



Centrum voor Wiskunde en Informatica

**REPORT***RAPPORT*

Dimensioning Large Call Centers

S.C. Borst, A. Mandelbaum, M.I. Reiman

Probability, Networks and Algorithms (PNA)

**PNA-R0015 November 30, 2000**

Report PNA-R0015  
ISSN 1386-3711

CWI  
P.O. Box 94079  
1090 GB Amsterdam  
The Netherlands

CWI is the National Research Institute for Mathematics and Computer Science. CWI is part of the Stichting Mathematisch Centrum (SMC), the Dutch foundation for promotion of mathematics and computer science and their applications.

SMC is sponsored by the Netherlands Organization for Scientific Research (NWO). CWI is a member of ERCIM, the European Research Consortium for Informatics and Mathematics.

Copyright © Stichting Mathematisch Centrum  
P.O. Box 94079, 1090 GB Amsterdam (NL)  
Kruislaan 413, 1098 SJ Amsterdam (NL)  
Telephone +31 20 592 9333  
Telefax +31 20 592 4199

# Dimensioning Large Call Centers

Sem Borst<sup>†,‡</sup>, Avi Mandelbaum<sup>§,\*</sup>, Marty Reiman<sup>‡</sup>

<sup>†</sup>*CWI*

*P.O. Box 94079, 1090 GB Amsterdam, The Netherlands*

<sup>‡</sup>*Bell Labs, Lucent Technologies*

*Murray Hill, NJ 07974-0636, USA*

<sup>§</sup>*Faculty of Industrial Engineering & Management  
Technion*

*Haifa 32000, Israel*

## ABSTRACT

We develop a framework for asymptotic optimization of a queueing system. The motivation is the staffing problem of call centers with 100's of agents (or more). Such a call center is modeled as an M/M/N queue, where the number of agents  $N$  is large. Within our framework, we determine the asymptotically optimal staffing level  $N^*$  that trades off agents' costs with service quality: the higher the latter, the more expensive is the former. As an alternative to this optimization, we also develop a constraint satisfaction approach where one chooses the least  $N^*$  that adheres to a given constraint on waiting cost. Either way, the analysis gives rise to three regimes of operation: quality-driven, where the focus is on service quality; efficiency-driven, which emphasizes agents' costs; and a rationalized regime that balances, and in fact unifies, the other two. Numerical experiments reveal remarkable accuracy of our asymptotic approximations: over a wide range of parameters, from the very small to the extremely large,  $N^*$  is *exactly* optimal, or it is accurate to within a single agent. We demonstrate the utility of our approach by revisiting the square-root safety staffing principle, which is a long-existing rule-of-thumb for staffing the M/M/N queue. In its simplest form, our rule is as follows: if  $c$  is the hourly cost of an agent, and  $a$  is the hourly cost of customers' delay, then  $N^* = R + y^*(\frac{a}{c})\sqrt{R}$ , where  $R$  is the offered load, and  $y^*(\cdot)$  is a function that is easily computable.

*2000 Mathematics Subject Classification:* 60K25 (primary), 90B22 (secondary).

*Keywords & Phrases:* asymptotics, call centers, multi-server queues, optimal staffing.

*Note:* Work of the first author carried out in part under the project PNA2.1 "Communication and Computer Networks".

---

\*The research of A.M. was partially supported by the ISF (Israeli Science Foundation) grant 388/99-2 and by the Technion funds for the promotion of research and sponsored research.

# 1 Introduction

Worldwide, telephone-based services have been expanding dramatically in both volume and scope. This has given rise to a huge growth industry – the (telephone) call center industry. Indeed, some assess [4] that 70% of all customer-business interactions in the U.S. occur in call centers, which employ (conservatively speaking) 3% of the U.S. workforce (about 1.5 million agents). Marketing managers refer to call centers as the modern business frontier, being the focus of Customer Relationship Management (CRM); Operations managers are challenged with the fact that personnel costs, specifically staffing, account for over 65% of the cost of running the typical call center. The trade-off between service quality (marketing) and efficiency (operations) thus naturally arises, and a central goal of ours is to contribute to its understanding.

We argue that call centers typify an emerging business environment in which the traditional quality-efficiency trade-off paradigm could collapse: extremely high levels of both service quality and efficiency can coexist. Consider, for example, a best-practice U.S. sales call center that attends to an average of 15,000 phone callers daily; the average duration of a call is 4 minutes and the variability of calls is significant; agents are highly utilized (over 90%), yet customers essentially never encounter a busy signal, hardly anyone abandons while waiting, and the average wait for service is a mere few seconds. Prerequisites for sustaining such performance, to the best of our judgment, are technology-enabled economies-of-scale and scientifically-based managerial principles and laws. In this paper we develop an analytical framework (Sections 3 and 8) that supports such principles. It is based on asymptotic optimization, which yields insight that does not come out of exact analysis. A convincing example is the *square-root safety staffing principle*, described in Subsection 1.3 below. It supports simple useful rules-of-thumb for staffing large call centers, rules that so far have been justified only heuristically. Indeed, formal asymptotic justifications of such rules are not common in the Operations Research literature. Hence another goal here is to convince the reader of their benefits.

## 1.1 Costs, optimization and constraint satisfaction

The cost of staffing is the principal component in the operating expenses of a call center. The staffing level is also the dominant factor to determine service level, as measured in terms of delay statistics: poor service levels incur either opportunity losses due to deteriorating goodwill, or more direct revenue losses in case of abandonment and blocking (busy signals). While the need to balance service quality and staffing cost is universal, the weight placed on each may vary dramatically. In some call centers, providing maximal customer care is the primary drive whereas in others, handling a high traffic volume at minimal cost is the overriding goal. The challenge, so we argue, is to translate such strategically-articulated goals into concrete staffing

levels: simply put, *how many agents are to be staffed in order to provide acceptable service quality and operational efficiency?* In this paper we answer this question for the M/M/N queue, which is the simplest yet most prevalent model that supports call center staffing. (In future research we are hoping to add central features of call centers such as abandonment and retrials [13].)

Within the M/M/N model, we postulate a staffing cost function  $F(N)$  for employing  $N$  agents. Low costs (small  $N$ ) give rise to long waits, which we quantify in terms of a delay cost function  $D(t)$  for a customer being served after waiting  $t$  units of time. When  $F$  dominates  $D$  (or conversely  $D$  dominates  $F$ ), the least costs are achieved in an *efficiency-driven* (or conversely a *quality-driven*) operation. When  $F$  and  $D$  are comparable, optimization leads to a *rationalized* operation which, as it turns out, is robust enough to encompass most circumstances. Formally, the three regimes emerge from an asymptotic analysis of the M/M/N queue, as the arrival rate  $\lambda$ , and accordingly the optimal staffing level  $N_\lambda^*$ , both scale up to infinity. We refer to such responsive staffing, in response to increased load, as *dimensioning* the call center, which inspired our title. While the staffing levels that we recommend are only asymptotically optimal, they are nevertheless remarkably accurate – to within a *single* agent in the majority of cases. The asymptotics also provide insight, beyond that of exact analysis, about the dependence of the optimal  $N_\lambda^*$  on  $\lambda$ ,  $F$  and  $D$ .

In industry practice, staffing levels are rarely determined through optimization. (One reason is that there is no standard practice for quantifying waiting costs, let alone abandonment, busy-signal and retrial costs; see [1] for some attempts.) Thus, if not by mere experience-based guessing, common practice seeks the least number of agents  $N^*$  that satisfies a given constraint on service level. The latter is expressed in terms of some congestion measure, for example the industry-standard Total Service Factor (TSF) given by

$$TSF = \Pr\{\text{Wait} > T\}, \quad \text{for some } T \geq 0,$$

perhaps combined with 1-800 operating costs. We call this practice *constraint satisfaction* (Section 8). It is to be contrasted with our previous *optimization* practice, where  $N^*$  was determined by cost minimization.

## 1.2 Structure of the paper

The rest of the Introduction is devoted to an exposition of the square-root safety staffing principle, followed by a review of some related literature.

In Section 2 we set up our M/M/N model and its cost structure. The framework for asymptotic optimality is developed in Section 3. Its applications require some special functions which are introduced in Section 4. Of central importance to us is the Halfin-Whitt delay function  $P(\cdot)$  [8], plotted in Figure 3. It provides an approximation for the delay probability in the

M/M/N queue,  $N$  large, which operates in the rationalized regime. (Some useful properties of  $P(\cdot)$  and other functions are verified in Appendices A-C.) In Sections 5-7 we analyze, respectively, the rationalized, efficiency- and quality-driven regimes, under the optimization approach. While the analysis is abstract, each of these sections concludes with examples of specific cost structures, for concreteness. In Section 8 we introduce the constraint satisfaction approach, which gives rise to the same three regimes of operation as optimization. Section 9 describes numerical experiments that test the accuracy of our asymptotically-supported approximations. The findings are astounding – rarely do we miss by more than a single agent, as far as optimal staffing levels are concerned. In addition, even though the theory is asymptotic, our approximations are accurate with as few as 3 agents. In order to apply our approximations, guidelines are required for fitting a given call center, represented by its parameters and costs, to one of the three operational regimes. This turns out simpler than expected. Indeed, our numerical experiments, backed up by some theory, clearly establish the robustness of the rationalized approximation, as it covers accurately both the efficiency- and quality-driven regimes. Thus, except for extreme settings, the rationalized approximation is the one to use, as we now do in the following example. We conclude in Section 10 with a few worthy directions for future research.

### 1.3 The square-root safety staffing principle

To recapitulate, we determine asymptotically optimal staffing levels in accordance with the relative importance of agents' costs and efficiency versus customers' service quality. This leads to the (re)discovery, as well as a deeper understanding, of a remarkably robust rule-of-thumb, the *square-root safety staffing* rule. It reads as follows: Suppose that the arrival rate is  $\lambda$  customers per hour, and service rate is  $\mu$ , which implies that the system's *offered load* is given by  $R = \frac{\lambda}{\mu}$ ; if the staffing cost is \$  $c$  per agent per hour, and waiting cost is \$  $a$  per customer per hour, our recommended number of servers  $N^*$  is given by

$$N^* = R + y^*\left(\frac{a}{c}\right)\sqrt{R}, \quad (1)$$

where the function  $y^*(r)$ ,  $r \geq 0$ , is plotted in Figures 1 and 2. In simple words, at least  $R$  agents ( $\lfloor R \rfloor + 1$  to be exact) are required to guarantee stability; however, safety staffing must be added to the minimum as a protection against stochastic variability. This number of additional agents is proportional to  $\sqrt{R}$ , and the proportionality coefficient  $y^*\left(\frac{a}{c}\right)$  is determined through the optimization (22), by the relative importance of customers' delay ( $a$ ) to agents' salary ( $c$ ).

Note that the right-hand side of (1) need not be an integer, in which case  $N^*$  is obtained by rounding it off. We demonstrate in Section 9, below (39), that this yields the staffing level that minimizes waiting plus staffing costs, exactly in most cases and off by a single agent in the other rare ones.

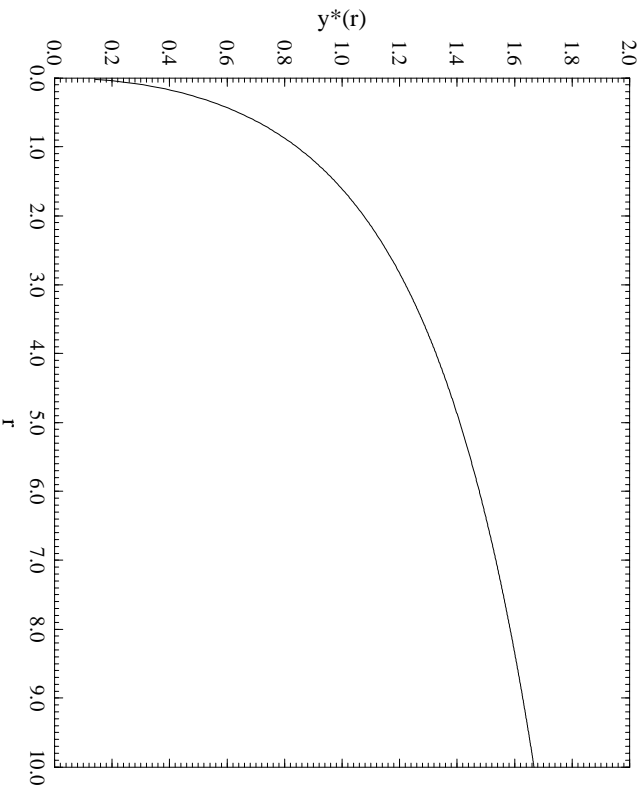


Figure 1:  $y^*(r)$  as function of  $r$ ,  $0 \leq r \leq 10$ .

The *form* of (1) already carries with it important insight. Let  $\Delta = y^*(\frac{a}{c})\sqrt{R}$  denote the *safety staffing* level (the excess number of servers above the minimum  $R = \frac{\lambda}{\mu}$ ). Then, with  $a$  and  $c$  fixed, an  $n$ -fold increase in the offered load  $R$  requires that the safety staffing  $\Delta$  increases by only  $\sqrt{n}$ -fold, which constitutes significant economies of scale [15].

Now suppose that  $R$  is measured in 100's, as it is in large call centers. Then  $\sqrt{R}$  is in the low 10's, hence  $\Delta$  is as well (since  $y^*$  grows so slowly:  $y^*(100) \approx 2.5$ ). It follows that the bulk of the agents, namely  $R$ , must be present for stability, and only a small fraction  $\Delta/R = y^*/\sqrt{R}$  of these must be added as safety against stochastic variability (typically up to 10%, and in fact significantly less for large call centers). This results in high agent utilization levels  $R/N^*$ , around 90% and up. Nevertheless, as shown in Examples 1.1-1.3 below, operational service quality ranges from the acceptable to the extremely high. (Indeed, small changes in  $\Delta$ , which amounts to small changes in agents' utilization, have noticeable effects on performance.) We are thus operating in a regime where high resource utilization and service level coexist, which is due to economies of scale that dominate stochastic variability. It is important and interesting to note that data from large call centers confirms these observations [5].

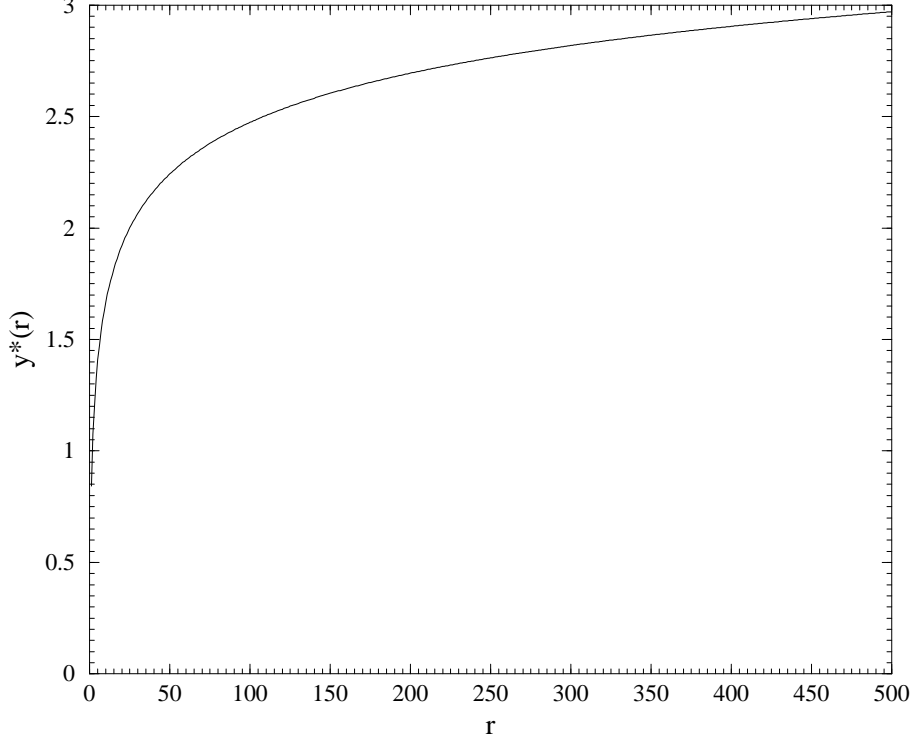


Figure 2:  $y^*(r)$  as function of  $r$ ,  $0 \leq r \leq 500$ .

