $\tilde{\omega}^\pi$ in (12) to obtain

$$\text{API}_\pi(\mathbf{n}) = \arg \max_a \left\{ \frac{R_a \mu_a}{\Delta} + \sum_{\mathbf{n}' \in S} p(\mathbf{n}' | \mathbf{s}, a) h(\mathbf{n}') \right\}.$$

In principle, the preceding procedure can be repeated multiple times. Although obtaining progressively better policies—a feature of exact PI—can no longer be guaranteed, we have found that, in practice, improvement in policy performance is indeed achieved. To highlight key design choices, we denote the preceding procedure by $\text{API}(\pi, n, m, r, t)$. The parametrising arguments offer great flexibility and are as follows: the initial policy π , n the number of states selected for bias estimation via simulation, m the number of replicated simulations at each selected state, r the fraction of selected states in the anchor set (so $1 - r$ the fraction in the support set), and t the number of iterations of the algorithm. In what follows, we write API_π for the best-performing policy from t iterations of the algorithm, including the initial policy, which ensures that we only consider policies which improve as t increases. The trade-off between different choices of the parameters will be explored in Section 6.2, where we will give a recommendation for their selection.

We now present an example to illustrate the algorithm. Consider a $k = 2$ example with the parameters: $\lambda_1 = 2.5$, $\lambda_2 = 3$, $\mu_1 = 3.5$, $\mu_2 = 4$, $\theta_1 = 0.75$, $\theta_2 = 2.5$, $R_1 = 2.5$, $R_2 = 1.7$. We use truncation levels $N_1 = N_2 = 20$ throughout. Please note that for this example the $R\mu$ rule gives priority to class 1, while the $R\mu\theta$ rule gives priority to class 2. We use algorithms of the form $\text{API}(R\mu\theta, 45, m, \frac{32}{45}, 1)$ to construct policies. Figure 1 illustrates the selection and interpolation stages of the algorithm for the case $m = 10^5$. Figure 1(a) shows S_{sel} , with anchor states shown as diamonds and support states as circles. Figure 1(b) shows the interpolated bias estimates over the entire state space. Although not shown here, the surfaces of exact biases ω^π and simulated bias estimates throughout S closely resemble the interpolated surface, capturing its shape and curvature well, especially so around the anchor set. Figure 2 shows the actions taken in each state by the optimal policy for this example, along with the actions resulting from use of the above algorithm with m set at 10^3 , 10^4 , and 10^5 . We observe that as m increases, the corresponding policies approach more closely the switching curve structure of the optimal policy.

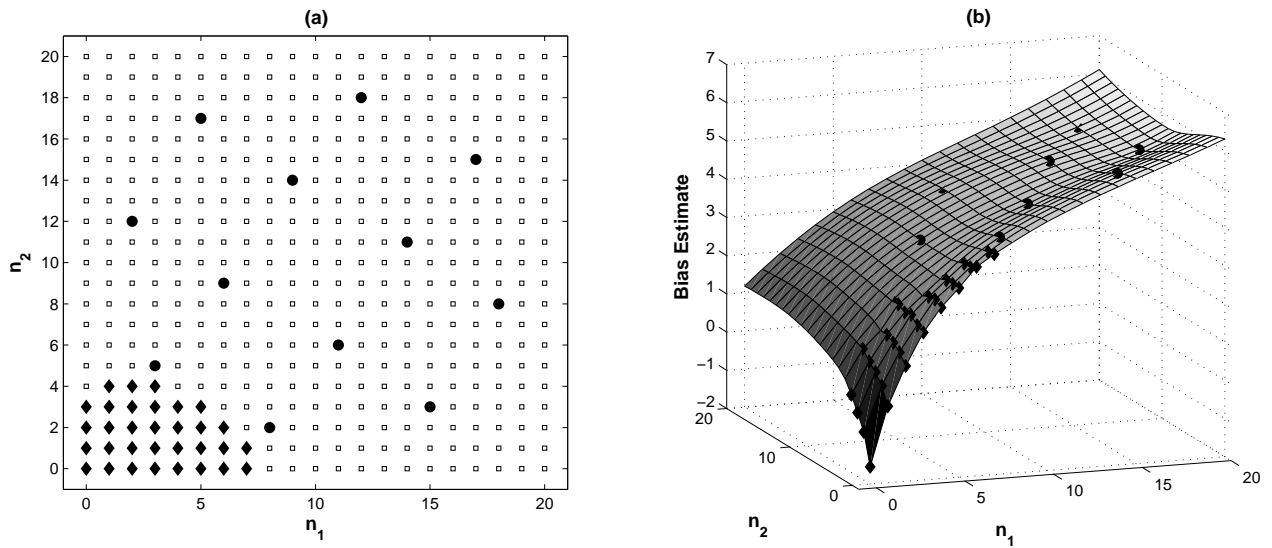


Figure 1 Illustration of the selection and interpolation stages of the algorithm in the example system.

