# Stationary Birth-and-Death Processes Fit to Queues with Periodic Arrival Rate Functions

James Dong    and   Ward Whitt

School of Operations Research and Information Engineering,
Cornell University, Ithaca, NY 14850 jd748@cornell.edu

Industrial Engineering and Operations Research
Columbia University, New York, NY, 10027 ww2040@columbia.edu

January 17, 2015

### Abstract

To better understand how to interpret birth-and-death (BD) processes fit to service system data, we investigate the consequences of fitting a BD process to a multi-server queue with a periodic time-varying arrival rate function. We study how this fitted BD process is related to the original queue-length process. If a BD process is fit to a segment of the sample path of the queue-length process, with the birth (death) rates in each state estimated by the observed number of arrivals (departures) in that state divided by the total time spent in that state, then under minor regularity conditions that BD process has the steady-state distribution of the queue length process in the original $M_t/GI/s$ queueing model as the sample size increases. The steady-state distribution can be estimated efficiently by fitting a parametric function to the observed birth and death rates.

*Keywords:* birth-and-death processes; grey-box stochastic models; fitting stochastic models to data; queues with time-varying arrival rate; speed ratio; transient behavior.

1

# 1   Introduction

A commonly applied queueing model to analyze the performance of service systems is the $M/M/s+$ $M$ Erlang A model; see [16, 32] and references therein. It is a stationary birth-and-death (BD) process with four parameters: the arrival rate $\lambda$, the service rate $\mu$, the number of servers $s$ and the individual customer abandonment rate from queue $\theta$. The familiar $M/M/s/0$ Erlang B (loss) and $M/M/s \equiv M/M/s/\infty$ Erlang C (delay) models are the special cases in which $\theta = \infty$ and $\theta = 0$. These models are convenient because there are so few parameters. The arrival rate $\lambda$ and service rate $\mu$ can quickly be estimated as the reciprocals of the average interarrival time and service time, respectively, but the abandonment rate is more complicated because of censoring; it is often better to estimate the hazard rate; see [2].

For successful applications, it is important to investigate to what extent the model is consistent with service system data. This is most often done by estimating the distributions of the interarrival times and service times to see if they are nearly exponential, but there are many other ways the system can differ from the model. Service systems typically have time-varying arrival rates and there may be significant dependence among interarrival times and service times. The number of servers may vary over time as well and the servers are often actually heterogeneous [15]. Indeed, careful statistical analysis of service system data can be quite complicated, e.g., see [2, 21, 22, 23].

In this paper we investigate an alternative way to fit the Erlang A model to data: We may do that by directly fitting a state-dependent BD process. Following common practice [37], we can estimate the birth rate in state $k$ from data over an interval $[0, t]$ by $\bar{\lambda}_k \equiv \bar{\lambda}_k(t)$, the number of arrivals observed in that state, divided by the total time spent on that state, while the death rate in state $k$ is estimated by $\bar{\mu}_k \equiv \bar{\mu}_k(t)$, the number of departures observed in that state, divided by the total time spent on that state. For a BD process, those are the maximum likelihood estimators of the actual birth and death rates.

Given that the data are from the Erlang A model, we will see simple linear structure in the estimated birth and death rates. With enough data, we will see that

$$\bar{\lambda}_k = \lambda, \quad k \geq 0, \quad \text{and} \quad \bar{\mu}_k = (k \wedge s)\mu + (k - s)^+\theta, \quad k \geq 1, \tag{1}$$

where $a \wedge b \equiv \min\{a, b\}$ and $(a)^+ \equiv \max\{a, 0\}$. By this procedure, we can estimate all four parameters and test if the model is appropriate. A direct BD fit of the form (1) may indicate that the model should be effective even though some other tests fail. For example, experience indicates that a good model fit can occur by this BD rate fit even though the servers are heterogeneous and the service-time distribution is not exponential. Moreover, in those cases we may find that the Erlang A model works well in setting staffing levels.

However, what do we conclude if the BD fit does not yield the birth and death rate functions in (1)? Some insights are relatively obvious. For example, if we do not see death rates with two linear pieces joined at some level $s$, then we can judge that the number of servers probably was not constant during the measurement period. But it remains to carefully evaluate how to interpret departures from the simple Erlang structure in (1).

We might also consider directly applying the fitted BD process even if we do not see the Erlang A structure in (1), because BD processes are remarkably tractable. If we happen to find piecewise-linear fits, then we may find diffusion approximations with large scale, as in [3], which is not limited to the classical Erlang models in [16, 19]. It is well known that we can calculate the steady-state distribution of a general BD process by solving local balance equations. Less well known is the fact that we can efficiently calculate first passage time distributions in general BD processes [1]. But we should remember that the actual process may not be a general BD process. Our purpose here is to gain further insight into what the fitted BD rates do imply for the original process.

We started in [10] by looking carefully at BD fits to the number in system in $GI/GI/s$ queues. We continue here by looking carefully at BD fits to the number in system in $M_t/GI/s$ queues, having nonhomogeneous Poisson processes (NHPP's) as arrival processes with sinusoidal arrival rate functions, paying especial attention to the case of $s = \infty$ servers.

## 1.1 The Steady-State Distribution

As usual, the steady-state distribution of the fitted BD model, denoted by $\bar{\alpha}_k^e \equiv \bar{\alpha}_k^e(t)$ (with superscript $e$ indicating the estimated rates), is well defined (under regularity conditions [33]) and characterized as the unique probability vector satisfying the local balance equations,

$$\bar{\alpha}_k^e \bar{\lambda}_k = \bar{\alpha}_{k+1}^e \bar{\mu}_{k+1}, \quad k \geq 0. \tag{2}$$

To obtain reasonable rate estimates for which $\bar{\alpha}_k^e$ is indeed well defined and unique, we truncate the state space to a region of states that are visited relatively frequently. Throughout this paper, we assume that the limiting values of the rates as $t \to \infty$ exist so we omit the $t$. We use large sample sizes in our simulations to justify this assumption.