Small values of  $r$  correspond to efficiency-driven staffing. In this range, the function  $y^*(\cdot)$  is reasonably approximated by

$$y^*(r) \approx \sqrt{\frac{r}{1 + r(\sqrt{\frac{\pi}{2}} - 1)}}, \quad 0 \leq r < 10.$$

Large values of  $r$  correspond to quality-driven staffing. In this range, a close lower bound is  $y^*(r) \approx \sqrt{s - \ln s}$ , where  $s = 2 \ln \frac{r}{\sqrt{2\pi}}$ ,  $r \uparrow \infty$ . (See Remark 5.4 for some details on these asymptotic expansions.)

Under our square-root safety staffing, it is anticipated that service level, as expressed by the industry-standard TSF, equals

$$TSF = \Pr\{\text{Wait} > T\} \approx P(y^*) \times e^{-Ty^* \sqrt{\lambda\mu}}; \quad y^* = y^*\left(\frac{a}{c}\right);$$

in which

$$P(y^*) = \left[1 + \frac{y^* \Phi(y^*)}{\phi(y^*)}\right]^{-1} \approx \Pr\{\text{Wait} > 0\}, \quad (2)$$

is the *Halfin-Whitt delay function* [8] (see Figure 3 and Section 4);  $\phi(\cdot)$  and  $\Phi(\cdot)$  are the density and cumulative distribution function of the standard normal distribution, respectively. A more management-friendly representation of TSF is

$$TSF = \Pr\{\text{Wait} > T \times \mathbb{E}[\text{Service Time}]\} \approx \Pr\{\text{Wait} > 0\} \times e^{-T\Delta}. \quad (3)$$



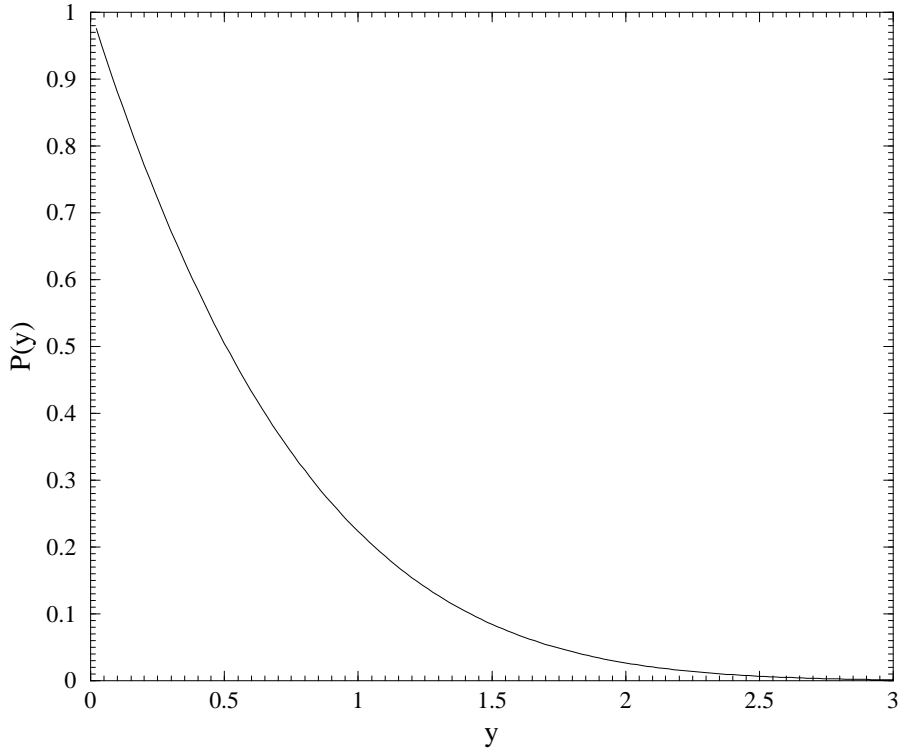


Figure 3: The Halfin-Whitt delay function  $P(y)$ .

Here delay is measured in units of average service time ( $E[\text{Service Time}] = \frac{1}{\mu}$ ), and  $\Delta = y^* \sqrt{R}$  is the safety staffing level. Another service level standard is the average waiting time, often referred to as Average Speed of Answer (ASA). With  $N^*$  as in (1), and again naturally quantified in units of service durations, it is given by

$$\frac{\text{ASA}}{1/\mu} = \frac{E[\text{Wait}]}{E[\text{Service Time}]} \approx \frac{P(y^*)}{\Delta}. \quad (4)$$

The industry standard for measuring operational efficiency is agent utilization, namely  $\frac{R}{N}$ , which is traded off against service level. Agents are thus idle, or more appropriately described as being available for service, a fraction  $\frac{\Delta}{N}$  of their time.

### Example 1.1

Consider, for example, the best-practice call center, described in the second paragraph of our Introduction. Assuming 1800 calls per busy hour, the offered load equals  $R = 120$ . With 90% utilization, one expects that about  $N^* \approx 133$  agents share the load ( $\Delta = 13$ ), hence the center operates with  $y^* \approx 1.22$ . Inverting  $y^*(\cdot)$  in Figure 1 shows that, in this call center, an hour wait of customers is valued as 3 times the hourly wage of an agent. With this staffing level, it is expected that about 15% of the customers ( $P(1.22) = 0.15$ ) are delayed; that 5% of the

customers are delayed over 20 seconds (using (3) with  $T = \frac{1}{12}$ ); and that, by (4), ASA equals 2.7 seconds (while those who were delayed actually averaged 18 seconds waiting).

□

But the staffing level in the example can be interpreted differently. To this end, recall that the prevalent alternative to the above optimization approach is constraint satisfaction. Specifically, in Example 8.5 it is shown that the least  $N$  that guarantees  $\Pr\{\text{Wait} > 0\} < \epsilon$  is closely approximated by rounding up

$$N^* = R + P^{-1}(\epsilon)\sqrt{R}, \quad (5)$$

where  $P(\cdot)$  is the Halfin-Whitt delay function introduced in (2). Returning to the above best-practice call center,  $P^{-1}(\epsilon) = 1.22$  yields, as expected,  $\epsilon = 0.15$ .

### Example 1.2

One should note that a constraint on the fraction of delayed customers is severe, hence it fits call centers that cater to say emergency calls. This can be nicely explained within our framework. For example, requiring that  $\epsilon = 0.01$ , namely 1 customer out of 100 delayed on average, corresponds to  $y^* = P^{-1}(0.01) = 2.38$  (see Figure 3), which *could* be interpreted as saying that  $\frac{a}{c} = (y^*)^{-1}(2.38) = 75$ ! An evaluation of customers' time as being worth 75-fold of agents' time seems reasonable only under extreme circumstances: for example, if the 'servers' are 'cheap' being say Interactive Voice Reponse (IVR) units, or customers' time is highly valued as with emergency call centers.

□

### Example 1.3

Most call centers define TSF with a positive  $T$ , and then requiring  $\epsilon = 0.01$  need not be extreme. We now illustrate this in terms of the following scenario. A prevalent standard is to aspire that no more than 80% of the callers are delayed over  $T = 20$  seconds. Incidentally, we believe that the source of this standard is the familiar 20:80 managerial rule-of-thumb, stating in great generality and vagueness that "only 20% of the reasons already give rise to 80% of the problems". While there is no apparent reason for connecting this rule-of-thumb with any staffing standard, it is nevertheless worthwhile to note that our framework provides some interesting implications for using this rule.

Consider a large call center with  $\lambda = 100$  calls per minute, and 4 minutes average call duration. Thus  $R = 400$ , and adhering to the 20:80 rule implies that  $y^* = 0.53$ , hence  $N^* = 411$ . By Figure 1, this translates into  $\frac{a}{c} = 0.32$ . It follows that, while customers are not highly valued, the 20:80 rule is 'easy' to adhere to because of the call center's size. To wit, increasing  $N^*$  to 429 amounts to  $y^* = 1.4$ , or  $\frac{a}{c} = 4.9$ , reflecting a significant yet reasonable increase of the relative value of customers' to agents' time. This is accompanied by an increase in server availability

(idleness), from 3% to 7%, which enables an order-of-magnitude reduction in TSF, from 0.2 to little less than 0.01: about 1 out of 100 customers is delayed for more than 20 seconds.

To underscore the role of scale in the above scenario, consider a call center with the same offered load parameters as Example 1.1: 30 calls per minute, and again 4 minutes average call duration. Now  $R = 120$ , but it takes  $N^* = 140$  to achieve  $\text{TSF} = 0.01$ , with  $T = 20$  seconds. This corresponds to  $y^* = 1.75$ , or  $\frac{a}{c} = 12.5$ , a 2.5-fold increase over the large call center. It is interesting to note that with an average call duration of 30 seconds (as in 411 services), with  $T$  held at 20 seconds,  $N^* = 126$  would suffice, which amounts to  $y^* = 0.53$  and  $\frac{a}{c} = 0.32$ . This is identical to the large call center with the 20:80 rule operation, but the latter accommodates mean service time of 4 minutes, in contrast to the 30 seconds here.

□

The square-root safety staffing principle emerged from the simplest cost structure (linear staffing and waiting costs). While our framework accommodates general costs, the corresponding safety staffing levels are nevertheless *always* proportional to  $\sqrt{R}$ ; it is only the proportionality coefficient that varies with the cost.

## 1.4 Related Literature

The square-root safety staffing principle has been part of the queueing-theory folklore for a long time. This is well documented by Grassmann [6, 7], and recently revisited by Kolesar & Green [12], where both its accuracy and applicability have been convincingly confirmed. The principle was substantiated by Whitt [15], then adapted in Jennings *et.al.* [11] to non-stationary models. All of this works applies infinite-server heuristics, grounded in the fact that the steady-state number of customers in the M/M/ $\infty$  queue, say  $Q^\infty$ , is Poisson distributed with mean  $R = \frac{\lambda}{\mu}$ . It follows that  $Q^\infty$  is approximately normally distributed, with mean  $R$  and standard deviation  $\sqrt{R}$ , when  $R$  is not too small. To relate this to staffing in the M/M/N model, one approximates the latter's probability of delay by

$$\Pr\{Q^\infty \geq N\} \approx 1 - \Phi\left(\frac{N - R}{\sqrt{R}}\right).$$

Then, the staffing level  $N^*$  that guarantees  $\epsilon$  delay probability is chosen to be

$$N^* = R + \bar{\Phi}^{-1}(\epsilon)\sqrt{R}, \tag{6}$$

where  $\bar{\Phi} = 1 - \Phi$ .

The square-root principle has two parts to it: first, the conceptual observation that the safety staffing level is proportional to the square-root of the offered load; and second, the explicit calculation of the proportionality coefficient  $y^*$ . Our framework accommodates *both* of these two needs, while in all previous works, to the best of our understanding, at least one of them

is treated in a heuristic fashion or simply ignored. (We shall be specific momentarily.) More important, however, is the fact that our approach and framework allow an arbitrary cost structure, and they have the potential to generalize beyond Erlang-C. For a concrete example, Garnet *et.al.* [5] accommodate impatient customers: in their main result, the square-root rule arises conceptually, but the determination of the value of  $y^*$  is left open.

Being specific now, [15] and [11] refer to  $y^*$  as a measure of service level, but leave out any explicit calculation of it. Grassmann [7], taking the optimization approach, leads the reader through an instructive progression of increasingly complex staffing models, culminating in his ‘equilibrium model’ (Erlang-C), for which no ‘square-root’ justification is provided. (It is justified for his less complex model, under the ‘Independence Assumption’, but this amounts to using (6).) Some numerical experiments, inspired by [7], are reported at the beginning of Section 9. Finally Kolesar & Green [12] advocate the use of (6), in order to support constraint satisfaction that achieves  $\Pr\{\text{Wait} > 0\} \leq \epsilon$ . We, on the other hand, recommend the use of (5) for constraint satisfaction, which is proven asymptotically accurate in Example 8.5. The approximations (5) and (6) essentially coincide for small  $\epsilon$ ’s, but (5) is uniformly more accurate. We refer to the beginning of Section 9 for more details.

## 2 Model description

We consider the classical M/M/N (Erlang-C) model with  $N$  servers and infinite-capacity waiting room. Customers arrive as a Poisson process of rate  $\lambda$ , and have independent exponentially distributed service times with mean  $1/\mu$ . The service rate  $\mu$  will be arbitrary but fixed, whereas the arrival rate  $\lambda$  will grow large in order to obtain asymptotic scaling results. We assume  $\lambda/N\mu < 1$  for stability. Customers are served in order of arrival; then (see for instance [3]) the waiting-time distribution is given by

$$\Pr\{\text{Wait} > t\} = \pi(N, \lambda/\mu) e^{-(N\mu - \lambda)t},$$

where the probability of waiting  $\pi(N, \lambda/\mu) = \Pr\{\text{Wait} > 0\}$  is determined by

$$\pi(N, \nu) = \frac{\nu^N}{N!} \left\{ (1 - \nu/N) \sum_{n=0}^{N-1} \frac{\nu^n}{n!} + \frac{\nu^N}{N!} \right\}^{-1}.$$

We consider the problem of determining the staffing level  $N$  that optimally balances staffing cost against quality-of-service. To this end, a staffing cost  $F(N)$  per unit of time is associated with staffing  $N$  servers. We assume that  $F(N)$  is also defined for all non-integer values  $N > \lambda/\mu$ , and that this extended function  $F(\cdot)$  is convex and strictly increasing.

Quality-of-service is quantified in terms of a waiting-cost function  $D_\lambda(\cdot)$ : a cost  $D_\lambda(t)$  is incurred when a customer waits for  $t$  time units. (The subscript  $\lambda$  is attached to allow for the possibility that the primitives vary with the arrival intensity.) We assume that  $D_\lambda(\cdot)$  is

strictly increasing. Without loss of generality, we may take  $D_\lambda(0) = 0$ . The expected total cost per unit of time is then given by

$$C(N, \lambda) = F(N) + \lambda \mathbb{E}D_\lambda(\text{Wait}) = F(N) + \lambda \pi(N, \lambda/\mu)G(N, \lambda),$$

where

$$G(N, \lambda) = \mathbb{E}D_\lambda(\text{Wait}|\text{Wait} > 0) = (N\mu - \lambda) \int_0^\infty D_\lambda(t) e^{-(N\mu - \lambda)t} dt.$$

Notice that  $G(N, \lambda)$  is also defined for all non-integer values  $N > \lambda/\mu$ . We assume that  $D_\lambda(\cdot)$  is such that  $G(N, \lambda)$  is finite for all  $\lambda/\mu < N$ .

We are interested in determining the optimum staffing level

$$N_\lambda^* := \arg \min_{N > \lambda/\mu} C(N, \lambda) \tag{7}$$

(the minimization being over integer values). To see that  $N_\lambda^*$  is well-defined, notice that  $\lim_{N \rightarrow \infty} F(N) = \infty$ , and thus  $\lim_{N \rightarrow \infty} C(N, \lambda) = \infty$ . Hence,  $C(N, \lambda)$  indeed achieves a minimum value.

### 3 Framework for asymptotic optimality

In principle, the optimum staffing level  $N_\lambda^*$  in equation (7) may be obtained through brute-force enumeration. Rather than determining the optimum staffing level numerically, however, we are primarily interested in gaining insight into how  $N_\lambda^*$  grows with the arrival intensity  $\lambda$ , and how it depends on the staffing and waiting cost functions  $F(N)$  and  $G(N, \lambda)$ . In order to do so, we develop an approximate analytical approach for determining the optimum staffing level. As a first step, we translate the discrete optimization problem (7) into a continuous one. The next step is to approximate the latter problem by a related continuous version which is easier to solve. To validate the approach, we then prove that the optimal solution to the approximating continuous problem provides an asymptotically optimal solution to the original discrete problem.

