

## Accepted Manuscript

Rare-event analysis of mixed Poisson random variables, and applications in staffing

Mariska Heemskerk, Julia Kuhn, Michel Mandjes

PII: S0166-5316(16)30215-2

DOI: <http://dx.doi.org/10.1016/j.peva.2017.03.006>

Reference: PEVA 1910

To appear in: *Performance Evaluation*



Please cite this article as: M. Heemskerk, et al., Rare-event analysis of mixed Poisson random variables, and applications in staffing, *Performance Evaluation* (2017), <http://dx.doi.org/10.1016/j.peva.2017.03.006>

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

## RARE-EVENT ANALYSIS OF MIXED POISSON RANDOM VARIABLES, AND APPLICATIONS IN STAFFING

MARISKA HEEMSKERK, JULIA KUHN, MICHEL MANDJES

ABSTRACT. A common assumption when modeling queuing systems is that arrivals behave like a Poisson process with constant parameter. In practice, however, call arrivals are often observed to be significantly overdispersed. This motivates that in this paper we consider a *mixed* Poisson arrival process with arrival rates that are resampled every  $N^{-\alpha}$  time units, where  $\alpha > 0$  and  $N$  a scaling parameter.

In the first part of the paper we analyse the asymptotic tail distribution of this doubly stochastic arrival process. That is, for large  $N$  and i.i.d. arrival rates  $X_1, \dots, X_N$ , we focus on the evaluation of the probability that the scaled number of arrivals exceeds  $Na$ ,

$$P_N(a) := \mathbb{P}(\text{Pois}(N\bar{X}_{N^\alpha}) \geq Na), \quad \text{with } \bar{X}_N := \frac{1}{N} \sum_{i=1}^N X_i.$$

The logarithmic asymptotics of  $P_N(a)$  are easily obtained from previous results; we find constants  $r_P$  and  $\gamma$  such that  $N^{-\gamma} \log P_N(a) \rightarrow -r_P$  as  $N \rightarrow \infty$ . Relying on elementary techniques, we then derive the exact asymptotics of  $P_N(a)$ : For  $\alpha < \frac{1}{3}$  and  $\alpha > 3$  we identify (in closed-form) a function  $\tilde{P}_N(a)$  such that  $P_N(a)/\tilde{P}_N(a)$  tends to 1 as  $N \rightarrow \infty$ . For  $\alpha \in [\frac{1}{3}, \frac{1}{2})$  and  $\alpha \in [2, 3)$  we find a partial solution in terms of an asymptotic lower bound. For the special case that the  $X_i$ s are gamma distributed, we establish the exact asymptotics across all  $\alpha > 0$ . In addition, we set up an asymptotically efficient importance sampling procedure that produces reliable estimates at low computational cost.

The second part of the paper considers an infinite-server queue assumed to be fed by such a mixed Poisson arrival process. Applying a scaling similar to the one in the definition of  $P_N(a)$ , we focus on the asymptotics of the probability that the number of clients in the system exceeds  $Na$ . The resulting approximations can be useful in the context of staffing. Our numerical experiments show that, astoundingly, the required staffing level can actually *decrease* when service times are more variable.

## 1. INTRODUCTION

In communications engineering it is increasingly accepted that traditional Poisson processes do not succeed in capturing the variability that is typically observed in real call arrival processes [12, 19]. This led to the idea to instead use Cox processes [5] to model arrivals, i.e., Poisson processes in which the arrival rate follows some (non-negative) stochastic process. Perhaps the simplest choice, advocated in [9], is to *resample* the arrival rate (in an i.i.d. manner) every  $\Delta$  units of time; during the resulting time intervals the arrival rate is assumed constant. We denote these i.i.d. arrival rates by  $(X_i)_{i \in \mathbb{N}}$ . This paper studies two settings in which such an overdispersed arrival process is featured.

1. *Number of arrivals.* We start by studying the tail asymptotics of the total number of arrivals in a time interval of given length. We do so in a scaling regime that was proposed in [9], in which the arrival rates and sampling frequency are jointly inflated as follows. In the first place, it is natural to assume that arrival rates are large, as these represent the contributions of many potential clients; this can be achieved by letting these arrival rates be  $NX_1, NX_2, \dots$  for i.i.d.  $(X_i)_{i \in \mathbb{N}}$  and some large  $N$ . In addition, the sampling frequency is set to  $N^\alpha$  (assumed to be integer) and hence the size of each time slot is assumed to be  $\Delta = N^{-\alpha}$ . Evidently, the larger  $\alpha$ , the more frequently the arrival rate is resampled.

The focus is on the probabilities  $P_N(a)$  and  $p_N(a)$ , where

$$P_N(a) := \mathbb{P}(\text{Pois}(N\bar{X}_{N^\alpha}) \geq Na), \quad \text{with} \quad \bar{X}_N := \frac{1}{N} \sum_{i=1}^N X_i,$$

and  $p_N(a)$  denotes the corresponding probability that the mixed Poisson random variable equals  $Na$  (assumed to be integer). We consider the situation that  $a$  is larger than  $\nu := \mathbb{E}X_i$ , which entails that the event under consideration is rare and that we are in the framework of large deviations theory.

We would like to stress the important role that is played by the time-scale parameter  $\alpha > 0$ . One could imagine that in a rapidly changing environment, the inherent overdispersion of the arrival process hardly plays a role, whereas in a slowly changing random environment, overdispersion is expected to be more dominant. Hence the parameter  $\alpha$  can be tweaked in order to match any real-world scenario in that sense. That is, if  $\alpha$  is large, since the arrival rate is resampled relatively frequently, it is anticipated that the mixed Poisson random variable behaves Poissonian with parameter  $N\nu$ . If on the contrary  $\alpha$  is small, one would expect that detailed characteristics of the distribution of the  $X_i$  matter. For  $\alpha = 1$  both effects play a role. This intuition underlies nearly all results presented in this paper.

2. *Number of customers in an infinite-server queue.* In the second part of this paper we focus on a cornerstone model in the design and performance evaluation of communication networks: the *infinite-server queue*. This model can be used to produce approximations for many-server systems. In our paper, the arrival process is the overdispersed process we introduced above, and the service times are i.i.d. samples from a (non-negative) distribution with distribution function  $F(\cdot)$ . The number of clients in this infinite-server queue, under the arrival process described above, is studied in [9]. As it turns out, one can prove the (conceivable) property that the number of clients in the system at time  $t$  (which we, for simplicity, assume to be a multiple of  $\Delta$ ), has a *mixed Poisson* distribution, i.e., a Poisson distribution with random parameter. This parameter is given by

$$\sum_{i=1}^{t/\Delta} X_i \Delta f_i(t, \Delta),$$

where  $f_i(t, \Delta)$  denotes the probability that a call arriving at a uniformly distributed epoch in the interval  $[(i-1)\Delta, i\Delta)$  is still in the system at time  $t$ . Evidently, for small  $\Delta$  this probability essentially behaves as  $\bar{F}(t - i\Delta)$ , with  $\bar{F}(\cdot) := 1 - F(\cdot)$  denoting the complementary distribution function.

We renormalize time such that  $t \equiv 1$  (which can be done without loss of generality), and again impose the scaling along the lines of [9]: the arrival rates are  $NX_i$  and the interval width  $N^{-\alpha}$ . Then the number of clients in the system is Poisson with random parameter

$$\sum_{i=1}^{N^\alpha} (NX_i) N^{-\alpha} f_i(1, N^{-\alpha}) = N^{1-\alpha} \sum_{i=1}^{N^\alpha} X_i \omega_i(N^\alpha), \quad (1)$$

where  $\omega_i(N) := f_i(1, N^{-1}) \approx \bar{F}(1 - i/N)$ . A clearly relevant object of study concerns the probability that the number of clients in the system exceeds some threshold  $Na$ :

$$Q_N(a) := \mathbb{P} \left( \text{Pois} \left( N^{1-\alpha} \sum_{i=1}^{N^\alpha} X_i \omega_i(N^\alpha) \right) \geq Na \right); \quad (2)$$

$q_N(a)$  denotes the corresponding probability that the mixed Poisson random variable equals  $Na$ . To ensure that the event under consideration is rare,  $a$  is assumed to be larger than

$$\frac{\nu}{N^\alpha} \sum_{i=1}^{N^\alpha} \omega_i(N^\alpha) \approx \frac{\nu}{N^\alpha} \sum_{i=1}^{N^\alpha} \bar{F}(1 - i/N^\alpha) \approx \nu \int_0^1 \bar{F}(x) dx.$$

A related question of practical interest concerns *staffing*: how many servers should be allocated to ensure a given service level for customers or jobs arriving according to a mixed Poisson process in a random environment? Approximating the many-server model by its infinite-server counterpart, we approach this classical problem in queueing theory as an asymptotic dimensioning problem: we want to find the smallest  $a$  such that  $Q_N(a)$  (or  $q_N(a)$ ) is below some desired (typically small)  $\varepsilon$  as  $N$  tends to infinity (cf. [4]). The resulting procedure has applications in the context of call centers, cloud computing or in the design of data centers [14, 17]. Related literature on (dynamic) staffing procedures in such settings is, e.g., [10, 18, 19]; see also the recent review [6] and the references therein. Previous work that addresses overdispersion in the arrival process includes [2, 8, 11]. Our approach in this paper is different to earlier work in that it uses exact asymptotics to approximate the objective function that we want to minimize in the staffing problem. As we focus on the large-deviations setting, the technique we develop is specifically useful in the regime in which the performance requirements are strict (i.e., the probability of service degradation should be kept low).

We now comment on the type of results we establish in this paper. As is common in the literature, we first consider *logarithmic asymptotics*, that is, we identify a constant  $r_Q > 0$  (that depends on  $a$ ) such that, for  $\gamma := \min\{\alpha, 1\}$ ,

$$\lim_{N \rightarrow \infty} \frac{1}{N^\gamma} \log Q_N(a) = -r_Q. \quad (3)$$

This is easily done by using the techniques from [9].

These logarithmic asymptotics provide useful insight into the decay of the probabilities of interest, but it should be noted that they are inherently imprecise. More specifically, they suggest that one could use ‘naive’ approximations of the form

$$P_N(a) \approx e^{-r_P N^\gamma}, \quad Q_N(a) \approx e^{-r_Q N^\gamma}$$

for  $N$  large. It is important to notice, however, that (3) only entails that  $P_N(a) = \xi(N) \exp(-r_P N^\gamma)$ , with  $\xi(\cdot)$  being subexponential in the sense that

$$\lim_{N \rightarrow \infty} \frac{1}{N^\gamma} \log \xi(N) = 0.$$

In other words, it does not rule out that, for instance,  $\xi(N) = 10^{10}$ , or  $N^M$  for some given  $M$ , or even  $\exp(N^{0.99\gamma})$ . This motivates the interest in *exact asymptotics*. Here, the objective is to identify a function  $\tilde{P}_N(a)$  such that  $\tilde{P}_N(a)/P_N(a) \rightarrow 1$  as  $N \rightarrow \infty$  (which we denote throughout the paper by  $P_N(a) \sim \tilde{P}_N(a)$ ), leading to the evident approximation  $P_N(a) \approx \tilde{P}_N(a)$ . Along the same lines we would like to find the exact asymptotics  $\tilde{Q}_N(a)$  for the probability  $Q_N(a)$ .

The contributions and organization of our paper are as follows. In Sections 2, 3 and 4 we focus on the evaluation of the probabilities  $P_N(a)$  and  $p_N(a)$ . After having introduced the notation, in Section 2 we first briefly present the logarithmic asymptotics. We then use elementary techniques to derive the exact asymptotics, however, as it turns out, these only apply when the time scales of the arrival process and the resampling are sufficiently separated: we address the cases  $\alpha < \frac{1}{3}$  and  $\alpha > 3$  (with a partial solution for  $\alpha \in [\frac{1}{3}, \frac{1}{2})$  and  $\alpha \in [2, 3)$  in terms of an asymptotic lower bound). In Section 3 it becomes clear why such elementary techniques do not work across all values of  $\alpha$ : for the important special case of the  $X_i$  corresponding to i.i.d. gamma distributed random variables [11] we find the exact asymptotics for all  $\alpha > 0$ , and in the range  $(\frac{1}{2}, 2) \setminus \{1\}$  these turn out to have a rather intricate shape.

Section 4 focuses on rare-event simulation as a means to find an accurate approximation at relatively low computational cost: we propose an importance-sampling based technique, which we prove to be asymptotically efficient.

In Section 5 we shift our attention to the probabilities  $Q_N(a)$  and  $q_N(a)$ . Again, logarithmic asymptotics can be found, and in addition we manage to identify the exact asymptotics for the case  $\alpha = 1$ . By a series of numerical examples it is illustrated how the resulting approximation can be used for staffing purposes. We performed extensive experiments, and make the striking observation that increasing the variability of the service times (e.g. Pareto service times rather than exponential ones) often leads to less conservative staffing rules.