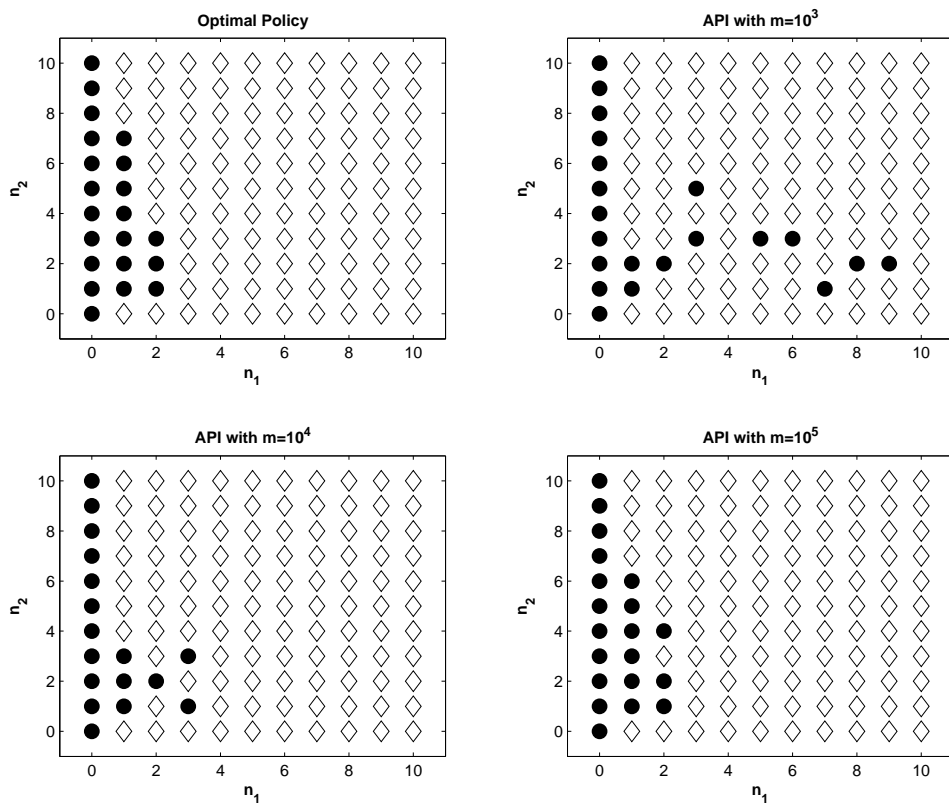


Figure 2 Class 1 (diamonds) and class 2 (circles) actions in each state under various policies in the example system.

5. An Upper Bound on Achievable Rewards

In order to evaluate heuristic policies when the optimal solution in (7) is not available, we derive an upper bound for the long-run reward rate. For a given feasible policy, if x_i

represents the implied fraction of time the server spends on class i customers, then

$$\sum_{i=1}^k R_i \mu_i x_i \quad (16)$$

is the long-run reward rate for the feasible policy. To compute an upper bound for the optimal long-run reward rate, we formulate a linear program with the variables $x_i \geq 0$, $1 \leq i \leq k$, and the objective function to maximize (16), subject to the constraint $\sum_{i=1}^k x_i \leq 1$. The key to get a tight upper bound is to impose additional constraints on the x_i so that the resulting optimal policies come as close as possible to those implied by a feasible policy.

First, denote by $A_{\{i\}}$ the long-run fraction of time the server is busy if he serves only class i customers and ignores all other classes, $1 \leq i \leq k$. Taking the number of class i customers as the state, we have a birth-and-death process, so it is straightforward to compute

$$A_{\{i\}} = 1 - \left[\sum_{n=0}^{\infty} (\lambda_i)^n \left\{ \prod_{m=1}^n (\mu_i + m\theta_i) \right\}^{-1} \right]^{-1}, \quad 1 \leq i \leq k.$$

We can add $x_i \leq A_{\{i\}}$ as a constraint in the aforementioned linear program, $1 \leq i \leq k$, or a total of k constraints.

To extend this idea, for $T \subseteq \{1, \dots, k\}$, we can add a constraint $\sum_{i \in T} x_i \leq A_T$, where A_T denotes the maximal long-run fraction of time that the server serves customer classes in T by ignoring all other classes. To compute A_T , consider the same MDP model in Section 2 with customer class set T , and substitute $R_i = \mu_i^{-1}$, $i \in T$, so that the long-run reward rate becomes equivalent to the long-run fraction of time that the server is busy. Using DP value iteration to compute the optimal solution when $|T| = 2$ or $|T| = 3$ is computationally viable, resulting in $\binom{k}{2} + \binom{k}{3}$ additional constraints.