We first transform the discrete optimization problem into a continuous one. Let

$$N_\lambda(x) = \lambda/\mu + x\sqrt{\lambda/\mu},$$

so that the variable  $x = (N - \lambda/\mu)/\sqrt{\lambda/\mu}$  is the (normalized) number of servers in excess of the minimum number  $\lambda/\mu$  required for stability. In terms of  $x$ , we define

$$F_\lambda(x) := F(N_\lambda(x)) - F(\lambda/\mu);$$

$$G_\lambda(x) := \lambda G(N_\lambda(x), \lambda);$$

$$C_\lambda(x) := C(N_\lambda(x), \lambda) - F(\lambda/\mu);$$

$$\pi_\lambda(x) := H(N_\lambda(x), \lambda/\mu),$$

with

$$H(M, \alpha) = \left\{ \alpha \int_0^\infty e^{-\alpha t} t(1+t)^{M-1} dt \right\}^{-1}.$$

It can be verified ([9], [10]) that  $\pi_\lambda(x) = \pi(N_\lambda(x), \lambda/\mu)$  for integer values of  $N_\lambda(x)$ . The total cost per unit of time (up to the additive constant factor  $F(\lambda/\mu)$ ) can thus be rewritten

$$C_\lambda(x) = F_\lambda(x) + \pi_\lambda(x)G_\lambda(x).$$

Denote

$$x_\lambda^* := \arg \min_{x>0} C_\lambda(x). \quad (8)$$

To see that  $x_\lambda^*$  is well-defined, first notice that the function  $C_\lambda(\cdot)$  is strictly convex. This follows from the assumption that  $F(\cdot)$  is convex and the fact that  $\pi_\lambda(\cdot)$  is convex ([9], [10]) and  $G_\lambda(\cdot)$  is strictly convex (Appendix C). In addition,  $\lim_{x \downarrow 0} C_\lambda(x) = \infty$ , since  $\lim_{N \downarrow \lambda/\mu} G(N, \lambda) = \infty$ . Also,  $\lim_{x \rightarrow \infty} C_\lambda(x) = \infty$ , because  $\lim_{N \rightarrow \infty} F(N) = \infty$ . Hence,  $C_\lambda(\cdot)$  is unimodal, implying that it indeed achieves a unique minimum value at  $x_\lambda^* \in (0, \infty)$ . Further notice that either  $N_\lambda^* = \lfloor N_\lambda(x_\lambda^*) \rfloor$  or  $N_\lambda^* = \lceil N_\lambda(x_\lambda^*) \rceil$ , which establishes the link between the discrete problem and the corresponding continuous problem. (Here  $\lfloor u \rfloor$  and  $\lceil u \rceil$  denote the largest integer smaller than or equal to  $u$ , and the smallest integer larger than or equal to  $u$ , respectively.)

Next, we approximate  $x_\lambda^*$  in (8) by

$$z_\lambda^* := \arg \min_{z>0} C[z; \hat{F}_\lambda, \hat{\pi}_\lambda, \hat{G}_\lambda], \quad (9)$$

where

$$C[z; \hat{F}_\lambda, \hat{\pi}_\lambda, \hat{G}_\lambda] := \hat{F}_\lambda(z) + \hat{\pi}_\lambda(z)\hat{G}_\lambda(z),$$

with the functions  $\hat{F}_\lambda(\cdot)$ ,  $\hat{\pi}_\lambda(\cdot)$ ,  $\hat{G}_\lambda(\cdot)$  ‘approximating’  $F_\lambda(\cdot)$ ,  $\pi_\lambda(\cdot)$ ,  $G_\lambda(\cdot)$ , respectively. (Note that with this notation,  $x_\lambda^* = \arg \min_{x>0} C[x; F_\lambda, \pi_\lambda, G_\lambda]$ .) The approximating functions  $\hat{F}_\lambda(\cdot)$ ,  $\hat{G}_\lambda(\cdot)$ , and  $\hat{\pi}_\lambda(\cdot)$  that we consider will always be such that  $z_\lambda^*$  exists and is unique. If  $\hat{F}_\lambda(\cdot)$ ,  $\hat{G}_\lambda(\cdot)$ ,  $\hat{\pi}_\lambda(\cdot)$  have a simple form, then solving for  $z_\lambda^*$  will be easier than determining  $x_\lambda^*$ . At the same time, if  $\hat{F}_\lambda(\cdot)$ ,  $\hat{G}_\lambda(\cdot)$ , and  $\hat{\pi}_\lambda(\cdot)$  approximate  $F_\lambda(\cdot)$ ,  $G_\lambda(\cdot)$ , and  $\pi_\lambda(\cdot)$  well, then it is reasonable to expect that  $z_\lambda^*$  provides a good approximation to  $x_\lambda^*$  and, moreover,  $N_\lambda(z_\lambda^*)$  yields a good approximation to  $N_\lambda^*$ .

Before formalizing the above approximation principle, we first introduce the following notational conventions: for any pair of functions  $a_\lambda$  and  $b_\lambda$  (implicitly assuming existence of the limits), denote

$$\begin{aligned}
a_\lambda &\stackrel{\infty}{\sim} b_\lambda : \lim_{\lambda \rightarrow \infty} \frac{a_\lambda}{b_\lambda} = 1; & a_\lambda &\stackrel{\infty}{\approx} b_\lambda : \lim_{\lambda \rightarrow \infty} \frac{a_\lambda}{b_\lambda} = \gamma, \quad 0 < \gamma < \infty; \\
a_\lambda &\stackrel{\infty}{\ll} b_\lambda : \lim_{\lambda \rightarrow \infty} \frac{a_\lambda}{b_\lambda} = 0; & a_\lambda &\stackrel{\infty}{\gg} b_\lambda : \lim_{\lambda \rightarrow \infty} \frac{a_\lambda}{b_\lambda} = \infty; \\
a_\lambda &\stackrel{\sup}{\leq} b_\lambda : \limsup_{\lambda \rightarrow \infty} \frac{a_\lambda}{b_\lambda} \leq 1; & a_\lambda &\stackrel{\inf}{\geq} b_\lambda : \liminf_{\lambda \rightarrow \infty} \frac{a_\lambda}{b_\lambda} \geq 1; \\
a_\lambda &\stackrel{\inf}{<} b_\lambda : \liminf_{\lambda \rightarrow \infty} \frac{a_\lambda}{b_\lambda} < 1; & a_\lambda &\stackrel{\sup}{>} b_\lambda : \limsup_{\lambda \rightarrow \infty} \frac{a_\lambda}{b_\lambda} > 1; \\
a_\lambda &\stackrel{\inf}{\ll} b_\lambda : \liminf_{\lambda \rightarrow \infty} \frac{a_\lambda}{b_\lambda} = 0; & a_\lambda &\stackrel{\sup}{\gg} b_\lambda : \limsup_{\lambda \rightarrow \infty} \frac{a_\lambda}{b_\lambda} = \infty.
\end{aligned}$$

**Lemma 3.1**

Denote  $\hat{C}_\lambda(z) = C[z; \hat{F}_\lambda, \hat{\pi}_\lambda, \hat{G}_\lambda]$ .

Then  $C_\lambda(z_\lambda^*) \stackrel{\infty}{\sim} C_\lambda(x_\lambda^*)$  if both  $C_\lambda(x_\lambda^*) \stackrel{\infty}{\sim} \hat{C}_\lambda(x_\lambda^*)$  and  $C_\lambda(z_\lambda^*) \stackrel{\infty}{\sim} \hat{C}_\lambda(z_\lambda^*)$ .

**Proof**

By definition of  $x_\lambda^*$ ,  $C_\lambda(z_\lambda^*) \geq C_\lambda(x_\lambda^*)$ , so it suffices to show that  $C_\lambda(z_\lambda^*) \stackrel{\sup}{\leq} C_\lambda(x_\lambda^*)$ , which follows directly from

$$C_\lambda(z_\lambda^*) \stackrel{\infty}{\sim} \hat{C}_\lambda(z_\lambda^*) \leq \hat{C}_\lambda(x_\lambda^*) \stackrel{\infty}{\sim} C_\lambda(x_\lambda^*).$$

□

Define

$$S_\lambda(x) := \min\{C(\lfloor N_\lambda(x) \rfloor, \lambda), C(\lceil N_\lambda(x) \rceil, \lambda)\}. \quad (10)$$

**Lemma 3.2**

If  $C_\lambda(z_\lambda^*) \stackrel{\infty}{\sim} C_\lambda(x_\lambda^*)$ , then  $S_\lambda(z_\lambda^*) - F(\lambda/\mu) \stackrel{\infty}{\sim} C(N_\lambda^*, \lambda) - F(\lambda/\mu)$ .

**Proof**

By definition,  $S_\lambda(z_\lambda^*) \geq C(N_\lambda^*, \lambda)$ , so it suffices to show that  $S_\lambda(z_\lambda^*) - F(\lambda/\mu) \stackrel{\sup}{\leq} C(N_\lambda^*, \lambda) - F(\lambda/\mu)$ .

For fixed  $\lambda$ , we distinguish between four cases.

- i.  $N_\lambda^* - 1 < N_\lambda(z_\lambda^*) \leq N_\lambda^*$ . Then  $\lceil N_\lambda(z_\lambda^*) \rceil = N_\lambda^*$ , and  $S_\lambda(z_\lambda^*) = C(N_\lambda^*, \lambda)$ .
- ii.  $N_\lambda^* \leq N_\lambda(z_\lambda^*) < N_\lambda^* + 1$ . Then  $\lfloor N_\lambda(z_\lambda^*) \rfloor = N_\lambda^*$ , and  $S_\lambda(z_\lambda^*) = C(N_\lambda^*, \lambda)$ .
- iii.  $N_\lambda(z_\lambda^*) \leq N_\lambda^* - 1$ . Then

$$z_\lambda^* \leq \frac{\lceil N_\lambda(z_\lambda^*) \rceil - \lambda/\mu}{\sqrt{\lambda/\mu}} \leq \frac{N_\lambda^* - 1 - \lambda/\mu}{\sqrt{\lambda/\mu}} \leq \frac{\lfloor N_\lambda(x_\lambda^*) \rfloor - \lambda/\mu}{\sqrt{\lambda/\mu}} \leq x_\lambda^*,$$

so that

$$S_\lambda(z_\lambda^*) - F(\lambda/\mu) \leq C(\lceil N_\lambda(z_\lambda^*) \rceil, \lambda) - F(\lambda/\mu) = C_\lambda\left(\frac{\lceil N_\lambda(z_\lambda^*) \rceil - \lambda/\mu}{\sqrt{\lambda/\mu}}\right) \leq C_\lambda(z_\lambda^*)$$

because of the unimodality of  $C_\lambda(\cdot)$ .

iv.  $N_\lambda(z_\lambda^*) \geq N_\lambda^* + 1$ . Then

$$z_\lambda^* \geq \frac{\lfloor N_\lambda(z_\lambda^*) \rfloor - \lambda/\mu}{\sqrt{\lambda/\mu}} \geq \frac{N_\lambda^* + 1 - \lambda/\mu}{\sqrt{\lambda/\mu}} \leq \frac{\lceil N_\lambda(x_\lambda^*) \rceil - \lambda/\mu}{\sqrt{\lambda/\mu}} \leq x_\lambda^*,$$

so that

$$S_\lambda(z_\lambda^*) - F(\lambda/\mu) \leq C(\lfloor N_\lambda(z_\lambda^*) \rfloor, \lambda) - F(\lambda/\mu) = C_\lambda\left(\frac{\lfloor N_\lambda(z_\lambda^*) \rfloor - \lambda/\mu}{\sqrt{\lambda/\mu}}\right) \leq C_\lambda(z_\lambda^*).$$

Thus, for all  $\lambda$ ,

$$\begin{aligned} S_\lambda(z_\lambda^*) - F(\lambda/\mu) &\leq \max\{C(N_\lambda^*, \lambda) - F(\lambda/\mu), C_\lambda(z_\lambda^*)\} \\ &\approx \max\{C(N_\lambda^*, \lambda) - F(\lambda/\mu), C_\lambda(x_\lambda^*)\} \\ &= C(N_\lambda^*, \lambda) - F(\lambda/\mu). \end{aligned}$$

□

Combining Lemmas 3.1 and 3.2, we obtain the fundamental approximation principle underlying our approach:

**Corollary 3.3** (*Asymptotic Optimality*)

Denote  $\hat{C}_\lambda(z) = C[z; \hat{F}_\lambda, \hat{\pi}_\lambda, \hat{G}_\lambda]$ . Let  $x_\lambda^*$  and  $z_\lambda^*$  be as in (8) and (9) respectively. If  $C_\lambda(x_\lambda^*) \approx \hat{C}_\lambda(x_\lambda^*)$  and  $C_\lambda(z_\lambda^*) \approx \hat{C}_\lambda(z_\lambda^*)$ , then the staffing function  $z_\lambda^*$  is asymptotically optimal in the sense that, as  $\lambda \rightarrow \infty$ ,

$$S_\lambda(z_\lambda^*) - F(\lambda/\mu) \approx C(N_\lambda^*, \lambda) - F(\lambda/\mu),$$

with  $S_\lambda(x)$  given in (10).

Note that the quantities  $S_\lambda(z_\lambda^*) - F(\lambda/\mu)$  and  $C(N_\lambda^*, \lambda) - F(\lambda/\mu)$  may be interpreted as the total cost in excess of the minimum required staffing cost  $F(\lambda/\mu)$  for the approximately optimal staffing level  $N_\lambda(z_\lambda^*)$  and for the truly optimal level  $N_\lambda^*$ , respectively. The above corollary identifies conditions under which these two quantities are asymptotically equal, implying that the approximate solution is asymptotically optimal in a certain sense.

In the next sections, we will identify ‘simple’ functions  $\hat{F}_\lambda(\cdot)$ ,  $\hat{G}_\lambda(\cdot)$ , and  $\hat{\pi}_\lambda(\cdot)$ , such that  $F_\lambda(x_\lambda^*) \approx \hat{F}_\lambda(x_\lambda^*)$  and  $F_\lambda(z_\lambda^*) \approx \hat{F}_\lambda(z_\lambda^*)$ ,  $G_\lambda(x_\lambda^*) \approx \hat{G}_\lambda(x_\lambda^*)$  and  $G_\lambda(z_\lambda^*) \approx \hat{G}_\lambda(z_\lambda^*)$ ,  $\pi_\lambda(x_\lambda^*) \approx \hat{\pi}_\lambda(x_\lambda^*)$  and  $\pi_\lambda(z_\lambda^*) \approx \hat{\pi}_\lambda(z_\lambda^*)$ . This implies  $C_\lambda(x_\lambda^*) \approx \hat{C}_\lambda(x_\lambda^*)$  and  $C_\lambda(z_\lambda^*) \approx \hat{C}_\lambda(z_\lambda^*)$  as required in the above corollary, which then will enable us to gain insight into the behavior of  $N_\lambda^*$ , as a function of  $\lambda$ .



## 4 Some special functions

In this section we introduce some functions that will play a central role in our analysis.

For any  $x > 0$ , define

$$P(x) := \frac{1}{1 + \frac{x}{h(-x)}}, \quad (11)$$

where  $h(\cdot)$  is the ‘hazard rate’ function of the standard normal distribution, namely

$$h(x) := \frac{\phi(x)}{1 - \Phi(x)},$$

with

$$\phi(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}; \quad \Phi(x) = \int_{-\infty}^x \phi(y) dy.$$

In Lemma B.1 we prove that  $P(\cdot)$  is strictly convex decreasing.

Also define

$$Q_\lambda(x) := \frac{\exp\{N_\lambda(x)[1 - r_\lambda(x) + \log r_\lambda(x)]\}}{\sqrt{2\pi N_\lambda(x)}(1 - r_\lambda(x))},$$

with

$$r_\lambda(x) := \frac{\lambda/\mu}{N_\lambda(x)},$$

and let

$$Q(x) := \frac{\phi(x)}{x} = \frac{e^{-x^2/2}}{x\sqrt{2\pi}}. \quad (12)$$

The following two lemmas characterize the asymptotic behavior of  $\pi_\lambda(\cdot)$ , as  $\lambda \rightarrow \infty$ .