## 2. ASYMPTOTICS OF $P_N(a)$

We start by introducing the framework that we consider throughout the paper. In our setup we let  $(X_i)_{i \in \mathbb{N}}$  be a sequence of i.i.d. random variables distributed as a generic random variable  $X$ , where  $\nu := \mathbb{E}X_i$ . Assume that the moment-generating function of  $X$ , denoted by  $M_X(\vartheta) := \mathbb{E}[e^{\vartheta X}]$ , is finite in an open set containing the origin. The *Fenchel-Legendre transform* (or convex conjugate) of the cumulant-generating function  $\Lambda_X(\vartheta) := \log M_X(\vartheta)$  is defined as

$$I_X(a) := \sup_{\vartheta \in \mathbb{R}} \{\vartheta a - \Lambda_X(\vartheta)\}. \quad (4)$$

We assume that the optimizing  $\vartheta$  in (4) indeed exists, and we denote it by  $\vartheta_X^*$  (thus suppressing that  $\vartheta_X^*$  actually depends on  $a$ ). Under these conditions, it is known that the sample mean  $\bar{X}_N := N^{-1} \sum_{i=1}^N X_i$  satisfies a large deviations principle with *rate function*  $I_X(\cdot)$  (see, e.g., [7]). Furthermore, a result by Bahadur and Rao [1] states that we have the following exact asymptotics for  $\bar{X}_N$ : when  $a > \nu$ ,

$$\lim_{N \rightarrow \infty} \mathbb{P}(\bar{X}_N \geq a) e^{N I_X(a)} \sqrt{N} = C_X(a). \quad (5)$$

We assume that  $X$  is non-lattice, in which case  $C_X(\cdot)$  takes the form

$$C_X(a) = \frac{1}{\vartheta_X^* \sqrt{2\pi \Lambda_X''(\vartheta_X^*)}}, \quad (6)$$

where  $\Lambda_X''(\vartheta_X^*)$  denotes the second derivative of  $\Lambda_X(\vartheta)$  evaluated at  $\vartheta_X^*$ ; if  $X$  is lattice, the constant  $C_X(a)$  should be defined slightly differently [7, Thm. 3.7.4]. There is also a local limit version of (5): with  $\xi_N(\cdot)$  the density of  $\sum_{i=1}^N X_i$ , from [15],

$$\lim_{N \rightarrow \infty} \xi_N(Na) e^{NI_X(a)} \sqrt{N} = C_X(a) I_X'(a). \quad (7)$$

In our analysis the tail asymptotics of Poisson random variables play a crucial role. We note that the Bahadur-Rao asymptotics entail that for the probabilities

$$\psi_N(a|x) := \mathbb{P}(\text{Pois}(Nx) \geq Na), \quad (8)$$

it holds that

$$\lim_{N \rightarrow \infty} \psi_N(a|x) e^{NI(a|x)} \sqrt{N} = C(a|x), \quad (9)$$

for  $a > x$ . Here,  $I(\cdot|x)$  is the rate function associated with a Poisson random variable with parameter  $x$ , that is,  $I(\cdot|x)$  is the Fenchel-Legendre transform of the cumulant-generating function  $\Lambda(\vartheta) = x(e^\vartheta - 1)$  of the Poisson random variable. Inserting the optimizer  $\vartheta^* = \log(a/x)$ , this yields  $I(a|x) = a \log(a/x) - a + x$ . Bearing in mind that the Poisson variable is lattice, it turns out that the function  $C(a|x)$  takes the form (cf. (6))

$$C(a|x) := \frac{1}{1 - \exp(\vartheta^*)} \frac{1}{\sqrt{2\pi\Lambda''(\vartheta^*)}} = \frac{1}{1 - a/x} \frac{1}{\sqrt{2\pi a}}.$$

Let us first present the logarithmic asymptotics of  $P_N(a)$  (the same logarithmic asymptotics hold for  $p_N(a)$ ). Here we merely state the results as the proof is exactly as in [9, Section 4.1]. We distinguish between the cases  $\alpha > 1$  and  $\alpha < 1$ ; the former case we refer to as the *fast* regime as the  $X_i$ 's are sampled relatively frequently, whereas the latter case is the *slow* regime. For completeness, the logarithmic asymptotics for the intermediate case  $\alpha = 1$ , though standard, are included as well.

- In the fast regime  $N^\alpha$  is substantially larger than  $N$ , and hence the rare event will be essentially due to  $\bar{X}_{N^\alpha}$  being close to  $\nu$ , and the Poisson random variable with parameter (roughly)  $N\nu$  exceeding  $Na$ . Accordingly, following the argumentation in [9], one obtains

$$\lim_{N \rightarrow \infty} \frac{1}{N} \log P_N(a) = -I(a|\nu).$$

This result entails that  $P_N(a)$  decays essentially exponentially.

- In the slow regime, assuming the support of  $X_i$  is unbounded, the rare event will be a consequence of the joint effect of (i)  $\bar{X}_{N^\alpha}$  being close to  $a$ , and (ii) the Poisson variable with parameter (roughly)  $Na$  attaining a typical value; the first event is rare, but the second is not. In this regime, we thus have

$$\lim_{N \rightarrow \infty} \frac{1}{N^\alpha} \log P_N(a) = -I_X(a);$$

observe that this corresponds to subexponential decay.

- For  $\alpha = 1$ , the random variable  $\text{Pois}(N\bar{X}_{N^\alpha})$  can be written as the sum of  $N$  i.i.d. contributions, each of them distributed as  $Z := \text{Pois}(X)$ . Noting that

$$\log \mathbb{E} \exp(\vartheta Z) = \Lambda_X(e^\vartheta - 1),$$

a straightforward application of Cramér's theorem [7] yields that the decay is exponential:

$$\lim_{N \rightarrow \infty} \frac{1}{N} \log P_N(a) = -\sup_{\vartheta} \left( \vartheta a - \Lambda_X(e^\vartheta - 1) \right) =: I_Z(a). \quad (10)$$

In the remainder of this section we show that for a range of values of  $\alpha$  the exact asymptotics of  $P_N(a)$  and  $p_N(a)$  can be found relying on elementary probabilistic techniques. We focus on the fast regime in Section 2.1, and on the slow regime in Section 2.2. We conclude with the exact asymptotics for the intermediate case  $\alpha = 1$ , which follow directly from the Bahadur-Rao result; see Section 2.3.

**2.1. Fast regime.** In this section we assume that  $\alpha > 1$ . We start by proving an upper bound for  $P_N(a)$ . In self-evident notation, we have

$$P_N(a) = \int_0^\infty \psi_N(a|x) \mathbb{P}(\bar{X}_{N^\alpha} \in dx), \quad (11)$$

with  $\psi_N(a|x)$  as defined in (8). For any  $\delta$ , Eqn. (11) is majorized by

$$\int_{\nu-N^\delta}^{\nu+N^\delta} \psi_N(a|x) \mathbb{P}(\bar{X}_{N^\alpha} \in dx) + \mathbb{P}(|\bar{X}_{N^\alpha} - \nu| \geq N^\delta); \quad (12)$$

we determine an appropriate value for  $\delta$  later on. The first term in (12) is evidently bounded from above by  $\psi_N(a|\nu + N^\delta)$ . Motivated by (9), we will show that, as  $N \rightarrow \infty$ ,

$$\psi_N(a|\nu + N^\delta) e^{NI(a|\nu)} \sqrt{N} \rightarrow C(a|\nu), \quad (13)$$

whereas the second term in (12) turns out to be asymptotically negligible.

To verify that (13) holds, note that  $C(a|\nu)/C(a|\nu + N^\delta) \rightarrow 1$  when  $\delta < 0$ , which follows by a standard continuity argument. We therefore proceed by considering  $NI(a|\nu) - NI(a|\nu + N^\delta)$ , which behaves as

$$\begin{aligned} & N \left( a \log \frac{a}{\nu} + a - \nu \right) - N \left( a \log \frac{a}{\nu + N^\delta} + a - (\nu + N^\delta) \right) \\ &= Na \log \left( 1 + \frac{N^\delta}{\nu} \right) + N^{1+\delta} = \left( \frac{a}{\nu} + 1 \right) N^{1+\delta} + O(N^{1+2\delta}) \rightarrow 0 \end{aligned}$$

if  $\delta < -1$ . Thus, for such  $\delta$  we have established (13).

Now consider the second term of (12), and, more specifically,

$$\mathbb{P}(|\bar{X}_{N^\alpha} - \nu| \geq N^\delta) e^{NI(a|\nu)} \sqrt{N}, \quad (14)$$

for  $N \rightarrow \infty$ . Due to a Chernoff bound, we have

$$\mathbb{P}(\bar{X}_{N^\alpha} \geq \nu + N^\delta) \leq \exp \left( -N^\alpha \sup_{\vartheta} \left( \vartheta(\nu + N^\delta) - \log \mathbb{E} e^{\vartheta X_i} \right) \right) = e^{-N^\alpha I_X(\nu + N^\delta)},$$

and hence (14) is majorized by

$$e^{-N^\alpha I_X(\nu + N^\delta)} e^{NI(a|\nu)} \sqrt{N} + e^{-N^\alpha I_X(\nu - N^\delta)} e^{NI(a|\nu)} \sqrt{N}.$$

Now realize that  $I_X(\nu + N^\delta) = \frac{1}{2} I_X''(\nu) N^{2\delta} + O(N^{3\delta})$  and similarly for  $I_X(\nu - N^\delta)$ . Thus, the expression from the previous display vanishes when  $\alpha + 2\delta > 1$ , or, equivalently,  $\delta > (1 - \alpha)/2$ , where  $(1 - \alpha)/2 < 0$  since  $\alpha > 1$ .

We note that the requirements  $\delta < -1$  (corresponding to the first term) and  $\delta > (1 - \alpha)/2$  (corresponding to the second term) are both fulfilled when  $\alpha > 3$ . Thus, we have shown that for  $\alpha > 3$  an asymptotic upper bound for  $P_N(a)$  is given by (13).

Let us now turn to the corresponding lower bound. The probability of interest majorizes

$$\psi_N(a|\nu - N^\delta) \int_{\nu-N^\delta}^{\nu+N^\delta} \mathbb{P}(\bar{X}_{N^\alpha} \in dx).$$

As above, we can check that for  $\delta < -1$ ,

$$\psi_N(a|\nu - N^\delta) e^{NI(a|\nu)} \sqrt{N} \rightarrow C(a|\nu),$$

and, by the Bahadur-Rao result (5),

$$\int_{\nu-N^\delta}^{\nu+N^\delta} \mathbb{P}(\bar{X}_{N^\alpha} \in dx) \sim 1 - 2 \exp\left(-\frac{1}{2} I_X''(\nu) N^\alpha N^{2\delta}\right) \rightarrow 1,$$

when  $\delta > -\alpha/2$ . This can be realized when  $\alpha > 2$  (and is hence fulfilled when  $\alpha > 3$  as well). This proves the lower bound.

Combining the upper and lower bounds, we thus find the following result.

**Proposition 2.1.** *For  $\alpha > 3$ , as  $N \rightarrow \infty$ ,*

$$P_N(a) \sim e^{-NI(a|\nu)} \frac{C(a|\nu)}{\sqrt{N}}.$$

For  $\alpha \in (2, 3]$ ,

$$\liminf_{N \rightarrow \infty} P_N(a) e^{NI(a|\nu)} \sqrt{N} \geq C(a|\nu).$$

*Remark 1.* This result is in accordance with the intuition we gave at the beginning of the section – in the fast regime the asymptotics of  $P_N(a)$  depend on the distribution of the  $X_i$  only through their mean  $\nu$ . This also gives an indication as to why the asymptotics for  $\alpha$  closer to 1 may be more delicate to deal with. One can imagine that for more moderate values of  $\alpha$  the result may not be precise enough, and that also large deviations coming from  $\bar{X}_{N^\alpha}$  may play a role in that regime. This is confirmed in Section 3, where we consider an example with  $X_i \sim \text{Exp}(\lambda)$ . It turns out that the exact asymptotic expression for  $\alpha \in (1, 2)$  is indeed more intricate than the expression provided in Thm. 2.1.  $\diamond$

*Remark 2.* Along the same lines the asymptotics for  $p_N(a)$  can be found. They turn out to be, for  $\alpha > 3$ , as  $N \rightarrow \infty$ ,

$$p_N(a) \sim e^{-NI(a|\nu)} \frac{C(a|\nu)}{\sqrt{N}} \left(1 - e^{-I'(a|\nu)}\right).$$

This is in line with the result of Prop. 2.1: informally,

$$\begin{aligned} p_N(a) &= P_N(a) - P_N(a + 1/N) \\ &\approx \frac{C(a|\nu)}{\sqrt{N}} e^{-NI(a|\nu)} - \frac{C(a + 1/N|\nu)}{\sqrt{N}} e^{-NI(a+1/N|\nu)} \\ &\approx \frac{C(a|\nu)}{\sqrt{N}} e^{-NI(a|\nu)} \left(1 - e^{-I'(a|\nu)}\right), \end{aligned}$$

for large  $N$ , based on elementary Taylor arguments.  $\diamond$

**2.2. Slow regime.** We now consider the slow regime, i.e.,  $\alpha < 1$ . We have to distinguish between two cases.

- In Case I we assume that  $X_i$  may have outcomes larger than  $a$  with positive probability:

$$b_+ := \sup\{b : \mathbb{P}(X_i > b) > 0\} > a;$$

as a consequence  $I_X(a) < \infty$ . Recall that in this case, for  $N^\alpha$  substantially smaller than  $N$ , it can be argued that  $P_N(a)$  essentially behaves as  $\mathbb{P}(\bar{X}_{N^\alpha} \geq a)$ .

- In Case II we consider the opposite situation:  $b_+ < a$ . Then the intuition is that the rare event under consideration is the consequence of large deviations of both random components: of (i)  $\bar{X}_{N^\alpha}$  being close to  $b_+$ , and (ii) the Poisson variable with parameter (roughly)  $Nb_+$  attaining the atypical value  $Na$ .



Case I. We start by establishing an upper bound. Note that  $P_N(a)$  is majorized by

$$\mathbb{P}\left(\bar{X}_{N^\alpha} \geq a - N^\delta\right) + \psi_N(a | a - N^\delta).$$

Due to the Bahadur-Rao result stated in (5), the first term is asymptotically equivalent to

$$N^{-\alpha/2} C_X(a - N^\delta) e^{-N^\alpha I_X(a - N^\delta)},$$

which behaves as  $N^{-\alpha/2} C_X(a) e^{-N^\alpha I_X(a)}$  when  $\delta < -\alpha$  (as a direct consequence of the standard expansion  $I_X(a - N^\delta) = I_X(a) - N^\delta I'_X(a) + O(N^{2\delta})$ ). In addition, again using the Chernoff bound, we have

$$e^{N^\alpha I_X(a)} \psi_N(a | a - N^\delta) \leq e^{N^\alpha I_X(a)} \exp\left(-N \left(a \log \frac{a}{a - N^\delta} + N^\delta\right)\right). \quad (15)$$

Observe that the exponent in the second factor of the right hand side of (15) behaves as  $N^{2\delta+1}$ . We conclude that (15) vanishes if  $2\delta + 1 > \alpha$ , or, equivalently,  $\delta > (\alpha - 1)/2$  (note that  $(\alpha - 1)/2 < 0$ ). In order to simultaneously meet  $\delta < -\alpha$  and  $\delta > (\alpha - 1)/2$ , we need to have  $\alpha < \frac{1}{3}$ .

We now turn to the lower bound. The probability of interest is bounded from below by

$$\psi_N(a | a + N^\delta) \mathbb{P}\left(\bar{X}_{N^\alpha} \geq a + N^\delta\right).$$

The first factor is bounded from below by 1 minus a term that decays as  $\exp(-N^{1+2\delta})$  (which goes to 1 when  $\delta > -\frac{1}{2}$ ), whereas the second behaves as  $N^{-\alpha/2} C_X(a) e^{-N^\alpha I_X(a)}$  when  $\delta < -\alpha$ . In other words, there is an appropriate  $\delta$  for all  $\alpha < \frac{1}{2}$ . We have thus arrived at the following result.

