

# MAD dispersion measure makes extremal queue analysis simple

Wouter van Eekelen

Department of Econometrics and Operations Research, Tilburg University, w.j.e.c.vaneekelen@tilburguniversity.edu

Dick den Hertog

Amsterdam Business School, University of Amsterdam, d.denhertog@uva.nl

Johan S.H. van Leeuwen

Department of Econometrics and Operations Research, Tilburg University, and Department of Mathematics and Computer Science, Eindhoven University of Technology, j.s.h.vanleeuwen@tilburguniversity.edu

A notorious problem in queueing theory is to compute the worst possible performance of the GI/G/1 queue under mean-dispersion constraints for the interarrival and service time distributions. We address this extremal queue problem by measuring dispersion in terms of Mean Absolute Deviation (MAD) instead of variance, making available recently developed techniques from Distributionally Robust Optimization (DRO). Combined with classical random walk theory, we obtain explicit expressions for the extremal interarrival time and service time distributions, and hence the best possible upper bounds, for all moments of the waiting time. We also apply the DRO techniques to obtain tight lower bounds that together with the upper bounds provide robust performance intervals. We show that all bounds are computationally tractable and remain sharp, also when the mean and MAD are not known precisely, but estimated based on available data instead.

*Key words:* extremal queue problem, GI/G/1 queue, random walk theory, tight bounds, distributionally robust optimization

---

## 1. Introduction

Queueing theory exists for more than a century with throughout a central role for the GI/G/1 queue with i.i.d. interarrival times  $\{U_n\}$  distributed as  $U$  and i.i.d. service times  $\{V_n\}$  distributed as  $V$ . The waiting times in the GI/G/1 queue can be expressed as the maxima of a random walk with step size  $X = V - U$ , the subject of an enormous literature: Chung (2001), Feller (1971), Asmussen (2003). For all moments of the maxima (i.e., waiting times), general expressions are available that involve convolutions of the distribution of  $X$ . To use these general expressions, one thus needs to specify the precise distribution of  $X$ , and in the case of the GI/G/1 queue the distributions of both  $U$  and  $V$ .

Special cases of the GI/G/1 queue can be studied with dedicated techniques for Markov chains. For instance, the M/G/1 queue with Poisson arrivals and the GI/M/1 queue with exponential services have explicit solutions that are more insightful than the general random walk results: Asmussen (2003), Cohen (1982). Another large, somewhat opposite branch of queueing theory

concerns finding approximations and bounds. For the steady-state waiting time  $W$  in the GI/G/1 queue, the arguably most famous upper bound for  $\mathbb{E}[W]$  was obtained by Kingman (1962) in terms of the first two moments of both  $U$  and  $V$ . While Kingman’s bound is sharp in situations of heavy traffic, when  $\mathbb{E}[U]/\mathbb{E}[V]$  approaches 1, it leaves room for improvement for all other values of  $\mathbb{E}[U]/\mathbb{E}[V]$ .

In search for that sharpest possible (tight) upper bound under the first two moments constraints, foundational work was done by Rolski (1972), Eckberg Jr (1977), and Whitt (1984) in the context of the GI/M/1 queue. Whitt (1984) considered the GI/M/1 queue with given mean and variance of  $U$ , and showed that  $\mathbb{E}[W]$  is maximized when the interarrivals follow a specific two-point distribution. It also led to the conjecture that the overall worst case behavior (in terms of  $\mathbb{E}[W]$ ) would be caused by two-point distributions, for both  $U$  and  $V$ . That conjecture was proved invalid by counterexamples in Whitt (1984) when fixing either  $U$  or  $V$ , but the conjecture remained standing for the case when both  $U$  and  $V$  are unspecified, except for their first two moments. After that it remained silent for a while, until Chen and Whitt (2019) showed recently, for distributions with finite support, that the extremal distributions of  $U$  and  $V$  both have supports on at most three points. While existence is thus proved, the exact form of the extremal three-or-fewer-points distributions can only be determined numerically, as the solution of a hard non-convex nonlinear optimization problem. Extensive numerical experiments led Chen and Whitt to conjecture that the worst case is formed by two-point distributions for both  $U$  and  $V$ , in line with the conjecture postulated several decades ago. Finding the extremal queue for given mean-variance information is therefore one of the longest standing problems in the field. That problem remains open, also after publication of the present paper.

We do consider the same problem of finding the sharpest possible bounds for GI/G/1 queue performance measures, but take a radical turn by quantifying dispersion in terms of mean absolute deviation (MAD) instead of variance. That may appear a bold decision, because MAD is hardly used in queueing theory, or random walk theory for that matter. We can only speculate about the historical reasons for variance preference, but the random walk and GI/G/1 queue are intrinsically linked with i.i.d. sums of random variables, and variance then enters naturally (e.g., variance of the sum, central limit theorem). The variance and MAD, however, are equally adequate descriptors of dispersion, and are both easily calibrated on data using basic statistical estimators.

The MAD perspective offered in this paper departs from the variance-based formulations of the past (see Rolski (1972), Eckberg Jr (1977), Whitt (1984), Chen and Whitt (2019) and the references therein), and brings to bear the rich theory of robust optimization, in particular the rapidly expanding theory of distributionally robust optimization (DRO). The exact expressions for the random walk maxima form a crucial ingredient for our proof methodology. These expressions

are convex functions of the driving random variables, a prerequisite for the mean-MAD approach. Indeed, recent advances in DRO, see Postek et al. (2018), show that knowledge on the support, mean and MAD can lead to closed-form expressions for stochastic quantities such as the minimum and maximum expectation of a convex function.

Using the MAD instead of the variance as dispersion measure has several important advantages for, e.g., analyzing the waiting times in GI/G/1 queues. First, not only simple explicit expressions for the worst-case distributions can be obtained, but also for the best-case ones. Hence, a sharp upper bound and a sharp lower bound for the expected waiting time can be obtained. Second, our approach is for i.i.d. sums of random variables, while existing DRO approaches have to tolerate possible dependence structures between the random variables. Third, our approach is suitable for analyzing both transient behavior and the steady state. Fourth, because of its computational tractability our approach can also be extended to many optimization variants.

The contributions of this paper can be summarized as follows:

1. We suggest to use MAD instead of variance, and obtain by concise mathematical proof the worst-case three-point distribution for a rich class of extremal problems. This proof for MAD gives insight into why the traditional moment constraints, although a popular choice, may not necessarily yield tractable counterparts.
2. We leverage this result to obtain tight upper and lower bounds for performance measures, including transient and steady-state queue length moments. Under mean-MAD constraints, these bounds are the sharpest possible (and thus cannot be improved). The mean-MAD approach in this paper is a new quantitative method applicable to random walks, queues and related stochastic processes. This generic approach is a computationally tractable way to analyze key performance measures of such processes.
3. We present guidelines that describe how to compute the novel tight bounds efficiently. Moreover, we demonstrate our approach when the mean and MAD are not known precisely and need to be estimated from data. Also in these more realistic settings, the bounds remain sharp.

**Outline.** The remainder of the paper is organized as follows. Section 2 presents the MAD perspective. Section 3 discusses methods to obtain upper and lower bounds for both best and worst-case performance. Section 4 presents a full solution of the extremal queue problem with mean-MAD constraints, and draws a comparison with the traditional mean-variance setting. We conclude in Section 5, also mentioning possibilities for follow-up research.

**Notation.** Boldfaced characters represent vectors, and  $x_i$  denotes the  $i$ -th element of vector  $\mathbf{x}$ . For a random variable  $X$ , we use  $X \sim \mathbb{P} \in \mathcal{P}$  to say that  $X$  is a random variable with probability distribution  $\mathbb{P}$  from the set of probability distributions  $\mathcal{P}$ . We denote  $\mathbb{E}_{\mathbb{P}}[\cdot]$  as the expectation over the probability distribution  $\mathbb{P}$ . When we consider  $\mathbb{E}_{\mathbb{P}}[f(\mathbf{X})]$  with  $\mathbf{X} = (X_1, \dots, X_n)$ , it is tacitly assumed that  $f(\cdot)$  is a measurable function from  $\mathbb{R}^n$  to  $\mathbb{R}$ , and such that  $\mathbb{E}_{\mathbb{P}}[f(\mathbf{X})]$  exists.

## 2. Extremal random walk

Consider the partial sums  $S_n := X_1 + \dots + X_n$  ( $S_0 := 0$ ) of i.i.d. random variables  $X_1, X_2, \dots$  distributed as  $X$ . The random walk  $(S_n, n \geq 0)$  arises in many application domains, including queueing theory, inventory management and risk theory. If  $(S_n, n \geq 0)$  indeed models congestion, shortfall or capital position, large values of  $S_n$  are of particular interest, and it is natural to consider the maxima sequence  $M_n := \max\{S_0, S_1, \dots, S_n\}$ . The random walk and its maxima can be studied with mathematical techniques for sums of random variables, covered in many standard texts on probability theory, e.g., Asmussen (2003), Chung (2001), Cohen (1982), Feller (1971). For the distribution and moments of  $M_n$  there exist general formulas in terms of finitely many convolutions. However, applying these exact formula requires full specification of the distribution of  $X$ . This paper searches for the sharpest possible bounds on  $\mathbb{E}[M_n]$  and related quantities, when only information is available on the mean and dispersion of  $X$ . We now present such bounds when the partial information consists of the mean, range and MAD of  $X$ .

### 2.1. Extremal distribution

Notice that  $M_n$  can be expressed as  $h_n(X_1, \dots, X_n)$ , with

$$h_n(x_1, \dots, x_n) = \max\{0, x_1, \dots, x_1 + \dots + x_n\}, \quad (1)$$

and the expected maximum can be expressed as  $\mathbb{E}[M_n] = \mathbb{E}[h_n(\mathbf{X})]$  with  $\mathbf{X} = (X_1, \dots, X_n)$ . For now assume that  $X_1, \dots, X_n$  are independent, but that each  $X_i$  can have a different distribution. Assuming we only have partial information consisting of means and dispersion measures of the random variables  $X_1, \dots, X_n$ , the first question we ask and answer in this paper is: What *extremal* distributions of  $X_i$  result in the worst-case expected maxima? Extremal distributions have been studied in many contexts, and in the literature variance is predominantly used as the dispersion measure. Here we shall use the MAD. To describe all considered distributions we define an ambiguity set that consists of all distributions of componentwise independent  $\mathbf{X}$  with known supports, means, and MADs. The partial information for  $(X_1, \dots, X_n)$  consists of (i)  $X_i$  has support  $\text{supp}(X_i) = [a_i, b_i]$  with  $-\infty < a_i \leq b_i < \infty, i = 1, \dots, n$ , (ii)  $\mathbb{E}_{\mathbb{P}}(X_i) = \mu_i$  and (iii)  $\mathbb{E}_{\mathbb{P}}|X_i - \mu_i| = d_i$ . This defines the *ambiguity set*

$$\mathcal{P}_{(\mu, d)} = \{\mathbb{P} : \text{supp}(X_i) \subseteq [a_i, b_i], \mathbb{E}_{\mathbb{P}}(X_i) = \mu_i, \mathbb{E}_{\mathbb{P}}|X_i - \mu_i| = d_i, \forall i, X_i \perp\!\!\!\perp X_j, \forall i \neq j\}, \quad (2)$$

where  $X_i \perp\!\!\!\perp X_j, \forall i \neq j$ , denotes stochastic independence of the components  $X_1, \dots, X_n$ . In what follows,  $\mathbf{X}$  is a vector of random variables whose distribution  $\mathbb{P}$  belongs to the set  $\mathcal{P}_{(\mu, d)}$ .