**Lemma 4.1** (*Halfin & Whitt [8]*)

For any function  $x_\lambda$  with  $\limsup_{\lambda \rightarrow \infty} x_\lambda < \infty$ ,

$$\pi_\lambda(x_\lambda) \stackrel{\infty}{\sim} P(x_\lambda).$$

If, moreover,  $\lim_{\lambda \rightarrow \infty} x_\lambda = x$ , then  $\pi_\lambda(x_\lambda) \stackrel{\infty}{\sim} P(x)$ ,  $x \geq 0$ .

In particular, if  $\lim_{\lambda \rightarrow \infty} x_\lambda = 0$ , then  $\pi_\lambda(x_\lambda) \stackrel{\infty}{\sim} 1$ .

**Proof**

Suppose to the contrary. Then there must be a subsequence  $\{\lambda_n\}$  with  $\lim_{n \rightarrow \infty} \lambda_n = \infty$  such that  $\lim_{n \rightarrow \infty} x_{\lambda_n} = \beta$  and  $\lim_{n \rightarrow \infty} \pi_{\lambda_n}(x_{\lambda_n}) = \alpha$ , where  $0 < \beta < \infty$  and  $\alpha \neq P(\beta)$ . This is in contradiction with Proposition 1 of Halfin & Whitt [8], which asserts that  $\alpha = P(\beta)$  must prevail for such a sequence  $\{\lambda_n\}$ .

□

**Lemma 4.2** (*Appendix A*)

For any function  $x_\lambda$  with  $\lim_{\lambda \rightarrow \infty} x_\lambda = \infty$ ,

$$\pi_\lambda(x_\lambda) \stackrel{\infty}{\sim} Q_\lambda(x_\lambda).$$

If also  $x_\lambda \stackrel{\sup}{\leq} \lambda^{1/6}$ , then

$$\pi_\lambda(x_\lambda) \stackrel{\infty}{\sim} Q(x_\lambda).$$

If specifically  $x_\lambda = \kappa \sqrt{\lambda/\mu}$  for some constant  $\kappa > 0$ , then

$$\pi_\lambda(x_\lambda) \stackrel{\infty}{\sim} \frac{1}{\kappa \sqrt{2\pi\lambda/\mu(1+\kappa)}} \left( \frac{e^\kappa}{(1+\kappa)^{1+\kappa}} \right)^{\lambda/\mu}.$$

□

We conclude the section with some observations on the behavior of the functions  $F_\lambda(\cdot)$ ,  $G_\lambda(\cdot)$ , and  $\pi_\lambda(\cdot)$ , as defined at the outset of Section 3.

Recall that the staffing cost function  $F(\cdot)$  is convex increasing, which implies that the function  $F_\lambda(\cdot)$  is convex increasing as well. In addition,  $F_\lambda(0) = 0$ . Hence,  $F_\lambda(x)/F_\lambda(y) \leq x/y$  for any pair of numbers  $x \leq y$ . Thus,

$$a_\lambda \stackrel{\sup}{>} b_\lambda \implies F_\lambda(a_\lambda) \stackrel{\sup}{>} F_\lambda(b_\lambda), \quad (13)$$

and

$$a_\lambda \stackrel{\sup}{\gg} b_\lambda \implies F_\lambda(a_\lambda) \stackrel{\sup}{\gg} F_\lambda(b_\lambda). \quad (14)$$

Also, from Lemmas 4.1 and B.1, for fixed  $b \geq 0$ ,

$$a_\lambda \stackrel{\inf}{<} b \implies P(a_\lambda) \stackrel{\sup}{>} P(b), \pi_\lambda(a_\lambda) \stackrel{\sup}{>} \pi_\lambda(b), \quad (15)$$

$$a_\lambda \stackrel{\sup}{>} b \implies P(a_\lambda) \stackrel{\inf}{<} P(b), \pi_\lambda(a_\lambda) \stackrel{\inf}{<} \pi_\lambda(b), \quad (16)$$

and noting that  $\lim_{x \rightarrow \infty} \pi_\lambda(x) = 0$ ,

$$a_\lambda \stackrel{\sup}{\gg} b \implies P(a_\lambda) \stackrel{\inf}{\ll} P(b), \pi_\lambda(a_\lambda) \stackrel{\inf}{\ll} \pi_\lambda(b). \quad (17)$$

## 5 Case I: Rationalized regime

In this section we consider what we call a *rationalized* scenario, by which we mean that, for some  $\kappa > 0$ ,

$$F_\lambda(\kappa) \stackrel{\infty}{\approx} G_\lambda(\kappa), \quad (18)$$

or in words, the staffing cost  $F_\lambda(\cdot)$  is comparable to the waiting cost  $G_\lambda(\cdot)$ , as  $\lambda \rightarrow \infty$ .

For any  $\lambda > 0$ , define

$$y_\lambda^* := \arg \min_{y > 0} C[y; F_\lambda, P, G_\lambda], \quad (19)$$

with  $P(\cdot)$  as in (11).

**Theorem 5.1**

*The staffing function  $y_\lambda^*$  is asymptotically optimal in the sense of Corollary 3.3.*

**Proof**

We start with showing that  $\limsup_{\lambda \rightarrow \infty} x_\lambda^* < \infty$ . Suppose to the contrary. Then  $x_\lambda^* \xrightarrow{\sup} \infty$ , so that

$$F_\lambda(x_\lambda^*) \xrightarrow{\sup} \infty$$

from (14). By definition,

$$C_\lambda(x_\lambda^*) = F_\lambda(x_\lambda^*) + \pi_\lambda(x_\lambda^*)G_\lambda(x_\lambda^*) \geq F_\lambda(x_\lambda^*).$$

Using (18),

$$F_\lambda(\kappa) \stackrel{\infty}{\approx} F_\lambda(\kappa) + G_\lambda(\kappa) \geq F_\lambda(\kappa) + \pi_\lambda(\kappa)G_\lambda(\kappa) = C_\lambda(\kappa).$$

Combining the above relations, we deduce

$$C_\lambda(x_\lambda^*) \xrightarrow{\sup} \infty \geq C_\lambda(\kappa),$$

contradicting the optimality of  $x_\lambda^*$ .

Thus,  $\limsup_{\lambda \rightarrow \infty} x_\lambda^* < \infty$ .

By similar arguments,  $\limsup_{\lambda \rightarrow \infty} y_\lambda^* < \infty$ .

Hence, according to Lemma 4.1,  $\pi_\lambda(x_\lambda^*) \stackrel{\infty}{\approx} \hat{\pi}_\lambda(x_\lambda^*) = P(x_\lambda^*)$  and  $\pi_\lambda(y_\lambda^*) \stackrel{\infty}{\approx} \hat{\pi}_\lambda(y_\lambda^*) = P(y_\lambda^*)$ .

Applying Corollary 3.3 then completes the proof.  $\square$

**Proposition 5.2**

*Assume that there exist functions  $f(\cdot)$ ,  $g(\cdot)$ , and  $H_\lambda$  such that, for any function  $k_\lambda > 0$ ,*

$$F_\lambda(k_\lambda) \stackrel{\infty}{\approx} f(k_\lambda)H_\lambda, \quad (20)$$

*and*

$$G_\lambda(k_\lambda) \stackrel{\infty}{\approx} g(k_\lambda)H_\lambda, \quad (21)$$

*so certainly  $F_\lambda(\kappa) \stackrel{\infty}{\approx} G_\lambda(\kappa)$  for all  $\kappa > 0$ .*

Define

$$y^* = \arg \min_{y>0} C[y; f, P, g].$$

The staffing function  $y^*$  is then asymptotically optimal in the sense of Corollary 3.3.

### Proof

Similar to that of Theorem 5.1. Note that  $F_\lambda(x_\lambda^*) \approx \hat{F}_\lambda(x_\lambda^*) = f(x_\lambda^*)H_\lambda$ ,  $F_\lambda(y^*) \approx \hat{F}_\lambda(y^*) = f(y^*)H_\lambda$ ,  $G_\lambda(x_\lambda^*) \approx \hat{G}_\lambda(x_\lambda^*) = g(x_\lambda^*)H_\lambda$ , and  $G_\lambda(y^*) \approx \hat{G}_\lambda(y^*) = g(y^*)H_\lambda$ . □

### Example 5.3

Assume there is a staffing cost  $c_\lambda$  per server per time unit, and a waiting cost  $a_\lambda$  per customer per time unit, as well as a fixed penalty cost  $b_\lambda$  when the waiting time exceeds  $d_\lambda$  time units, i.e.,  $F(N) := Nc_\lambda$  and  $D_\lambda(t) = a_\lambda t + b_\lambda \mathbf{I}_{\{t>d_\lambda\}}$ , so that  $G(N, \lambda) = \frac{a_\lambda}{N\mu - \lambda} + b_\lambda e^{-(N\mu - \lambda)d_\lambda}$ .

Thus  $F_\lambda(\kappa) = c_\lambda \kappa \sqrt{\lambda/\mu}$  and  $G_\lambda(\kappa) = \frac{a_\lambda \sqrt{\lambda/\mu}}{\kappa} + b_\lambda \lambda e^{-\mu \kappa d_\lambda \sqrt{\lambda/\mu}}$ .

First suppose that  $a_\lambda \approx a$ ,  $b_\lambda \sqrt{\lambda} e^{-\mu \kappa d_\lambda \sqrt{\lambda/\mu}} \approx 1$ , and  $c_\lambda \approx c$ . Then (20)-(21) are satisfied for  $f(\kappa) = c\kappa$ ,  $g(\kappa) = \frac{a}{\kappa}$ , and  $H_\lambda = \sqrt{\lambda/\mu}$ . Proposition 5.2 says that the staffing function

$$y^* = \arg \min_{y>0} \left\{ cy + \frac{aP(y)}{y} \right\} \quad (22)$$

is asymptotically optimal in the sense of Corollary 3.3. This is exactly the  $y^*(\frac{a}{c})$  that arose in (1) of the Introduction. The numerical search for  $y^*$  is straightforward, since the function that it minimizes is unimodal. It is important enough for our purposes, as demonstrated in the Introduction, that we plotted it in Figures 1 and 2.

Now suppose that  $a_\lambda \approx 1$ ,  $b_\lambda \approx b/\sqrt{\lambda}$ ,  $c_\lambda \approx c$ , and  $d_\lambda \approx d/\sqrt{\lambda}$  with  $bd\mu > c$ . Then (20)-(21) are satisfied for  $f(\kappa) = c\kappa$ ,  $g(\kappa) = b\sqrt{\mu} e^{-d\kappa\sqrt{\mu}}$ , and  $H_\lambda = \sqrt{\lambda/\mu}$ . The asymptotically optimal staffing function is

$$y^* = \arg \min_{y>0} \left\{ cy + bP(y)\sqrt{\mu} e^{-dy\sqrt{\mu}} \right\}.$$

Finally, consider the ‘combined’ case where all costs show up in the limit, i.e., suppose that  $a_\lambda \approx a$ ,  $b_\lambda \approx b/\sqrt{\lambda}$ ,  $c_\lambda \approx c$ , and  $d_\lambda \approx d/\sqrt{\lambda}$ . Then (20)-(21) are satisfied for  $f(\kappa) = c\kappa$ ,  $g(\kappa) = \frac{a}{\kappa} + b\sqrt{\mu} e^{-d\kappa\sqrt{\mu}}$ , and  $H_\lambda = \sqrt{\lambda/\mu}$ . The asymptotically optimal staffing function is

$$y^* = \arg \min_{y>0} \left\{ cy + P(y) \left[ \frac{a}{y} + b\sqrt{\mu} e^{-dy\sqrt{\mu}} \right] \right\}.$$

□

**Remark 5.4** (*Asymptotic expansions for  $y^*(\cdot)$* )

Two asymptotic expansions for  $y^*(r)$  were quoted in the Introduction, one for small  $r$  and the other for large. To derive the former, one simply replaces  $P(y)$  in (22) by  $P(0) + yP'(0) + \frac{1}{2}y^2P''(0)$ , then uses the values for the derivatives from Appendix B, and finally minimizes the resulting simple function. As for large values, one uses the well-known approximation  $1 - \Phi(y) \approx \phi(y)/y$  to get that also  $P(y) \approx \phi(y)/y$ . Substituting this approximation for  $P(y)$  into (22) identifies  $y^*$  as the solution  $y$  of

$$ye^{y^2/2} = \frac{a}{\sqrt{2\pi}}.$$

Changing variables to  $x = y^2$ , then squaring, gives  $x(t)e^{x(t)} = t$ , where  $t = \frac{a^2}{2\pi}$ . A complete asymptotic expansion of  $x(t)$ , for large  $t$ , is calculated in [2], pages 25–28. Only its first two terms were used by the approximation in the Introduction, namely

$$x(t) = \log t - \log \log t + O\left(\frac{\log \log t}{\log t}\right).$$

Our qualitative assessments in the Introduction, about the above two approximations, are based on plotting them against the true  $y^*$ . □

## 6 Case II: Efficiency-driven regime

In this section we consider an *efficiency-driven* scenario, meaning that, for all  $\kappa > 0$ ,

$$F_\lambda(\kappa) \gg^\infty G_\lambda(\kappa), \tag{23}$$

( $\limsup_{\lambda \rightarrow \infty} \frac{G_\lambda(\kappa)}{F_\lambda(\kappa)} < \infty$  is actually sufficient) or in words, the staffing cost dominates the waiting cost, as  $\lambda \rightarrow \infty$ .

For any  $\lambda > 0$ , define

$$y_\lambda^* := \arg \min_{y > 0} C[y; F_\lambda, 1, G_\lambda]. \tag{24}$$

### Theorem 6.1

*The staffing level  $y_\lambda^*$  is asymptotically optimal in the sense of Corollary 3.3.*

### Proof

We start with showing that  $\lim_{\lambda \rightarrow \infty} x_\lambda^* = 0$ . Suppose to the contrary. Then  $x_\lambda^* \sup > u$  for some  $u > 0$ , so that

$$F_\lambda(x_\lambda^*) \sup > F_\lambda(u)$$

from (13).

By definition,

$$C_\lambda(x_\lambda^*) = F_\lambda(x_\lambda^*) + \pi_\lambda(x_\lambda^*)G_\lambda(x_\lambda^*) \geq F_\lambda(x_\lambda^*).$$

Using (23),

$$F_\lambda(u) \approx F_\lambda(u) + G_\lambda(u) \geq F_\lambda(u) + \pi_\lambda(u)G_\lambda(u) = C_\lambda(u).$$

Combining the above relations, we obtain

$$C_\lambda(x_\lambda^*) \stackrel{\sup}{>} C_\lambda(u),$$

contradicting the optimality of  $x_\lambda^*$ .

Thus,  $\lim_{\lambda \rightarrow \infty} x_\lambda^* = 0$ .

By similar arguments,  $\lim_{\lambda \rightarrow \infty} y_\lambda^* = 0$ .

Hence, according to Lemma 4.1,  $\pi_\lambda(x_\lambda^*) \approx \hat{\pi}_\lambda(x_\lambda^*) = P(x_\lambda^*) \approx 1$ , and  $\pi_\lambda(y_\lambda^*) \approx \hat{\pi}_\lambda(y_\lambda^*) = P(y_\lambda^*) \approx 1$ .

Applying Corollary 3.3 then completes the proof.  $\square$

### Proposition 6.2

Assume that there exist functions  $f(\cdot)$ ,  $g(\cdot)$ ,  $H_\lambda$ , and  $J_\lambda$  such that, for any function  $k_\lambda > 0$ ,

$$F_\lambda(k_\lambda) \approx f(k_\lambda)H_\lambda, \tag{25}$$

and

$$G_\lambda(k_\lambda) \approx g(k_\lambda)H_\lambda J_\lambda, \tag{26}$$

with  $J_\lambda \ll 1$ , so certainly  $F_\lambda(\kappa) \gg G_\lambda(\kappa)$  for all  $\kappa > 0$ .