Computing A_T when $|T| \geq 4$ is computationally infeasible, but we can still impose constraints derived from relaxed systems. To do so, we create a single fictitious class by aggregation and relaxation of the customer classes in T . Denote the arrival, service, and abandonment rates of this fictitious class by $\lambda = \sum_{i \in T} \lambda_i$, $\mu = \min_{i \in T} \{\mu_i\}$, and $\theta = \min_{i \in T} \{\theta_i\}$, respectively. Since the server can only be busier with this fictitious class, the long-run fraction of time that the server is busy in this relaxed system is a legitimate upper bound for $\sum_{i \in T} x_i$.

Taking this idea further, we could improve this upper bound by formulating a number of two-class MDPs. Divide customer classes in T into two groups and aggregate the classes in

each group into a fictitious class, as before. We then use DP value iteration to compute the maximal fraction of time that the server is busy dealing with these two fictitious classes. When $|T| \geq 4$, we write B_T for the tightest upper bound for $\sum_{i \in T} x_i$ derived with this method, and add it as one constraint. Although it is possible to divide T into three or more groups, the marginal benefit is outweighed by the increased computational burden.

To formulate a linear program to compute an upper bound for the optimal long-run reward rate, define $\mathcal{S}_{k'} \equiv \{T \subseteq \{1, \dots, k\} : |T| = k'\}$, which is the set of all subsets of $\{1, \dots, k\}$ whose cardinality is k' . This linear program is thus given by

$$\begin{aligned} \max \quad & \sum_{i=1}^k R_i \mu_i x_i \\ \text{subject to} \quad & \sum_{i=1}^k x_i \leq 1, \quad x_i \geq 0, \quad 1 \leq i \leq k; \\ & \sum_{i \in T} x_i \leq A_T \quad \text{for all } T \in \mathcal{S}_{k'}, \quad k' = 1, 2, 3; \\ & \sum_{i \in T} x_i \leq B_T \quad \text{for all } T \in \mathcal{S}_{k'}, \quad 4 \leq k' \leq k. \end{aligned}$$

We would expect the upper bound to come close to the optimal long-run reward rate in smaller systems $k \leq 3$, mainly because of the optimized upper bounds A_T . The upper bounds B_T in subsystems of size $k' > 3$ will worsen as k' increases, due to a greater relaxation when creating more customer classes. Consequently, the quality of the upper bound tends to degrade as the size of the system k increases.

6. Numerical Study

In this section, we conduct extensive numerical experiments to assess the impact and design of our API method, as well as the performance of a range of heuristics which includes our suite of priority policies. Section 6.1 uses a numerical study based on cases with two customer classes to explore design choices for our API heuristics. We assess, *inter alia*, the relative performance of the candidate initialising priority rules $R\mu$ and $R\mu\theta$ as well as testing different choices of parameters for our API method. This test yields a recommended API policy which we denote rAPI. Using numerical studies based on cases with three and five customer classes, Section 6.2 compares the performance of rAPI with that of other heuristics. Section 6.3 contains a brief discussion of the computational burden of developing rAPI and the upper bound discussed in Section 5.

6.1. Selecting Parameters for the API Algorithm

To explore the trade-off between different choices of parameters for our API algorithm, we test the algorithm on systems with $k = 2$ customer classes. Problems were randomly generated to reflect a wide range of conditions with regard to (1) the length of job lifetimes in relation to service times (reflected in the categorisation A, B, C in (17c)–(17e) below); and (2) the traffic intensity or workload in the corresponding system without abandonments. There are three categories of traffic—namely light, moderate, and heavy—as determined by the value of $\rho = \sum_{i=1}^k \lambda_i / \mu_i$; see (17f)–(17h) below. For all nine combinations of A, B, C with the traffic categorisation light, moderate, heavy, 500 problems were generated at random. Parameters were sampled as follows:

$$\mu_i \sim U[0.2, 5] \quad (\text{all cases}); \quad (17a)$$

$$\lambda_i \sim U[0.2, 5] \quad (\text{all cases}); \quad (17b)$$

$$\theta_i^{-1} \mu_i | \mu_i \sim U[0.5, 2] \quad (\text{short lifetimes, A}); \quad (17c)$$

$$\theta_i^{-1} \mu_i | \mu_i \sim U[5, 10] \quad (\text{moderate lifetimes, B}); \quad (17d)$$

$$\theta_i^{-1} \mu_i | \mu_i \sim U[20, 200] \quad (\text{long lifetimes, C}); \quad (17e)$$