**Proposition 2.2.** *Assume  $b_+ > a$ . For  $\alpha < \frac{1}{3}$ , as  $N \rightarrow \infty$ ,*

$$P_N(a) \sim e^{-N^\alpha I_X(a)} \frac{C_X(a)}{N^{\alpha/2}}.$$

For  $\alpha \in [\frac{1}{3}, \frac{1}{2})$ ,

$$\liminf_{N \rightarrow \infty} P_N(a) e^{N^\alpha I_X(a)} N^{\alpha/2} \geq C_X(a).$$

*Remark 3.* Note that here, in contrast with Prop. 2.1, the rate function is that of  $X$  rather than the Poisson random variable. As expected, when  $\alpha$  is small, the rare event is typically a result of a large deviation of  $\bar{X}_{N^\alpha}$ . However, for values of  $\alpha$  closer to 1 the same reasoning as in Remark 1 applies, and we do not expect a simple asymptotic expression as given in Prop. 2.2 to hold for all  $\alpha \in (\frac{1}{3}, 1)$  (as will be confirmed in Section 3, which covers the special case in which the  $X_i$  are exponentially distributed).  $\diamond$

*Remark 4.* As in Remark 2, the asymptotics for  $p_N(a)$  can be found as well. As it turns out, as  $N \rightarrow \infty$ ,

$$p_N(a) \sim e^{-N^\alpha I_X(a)} \frac{C_X(a) I'_X(a)}{N^{1-\alpha/2}}.$$

This is consistent with the result stated in Prop. 2.2:

$$\begin{aligned} p_N(a) &= P_N(a) - P_N(a + 1/N) \\ &\approx \frac{C_X(a)}{N^{\alpha/2}} e^{-N^\alpha I_X(a)} - \frac{C_X(a + 1/N)}{N^{\alpha/2}} e^{-N^\alpha I_X(a + 1/N)} \\ &\approx \frac{C_X(a)}{N^{\alpha/2}} e^{-N^\alpha I_X(a)} \left(1 - e^{-N^{\alpha-1} I'_X(a)}\right) \approx C_X(a) I'_X(a) e^{-N^\alpha I_X(a)} N^{\alpha/2-1}, \end{aligned}$$

for large  $N$ . Note that the asymptotic expansion of  $P_N(a)$  has a polynomial factor  $N^{-\alpha/2}$ , whereas  $p_N(a)$  has a polynomial factor  $N^{\alpha/2-1}$ . So in this case  $P_N(a)$  and  $p_N(a)$  are not (asymptotically) off by a constant, but by a constant multiplied by  $N^{\alpha-1}$ .  $\diamond$

*Case II.* In the above arguments for the slow regime, it is crucial that we assumed that  $X_i$  can exceed  $a$  with positive probability (which entails that  $I_X(a) < \infty$ ). We now consider the situation that  $b_+ < a$ . We derive the exact asymptotics of  $P_N(a)$  by separately considering a lower bound and an upper bound. We throughout assume that both  $I_X(b_+)$  and  $I'_X(b_+)$  are finite. The proof essentially follows that of [3], in which exact asymptotics of the Markov-modulated infinite-server queue are addressed.

We start with the lower bound. Let  $K$  the smallest value in  $\{2, 3, \dots\}$  such that  $-1/K$  is strictly larger than  $\alpha - 1$ . Fix  $\delta \in (\alpha - 1, -1/K)$ . We have, for  $\alpha < 1$ ,

$$P_N(a) \geq \int_{b_+ - N^\delta}^{b_+} \psi_N(a | x) \mathbb{P}(\bar{X}_{N^\alpha} \in dx) = \int_{b_+ - N^\delta}^{b_+} \psi_N(a | x) N^\alpha \xi_{N^\alpha}(N^\alpha x) dx,$$

recalling that  $\xi_N(\cdot)$  denotes the density of  $\sum_{i=1}^N X_i$ . Fix an arbitrary  $\zeta > 0$ . The right-hand side of the previous display majorizes, by Petrov's local limit version of the Bahadur-Rao result (5), in combination with (9), for  $N$  sufficiently large,

$$(1 - \zeta) \int_{b_+ - N^\delta}^{b_+} \frac{C(a | x)}{\sqrt{N}} e^{-N I(a | x)} \cdot C_X(x) I'_X(x) N^{\alpha/2} e^{-N^\alpha I_X(x)} dx.$$

This is in turn asymptotically equal to, with  $\gamma(a) := C(a | b_+) C_X(b_+) I'_X(b_+)$ , using the transformation  $y := b_+ - x$ ,

$$(1 - \zeta) \gamma(a) N^{(\alpha-1)/2} e^{-N I(a | b_+)} e^{-N^\alpha I_X(b_+)} \int_0^{N^\delta} e^{N \phi_1(y)} e^{N^\alpha \phi_2(y)} dy, \quad (16)$$

where

$$\phi_1(y) := -I(a | b_+ - y) + I(a | b_+) = a \log \left( 1 - \frac{y}{b_+} \right) + y, \quad \phi_2(y) := -I_X(b_+ - y) + I_X(b_+).$$

For all  $y \in [0, N^\delta]$ , there exist  $\ell_i$  and  $u_i$  ( $i = 1, 2$ ) such that

$$\begin{aligned} \ell_1 N^{1+K\delta} + \sum_{k=1}^{K-1} \beta_{1,k} N y^k &\leq N \phi_1(y) \leq u_1 N^{1+K\delta} + \sum_{k=1}^{K-1} \beta_{1,k} N y^k, \quad \beta_{1,1} := 1 - \frac{a}{b_+}, \\ \ell_2 N^{\alpha+K\delta} + \sum_{k=1}^{K-1} \beta_{2,k} N^\alpha y^k &\leq N^\alpha \phi_2(y) \leq u_2 N^{\alpha+K\delta} + \sum_{k=1}^{K-1} \beta_{2,k} N^\alpha y^k; \end{aligned}$$

observe that  $\beta_{1,1} < 0$ . We now further analyze the integral in (16). We find, using the above inequalities and the fact that both  $1 + K\delta < 0$  and  $\alpha + K\delta < 0$  (as we have chosen  $\delta < -1/K$ ),

$$\begin{aligned} \int_0^{N^\delta} e^{N \phi_1(y)} e^{N^\alpha \phi_2(y)} dy &\geq e^{\ell_1 N^{1+K\delta} + \ell_2 N^{\alpha+K\delta}} \int_0^{N^\delta} \exp \left( \sum_{k=1}^{K-1} \beta_{1,k} N y^k + \sum_{k=1}^{K-1} \beta_{2,k} N^\alpha y^k \right) dy \\ &\sim \int_0^{N^\delta} \exp \left( \sum_{k=1}^{K-1} \beta_{1,k} N y^k + \sum_{k=1}^{K-1} \beta_{2,k} N^\alpha y^k \right) dy. \end{aligned}$$

Applying the transformation  $z := N y$ , and using that  $\delta > -1$ , this integral can be evaluated as

$$\frac{1}{N} \int_0^{N^{\delta+1}} \exp \left( \sum_{k=1}^{K-1} \beta_{1,k} N^{1-k} z^k + \sum_{k=1}^{K-1} \beta_{2,k} N^{\alpha-k} z^k \right) dz \sim \frac{1}{N} \int_0^\infty e^{\beta_{1,1} z} dz = \frac{1}{N} \cdot \frac{b_+}{a - b_+}.$$

Letting  $\zeta \downarrow 0$ , we have thus found

$$\liminf_{N \rightarrow \infty} P_N(a) N^{(\alpha+1)/2} e^{N I(a | b_+)} e^{N^\alpha I_X(b_+)} \geq \gamma(a) \cdot \frac{b_+}{a - b_+}.$$

We proceed by the upper bound. Evidently,

$$P_N(a) = \int_{b_+ - N^\delta}^{b_+} \psi_N(a | x) \mathbb{P}(\bar{X}_{N^\alpha} \in dx) + \int_0^{b_+ - N^\delta} \psi_N(a | x) \mathbb{P}(\bar{X}_{N^\alpha} \in dx).$$

The first integral in the previous display can be dealt with as in the upper bound (*mutatis mutandis*; e.g. the factor  $1 - \zeta$  becomes  $1 + \zeta$ , and the  $u_i$  need to be used rather than the  $\ell_i$ ). We therefore focus on the second integral, which is clearly bounded above by  $\psi_N(a | b_+ - N^\delta)$ . Now observe that

$$I(a | b_+) - I(a | b_+ - N^\delta) = a \log \left( \frac{b_+ - N^\delta}{b_+} \right) + N^\delta \leq \left( 1 - \frac{a}{b_+} \right) N^\delta = \beta_{1,1} N^\delta.$$

As a consequence, as  $N \rightarrow \infty$ , recalling that  $\delta > \alpha - 1$  and  $\beta_{1,1} < 0$ ,

$$e^{NI(a | b_+)} e^{N^\alpha I_X(b_+)} \psi_N(a | b_+ - N^\delta) \leq e^{\beta_{1,1} N^{\delta+1}} e^{N^\alpha I_X(b_+)} \rightarrow 0.$$

We conclude the interval  $[0, b_+ - N^\delta)$  does not contribute to the asymptotics. We thus have established the upper bound, leading to the following result.

**Proposition 2.3.** *Assume  $\alpha < 1$  and  $b_+ < a$ . Then*

$$\lim_{N \rightarrow \infty} P_N(a) \sim e^{-NI(a | b_+)} e^{-N^\alpha I_X(b_+)} N^{-(\alpha+1)/2} \gamma(a) \frac{b_+}{a - b_+},$$

where  $\gamma(a) := C(a | b_+) C_X(b_+) I'_X(b_+)$ .

*Remark 5.* As before, we can identify the asymptotics of  $p_N(a)$  as well:

$$p_N(a) \sim e^{-NI(a | b_+)} e^{-N^\alpha I_X(b_+)} N^{-(\alpha+1)/2} \cdot \gamma(a) \cdot \frac{b_+}{a - b_+} \left( 1 - e^{-I'(a | b_+)} \right)$$

as  $N \rightarrow \infty$ . ◇

**2.3. Intermediate range.** We finally consider the case  $\alpha = 1$ . The random variable  $\text{Pois}(N\bar{X}_{N^\alpha})$  is distributed as the sum of  $N$  i.i.d. contributions, each of them distributed as  $Z := \text{Pois}(X)$ . Assuming that maximum in the definition (10) of  $I_Z(a)$  is attained at  $\vartheta_Z^*$ , the Bahadur-Rao result yields, as  $N \rightarrow \infty$ ,

$$P_N(a) \sim e^{-NI_Z(a)} \frac{C_Z(a)}{\sqrt{N}},$$

where now

$$\begin{aligned} C_Z(a) &:= \frac{1}{1 - e^{\vartheta_Z^*}} \frac{1}{\sqrt{2\pi \Lambda_Z''(\vartheta_Z^*)}} = \frac{1}{1 - e^{\vartheta_Z^*}} \frac{1}{\sqrt{2\pi (e^{\vartheta_Z^*} \Lambda'_X(e^{\vartheta_Z^*} - 1) + e^{2\vartheta_Z^*} \Lambda_X''(e^{\vartheta_Z^*} - 1))}} \\ &= \frac{1}{1 - e^{\vartheta_Z^*}} \frac{1}{\sqrt{2\pi (a + e^{2\vartheta_Z^*} \Lambda_X''(e^{\vartheta_Z^*} - 1))}}. \end{aligned}$$

Based on the same arguments as in Remark 2 we infer that

$$p_N(a) \sim \frac{1}{\sqrt{2\pi N (a + e^{2\vartheta_Z^*} \Lambda_X''(e^{\vartheta_Z^*} - 1))}} e^{-NI_Z(a)}.$$

### 3. ASYMPTOTICS OF $P_N(a)$ : SPECIAL CASE OF GAMMA $X_i$ 'S

In this section we consider the special case that the  $X_i$ s are i.i.d. samples from the gamma distribution. The use of this specific mixed Poisson distribution for call center staffing purposes is advocated in e.g. [11]. In the analysis, this can be reduced to the case where the  $X_i$ s are exponentially distributed with parameter  $\lambda$  (i.e., mean  $\lambda^{-1}$ ), see Remark 7.

To start the exposition, we note that if the  $X_i$ s are exponential with parameter  $\lambda$ , then  $\sum_{i=1}^{N^\alpha} X_i$  has a gamma distribution with parameters  $N^\alpha$  and  $\lambda$ . The objective of this section is to evaluate the asymptotics of  $p_N(a)$  across all values of  $\alpha$ ; later we comment on what the corresponding  $P_N(a)$  looks like. We assume throughout that  $a$  is larger than  $\lambda^{-1}$ . The computations are facilitated by

the fact that an exact expression for  $p_N(a)$  is available. It takes a routine calculation, which we include for completeness, to compute  $p_N(a)$ :

$$\begin{aligned}
 p_N(a) &= \int_0^\infty \frac{(N^{1-\alpha}x)^{Na}}{(Na)!} e^{-(\lambda+N^{1-\alpha})x} \frac{\lambda^{N^\alpha}}{(N^\alpha-1)!} x^{N^\alpha-1} dx \\
 &= \frac{(N^{1-\alpha})^{Na}}{(Na)!} \frac{\lambda^{N^\alpha}}{(N^\alpha-1)!} \int_0^\infty e^{-(\lambda+N^{1-\alpha})x} x^{Na+N^\alpha-1} dx \\
 &= \frac{(Na+N^\alpha-1)!}{(Na)!(N^\alpha-1)!} \frac{(N^{1-\alpha})^{Na} \lambda^{N^\alpha}}{(\lambda+N^{1-\alpha})^{Na+N^\alpha}} \int_0^\infty \frac{(\lambda+N^{1-\alpha})^{N^\alpha}}{(Na+N^\alpha-1)!} e^{-(\lambda+N^{1-\alpha})x} x^{Na+N^\alpha-1} dx \\
 &= \binom{Na+N^\alpha-1}{Na} \left( \frac{N^{1-\alpha}}{\lambda+N^{1-\alpha}} \right)^{Na} \left( \frac{\lambda}{\lambda+N^{1-\alpha}} \right)^{N^\alpha}.
 \end{aligned}$$

*Remark 6.* We recognize here the probability that a *negative binomially distributed* random variable with success probability  $p := N^{1-\alpha}/(\lambda+N^{1-\alpha})$  attains  $Na$  successes before  $N^\alpha$  failures have occurred. This can be understood as follows. Note that a Poisson random variable with parameter  $xT$  represents the number of  $\text{Exp}(x)$  “success clocks” expiring within a period of length  $T$ . In our case the rate of the success clocks is  $x = N^{1-\alpha}$  and the length of the period corresponds to the time it takes for  $N$  exponential “failure clocks” of rate  $\lambda$  to expire, that is, we have  $T = \sum_{i=1}^{N^\alpha} X_i$ . Thus,  $p_N(a)$  is the probability that  $Na$  success clocks expire before the  $N^\alpha$ th failure clock expires and the period ends. The success probability is indeed given by  $p$  as it is the probability that the next  $\text{Exp}(N^{1-\alpha})$  success clock expires before a  $\text{Exp}(\lambda)$  failure clock.  $\diamond$

*Remark 7.* In the above setup we considered exponentially distributed  $X_i$ s. Note, however, that our analysis only relies on  $\sum_{i=1}^{N^\alpha} X_i$  having a gamma distribution, and thus can easily be extended to the practically relevant case [11] that the  $X_i$ s are i.i.d. samples from a gamma distribution. It is noted that the gamma distribution has two parameters (as opposed to the exponential distribution), and therefore allows for more modelling flexibility (e.g., the mean and variance can be fitted).  $\diamond$