As the title says, with MAD as dispersion measure, the extremal problem becomes simple. Observe that the function  $h_n$  is convex in the vector  $(x_1, \dots, x_n)$ . We can thus apply the general upper bound in Ben-Tal and Hochman (1972) on the expectation of a convex function of independent random variables with mean-MAD ambiguity, which gives the following result:

THEOREM 1. *The extremal distribution that solves*

$$\max_{\mathbb{P} \in \mathcal{P}(\mu, d)} \mathbb{E}_{\mathbb{P}}[h_n(\mathbf{X})] \quad (3)$$

*consists for each  $X_i$  of a three-point distribution with values  $\tau_1^{(i)} = a_i$ ,  $\tau_2^{(i)} = \mu_i$ ,  $\tau_3^{(i)} = b_i$  and probabilities*

$$p_1^{(i)} = \frac{d_i}{2(\mu_i - a_i)}, \quad p_2^{(i)} = 1 - \frac{d_i}{2(\mu_i - a_i)} - \frac{d_i}{2(b_i - \mu_i)}, \quad p_3^{(i)} = \frac{d_i}{2(b_i - \mu_i)}. \quad (4)$$

Ben-Tal and Hochman (1972) prove Theorem 1 (for general convex functions) by introducing a piecewise linear function on the interval  $[a, b]$  that intersects the convex function in  $a$ ,  $\mu$  and  $b$ , and then applying the classic Jensen bound to the subintervals  $[a, \mu]$  and  $[\mu, b]$ . In the next section, we give another proof of Theorem 1 that also gives insight into why using as dispersion measure MAD instead of variance makes the analysis so simple.

## 2.2. Novel primal-dual proof of Theorem 1

Our proof will crucially rely on the fact that the univariate case of Theorem 1 is tractable, and can be straightforwardly extended to the multivariate case. We thus start by considering some univariate measurable function  $f(x)$  (with the univariate function  $h_1(x_1)$  as an example) that has finite values on  $[a, b]$ , the support of the distribution  $\mathbb{P}(x)$ . Under mean-MAD ambiguity of one random variable  $X$  we thus need to solve

$$\begin{aligned} & \max_{\mathbb{P}(x) \geq 0} \int_x f(x) d\mathbb{P}(x) \\ \text{s.t.} \quad & \int_x |x - \mu| d\mathbb{P}(x) = d, \int_x x d\mathbb{P}(x) = \mu, \int_x d\mathbb{P}(x) = 1, \end{aligned} \quad (5)$$

a semi-infinite linear program (LP) with three equality constraints. A perhaps surprising, yet classical fact, is that the semi-infinite LP (5) can be reduced to an equivalent finite LP that yields the same optimal value. Indeed, the Richter-Rogosinski Theorem (e.g., Rogosinski (1958), Shapiro et al. (2009), Han et al. (2015)) states that there exists an extremal distribution for problem (5) with at most three support points. While finding these points in closed form is typically not possible (for general semi-infinite problems), we next show that this is possible for the problem at hand, by resorting to the dual problem and exploiting both the specific shape of the MAD constraint  $\int_x |x - \mu| d\mathbb{P}(x) = d$  and convexity of  $f$ .

Consider the dual of (5),

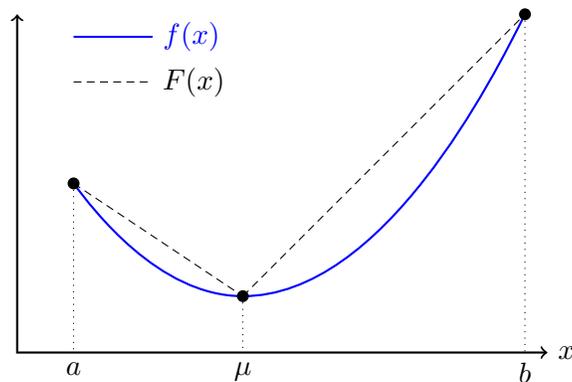
$$\begin{aligned} & \min_{\lambda_1, \lambda_2, \lambda_3} \lambda_1 d + \lambda_2 \mu + \lambda_3 \\ \text{s.t.} \quad & f(x) - \lambda_1 |x - \mu| - \lambda_2 x - \lambda_3 \leq 0, \quad \forall x \in [a, b]. \end{aligned} \quad (6)$$

Define  $F(x) = \lambda_1 |x - \mu| + \lambda_2 x + \lambda_3$ . Then the inequality in (6) can be written as  $f(x) \leq F(x)$ ,  $\forall x$ , i.e.  $F(x)$  majorizes  $f(x)$ . Note that  $F(x)$  has a ‘kink’ at  $x = \mu$ . Since the dual problem (6) has three

variables, the tightest majorant  $F(x)$  touches  $f(x)$  at three points:  $x = a, \mu$  and  $b$ , as illustrated in Figure 1. The optimal probabilities of (5) can now easily be obtained by solving the linear system resulting from the equations of (5). This is a linear system of three unknown probabilities and three equations, with a solution as stated in Theorem 1.

To deal with the multivariate case, we recursively apply the univariate result. Suppose we first apply this result to  $x_1$ , then the worst-case distribution is as in Theorem 1, independent of the values for  $x_2, \dots, x_n$ . Moreover, the worst-case expectation becomes a convex function in  $x_2, \dots, x_n$ , since the worst-case probabilities for  $x_1$  are nonnegative. Hence, we can apply the result above for the univariate case to  $x_2$ , etc. This completes the proof. Note that for multiperiod problems that involve multivariate optimization, such as the waiting time in the GI/G/1 queue, determining the extremal distribution for period  $n$  is unaffected by all previous periods.

To the best of our knowledge, our proof is the first to exploit the specific shape of the kink-majorant to find an analytic solution for the semi-infinite LP. While the dual problems are often solvable as semi-definite or second-order conic programs, analytic solutions as in our case are typically hard to attain, and require special structural properties of the LP's objective function or its constraints. Notice that in the univariate case, this proof method does not require convexity of  $f(x)$  and in fact could work for an arbitrary measurable function  $f$ . Convexity is needed, however, in the proof of Theorem 1 to extend the univariate case to the multivariate case. The proof method is of independent interest, and can for instance be applied to study the mean-MAD counterparts of the mean-variance analyses in e.g. Xin and Goldberg (2013), Natarajan and Zhou (2007), Perakis and Roels (2008), Natarajan et al. (2017), and Das et al. (2018).



**Figure 1** Some convex function  $f(x)$  and its piecewise linear majorant  $F(x)$ .

### 2.3. Why is MAD computationally easier than variance?

Now that we fully grasp why and how the proof of Theorem 1 relies on the specific structural properties of the mean-MAD constraints, and in particular the univariate result seamlessly passes into the multivariate counterpart, we can also explain why the comparable challenge with mean-variance constraints becomes much more difficult if not impossible. Observe that for the univariate case, the same proof argument works when  $\sigma^2$  is given instead of  $d$ , i.e., when  $|x - \mu|$  in (5) is replaced by  $(x - \mu)^2$ . Hence, irrespective of whether MAD or variance is used as dispersion measure, for determining the tight upper bound of  $f(x)$ , it suffices to consider distributions with support on at most three points. There is however a crucial complication when extending to the multivariate case.

To see this, observe that when  $\sigma^2$  is used as dispersion measure, the end points and kink point do *not* necessarily span the support of the extremal distribution. That is, upon replacing  $|x - \mu|$  with  $(x - \mu)^2$ , the tightest majorant  $F(x)$  does not necessarily touch  $f(x)$  in  $a$ ,  $b$  and  $\mu$ .

Hence, if the variance is used as dispersion measure, then the worst-case distribution depends on the function  $f(x)$ . This has severe consequences for the multivariate case, i.e., when we consider  $h_n(x_1, \dots, x_n)$ . In that case, the worst-case distribution depends on the values of  $x_2, \dots, x_n$ , and calculating (in closed form) the worst-case distribution as a function of  $x_2, \dots, x_n$  seems to be impossible. Moreover, even if we would be able to derive such a worst-case distribution, substituting this distribution in the worst-case expectation would result in an extremely difficult function in  $x_2, \dots, x_n$  that is likely non-convex, and hence applying the univariate result to  $x_2$  is no longer possible. Our duality proof thus reveals that the complicating feature of the mean-variance framework applied to multiperiod problems is the fact that the extremal distribution in period  $n$  is affected by all previous periods.

## 3. Sharpest possible bounds

A direct consequence of Theorem 1 is that the worst-case expectation of  $h_n(\mathbf{X})$  is obtained by enumerating over all  $3^n$  permutations of outcomes  $a_i, \mu_i, b_i$  of components  $X_i$ .

COROLLARY 1.

$$\max_{\mathbb{P} \in \mathcal{P}_{(\mu, d)}} \mathbb{E}_{\mathbb{P}}[h_n(\mathbf{X})] = \sum_{\alpha \in \{1, 2, 3\}^n} h_n(\tau_{\alpha_1}^{(1)}, \dots, \tau_{\alpha_n}^{(n)}) \prod_{i=1}^n p_{\alpha_i}^{(i)}. \quad (7)$$

Thus, under the partial information contained in  $\mathcal{P}_{(\mu, d)}$ , (7) is an upper bound on  $\mathbb{E}[M_n]$  that cannot be improved. We next specialize to the random walk setting with  $X_1, X_2, \dots$  independent and distributed as  $X$ , obtain representations for the tight upper bound that are computationally less cumbersome than (7), and extend to all moments of the all-time maximum (when  $n \rightarrow \infty$ ).

### 3.1. Random walk upper bounds

We recall that Spitzer (1956) used combinatorial arguments to establish for  $\mathbb{E}[M_n]$  the alternative expression (which strictly requires i.i.d. increments)

$$\mathbb{E}[M_n] = \sum_{k=1}^n \frac{1}{k} \mathbb{E}[S_k^+], \quad (8)$$

with  $x^+ = \max\{0, x\}$ . This can be written as  $\mathbb{E}[M_n] = \mathbb{E}[f_n(\mathbf{X})]$  with

$$f_n(x_1, \dots, x_n) = \sum_{k=1}^n \frac{1}{k} \max\{0, x_1 + \dots + x_k\}. \quad (9)$$

A first usage of Spitzer's formula (8) is a considerable improvement, in terms of computational complexity, of the tight bound for  $\mathbb{E}[M_n]$  in (7). To state the result and for later reference, let  $\Omega(\mu, d, a, b)$  denote a three-point distribution on the values  $\{a, \mu, b\}$  with probabilities

$$p_1 = \frac{d}{2(\mu - a)}, \quad p_2 = 1 - \frac{d}{2(\mu - a)} - \frac{d}{2(b - \mu)}, \quad p_3 = \frac{d}{2(b - \mu)}. \quad (10)$$

Let  $X_{(3)}$  denote the random variable with the extremal three-point distribution, identified in Theorem 1 for the special case when  $X_1, X_2, \dots$  are i.i.d., hence  $X_{(3)} \sim \Omega(\mu, d, a, b)$ .

COROLLARY 2.

$$\max_{\mathbb{P} \in \mathcal{P}(\mu, d)} \mathbb{E}_{\mathbb{P}}[f_n(\mathbf{X})] = \sum_{k=1}^n \frac{1}{k} \sum_{\sum_i k_i = k} \max\{0, k_1 a + k_2 \mu + k_3 b\} \cdot \frac{k!}{k_1! k_2! k_3!} p_1^{k_1} p_2^{k_2} p_3^{k_3}. \quad (11)$$

Note that for each fixed  $k$ , (11) contains a multinomial distribution with support set  $\{(k_1, k_2, k_3) \in \mathbb{N}^3 : k_1 + k_2 + k_3 = k\}$  with cardinality  $\binom{k+2}{2}$ . This implies that the sum over  $k$  in (11) is over roughly  $n^3$  terms, which is way better than the  $3^n$  terms in (7).

For  $\mathbb{E}[X] < 0$  the all-time maximum  $M := \lim_{n \rightarrow \infty} M_n$  is a proper random variable ( $M_n$  converges in distribution to  $M$ , which will be finite with probability one if  $\mathbb{E}[X] < 0$ ). Let  $c_m(M)$  denote the  $m$ -th cumulant of  $M$ . Recall that  $c_1(M)$  is the mean,  $c_2(M)$  is the variance, and  $c_3(M)$  is the central moment  $\mathbb{E}[(M - \mathbb{E}[M])^3]$ . From general random walk theory we know that (see e.g., Abate et al. (1993))

$$c_m(M) = \sum_{k=1}^{\infty} \frac{1}{k} \mathbb{E}[(S_k^+)^m]. \quad (12)$$

We can now prove results similar as for  $\mathbb{E}[M_n]$ , regarding the extremal distribution and tight upper bound.