In [33] we cautioned against drawing unwarranted positive conclusions if the fitted BD steady-state distribution $\{\bar{\alpha}_k^e : k \geq 0\}$ in (2) closely matches the empirical steady-state distribution, $\{\bar{\alpha}_k : k \geq 0\}$, where $\bar{\alpha}_k \equiv \bar{\alpha}_k(t)$ is the proportion of total time spent in each state, because these two distribution are automatically closely related. Indeed, as has been known for some time (e.g., see Chapter 4 of [13]), under regularity conditions, these two distributions coincide asymptotically as $t$ (and thus the sample size) increases, even if the actual system evolves in a very different way from the fitted BD process. For example, the actual process $\{Q(t) : t \geq 0\}$ might be non-Markovian (as in [10]) or have a time-varying arrival rate (as here). Stochastic comparisons between the two distributions, depending on the beginning and ending states, were also derived in [33]. If the ending state coincides with the initial state, then these two empirical distributions are identical for any sample size!

## 1.2 Grey-Box Stochastic Modeling

Even though a close match between the empirical steady-state distribution, $\{\bar{\alpha}_k\}$, and the steady-state distribution of the fitted BD model, $\{\bar{\alpha}_k^e\}$, does not nearly imply that the actual system evolves as a BD process, we think that the fitted BD model has the potential to become a useful modeling and analysis tool, providing insight into the actual system. Of course, if the actual system can be well modeled by a standard BD model, such as one of the classical Erlang models, then we will see a good fit to that model with enough data. Of primary interest here is to be able to see deviations from classical models through the fitted birth and death rates. Actual service systems may have complex time-dependence and stochastic dependence that may be difficult to assess directly. Fitting a BD process may be a useful way to probe into system data. In [10] we referred to this as "grey-box stochastic modeling."

In [10] we applied this analysis to various conventional $GI/GI/s$ queueing models. We saw how the fitted rates differ from the corresponding $M/M/s$ model. We also saw that they differed in systematic ways that enabled us to see a "signature" of the $G/G/s$ model. Here we consider many-server $M_t/GI/s$ queueing models with sinusoidal periodic arrival rate functions. Now we find significant differences in the fitted birth rates from what we saw before for the $GI/GI/s$ models. And we see a signature of the $M_t/GI/s$ model with sinusoidal arrival rates. The results for the basic stochastic model with periodic arrival rate functions here should be useful to compare to similar analyses of service system data, such as hospital occupancy levels, where the arrival rates

have periodic structure over the days of each week and over the hours of each day; see [36]. Indeed, preliminary analysis of the hospital data shows striking differences, which should not be surprising, because hospitals tend not to be well modeled as standard queueing models. Just as in [10], we see telling structure in the fitted birth and death rates. From such empirical plots, we can recognize both consistency and deviations from basic models, such as the $M/M/s$ Erlang delay model, its $GI/GI/s$ extension and the associated $M_t/GI/s$ model with a periodic time-varying arrival rate function.

## 1.3 Operational Analysis

The present study is related to early work on operational analysis. In early performance analysis of computer systems, Buzen and Denning [4, 5, 9] advocated working with BD processes fit directly to data as part of a general operational analysis directly. The goal was to understand performance empirically, directly from data, without using customary stochastic models. Key support for this approach was provided by conservation laws that must hold among the statistics collected, as in Little's law.

However, we prefer to think of there actually being an underlying stochastic model. With that in mind, the fitted BD process provides partial information about the underlying model. Problems with a direct application of operational analysis are discussed in §§4.6-4.7 in [13]. In that context, though, [10] and this paper provides the first comparison between an underlying stochastic process model and the operational analysis BD model fit to data. For either to be useful in prediction, the future system of interest should be like the current system being measured. To judge whether candidate models are appropriate, we think that it is appropriate to apply statistical analysis to analyze the measurements. Sound statistical analysis, as in [2, 22, 23], can strongly support an underlying stochastic model, which will behave differently from the fitted BD model if the data are inconsistent with the BD model, as we show here.

## 1.4 Periodic Queues

Our goal in the present paper is to consider many-server queues with periodic arrival rates. These have been studied in [8, 11, 12, 14, 20, 24, 25, 27, 28, 34] and references therein. As in [10], we want to understand how the fitted birth and death rates depend on the model structure. We find that the fitted birth and death rates provide very useful information about the structure of the actual model. In this paper we concentrate on $M_t/GI/s$ multi-server queues, where the arrival process is a nonhomogeneous Poisson process (NHPP) with a periodic arrival rate function, emphasizing the tractable limiting case of the infinite-server (IS) model [11, 12]. For these models, there is a proper steady-state distribution, which is the time average of the time-dependent distributions over each periodic cycle. For the special case of the $M_t/M/\infty$ model with a sinusoidal arrival rate function, the steady-state distribution is studied in [35].

There are very few available results for actually computing the steady-state distribution in periodic queues. For Markovian models, the steady-state distribution may be calculated by numerically solving ordinary differential equations, possibly simplified by closure approximations [29], However, simulation seems to be the only available method for non-Markovian models. Thus, a significant contribution in this paper is to provide a new way to estimate the steady-state distribution; see §2.8. We suggest fitting parametric functions to estimated birth and death rates and then solving the local balance equations in (2). This approach has potential because the fitted birth rates and death rates often have more elementary structure, such as linearity.