$$\rho \in [0.6, 0.8] \quad (\text{light traffic}); \quad (17f)$$

$$\rho \in [0.9, 1.1] \quad (\text{moderate traffic}); \quad (17g)$$

$$\rho \in [1.2, 1.4] \quad (\text{heavy traffic}); \quad (17h)$$

In the parameter generation, μ_i and λ_i were sampled according to (17a) and (17b) by means of a rejection algorithm until a desired ρ condition (17f)–(17h) was met. An additional rejection step ensured that the $R\mu\theta$ and $R\mu$ rules of each parameter set were distinct; otherwise, all parameters were resampled. In all cases, rewards were sampled as follows: $R_2 \sim U[1, 3]$ and $R_1 R_2^{-1} | R_2 \sim U[1.25, 2]$. To compute the optimal policy, we use DP value iteration by truncating the state space at $N_i = 40$ for each class i with case A, and $N_i = 60$ with cases B and C, as discussed in Section 2.

Table 1 reports the numerical results with $k = 2$ customer classes. In comparing the $R\mu\theta$ and $R\mu$ rules, please recall that we use the descriptors ‘light’, ‘moderate’, and ‘heavy’ as a shorthand for ranges of the traffic intensity ρ . The actual volume of traffic in the system will also be strongly influenced by the abandonment rate θ , with case C (small θ) yielding higher volumes than case A (large θ). Hence, while the $R\mu\theta$ rule performs very well in the

case $\{C, \text{light}\}$, as is consistent with Theorem 1, it performs poorly in the heaviest traffic case of all, namely $\{C, \text{heavy}\}$. Its performance under A is less variable than under C, as larger abandonment rates act as a moderator on traffic levels, though it still performs best under A when ρ is small. The $R\mu\theta$ rule clearly outperforms the $R\mu$ rule when ρ is small and θ not too large, while $R\mu$ is the better policy when ρ is large, increasingly so as θ declines in value and the traffic levels increase. It is worth noting that at least one of these two priority rules delivers a median performance less than 1% suboptimal across all cases, so they complement each other well.

Table 1 also reports the performance of the policy PI- $R\mu\theta$, which is derived from exact application of a single PI step to the $R\mu\theta$ rule. The fact that the PI- $R\mu\theta$ is nearly optimal shows promise of the API method, if the bias function can be approximated satisfactorily. The policy GADL—due to Glazebrook et al. (2004), with GADL referring to the paper’s four coauthors—requires much computational effort but is typically not as good as the better between the $R\mu\theta$ and $R\mu$ rules.

The next seven columns in Table 1 explore the trade-off between different choices of parameters of the $\text{API}(R\mu\theta, n, m, r, t)$ proposed in Section 4.2. We use $n = 45$ selected states throughout, with $r = \frac{32}{45}$. As one would expect, increasing m and t improves performance whilst increasing computational effort. For a given level of computational effort, the strong performance of PI- $R\mu\theta$ suggests that the policy $\text{API}(\pi, n, m, r, t)$ may perform better with a single, more detailed iteration. Based on our wider experimentation, there appears to be a degree of indifference in performance between multiple, less detailed iterations and a single, more detailed iteration. Further, and unsurprisingly, a strongly performing initial policy π usually improves performance. Based on these observations, to choose parameters in $\text{API}(\pi, n, m, r, t)$ for general k class systems, we recommend $t = 1$, a large value of m (10^5 , say) and allow n , the number of selected states, to scale roughly linearly with k , so that $20k \leq n \leq 25k$. To choose the initial policy π in general, we first run the pairwise-swapping mechanism in Section 3.3 on the $R\mu\theta$ rule, and on the $R\mu$ rule, separately. It turns out that in all numerical tests in this section, the final orderings are the same, which we label PaS. The initial policy π in the API method is thus set to be the best performing among $R\mu\theta$, $R\mu$, and PaS. We shall denote our recommended API policy by rAPI. As seen in Table 1, the rAPI is nearly optimal in all our $k = 2$ class cases. The final column in Table 1 reports the quality of the upper bound presented in Section 5. The upper bound typically sits about 1–2% above the optimal value.

6.2. Comparing the rAPI and Other Heuristics

This section compares the rAPI and other heuristics in systems with $k = 3$ and 5 customer classes. Problem parameters were again generated according to (17a)–(17h), along with suitable rejection algorithms. We now use $R_i \sim U[1, 4]$ for sampled rewards. For each lifetime/traffic combination, 100 problems were generated at random.