**THEOREM 2.** *Consider the random walk with generic step size  $X$  contained in the ambiguity set  $\mathcal{P}(\mu, d)$ . The tight upper bounds for all cumulants  $c_m(M)$  of the all-time maximum  $M$  are the cumulants of the random walk with extremal step size  $X_{(3)}$ .*

*Proof.* Consider the function

$$f_n^m(x_1, \dots, x_n) = \sum_{k=1}^n \frac{1}{k} (\max\{0, x_1 + \dots + x_k\})^m, \quad (13)$$

which is convex in the vector  $(x_1, \dots, x_n)$ . Hence, for i.i.d. increments with generic  $X$ ,

$$\max_{\mathbb{P} \in \mathcal{P}(\mu, d)} \mathbb{E}_{\mathbb{P}}[f_n^m(\mathbf{X})] \quad (14)$$

is solved by the extremal random variable  $X_{(3)}$ . This gives the bound, with  $X_1^*, X_2^*, \dots$  i.i.d. as  $X_{(3)}$ ,

$$l_n := \sum_{k=1}^n \frac{1}{k} \mathbb{E}[(S_k^+)^m] \leq \mathbb{E}f_n^m(X_1^*, \dots, X_n^*) =: u_n. \quad (15)$$

The result follows by observing that the sequences  $\{l_n\}$  and  $\{u_n\}$  are both monotone, and converging to well-defined limits.  $\square$

We conclude that the extremal three-point distribution for  $\mathbb{E}[M_n]$  in Theorem 1 is also the extremal distribution for all cumulants of  $M$ . When calculating the associate tight upper bounds for  $c_m(M)$ , (12) shows that we are confronted with an infinite summation of increasingly complex summands. Here, another line of classical random walk theory can help, which transforms such infinite sums into complex contour integrals.

Consider the random walk with generic step size  $X$ . It is known that formal solutions of the distribution of  $M_n$  and  $M$  can be expressed in terms of complex contour integrals (see Abate et al. (1993), Janssen et al. (2015) for the algorithmic aspects of these contour integrals). Assume that  $\phi_X(s) = \mathbb{E}[e^{sX}]$  is analytic for complex  $s$  in the strip  $|\operatorname{Re}(s)| < \delta$  for some  $\delta > 0$ . A sufficient condition is that the moment generating function  $\phi_X(s)$  is finite in a neighborhood of the origin, and hence all moments of  $X$  exist. Then

$$\mathbb{E}[e^{-sM}] = \exp \left\{ \frac{-1}{2\pi i} \int_{\mathcal{C}} \frac{s}{u(s-u)} \log(1 - \phi_X(-u)) du \right\}, \quad (16)$$

where  $s$  is a complex number with  $\operatorname{Re}(s) \geq 0$ ,  $\mathcal{C}$  is a contour to the left of, and parallel to, the imaginary axis, and to the right of any singularities of  $\log(1 - \phi_X(-u))$  in the left half plane. From (16) contour integral expressions for the cumulants follow by differentiation:

$$c_m(M) = \frac{(-1)^m}{2\pi i} \int_{\mathcal{C}} \frac{\log(1 - \phi_X(-u))}{u^{m+1}} du. \quad (17)$$

Consider  $X = X_{(3)}$  with a three-point distribution on values  $\{a, b, c\}$  with probabilities  $p_a, p_b, p_c$  and moment generating function

$$\phi_{X_{(3)}}(s) = p_a e^{sa} + p_b e^{sb} + p_c e^{sc}. \quad (18)$$

Notice that all moments of  $X_{(3)}$  exist, and hence  $\phi_{X_{(3)}}(s)$  satisfies the assumption required for representation (16) to hold. Since  $X_{(3)}$  follows the extremal three-point distribution associated with the tight upper bounds for  $c_m(M)$ , we obtain the following result:

COROLLARY 3. Let  $\phi_{X_{(3)}}(s) := \mathbb{E}[e^{sX_{(3)}}] = p_1 e^{sa} + p_2 e^{s\mu} + p_3 e^{sb}$ . The tight upper bounds on  $c_m(M)$  identified in Theorem 2 are given by

$$\frac{(-1)^m}{2\pi i} \int_{\mathcal{C}} \frac{\log(1 - \phi_{X_{(3)}}(-u))}{u^{m+1}} du, \quad m = 1, 2, \dots, \quad (19)$$

where  $\mathcal{C}$  is a contour to the left of, and parallel to, the imaginary axis, and to the right of any singularities of  $\log(1 - \phi_{X_{(3)}}(-u))$  in the left half plane.

Observe that (19) bypasses the cumbersome calculations with convolutions in (12). In EC.3 we demonstrate that this is a numerically efficient way of computing the tight bounds.

### 3.2. Random walk lower bounds

The tight upper bounds correspond to worst-case scenarios. We next show how the same MAD approach can identify best-case scenarios and hence tight lower bounds. For each  $X_i$ , define a second ambiguity set, which is a subset of  $\mathcal{P}_{(\mu,d)}$ :

$$\mathcal{P}_{(\mu,d,\beta)} = \{\mathbb{P} : \mathbb{P} \in \mathcal{P}_{(\mu,d)}, \mathbb{P}(X_i \geq \mu_i) = \beta_i, \forall i\}. \quad (20)$$

Hence, for obtaining a lower bound, we include the additional information  $\mathbb{P}(X_i \geq \mu_i) = \beta_i$  in the ambiguity set. Now, instead of finding the worst-case distribution, we want to identify the best-case distribution and corresponding tight lower bound. The following result is a direct consequence of the general lower bound in Ben-Tal and Hochman (1972) on the expectation of a convex function of independent random variables with  $\mathcal{P}_{(\mu,d,\beta)}$  ambiguity. In Section EC.2 we present a novel proof using the primal-dual method developed earlier for proving Theorem 1.

THEOREM 3.

$$\min_{\mathbb{P} \in \mathcal{P}_{(\mu,d,\beta)}} \mathbb{E}_{\mathbb{P}}[h_n(\mathbf{X})] = \sum_{\alpha \in \{1,2\}^n} h_n(v_{\alpha_1}^{(1)}, \dots, v_{\alpha_n}^{(n)}) \prod_{i=1}^n q_{\alpha_i}^{(i)}, \quad (21)$$

where

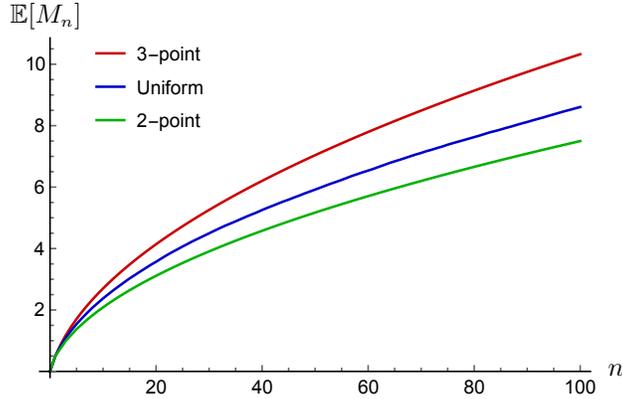
$$q_1^{(i)} = \beta_i, \quad q_2^{(i)} = 1 - \beta_i, \quad v_1^{(i)} = \mu_i + d_i/2\beta_i, \quad v_2^{(i)} = \mu_i - d_i/2(1 - \beta_i). \quad (22)$$

Again specialize to the i.i.d. setting, and denote by  $Y$  the random variable with two-point distribution on values

$$v_1 = \mu + \frac{d}{2\beta}, \quad v_2 = \mu - \frac{d}{2(1-\beta)},$$

with probabilities  $\beta$  and  $1 - \beta$ , respectively. Using similar reasonings as for the upper bound, we obtain for the tight lower bound for  $\mathbb{E}[M_n]$  an expression that sums over  $O(n^2)$  terms:

$$\sum_{k=1}^n \frac{1}{k} \sum_{k_1+k_2=k} \frac{k!}{k_1!k_2!} \beta^{k_1} (1-\beta)^{k_2} \max\{0, k_1 v_1 + k_2 v_2\}. \quad (23)$$



**Figure 2** Expected random walk maximum  $\mathbb{E}[M_n]$  for  $U(-b, b)$  and  $b=2$  distributed step sizes with MAD  $b/2$  (middle curve, obtained by simulation). The upper curve corresponds to the extremal three-point distribution within the ambiguity set with  $\mu=0$ ,  $d=b/2$  and range  $[-b, b]$ , and the lower curve is the bound (23) from the two-point distribution with  $\beta=1/2$ .

The tight lower bound for  $c_m(M)$  can be expressed in terms of the integral

$$\frac{(-1)^m}{2\pi i} \int_{\mathcal{C}} \frac{\log(1 - \phi_Y(-u))}{u^{m+1}} du, \quad (24)$$

where  $\phi_Y(s) = \beta e^{sv_1} + (1-\beta)e^{sv_2}$ ,  $\mathcal{C}$  is a contour to the left of, and parallel to, the imaginary axis, and to the right of any singularities of  $\log(1 - \phi_Y(-u))$  in the left half plane.

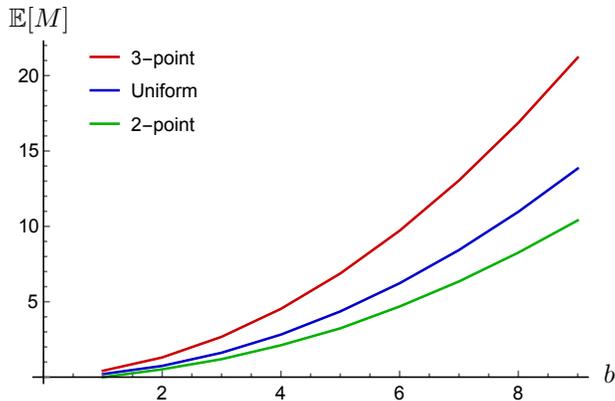
We illustrate the lower bound (21) (calculated using (23)) in Figure 2 for the random walk with step size  $X$  having a uniform distribution on  $[a, b]$ . Here we assume a specific distribution just for illustration purposes. The MAD of  $X$  can be shown to be  $(b-a)/4$ . In Figure 2 we choose  $b=-a=2$  so that  $\mu=0$  and  $d=1$ . Observe that upper and lower bound together provide a tight interval for all possible distributions in the ambiguity set  $\mathcal{P}_{(0,1,1/2)}$ .

Figure 3 shows the tight upper bound (19) and the lower bound (24) for  $\mathbb{E}[W]$  with ambiguity set with  $\mu=-1$ ,  $d=b/2$  and range  $[-b-2, b]$ . Observe that the bounds increase with the range and the MAD (which can be shown to hold in general). For a point of reference, we also plot the exact results for one member of the ambiguity set, when generic increment having a uniform distribution on  $[-b-2, b]$ .

### 3.3. Comparison with mean-variance ambiguity

As explained earlier, mean-variance ambiguity appears less computationally tractable than mean-MAD ambiguity. We now show how the key result for mean-MAD ambiguity, Theorem 1, can be used to obtain results for mean-variance ambiguity. Let  $\mathcal{P}_{(\mu, \sigma)}^*$  denote the ambiguity set that contains all distributions with known range, mean and variance, i.e.

$$\mathcal{P}_{(\mu, \sigma)}^* = \left\{ \mathbb{P} : \text{supp}(X_i) \subseteq [a_i, b_i], \mathbb{E}_{\mathbb{P}}(X_i) = \mu, \mathbb{E}_{\mathbb{P}}(X_i - \mu)^2 = \sigma^2, \forall i, X_i \perp X_j, \forall i \neq j \right\}. \quad (25)$$



**Figure 3** Expected all-time maximum  $\mathbb{E}[M]$  for  $U(-b-2, b)$  and  $b \in (1, 10)$  (middle curve, obtained by simulation). The upper curve corresponds to the extremal three-point distribution within the ambiguity set with  $\mu = -1$ ,  $d = (b+1)/2$  and range  $[-b-2, b]$ , and the lower curve is the bound (23) from the two-point distribution with  $\beta = 1/2$ .