We start in §2 by reporting results of simulation experiments for $M_t/GI/s$ queueing models

with sinusoidal arrival rates, which are directly of interest and serve to modify theoretical results that follow. In §3 and §4 we develop supporting theory. In §5 we draw conclusions.

## 2    Simulation Experiments

All the models considered in this paper will be $M_t/GI/s$ queueing models, having a nonhomogeneous Poisson process (NHPP, the $M_t$) as an arrival process, independent of i.i.d. service times distributed as a random variable $S$ with mean $E[S] = 1/\mu = 1$, $s$ servers, $1 \leq s \leq \infty$, and unlimited waiting space. Moreover, we consider the stylized sinusoidal arrival rate function

$$\lambda(t) \equiv \bar{\lambda}\left(1 + \beta \sin\left(\gamma t\right)\right), \tag{3}$$

where the cycle is $c = 2\pi/\gamma$. There are three parameters: (i) the average arrival rate $\bar{\lambda}$, (ii) the relative amplitude $\beta$ and (iii) the time scaling factor $\gamma$ or, equivalently the cycle length $c = 2\pi/\gamma$. Our base model is the $M_t/M/\infty$ model, which is the special case of the $M_t/GI/s$ model in which $s = \infty$, $S$ has an exponential distribution and $\beta = 10/35$.

### 2.1    Designing the Simulation Experiments

The simulation experiments were conducted much as in the prequel to this paper [10]. We generated the NHPP arrival process by thinning a Poisson process with rate equal to the maximum arrival rate over a sine cycle. Since we use relative amplitude $\beta = 10/35$, with $\bar{\lambda} = 35$ a proportion $10/(35+10) = 10/45 = 0.222$ of the potential arrivals were not actual arrivals. The fitted birth and death rates as well as the empirical mass function were estimated using 30 independent replications of 1.5 million potential arrivals before thinning. Overall, that means about 35 million arrivals in each experiment. Multiple i.i.d. repetitions were performed to confirm high accuracy within the regions shown. In order to compare the transient behavior of the fitted BD process to the original process, we simulated a separate version of the fitted BD process in a similar manner. To compute the first passage times starting from steady state (see §2.6), the process is initialized in steady state by choosing the initial state from the estimated steady-state distribution.

### 2.2    Comparing the Fitted Rates in the $M_t/M/\infty$ and $GI/M/\infty$ Models

Our main hypothesis is that the fitted birth and death rates can reveal features of the underlying model. To compare the impact of predictable deterministic variability in the arrival process, as manifested in a time-varying arrival rate function, to stochastic variability, we see how the fitted birth rates differ in the $M_t/M/\infty$ infinite-server (IS) model with a sinusoidal arrival rate function and the stationary $GI/M/\infty$ model with a renewal process having an interarrival time more variable than the exponential distribution. (When the service-time distribution is exponential, the fitted death rates coincide with the exact death rates in both cases, i.e., $\bar{\mu}_k^e(\infty) = k$; see Theorem 3.1 of [10] and Theorem 3.3 here.) However, the fitted birth rates are revealing.

In [10] we found that, when the actual arrival rate is $n$ (provided that $n$ is not too small), with the service rate fixed at $\mu = 1$, the fitted birth rates in state $k$, denoted by $\lambda_{n,k}$, tended to have the form

$$\bar{\lambda}_{n,k}^e = (n + b(k - n)) \vee 0, \tag{4}$$

where $b$ is a constant such that $-1 < b < 1$ and

$$b \approx 1 - \frac{2}{1 + c_a^2}, \tag{5}$$

with $c_a^2$ being the *squared coefficient of variation* (scv, variance divided by the square of the mean) of the interarrival-time distribution of the renewal arrival process. This is illustrated in Figure 1, which shows the fitted birth rates and death rates in five $GI/M/\infty$ models with arrival rate $\lambda = 39$ and service rate $\mu = 1$. The five interarrival-time distributions are Erlang $E_4$, $E_2$, $M$, and hyperexponential, $H_2$ with $c_a^2 = 2$ and $c_a^2 = 4$.

Figure 1 shows that the fitted birth rates tend to be approximately linear (over the region where the process visits relatively frequently, so that there are ample data for the estimation), with $\lambda_{n,n} = n$ and slope increasing as the variability increases. This is consistent with greater variability in the arrival process leading to a a larger steady-state number in system. For $c_a^2 < 1$, the slope is negative; for $c_a^2 > 1$, the slope is positive. As $c_a^2$ increases to $\infty$, the slope approaches 1.

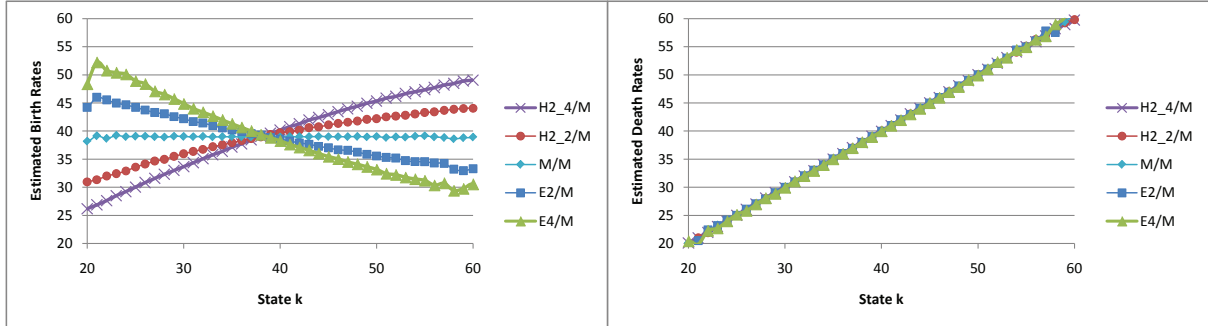

Figure 1: Fitted birth rates and death rates for five $G/M/\infty$ models with $\lambda = 39$ and $\mu = 1$.