Define

$$y_\lambda^* = \arg \min_{y > 0} C[y; f, 1, gJ_\lambda].$$

The staffing function  $y_\lambda^*$  is then asymptotically optimal in the sense of Corollary 3.3.

### Proof

Similar to that of Theorem 6.1. Note that  $F_\lambda(x_\lambda^*) \approx \hat{F}_\lambda(x_\lambda^*) = f(x_\lambda^*)H_\lambda$ ,  $F_\lambda(y_\lambda^*) \approx \hat{F}_\lambda(y_\lambda^*) = f(y_\lambda^*)H_\lambda$ ,  $G_\lambda(x_\lambda^*) \approx \hat{G}_\lambda(x_\lambda^*) = g(x_\lambda^*)H_\lambda J_\lambda$ , and  $G_\lambda(y_\lambda^*) \approx \hat{G}_\lambda(y_\lambda^*) = g(y_\lambda^*)H_\lambda J_\lambda$ .  $\square$

**Remark 6.3** The rationalized staffing function (19) is in fact also asymptotically optimal in the efficiency-driven regime as defined by (23). However, the staffing function (24), where  $P(\cdot)$  is replaced by 1, is asymptotically optimal as well, while considerably simpler.  $\square$

### Example 6.4

Assume  $F(N) := Nc_\lambda$  and  $D_\lambda(t) = a_\lambda t + b_\lambda \mathbf{I}_{\{t > d_\lambda\}}$ , so that  $G(N, \lambda) = \frac{a_\lambda}{N\mu - \lambda} + b_\lambda e^{-(N\mu - \lambda)d_\lambda}$ .

Thus  $F_\lambda(\kappa) = c_\lambda \kappa \sqrt{\lambda/\mu}$  and  $G_\lambda(\kappa) = \frac{a_\lambda \sqrt{\lambda/\mu}}{\kappa} + b_\lambda \lambda e^{-\mu \kappa d_\lambda \sqrt{\lambda/\mu}}$ .

First suppose that  $a_\lambda \stackrel{\infty}{\sim} aJ_\lambda$ ,  $b_\lambda \sqrt{\lambda} e^{-\mu \kappa d_\lambda \sqrt{\lambda/\mu}} \stackrel{\infty}{\ll} J_\lambda$ , and  $c_\lambda \stackrel{\infty}{\sim} c$ , with  $J_\lambda \stackrel{\infty}{\ll} 1$ . Then (25)-(26) are satisfied for  $f(\kappa) = c\kappa$ ,  $g(\kappa) = \frac{a}{\kappa}$ , and  $H_\lambda = \sqrt{\lambda/\mu}$ . Proposition 6.2 says that the staffing function

$$y_\lambda^* = \arg \min_{y>0} \left\{ cy + \frac{a}{y} J_\lambda \right\}$$

is asymptotically optimal in the sense of Corollary 3.3.

Next, consider the ‘combined’ case where all costs show up in the limit, i.e., suppose that  $a_\lambda \stackrel{\infty}{\sim} aJ_\lambda$ ,  $b_\lambda \stackrel{\infty}{\sim} bJ_\lambda/\sqrt{\lambda}$ ,  $c_\lambda \stackrel{\infty}{\sim} c$ ,  $d_\lambda \stackrel{\infty}{\sim} d/\sqrt{\lambda}$ , and  $J_\lambda \stackrel{\infty}{\ll} 1$ . Then (25)-(26) are satisfied for  $f(\kappa) = c\kappa$ ,  $g(\kappa) = \frac{a}{\kappa} + b\sqrt{\mu} e^{-d\kappa\sqrt{\mu}}$ , and  $H_\lambda = \sqrt{\lambda/\mu}$ . The asymptotically optimal staffing function is

$$y_\lambda^* = \arg \min_{y>0} \left\{ cy + \left[ \frac{a}{y} + b\sqrt{\mu} e^{-dy\sqrt{\mu}} \right] J_\lambda \right\}.$$

□

## 7 Case III: Quality-driven regime

In this section we consider a *quality-driven* regime, meaning that, for all  $\kappa > 0$ ,

$$F_\lambda(\kappa) \stackrel{\infty}{\ll} G_\lambda(\kappa), \tag{27}$$

or in words, the staffing cost is negligible compared to the waiting cost, as  $\lambda \rightarrow \infty$ .

For any  $\lambda > 0$ , define

$$y_\lambda^* := \arg \min_{y>0} C[y; F_\lambda, Q_\lambda, G_\lambda].$$

### Theorem 7.1

*The staffing function  $y_\lambda^*$  is asymptotically optimal in the sense of Corollary 3.3.*

### Proof

We start with showing that  $\lim_{\lambda \rightarrow \infty} x_\lambda^* = \infty$ . Suppose to the contrary. Then  $x_\lambda^* \stackrel{\inf}{<} u$  for some  $u > 0$ , so that

$$\pi_\lambda(x_\lambda^*) G_\lambda(x_\lambda^*) \stackrel{\sup}{>} \pi_\lambda(u) G_\lambda(u)$$

from (15) and the fact that  $G_\lambda(\cdot)$  is decreasing.

By definition,

$$C_\lambda(x_\lambda^*) = F_\lambda(x_\lambda^*) + \pi_\lambda(x_\lambda^*)G_\lambda(x_\lambda^*) \geq \pi_\lambda(x_\lambda^*)G_\lambda(x_\lambda^*).$$

Using Lemma 4.2 and (27),

$$\pi_\lambda(u)G_\lambda(u) \approx P(u)G_\lambda(u) \approx F_\lambda(u) + P(u)G_\lambda(u) \approx F_\lambda(u) + \pi_\lambda(u)G_\lambda(u) = C_\lambda(u).$$

Combining the above relations, we deduce

$$C_\lambda(x_\lambda^*) \stackrel{\sup}{>} C_\lambda(u),$$

contradicting the optimality of  $x_\lambda^*$ .

Thus,  $\lim_{\lambda \rightarrow \infty} x_\lambda^* = \infty$ .

By similar arguments,  $\lim_{\lambda \rightarrow \infty} y_\lambda^* = \infty$ .

Hence, according to Lemma 4.2,  $\pi_\lambda(x_\lambda^*) \approx \hat{\pi}_\lambda(y_\lambda^*) = Q_\lambda(x_\lambda^*)$ , and  $\pi_\lambda(y_\lambda^*) \approx \hat{\pi}_\lambda(y_\lambda^*) = Q_\lambda(y_\lambda^*)$ .

Applying Corollary 3.3 then completes the proof.  $\square$

### Proposition 7.2

Assume that there exist functions  $f(\cdot)$ ,  $g(\cdot)$ ,  $H_\lambda$ , and  $J_\lambda$  such that, for any function  $k_\lambda > 0$ ,

$$F_\lambda(k_\lambda) \approx f(k_\lambda)H_\lambda \tag{28}$$

and

$$G_\lambda(k_\lambda) \approx g(k_\lambda)H_\lambda J_\lambda, \tag{29}$$

with  $J_\lambda \gg 1$ , and  $g(\lambda^\alpha)J_\lambda Q(\lambda^\alpha) \ll f(\lambda^\alpha)$  for some  $\alpha < 1/6$ , so that certainly  $F_\lambda(\kappa) \ll G_\lambda(\kappa)$  for all  $\kappa > 0$ .

Define

$$y_\lambda^* = \arg \min_{y > 0} C[y; f, Q, gJ_\lambda]. \tag{30}$$

The staffing function  $y_\lambda^*$  is then asymptotically optimal in the sense of Corollary 3.3.

### Proof

We start with showing that  $x_\lambda^* \stackrel{\sup}{\leq} \lambda^\alpha$ . Suppose to the contrary. Then  $x_\lambda^* \stackrel{\sup}{>} \lambda^\alpha$ , so that

$$F_\lambda(x_\lambda^*) \stackrel{\sup}{>} F_\lambda(\lambda^\alpha)$$

from (13).

By definition,

$$C_\lambda(x_\lambda^*) = F_\lambda(x_\lambda^*) + \pi_\lambda(x_\lambda^*)G_\lambda(x_\lambda^*) \geq F_\lambda(x_\lambda^*).$$



Using Lemma 4.2 and (28), (29),

$$F_\lambda(\lambda^\alpha) \approx F_\lambda(\lambda^\alpha) + Q(\lambda^\alpha)G_\lambda(\lambda^\alpha) \approx F_\lambda(\lambda^\alpha) + \pi_\lambda(\lambda^\alpha)G_\lambda(\lambda^\alpha) = C_\lambda(\lambda^\alpha).$$

Combining the above relations, we obtain

$$C_\lambda(x_\lambda^*) \stackrel{\sup}{>} C_\lambda(\lambda^\alpha),$$

contradicting the optimality of  $x_\lambda^*$ .

Thus,  $x_\lambda^* \stackrel{\sup}{\leq} \lambda^\alpha$ .

By similar arguments,  $y_\lambda^* \stackrel{\sup}{\leq} \lambda^\alpha$ .

It may further be shown that  $\lim_{\lambda \rightarrow \infty} x_\lambda^* = \infty$  and  $\lim_{\lambda \rightarrow \infty} y_\lambda^* = \infty$  in a similar fashion as in the proof of Theorem 7.1. Hence, according to Lemma 4.2,  $\pi_\lambda(x_\lambda^*) \approx \hat{\pi}_\lambda(x_\lambda^*) = Q(x_\lambda^*)$ , and  $\pi_\lambda(y_\lambda^*) \approx \hat{\pi}_\lambda(y_\lambda^*) = Q(y_\lambda^*)$ .

Applying Corollary 3.3 then completes the proof. Note that  $F_\lambda(x_\lambda^*) \approx \hat{F}_\lambda(x_\lambda^*) = f(x_\lambda^*)H_\lambda$ ,  $F_\lambda(y_\lambda^*) \approx \hat{F}_\lambda(y_\lambda^*) = f(y_\lambda^*)H_\lambda$ ,  $G_\lambda(x_\lambda^*) \approx \hat{G}_\lambda(x_\lambda^*) = g(x_\lambda^*)H_\lambda J_\lambda$ , and  $G_\lambda(y_\lambda^*) \approx \hat{G}_\lambda(y_\lambda^*) = g(y_\lambda^*)H_\lambda J_\lambda$ . □

**Remark 7.3** The rationalized staffing function (19) is in fact also asymptotically optimal in the quality-driven regime as defined by (27) when  $x_\lambda^* \ll \lambda^{1/6}$ , since the proof of Theorem 7.2 then shows that  $\pi_\lambda(x_\lambda^*) \approx Q(x_\lambda^*)$ , while  $Q(x_\lambda^*) \approx P(x_\lambda^*)$ . However, the staffing function (30) is asymptotically optimal as well, while simpler. □

#### Example 7.4

Assume  $F(N) := Nc_\lambda$  and  $D_\lambda(t) = a_\lambda t + b_\lambda \mathbf{I}_{\{t > d_\lambda\}}$ , so that  $G(N, \lambda) = \frac{a_\lambda}{N\mu - \lambda} + b_\lambda e^{-(N\mu - \lambda)d_\lambda}$ .

Thus  $F_\lambda(\kappa) = c_\lambda \kappa \sqrt{\lambda/\mu}$  and  $G_\lambda(\kappa) = \frac{a_\lambda \sqrt{\lambda/\mu}}{\kappa} + b_\lambda \lambda e^{-\mu \kappa d_\lambda} \sqrt{\lambda/\mu}$ .

First suppose that  $a_\lambda \approx aJ_\lambda$ ,  $b_\lambda \sqrt{\lambda} e^{-\mu \kappa d_\lambda} \sqrt{\lambda/\mu} \ll J_\lambda$ , and  $c_\lambda \approx c$ , with  $J_\lambda \gg 1$ . Then (28)-(29) are satisfied for  $f(\kappa) = c\kappa$ ,  $g(\kappa) = \frac{a}{\kappa}$ , and  $H_\lambda = \sqrt{\lambda/\mu}$ . Proposition 7.2 says that the staffing function

$$y_\lambda^* = \arg \min_{y > 0} \left\{ cy + \frac{aQ(y)}{y} J_\lambda \right\}$$

is asymptotically optimal in the sense of Corollary 3.3.

Now suppose that  $a_\lambda \ll J_\lambda$ ,  $b_\lambda \approx bJ_\lambda/\sqrt{\lambda}$ ,  $c_\lambda \approx c$ , and  $d_\lambda \approx d/\sqrt{\lambda}$ , with  $J_\lambda \gg 1$ . Then (28)-(29) are satisfied for  $f(\kappa) = c\kappa$ ,  $g(\kappa) = b\sqrt{\mu} e^{-d\kappa\sqrt{\mu}}$ , and  $H_\lambda = \sqrt{\lambda/\mu}$ . The asymptotically optimal staffing function is

$$y_\lambda^* = \arg \min_{y > 0} \left\{ cy + bQ(y)\sqrt{\mu} e^{-dy\sqrt{\mu}} J_\lambda \right\}.$$

Finally, consider the ‘combined’ case where all costs show up in the limit, i.e., suppose that  $a_\lambda \approx a J_\lambda$ ,  $b_\lambda \approx b J_\lambda / \sqrt{\lambda}$ ,  $c_\lambda \approx c$ , and  $d_\lambda \approx d / \sqrt{\lambda}$ , with  $J_\lambda \gg 1$ . Then (28)-(29) are satisfied for  $f(\kappa) = c\kappa$ ,  $g(\kappa) = \frac{a}{\kappa} + b\sqrt{\mu} e^{-d\kappa\sqrt{\mu}}$ , and  $H_\lambda = \sqrt{\lambda/\mu}$ . The asymptotically optimal staffing function is

$$y_\lambda^* = \arg \min_{y>0} \left\{ cy + Q(y) \left[ \frac{a}{y} + b\sqrt{\mu} e^{-dy\sqrt{\mu}} \right] J_\lambda \right\}.$$

□

## 8 Constraint satisfaction

In the previous sections, we have considered the problem of determining the staffing level so as to minimize the total staffing and waiting cost. A closely related problem, which is in fact motivated by actual practice, is to minimize the staffing level subject to a constraint  $M_\lambda > 0$  on the waiting cost.

So we are now interested in determining

$$N_\lambda^* := \min_{N>\lambda/\mu} \{N : K(N, \lambda) \leq M_\lambda\}, \quad (31)$$

with  $K(N, \lambda) := \lambda\pi(N, \lambda/\mu)G(N, \lambda)$  denoting the waiting cost. Notice that  $\lim_{N \rightarrow \infty} K(N, \lambda) = 0$ , so  $N_\lambda^*$  is well-defined.

As for the cost minimization problem, our approach is to first translate the discrete problem (31) into a continuous one, and then approximate the latter problem by a related continuous problem which is easier to solve. Denote

$$x_\lambda^* := \min_{x>0} \{x : K_\lambda(x) \leq M_\lambda\},$$

with  $K_\lambda(x) := \pi_\lambda(x)G_\lambda(x)$ . Since  $\pi_\lambda(\cdot)$  and  $G_\lambda(\cdot)$  are both continuous and strictly decreasing,  $x_\lambda^*$  is the unique solution to the equation  $K_\lambda(x) = M_\lambda$ . Further notice that  $N_\lambda^* = \lceil N_\lambda(x_\lambda^*) \rceil$ , which establishes the link between the discrete problem and the corresponding continuous problem.

To approximate  $x_\lambda^*$ , define  $z_\lambda^*$  as the solution to the equation  $\hat{\pi}_\lambda(z)\hat{G}_\lambda(z) = M_\lambda$ . The functions  $\hat{\pi}_\lambda(\cdot)$  and  $\hat{G}_\lambda(\cdot)$  that we consider will always be such that  $z_\lambda^*$  exists and is unique.

We now formulate the approximation principle underlying our approach, in parallel to that for the cost minimization problem.