PROPOSITION 1. Let  $d_{\min} = 2\sigma^2/(b-a)$  and  $d_{\max} = \sigma$ . Then,

$$\max_{\mathbb{P} \in \mathcal{P}(\mu, d_{\min})} \mathbb{E}_{\mathbb{P}}[h_n(\mathbf{X})] \leq \max_{\mathbb{P} \in \mathcal{P}(\mu, \sigma)} \mathbb{E}_{\mathbb{P}}[h_n(\mathbf{X})] \leq \max_{\mathbb{P} \in \mathcal{P}(\mu, d_{\max})} \mathbb{E}_{\mathbb{P}}[h_n(\mathbf{X})] \quad (26)$$

*Proof.* From Ben-Tal and Hochman (1985), we know that

$$\frac{2\sigma^2}{b-a} \leq d \leq \sigma.$$

Hence,  $\max_{\mathbb{P} \in \mathcal{P}(\mu, \sigma)} \mathbb{E}_{\mathbb{P}}[h_n(\mathbf{X})] = \max_{\mathbb{P} \in \mathcal{P}(\mu, d^*)} \mathbb{E}_{\mathbb{P}}[h_n(\mathbf{X})]$  for some  $d^* \in [2\sigma^2/(b-a), \sigma]$ . Since  $\max_{\mathbb{P} \in \mathcal{P}(\mu, d)} \mathbb{E}_{\mathbb{P}}[h_n(\mathbf{X})]$  is non-decreasing in  $d$ , see Postek et al. (2018), the result follows.  $\square$

Notice that Proposition 1 presents a way to delimit the upper bounds of all stationary cumulants  $c_m(M)$  and the transient mean  $\mathbb{E}[M_n]$  under mean-variance ambiguity. The mean-MAD bounds are specified in terms of specific three-point distributions.

We next show that the lower bound in Proposition 1 can lead to a result for infinite-support distributions. Set  $b = a + \xi(\mu - a)$  with  $\xi \geq 1$ , and observe that the lower bound in (26) comes with the extremal three-point distribution

$$X_{(3)}^{\xi} = \begin{cases} a & \text{w.p. } \frac{\sigma^2}{(\mu-a)^2\xi}, \\ \mu & \text{w.p. } 1 - \frac{\sigma^2}{(\mu-a)^2\xi} - \frac{\sigma^2}{(\mu-a)^2\xi(\xi-1)}, \\ a + \xi(\mu - a) & \text{w.p. } \frac{\sigma^2}{(\mu-a)^2\xi(\xi-1)}. \end{cases}$$

This distribution has mean  $\mu$  and variance  $\sigma^2$ , irrespective of the range  $[a, b]$ . We can thus let  $\xi$  grow to infinity to investigate what happens for infinite-support distributions.

For the expected all-time maximum, we can exploit an argument very similar to Chen and Whitt (2019), Theorem EC.3. A classic result from regenerative analysis says that the expected all-time

maximum is the expected sum of the random walk position over one cycle, denoted by  $\mathbb{E}[\text{integral}]$ , divided by the expected length of one cycle, i.e.  $\mathbb{E}[\text{cycle length}]$ . This cycle will consists of a period during which the queue remains empty, corresponding of consecutive (negative) steps of size  $a$  or  $\mu$ . As  $\xi$  increases, the three-point distribution places probabilities of order  $O(1/\xi^2)$  on  $a$  and  $a + \xi(\mu - a)$ , and the rest of the mass on point  $\mu$ . As  $\xi$  grows large, only rarely with probability  $O(1/\xi^2)$ , a large positive step occurs. The impact of the very large step of size  $a + \xi(\mu - a)$  is roughly the area of the triangle with height  $a + \xi(\mu - a)$  and width  $(a + \xi(\mu - a))/(-\mu)$ , and hence  $\mathbb{E}[\text{integral}] = (a + \xi(\mu - a))^2/(-2\mu) \sim (\xi(\mu - a))^2/(-2\mu)$  as  $\xi \rightarrow \infty$ . The cycle then consists of an empty period of expected length  $(1 - p_b)/p_b \sim (\xi(\mu - a))^2/\sigma^2$  and the positive period due to the large step of expected length  $(a + \xi(\mu - a))/(-\mu)$ , so that  $\mathbb{E}[\text{cycle length}] \sim (\xi(\mu - a))^2/\sigma^2$ , and the expected all-time maximum converges to  $\sigma^2/(-2\mu)$  as  $\xi \rightarrow \infty$ . Since this is a lower bound for  $\max_{\mathbb{P} \in \mathcal{P}_{(\mu, \sigma)}^*} \mathbb{E}[M]$ , we know that for the random walk with generic step size  $X$  it holds that  $\max_{\mathbb{P} \in \mathcal{P}_{(\mu, \sigma)}^*} \mathbb{E}[M] \geq \sigma^2/(-2\mu)$ . This lower bound matches Kingman's upper bound  $\mathbb{E}[M] \leq \sigma^2/(-2\mu)$ , which proves that Kingman's upper bound is tight. Tightness of Kingman's bound was already proven in Daley et al. (1992) by identifying a two-point distribution with mean  $\mu$ , variance  $\sigma^2$  such that  $\mathbb{E}[M]$  approaches the upper limit as one of the two points goes to infinity.

### 3.4. Degenerate behavior for infinite range

Compared to variance, MAD may be more appropriate in case of real-life empirical data that display non-Gaussian features and outliers. Indeed, unlike standard deviation, MAD does not require existence of second moments, and is not so much affected by large deviations from the mean. This feature, however, has major consequences when we let the range  $[a, b]$  grow large in which case conditioning on the MAD being  $d$  thus allows for distributions with relatively heavy tails. In particular, in the limit  $b \rightarrow \infty$ , this will lead to overly pessimistic scenarios as heavy-tailed distributions with infinite second moments would still have a finite  $d$  and hence be member of the ambiguity set. While for large but finite  $b$  a truly heavy-tailed distribution with infinite second moment is ruled out, the dispersion allowed by the ambiguity set might become too loose for practical purposes. An effective usage of the robust mean-MAD framework therefore requires a careful selection of the range, for which we now present some guidelines.

Observe that the variance of  $X_{(3)}$  is  $\frac{d}{2}(b - a)$ , the maximal variance for distributions in the ambiguity set  $\mathcal{P}_{(\mu, d)}$ . Hence, for fixed  $d$ , the variance becomes unbounded when  $b \rightarrow \infty$ . As a consequence, this results in fairly crude bounds:

PROPOSITION 2. *As  $b \rightarrow \infty$ , the bound  $\max_{\mathbb{P} \in \mathcal{P}_{(\mu, d)}} \mathbb{E}_{\mathbb{P}}[f_n(\mathbf{X})]$  converges to*

$$n \cdot \frac{d}{2} + \sum_{k=1}^n \frac{1}{k} \sum_{k_1+k_2=k} \max\{0, k_1 a + k_2 \mu\} \cdot \frac{k!}{k_1! k_2!} p_1^{k_1} p_2^{k_2} \quad (27)$$

with  $p_1 = \frac{d}{2(\mu-a)}$  and  $p_2 = 1 - \frac{d}{2(\mu-a)}$ .

*Proof.* Split the inner summation in (11) into three parts. First consider the summation over  $\sum_i k_i = k : k_3 \geq 2$ , hence those instances for which the value  $b$  occurs multiple times. Taking the limit  $b \rightarrow \infty$  inside of the summation and recognizing the fact that the probability mass on the third point is  $O(\frac{1}{b^{k_3}})$  gives

$$\lim_{b \rightarrow \infty} \sum_{\sum_i k_i = k : k_3 \geq 2} \frac{d^{k_3} \max\{0, k_1 a + k_2 \mu + k_3 b\}}{2^{k_3} (b - \mu)^{k_3}} \cdot \frac{k!}{k_1! k_2! k_3!} p_1^{k_1} p_2^{k_2} = 0. \quad (28)$$

Next consider  $\sum_i k_i = k : k_3 = 1$ , describing the instances for which the extremal point  $b$  occurs precisely once. Taking the limit  $b \rightarrow \infty$  inside the sum and using that the probability mass on the point  $b$  is  $O(\frac{1}{b})$  gives

$$\lim_{b \rightarrow \infty} \sum_{\sum_i k_i = k : k_3 = 1} \frac{d \max\{0, k_1 a + k_2 \mu + b\}}{2(b - \mu)} \cdot \frac{k!}{k_1! k_2!} p_1^{k_1} p_2^{k_2} = k \cdot \frac{d}{2} \cdot \sum_{\sum_i k_1 + k_2 = k-1} \frac{(k-1)!}{k_1! k_2!} p_1^{k_1} p_2^{k_2} = k \cdot \frac{d}{2}. \quad (29)$$

The third part is then  $\sum_i k_i = k : k_3 = 0$ , representing the instances without occurrence of the point  $b$ . Taking the limit inside of the summation we get

$$\lim_{b \rightarrow \infty} \sum_{\sum_i k_i = k : k_3 = 0} \max\{0, k_1 a + k_2 \mu\} \cdot \frac{k!}{k_1! k_2!} p_1^{k_1} p_2^{k_2} = \sum_{k_1 + k_2 = k} \max\{0, k_1 a + k_2 \mu\} \cdot \frac{k!}{k_1! k_2!} p_1^{k_1} p_2^{k_2}. \quad (30)$$

This completes the proof.  $\square$

The proof reflects that large running maxima are likely due to a single large step. The feature is caused by heavy-tailed distributions, and in queueing theory dubbed the single big jump principle (see e.g., Foss et al. (2007)). This dominance of one step sharply contrasts intuition for light-tailed distributions, where typically all steps together lead to large sums or maxima. The bound (27) for  $\mathbb{E}[M_n]$  grows to infinity as  $n \rightarrow \infty$ , rendering the bound useless for the expected all-time maximum  $\mathbb{E}[M]$ . This is indeed anticipated, and can be understood as follows. Define a sequence of random walks indexed by  $b$  with the extremal three-point distribution. Consider the limiting all-time maximum  $M$  as  $b \rightarrow \infty$ . Assume that the random walk has negative drift (i.e.,  $\mathbb{E}[X] < 0$ ). Then the associated sequence of distributions of  $M = M_{(b)}$  will converge to a proper limit  $M_{(\infty)}$ . However, as  $\lim_{b \rightarrow \infty} \mathcal{P}_{(\mu, d)}$  contains distributions with infinite second moment, Asmussen (2003), Theorem X.2.1, says that  $\mathbb{E}[M_{(\infty)}]$  will be infinite.

### 3.5. Setting the range to construct adequate bounds

We now present some guidelines for setting the range, based on the observation that many distributions come with a MAD and standard deviation of comparable size. For the Pearson family of

distributions (which includes the gamma and normal distribution) with mean  $\mu$  and variance  $\sigma^2$ , the MAD  $d$  and variance are related as

$$d = 2\alpha\sigma^2 p(\mu) \tag{31}$$

with  $\alpha$  a constant depending on skewness and kurtosis and  $p(\mu)$  the density in  $\mu$ . For the exponential distribution this relation gives  $d = (2/e)\sigma$  and for the normal distribution  $d = (\sqrt{2/\pi})\sigma$ . Other distributions for which the ratio  $d/\sigma$  is constant include the uniform distribution and discrete distributions such as the Poisson, binomial, and negative binomial distribution. With this in mind, in a way similar to constructing confidence intervals in statistical estimation, we then choose to set the range as the mean plus or minus a constant times the MAD:

$$a = \mu - k \cdot d, \quad b = \mu + k \cdot d. \tag{32}$$

Here we regard  $d$  as the natural scale of deviation, and  $k$  as a free parameter that sets the robustness level. So we take the mean and MAD as given, and regard the range as tunable (using common sense or statistical evidence) by the decision maker. We should stress that, while intuitive from a probabilistic perspective, the rule (32) is only one of many ways to choose the parameters  $a, b$ .

We demonstrate (32) for a setting where we take the M/M/1 queue as the ‘true’ model. The increment  $X$  now becomes the difference of two exponential random variables for which we have a closed-form MAD expression in terms of the mean value of  $X$  (see the caption of Table 1). We thus have reference values for  $\mu$  and  $d$ , and can investigate the impact of  $k$ . Observe that the bound grows almost linearly with  $k$ , in particular in heavy-traffic scenarios, and this underlines the need for careful selection of the range. While the actual range of the M/M/1 queue spans all real numbers, we see that restricting deviations to twice the MAD ( $k = 2$ ) gives comparable model performance. When reading Table 1, keep in mind that the overall goal in this paper is not to approximate specific models (such as the M/M/1 queue), but rather to come with conservative, robust estimates for an entire class of models that share the same mean-MAD-range properties. In that sense,  $k = 2$  is not better than  $k = 1.5$  or  $k = 2.5$ , but rather expresses a different ambiguity assessment or robustness level.