We now consider the $M_t/M/\infty$ IS model with the sinusoidal arrival rate function in (3). Very roughly, we expect the predictable variability of a nonhomogeneous Poisson arrival process with a periodic arrival rate function to correspond approximately to a stationary model with a renewal arrival process having an interarrival-time distribution that is more variable than an exponential distribution [26]. That means we expect to see something like the fitted birth rates with increasing linear slopes in Figure 1. And indeed that is exactly what we do see, but restricted to a subinterval centered at the long-run average $\lambda_{n,n} = n$, as illustrated in Figure 2.
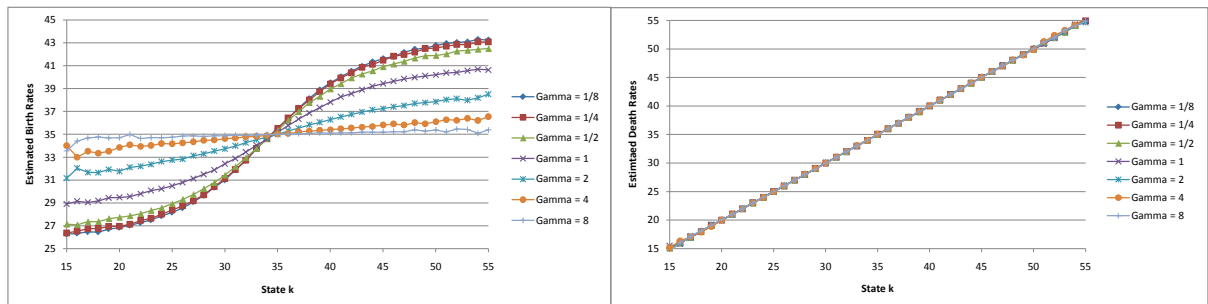


Figure 2: Fitted birth rates (left) and fitted death rates (right) for the $M_t/M/\infty$ model with the sinusoidal arrival rate function in (3) having parameters $\bar{\lambda} = 35$ and $\beta\bar{\lambda} = 10$ and 7 values of $\gamma$ ranging from 1/8 to 8.

The evolution of a BD queue primarily depends on the birth and death rates $\lambda_k$ and $\mu_k$ through their difference, the drift $\delta_k \equiv \lambda_k - \mu_k$, $k \geq 0$. Thus, we plot the drift functions associated with the $G/M/\infty$ and $M_t/M/\infty$ models in Figures 1 and 2 in Figure 3. These show that there is drift toward the overall mean in all cases, which is stronger when there is less variability.
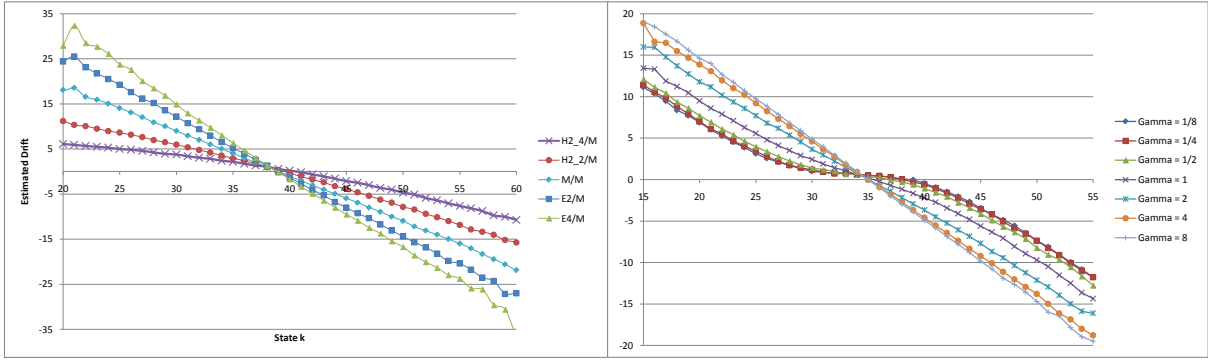
Figure 3: The estimated drift functions (birth rates minus death rates) for the $G/M/\infty$ model in Figure 1 (left) and the $M_t/M/\infty$ model in Figure 2 (right).

Similar results hold for models with finitely many servers. We show the results paralleling Figure 2 for the case of 40 servers in Figure 4.
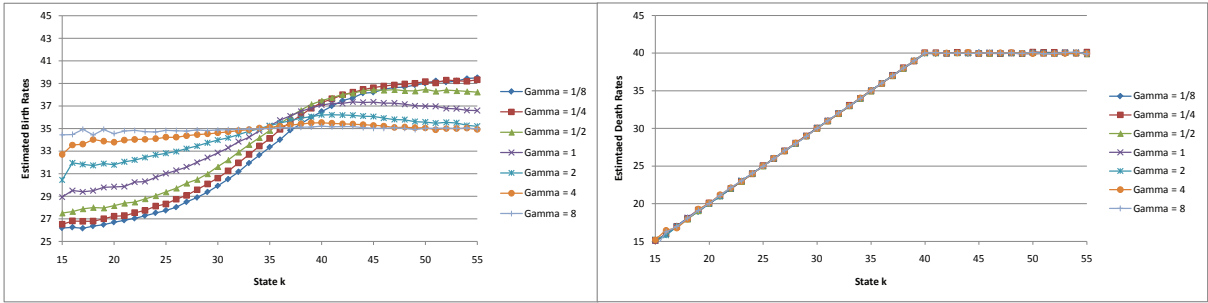


Figure 4: Fitted birth rates (left) and fitted death rates (right) for the $M_t/M/40$ queue with the sinusoidal arrival rate function in (3) having parameters $\bar{\lambda} = 35$ and $\beta\bar{\lambda} = 10$ and 7 values of $\gamma$ ranging from 1/8 to 8.