Define

$$T_\lambda(x) := \min\{ |K(\lfloor N_\lambda(x) \rfloor, \lambda) - M_\lambda|, |K(\lceil N_\lambda(x) \rceil, \lambda) - M_\lambda|, |K(\lceil N_\lambda(x) \rceil, \lambda) - K(N_\lambda^*, \lambda)| \}. \quad (32)$$

### **Lemma 8.1** (*Asymptotic Optimality*)

Denote  $\hat{K}_\lambda(y) = \hat{\pi}_\lambda(y)\hat{G}_\lambda(y)$ . Let  $z_\lambda^*$  be as defined above. If  $K_\lambda(z_\lambda^*) \approx \hat{K}_\lambda(z_\lambda^*)$ , then the staffing function  $z_\lambda^*$  is asymptotically optimal in the sense that, as  $\lambda \rightarrow \infty$ ,  $T_\lambda(z_\lambda^*) \ll^\infty M_\lambda$ , with  $T_\lambda(\cdot)$  given in (32).

## Proof

For fixed  $\lambda$ , we distinguish between three cases.

- i.  $N_\lambda^* - 1 < N_\lambda(z_\lambda^*) \leq N_\lambda^*$ . Then  $\lceil N_\lambda(z_\lambda^*) \rceil = N_\lambda^*$ , so that  $T_\lambda(z_\lambda^*) = 0$ .
- ii.  $N_\lambda(z_\lambda^*) \leq N_\lambda^* - 1$ . Then

$$M_\lambda = K_\lambda(x_\lambda^*) \leq K(N_\lambda^* - 1, \lambda) \leq K(\lceil N_\lambda(z_\lambda^*) \rceil, \lambda) \leq K_\lambda(z_\lambda^*),$$

so that

$$T_\lambda(z_\lambda^*) \leq |K(\lceil N_\lambda(z_\lambda^*) \rceil, \lambda) - M_\lambda| \leq K_\lambda(z_\lambda^*) - M_\lambda.$$

- iii.  $N_\lambda(z_\lambda^*) > N_\lambda^*$ . Then

$$K_\lambda(z_\lambda^*) \leq K(\lfloor N_\lambda(z_\lambda^*) \rfloor, \lambda) \leq K(N_\lambda^*, \lambda) \leq K_\lambda(x_\lambda^*) = M_\lambda,$$

so that

$$T_\lambda(z_\lambda^*) \leq |K(\lfloor N_\lambda(z_\lambda^*) \rfloor, \lambda) - M_\lambda| \leq M_\lambda - K_\lambda(z_\lambda^*).$$

Thus, for all  $\lambda$ ,

$$T_\lambda(z_\lambda^*) \leq |K_\lambda(z_\lambda^*) - M_\lambda| \stackrel{\infty}{\ll} M_\lambda,$$

as  $\hat{K}_\lambda(z_\lambda^*) = M_\lambda$  by definition.

□

In full generality, it seems difficult to establish a stronger optimality property than indicated in the above lemma. Under additional conditions, however, it is possible to make sharper statements. For example, a more desirable criterion for asymptotic optimality would be

$$|K(\lceil N_\lambda(z_\lambda^*) \rceil, \lambda) - K(N_\lambda^*, \lambda)| \stackrel{\infty}{\ll} M_\lambda.$$

And indeed, following the same reasoning of the proof of Lemma 8.1, it can be guaranteed to hold but under additional constraints on the oscillation of our costs.

## 8.1 Rationalized regime

We first consider a rationalized scenario, by which we mean that, for some  $\kappa > 0$ ,

$$G_\lambda(\kappa) \stackrel{\infty}{\approx} M_\lambda, \tag{33}$$

( $\limsup_{\lambda \rightarrow \infty} \frac{G_\lambda(\kappa)}{M_\lambda} < \infty$  is actually sufficient) or in words, the waiting cost  $G_\lambda(\cdot)$  is comparable to the constraint  $M_\lambda$ , as  $\lambda \rightarrow \infty$ .

For any  $\lambda > 0$ , define  $y_\lambda^*$  as the solution to the equation  $P(y)G_\lambda(y) = M_\lambda$ .

Note that  $P(\cdot)$  and  $G_\lambda(\cdot)$  are both continuous and strictly decreasing with  $\lim_{x \downarrow 0} P(x)G_\lambda(x) = \infty$  and  $\lim_{x \rightarrow \infty} P(x)G_\lambda(x) = 0$ , so that  $y_\lambda^*$  exists and is unique.

### Theorem 8.2

The staffing function  $y_\lambda^*$  is asymptotically optimal in the sense of Lemma 8.1.

#### Proof

We start with showing that  $\limsup_{\lambda \rightarrow \infty} y_\lambda^* < \infty$ . Suppose to the contrary. Then  $y_\lambda^* \xrightarrow{\sup} \infty$ , so that

$$P(y_\lambda^*)G_\lambda(y_\lambda^*) \xrightarrow{\inf} P(\kappa)G_\lambda(\kappa)$$

from (17) and the fact that  $G_\lambda(\cdot)$  is decreasing.

Using (33),

$$P(\kappa)G_\lambda(\kappa) \leq G_\lambda(\kappa) \approx M_\lambda.$$

Combining the above relations, we deduce

$$P(y_\lambda^*)G_\lambda(y_\lambda^*) \xrightarrow{\inf} M_\lambda,$$

contradicting the definition of  $y_\lambda^*$ .

Thus,  $\limsup_{\lambda \rightarrow \infty} y_\lambda^* < \infty$ .

Hence, according to Lemma 4.1,  $\pi_\lambda(y_\lambda^*) \approx \hat{\pi}_\lambda(y_\lambda^*) = P(y_\lambda^*)$ .

Applying Lemma 8.1 then completes the proof. □

### Proposition 8.3

Assume that there exists a function  $g(\cdot)$ , such that, for any function  $k_\lambda > 0$ ,

$$G_\lambda(\kappa) \approx g(\kappa)M_\lambda, \tag{34}$$

so certainly  $G_\lambda(\kappa) \approx M_\lambda$  for all  $\kappa > 0$ .

Define  $y^*$  as the solution to the equation  $P(y)g(y) = 1$ .

Assume that  $g(\cdot)$  is continuous and decreasing, with  $\lim_{x \downarrow 0} g(x) > 1$  so that  $y^*$  exists and is unique.

The staffing function  $y^*$  is then asymptotically optimal in the sense of Lemma 8.1.

#### Proof

Similar to that of Theorem 8.2. Note that  $G_\lambda(y^*) \approx \hat{G}_\lambda(y^*) = g(y^*)M_\lambda$ . □

### Example 8.4

Assume there is a waiting cost  $a_\lambda$  per customer per time unit, as well as a fixed penalty cost  $b_\lambda$  when the waiting time exceeds  $d_\lambda$  time units, i.e.,  $D_\lambda(t) = a_\lambda t + b_\lambda \mathbf{I}_{\{t > d_\lambda\}}$ , so that

$$G(N, \lambda) = \frac{a_\lambda}{N\mu - \lambda} + b_\lambda e^{-(N\mu - \lambda)d_\lambda}. \text{ Thus } G_\lambda(\kappa) = \frac{a_\lambda \sqrt{\lambda/\mu}}{\kappa} + b_\lambda \lambda e^{-\mu \kappa d_\lambda \sqrt{\lambda/\mu}}.$$

First suppose that  $a_\lambda \approx a\sqrt{\lambda/\mu}$ ,  $b_\lambda e^{-\mu\kappa d_\lambda\sqrt{\lambda/\mu}} \ll 1$ , and  $M_\lambda = M\lambda$ . Then (34) is satisfied for  $g(\kappa) = \frac{a}{\kappa\mu M}$ . Proposition 8.3 says that the staffing function  $y^*$  determined as the unique solution to the equation

$$\frac{aP(y)}{\mu y} = M$$

is asymptotically optimal in the sense of Lemma 8.1.

Now suppose that  $a_\lambda/\sqrt{\lambda} \ll 1$ ,  $b_\lambda \approx b$ ,  $d_\lambda \approx d/\sqrt{\lambda}$ , and  $M_\lambda = M\lambda$  with  $b > M$ . Then (34) is satisfied for  $g(\kappa) = b e^{-d\kappa\sqrt{\mu}}/M$ . The asymptotically optimal staffing function is the unique solution to the equation

$$bP(y) e^{-dy\sqrt{\mu}} = M.$$

Finally, consider the ‘combined’ case where all costs show up in the limit, i.e., suppose that  $a_\lambda \approx a\sqrt{\lambda/\mu}$ ,  $b_\lambda \approx b$ ,  $d_\lambda \approx d/\sqrt{\lambda}$ , and  $M_\lambda = M\lambda$ . Then (34) is satisfied for  $g(\kappa) = \frac{a}{\kappa\mu M} + b e^{-d\kappa\sqrt{\mu}}/M$ . The asymptotically optimal staffing function is the unique solution to the equation

$$P(y)\left[\frac{a}{\mu y} + b e^{-dy\sqrt{\mu}}\right] = M.$$

□

### Example 8.5

An important special case is  $a_\lambda = 0$ ,  $b_\lambda = 1$ ,  $d_\lambda = 0$ ,  $M_\lambda = \epsilon\lambda$ , which corresponds to a target waiting probability  $\epsilon$ . Case ii) of the above example then shows that the staffing function  $y^* = P^{-1}(\epsilon)$  is asymptotically optimal. We described this example in our Introduction – see (5). It is to be compared with (6), used in [15] and [12]. On the differences and similarities between the two approximations, see Section 9.

□

## 8.2 Efficiency-driven regime

We now consider an efficiency-driven scenario, meaning that, for all  $\kappa > 0$ ,

$$G_\lambda(\kappa) \ll M_\lambda, \tag{35}$$

(in fact  $G_\lambda(\kappa) \leq^{\sup} M_\lambda$  would be sufficient for the results below to hold) or in words, the waiting cost is dominated by the target upper bound, as  $\lambda \rightarrow \infty$ .

For any  $\lambda > 0$ , define  $y_\lambda^*$  as the solution to the equation  $G_\lambda(y) = M_\lambda$ .

Note that  $G_\lambda(\cdot)$  is continuous and strictly decreasing with  $\lim_{x \downarrow 0} G_\lambda(x) = \infty$  and  $\lim_{x \rightarrow \infty} G_\lambda(x) = 0$ , so that  $y_\lambda^*$  exists and is unique.

### Theorem 8.6

*The staffing function  $y_\lambda^*$  is asymptotically optimal in the sense of Lemma 8.1.*

**Proof**

We start with showing that  $\lim_{\lambda \rightarrow \infty} y_\lambda^* = 0$ . Suppose to the contrary. Then  $y_\lambda^* \stackrel{\sup}{>} u$  for some  $u > 0$ , so that

$$P(y_\lambda^*)G_\lambda(y_\lambda^*) \stackrel{\inf}{<} P(u)G_\lambda(u),$$

from (16) and the fact that  $G_\lambda(\cdot)$  is decreasing.

Using (35),

$$P(u)G_\lambda(u) \leq G_\lambda(u) \stackrel{\sup}{\leq} M_\lambda,$$

Combining the above relations, we deduce

$$P(y_\lambda^*)G_\lambda(y_\lambda^*) \stackrel{\inf}{<} M_\lambda,$$

contradicting the definition of  $y_\lambda^*$ .

Thus,  $\lim_{\lambda \rightarrow \infty} y_\lambda^* = 0$ .

Hence, according to Lemma 4.1,  $\pi_\lambda(y_\lambda^*) \stackrel{\infty}{\approx} \hat{\pi}_\lambda(y_\lambda^*) = P(y_\lambda^*) \stackrel{\infty}{\approx} 1$ .

Applying Lemma 8.1 then completes the proof. □

**Proposition 8.7**

Assume that there exist functions  $g(\cdot)$  and  $J_\lambda$  such that, for any function  $k_\lambda > 0$ ,

$$G_\lambda(k_\lambda) \stackrel{\infty}{\approx} g(k_\lambda)J_\lambda M_\lambda, \tag{36}$$

with  $J_\lambda \stackrel{\infty}{\ll} 1$ , so certainly  $G_\lambda(\kappa) \stackrel{\infty}{\ll} M_\lambda$  for all  $\kappa > 0$ .

Define  $y_\lambda^*$  as the unique solution to the equation  $g(y)J_\lambda = 1$ .

Assume that  $g(\cdot)$  is continuous and strictly decreasing, with  $\lim_{x \downarrow 0} g(x) = \infty$  so that  $y_\lambda^*$  exists and is unique.

The staffing function  $y_\lambda^*$  is then asymptotically optimal in the sense of Lemma 8.1.

**Proof**

Similar to that of Theorem 6.1. Note that  $G_\lambda(y_\lambda^*) \stackrel{\infty}{\approx} g(y_\lambda^*)J_\lambda M_\lambda$ . □

**Example 8.8**

Assume  $D_\lambda(t) = a_\lambda t + b_\lambda \mathbf{I}_{\{t > d_\lambda\}}$ , so that  $G(N, \lambda) = \frac{a_\lambda}{N\mu - \lambda} + b_\lambda e^{-(N\mu - \lambda)d_\lambda}$ . Thus  $G_\lambda(\kappa) = \frac{a_\lambda \sqrt{\lambda/\mu}}{\kappa} + b_\lambda \lambda e^{-\mu \kappa d_\lambda} \sqrt{\lambda/\mu}$ .

First suppose that  $a_\lambda \stackrel{\infty}{\approx} a \sqrt{\lambda/\mu}$ ,  $b_\lambda e^{-\mu \kappa d_\lambda} \sqrt{\lambda/\mu} \stackrel{\infty}{\ll} 1$ , and  $M_\lambda = M\lambda/J_\lambda$ , with  $J_\lambda \stackrel{\infty}{\ll} 1$ . Then (36) is satisfied for  $g(\kappa) = \frac{a}{\kappa \mu M}$ . Proposition 8.7 says that the staffing function

$$y_\lambda^* = \frac{a}{\mu M} J_\lambda$$

is asymptotically optimal in the sense of Lemma 8.1.

Next, consider the ‘combined’ case where all costs show up in the limit, i.e., suppose that  $a_\lambda \approx a\sqrt{\lambda/\mu}$ ,  $b_\lambda \approx b$ ,  $d_\lambda \approx d/\sqrt{\lambda}$ , and  $M_\lambda = M\lambda/J_\lambda$ , with  $J_\lambda \ll 1$ . Then (36) is satisfied for  $g(\kappa) = \frac{a}{\kappa\mu M} + b e^{-d\kappa\sqrt{\mu}}/M$ . The asymptotically optimal staffing function is the unique solution to the equation

$$\left[ \frac{a}{\mu y} + b e^{-dy\sqrt{\mu}} \right] J_\lambda = M.$$

□

### 8.3 Quality-driven regime

We finally consider a quality-driven scenario, meaning that, for all  $\kappa > 0$ ,

$$G_\lambda(\kappa) \gg M_\lambda, \tag{37}$$

or in words, the waiting cost dominates the target upper bound, as  $\lambda \rightarrow \infty$ .

For any  $\lambda > 0$ , define  $y_\lambda^*$  as the solution to the equation  $G_\lambda(y)Q_\lambda(y) = M_\lambda$ .

Note that  $G_\lambda(\cdot)$  and  $Q_\lambda(\cdot)$  are continuous and strictly decreasing with  $\lim_{x \downarrow 0} G_\lambda(x)Q_\lambda(x) = \infty$  and  $\lim_{x \rightarrow \infty} G_\lambda(x)Q_\lambda(x) = 0$ , so that  $y_\lambda^*$  exists and is unique.

#### Theorem 8.9

*The staffing function  $y_\lambda^*$  is asymptotically optimal in the sense of Lemma 8.1.*

#### Proof

We start with showing that  $\lim_{\lambda \rightarrow \infty} y_\lambda^* = \infty$ . Suppose to the contrary. Then  $y_\lambda^* < u$  for some  $u > 0$ , so that

$$P(y_\lambda^*)G_\lambda(y_\lambda^*) \stackrel{\sup}{>} P(u)G_\lambda(u)$$

from (15) and the fact that  $G_\lambda(\cdot)$  is decreasing.