As a first step in deriving the exact asymptotics of  $p_N(a)$ , we approximate the binomial coefficients by applying Stirling’s formula, which says that  $n! \sim \sqrt{2\pi n} n^n e^{-n}$ . As a consequence we find that

$$\binom{Na+N^\alpha-1}{Na} \sim \frac{1}{\sqrt{2\pi}} \frac{\sqrt{Na+N^\alpha-1}}{\sqrt{Na}\sqrt{N^\alpha-1}} \frac{(Na+N^\alpha-1)^{Na+N^\alpha-1}}{(Na)^{Na}(N^\alpha-1)^{N^\alpha-1}}$$

Applying this in the expression for  $p_N(a)$  then yields

$$\begin{aligned}
 p_N(a) &= \binom{Na+N^\alpha-1}{Na} \left( \frac{N^{1-\alpha}}{\lambda+N^{1-\alpha}} \right)^{Na} \left( \frac{\lambda}{\lambda+N^{1-\alpha}} \right)^{N^\alpha} \\
 &\sim \frac{1}{\sqrt{2\pi}} \frac{\sqrt{Na+N^\alpha-1}}{\sqrt{Na}\sqrt{N^\alpha-1}} \frac{(Na+N^\alpha-1)^{Na+N^\alpha-1}}{(Na)^{Na}(N^\alpha-1)^{N^\alpha-1}} \left( \frac{N^{1-\alpha}}{\lambda+N^{1-\alpha}} \right)^{Na} \left( \frac{\lambda}{\lambda+N^{1-\alpha}} \right)^{N^\alpha} \\
 &= \frac{1}{\sqrt{2\pi}} \frac{\sqrt{N^\alpha-1}}{\sqrt{Na}\sqrt{Na+N^\alpha-1}} \cdot \left( \frac{Na+N^\alpha-1}{a\lambda(N^\alpha+\frac{N}{\lambda})} \right)^{Na} \cdot \left( \frac{\lambda(Na+N^\alpha-1)}{(\lambda+N^{1-\alpha})(N^\alpha-1)} \right)^{N^\alpha}. \quad (17)
 \end{aligned}$$

In order to determine the asymptotic behavior of this expression for large  $N$ , we again consider the three regimes separately. We do so by evaluating the three factors in (17).

**3.1. Fast regime.** We start by examining the case  $\alpha > 1$ . For the first factor we have

$$\frac{1}{\sqrt{2\pi}} \frac{\sqrt{N^\alpha-1}}{\sqrt{Na}\sqrt{Na+N^\alpha-1}} \sim \frac{1}{\sqrt{2\pi}} \frac{1}{\sqrt{Na}}.$$

The middle factor can be addressed as follows. For ease we analyze its logarithm:

$$Na \log \left( \frac{Na + N^\alpha - 1}{a\lambda(N^\alpha + N/\lambda)} \right) = -Na \log(a\lambda) + Na \log(1 + N^{1-\alpha}a - N^{-\alpha}) - Na \log(1 + N^{1-\alpha}/\lambda) \quad (18)$$

For the last factor we similarly obtain

$$N^\alpha \log \left( \frac{\lambda(Na + N^\alpha - 1)}{(\lambda + N^{1-\alpha})(N^\alpha - 1)} \right) = N^\alpha \log(1 + N^{1-\alpha}a - N^{-\alpha}) - N^\alpha \log \left( 1 + \frac{1}{\lambda}N^{1-\alpha} - \frac{1}{\lambda}N^{1-2\alpha} - N^{-\alpha} \right) \quad (19)$$

Define  $\bar{k} := (\alpha - 1)^{-1}$  and  $k_+ := \lceil \bar{k} \rceil$ . The logarithms can be expanded relying on their standard Taylor series form, but it can be argued that the resulting infinite series can be truncated. For instance,

$$Na \log(1 + N^{1-\alpha}a - N^{-\alpha}) = Na \sum_{k=1}^{\infty} \frac{(-1)^{k+1}}{k} (N^{(1-\alpha)}a - N^{-\alpha})^k \sim Na \sum_{k=1}^{k_+} \frac{(-1)^{k+1}a^k}{k} N^{(1-\alpha)k}.$$

Likewise,

$$Na \log(1 + N^{1-\alpha}/\lambda) \sim Na \sum_{k=1}^{k_+} \frac{(-1)^{k+1}\lambda^{-k}}{k} N^{(1-\alpha)k}.$$

We thus find that (18) asymptotically equals

$$-Na \log(a\lambda) + Na \sum_{k=1}^{k_+} \frac{(-1)^k(\lambda^{-k} - a^k)}{k} N^{(1-\alpha)k}.$$

For the last factor, note that from  $k_+ + 1$  on all terms vanish, leaving us with

$$N^\alpha \log(1 + N^{1-\alpha}a - N^{-\alpha}) \sim N^\alpha \sum_{k=1}^{k_++1} \frac{(-1)^{k+1}a^k}{k} N^{(1-\alpha)k} - 1,$$

$$N^\alpha \log(1 + N^{1-\alpha}/\lambda - N^{1-2\alpha}/\lambda - N^{-\alpha}) \sim N^\alpha \sum_{k=1}^{k_++1} \frac{(-1)^{k+1}\lambda^{-k}}{k} N^{(1-\alpha)k} - 1.$$

After a bit of rewriting, we conclude that (19) equals

$$N \sum_{k=0}^{k_+} \frac{(-1)^k(a^{k+1} - \lambda^{-(k+1)})}{k+1} N^{(1-\alpha)k}.$$

Defining

$$\xi_0 := -a \log(\lambda a) + a - \frac{1}{\lambda}, \quad \xi_k := (-1)^k \left( \lambda^{-k} \left( \frac{a}{k} - \frac{1/\lambda}{k+1} \right) - a^{k+1} \left( \frac{1}{k} - \frac{1}{k+1} \right) \right),$$

we conclude that in case  $\alpha > 1$ ,

$$p_N(a) \sim \frac{1}{\sqrt{2\pi a N}} e^{\xi_0 N} \exp \left( \sum_{k=1}^{k_+} \xi_k N^{(1-\alpha)k+1} \right).$$

In particular, if  $\alpha > 2$ , then the last factor equals 1 (the empty sum being defined as 0). It is not hard to check that this result agrees with what has been found for  $\alpha > 3$  in Section 2.

3.2. **Slow regime.** If  $\alpha < 1$ , then the first factor behaves as

$$\frac{1}{\sqrt{2\pi}} \frac{\sqrt{N^\alpha - 1}}{\sqrt{Na}\sqrt{Na + N^\alpha - 1}} \sim \frac{1}{\sqrt{2\pi}} \frac{1}{a} N^{\alpha/2-1}.$$

For the logarithm of the middle factor we now have

$$Na \log \left( \frac{Na + N^\alpha - 1}{a\lambda(N^\alpha + N/\lambda)} \right) = Na \log \left( 1 + \frac{1}{a}(N^{\alpha-1} - N^{-1}) \right) - Na \log(1 + \lambda N^{\alpha-1}). \quad (20)$$

With  $\tilde{k} := \alpha(1 - \alpha)^{-1}$  and  $k_- := \lfloor \tilde{k} \rfloor$ , we obtain that this factor asymptotically equals

$$Na \sum_{k=1}^{k_-+1} \frac{(-1)^{k+1}(a^{-k} - \lambda^k)}{k} N^{(\alpha-1)k} - 1 = Na \sum_{k=0}^{k_-} \frac{(-1)^k(a^{-(k+1)} - \lambda^{k+1})}{k+1} N^{(\alpha-1)(k+1)} - 1.$$

For the last factor we find

$$\begin{aligned} N^\alpha \log \left( \frac{\lambda(Na + N^\alpha - 1)}{(\lambda + N^{1-\alpha})(N^\alpha - 1)} \right) &= N^\alpha \log(\lambda a) + N^\alpha \log \left( 1 + \frac{1}{a}N^{\alpha-1} - \frac{1}{a}N^{-1} \right) \\ &\quad - N^\alpha \log(1 + \lambda N^{\alpha-1} - \lambda N^{-1} - N^{-\alpha}) \end{aligned}$$

where

$$\begin{aligned} N^\alpha \log \left( 1 + \frac{1}{a}N^{\alpha-1} - \frac{1}{a}N^{-1} \right) &\sim N^\alpha \sum_{k=1}^{k_-} \frac{(-1)^{k+1}}{k} a^{-k} N^{(\alpha-1)k}, \\ N^\alpha \log(1 + \lambda N^{\alpha-1} - \lambda N^{-1} - N^{-\alpha}) &\sim N^\alpha \sum_{k=1}^{k_-} \frac{(-1)^{k+1}}{k} \lambda^k N^{(\alpha-1)k} - 1. \end{aligned}$$

Combining the above we conclude

$$p_N(a) \sim \frac{1}{\sqrt{2\pi a}} N^{\frac{\alpha}{2}-1} e^{\zeta_0 N^\alpha} \exp \left( \sum_{k=1}^{k_-} \zeta_k N^{(\alpha-1)k+\alpha} \right), \quad (21)$$

where

$$\zeta_0 := \log(\lambda a) + 1 - \lambda a, \quad \zeta_k := (-1)^k \left( \lambda^k \left( \frac{1}{k} - \frac{a\lambda}{k+1} \right) - a^{-k} \left( \frac{1}{k} - \frac{1}{k+1} \right) \right).$$

It can again be verified that this result coincides for  $\alpha < \frac{1}{3}$  with the one derived in Section 2.

3.3. **Intermediate regime.** For completeness, we also include the result for the case  $\alpha = 1$ . We find

$$p_N(a) \sim \frac{1}{\sqrt{2\pi}} \frac{1}{\sqrt{Na(a+1)}} \exp \left( -N \left( a \log \left( a \frac{1+\lambda}{1+a} \right) + \log \left( \frac{1}{\lambda} \frac{1+\lambda}{1+a} \right) \right) \right). \quad (22)$$

It is noted that the asymptotics of  $P_N(a)$  and  $p_N(a)$  could have been found by applying the Bahadur-Rao result directly, as noted in Section 2.3:

$$P_N(a) \sim \frac{1}{1 - e^{\vartheta_Z^*}} \frac{1}{\sqrt{2\pi N \Lambda_Z''(\vartheta_Z^*)}} e^{-NI_Z(a)} = \frac{1}{1 - a \frac{1+\lambda}{1+a}} \frac{1}{\sqrt{2\pi Na(a+1)}} e^{-NI_Z(a)}.$$

and

$$p_N(a) = P_N(a) - P_N \left( a + \frac{1}{N} \right) \sim \frac{1}{\sqrt{2\pi Na(a+1)}} e^{-NI_Z(a)},$$

where it can be verified that  $I_Z(a)$  coincides with the exponent found in (22).

**3.4. Example.** In Fig. 1 we illustrate the accuracy of the approximation, by displaying the ratio of the approximation  $\tilde{p}_N(a)$  and the exact expression for  $p_N(a)$ . We observe that this ratio tends to 1 as  $N$  grows, as expected.

Note that the naive approximation  $p_N(a) \approx \exp(-N^\alpha I_X(a))$  obtained from the logarithmic asymptotics is still very far off the true value for the small values  $N$  considered in this example. For example, with  $\alpha < 1$  the ratio is as high as  $\sqrt{2\pi a} N^{1-\alpha/2} \tilde{p}_N(a)/p_N(a)$  (cf. (21)). This shows how important it can be to consider exact asymptotics instead of logarithmic asymptotics.

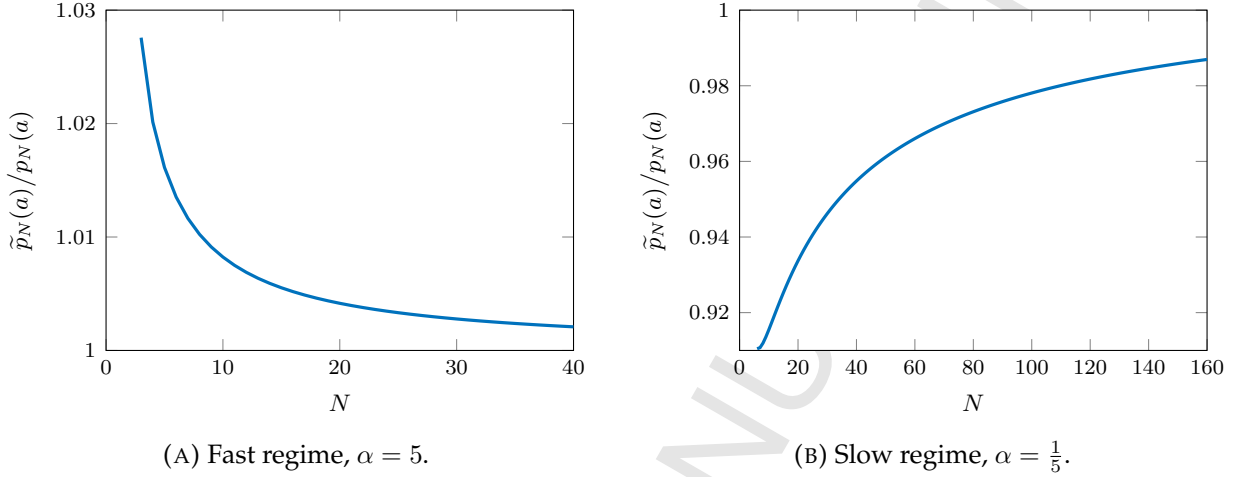


FIGURE 1. Ratio of approximation  $\tilde{p}_N(a)$  and exact value  $p_N(a)$ , where  $X_i$  is exponentially distributed with parameter  $\lambda = 2.5$  and  $a = 1$ .

#### 4. IMPORTANCE SAMPLING FOR $P_N(a)$

In the previous sections we found exact asymptotics for the rare-event probabilities  $p_N(a)$  and  $P_N(a)$  for (i) a specific range of  $\alpha$ , and (ii) for the specific case that the  $X_i$  are exponentially distributed. To facilitate numerical evaluation (which we need, for example, if (i) and (ii) do not apply), we propose in this section importance sampling estimators for  $p_N(a)$  and  $P_N(a)$ . We establish asymptotic efficiency properties, thus guaranteeing fast computation even for large  $N$ . As before, we distinguish the cases  $\alpha < 1$  and  $\alpha > 1$ ; the case  $\alpha = 1$  can be addressed by using a classical importance sampling procedure.