Figure 4 shows the piecewise-linear death rates, with two linear components, joined at the number of servers, that are characteristic of multi-server queues. Figure 2 of [10] displays similar plots for $GI/GI/s$ queues. However, the estimated birth rates in Figures 2 and 4 are unlike those of any $GI/GI/s$ queue. Theorems 4.4 and 4.5 establish finite bounds and heavy-traffic limits for the fitted birth rates, consistent with these figures.

## 2.3 The Steady-State Distribution of the $M_t/M/\infty$ Model

The estimated birth and death rates in §2.2 yield corresponding estimates of the steady-state distribution by solving the local balance equation (2). The estimated steady-state distributions for the $M_t/M/\infty$ model with the sinusoidal arrival rate function in (3) having parameters $\bar{\lambda} = 35$ and $\beta = 10/35$ for different ranges of $\gamma$ are shown in Figure 5. On the left (right) is shown different cases varying in a power of 10 (2). Many of the plots on the left coincide, so that we see convergence as $\gamma \uparrow \infty$ and as $\gamma \downarrow 0$. Indeed, the relevant ranges for intermediate behavior can be said to be $1/8 \le \gamma \le 8$ for these parameters $\bar{\lambda} = 35$ and $\beta = 10/35$, with the limits serving as effective approximations outside this interval.

The steady-state distribution of the number in system in the $M_t/M/\infty$ IS model with the sinusoidal arrival rate function in (3) is analyzed in [35] by applying [11]. By §5 of [11], the number
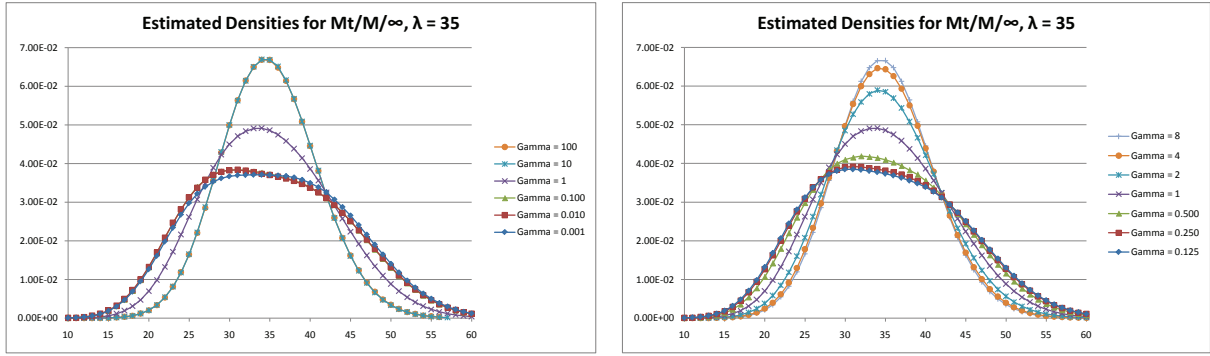
Figure 5: the estimated steady state number in the $M_t/M/\infty$ model with the sinusoidal arrival rate function in (3) having parameters $\bar{\lambda} = 35$ and $\beta = 10/35$ for different ranges of $\gamma$.

of customers in the system (or the number of busy servers), $Q(t)$, starting empty in the distant past, has a Poisson distribution at each time $t$ with mean

$$m(t) \equiv E[Q(t)] = \bar{\lambda}(1 + s(t)), \quad s(t) = \frac{\beta}{1 + \gamma^2}\left(\sin(\gamma t) - \gamma \cos(\gamma t)\right). \tag{6}$$

Moreover,

$$s^U \equiv \sup_{t \geq 0} s(t) = \frac{\beta}{\sqrt{1 + \gamma^2}} \tag{7}$$

and

$$s(t_0^m) = 0 \quad \text{and} \quad \dot{s}(t_0^m) > 0 \quad \text{for} \quad t_0^m = \frac{\cot^{-1}(1/\gamma)}{\gamma}. \tag{8}$$

The function $s(t)$ increases from 0 at time $t_0^m$ to its maximum value $s^U = \beta/\sqrt{1 + \gamma^2}$ at time $t_0^m + \pi/(2\gamma)$. The interval $[t_0^m, t_0^m + \pi/(2\gamma)]$ corresponds to its first quarter cycle.

Let $Z$ be a random variable with the steady-state probability mass function (pmf) of $Q(t)$; its pmf is a mixture of Poisson pmf's. In particular,

$$P(Z = k) = \frac{\gamma}{2\pi} \int_0^{2\pi/\gamma} P(Q(t) = k)\,dt, \quad k \geq 0, \tag{9}$$

The moments of $Z$ are given by the corresponding mixture

$$E[Z^k] = \frac{\gamma}{2\pi} \int_0^{2\pi/\gamma} E[Q(t)^k]\,dt, \quad k \geq 1, \tag{10}$$

so that $E[Z] = \bar{\lambda}$.