Using Lemma 4.2 and (37),

$$P(u)G_\lambda(u) \approx P(u)G_\lambda(u) \gg M_\lambda.$$

Combining the above relations, we deduce

$$P(y_\lambda^*)G_\lambda(y_\lambda^*) \stackrel{\sup}{>} M_\lambda,$$

contradicting the definition of  $y_\lambda^*$ .

Thus,  $\lim_{\lambda \rightarrow \infty} y_\lambda^* = \infty$ .

Hence, according to Lemma 4.2,  $\pi_\lambda(y_\lambda^*) \approx \hat{\pi}_\lambda(y_\lambda^*) = Q_\lambda(y_\lambda^*)$ .

Applying Lemma 8.1 then completes the proof.

□

**Proposition 8.10**

Assume that there exist functions  $g(\cdot)$  and  $J_\lambda$  such that, for any function  $k_\lambda > 0$ ,

$$G_\lambda(k_\lambda) \approx g(k_\lambda) J_\lambda M_\lambda, \quad (38)$$

with  $J_\lambda \gg 1$ , and  $g(\lambda^\alpha) J_\lambda Q(\lambda^\alpha) \ll 1$  for some  $\alpha < 1/6$ , so that certainly  $G_\lambda(\kappa) \gg M_\lambda$  for all  $\kappa > 0$ .

Define  $y_\lambda^*$  as the unique solution to the equation

$$Q(y)g(y)J_\lambda = 1.$$

Assume that  $g(\cdot)$  is continuous and strictly decreasing, with  $\lim_{x \downarrow 0} g(x) > 0$  so that  $y_\lambda^*$  exists and is unique.

The staffing function  $y_\lambda^*$  is then asymptotically optimal in the sense of Lemma 8.1.

**Proof**

It may easily be shown that  $y_\lambda^* \leq \sup \lambda^\alpha$ . Hence, according to Lemma 4.2,  $\pi_\lambda(y_\lambda^*) \approx \hat{\pi}_\lambda(y_\lambda^*) = Q(y_\lambda^*)$ .

The proof is further similar to that of Theorem 8.9. Note that  $G_\lambda(y_\lambda^*) \approx \hat{G}_\lambda(y_\lambda^*) = g(y_\lambda^*) M_\lambda J_\lambda$ .  $\square$

**Example 8.11**

Assume  $D_\lambda(t) = a_\lambda t + b_\lambda \mathbf{I}_{\{t > d_\lambda\}}$ , so that  $G(N, \lambda) = \frac{a_\lambda}{N\mu - \lambda} + b_\lambda e^{-(N\mu - \lambda)d_\lambda}$ . Thus  $G_\lambda(\kappa) = \frac{a_\lambda \sqrt{\lambda/\mu}}{\kappa} + b_\lambda \lambda e^{-\mu \kappa d_\lambda \sqrt{\lambda/\mu}}$ .

First suppose that  $a_\lambda \approx a \sqrt{\lambda/\mu}$ ,  $b_\lambda e^{-\mu \kappa d_\lambda \sqrt{\lambda/\mu}} \ll 1$ , and  $M_\lambda = M\lambda/J_\lambda$ , with  $J_\lambda \gg 1$ . Then (38) is satisfied for  $g(\kappa) = \frac{a}{\kappa \mu M}$ . Proposition 8.10 says that the staffing function  $y_\lambda^*$  determined by the unique solution of the equation

$$\frac{aQ(y)}{\mu y} J_\lambda = M$$

is asymptotically optimal in the sense of Lemma 8.1.

Now suppose that  $a_\lambda \ll \sqrt{\lambda}$ ,  $b_\lambda \approx b$ ,  $d_\lambda \approx d/\sqrt{\lambda}$ , and  $M_\lambda = M\lambda/J_\lambda$ , with  $J_\lambda \gg 1$ . Then (38) is satisfied for  $g(\kappa) = b e^{-d\kappa\sqrt{\mu}}/M$ . The asymptotically optimal staffing function is the unique solution of the equation

$$bQ(y) e^{-dy\sqrt{\mu}} J_\lambda = M.$$

Finally, consider the ‘combined’ case where all costs show up in the limit, i.e., suppose that  $a_\lambda \approx a \sqrt{\lambda/\mu}$ ,  $b_\lambda \approx b$ ,  $d_\lambda \approx d/\sqrt{\lambda}$ , and  $M_\lambda = M\lambda/J_\lambda$ , with  $J_\lambda \gg 1$ . Then (38) is satisfied for  $g(\kappa) = \frac{a}{\kappa \mu M} + b e^{-d\kappa\sqrt{\mu}}/M$ . The asymptotically optimal staffing function is the unique solution of the equation

$$Q(y) \left[ \frac{a}{\mu y} + b e^{-dy\sqrt{\mu}} \right] J_\lambda = M.$$

$\square$



**Remark 8.12** Notice that the results for the constraint satisfaction problem closely mirror those for the cost minimization problem. In fact, the two problems may be formally related as follows. Consider a strictly decreasing function  $M(\cdot)$  on  $(0, \infty)$  with  $\lim_{x \downarrow 0} M(x) = \infty$  and  $\lim_{x \rightarrow \infty} M(x) = 0$ . Then the (unique) solution to the equation  $M(x) = 1$  is  $\arg \min_{x > 0} \{M(x) + 1/M(x)\}$ . Thus, the solution to the constraint satisfaction problem  $\pi_\lambda(x)G_\lambda(x) = M_\lambda$  may also be represented as  $\arg \min_{x > 0} \{\frac{M_\lambda^2}{\pi_\lambda(x)G_\lambda(x)} + \pi_\lambda(x)G_\lambda(x)\}$ , which has the form of the cost minimization problem. The relation thus established is only formal. We could not utilize it to derive the results for constraint satisfaction from the optimization results.  $\square$

## 9 Numerical experiments

In this section we present the results of some numerical experiments that we carried out. The main purpose of the numerical experiments was to test the accuracy of the approximations that arise from our asymptotically optimal staffing levels. The numerical results indicate that the rationalized approximation performs exceptionally well in all regimes. By Remark 6.3 we know that the rationalized approximation is in fact asymptotically optimal in all regimes. On the other hand, the accuracy displayed by the rationalized approximation is astonishingly better than our rigorous results lead us to believe.

Our first two experiments address [7] and [12], which correspond to Examples 5.3 and 8.5, respectively.

*Grassmann [7], Table 3:* Grassmann calculates the optimal staffing level  $N^*$  for the M/M/N queue, with offered loads  $R = 1, 3, 10, 30, 100$  and costs  $r = \frac{a}{c} = 10, 20, 100, 200$ . While the latter are rather extreme values, our approximation (1) is nevertheless accurate: it is exact in 7 cases and off by only 1 agent in the other 13 cases.  $\square$

*Kolesar & Green [12], Table 1:* Kolesar & Green use (6) in order to calculate  $N^*$  that achieves  $\Pr\{\text{Wait} > 0\} = \epsilon$ , for  $\epsilon = 0.2, 0.1, 0.05, 0.025, 0.01, 0.001$ , and offered loads  $R = 2^m, m = 0, 1, \dots, 10$ . The approximation (1) is superior to (6). This is to be expected in view of the theory that supports the former, while the latter is heuristically based. Indeed, for  $\epsilon = 0.2$ , (1) is exact for 9 cases and misses by 1 agent for the other 2. In contrast, (6) misses by up to 7 agents. But more significantly, the misses in staffing levels lead to misses in the target delay probabilities, off by 25-75% in 8 out of the 11 cases.

The approximation (6) improves as  $\epsilon$  decreases, until eventually it coincides with (1). This is understood as follows: small values of  $\epsilon$  give rise to large  $y^*$  (quality-driven), for which

$$\bar{\Phi}(y) \approx \frac{\phi(y)}{y} \approx P(y).$$

□

We now turn to numerical experiments related to Examples 5.3, 6.4, and 7.4. In all cases we compared an approximation to the optimal staffing level obtained from the asymptotics with the exact optimal staffing level, which we obtained through a simple search procedure. (The unimodality of  $C(N, \lambda)$  in  $N$  makes such a search simple.)

In all of the examples we use  $c_\lambda = c = 1$  and  $\mu = 1$ . Once we set  $c_\lambda = c$ , taking  $c = 1$  is without loss of generality because we can take this as the definition of the monetary unit. Similarly, taking  $\mu = 1$  is also without loss of generality because we can take  $1/\mu$  as the definition of the time unit.

We first describe the numerical results related to Example 5.3. We considered three cases:

- i)  $a_\lambda = a, \quad b_\lambda = 0$
- ii)  $a_\lambda = 0, \quad b_\lambda = b/\sqrt{\lambda}, \quad d_\lambda = d/\sqrt{\lambda}$
- iii)  $a_\lambda = a, \quad b_\lambda = b\sqrt{\lambda}, \quad d_\lambda = d/\sqrt{\lambda}.$

These correspond to the three cases in Example 5.3. First consider case i). For  $r > 0$ , let

$$y^*(r) = \arg \min_{y>0} \left\{ y + \frac{rP(y)}{y} \right\}. \quad (39)$$

We plot  $y^*(r)$ ,  $0 \leq r \leq 10$  in Figure 1.

Let

$$n^* = \lambda/\mu + y^*(r)\sqrt{\lambda/\mu}.$$

The rationalized approximation for case i) of Example 5.3 is obtained by rounding  $n^*$  to the nearest integer. (The asymptotic analysis gives no guidance of how to go from  $n^*$  to an integer staffing level. Preliminary numerical calculations comparing rounding up, rounding down, and rounding off showed that rounding off is generally superior. So all of our numerical results involve rounding off. Of course, if rounding off  $n^*$  yields a value smaller than  $\lambda/\mu$ , the staffing level must be increased by 1 to avoid an unstable system.)

To check this approximation we first tried  $\lambda = 100$  and the seven different values  $a = 0.1, 0.25, 0.5, 1, 2, 4, 10$ . In all of these cases rounding off  $n^*$  gave the exact optimal staffing level. We next set  $a = 2$  and tried all integer values of  $\lambda$  between 5 and 100. Here rounding off  $n^*$  is never off by more than 1, and is usually exact (83 out of 96 cases).

For case ii) we let

$$n^* = \lambda/\mu + y^*(b, d)\sqrt{\lambda/\mu},$$

where

$$y^*(b, d) = \arg \min_{y>0} \left\{ y + bP(y) e^{-dy} \right\}.$$

To check this approximation we first tried  $\lambda = 100$  and the seven different values  $b = 0.1, 0.25, 0.5, 1, 2, 4, 10$ . Again we found that in all of these cases rounding off  $n^*$  gave the exact optimal staffing level. We next set  $b = 5$  and  $d = 1$  and tried all integer values between 5 and 100 for  $\lambda$ . Rounding off  $n^*$  is almost always exact (84 out of 96 times), and is never off by more than 1.

For case iii) we let

$$n^* = \lambda/\mu + y^*(a, b, d)\sqrt{\lambda/\mu},$$

where

$$y^*(a, b, d) = \arg \min_{y>0} \left\{ y + P(y) \left[ \frac{a}{y} + b e^{-dy} \right] \right\}.$$

Here we set  $a = 2$ ,  $b = 2.5$ , and  $d = 0.1$ , and tried all integer values between 5 and 100 for  $\lambda$ . Here again, rounding off  $n^*$  is almost always exact (80 out of 96 cases), and is never off by more than 1.

For Example 6.4 we restricted our attention to  $b_\lambda = 0$ . We initially set

$$a_\lambda = a\lambda^{-1/2}. \tag{40}$$

Defining

$$y_\lambda^*(a) = \arg \min_{y>0} \left\{ y + \frac{a}{y\sqrt{\lambda}} \right\},$$

we can solve explicitly to obtain  $y_\lambda^*(a) = \sqrt{a}\lambda^{-1/4}$ . We thus let

$$n^* = \lambda/\mu + \sqrt{a/\mu}\lambda^{1/4}.$$

For the numerical test of this approximation we set  $a = 1$  and tried all integer multiples of 10 between 10 and 200 for  $\lambda$ , using (40) to determine  $a_\lambda$ . Rounding off  $n^*$  to the nearest integer is almost always exact (95 out of 96 times), and is never off by more than 1.

For Example 7.4 we restricted our attention to  $b_\lambda = 0$ . We took

$$a_\lambda = a\sqrt{\lambda}. \tag{41}$$

Let

$$n^* = \lambda/\mu + y_\lambda^*(a)\sqrt{\lambda/\mu},$$

where

$$y_\lambda^*(a) = \arg \min_{y>0} \left\{ y + \frac{aQ(y)}{\sqrt{\lambda}}y \right\}, \tag{42}$$

and  $Q(y)$  is given by (12).

For the numerical test of this approximation we set  $a = 1$  and tried all integer multiples of 10 between 10 and 200 for  $\lambda$ , using (41) to determine  $a_\lambda$ . Rounding off  $n^*$  to the nearest integer is almost always exact (92 out of the 96 times), and is never off by more than 1.

The above tests were run under favorable conditions: The asymptotic scaling of the parameters was known, and the approximation used was that associated with the regime corresponding to the parameter scaling. On the other hand, the results of the test are astoundingly good: Rounding  $n^*$  to the nearest integer is almost always exact, and is never off by more than 1. Note that these tests include values of  $\lambda$  that do not appear to be very large.

It is clear that more numerical testing is in order. There are two aims to this testing: 1) Find parameter values that ‘break’ the approximations, and 2) Determine if any of the asymptotic approximations is robust enough to work outside of its regime and/or determine rules of thumb for when each approximation should be used. (Indeed, this last point is central for obtaining a practically useful approximation.)

The additional testing takes the form of ‘wrong-regime’ testing. The first ‘wrong-regime’ test we conducted involved scaling the parameters as in the efficiency-driven regime with  $c_\lambda = c = 1$ ,  $\mu = 1$ ,  $b_\lambda = 0$ , and  $a_\lambda = a\lambda^{-1/2}$ , and using the approximation from the rationalized regime. Thus our approximation rounded off

$$n^* = \lambda/\mu + y^*(a\lambda^{-1/2})\sqrt{\lambda/\mu}$$

to obtain the staffing level. For our numerical results we set  $a = 1$  and tried all integer multiples of 10 between 10 and 200 for  $\lambda$ . The approximation was exact in all but one case ( $\lambda = 160$ ), where it was off by 1.

Using the same parameters as in the preceding example we used the approximation from the quality-driven regime (as if  $a_\lambda = a\sqrt{\lambda}$ ). Thus our approximation rounded off

$$n^* = \lambda/\mu + y_\lambda^*\left(\frac{a}{\lambda}\right)\sqrt{\lambda/\mu},$$

where  $y_\lambda^*$  is given by (42), to obtain the staffing level. For our numerical results we set  $a = 1$  and tried all integer multiples of 10 between 10 and 200 for  $\lambda$ . The approximation was good, but not as good as the rationalized regime: it was off by 1 in 14 cases and off by 2 in 4 cases. We next scaled the parameters as in the quality-driven regime with  $c_\lambda = c = 1$ ,  $\mu = 1$ ,  $b_\lambda = 0$ , and  $a_\lambda = a\sqrt{\lambda}$ , and used the approximation from the rationalized regime, rounding off

$$n^* = \lambda/\mu + y^*(a\sqrt{\lambda})\sqrt{\lambda/\mu},$$

where  $y^*(\cdot)$  is given by (39), to obtain the staffing level. For our numerical results we set  $a = 1$  and tried all integer multiples of 10 between 10 and 200 for  $\lambda$ . The approximation was off by 1 in 6 cases and exact in the others.

Using the same parameters as in the preceding example we used the approximation from the efficiency-driven regime (as if  $a_\lambda = a\lambda^{-1/2}$ ). Thus our approximation rounded off

$$n^* = \lambda/\mu + \sqrt{a/\mu}\lambda^{3/4}$$

to obtain the staffing level. For our numerical results we set  $a = 1$  and tried all integer multiples of 10 between 10 and 200 for  $\lambda$ . This approximation did not perform well, leading to overstaffing of 10–15%.