#### 4. Extremal GI/G/1 queue

Let us now turn to the extremal GI/G/1 queue problem, as described in the introduction. Let  $W_n$  be the waiting time of customer  $n$ . The sequence  $(W_n, n \geq 0)$  with  $W_0 = 0$  satisfies the Lindley recursion

$$W_{n+1} = (W_n + V_n - U_n)^+, \quad n \geq 0. \tag{33}$$

**Table 1** The actual value and bounds of the expected steady-state waiting time  $\mathbb{E}[W]$  of the M/M/1 queue with unit mean exponential interarrival times and exponential service times with mean  $\rho$ , where the increment  $X$  has mean  $\mu = \rho - 1$  and MAD  $d = \frac{2e^{\rho-1}}{\rho+1}$ , with the range  $[a, b]$  set through the rule (32).

$\rho$	$\mathbb{E}[W]$	$k$					
		1.5	1.75	2	2.25	2.5	3
0.1	0.01111	0.10497	0.16434	0.21535	0.25915	0.30116	0.40782
0.5	0.50000	0.56329	0.67919	0.79663	0.91459	1.02840	1.26462
0.6	0.90000	0.86690	1.03323	1.19770	1.36332	1.52804	1.85818
0.7	1.63333	1.41436	1.66589	1.91885	2.17142	2.42373	2.92850
0.8	3.20000	2.57273	3.01339	3.45454	3.89573	4.33672	5.21866
0.9	8.10000	6.21057	7.25250	8.29428	9.33642	10.37811	12.46184
0.99	98.01000	73.55537	85.81540	98.07540	110.33542	122.59543	147.11548

Let  $W$  be the steady-state waiting time. Since  $W_n \stackrel{d}{=} M_n$  and  $W \stackrel{d}{=} M$  the results for the random walk maxima likely carry over to the waiting times. The main difference is that the step size  $X$  is now interpreted as the difference  $V - U$  between the generic service time and generic interarrival time. If one has mean-MAD information about both  $V$  and  $U$  this is more informative than mean-MAD information about  $V - U$ , and this additional information should lead to even sharper bounds.

#### 4.1. A complete picture

The GI/G/1 queue assumes that interarrival times and service times are independent, so it is natural to assume that  $V$  has ambiguity set  $\mathcal{P}_{(\mu_V, d_V)}$  and  $U$  has ambiguity set  $\mathcal{P}_{(\mu_U, d_U)}$ , where the ambiguity sets now contain all distributions for *univariate*  $V$  and  $U$ , that is,

$$\mathcal{P}_{(\mu_V, d_V)} = \{\mathbb{P} : \text{supp}(V) \subseteq [a_V, b_V], \mathbb{E}_{\mathbb{P}}(V) = \mu_V, \mathbb{E}_{\mathbb{P}}|V - \mu_V| = d_V\}$$

and

$$\mathcal{P}_{(\mu_U, d_U)} = \{\mathbb{P} : \text{supp}(U) \subseteq [a_U, b_U], \mathbb{E}_{\mathbb{P}}(U) = \mu_U, \mathbb{E}_{\mathbb{P}}|U - \mu_U| = d_U\}.$$

The extremal queue problem with mean-MAD dispersion information can then be phrased as

$$\max_{\mathbb{P} \in \mathcal{P}_{(\mu_V, d_V)} \times \mathcal{P}_{(\mu_U, d_U)}} \mathbb{E}[f(\mathbf{X})], \quad (34)$$

where  $\mathbb{E}[f(\mathbf{X})]$  describes  $\mathbb{E}[W_n]$  or  $c_m(W)$  and  $\mathbf{X}$  is the random vector with elements  $U_1, V_1, U_2, V_2, \dots$ . This is the classical setting of the extremal GI/G/1 queue treated in Rolski (1972), Eckberg Jr (1977), Whitt (1984), Chen and Whitt (2019), but with MADs instead of variances describing the ambiguity set. Let the random variables  $V_{(3)}$  and  $U_{(3)}$  follow the extremal three-point distributions  $\Omega(\mu_V, d_V, a_V, b_V)$  and  $\Omega(\mu_U, d_U, a_U, b_U)$ , respectively.

**THEOREM 4.** *Consider the GI/G/1 queue with generic interarrival time  $U$  with ambiguity set  $\mathcal{P}_{(\mu_U, d_U)}$  and generic service times  $V$  with ambiguity set  $\mathcal{P}_{(\mu_V, d_V)}$ . Consider the tight upper bounds for the transient mean waiting time  $\mathbb{E}[W_n]$  and all cumulants of the steady-state waiting time  $W$ .*

- (i) For given interarrival time  $U$ , the tight upper bounds follow from the service time  $V_{(3)}$ .
- (ii) For given service time  $V$ , the tight upper bounds follow from the interarrival time  $U_{(3)}$ .
- (iii) The overall tight upper bounds follow from interarrival time  $U_{(3)}$  and service time  $V_{(3)}$ .

*Proof.* Like Theorem 1, the tight bounds for  $\mathbb{E}[W_n]$  follow from the general upper bound in Ben-Tal and Hochman (1972) on the expectation of a convex function of the random vector  $(X_1, \dots, X_n)$  with mean-MAD ambiguity, but now with  $X_i$  replaced by  $V_i - U_i$ . The function describing  $\mathbb{E}[W_n]$  (see Theorem 1) is indeed convex in *both*  $V_i$  and  $U_i$ , and hence the result follows. Similarly, Spitzer's formula for  $c_m(W)$  (see Theorem 2) is also convex in both  $V_i$  and  $U_i$ , and hence the tight bounds for  $c_m(W)$  follow from our proof of Theorem 2.  $\square$

Using the earlier results for the random walk, we present in EC.3 expressions that are helpful in evaluating the tight bounds. Table 2 shows an example of the tight bound for  $\mathbb{E}[W]$  associated with  $(U_{(3)}, V_{(3)})$ , also compared with other known bounds that require variance information (see EC.4). The variance of the extremal three-point distribution  $\Omega(\mu, d, a, b)$  is  $\frac{d}{2}(b-a)$ , the maximal variance for distributions in the ambiguity set  $\mathcal{P}_{(\mu, d)}$ . We thus know the variances of  $U_{(3)}$  and  $V_{(3)}$ , and can calculate the other three bounds. In heavy traffic, Kingman's bound is known to be asymptotically correct, and hence the other three (sharper) bounds also converge to the heavy-traffic limit as  $\rho \uparrow 1$ . See EC.5 for more numerical results. Notice that Table 2 is not meant to compare mean-MAD with mean-variance bounds. The displayed differences merely express different ways of dealing with ambiguity. Also remember that the mean-MAD bounds in Theorem 4 are crucially influenced by the choice of range, in this example set to  $[0, 10]$  for both the interarrival and service time distributions.

**Table 2**    **Bounds for  $(1-\rho)\mathbb{E}[W]/\rho$  for  $(\mu_U, d_U, a_U, b_U) = (1, 1, 0, 10)$  and  $(\mu_V, d_V, a_V, b_V) = (\rho, 0.1, 0, 10)$ .**

$\rho$	Thm. 4	C & W (EC.14)	Daley (EC.13)	Kingman (EC.12)
0.1	4.06613	7.00020	7.25000	27.50000
0.2	2.52306	5.27810	5.75000	13.75000
0.5	2.03141	3.63750	4.25000	5.50000
0.7	2.49160	3.17138	3.60714	3.92857
0.8	2.61932	3.00523	3.31250	3.43750
0.9	2.69802	2.86711	3.02778	3.05556
0.95	2.72609	2.80627	2.88816	2.89474
0.99	2.74547	2.76091	2.77753	2.77778

#### 4.2. Further comparison between MAD and variance

For the variance counterpart, Chen and Whitt (2019) also formulate a semi-infinite linear optimization problem. The crucial difference is that they cannot use the univariate function extension

(as explained in Section 2), and hence should work directly with the multivariate function. This in turn implies that the dual problem cannot be solved explicitly (like in the univariate case), let alone that there is a zero duality gap. Another complication is that the multivariate function based on Spitzer’s formulas (8) and (11) cannot be expressed directly in  $V$  and  $U$ , but rather in terms of convolutions of the distributions of  $V$  and  $U$ . Chen and Whitt (2019) resolve these considerable challenges by several ingenious arguments, a.o. exploiting the description of  $W$  as a fixed point of the stochastic equation  $W \stackrel{d}{=} (W + V - U)^+$ , and by imposing additional regularity conditions on  $V$ . In this way, Chen and Whitt (2019) prove a similar but weaker result than Theorem 4 for the exact same setting, but with variance as dispersion measure. They show that the extremal distributions of  $U$  and  $V$  both have supports on at most three points.

An important message of this paper is that with MAD the extremal distribution remains unaltered going from the univariate to the multivariate setting, and that with variance this reasoning fails. In fact, one intuitively expects formidable challenges when seeking for extremal distributions under variance constraints. This intuition is confirmed by Chen and Whitt’s formulation of the extremal distribution as the solution of a non-convex nonlinear optimization problem. While this optimization problem can be solved numerically, a closed-form solution and hence identification of the extremal distribution remains out of reach.

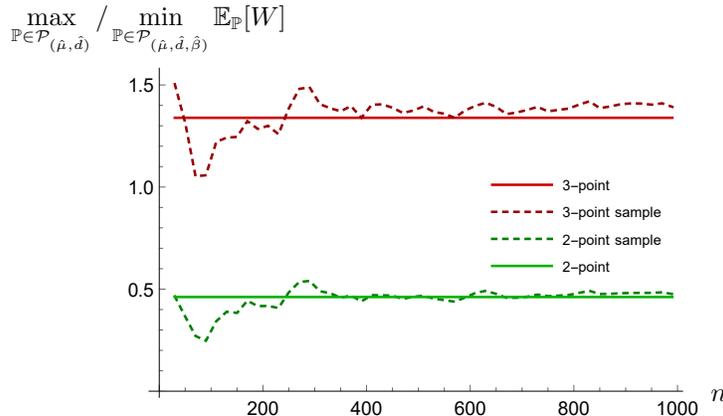
Under variance constraints, it is conjectured that the tight bound comes from specific two-point distributions for both  $U$  and  $V$ . In fact, the bound (EC.14) in Table 2 holds under the assumption that this conjecture is true, and was shown by Chen and Whitt (2019) to be very close to the tight upper bound. Theorem 4 rules out a similar conjecture in the MAD setting. The tight bounds in Theorem 4 always involve three-point distributions.

### 4.3. Data-driven setting

In applications, you may only have a limited number  $n$  of observed interarrival and service times. We consider this realistic setting where knowledge of the stochastic nature is restricted to a set of samples generated independently and randomly according to an unknown distribution  $\mathbb{P}$ . To apply the mean-MAD framework in this context, we need to construct the ambiguity set that is supposed to contain this unknown  $\mathbb{P}$ . We will show that we can efficiently estimate the mean, MAD, and  $\beta$ , and hence compute robust bounds that are useful in realistic settings.

Let  $\mu_n^{(V)}$ ,  $d_n^{(V)}$  and  $\beta_n^{(V)}$  denote the consistent estimators of  $\mu_V$ ,  $d_V$  and  $\beta_V = \mathbb{P}(V \geq \mu_V)$ , respectively, based on  $n$  observed service times  $v_1, \dots, v_n$ , and defined as  $\mu_n^{(V)} = \bar{v} = \frac{1}{n} \sum_{i=1}^n v_i$ ,  $d_n^{(V)} = \frac{1}{n} \sum_{i=1}^n |v_i - \bar{v}|$  and  $\beta_n^{(V)} = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{[\bar{v}, \infty)}(v_i)$  as consistent estimators. We define similar estimators based on  $n$  observed interarrival times. Next, we demonstrate the mean-MAD bounds in this data-driven setting. Since statistical accuracy of the estimators increases with the number of samples,

we expect the bounds to converge as  $n$  increases. Figure 4 illustrates two sample paths representing the estimates for the upper and lower bound and their convergence to the tight mean-MAD bounds, where  $V$  and  $U$  both follow a uniform distribution on the intervals  $[0, 5]$  and  $[0, 10]$ , respectively. Observe that convergence settles in quickly.