## 2.4   Transient Behavior

It should be evident that the transient behavior of the fitted BD process and the original process have significant differences. In particular, there is no periodicity in the fitted BD process. The differences are particularly striking with small $\gamma$, i.e., for long cycles $c(\gamma) = 2\pi/\gamma$. That is dramatically illustrated in Figure 6, which compares the sample paths of the number in system of the two processes for the $M_t/M/\infty$ queue with the sinusoidal arrival rate function in (3) having parameters $\bar{\lambda} = 100$ and $\beta = 10/35$ and $\gamma = 0.01$. Since $\gamma = 0.01$, the cycle length is 628. Hence in the time
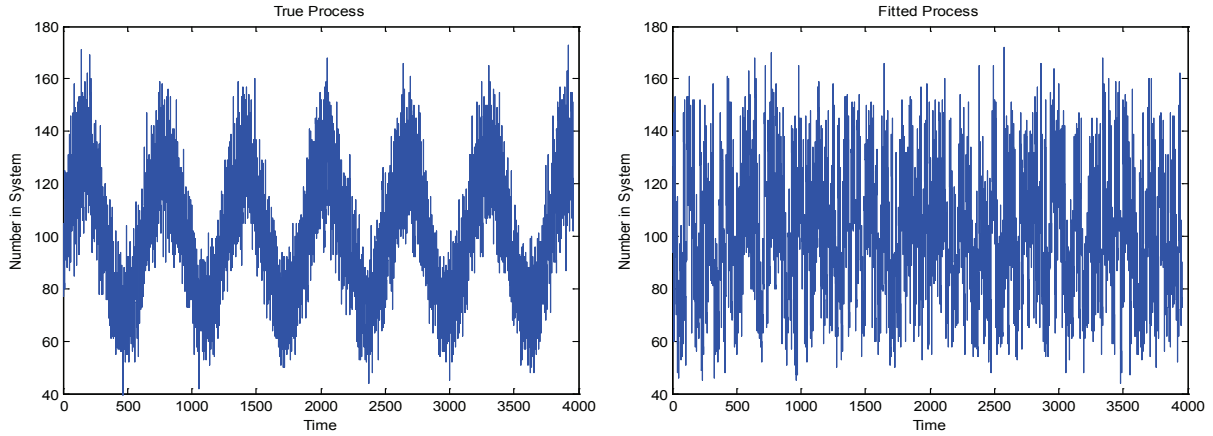
Figure 6: sample paths of the number in system for the original process (left) and the fitted BD process (right) for the $M_t/M/\infty$ queue with the sinusoidal arrival rate function in (3) having parameters $\bar{\lambda} = 100$ and $\beta = 10/35$ and $\gamma = 0.01$.

interval $[0, 4000]$ we see a bit more than six cycles, but there is no periodic behavior in the fitted BD process.

However, the sample paths are not always so strikingly different. Indeed, the sample paths get less different as $\gamma$ increases. Figures 7 and 8 illustrate by showing the sample paths for $\gamma = 1$ and $\gamma = 10$ over the interval $[0, 40]$. For $\gamma = 1$, there are again 6.28 sine cycles, but for $\gamma = 10$, there are 62.8 cycles. In these cases, the sample paths look much more similar. From Figures 7 and 8, we conclude that we might well use the fitted BD process to describe the transient behavior as well as the steady-state behavior for $\gamma \geq 1$, i.e., for relatively short cycles. Periodic arrival rates with short cycles often arise in practice in appointment-generated arrivals, where the actual arrivals are randomly distributed about the scheduled appointment times.
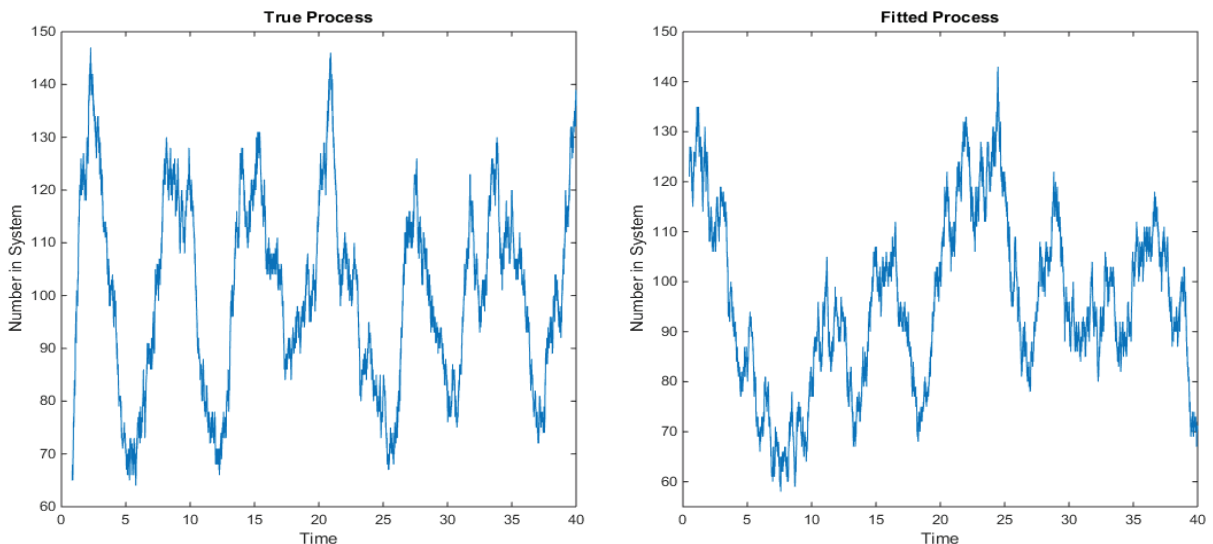