Table 2 reports the performance of various service policies against the optimal solution for systems with $k = 3$ customer classes. The rAPI was constructed with $t = 1$, $m = 10^5$, $n = 75$, and $r = \frac{52}{75}$. As seen in the table, the rAPI delivers near-optimal performance in all cases, which reaffirms the strength of policies based on a single, well-estimated (but nonetheless approximate) PI step applied to a well-chosen priority policy. Table 2 also shows that a naive heuristic that always serves the longest queue (labeled SLQ) can perform poorly. In most cases, PaS improves on both $R\mu\theta$ and $R\mu$, though this is not universal. The quality of the upper bound for $k = 3$ customer classes is similar to that for $k = 2$ customer classes.

Table 3 reports the performance of various service policies against an upper bound on the optimal solution, as discussed in Section 5, for systems with $k = 5$ customer classes. Since value iteration is not computationally feasible, the gain of each heuristic is estimated as the mean of 1000 Monte Carlo realisations, which is then compared with the upper bound presented in Section 5. The policy rAPI was constructed with $t = 1$, $m = 10^5$, $n = 100$, and $r = \frac{69}{100}$. As seen in Table 3, the relative quality among $R\mu\theta$, $R\mu$, PaS, and rAPI, is consistent with that in Table 2. The PaS typically improves $R\mu\theta$ and $R\mu$, and then the rAPI further improves the PaS, although the improvement, on average, is rather marginal. The rAPI is the best-performing policy in all cases, and its median performance is within 4% of the upper bound derived in Section 5. Although it is difficult to judge how the rAPI compares with the optimal policy, the fact that the rAPI is much closer to the optimal value than it is to the upper bound in Tables 1 and 2 suggests that the figures in Table 3 are a conservative statement of where the policies stand in relation to the optimal value.

Whereas our numerical experiments in Tables 1–3 show that the suite of priority policies ($R\mu$, $R\mu\theta$, and PaS) generally perform very well and, in several cases the API method offers only marginal improvement on average, it is not always the case. To conclude our numerical study, we offer one example where the improvement of the rAPI method is substantial. Consider a $k = 3$ example in which the class parameters $(\lambda_j, \mu_j, \theta_j, R_j)$ are

Table 2 Percentage suboptimalities in $k = 3$ class systems of various traffic and abandonment level combinations. In the last column, we report the percentage above the optimal policy of the upper bound. Median, 75th, and 90th percentiles are shown.

Case	Workload		$R\mu\theta$	$R\mu$	PaS	SLQ	rAPI	UB
A	Light	90th	0.76	0.46	0.01	5.54	0.00	1.03
		75th	0.23	0.23	0.00	3.29	0.00	0.70
		Median	0.03	0.07	0.00	1.20	0.00	0.34
	Moderate	90th	1.30	0.75	0.01	6.64	0.00	1.25
		75th	0.62	0.39	0.00	3.25	0.00	0.94
		Median	0.11	0.03	0.00	1.50	0.00	0.51
	Heavy	90th	1.32	0.86	0.02	8.52	0.00	1.54
		75th	0.62	0.32	0.00	5.28	0.00	0.99
		Median	0.11	0.02	0.00	2.41	0.00	0.61
B	Light	90th	0.26	0.69	0.04	3.39	0.01	1.15
		75th	0.07	0.24	0.01	2.32	0.00	0.76
		Median	0.01	0.09	0.00	1.43	0.00	0.45
	Moderate	90th	0.85	0.89	0.08	6.03	0.02	1.58
		75th	0.38	0.30	0.01	4.27	0.00	0.92
		Median	0.10	0.05	0.00	2.58	0.00	0.55
	Heavy	90th	1.52	0.86	0.16	10.10	0.03	1.65
		75th	0.84	0.32	0.04	7.15	0.01	1.03
		Median	0.13	0.02	0.00	3.70	0.00	0.58
C	Light	90th	0.01	0.97	0.00	1.44	0.00	0.63
		75th	0.00	0.50	0.00	0.90	0.00	0.29
		Median	0.00	0.21	0.00	0.50	0.00	0.14
	Moderate	90th	0.67	1.52	0.29	4.76	0.10	1.51
		75th	0.24	0.70	0.07	3.28	0.02	0.93
		Median	0.02	0.22	0.00	1.88	0.00	0.56
	Heavy	90th	6.62	0.45	3.02	13.48	0.09	0.80
		75th	2.78	0.13	1.03	8.83	0.01	0.38
		Median	0.76	0.00	0.09	5.23	0.00	0.17

given for classes 1, 2, and 3 by $(\lambda, 3, 0.1, 5)$, $(5\lambda/3, 5, 1, 2)$, and $(4\lambda/3, 4, 5, 1)$, respectively. With these parameters we have $\lambda = \rho$, and the $R\mu\theta$ rule gives class ordering 321, while the $R\mu$ rule gives 123. As seen in Figure 3, the $R\mu\theta$ rule performs well for small ρ , while the $R\mu$ rule performs well for large ρ , which coincides with intuition. For intermediate ρ values, however, there is a substantial gap between the suite of priority policies ($R\mu$, $R\mu\theta$, and PaS) and the rAPI method (with $t = 1$, $m = 10^5$, $n = 75$). In particular, when $\rho = 1.7$, the $R\mu$ rule is 4.26% suboptimal, while the rAPI is 0.04% suboptimal.