**4.1. Fast regime.** Recall that in this regime a rare event is typically the result of a large deviation of the Poisson random variable, while the sample mean  $X_1, \dots, X_{N^\alpha}$  will typically be close to  $\nu$  (under their true distribution, which we shall indicate by a subscript  $\nu$ ). In view of this, we propose a somewhat unconventional importance sampling estimator (cf. the more classical estimator (26) that we will come across in the slow regime). Based on  $n \in \mathbb{N}$  runs,  $P_N^{(n)}(a)$  can be unbiasedly estimated by

$$\hat{P}_N^{(n)}(a) = \frac{1}{n} \sum_{i=1}^n \frac{\mathbb{P}(\text{Pois}(N\bar{X}_{N^\alpha, i}) = Z_i)}{\mathbb{P}(\text{Pois}(Na) = Z_i)} \mathbb{1}\{Z_i \geq Na\}, \quad (23)$$

where (i)  $Z_1, \dots, Z_n \sim \text{Pois}(Na)$  (independently sampled), and (ii)  $\bar{X}_{N^\alpha, 1}, \dots, \bar{X}_{N^\alpha, n}$  independently sampled under the original measure.

Observe that the contribution from the  $i$ th run depends on  $\bar{X}_{N^\alpha, i}$  as well as  $Z_i$ . It is therefore easier to analyze the corresponding estimator for  $p_N(a)$ ,

$$\hat{p}_N^{(n)}(a) := \frac{1}{n} \sum_{i=1}^n \frac{\mathbb{P}(\text{Pois}(N\bar{X}_{N^\alpha, i}) = Z_i)}{\mathbb{P}(\text{Pois}(Na) = Z_i)} \mathbb{1}\{Z_i = Na\},$$

which does not depend on the specific value of  $Z_i$  (as it is  $Na$  with certainty). We later comment on efficient estimation of  $P_N(a)$ .

The contribution due to the likelihood ratio of the  $i$ th run is

$$L(\bar{X}_{N^\alpha, i}) := \left( \frac{\bar{X}_{N^\alpha, i}}{a} \right)^{Na} e^{N(a - \bar{X}_{N^\alpha, i})}.$$

The variance of the estimator (with respect to the joint distribution of  $Z \sim \text{Pois}(Na)$  and  $\bar{X}_{N^\alpha}$ ) can be evaluated to be

$$\frac{1}{n} \mathbb{E} \left[ (L(\bar{X}_{N^\alpha}) \mathbb{1}\{Z = Na\})^2 \right] - p_N(a)^2 = \frac{1}{n} \mathbb{E} [L^2(\bar{X}_{N^\alpha}) \mathbb{1}\{Z = Na\}] - p_N(a)^2, \quad (24)$$

with  $Z$  distributed as each of the  $Z_i$ , and  $\bar{X}_{N^\alpha}$  as each of the  $\bar{X}_{N^\alpha, i}$ . As we have seen in the introduction of Section 2, the logarithmic decay rate of  $p_N(a)^2$  is  $-2I(a|\nu)$ . Since the variance is always non-negative, this implies that the first term in (24) vanishes no faster than with exponential rate  $-2I(a|\nu)$ . This motivates the following notion of *asymptotic efficiency* (or logarithmical efficiency), as suggested in e.g. [16].

**Proposition 4.1.** *The estimator  $\hat{p}_N^{(n)}(a)$  is asymptotically efficient for estimating  $p_N(a)$ ; that is*

$$\limsup_{N \rightarrow \infty} \frac{1}{N} \log \mathbb{E} [L^2(\bar{X}_{N^\alpha}) \mathbb{1}\{Z = Na\}] \leq -2I(a|\nu).$$

*Proof.* First, note that

$$\begin{aligned} \mathbb{E} [L^2(\bar{X}_{N^\alpha}) \mathbb{1}\{Z = Na\}] &= \mathbb{E}_\nu \left[ \left( \frac{\bar{X}_{N^\alpha}}{a} \right)^{2Na} e^{2N(a - \bar{X}_{N^\alpha})} \right] \mathbb{P}(Z = Na) \\ &\leq \mathbb{E}_\nu \left[ \left( \frac{\bar{X}_{N^\alpha}}{a} \right)^{2Na} e^{2N(a - \bar{X}_{N^\alpha})} \right]. \end{aligned}$$

Define  $\mathcal{F}_\varepsilon^{(N)} := \{\bar{X}_{N^\alpha} \in (\nu - \varepsilon, \nu + \varepsilon)\}$ , where  $\varepsilon > 0$ . Then

$$\mathbb{E}_\nu \left[ \left( \frac{\bar{X}_{N^\alpha}}{a} \right)^{2Na} e^{2N(a - \bar{X}_{N^\alpha})} \mathbb{1}\{\mathcal{F}_\varepsilon^{(N)}\} \right] \leq \left( \frac{\nu + \varepsilon}{a} \right)^{2Na} e^{2N(a - \nu + \varepsilon)}. \quad (25)$$

On the other hand, we have

$$\mathbb{E}_\nu \left[ \left( \frac{\bar{X}_{N^\alpha}}{a} \right)^{2Na} e^{2N(a - \bar{X}_{N^\alpha})} \mathbb{1}\{(\mathcal{F}_\varepsilon^{(N)})^c\} \right] = \mathbb{E}_\nu \left[ e^{-2NI(a|\bar{X}_{N^\alpha})} \mathbb{1}\{(\mathcal{F}_\varepsilon^{(N)})^c\} \right] \leq \mathbb{P} \left( [\mathcal{F}_\varepsilon^{(N)}]^c \right).$$

where the last inequality is due to  $I(a|x) \geq 0$  for any  $x$ . Invoking Chernoff's bound, we note that

$$\mathbb{P} \left( [\mathcal{F}_\varepsilon^{(N)}]^c \right) \leq 2 \exp(-N^\alpha j_\varepsilon), \quad \text{where } j_\varepsilon := \inf_{x \notin (\nu - \varepsilon, \nu + \varepsilon)} I_X(x) > 0.$$

We conclude that for  $\alpha > 1$ ,

$$\limsup_{N \rightarrow \infty} \frac{N^\alpha}{N} \frac{1}{N^\alpha} \log \mathbb{P} \left( [\mathcal{F}_\varepsilon^{(N)}]^c \right) \leq \limsup_{N \rightarrow \infty} -\frac{N^\alpha}{N} j_\varepsilon = -\infty.$$

Combining this with (25), we conclude that

$$\limsup_{N \rightarrow \infty} \frac{1}{N} \log \mathbb{E} \left[ (L(\bar{X}_{N^\alpha}) \mathbb{1}\{Z = Na\})^2 \right] \leq 2a \log \left( \frac{\nu + \varepsilon}{a} \right) + 2(a - \nu + \varepsilon).$$

The desired result follows when taking  $\varepsilon \downarrow 0$ . □



Formally, this result on asymptotic efficiency for  $\hat{p}_N^{(n)}(a)$  does not imply asymptotic efficiency for  $\hat{P}_N^{(n)}(a)$ . In practice, however, we can use

$$\hat{P}_N^{(n)}(a) = \sum_{k=Na}^K \hat{p}_N^{(n)}(k/N),$$

with  $K$  sufficiently large, to estimate  $P_N^{(n)}(a)$ .

**4.2. Slow regime.** In the slow regime, assuming that  $b_+ \geq a$ , the rare event is typically caused by a large deviation of  $\bar{X}_{N^\alpha}$ . Suppose that  $\bar{X}_{N^\alpha,1}, \dots, \bar{X}_{N^\alpha,n}$  are independently sampled according to the original measure  $\mathbb{P}_\nu$  (where the subscript indicates that the expectation of each of the sample means  $\bar{X}_{N^\alpha,i}$  involved is  $\nu$ ). In this case we suggest the estimator

$$\hat{P}_N^{(n)}(a) = \frac{1}{n} \sum_{i=1}^n \frac{\mathbb{P}_\nu(\bar{X}_{N^\alpha,i} \in dY_i)}{\mathbb{P}_a(\bar{X}_{N^\alpha,i} \in dY_i)} \mathbb{1} \{ \text{Pois}(NY_i) \geq Na \}, \quad (26)$$

where  $Y_1, \dots, Y_n \sim \mathbb{P}_a$ . The measure  $\mathbb{P}_a$  corresponds to the exponentially twisted version such that the mean becomes  $a$  (rather than  $\nu$ ).

For each run we have the likelihood ratio, with  $\vec{y} = (y_1, \dots, y_{N^\alpha})$ ,

$$L(\vec{y}) = \prod_{i=1}^{N^\alpha} M_X(\vartheta_a) e^{-\vartheta_a y_i},$$

where we recall that  $M_X(\cdot)$  is the moment-generating function of  $X$  and  $\vartheta_a$  is the unique solution to

$$\mathbb{E}_a[X] = \mathbb{E}_\nu \left[ X \frac{e^{\vartheta X}}{M_X(\vartheta)} \right] = \frac{M'_X(\vartheta)}{M_X(\vartheta)} = a.$$

In this case we have seen before that  $N^{-\alpha} \log P_N(a) \rightarrow -I_X(a)$  as  $N \rightarrow \infty$ .

**Proposition 4.2.** *The estimator  $\hat{P}_N^{(n)}(a)$  is asymptotically efficient for estimating  $P_N(a)$ ; that is*

$$\limsup_{N \rightarrow \infty} \frac{1}{N^\alpha} \log \mathbb{E}_a \left[ \left( L(\vec{X}) \mathbb{1} \{ \text{Pois}(N\bar{X}_{N^\alpha}) \geq Na \} \right)^2 \right] \leq -2I_X(a).$$

*Proof.* Note that

$$\mathbb{E}_a \left[ \left( L(\vec{X}) \mathbb{1} \{ \text{Pois}(N\bar{X}_{N^\alpha}) \geq Na \} \right)^2 \right] = M(\vartheta_a)^{2N^\alpha} \mathbb{E}_a \left[ e^{-2\vartheta_a N^\alpha \bar{X}_{N^\alpha}} \mathbb{1} \{ \text{Pois}(N\bar{X}_{N^\alpha}) \geq Na \} \right].$$

On  $\mathcal{F}_\varepsilon^{(N)} := \{ \bar{X}_{N^\alpha} \in (a - \varepsilon, \infty) \}$  we have

$$\mathbb{E}_a \left[ e^{-2\vartheta_a N^\alpha \bar{X}_{N^\alpha}} \mathbb{1} \{ \text{Pois}(N\bar{X}_{N^\alpha}) \geq Na \} \mathbb{1} \{ \mathcal{F}_\varepsilon^{(N)} \} \right] \leq e^{-2\vartheta_a N^\alpha (a - \varepsilon)},$$

while outside of  $\mathcal{F}_\varepsilon^{(N)}$  we have

$$\mathbb{E}_a \left[ e^{-2\vartheta_a N^\alpha \bar{X}_{N^\alpha}} \mathbb{1} \{ \text{Pois}(N\bar{X}_{N^\alpha}) \geq Na \} \mathbb{1} \left\{ \left[ \mathcal{F}_\varepsilon^{(N)} \right]^c \right\} \right] \leq \mathbb{P}_a(\text{Pois}(N(a - \varepsilon)) \geq Na),$$

where we used that  $\vartheta_a > 0$  because  $a > \nu$  [7, Lemma 2.2.5]. By virtue of the Chernoff bound,

$$\mathbb{P}_a(\text{Pois}(N(a - \varepsilon)) \geq Na) \leq e^{-NI(a|a - \varepsilon)}, \text{ where } I(a|a - \varepsilon) > 0.$$

This implies that

$$\begin{aligned} \limsup_{N \rightarrow \infty} \frac{1}{N^\alpha} \log \mathbb{E}_a \left[ e^{-2\vartheta_a N^\alpha \bar{X}_{N^\alpha}} \mathbb{1} \{ \text{Pois}(N\bar{X}_{N^\alpha}) \geq Na \} \mathbb{1} \left\{ \left[ \mathcal{F}_\varepsilon^{(N)} \right]^c \right\} \right] \\ \leq \limsup_{N \rightarrow \infty} -\frac{N}{N^\alpha} I(a|a - \varepsilon) = -\infty. \end{aligned}$$

We let first  $N \rightarrow \infty$  and then  $\varepsilon \downarrow 0$ , to conclude that

$$\limsup_{N \rightarrow \infty} \frac{1}{N^\alpha} \log \mathbb{E}_a \left[ \left( L(\vec{X}) \mathbb{1} \{ \text{Pois}(N\bar{X}_{N^\alpha}) \geq Na \} \right)^2 \right] \leq 2 \log M_X(\vartheta_a) - 2\vartheta_a a = -2I_X(a),$$

as claimed.  $\square$

**4.3. Numerical example.** We provide a numerical example with exponentially distributed  $X_i$ . Specifically, we consider  $X_i \sim \text{Exp}(1)$ ,  $a = 2$ , and  $\alpha \in \{0.5, 2\}$ . Fig. 2 shows the logarithm of  $\hat{P}_N(a)$  as well as the corresponding crude Monte Carlo estimators, as a function of  $N$ . We generated  $\sum_{i=1}^{N^\alpha} X_i$  by drawing from the gamma distribution with parameters  $N^\alpha$  and  $1/\lambda$ . This allowed us to include values of  $N$  for which  $N^\alpha \notin \mathbb{N}$  in Fig. 2.(B). The dotted lines in the figures indicate the upper bounds of the standard normal 95% confidence intervals evaluated using sample standard deviations (multiplied by a factor  $10^3$  to make them visible). It can be seen that for the importance sampling estimator the width of the confidence interval hardly depends on  $N$ . In contrast, for the Monte Carlo estimator the width of the confidence interval increases significantly.

## 5. ASYMPTOTICS FOR INFINITE-SERVER SYSTEM, AND IMPLICATIONS FOR STAFFING

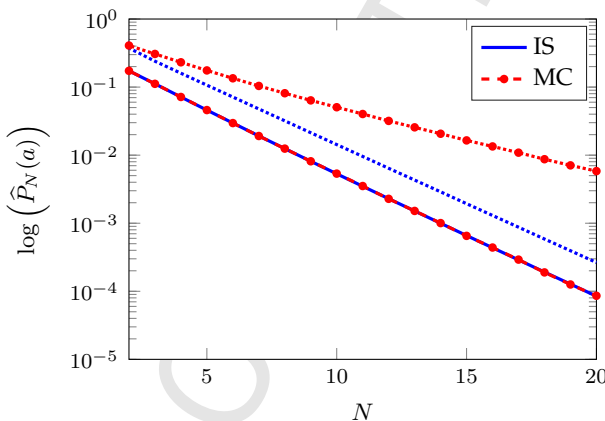
In this section we investigate the asymptotic behavior of  $Q_N(a)$  ( $q_N(a)$ ), the probability that the number of clients in the system exceeds (equals) some threshold  $Na$ . We consider the scaled system previously studied in [9].

We start by presenting the logarithmic asymptotics, which can be identified with exactly the same techniques as in [9, Section 4.1]. As before, we distinguish three cases (where it is noted that  $q_N(a)$  has the same logarithmic asymptotics as  $Q_N(a)$ ); the intuition behind the three regimes is as before.