Figure 7: sample paths of the number in system for the original process (left) and the fitted BD process (right) for the $M_t/M/\infty$ queue with the sinusoidal arrival rate function in (3) having parameters $\bar{\lambda} = 100$ and $\beta = 10/35$ and $\gamma = 1.0$.
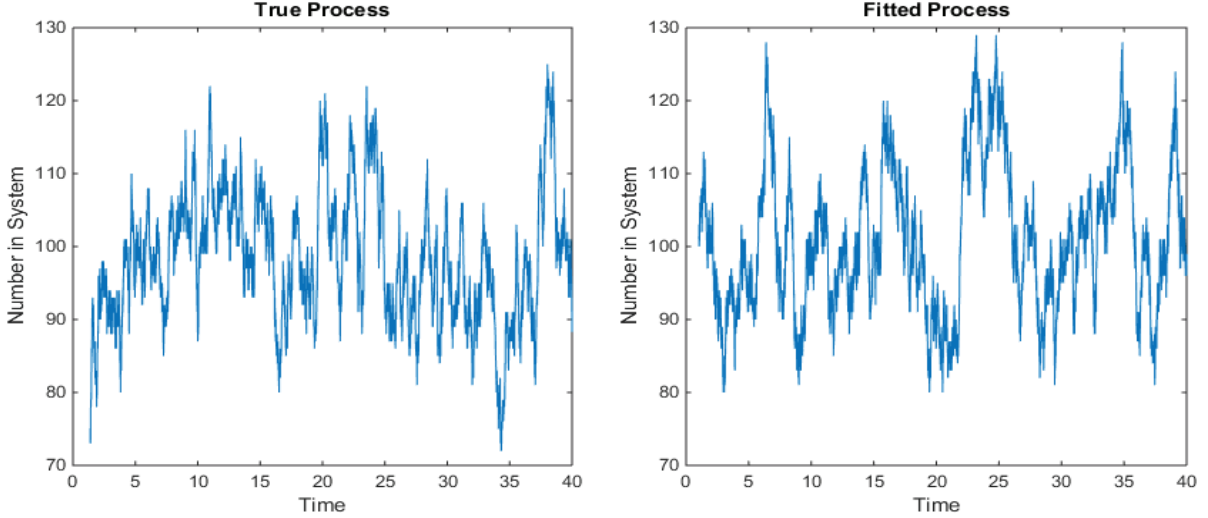
9

Figure 8: sample paths of the number in system for the original process (left) and the fitted BD process (right) for the $M_t/M/\infty$ queue with the sinusoidal arrival rate function in (3) having parameters $\bar{\lambda} = 100$ and $\beta = 10/35$ and $\gamma = 10$.

## 2.5 Limits for Small and Large $\gamma$

The behavior of the fitted BD process can be better understood by limits for the steady-state distribution of the $M_t/M/\infty$ model as $\gamma \uparrow \infty$ and as $\gamma \downarrow 0$. First, as $\gamma \uparrow \infty$, even though the arrival rate function oscillates more and more rapidly, the cumulative arrival rate function $\Lambda(t) \equiv \int_0^t \lambda(s)\, ds$ converges to the linear function $\bar{\lambda}t$. Consequently, the arrival process converges to a stationary Poisson process ($M$) with the average arrival rate $\bar{\lambda}$ and the steady-state number in system converges to the Poisson steady state distribution in associated the stationary $M/M/\infty$ model with mean $\bar{\lambda}$. That follows from Theorem 1 of [30] and references therein. As a consequence, as $\gamma \uparrow \infty$ we must have the fitted birth rates in the fitted BD process converge to the constant birth rates of a Poisson process, and that is precisely what we see as $\gamma$ increases in Figure 2.

Second, as $\gamma \downarrow 0$, the cycles get longer and longer, so that the system behaves at each time as a stationary model with the instantaneous arrival rate at that particular time. That is the perspective of the pointwise stationary approximation for queues with time-varying arrival rates [18], which is asymptotically correct for the $M_t/M/\infty$ model as $\gamma \downarrow 0$. That follows from Theorem 1 of [31]. As a consequence, as $\gamma \downarrow 0$ we must have the fitted birth rates in the fitted BD process converge to a proper limit, and that is precisely what we see as $\gamma$ increases in Figure 2.

The limit $Z_0$ of the steady-state variable $Z \equiv Z_\gamma$ as $\gamma \downarrow 0$ is the mixture of the steady-state distributions. That is, by combining the PSA limit with (9), we see that

$$P(Z_0 = k) = \frac{\gamma}{2\pi} \int_0^{2\pi} P(Q_0(t) = k)\, dt, \quad k \geq 0, \tag{11}$$

where $Q_0(t)$ has a Poisson distribution with mean $m_0(t) = \lambda_1(t)$, where we let $\gamma = 1$. In particular, this limit as $\gamma \downarrow 0$ becomes independent of $\gamma$.

These two limits can be seen by comparing the sample paths of the fitted BD processes for different $\gamma$. This is especially interesting for the long-cycle case. Figure 9 illustrates by showing the sample paths of the number in system for the fitted BD process in the $M_t/M/\infty$ queue with the