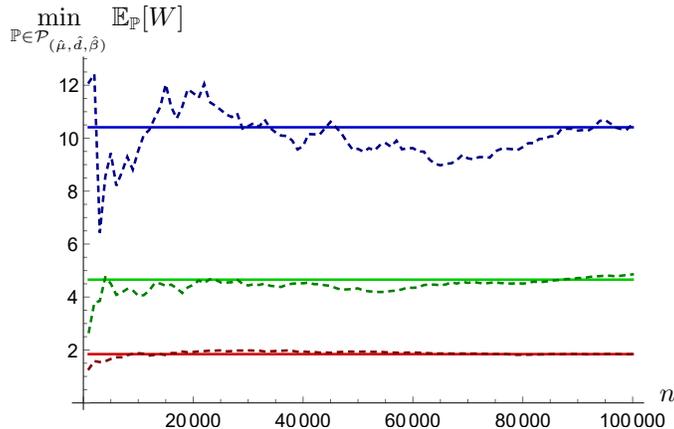
**Figure 4** Estimation of the mean-MAD ambiguity upper and lower bound. The red and green line represent the true upper and lower bound, respectively, and the dashed lines represent bound estimates which are computed using the realizations  $v_1, \dots, v_n$  drawn from a  $U(0, 5)$  distribution and  $u_1, \dots, u_n$  sampled from a  $U(0, 10)$  distribution.

We have also performed extensive simulations to investigate the error between the estimated and true bounds for several values of the sample size  $n$ . We generate 1,000 sample paths of sample size 10,000 and compute the corresponding mean relative error. Table 3 displays the mean absolute percentage error (MAPE) for both the upper and lower bound estimates, where the interarrival time is  $U(0, 10)$  distributed and we differentiate between a 50% and 90% utilization level. Observe that estimating the lower bound is slightly harder than estimating the upper bound. Indeed, the lower bound requires estimating the additional parameters  $\beta_V$  and  $\beta_U$ . Also observe that the relative error increases with the system utilization.

**Table 3** MAPE of the bound estimates for  $n \in \{150, 200, 500, 1000, 2000, 5000, 10000\}$ . The interarrival times are  $U(0, 10)$  distributed and the results differentiate between two service time distributions and the upper and lower mean-MAD bounds. Sample paths resulting in instable systems were removed and done over.

Service times	Bound	MAPE with sample size $n$						
		150	200	500	1000	2000	5000	10000
$U(0, 5)$	UB	15.44%	13.22%	8.31%	5.84%	4.28%	2.72%	1.89%
	LB	25.51%	22.30%	13.86%	9.75%	7.08%	4.53%	3.16%
$U(0, 9)$	UB	33.35%	30.93%	21.93%	16.35%	13.29%	8.92%	6.41%
	LB	36.27%	35.01%	28.72%	22.30%	17.35%	10.77%	7.58%

To further highlight the role of system utilization, we perform a similar data-driven experiment, but now with ground truth a single trace of  $n$  customers in an M/M/1 queue. The results are shown in Figure 5. Indeed, as  $\rho$  increases, more observations are required for accurate parameter estimates and hence accurate bounds.



**Figure 5** Estimation of the mean-MAD ambiguity lower bound for the M/M/1 queue. The solid red, green, and blue lines depict the bounds for  $\rho = 0.8, 0.9, 0.95$ , respectively. The dashed lines represent the corresponding estimates of the bounds, where the  $U_i$  are sampled from a unit mean exponential distribution and the  $V_i$  are exponentially distributed with mean  $\rho$ .

Taken together, we conclude that the robust bounds are useful for realistic data-driven settings that require statistical estimation of the summary statistics such as the mean and MAD.

## 5. Conclusions

This paper explains why MAD simplifies comparable variance-based optimization problems, in a way that is almost unreasonably effective, resulting in a full solution to the extremal queue problem with mean-MAD constraints. When partial information is available in the form of mean, range and MAD, we have obtained the sharpest possible bounds. Through basic statistical estimation of this partial information, the GI/G/1 queue becomes a data-driven model that adjusts to available training data, for which this paper presents tight performance guarantees.

The key idea of using MAD instead of variance as dispersion measure, is likely applicable to many other queueing systems. Examples are queues with dependency and correlation structures in the series  $\{U_n\}$  and  $\{V_n\}$ , the multi-server GI/G/c queue and networks of queues. Indeed, most of the key performance measures for such systems are expectations of functions that are convex in the random variables (see e.g., Shaked and Shanthikumar (1988)), and therefore the mean-MAD approach can be used. The MAD perspective is of interest beyond queueing theory, because the search for extremal distributions of convex functions is relevant in many other settings. Moreover,

whenever a performance measure can be viewed as a convex function of i.i.d. random variables with mean-MAD ambiguity (e.g., nested max-operators in production systems; see Glasserman (1997), Bradley and Glynn (2002)), our approach will identify the extremal distribution and tight bounds.

The MAD approach stays close to the common practice in the stochastic field, namely to use probability distributions to model uncertainty. The nucleus of the MAD approach consists of the explicitly solvable dual LP described in Section 2. A simple reasoning then showed that this solution is independent of the precise objective function (in this paper describing waiting time moments of the GI/G/1 queue). Hence, the MAD approach is a generic, computationally tractable way to analyze stochastic processes, such as random walks and queues.

Let us conclude with a broader robust optimization perspective. It is well-known that the use of probability distributions in stochastic systems often leads to computationally intractability (e.g., calculation of high dimensional convolutions). Therefore, Bandi and Bertsimas (2012), Bandi et al. (2015), Whitt and You (2018) suggest to use uncertainty sets instead of probability distributions. The MAD approach described in this paper can serve in many situations as an alternative (not per se better), bringing new opportunities. The uncertainty set approach yields a worst-case scenario. Our approach yields both worst-case and best-case distributions, i.e., both upper and lower bounds. In stochastic systems one often studies convex functions in the stochastic variables. In the uncertainty set approach it is in general hard (in fact, NP-hard) to find worst-case scenarios for such convex functions. Our approach can easily find worst-case distributions as shown in this paper.

## Acknowledgments

The authors would like to thank Daniel Kuhn and Krzysztof Postek for pointing out the primal-dual reasoning that gives the intuitive proof of Theorem 1, and Marko Boon for helping with the experiments in Section EC.3.

## References

- Abate J, Choudhury G, Whitt W (1993) Calculation of the GI/G/1 waiting-time distribution and its cumulants from Pollaczek’s formulas. *Archiv fur Elektronik und Ubertragungstechnik (International Journal of Electronics and Communication)* 47(5/6):311–321.
- Abate J, Whitt W (1992) The Fourier-series method for inverting transforms of probability distributions. *Queueing Systems* 10(1-2):5–87.
- Asmussen S (2003) *Applied Probability and Queues* (New York: Springer-Verlag), second edition.
- Bandi C, Bertsimas D (2012) Tractable stochastic analysis in high dimensions via robust optimization. *Mathematical Programming* 134(1):23–70.
- Bandi C, Bertsimas D, Youssef N (2015) Robust queueing theory. *Operations Research* 63(3):676–700.

- Ben-Tal A, Hochman E (1972) More bounds on the expectation of a convex function of a random variable. *Journal of Applied Probability* 9:803–812.
- Ben-Tal A, Hochman E (1985) Approximation of expected returns and optimal decisions under uncertainty using mean and mean absolute deviation. *Zeitschrift für Operations Research* 29(7):285–300.
- Bradley J, Glynn P (2002) Managing capacity and inventory jointly in manufacturing systems. *Management Science* 48(2):273–288.
- Chen Y, Whitt W (2019) Extremal GI/GI/1 queues given two moments. *Submitted to Operations Research*, Preprint.
- Chen Y, Whitt W (2020) Algorithms for the upper bound mean waiting time in the GI/GI/1 queue. *Queueing Systems* 94:327–356.
- Chung K (2001) *A Course in Probability Theory* (London: Academic Press).
- Cohen J (1982) *The Single Server Queue* (Amsterdam: North-Holland Publishing Co.), second edition.
- Daley DJ, Kreinin AY, Trengove CD (1992) Inequalities concerning the waiting-time in single-server queues: a survey. Bhat UN, Basawa IV, eds., *Queueing and Related Models*, 177–223 (Oxford: Clarendon Press).
- Das B, Dhara A, Natarajan K (2018) On the heavy-tail behavior of the distributionally robust newsvendor. *arXiv preprint arXiv:1806.05379*.
- Eckberg Jr A (1977) Sharp bounds on Laplace-Stieltjes transforms, with applications to various queueing problems. *Mathematics of Operations Research* 2(2):135–142.
- Feller W (1971) *An Introduction to Probability Theory and its Applications. Vol. II.* (New York: John Wiley & Sons Inc.), second edition.
- Foss S, Konstantopoulos T, Zachary S (2007) Discrete and continuous time modulated random walks with heavy-tailed increments. *Journal of Theoretical Probability* 20(3):581–612.
- Glasserman P (1997) Bounds and asymptotics for planning critical safety stocks. *Operations Research* 45(2):244–257.
- Han S, Tao M, Topcu U, Owhadi H, Murray RM (2015) Convex optimal uncertainty quantification. *SIAM Journal on Optimization* 25(3):1368–1387.
- Janssen A, van Leeuwen J, Mathijssen B (2015) Novel heavy-traffic regimes for large-scale service systems. *SIAM Journal on Applied Mathematics* 75(2):787–812.
- Kingman JF (1962) Some inequalities for the queue GI/G/1. *Biometrika* 49(3/4):315–324.
- Natarajan K, Sim M, Uichanco J (2017) Asymmetry and ambiguity in newsvendor models. *Management Science* 64(7):3146–3167.
- Natarajan K, Zhou L (2007) A mean–variance bound for a three-piece linear function. *Probability in the Engineering and Informational Sciences* 21(4):611–621.

- Perakis G, Roels G (2008) Regret in the newsvendor model with partial information. *Operations Research* 56(1):188–203.
- Postek K, Ben-Tal A, Den Hertog D, Melenberg B (2018) Robust optimization with ambiguous stochastic constraints under mean and dispersion information. *Operations Research* 66(3):814–833.
- Rogosinski WW (1958) Moments of non-negative mass. *Proceedings of the Royal Society of London. Series A. Mathematical and Physical Sciences* 245(1240):1–27.
- Rolski T (1972) Some inequalities for GI/M/n queues. *Applicaciones Mathematicae* 1(13):42–47.
- Shaked M, Shanthikumar JG (1988) Stochastic convexity and its applications. *Advances in Applied Probability* 20(2):427–446.
- Shapiro A, Dentcheva D, Ruszczyński A (2009) *Lectures on Stochastic Programming: Modeling and Theory* (Philadelphia: SIAM).
- Spitzer F (1956) A combinatorial lemma and its application to probability theory. *Transactions of the American Mathematical Society* 82:323–339, ISSN 0002-9947.
- Whitt W (1984) On approximations for queues, I: Extremal distributions. *AT&T Bell Laboratories Technical Journal* 63(1):115–138.
- Whitt W, You W (2018) Using robust queueing to expose the impact of dependence in single-server queues. *Operations Research* 66(1):184–199.
- Xin L, Goldberg DA (2013) Time (in)consistency of multistage distributionally robust inventory models with moment constraints. *arXiv preprint arXiv:1304.3074*.

# E-Companion to “MAD dispersion measure makes extremal queue analysis simple”

## EC.1. Properties of MAD

We recall some well known properties of the MAD, see e.g. Ben-Tal and Hochman (1985). Denote by  $\sigma^2$  the variance of the random variable  $X$ , whose distribution is known to belong to the set  $\mathcal{P}_{(\mu,d)}$ . Then

$$\frac{d^2}{4\beta(1-\beta)} \leq \sigma^2 \leq \frac{d(b-a)}{2}.$$

In particular, since

$$d^2 \leq 4\beta(1-\beta)\sigma^2 \leq \sigma^2,$$

it holds that  $d \leq \sigma$ . For a proof, we refer the reader to Ben-Tal and Hochman (1985). For some distributions, an explicit formula for  $d$  is available:

- Uniform distribution on  $[a, b]$ :

$$d = \frac{1}{4}(b-a)$$

- Normal distribution  $N(\mu, \sigma^2)$ :

$$d = \sqrt{\frac{2}{\pi}}\sigma$$

- Gamma distribution with parameters  $\lambda$  and  $k$  (for which  $\mu = k/\lambda$ ):

$$d = \frac{2k^k}{\Gamma(k)\exp(k)} \frac{1}{\lambda}.$$

The MAD is known to satisfy the bound

$$0 \leq d \leq \frac{2(b-\mu)(\mu-a)}{b-a}. \tag{EC.1}$$

Let  $\beta = \mathbb{P}(X \geq \mu)$ . For example, in the case of continuous symmetric distribution of  $X$  we know that  $\beta = 0.5$ . This quantity is known to satisfy the bounds:

$$\frac{d}{2(b-\mu)} \leq \beta \leq 1 - \frac{d}{2(\mu-a)}. \tag{EC.2}$$

## EC.2. Primal-dual proof of Theorem 3

In a similar manner as for the upper bound, we will show that the best-case distribution is a two-point distribution. We again consider the convex univariate measurable function  $f(x)$  that has finite values on  $[a, b]$ . Under  $\mathcal{P}_{(\mu,d,\beta)}$  ambiguity of the random variable  $X$  we now need to solve

$$\begin{aligned} & \min_{\mathbb{P}(x) \geq 0} \int_x f(x) d\mathbb{P}(x) \\ \text{s.t.} \quad & \int_x \mathbb{1}_{\{x \geq \mu\}} d\mathbb{P}(x) = \beta, \int_x |x - \mu| d\mathbb{P}(x) = d, \int_x x d\mathbb{P}(x) = \mu, \int_x d\mathbb{P}(x) = 1, \end{aligned} \tag{EC.3}$$

which is a semi-infinite linear program with four equality constraints.