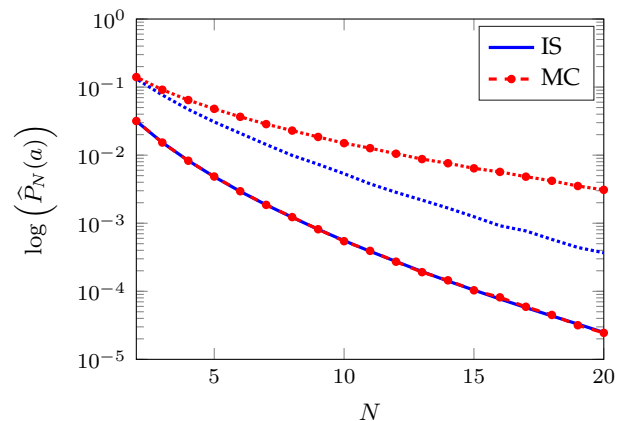
- For  $\alpha > 1$ ,

$$\lim_{N \rightarrow \infty} \frac{1}{N} \log Q_N(a) = -I \left( a \mid \nu \int_0^1 \bar{F}(x) dx \right),$$

where  $\bar{F}(\cdot)$  denotes the complementary distribution function of the service times.



(A) Fast regime,  $\alpha = 2$ .



(B) Slow regime,  $\alpha = 0.5$ .

FIGURE 2. Logarithmic importance sampling (IS) and crude Monte Carlo (MC) estimators for  $P_N(a)$ , where  $X_i$  is exponentially distributed with parameter  $\lambda_\alpha$  (where  $\lambda_2 = 1$ ,  $\lambda_{0.5} = 2.5$ ) and  $a = 2$ , averaged over  $n = 10^7$  samples. The upper bounds of the sample confidence intervals are indicated by dashed lines; the width of the intervals is inflated by a factor  $10^3$  for better visibility.

- For  $\alpha < 1$ , assuming the support of  $X_i$  is unbounded,

$$\lim_{N \rightarrow \infty} \frac{1}{N^\alpha} \log Q_N(a) = - \sup_{\vartheta} \left( \vartheta a - \int_0^1 \Lambda_X (\vartheta \bar{F}(x)) dx \right).$$

- For  $\alpha = 1$ ,

$$\lim_{N \rightarrow \infty} \frac{1}{N} \log Q_N(a) = - \sup_{\vartheta} \left( \vartheta a - \int_0^1 \Lambda_X \left( (e^\vartheta - 1) \bar{F}(x) \right) dx \right). \quad (27)$$

In the remainder of this section, we first determine the exact asymptotics for the special case  $\alpha = 1$ . That is, we assume that the arrival rates are resampled every  $1/N$  time units, and we are interested in the number of customers present at time 1 (that is, after  $N$  time periods of length  $1/N$ ). As it turns out, the case  $\alpha \neq 1$  is considerably harder to deal with, and therefore left for future research. We conclude this section by a set of numerical experiments.

**5.1. Exact asymptotics.** As mentioned in the introduction (*viz.* Eqn. (1)), under the scaling of [9] the number of clients in the system at time 1 is distributed as the sum of  $N$  Poisson random variables, say,  $Z_1$  up to  $Z_N$ , where  $Z_i$  can be interpreted as the contribution due to arrivals in the interval  $[(i-1)/N, i/N]$ ; for details we refer to [9]. Then it can be argued that

$$Z_i \stackrel{d}{=} \text{Pois} (NX_i \cdot N^{-1} \omega_i(N)) = \text{Pois} (X_i \omega_i(N)),$$

where we defined  $\omega_i(N)$  as the probability that a call that arrived at a uniform epoch in the interval  $[(i-1)/N, i/N]$  is still present at time 1. It can be verified that

$$\omega_i(N) = N \int_{(i-1)/N}^{i/N} \bar{F}(1-x) dx;$$

because the  $X_i$  are i.i.d., we can reverse time, and hence replace  $\bar{F}(1-x)$  in the previous display by  $\bar{F}(x)$ .

We now wish to evaluate

$$Q_N(a) = \mathbb{P} \left( \text{Pois} \left( \sum_{i=1}^N X_i \omega_i(N) \right) \geq Na \right), \quad q_N(a) = \mathbb{P} \left( \text{Pois} \left( \sum_{i=1}^N X_i \omega_i(N) \right) = Na \right).$$

Let  $S_N = \sum_{i=1}^N Z_i$ , where  $Z_i \stackrel{d}{=} \text{Pois} (X_i \omega_i(N))$ , with the  $Z_i$  independent; hence  $q_N(a) = \mathbb{P}(S_N = Na)$ . It is immediately verified that, with  $M_X(\cdot)$  the moment generating function of the  $X_i$ ,

$$\mathbb{E} \left[ e^{\vartheta S_N} \right] = \prod_{i=1}^N M_X \left( \omega_i(N) (e^\vartheta - 1) \right). \quad (28)$$

Bearing in mind (27), we define

$$\vartheta^* := \arg \sup_{\vartheta} \left\{ \vartheta a - \int_0^1 \Lambda_X \left( \bar{F}(x) (e^\vartheta - 1) \right) dx \right\}.$$

The idea is now that we construct a measure  $\mathbb{Q}$ , under which the event of interest is not rare so that a central limit theorem applies. Concretely, we choose  $\mathbb{Q}$  to be an  $\vartheta^*$ -twisted version of the original measure such that  $S_N$  has moment generating function (*cf.* (28))

$$\mathbb{E}_{\mathbb{Q}} \left[ e^{\vartheta S_N} \right] = \prod_{i=1}^N M_X \left( \omega_i(N) (e^{\vartheta + \vartheta^*} - 1) \right) \Big/ \prod_{i=1}^N M_X \left( \omega_i(N) (e^{\vartheta^*} - 1) \right). \quad (29)$$

As a consequence,  $q_N(a) = \mathbb{E}_{\mathbb{Q}} LI_N$ , with the indicator function  $I_N := 1_{\{S_N = Na\}}$  and the likelihood ratio

$$L := e^{-\vartheta^* S_N} \prod_{i=1}^N M_X \left( \omega_i(N) (e^{\vartheta^*} - 1) \right).$$

It thus follows that

$$q_N(a) = e^{-\vartheta^* Na} \left( \prod_{i=1}^N M_X \left( \omega_i(N)(e^{\vartheta^*} - 1) \right) \right) \mathbb{Q}(S_N = Na).$$

We now point out how to evaluate the middle factor in the previous display (i.e., the product), namely, we check that asymptotically this middle factor behaves as

$$\exp \left( N \int_0^1 \Lambda_X(\tau \bar{F}(x)) dx \right), \quad (30)$$

with  $\tau := e^{\vartheta^*} - 1$ . The logarithm of the middle factor is

$$\sum_{i=1}^N \Lambda_X(\tau \omega_i(N)) = \sum_{i=1}^N \Lambda_X \left( \tau N \int_{(i-1)/N}^{i/N} \bar{F}(x) dx \right),$$

where, by a Taylor expansion of  $\bar{F}$ ,

$$N \int_{(i-1)/N}^{i/N} \bar{F}(x) dx = \bar{F} \left( \frac{i-1}{N} \right) + \frac{1}{2N} \bar{F}' \left( \frac{i-1}{N} \right) + O \left( \frac{1}{N^2} \right).$$

As a consequence, from a Taylor expansion of  $\Lambda_X(\cdot)$  we have

$$\sum_{i=1}^N \Lambda_X(\tau \omega_i(N)) = \sum_{i=1}^N \Lambda_X \left( \tau \bar{F} \left( \frac{i-1}{N} \right) \right) + \frac{\tau}{2N} \sum_{i=1}^N \bar{F}' \left( \frac{i-1}{N} \right) \Lambda'_X \left( \tau \bar{F} \left( \frac{i-1}{N} \right) \right) + O \left( \frac{1}{N} \right),$$

where, as  $N \rightarrow \infty$ ,

$$\begin{aligned} \frac{\tau}{2N} \sum_{i=1}^N \bar{F}' \left( \frac{i-1}{N} \right) \Lambda'_X \left( \tau \bar{F} \left( \frac{i-1}{N} \right) \right) &\rightarrow \frac{\tau}{2} \int_0^1 \bar{F}'(x) \Lambda'_X(\tau \bar{F}(x)) dx \\ &= \frac{1}{2} (\Lambda_X(\tau \bar{F}(1)) - \Lambda_X(\tau \bar{F}(0))), \end{aligned}$$

provided that  $\bar{F}(\cdot)$  is twice differentiable on  $[0, 1]$  (recognize the *left Riemann sum* approximation). Now recall the *trapezoidal rule* version of the Riemann sum approximation, that holds for any Riemann-integrable  $G(\cdot)$ :

$$\frac{1}{N} \sum_{i=1}^N G(i/N) = \int_0^1 G(x) dx + \frac{1}{2N} (G(1) - G(0)) + O \left( \frac{1}{N^2} \right).$$

Since  $\Lambda_X$  is Riemann integrable on  $[0, 1]$ , this can be applied to yield

$$\begin{aligned} N \int_0^1 \Lambda_X(\tau \bar{F}(x)) dx &= \sum_{i=1}^N \Lambda_X \left( \tau \bar{F} \left( \frac{i}{N} \right) \right) - \frac{1}{2} (\Lambda_X(\tau \bar{F}(1)) - \Lambda_X(\tau \bar{F}(0))) + O \left( \frac{1}{N} \right) \\ &= \sum_{i=1}^N \Lambda_X \left( \tau \bar{F} \left( \frac{i-1}{N} \right) \right) + \frac{1}{2} (\Lambda_X(\tau \bar{F}(1)) - \Lambda_X(\tau \bar{F}(0))) + O \left( \frac{1}{N} \right). \end{aligned}$$

We have thus arrived at

$$q_N(a) \sim e^{-\vartheta^* Na} \exp \left( N \int_0^1 \Lambda_X(\bar{F}(x)(e^{\vartheta^*} - 1)) dx \right) \mathbb{Q}(S_N = Na).$$

We are left to evaluate  $\mathbb{Q}(S_N = Na)$ . We do so by first proving the claim that, under  $\mathbb{Q}$ ,  $S_N$  obeys a central limit theorem: as  $N \rightarrow \infty$ ,

$$\frac{S_N - Na}{\sqrt{N}}$$

converges to a zero-mean Normal random variable. Recall from (29) that we have

$$\log \mathbb{E}_{\mathbb{Q}} e^{\vartheta S_N} = \sum_{i=1}^N \Lambda_X \left( \omega_i(N)(e^{\vartheta + \vartheta^*} - 1) \right) - \sum_{i=1}^N \Lambda_X \left( \omega_i(N)(e^{\vartheta^*} - 1) \right).$$

In order to establish that  $S_N$  satisfies the anticipated central limit theorem, we prove that  $\Psi_N(\vartheta) := \log \mathbb{E}_{\mathbb{Q}} e^{\vartheta S_N / \sqrt{N}} - \vartheta a \sqrt{N} \rightarrow \frac{1}{2} \sigma^2 \vartheta^2$ , for some  $\sigma^2 > 0$ . This is done as follows. Observe that we can write the logarithmic moment generating function  $\Psi_N(\vartheta)$  as

$$\sum_{i=1}^N \Lambda_X \left( \omega_i(N) \left( e^{\vartheta^*} - 1 + \left( e^{\vartheta^*} (e^{\vartheta/\sqrt{N}} - 1) \right) \right) \right) - \sum_{i=1}^N \Lambda_X \left( \omega_i(N) (e^{\vartheta^*} - 1) \right) - \vartheta a \sqrt{N}.$$

By applying a Taylor expansion to  $e^{\vartheta/\sqrt{N}} - 1$ , this can be written as (neglecting higher order terms)

$$\sum_{i=1}^N \Lambda_X \left( \omega_i(N) \left( e^{\vartheta^*} - 1 + \left( e^{\vartheta^*} \left( \frac{\vartheta}{\sqrt{N}} + \frac{\vartheta^2}{2N} \right) \right) \right) \right) - \sum_{i=1}^N \Lambda_X \left( \omega_i(N) (e^{\vartheta^*} - 1) \right) - \vartheta a \sqrt{N}.$$

This can be expanded to, up to terms that are  $o(1)$  as  $N \rightarrow \infty$ ,

$$\begin{aligned} & \sum_{i=1}^N \left[ \Lambda'_X \left( \omega_i(N) (e^{\vartheta^*} - 1) \right) \omega_i(N) e^{\vartheta^*} \left( \frac{\vartheta}{\sqrt{N}} + \frac{\vartheta^2}{2N} \right) \right. \\ & \left. + \frac{1}{2} \Lambda''_X \left( \omega_i(N) (e^{\vartheta^*} - 1) \right) \omega_i(N)^2 e^{2\vartheta^*} \frac{\vartheta^2}{N} \right] - \vartheta a \sqrt{N}. \end{aligned} \quad (31)$$

Now note that, similar to what we have seen before,

$$\frac{1}{N} \sum_{i=1}^N \Lambda'_X \left( \omega_i(N) (e^{\vartheta^*} - 1) \right) \omega_i(N) e^{\vartheta^*} = \int_0^1 \Lambda'_X(\bar{F}(x) (e^{\vartheta^*} - 1)) \bar{F}(x) e^{\vartheta^*} dx + O\left(\frac{1}{N}\right),$$

where the integral equals  $a$  by the definition of  $\vartheta^*$ . We conclude that (31) converges to  $\frac{1}{2} \sigma^2 \vartheta^2$  as  $N \rightarrow \infty$ , where the corresponding variance is given by

$$\begin{aligned} \sigma^2 & := \int_0^1 \Lambda'_X(\bar{F}(x) (e^{\vartheta^*} - 1)) \bar{F}(x) e^{\vartheta^*} dx + \int_0^1 \Lambda''_X(\bar{F}(x) (e^{\vartheta^*} - 1)) \bar{F}^2(x) e^{2\vartheta^*} dx \\ & = a + \int_0^1 \Lambda''_X(\bar{F}(x) (e^{\vartheta^*} - 1)) \bar{F}^2(x) e^{2\vartheta^*} dx. \end{aligned}$$

We have thus established that, under  $\mathbb{Q}$ ,  $S_N$  satisfies the claimed central limit theorem. It directly implies that, by applying the usual continuity correction idea,  $\mathbb{Q}(S_N = Na)$  behaves inversely proportionally to  $\sqrt{N}$  in the sense that

$$\sqrt{N} \mathbb{Q}(S_N = Na) \sim \sqrt{N} \mathbb{P} \left( \mathcal{N}(0, \sigma^2) \in \left( -\frac{1}{2\sqrt{N}}, \frac{1}{2\sqrt{N}} \right) \right) \rightarrow \frac{1}{\sqrt{2\pi}\sigma}.$$

Upon combining the above, we conclude that the following asymptotic relationship holds.