We next scaled the parameters as in the rationalized regime, with  $c_\lambda = c = 1$ ,  $\mu = 1$ ,  $b_\lambda = 0$ , and  $a_\lambda = a$ , using the approximation from the quality-driven regime (as if  $a_\lambda = a\sqrt{\lambda}$ ). Thus our approximation rounded off

$$n^* = \lambda/\mu + y_\lambda^*\left(\frac{a}{\sqrt{\lambda}}\right)\sqrt{\lambda/\mu},$$

where  $y_\lambda^*(\cdot)$  is given by (42), to obtain the staffing level. For our numerical results we set  $a = 1$  and tried all integer multiples of 10 between 10 and 200 for  $\lambda$ . The approximation here was exact in 8 cases and off by 1 in 12 cases.

Using the same parameters as in the preceding example we used the approximation from the efficiency-driven regime (as if  $a_\lambda = a\lambda^{-1/2}$ ). Thus our approximation rounded off

$$n^* = \lambda/\mu + \sqrt{a/\mu}\sqrt{\lambda}$$

to obtain the staffing level. For our numerical results we set  $a = 1$  and tried all integer multiples of 10 between 10 and 200 for  $\lambda$ . The approximation is not as good as that of the rationalized or quality-driven regime: 3 cases were exact, 8 were off by 1, and 9 were off by 2.

In the tests so far that involve scaling parameters for the efficiency-driven and quality-driven regimes we have used  $a_\lambda = a\lambda^{-1/2}$  and  $a_\lambda = a\sqrt{\lambda}$ , respectively. These regimes hold more generally with  $a_\lambda = a\lambda^{-\alpha}$  and  $a_\lambda = a\lambda^\alpha$  respectively, for  $\alpha > 0$ . In addition to  $\alpha = 1/2$ , we also tried values of  $\alpha = 1/4$  and 1. The runs for  $\alpha = 1$  involved multiples of 50 from 50 to 1000 for  $\lambda$ . (The approximations used the same value of  $\alpha$  as used to scale the actual parameters.) The results can be summarized as follows. The rationalized approximation was excellent: it was mostly exact, and when it was wrong it was never off by more than 1. The quality-driven approximation was good, but not as good as the rationalized approximation (even in the quality-driven regime!). The efficiency-driven solution was the worst of the three, and substantially overstaffed in the quality-driven regime. In the efficiency-driven regime with  $\alpha = 1$ , the efficiency-driven solution provided the infeasible solution of  $\lambda$  as the staffing level. Of course this can be corrected by simply requiring that the staffing level must be strictly greater than  $\lambda$ .

## 10 Future research

There are a few directions of research that suggest themselves. First, we would like to explain theoretically the extreme accuracy of our approximations. This cannot be anticipated from the corresponding asymptotic approximations.

The call center environment enjoys features that are not captured by the M/M/N (Erlang-C) model. Examples include non-exponential service times (M/G/N) and abandonment. The goal is to incorporate such features into our framework. Regarding abandonment, a first step was taken in [5]: in the M/M/N queue with exponential abandonment, the square-root safety staffing rule was structurally identified, but  $y^*$  was not calculated. (It is important to note that abandonment renders the queue always stable, hence  $y^*$  can also take negative values.) As for the M/G/N queue, it is not amenable to exact analysis, and letting  $N \rightarrow \infty$  does not make things easier. Indeed, even the accuracy of the standard multi-server numerical approximations is unclear as  $N$  becomes large, which is the relevant regime for call centers. A key challenge is the calculation of the Halfin-Whitt delay function for the M/G/N queue. To this end, one could first attempt the M/PH/N queue, following [14].

Also, in the present paper we assumed that the offered traffic parameters are exactly known, and focused on determining the amount of safety staffing needed to deal with stochastic variability only. In practice, the offered traffic forecasts are typically not completely accurate, and additional staffing may be required in view of the inherent uncertainty in the parameters. However, we feel that the analysis for known parameters is an essential step towards situations where the offered traffic estimates may contain inaccuracies.

Lastly, one could perhaps use duality theory from Mathematical Programming to relate the optimization approach to constraint satisfaction. This could perhaps add insight on the optimality criteria for constraint satisfaction, which should be sharpened.

## References

- [1] Andrews, B. and Parsons, H. (1993). Establishing telephone-agent staffing levels through economic optimizations. *Interfaces* **23:2**, 14–20.
- [2] De Bruijn, N.G. (1981). *Asymptotic Methods in Analysis*. Dover.
- [3] Cooper, R.B. (1981). *Introduction to Queueing Theory*. 2nd Edition, North Holland.
- [4] Call Center Statistics (WEB site), September 2000.  
<http://www.callcenternews.com/resources/statistics.shtml>.
- [5] Garnet, O., Mandelbaum, A., and Reiman, M.I. (1999). Designing a telephone call center with impatient customers. Submitted for publication in MSOM.

- [6] Grassmann, W.K. (1986). Is the fact that the emperor wears no clothes a subject worthy of publication? *Interfaces* **16:2**, 43–51.
- [7] Grassmann, W.K. (1988). Finding the right number of servers in real-world queueing systems. *Interfaces* **18:2**, 94–104.
- [8] Halfin, S., Whitt, W. (1981). Heavy-traffic limits for queues with many exponential servers. *Oper. Res.* **29**, 567–588.
- [9] Jagers, A.A., Van Doorn, E.A. (1991). Convexity of functions which are generalizations of the Erlang loss function and the Erlang delay function. *SIAM Review* **33**, 281–282.
- [10] Jagers, A.A., Van Doorn, E.A. (1986). On the continued Erlang loss function. *Oper. Res. Lett.* **5**, 43–46.
- [11] Jennings, O.B., Mandelbaum, A., Massey, W.A., and Whitt, W. (1996). Service staffing to meet time-varying demand. *Mgmt. Sc.* **42**, 1383–1394.
- [12] Kolesar, P.J., Green, L.V. (1998). Insights on service system design from a normal approximation to Erlang’s delay formula. *Prod. Oper. Mgmt.* **7**, 282–293.
- [13] Mandelbaum, A., Massey, W. A., Reiman, M. I., Rider, B., and Stolyar, A. (2000). Queue length and waiting times for multiserver queues with abandonment and retries. Selected Proceedings of the Fifth INFORMS Telecommunications Conference.
- [14] Puhalskii, A.A., Reiman, M.I. (2000). The multiclass GI/PH/N queue in the Halfin-Whitt regime. *Adv. Appl. Prob.* **32**, 564–595.
- [15] Whitt, W. (1992). Understanding the efficiency of multi-server service systems. *Mgmt. Sc.* **38**, 708–723.

## A Proof of Lemma 4.2

### Lemma 4.2

For any function  $x_\lambda$  with  $\lim_{\lambda \rightarrow \infty} x_\lambda = \infty$ ,

$$\pi_\lambda(x_\lambda) \stackrel{\infty}{\sim} Q_\lambda(x_\lambda).$$

If also  $x_\lambda \stackrel{\sup}{\leq} \lambda^{1/6}$ , then

$$\pi_\lambda(x_\lambda) \stackrel{\infty}{\sim} Q(x_\lambda).$$

If specifically  $x_\lambda = \kappa \sqrt{\lambda/\mu}$  for some constant  $\kappa > 0$ , then

$$\pi_\lambda(x_\lambda) \stackrel{\infty}{\sim} \frac{1}{\kappa \sqrt{2\pi\lambda/\mu(1+\kappa)}} \left( \frac{e^\kappa}{(1+\kappa)^{1+\kappa}} \right)^{\lambda/\mu}.$$

### Proof

The first statement follows after some manipulations from the proof of Proposition 1 of Halfin & Whitt [8].

If  $x_\lambda = \kappa\sqrt{\lambda/\mu}$  for some constant  $\kappa > 0$ , then  $N_\lambda(x_\lambda) = (1 + \kappa)\lambda/\mu$ ,  $r_\lambda(x_\lambda) = 1/(1 + \kappa)$ , and  $1 - r_\lambda(x_\lambda) = \kappa/(1 + \kappa)$ , so that

$$Q_\lambda(x_\lambda) = \frac{1}{\kappa\sqrt{2\pi\lambda/\mu(1 + \kappa)}} \left( \frac{e^\kappa}{(1 + \kappa)^{1+\kappa}} \right)^{\lambda/\mu}.$$

We now prove the third statement. Using the Taylor series expansion

$$\log u = \log(1 - (1 - u)) = - \sum_{m=1}^{\infty} \frac{(1 - u)^m}{m} = -(1 - u) - \sum_{m=2}^{\infty} \frac{(1 - u)^m}{m},$$

we obtain

$$\begin{aligned} \frac{\exp\{N_\lambda(x)[1 - r_\lambda(x) + \log r_\lambda(x)]\}}{\sqrt{2\pi N_\lambda(x)}(1 - r_\lambda(x))} &= \frac{\exp\left\{-N_\lambda(x) \sum_{m=2}^{\infty} \frac{(1 - r_\lambda(x))^m}{m}\right\}}{\sqrt{2\pi N_\lambda(x)}(1 - r_\lambda(x))} = \\ &= \frac{\exp\left\{-N_\lambda(x)(1 - r_\lambda(x))^2/2 - N_\lambda(x) \sum_{m=3}^{\infty} \frac{(1 - r_\lambda(x))^m}{m}\right\}}{\sqrt{2\pi N_\lambda(x)}(1 - r_\lambda(x))} = \\ Q(x) \frac{x}{\sqrt{N_\lambda(x)}(1 - r_\lambda(x))} \exp\left\{[x^2 - N_\lambda(x)(1 - r_\lambda(x))^2]/2\right\} \exp\left\{-N_\lambda(x) \sum_{m=3}^{\infty} \frac{(1 - r_\lambda(x))^m}{m}\right\}. \end{aligned}$$

Thus it remains to be shown that

$$\frac{x_\lambda}{\sqrt{N_\lambda(x_\lambda)}(1 - r_\lambda(x_\lambda))} \exp\{[x_\lambda^2 - N_\lambda(x_\lambda)(1 - r_\lambda(x_\lambda))^2]/2\} \exp\{-N_\lambda(x_\lambda) \sum_{m=3}^{\infty} \frac{(1 - r_\lambda(x_\lambda))^m}{m}\} \approx 1$$

if  $x_\lambda \ll \lambda^{1/6}$ .

Note that

$$\sqrt{N_\lambda(x_\lambda)}(1 - r_\lambda(x_\lambda)) = \frac{x_\lambda \sqrt{\lambda/\mu}}{\sqrt{\lambda/\mu + x_\lambda \sqrt{\lambda/\mu}}} \approx x_\lambda,$$

and

$$N_\lambda(x_\lambda)(1 - r_\lambda(x_\lambda))^2 = \frac{x_\lambda^2 \lambda/\mu}{\lambda/\mu + x_\lambda \sqrt{\lambda/\mu}},$$

so that

$$0 \leq x_\lambda^2 - N_\lambda(x_\lambda)(1 - r_\lambda(x_\lambda))^2 = \frac{x_\lambda^3 \sqrt{\lambda/\mu}}{\lambda/\mu + x_\lambda \sqrt{\lambda/\mu}} \leq \frac{x_\lambda^3}{\sqrt{\lambda/\mu}} \ll 1.$$

Finally,

$$\begin{aligned} 0 \leq N_\lambda(x_\lambda) \sum_{m=3}^{\infty} \frac{(1 - r_\lambda(x_\lambda))^m}{m} &\leq N_\lambda(x_\lambda) \sum_{m=3}^{\infty} (1 - r_\lambda(x_\lambda))^m = \\ \frac{N_\lambda(x_\lambda)(1 - r_\lambda(x_\lambda))^3}{r_\lambda(x_\lambda)} &\leq \frac{x_\lambda^3}{\sqrt{\lambda/\mu}} \ll 1, \end{aligned}$$

which completes the proof. □



## B Properties of $P(\cdot)$

### Lemma B.1

The function  $P(\cdot)$  is strictly convex decreasing.

### Proof

The function  $P(\cdot)$  may be written

$$P(x) = \frac{1}{1 + U(x)},$$

with

$$U(x) := x e^{x^2/2} V(x),$$

and

$$V(x) := \int_{-\infty}^x e^{-y^2/2} dy.$$

Differentiating,

$$P'(x) = \frac{-U'(x)}{(1 + U(x))^2},$$

and

$$P''(x) = \frac{2U'(x)^2 - U''(x)(1 + U(x))}{(1 + U(x))^3}.$$

Observing that  $V'(x) = e^{-x^2/2}$  and  $V''(x) = -x e^{-x^2/2}$ ,

$$U'(x) = x + (x^2 + 1) e^{x^2/2} V(x),$$

and

$$U''(x) = x^2 + 2 + (x^3 + 3x) e^{x^2/2} V(x).$$

Thus,  $P(\cdot)$  is decreasing since  $U'(x) > 0$  for any  $x > 0$ .

Also,  $P(\cdot)$  is strictly convex because for any  $x > 0$ ,

$$\begin{aligned} 2U'(x)^2 - U''(x)(1 + U(x)) &= \\ 2 \left[ x + (x^2 + 1) e^{x^2/2} V(x) \right]^2 - \left[ x^2 + 2 + (x^3 + 3x) e^{x^2/2} V(x) \right] \left[ 1 + x e^{x^2/2} V(x) \right] &= \\ x^2 - 2 + \left[ 2x^3 - x \right] e^{x^2/2} V(x) + \left[ x^4 + x^2 + 2 \right] \left[ e^{x^2/2} V(x) \right]^2 &> \\ x^2 - 2 + \left[ x^4 + 2x^3 + x^2 - x + 2e^{x^2/2} V(x) \right] e^{x^2/2} V(x) &> \end{aligned}$$

$$\begin{aligned}
& x^2 - 2 + \left[ x^4 + 2x^3 + x^2 - x + 2\sqrt{\pi/2} \right] e^{x^2/2} V(x) > \\
& x^2 + \left[ x^4 + 2x^3 + x^2 - x + 2(\sqrt{\pi/2} - 1) \right] e^{x^2/2} V(x) > \\
& \left[ x^2 - x + 2(\sqrt{\pi/2} - 1) \right] e^{x^2/2} V(x) = \\
& \left[ (x - 1/2)^2 + 2(\sqrt{\pi/2} - 9/8) \right] e^{x^2/2} V(x) > 0.
\end{aligned}$$

□

## C Properties of $G_\lambda(\cdot)$

### Lemma C.1

*The function  $G_\lambda(\cdot)$  is strictly convex decreasing.*

### Proof

It suffices to show that  $K_\lambda(\cdot)$  is strictly convex decreasing with

$$K_\lambda(\omega) := \omega \int_0^\infty D_\lambda(t) e^{-\omega t} dt.$$

Differentiating,

$$K'_\lambda(\omega) = \int_0^\infty [1 - \omega t] e^{-\omega t} D_\lambda(t) dt,$$

and

$$K''_\lambda(\omega) = \int_0^\infty [\omega t - 2] t e^{-\omega t} D_\lambda(t) dt,$$

for all  $\omega > 0$ .

Since  $D_\lambda(\cdot)$  is strictly increasing, we have

$$[1 - \omega t] D_\lambda(t) \leq [1 - \omega t] D_\lambda(\omega)$$

for all  $\omega > 0$ ,  $t > 0$ , with *strict* inequality for  $t \neq 1/\omega$ , so that

$$\int_0^\infty [1 - \omega t] e^{-\omega t} D_\lambda(t) dt < D_\lambda(1/\omega) \int_0^\infty [1 - \omega t] e^{-\omega t} dt = 0,$$

and

$$\int_0^\infty [\omega t - 2] t e^{-\omega t} D_\lambda(t) dt > D_\lambda(2/\omega) \int_0^\infty [\omega t - 2] t e^{-\omega t} dt = 0,$$

for all  $\omega > 0$ .

□