6.3. Computational Time for rAPI and the Upper Bound

Table 4 summarizes the time needed to compute the rAPI heuristic. Please note that the algorithm was coded in the C programming language and carried out on a High Performance Computing cluster, with a typical node specification of 2.26Ghz Intel Xeon E5520

Table 3 Percentage below the upper bound in $k = 5$ class systems of various traffic and abandonment level combinations. Median, 75th, and 90th percentiles are shown.

Case	Workload		$R\mu\theta$	$R\mu$	PaS	SLQ	rAPI
A	Light	90th	3.86	4.05	3.76	7.95	3.72
		75th	3.40	3.43	3.24	6.28	3.23
		Median	2.88	2.89	2.73	5.46	2.73
	Moderate	90th	5.47	4.91	4.91	11.63	4.89
		75th	4.08	3.92	3.86	9.86	3.86
		Median	3.32	3.19	3.13	7.38	3.13
	Heavy	90th	6.00	5.86	5.52	13.39	5.52
		75th	5.16	4.97	4.80	10.60	4.80
		Median	4.38	4.05	3.94	8.79	3.94
B	Light	90th	3.76	3.97	3.74	7.85	3.74
		75th	3.25	3.51	3.25	6.52	3.25
		Median	2.74	2.88	2.70	5.36	2.70
	Moderate	90th	5.76	5.96	5.73	13.30	5.73
		75th	4.77	5.15	4.71	11.34	4.67
		Median	3.40	3.41	3.37	9.30	3.37
	Heavy	90th	6.45	6.29	6.08	17.56	6.07
		75th	4.86	4.94	4.75	14.64	4.74
		Median	3.87	3.82	3.78	11.66	3.64
C	Light	90th	1.00	1.55	1.00	2.85	1.00
		75th	0.80	1.23	0.80	2.32	0.80
		Median	0.59	0.77	0.59	1.81	0.59
	Moderate	90th	3.85	4.09	3.85	9.21	3.65
		75th	2.53	3.29	2.53	7.23	2.51
		Median	2.03	2.25	2.03	5.82	2.02
	Heavy	90th	4.73	1.77	4.40	19.11	1.51
		75th	2.92	1.01	2.24	15.09	0.88
		Median	1.20	0.48	1.04	11.03	0.41

processor. Unsurprisingly, the computational burden grows with the number of customer classes k . Recall that the number of selected states n used in the approximate PI step grows roughly linearly in k . Further, as k increases, the balance of computational effort moves toward the sampling stage of the API algorithm and, within the sampling stage, toward the estimation of bias at the support states. These trends particularly reflect the nature of the growth in the mean times $t^\pi(\mathbf{n})$ for a single simulation run during the estimation of the bias $\omega^\pi(\mathbf{n})$. The mean computation times for the upper bound in each problem are in the order of 10 seconds, 400 seconds, and 3000 seconds for systems with 2, 3, and 5 customer classes, respectively. This growth in the computational burden reflects the growth in the number of MDP subproblems that must be solved through DP methods to generate the constraints for the linear program in Section 5, when the number of customer classes increases.