**Proposition 5.1.** *As  $N \rightarrow \infty$ , if  $\bar{F}(\cdot)$  is twice differentiable on  $[0, 1]$ ,*

$$q_N(a) \sim \tilde{q}_N(a) := e^{-\vartheta^* Na} \exp \left( N \int_0^1 \Lambda_X(\bar{F}(x) (e^{\vartheta^*} - 1)) dx \right) \frac{1}{\sqrt{2\pi N} \sigma}.$$

Similar to Remark 2, we can convert the asymptotics of  $q_N(a)$  into those of  $Q_N(a)$ . More precisely, it can be argued that  $Q_N(a)$  has the same asymptotics as  $q_N(a)$ , except that the expansion for  $q_N(a)$  should be divided by  $1 - e^{-\vartheta^*}$  (which is smaller than 1). Note also that for the case  $\bar{F}(\cdot) \equiv 1$  we indeed recover the expression that we provided in Section 2.3. Furthermore, it is easily verified that if  $\mathbb{P}(X_i = \lambda) = 1$  (so the arrival rates are deterministic), the approximation we obtained in Prop. 5.1 coincides with that of the transient distribution of an M/G/ $\infty$  queue. With  $\varrho(1) := \lambda \int_0^1 \bar{F}(x) dx$ , recall that the number of customers present at time 1 is Poisson with mean  $\varrho(1)$ . By applying Stirling's approximation, and using that  $\vartheta^* = \log(a/\varrho(1))$ ,

$$q_N(a) = (N\varrho(1))^{Na} e^{-N\varrho(1)} \frac{1}{(Na)!} \sim \left( \frac{\varrho(1)}{a} \right)^{Na} e^{N(a-\varrho(1))} \frac{1}{\sqrt{2\pi Na}} = \tilde{q}_N(a).$$

**5.2. Numerical example.** We consider the following numerical example, which illustrates how Prop. 5.1 can be useful in devising staffing rules with possible applications in cloud provisioning, call center staffing or the design of data centers. Per time slot of length 1 time unit (which we refer to as  $\Delta$ ) a new arrival rate is sampled from a given distribution with a mean such that on average  $\lambda$  clients arrive in the time slot of length  $\Delta$ . The service times have a fixed mean  $E$ .

Let us assume the system starts empty, say at 8 AM. Suppose we wish to determine an appropriate staffing rule for slot 100 (evidently, any other slot for which we wish to adapt staffing levels can be dealt with analogously). Then we choose  $N = 100$  (recall the way we normalized time), and after scaling we have  $\mathbb{E}[NX_i\Delta] = \lambda$  (as  $N\Delta = 1$ ). Suppose the service facility wishes to maintain a rather strict quality level; its objective is to choose the number of servers in slot 100 to be  $\lfloor Na \rfloor$  (or, alternatively,  $\lceil Na \rceil$ ), where  $a$  is the smallest number such that  $Q_N(a)$  drops below  $\varepsilon$ .

For the service times we consider the following three distributions:

- In the first place, we assume that the service times are exponential with mean service time  $E$ , that is,  $\bar{F}(x) = e^{-x/E}$ .
- A second choice is to assume that the service times are deterministically equal to  $E$ , that is we define  $\bar{F}(x) = \mathbb{1}\{x < E\}$ .
- A third choice is to assume that the service times have a Pareto(2) distribution with mean  $E$ , that is,  $\bar{F}(x) = (1 + x/E)^{-2}$ .

As indicated in the introduction, in practice arrival rates for modeling call centers are typically not constant over time, but may be fluctuating around some mean value [11]. We assume that arrival rates follow a Poisson distribution in Section 5.2.1. In Section 5.2.2 we consider discrete arrival rates alternating between two values (corresponding to busy and quiet periods), motivated by applications in cloud computing, where the workload of virtual machines exhibits such bursty behaviour [20].

**5.2.1. Poisson arrival rates.** In this example we take  $X_i \sim \text{Pois}(\lambda)$ . We then have

$$\Lambda_X(\vartheta) = \lambda(e^\vartheta - 1); \quad \Lambda'_X(\vartheta) = \Lambda''_X(\vartheta) = \lambda e^\vartheta.$$

To compute  $\vartheta^*$  and  $\sigma^2$ , we evaluate

$$\int_0^1 \Lambda_X(\bar{F}(x)(e^\vartheta - 1)) dx = \int_0^1 \lambda \left( \exp(\bar{F}(x)(e^\vartheta - 1)) - 1 \right) dx$$

and

$$\int_0^1 \lambda \exp(\bar{F}(x)(e^{\vartheta^*} - 1)) \bar{F}^2(x) e^{2\vartheta^*} dx$$

by numerical integration. Inserting the resulting quantities into the formula provided in Prop. 5.1, we can compute the approximation  $\tilde{Q}_N(a)$  as  $\tilde{q}_N(a)(1 - e^{-\vartheta^*})^{-1}$  for various  $a$ . Consider Fig. 3 for a comparison of  $\tilde{Q}_N(a)$  with the corresponding estimators  $\hat{Q}_N(a)$  that are obtained by crude Monte Carlo estimation of the probability  $Q_N(a)$  as defined in (2).

We then proceed to find the value of  $a$ , denoted by  $a(\varepsilon)$ , for which we have  $|\tilde{Q}_N(a) - \varepsilon| < 10^{-9}$  using a bisection method. The results are displayed in Table 1; together with  $M_1$ , the expected number of customers present at time 1; the Monte Carlo estimates  $\hat{Q}_N(a(\varepsilon))$ ; and the values of  $\tilde{Q}_N(\underline{a})$  and  $\tilde{Q}_N(\bar{a})$ , where  $\underline{a}$  and  $\bar{a}$  are such that the number of servers is integer-valued:  $N\underline{a} = \lfloor Na(\varepsilon) \rfloor$  and  $N\bar{a} = \lceil Na(\varepsilon) \rceil$ . Surprisingly, the results we obtain for  $a(\varepsilon)$  and  $M_1$  suggest that the number of servers required decreases as the variability of the service distribution increases: a relatively small number of servers suffices when service times are Pareto(2), whereas a large number of servers is required for deterministic service times.

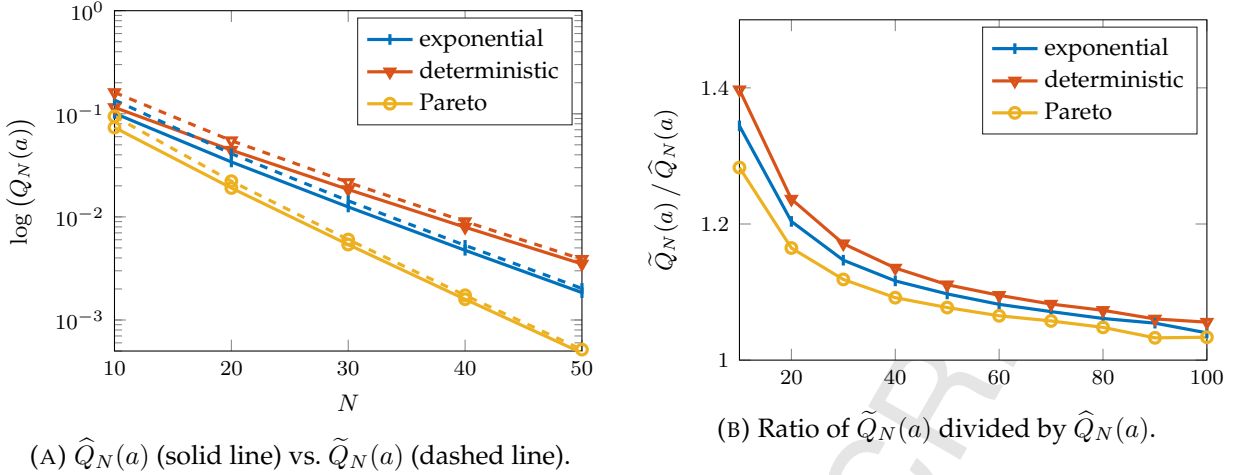
(A)  $\hat{Q}_N(a)$  (solid line) vs.  $\tilde{Q}_N(a)$  (dashed line).(B) Ratio of  $\tilde{Q}_N(a)$  divided by  $\hat{Q}_N(a)$ .

FIGURE 3. Comparison of crude Monte Carlo estimators  $\hat{Q}_N(a)$  and the approximation  $\tilde{Q}_N(a)$  as provided in Prop. 5.1. Parameters are chosen as  $a = 0.2$ ,  $\lambda = 0.1$ ,  $E = 1$ .

TABLE 1. Values of  $a(\varepsilon)$  needed to achieve  $|\tilde{Q}_N(a(\varepsilon)) - \varepsilon| < 10^{-9}$  with  $N = 100$ , expected arrival rate  $\lambda = 2$ , and mean service time  $E$ . The Monte Carlo estimates  $\hat{Q}_N(a(\varepsilon))$  are also provided (based on  $10^9$  runs) together with CI, the width of the standard normal 95% confidence interval, as well as the values of the approximation  $\tilde{Q}_N(a)$  with  $\underline{a}$  ( $\bar{a}$ , respectively) such that  $N\underline{a} = \lfloor Na \rfloor$  ( $N\bar{a} = \lceil Na \rceil$ , respectively). The inferred number of servers is  $N\bar{a}$ , which should be larger than the expected number of customers  $M_1$  at time 1.

$F$	$\varepsilon$	$E$	$a(\varepsilon)$	$N\bar{a}$	$\lceil M_1 \rceil$	$\frac{1}{\varepsilon} \left[ \hat{Q}_N(a(\varepsilon)) \pm \frac{\text{CI}}{2} \right]$	$\frac{1}{\varepsilon} \left( \tilde{Q}_N(\underline{a}), \tilde{Q}_N(\bar{a}) \right)$
Exponential	$10^{-3}$	0.05	0.2516	26	10	$0.5568 \pm 0.0015$	(1.1009, 0.6033)
		0.5	1.2602	127	87	$0.7215 \pm 0.0017$	(1.0053, 0.7802)
		1	1.7537	176	127	$0.8099 \pm 0.0018$	(1.0784, 0.8780)
	$10^{-4}$	0.05	0.2885	29	10	$0.8436 \pm 0.0057$	(1.7277, 0.9039)
		0.5	1.3460	135	87	$0.8380 \pm 0.0057$	(1.1858, 0.8921)
		1	1.8587	186	127	$0.9122 \pm 0.0059$	(1.2238, 0.9702)
Deterministic	$10^{-3}$	0.05	0.2782	28	10	$0.8382 \pm 0.0018$	(1.4983, 0.9133)
		0.5	1.4809	149	100	$0.7645 \pm 0.0017$	(1.0185, 0.8279)
		1	2.6636	267	200	$0.8353 \pm 0.0018$	(1.0565, 0.9070)
	$10^{-4}$	0.05	0.3223	33	10	$0.6146 \pm 0.0049$	(1.1319, 0.6547)
		0.5	1.5857	159	100	$0.8463 \pm 0.0057$	(1.1407, 0.9036)
		1	2.8048	281	200	$0.8590 \pm 0.0057$	(1.0869, 0.9136)
Pareto(2)	$10^{-3}$	0.05	0.2350	24	10	$0.6630 \pm 0.0016$	(1.3845, 0.7229)
		0.5	1.0074	101	67	$0.8559 \pm 0.0018$	(1.2375, 0.9268)
		1	1.4250	143	100	$0.8224 \pm 0.0018$	(1.1252, 0.8894)
	$10^{-4}$	0.05	0.2688	27	10	$0.5721 \pm 0.0057$	(1.8616, 0.9194)
		0.5	1.0818	109	67	$0.7223 \pm 0.0053$	(1.0613, 0.7633)
		1	1.5167	152	100	$0.8642 \pm 0.0058$	(1.1959, 0.9164)

At first sight, this outcome may seem counter-intuitive: one would perhaps have expected that unsteady service times would imply that more servers are needed. It is, however, easy to see that this conclusion is not necessarily valid (and in fact false for the example at hand). While it is true that customers arriving at an early slot can be served in time by the ‘deterministic servers’ with probability 1, customers arriving in later slots can never complete their service in time. For ‘random servers’ instead, customers arriving early may not finish their service in time but on the other hand customers arriving late still have a chance of completing their service.

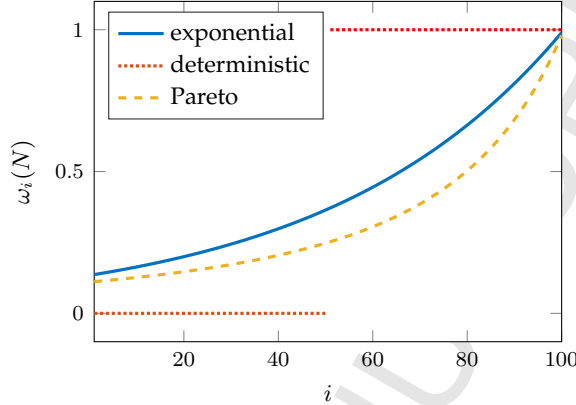


FIGURE 4. Values of  $\omega_i(N)$ , the probability that a customer arriving in the  $i$ -th time slot is still in the system at time 1, where  $N = 100$ ,  $E = 0.5$ .

In our example, this is reflected in the values of  $\omega_i(N)$ : bearing in mind that we fixed the value of the mean service time  $E$ , the arrival rates in the system with Pareto service times are thinned less in early slots but more in later slots, compared to deterministic service times (see Fig. 4). That Pareto service times turn out to be better is a result of the fact that the Pareto service times are smaller than  $E$  with large probability, and hence the regime in which the Pareto servers outperform the deterministic servers matters more than the regime in which the deterministic servers are better. Formally, we have that the sum of  $\omega_i(N)$  is smallest in the case of Pareto servers, and hence,  $S_N = \sum_{i=1}^N \text{Pois}(X_i \omega_i(N))$  has the smallest exceedance probability in that case.

To further investigate this issue, it is instructive to compute the variance of the steady-state number of clients in the system for the three models for the infinite-server queue. To this end, we can use the formulae that were provided in [9, Eqn. (2.31)] for the special case of exponential service times, noting that they can analogously be derived for more general service time distributions. We obtain

$$\text{Var} \left( \sum_{i=1}^N Z_i \right) = \text{Var} X \sum_{i=1}^N \omega_i^2(N) + \mathbb{E}X \sum_{i=1}^N \omega_i(N).$$

In case the service times are typically considerably smaller than 1, this behaves as

$$N \text{Var} X \int_0^1 \bar{F}^2(x) dx + N \mathbb{E}X \int_0^1 \bar{F}(x) dx \approx N \text{Var} X \int_0^\infty \bar{F}^2(x) dx + N \mathbb{E}X \int_0^\infty \bar{F}(x) dx. \quad (32)$$

In this decomposition the second part can be interpreted as the variance that one would obtain if the arrival process were Poisson with a constant (non-random) rate  $\mathbb{E}X$ , whereas the first part is the contribution due to overdispersion. In our example, because  $X$  has a Poisson distribution,  $\mathbb{E}X = \lambda = \text{Var} X$ .