Consider the dual of (EC.3),

$$\begin{aligned} \max_{\lambda_0, \lambda_1, \lambda_2, \lambda_3} \quad & \lambda_0 \beta + \lambda_1 d + \lambda_2 \mu + \lambda_3 \\ \text{s.t.} \quad & f(x) - \lambda_0 \mathbb{1}_{\{x \geq \mu\}} - \lambda_1 |x - \mu| - \lambda_2 x - \lambda_3 \geq 0, \quad \forall x \in [a, b]. \end{aligned} \quad (\text{EC.4})$$

Define  $F(x) = \lambda_0 \mathbb{1}_{\{x \geq \mu\}} + \lambda_1 |x - \mu| + \lambda_2 x + \lambda_3$ . Then the inequality in (EC.4) can be written as  $F(x) \leq f(x), \forall x$ , i.e.  $F(x)$  minorizes  $f(x)$ . Note that in our new situation  $F(x)$  has both a kink and a discontinuity at  $x = \mu$ , as depicted in Figure EC.1. The dual problem boils down to finding the tightest minorant that maximizes the dual problem's objective value. The minorant  $F(x)$  touches the epigraph of  $f(x)$  in at most two points on opposite sides of  $\mu$  (i.e.,  $x_1 \leq \mu \leq x_2$ ). This is a consequence of the supporting hyperplane theorem and the jump discontinuity at  $x = \mu$ . The dual problem now becomes

$$\begin{aligned} \max_{\lambda_0, \lambda_1, \lambda_2, \lambda_3} \quad & \lambda_0 \beta + \lambda_1 d + \lambda_2 \mu + \lambda_3 \\ \text{s.t.} \quad & \lambda_0 + \lambda_1(x_1 - \mu) + \lambda_2 x_1 + \lambda_3 = f(x_1), \\ & -\lambda_1(x_2 - \mu) + \lambda_2 x_2 + \lambda_3 = f(x_2). \end{aligned} \quad (\text{EC.5})$$

Now using Lagrange duality, we can show that the optimal solution satisfies

$$x_1 = \mu + \frac{d}{2\beta}, \quad x_2 = \mu - \frac{d}{2(1-\beta)},$$

which corresponds to the values of  $v_1$  and  $v_2$  stated in Theorem 3. Substituting this solution and solving for  $\lambda_0, \lambda_1, \lambda_2$ , and  $\lambda_3$  gives

$$\lambda_0 = f(v_1) - f(v_2) + \frac{\lambda_1 d}{(1-\beta)} - \frac{(\lambda_1 + \lambda_2)d}{2\beta(1-\beta)}, \quad \lambda_3 = f(v_2) + \frac{(\lambda_2 - \lambda_1)d}{2(1-\beta)} - \lambda_2 \mu,$$

and hence the objective value of the dual becomes  $\beta f(v_1) + (1-\beta)f(v_2)$ . Note that we have two free variables that can be chosen in a way that makes the solution dual feasible. The optimal probabilities of (EC.3) are obtained by solving the linear system resulting from (EC.3), which produces the solution stated in Theorem 3. Finally, one can verify that the primal and dual objective values are the same and that these results can be extended to the multivariate case in a manner analogous to that of Theorem 1.

### EC.3. Representations for the tight bounds

We will now present some efficient ways of calculating the tight bounds identified in this paper. But first we show a way to verify the contour integral representation.

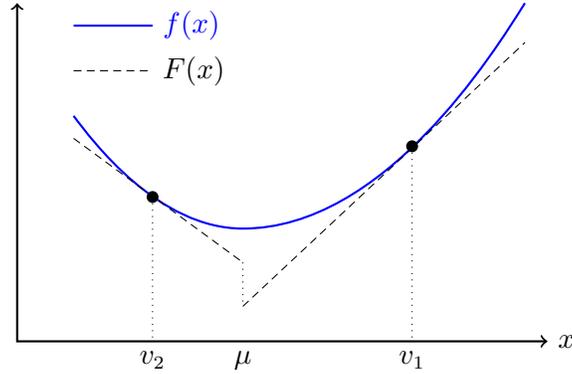


Figure EC.1 Some convex function  $f(x)$  and its non-continuous piecewise linear minorant  $F(x)$ .

### EC.3.1. Numerical experiments with contour integrals

Numerical aspects of integrals of the type (17) have been discussed in e.g., Abate et al. (1993), Janssen et al. (2015), Chen and Whitt (2020). For distributions with support on a finite set of points, potential numerical problems can arise, because  $|\operatorname{Re}(\phi_X(u))|$  does not converge to zero as  $|u| \rightarrow \infty$ ; see Abate and Whitt (1992), Chen and Whitt (2020). For the three-point distributions required in this paper we have performed extensive numerical experiments with (19). These experiments confirmed that the integrals can be calculated up to high accuracy with standard integration routines in Mathematica (our code is available upon request).

For many parameter values  $a, b, \mu, d$  such that (EC.1) holds, we have calculated  $\mathbb{E}[M]$  for generic increment  $X_{(3)}$  using (19), and compared this with results from extensive stochastic simulations. We also compared the results with a third numerical procedure, known to be extremely stable and accurate. Let us explain the third procedure, which might be of independent interest.

Choose the boundaries of the support as multiples of  $\beta = |\mu|$  by writing that  $a = -s\beta$  and  $b = m\beta$  with  $s, m$  positive integers. Denote by  $M_\beta = M/\beta$  the normalized steady-state waiting time. We then get

$$M_\beta \stackrel{d}{=} (M_\beta + X_\beta)^+,$$

with  $X_\beta = X/\beta$  a discrete random variable with support  $\{-s, -1, m\}$  and MAD

$$d_\beta := \mathbb{E}[|X_\beta - \mathbb{E}[X_\beta]|] = \frac{1}{\beta} \mathbb{E}[|X - \mathbb{E}[X]|] = d.$$

Define  $X_\beta = A_\beta - s$ , so that

$$M_\beta \stackrel{d}{=} (M_\beta + A_\beta - s)^+$$

for a discrete random variable  $A_\beta$  with support  $\{0, s-1, s+m\}$  and probability generating function

$$\mathbb{E}[z^{A_\beta}] = p_a + p_\mu z^{s-1} + p_b z^{m+s},$$

with

$$p_a = \frac{d_\beta}{2(s-1)}, \quad p_\mu = 1 - \frac{d_\beta}{2(s-1)} - \frac{d_\beta}{2(m+1)}, \quad p_b = \frac{d_\beta}{2(m+1)}.$$

Notice that  $\mathbb{E}[A_\beta] = s - 1$ . The resulting discrete queueing system is sometimes referred to as a bulk service queue. Let  $r_0$  be the unique zero of  $z^s - \mathbb{E}[z^{A_\beta}]$  with real  $z > 1$ . For any  $\varepsilon > 0$  with  $1 + \varepsilon < r_0$ ,

$$\mathbb{E}[w^{M_\beta}] = \exp\left(\frac{1}{2\pi i} \int_{|z|=1+\varepsilon} \ln\left(\frac{w-z}{1-z}\right) \frac{(z^s - \mathbb{E}[z^{A_\beta}])'}{z^s - \mathbb{E}[z^{A_\beta}]} dz\right) \quad (\text{EC.6})$$

holds when  $|w| < 1 + \varepsilon$ . Alternatively,

$$\mathbb{E}[w^{M_\beta}] = \frac{(s - \mathbb{E}[A_\beta])(w-1)}{w^s - A(w)} \prod_{k=1}^{s-1} \frac{w - z_k}{1 - z_k} \quad (\text{EC.7})$$

that holds for all  $w$ ,  $|w| < r_0$ , in which  $z_1, \dots, z_{s-1}$  are the  $s - 1$  zeros of  $z^s - \mathbb{E}[z^{A_\beta}]$  in  $|z| < 1$ . Upon differentiation, (EC.6) and (EC.7) provide expressions for all cumulants of  $M_\beta$  that are known to allow for accurate numerical evaluation, see Janssen et al. (2015). We have then performed for a wide range of parameters, the following experiment:

1. Fix  $\beta$ , and then choose integers  $s$  and  $m$ . In this way we create a standard bulk service queue with discrete-valued generic increment  $A_\beta$ .
2. For ranging  $d_\beta$ , calculate  $\mathbb{E}[M_\beta]$  using root-finding procedures and (EC.7) or using the contour integral (EC.6).
3. Calculate

$$\mathbb{E}[M] = \frac{-1}{2\pi i} \int_{\mathcal{C}} \frac{\log(1 - (p_a e^{-ua} + p_b e^{-ub} + p_c e^{-uc}))}{u^2} du.$$

4. Check whether  $\mathbb{E}[M] = \beta \mathbb{E}[M_\beta]$ .

### EC.3.2. Numerical procedures for the GI/G/1 queue

Calculations for  $\mathbb{E}[W_n]$  and  $c_n(W)$  in the GI/G/1 queue can be performed using similar expressions as for the random walk. Let the random variable  $V_{(3)}$  follow a three-point distribution on values  $\{s_1, s_2, s_3\}$  with probabilities

$$p_1 = \frac{d_V}{2(\mu_V - a_V)}, \quad p_2 = 1 - \frac{d_V}{2(\mu_V - a_V)} - \frac{d_V}{2(b_V - \mu_V)}, \quad p_3 = \frac{d_V}{2(b_V - \mu_V)}, \quad (\text{EC.8})$$

with  $0 \leq a_V < \mu_V < b_V$ , so that  $V_{(3)}$  has mean  $\mu_V$  and MAD  $d_V$ . Similarly, let  $U_{(3)}$  have a three-point distribution on values  $\{t_1, t_2, t_3\}$  with probabilities

$$r_1 = \frac{d_U}{2(\mu_U - a_U)}, \quad r_2 = 1 - \frac{d_U}{2(\mu_U - a_U)} - \frac{d_U}{2(b_U - \mu_U)}, \quad r_3 = \frac{d_U}{2(b_U - \mu_U)} \quad (\text{EC.9})$$

and  $0 \leq a_U < \mu_U < b_U$ , so that  $U_{(3)}$  has mean  $\mu_U$  and MAD  $d_U$ .

We then have the representation, see also Chen and Whitt (2019),

$$\mathbb{E}[W_n] = \sum_{k=1}^n \frac{1}{k} \sum_{\sum_i k_i=k, \sum_j l_j=k} \max\{0, \sum_{i=1}^3 k_i s_i - \sum_{j=1}^3 l_j t_j\} \cdot P(k_1, k_2, k_3) \cdot R(l_1, l_2, l_3) \quad (\text{EC.10})$$

with

$$P(k_1, k_2, k_3) = \frac{k!}{k_1!k_2!k_3!} p_1^{k_1} p_2^{k_2} p_3^{k_3}, \quad R(l_1, l_2, l_3) = \frac{k!}{l_1!l_2!l_3!} r_1^{l_1} r_2^{l_2} r_3^{l_3},$$

which requires summing  $O(n^5)$  terms.