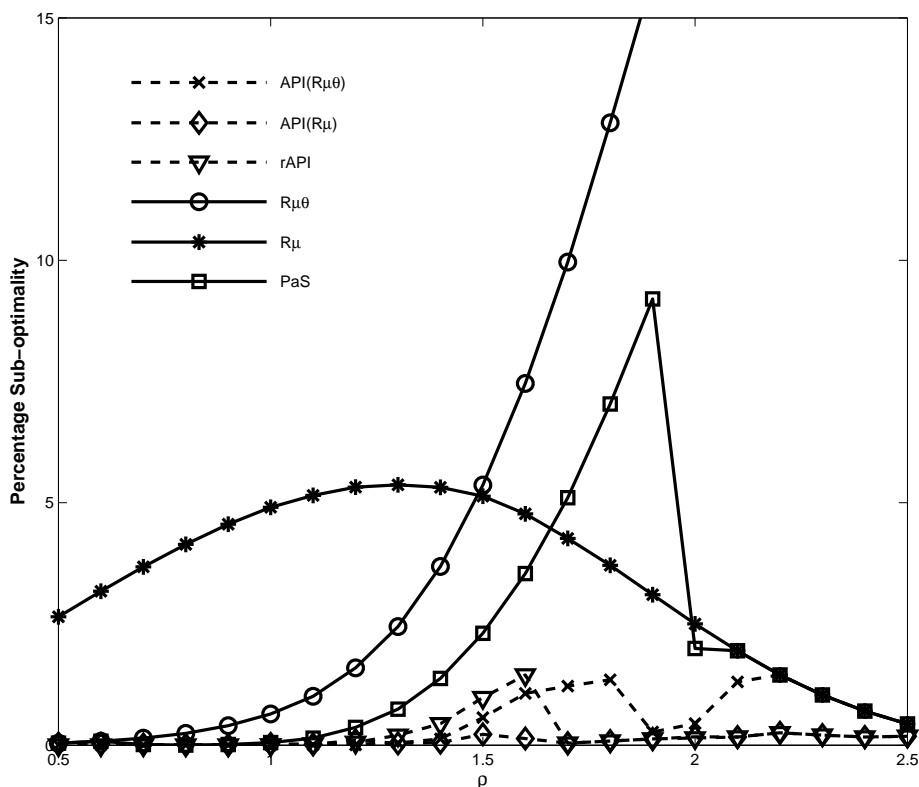


Figure 3 Percentage suboptimality for six heuristics for the example at the end of Section 6.2. The variations of API shown denote the initialising policy.

7. Conclusions

This paper studies the problem of developing effective service policies for multiclass queues with abandonment in a computationally efficient manner. Whereas it is known in the literature that the server can do well simply by maximizing the instantaneous reward rate using the R_{μ} rule in heavy-traffic systems, we show that in light-traffic systems it becomes important to take into account the abandonment rate using the $R_{\mu\theta}$ rule. We also consider an approximate policy improvement algorithm to improve a given service policy. Our numerical study shows that the $R_{\mu\theta}$ rule complements the R_{μ} rule, and applying a pairwise-swapping mechanism to each often yields an even stronger priority policy. The best priority policy that we compute is often nearly optimal; in the cases where it is not, the approximate policy improvement algorithm invariably substantially tightens up either the optimality gap, or the gap relative to the upper bound that we compute.

There are several interesting future research directions. If a customer's lifetime in the system does not follow an exponential distribution, then the server needs to take into account the arrival time of each customer when selecting which customer to serve. The

Table 4 Mean computation time (secs) needed to generate the rAPI policy in each problem of various k class systems. Also shown are the mean proportions of overall computation time spent on the pilot study, sampling of the anchor set, and sampling of the support set.

k	Case	Workload	Time	Proportion		
				Pilot	Anchor	Support
2	A	Light	25	0.07	0.54	0.39
		Moderate	26	0.07	0.54	0.39
		Heavy	29	0.06	0.54	0.39
	B	Light	39	0.05	0.43	0.52
		Moderate	49	0.04	0.44	0.52
		Heavy	73	0.03	0.47	0.50
	C	Light	69	0.03	0.36	0.61
		Moderate	233	0.01	0.45	0.54
		Heavy	291	0.01	0.50	0.48
3	A	Light	76	0.04	0.30	0.67
		Moderate	81	0.03	0.31	0.66
		Heavy	87	0.03	0.31	0.65
	B	Light	127	0.02	0.24	0.74
		Moderate	157	0.02	0.27	0.71
		Heavy	209	0.02	0.31	0.67
	C	Light	249	0.01	0.19	0.80
		Moderate	806	0.01	0.32	0.68
		Heavy	1822	0.00	0.45	0.54
5	A	Light	198	0.02	0.17	0.81
		Moderate	210	0.02	0.18	0.80
		Heavy	224	0.02	0.19	0.79
	B	Light	338	0.01	0.14	0.84
		Moderate	422	0.01	0.18	0.81
		Heavy	559	0.01	0.24	0.75
	C	Light	810	0.01	0.11	0.88
		Moderate	2371	0.00	0.29	0.71
		Heavy	8786	0.00	0.49	0.51

problem will be further complicated if the time needed for the server to switch to another queue cannot be ignored. Still another research direction is to allow the customers to be active decision makers. For example, in military applications, an adversarial customer may not select a queue to join at random as assumed in our model, but instead chooses the one that maximizes the expected gain from his standpoint.

Acknowledgments

The authors wish to thank the associate editor and two anonymous referees for their many helpful comments that substantially improve the presentation of this paper.