The mean number in the system in stationarity is

$$M_\infty := N \mathbb{E}X \int_0^\infty \bar{F}(x) dx = N \lambda E, \quad (33)$$



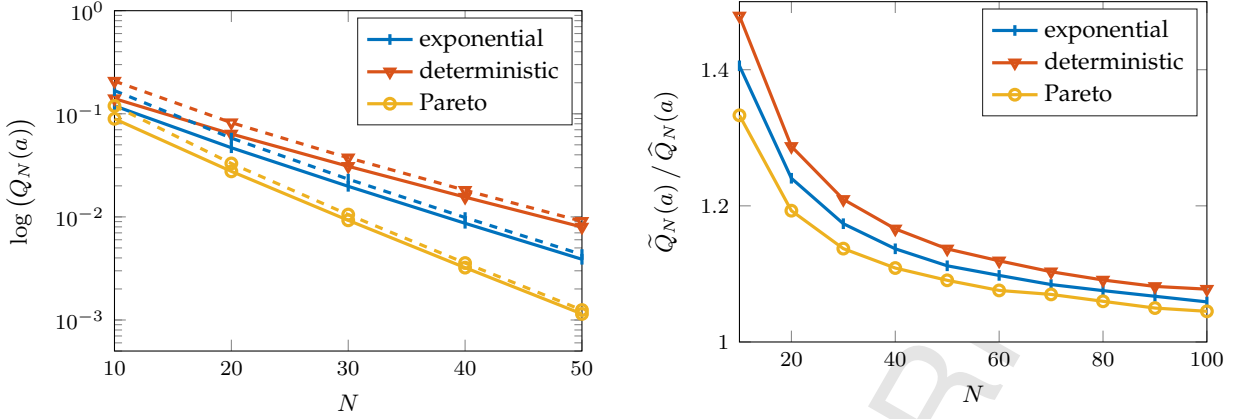
(A)  $\hat{Q}_N(a)$  (solid line) vs.  $\tilde{Q}_N(a)$  (dashed line).(B) Ratio of  $\tilde{Q}_N(a)$  divided by  $\hat{Q}_N(a)$ .

FIGURE 5. Comparison of crude Monte Carlo estimators  $\hat{Q}_N(a)$  and the approximation  $\tilde{Q}_N(a)$  as provided in Prop. 5.1. Parameters are chosen as  $E = 0.5$ ,  $p = 0.75$ ,  $\lambda_1 = 1$  and  $\lambda_2 = 5$ , with  $a = 1.6$  for deterministic,  $a = 1.4$  for exponential and  $a = 1.2$  for Pareto service times.

which shows that this term depends on the service-time distribution only through its mean  $E$ . It thus follows that the second term in the right-hand side of (32) equals  $N \lambda E$ . We now consider the first (overdispersion-related) term. In the exponential case,

$$\int_0^\infty \bar{F}^2(x) dx = \int_0^\infty e^{-2x/E} dx = \frac{E}{2};$$

in the deterministic case,

$$\int_0^\infty \bar{F}^2(x) dx = \int_0^E dx = E;$$

and in the Pareto(2) case,

$$\int_0^\infty \bar{F}^2(x) dx = \int_0^\infty (1 + x/E)^{-4} dx = \frac{E}{3}.$$

These computations confirm that the variability in the number of clients in the system is highest when the service times are deterministic, and lowest when they are Pareto(2). This entails that – as we saw from the results in Table 1 – if there is overdispersion (i.e.,  $\text{Var } X > 0$ ), the Pareto(2) case allows for a relatively conservative staffing policy, whereas in the deterministic case comparatively many servers are required.

The table also shows that the required number of servers given by  $N\bar{a}$  is, for obvious reasons, larger than  $M_1$ , the expected number of customers at time 1. At the same time,  $N\bar{a}$  can be substantially lower than the expected number of customers in the system in stationarity (i.e.,  $M_\infty$ , as defined in (33)), due to the fact that the system has not necessarily reached stationarity at time  $t = 1$  (recall that the system starts empty at time 0).

**5.2.2. Bursty arrival rate parameters.** In a second example we assume that the arrivals are Poisson and usually occur with a certain rate  $\lambda_1$ , but occasionally occur with some larger rate  $\lambda_2$  (corresponding to peak times in the network). Queueing networks with such ‘bursty’ arrival behaviour are of interest in the context of cloud computing, see for example [14, 20].

Specifically, we assume that  $\mathbb{P}(X_i = \lambda_1) = p$  and  $\mathbb{P}(X_i = \lambda_2) = 1 - p =: \bar{p}$ , where  $p$  is typically substantially larger than  $\frac{1}{2}$ . A routine calculation shows that

$$\Lambda_X(\vartheta) = \log\left(pe^{\vartheta\lambda_1} + \bar{p}e^{\vartheta\lambda_2}\right), \quad \Lambda'_X(\vartheta) = \frac{\lambda_1 p e^{\vartheta\lambda_1} + \lambda_2 \bar{p} e^{\vartheta\lambda_2}}{pe^{\vartheta\lambda_1} + \bar{p}e^{\vartheta\lambda_2}}, \quad \Lambda''_X(\vartheta) = \frac{p\bar{p}(\lambda_1 - \lambda_2)^2 e^{\vartheta(\lambda_1 + \lambda_2)}}{(pe^{\vartheta\lambda_1} + \bar{p}e^{\vartheta\lambda_2})^2}.$$

TABLE 2. Parameters are chosen as in Table 1, with arrival rate parameters  $p = 0.75$ ,  $\lambda_1 = 1$  and  $\lambda_2 = 5$  (so that the expected arrival rate is 2).

$F$	$\varepsilon$	$E$	$a(\varepsilon)$	$N\bar{a}$	$[M_1]$	$[\widehat{Q}_N(a(\varepsilon)) \pm \frac{CI}{2}] / \varepsilon$	$(\widetilde{Q}_N(\underline{a}), \widetilde{Q}_N(\bar{a})) / \varepsilon$
Exponential	$10^{-3}$	0.05	0.2662	27	10	$0.7501 \pm 0.0017$	(1.4061, 0.8115)
		0.5	1.2991	130	87	$0.9002 \pm 0.0019$	(1.2266, 0.9787)
		1	1.8061	181	127	$0.8576 \pm 0.0018$	(1.1182, 0.9307)
	$10^{-4}$	0.05	0.3056	31	10	$0.7199 \pm 0.0053$	(1.4107, 0.7615)
		0.5	1.3942	140	87	$0.8089 \pm 0.0056$	(1.1124, 0.8601)
		1	1.9234	193	127	$0.8230 \pm 0.0056$	(1.0742, 0.8717)
Deterministic	$10^{-3}$	0.05	0.3012	31	10	$0.6173 \pm 0.0015$	(1.0539, 0.6640)
		0.5	1.5438	155	100	$0.8215 \pm 0.0018$	(1.0708, 0.8934)
		1	2.7487	275	200	$0.9035 \pm 0.0019$	(1.1232, 0.9827)
	$10^{-4}$	0.05	0.3484	35	10	$0.8783 \pm 0.0058$	(1.5388, 0.9209)
		0.5	1.6632	167	100	$0.8187 \pm 0.0056$	(1.0669, 0.8690)
		1	2.9094	291	200	$0.9316 \pm 0.0060$	(1.1532, 0.9905)
Pareto(2)	$10^{-3}$	0.05	0.2461	25	10	$0.7264 \pm 0.0017$	(1.4490, 0.7888)
		0.5	1.0381	104	67	$0.8755 \pm 0.0018$	(1.2856, 0.7069)
		1	1.4671	147	100	$0.8651 \pm 0.0018$	(1.1606, 0.9393)
	$10^{-4}$	0.05	0.2817	29	10	$0.5315 \pm 0.0045$	(1.1255, 0.5649)
		0.5	1.1200	113	67	$0.6948 \pm 0.0052$	(1.0002, 0.7408)
		1	1.5688	157	100	$0.9138 \pm 0.0059$	(1.2335, 0.9709)

As before, we evaluate the approximation provided in Prop. 5.1 numerically. The obtained approximations and the corresponding Monte Carlo estimates are depicted in Fig. 5. The counterpart to Table 1 is Table 2, where the parameters are chosen as in Section 5.2.1 (we put  $\lambda_1 = 1$ ,  $\lambda_2 = 5$  and  $p = 0.75$  so that the mean arrival rate is 2 as before). Compared to the previous example, it seems that here the required number of servers is overall somewhat larger due to the greater variance of the  $X_i$ . The ordering of the service time distributions in terms of the required number of servers remains the same as before: the queuing system with deterministic service times requires the largest number of servers.

## 6. CONCLUSION

In this paper we considered an infinite-server queue with doubly stochastic Poisson arrivals, where the arrival rate is resampled every  $N^{-\alpha}$  time units. Among the main contributions of the paper are exact (non-logarithmic, that is) asymptotic expressions for  $P_N(a)$ , namely the tail distribution of the number of arrivals at a given time (for  $\alpha > 3$  or  $\alpha < \frac{1}{3}$ ), as well as for  $Q_N(a)$ , for which we consider the tail probability of having more than  $Na$  customers in the system (for the case  $\alpha = 1$ ).

As we saw for the specific example of exponentially distributed arrival rates, the asymptotic expression for  $P_N(a)$  can have a rather intricate shape for  $\alpha \in [\frac{1}{2}, 2]$ . We do, however, believe that it

is possible to derive the asymptotics for the cases  $\alpha \in [\frac{1}{3}, \frac{1}{2})$  and  $\alpha \in (2, 3]$  by using more precise bounds based on the Berry-Esseen inequality.

In numerical examples we showed how the approximation for  $Q_N(a)$  can be useful when determining the required number of servers such that at a specific time  $t$  (e.g. a certain time of the day) a specific performance target is met. This staffing rule could be extended to one that achieves the desired performance level during an extended period of time, rather than at a single time point. We expect that this requires more refined techniques, since the staffing level at a certain point in time affects the number of customers present in the subsequent time interval. However, we feel that the procedure developed in this paper may serve as a reasonably accurate proxy.

Finally, we believe that it is possible to extend the results of the paper by relaxing the assumption that the arrival rates are independent and identically distributed. Instead, one could consider the situation in which the arrival rates in subsequent time intervals depend on each other in a Markovian fashion. Another interesting topic relates to the infinite-server model in which the random rate of the arrival process changes continuously (rather than being redrawn periodically, and then being valid for the rest of the interval); in this context we could for instance consider a Coxian arrival process with a shot-noise rate [13].

#### ACKNOWLEDGMENTS AND AFFILIATIONS

The authors are with Korteweg-de Vries Institute for Mathematics, University of Amsterdam, Science Park 904, 1098 XH Amsterdam, the Netherlands. J. Kuhn is also with The University of Queensland, St Lucia, Queensland, Australia, and is supported by Australian Research Council (ARC) grant DP130100156. M. Mandjes is also with EURANDOM, Eindhoven University of Technology, Eindhoven, the Netherlands, and Amsterdam Business School, Faculty of Economics and Business, University of Amsterdam, Amsterdam, the Netherlands. The research of M. Heemskerk and M. Mandjes is partly funded by NWO Gravitation project NETWORKS, grant number 024.002.003.

EMAIL. {j.m.a.heemskerk|j.kuhn|m.r.h.mandjes}@uva.nl.

#### REFERENCES

- [1] R. R. Bahadur and R. R. Rao. On deviations of the sample mean. *Annals of Mathematical Statistics*, 31(4):1015–1027, 1960.
- [2] A. Bassamboo, R.S. Randhawa, and A. Zeevi. Capacity sizing under parameter uncertainty: safety staffing principles revisited. *Management Science*, 56(10):1668–1686, 2010.
- [3] J. Blom, M. Mandjes, and K. de Turck. Refined large deviations asymptotics for Markov-modulated infinite-server systems. *European Journal of Operational Research*; to appear, arXiv:1608.04250, 2016.
- [4] S. Borst, A. Mandelbaum, and M.I. Reiman. Dimensioning large call centers. *Operations Research*, 52(1):17–34, 2004.
- [5] D.R. Cox. Some statistical methods connected with series of events. *Journal of the Royal Statistical Society, Series B (Methodological)*, 17(2):129–164, 1955.
- [6] Mieke Defraeye and Inneke Van Nieuwenhuysse. Staffing and scheduling under nonstationary demand for service: A literature review. *Omega*, 58:4 – 25, 2016.
- [7] A. Dembo and O. Zeitouni. *Large Deviations Techniques and Applications*. Springer-Verlag, New York, 2 edition, 1998.
- [8] B. He, Y. Liu, and W. Whitt. Staffing a service system with non-Poisson nonstationary arrivals. *Probability in the Engineering and Information Sciences*, 30:593–621, 2016.
- [9] M. Heemskerk, J. van Leeuwen, and M. Mandjes. Scaling limits for infinite-server systems in a random environment. *Stochastic Systems*; to appear, arXiv:1602.00499, 2017.
- [10] O.B. Jennings, A. Mandelbaum, W.A. Massey, and W. Whitt. Server staffing to meet time-varying demand. *Management Science*, 42(10):1383–1394, 1996.
- [11] G. Jongbloed and G. Koole. Managing uncertainty in call centers using Poisson mixtures. *Applied Stochastic Models in Business and Industry*, 17(4):307–318, 2001.
- [12] S. Kim and W. Whitt. Are call center and hospital arrivals well modeled by nonhomogeneous Poisson processes? *Manufacturing & Service Operations Management*, 16(3):464–480, 2014.

- [13] D. Koops, M. Mandjes, and O. Boxma. Networks of  $\cdot/G/\infty$  server queues with shot-noise-driven arrival intensities. *Queueing Systems; to appear*, arXiv:1608.04924, 2017.
- [14] B. Patch and T. Taimre. Transient provisioning for cloud computing platforms. *Submitted*, arXiv:1612.01845, 2016.
- [15] V. V. Petrov. *Sums of Independent Random Variables*. Springer Verlag, New York, 1975.
- [16] J. S. Sadowsky and J. A. Bucklew. On large deviations theory and asymptotically efficient Monte Carlo estimation. *IEEE Transactions on Information Theory*, 36(3):579–588, 1990.
- [17] J.S.H. van Leeuwen, B.W.J. Mathijsen, and F. Sloothak. Cloud provisioning in the QED regime. In *Proceedings of the 9th EAI International Conference on Performance Evaluation Methodologies and Tools*, pages 180–187, 2016.
- [18] W. Whitt. Dynamic staffing in a telephone call center aiming to immediately answer all calls. *Operations Research Letters*, 24(5):205–212, 1999.
- [19] W. Whitt, L. V. Green, and P. J. Kolesar. Coping with time-varying demand when setting staffing requirements for a service system. *Production and Operations Management*, 16(1):13–39, 2007.
- [20] S. Zhang, Z. Qian, Z. Luo, J. Wu, and S. Lu. Burstiness-aware resource reservation for server consolidation in computing clouds. *IEEE Transactions on Parallel and Distributed Systems*, 27(4):964–977, 2016.