Let  $\phi_{V_{(3)}}(s)$  and  $\phi_{U_{(3)}}(s)$  denote the moment generating functions of  $V_{(3)}$  and  $U_{(3)}$ . The tight upper bounds on  $c_m(W)$  are given by

$$c_m(W) \leq \frac{(-1)^m}{2\pi i} \int_{\mathcal{C}} \frac{\log(1 - \phi_{V_{(3)}}(-u)\phi_{U_{(3)}}(u))}{u^{m+1}} du, \quad (\text{EC.11})$$

where  $\mathcal{C}$  is a contour to the left of, and parallel to, the imaginary axis, and to the right of any singularities of  $\log(1 - \phi_{V_{(3)}}(-u)\phi_{U_{(3)}}(u))$  in the left half plane. Again comparing with extensive simulation, we have found the expression (EC.11) accurate and hence suitable for calculating the tight bounds.

#### EC.4. Distribution-free upper bounds for the GI/G/1 queue

Consider the steady-state queue length  $W$  in the GI/G/1 queue, which satisfies  $W \stackrel{d}{=} (W + V - U)^+$ . Denote by  $\sigma_U^2$  and  $\sigma_V^2$  the variances of  $U$  and  $V$ , respectively. Let  $\rho = \mathbb{E}[V]/\mathbb{E}[U] < 1$ . The following bounds on  $\mathbb{E}[W]$  only require information about the first two moments of  $U$  and  $V$ :

- Kingman's upper bound:

$$\mathbb{E}[W] \leq \frac{\sigma_V^2 + \sigma_U^2}{2(\mathbb{E}[U] - \mathbb{E}[V])}. \quad (\text{EC.12})$$

- Daley's upper bound:

$$\mathbb{E}[W] \leq \frac{\sigma_V^2 + \rho(2 - \rho)\sigma_U^2}{2(\mathbb{E}[U] - \mathbb{E}[V])}. \quad (\text{EC.13})$$

- Upper bound of Chen and Whitt (2019) based on the two-point conjecture:

$$\mathbb{E}[W] \leq \frac{\sigma_V^2 + \kappa(\rho)\sigma_U^2}{2(\mathbb{E}[U] - \mathbb{E}[V])}, \quad (\text{EC.14})$$

with  $\kappa(\rho) = 2\rho(1 - \rho)/(1 - \delta)$  and  $\delta \in (0, 1)$  the solution of  $\delta = \exp(-(1 - \delta)/\rho)$ .

#### EC.5. Further numerical results for the bounds

We now complement Table 2 with some more numerical values for the bounds on  $\mathbb{E}[W]$ . Table EC.1 gives the unscaled values of  $\mathbb{E}[W]$  for the same parameter values as in Table 2.

The variance bounds are often reported in terms of the squared coefficient of variation (variance divided by the square of the mean), see Chen and Whitt (2019). For the extremal distributions with  $(\mu_V, d_V, a_V, b_V) = (\rho, d_V, 0, b_V)$  and  $(\mu_U, d_U, a_U, b_U) = (1, d_U, 0, b_U)$  this gives

$$c_V^2 = \frac{\sigma_V^2}{\mu_V^2} = \frac{d_V b_V}{2\rho^2}, \quad c_U^2 = \frac{\sigma_U^2}{\mu_U^2} = \frac{d_U b_U}{2}.$$

**Table EC.1** Bounds for  $\mathbb{E}[W]$  for  $(\mu_U, d_U, a_U, b_U) = (1, 1, 0, 10)$  and  $(\mu_V, d_V, a_V, b_V) = (\rho, 0.1, 0, 10)$ .

$\rho$	Tight (Thm. 4)	C & W (EC.14)	Daley (EC.13)	Kingman (EC.12)
0.1	0.45179	0.77780	0.80556	3.05556
0.2	0.63077	1.31953	1.43750	3.43750
0.5	2.03141	3.63750	4.25000	5.50000
0.7	5.81373	7.39989	8.41667	9.16667
0.8	10.47728	12.02090	13.25000	13.75000
0.9	24.28220	25.80400	27.25000	27.50000
0.95	51.79564	53.31910	54.87500	55.00000
0.99	271.80153	273.33100	274.97500	275.00000

Fixing the squared coefficient of variations  $c_V^2$  and  $c_U^2$  is equivalent with choosing the MADs as

$$d_V = \frac{2\rho^2 c_V^2}{b_V}, \quad d_U = \frac{2c_U^2}{b_U}. \quad (\text{EC.15})$$

We next present in Tables EC.2-EC.7 some further numerical results, for  $c_U^2 = c_V^2 = 0.5$  and  $c_U^2 = c_V^2 = 4$ .

**Table EC.2** Bounds for  $\mathbb{E}[W]$  for  $(\mu_V, d_V, a_V, b_V) = (\rho, d_V, 0, 10)$  and  $(\mu_U, d_U, a_U, b_U) = (1, d_U, 0, 10)$  with  $d_V, d_U$  as in (EC.15) and  $c_U^2 = c_V^2 = 0.5$ .

$\rho$	Tight (Thm. 4)	C & W (EC.14)	Daley (EC.13)	Kingman (EC.12)
0.1	0.00785	0.05278	0.05555	0.28055
0.2	0.02230	0.11320	0.12500	0.32500
0.5	0.14921	0.43875	0.50000	0.62500
0.7	0.48818	1.06499	1.16667	1.24167
0.8	0.99509	1.87709	2.00000	2.05000
0.9	2.85149	4.35540	4.50000	4.52500
0.95	7.29378	9.34441	9.50000	9.51250
0.99	46.78335	49.33560	49.50000	49.50250

**Table EC.3** Bounds for  $\mathbb{E}[W]$  for  $(\mu_V, d_V, a_V, b_V) = (\rho, d_V, 0, 10)$  and  $(\mu_U, d_U, a_U, b_U) = (1, d_U, 0, 10)$  with  $d_V, d_U$  as in (EC.15) and  $c_U^2 = c_V^2 = 4$ .

$\rho$	Tight (Thm. 4)	C & W (EC.14)	Daley (EC.13)	Kingman (EC.12)
0.1	0.09358	0.42224	0.44444	2.24444
0.2	0.26429	0.90562	1.00000	2.60000
0.5	2.05142	3.51000	4.00000	5.00000
0.7	6.76335	8.51991	9.33333	9.93333
0.8	13.18168	15.01670	16.00000	16.40000
0.9	32.95685	34.84320	36.00000	36.20000
0.95	72.84232	74.75520	76.00000	76.10000
0.99	392.74278	394.68400	396.00000	396.02000

**Table EC.4** Bounds for  $\mathbb{E}[W]$  for  $(\mu_V, d_V, a_V, b_V) = (\rho, d_V, 0, 10)$  and  $(\mu_U, d_U, a_U, b_U) = (1, d_U, 0, 10)$  with  $d_V, d_U$  as in (EC.15),  $c_U^2 = 4$  and  $c_V^2 = 0.5$ .

$\rho$	Tight (Thm. 4)	C & W (EC.14)	Daley (EC.13)	Kingman (EC.12)
0.1	0.07003	0.40280	0.42500	2.22500
0.2	0.15280	0.81812	0.91250	2.51250
0.5	0.91273	2.63500	3.12500	4.12500
0.7	3.73777	5.66158	6.47500	7.07500
0.8	7.53710	9.41674	10.40000	10.80000
0.9	18.82048	20.66820	21.82500	22.02500
0.95	41.31986	43.16770	44.41250	44.51250
0.99	221.30939	223.16700	224.48200	224.50200

**Table EC.5** Bounds for  $\mathbb{E}[W]$  for  $(\mu_V, d_V, a_V, b_V) = (\rho, d_V, 0, 10)$  and  $(\mu_U, d_U, a_U, b_U) = (1, d_U, 0, 10)$  with  $d_V, d_U$  as in (EC.15),  $c_U^2 = 0.5$  and  $c_V^2 = 4$ .

$\rho$	Tight (Thm. 4)	C & W (EC.14)	Daley (EC.13)	Kingman (EC.12)
0.1	0.02599	0.07222	0.07500	0.30000
0.2	0.10463	0.20070	0.21250	0.41250
0.5	1.00498	1.31375	1.37500	1.50000
0.7	3.39670	3.92332	4.02500	4.10000
0.8	6.81534	7.47709	7.60000	7.65000
0.9	17.72431	18.53040	18.67500	18.70000
0.95	40.05188	40.93190	41.08750	41.10000
0.99	219.91292	220.85300	221.01700	221.02000

**Table EC.6** Bounds for  $(1 - \rho)\mathbb{E}[W]/\rho$  for  $(\mu_V, d_V, a_V, b_V) = (\rho, d_V, 0, 10)$  and  $(\mu_U, d_U, a_U, b_U) = (1, d_U, 0, 10)$  with  $d_V, d_U$  as in (EC.15) and  $c_U^2 = c_V^2 = 0.5$ .

$\rho$	Tight (Thm. 4)	C & W (EC.14)	Daley (EC.13)	Kingman (EC.12)
0.1	0.07070	0.47502	0.50000	2.52500
0.2	0.08922	0.45281	0.50000	1.30000
0.5	0.14921	0.43875	0.50000	0.62500
0.7	0.20922	0.45642	0.50000	0.53214
0.8	0.24877	0.46927	0.50000	0.51250
0.9	0.31683	0.48393	0.50000	0.50277
0.95	0.38388	0.49181	0.50000	0.50065
0.99	0.47255	0.49833	0.50000	0.50002

**Table EC.7** Bounds for  $(1 - \rho)\mathbb{E}[W]/\rho$  for  $(\mu_V, d_V, a_V, b_V) = (\rho, d_V, 0, 10)$  and  $(\mu_U, d_U, a_U, b_U) = (1, d_U, 0, 10)$  with  $d_V, d_U$  as in (EC.15) and  $c_U^2 = c_V^2 = 4$ .

$\rho$	Tight (Thm. 4)	C & W (EC.14)	Daley (EC.13)	Kingman (EC.12)
0.1	0.84228	3.80016	4.00000	20.20000
0.2	1.05719	3.62248	4.00000	10.40000
0.5	2.05142	3.51000	4.00000	5.00000
0.7	2.89858	3.65139	4.00000	4.25714
0.8	3.29542	3.75418	4.00000	4.10000
0.9	3.66187	3.87146	4.00000	4.02222
0.95	3.83381	3.93449	4.00000	4.00526
0.99	3.96710	3.98671	4.00000	4.00020

**Table EC.8** Bounds for  $(1 - \rho)\mathbb{E}[W]/\rho$  for  $(\mu_V, d_V, a_V, b_V) = (\rho, d_V, 0, 10)$  and  $(\mu_U, d_U, a_U, b_U) = (1, d_U, 0, 10)$  with  $d_V, d_U$  as in (EC.15),  $c_U^2 = 4$  and  $c_V^2 = 0.5$ .

$\rho$	Tight (Thm. 4)	C & W (EC.14)	Daley (EC.13)	Kingman (EC.12)
0.1	0.63030	3.62516	3.82500	20.02500
0.2	0.61120	3.27248	3.65000	10.05000
0.5	0.91273	2.63500	3.12500	4.12500
0.7	1.60190	2.42639	2.77500	3.03214
0.8	1.88427	2.35418	2.60000	2.70000
0.9	2.09116	2.29646	2.42500	2.44722
0.95	2.17473	2.27199	2.33750	2.34276
0.99	2.23545	2.25421	2.26750	2.26770

**Table EC.9** Bounds for  $(1 - \rho)\mathbb{E}[W]/\rho$  for  $(\mu_V, d_V, a_V, b_V) = (\rho, d_V, 0, 10)$  and  $(\mu_U, d_U, a_U, b_U) = (1, d_U, 0, 10)$  with  $d_V, d_U$  as in (EC.15),  $c_U^2 = 0.5$  and  $c_V^2 = 4$ .

$\rho$	Tight (Thm. 4)	C & W (EC.14)	Daley (EC.13)	Kingman (EC.12)
0.1	0.23392	0.65002	0.67500	2.70000
0.2	0.41852	0.80281	0.85000	1.65000
0.5	1.00498	1.31375	1.37500	1.50000
0.7	1.45573	1.68142	1.72500	1.75714
0.8	1.70384	1.86927	1.90000	1.91250
0.9	1.96937	2.05893	2.07500	2.07778
0.95	2.10799	2.15431	2.16250	2.16316
0.99	2.22134	2.23084	2.23250	2.23253