References

- Argon, N. T., S. Ziya, R. Righter. 2008. Scheduling impatient jobs in a clearing system with insights on patient triage in mass casualty incidents. *Probability in the Engineering and Informational Sciences* **22** 301–322.
- Ata, B., M. H. Tongarlak. 2013. On scheduling a multiclass queue with abandonments under general delay costs. *Queueing Systems* **74** 65–104.
- Atar, R., C. Giat, N. Shimkin. 2010. The $c\mu/\theta$ -rule for many-server queues with abandonment. *Oper. Res.* **58** 1427–1439.
- Ayesta, U., P. Jacko, V. Novak. 2011. A nearly-optimal index rule for scheduling of users with abandonment. *INFOCOM, 2011 Proceedings IEEE*. 2849–2857.
- Bassamboo, A., J. M. Harrison, A. Zeevi. 2005. Dynamic routing and admission control in high-volume service systems: Asymptotic analysis via multi-scale fluid limits. *Queueing Systems* **51** 249–285.
- Bertsekas, D. P. 2012. *Dynamic Programming and Optimal Control: Approximate Dynamic Programming*, vol. 2. Athena Scientific.
- Down, D. G., G. Koole, M. E. Lewis. 2011. Dynamic control of a single-server system with abandonments. *Queueing Systems* **67** 63–90.
- Garnett, O., A. Mandelbaum, M. Reiman. 2002. Designing a call center with impatient customers. *Manufacturing & Service Operations Management* **4** 208–227.
- Gaver, D. P., P. A. Jacobs, G. Samorodnitsky, K. D. Glazebrook. 2006. Modeling and analysis of uncertain time-critical tasking problems. *Naval Research Logistics (NRL)* **53** 588–599.
- Gittins, John, Richard R. Weber, Kevin D. Glazebrook. 2011. *Multi-armed Bandit Allocation Indices*. John Wiley & Sons.
- Glazebrook, K. D., P. S. Ansell, R. T. Dunn, R. R. Lumley. 2004. On the optimal allocation of service to impatient tasks. *Journal of Applied Probability* **41** 51–72.
- Glazebrook, K. D., J. Niño-Mora. 2001. Parallel scheduling of multiclass $M/M/m$ queues: Approximate and heavy-traffic optimization of achievable performance. *Oper. Res.* **49** 609–623.
- Glazebrook, K. D., E. L. Punton. 2008. Dynamic policies for uncertain time-critical tasking problems. *Naval Research Logistics (NRL)* **55** 142–155.
- Harrison, M. J., A. Zeevi. 2004. Dynamic scheduling of a multiclass queue in the Halfin-Whitt heavy traffic regime. *Oper. Res.* **52** 243–257.
- Howard, R. 1960. *Dynamic Programming and Markov Processes*. MIT Press.

- Jouini, O., A. Pot, G. Koole, Y. Dallery. 2010. Online scheduling policies for multiclass call centers with impatient customers. *Eur. J. Oper. Res.* **207** 258–268.
- Kim, J., A. R. Ward. 2013. Dynamic scheduling of a $GI/GI/1 + GI$ queue with multiple customer classes. *Queueing Systems* **75** 339–384.
- Klimov, G. P. 1974. Time-sharing service systems I. *Theory of Probability & Its Applications* **19** 532–551.
- Klimov, G. P. 1978. Time-sharing service systems II. *Theory of Probability & Its Applications* **23** 314–321.
- Krishnan, K. R. 1987. Joining the right queue: A Markov decision-rule. *Proceedings of the 26th IEEE Conference on Decision and Control*, vol. 26. 1863–1868.
- Larrañaga, M., U. Ayesta, M. Verloop. 2014. Index policies for a multi-class queue with convex holding cost and abandonments. *The 2014 ACM international conference on Measurement and modeling of computer systems*. ACM, 125–137.
- Li, D., K. D. Glazebrook. 2010. An approximate dynamic programming approach to the development of heuristics for the scheduling of impatient jobs in a clearing system. *Naval Research Logistics (NRL)* **57** 225–236.
- Lin, K. Y., M. Kress, R. Szechtman. 2009. Scheduling policies for an antiterrorist surveillance system. *Naval Research Logistics (NRL)* **56** 113–126.
- Niederreiter, H. 1978. Existence of good lattice points in the sense of Hlawka. *Monatshefte für Mathematik* **86** 203–219.
- Powell, M. J. D. 1987. Radial basis functions for multivariable interpolation: a review. *Algorithms for approximation*. Clarendon Press, 143–167.
- Powell, M. J. D. 1999. Recent research at Cambridge on radial basis functions. *New Developments in Approximation Theory*, vol. 132. International Series of Numerical Mathematics, 215–232.
- Powell, W. B. 2011. *Approximate Dynamic Programming: Solving the Curses of Dimensionality*. 2nd ed. John Wiley & Sons.
- Puterman, M. L. 1994. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, Inc.
- Tijms, H. C. 1994. *Stochastic models: an algorithmic approach*. Wiley series in probability and mathematical statistics, John Wiley & Sons.
- Verloop, M. 2014. Asymptotic optimal control of multi-class restless bandits. Tech. rep., Université de Toulouse. CNRS technical report hal-00743781.