10

sinusoidal arrival rate function in (3) having parameters $\bar{\lambda} = 100$ and $\beta = 10/35$ and $\gamma = 0.1$ (left) and $\gamma = 0.01$ (right). The plots of different interval lengths show that the fitted BD processes are very similar.
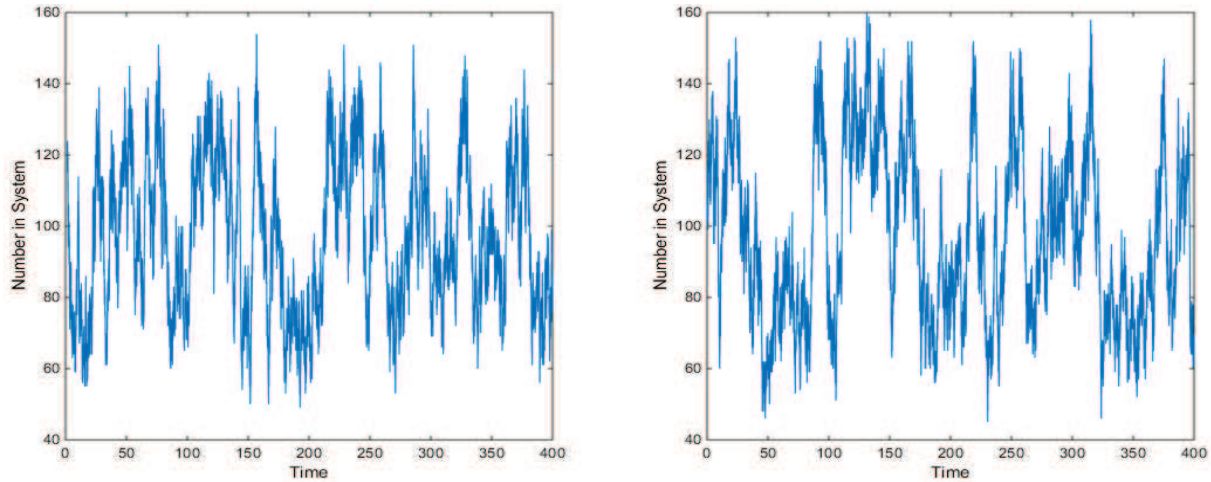


Figure 9: sample paths of the number in system for the fitted BD process in the $M_t/M/\infty$ queue with the sinusoidal arrival rate function in (3) having parameters $\bar{\lambda} = 100$ and $\beta = 10/35$ and $\gamma = 0.1$ (left) and $\gamma = 0.01$ (right).

## 2.6 Speed Ratios: Very Different Limits for the Finite-Server Models

The stationary Poisson limit as $\gamma \uparrow \infty$ is the same in $M_t/GI/s$ models with $s$ servers and general service times, but the limit as $\gamma \downarrow 0$ can be very different. Indeed, the limiting behavior will be very different if the finite-server model is overloaded with instantaneous traffic intensity $\rho(t) > 1$ at some time within its periodic cycle. If $\rho(t) > 1$ for some values of $t$ and if we make $\gamma$ very small, then these overload periods extend for longer and longer times, so that there can be a significant queue buildup. Indeed, proper limits as $\gamma \downarrow 0$ can only be obtained by adding additional scaling. This phenomenon is discussed in [7].
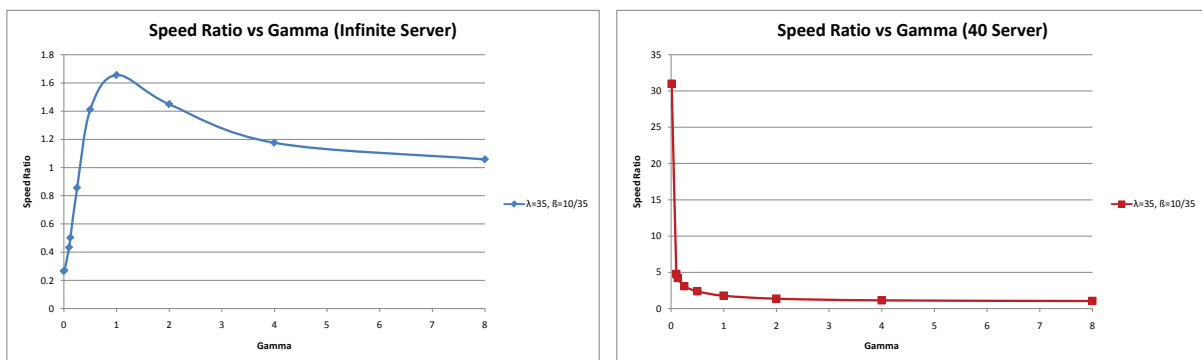


Figure 10: plots of the speed ratios in the $M_t/M/\infty$ (left) and $M_t/M/40$ models with the sinusoidal arrival rate function in (3) having parameters $\bar{\lambda} = 35$ and $\beta = 10/35$ as a function of the parameter $\gamma$.

The great difference as $\gamma \downarrow 0$ is illustrated by Figure 10, which plots the speed ratios for

11

the $M_t/M/\infty$ (left) and $M_t/M/40$ models with the sinusoidal arrival rate function in (3) having parameters $\bar{\lambda} = 35$ and $\beta = 10/35$ as a function of the parameter $\gamma$. The speed ratios were introduced in [10] to approximately characterize the transient behavior. Let $T(p,q)$ be the first passage time from the $p^{\text{th}}$ percentile of the steady-state distribution to the $q^{\text{th}}$ percentile of the steady-state distribution in the original process, and let $T_f(p,q)$ be the first passage time from the $p^{\text{th}}$ percentile of the steady-state distribution to the $q^{\text{th}}$ percentile of the steady-state distribution in the fitted BD process. These first passage times are fully specified for the fitted BD process because it is a Markov process, but they are not completely specified in the original model, because the stochastic process $\{Q(t) : t \geq 0\}$ is in general not Markov. Thus we need to specify the initial conditions. We understand the system to be in steady-state, so the initial condition is the steady-state distribution of the process conditional on starting at percentile $p$.

We in fact estimate the expected first passage times for the original process from simulations, by considering successive alternating visits to the $p^{\text{th}}$ and $q^{\text{th}}$ percentiles of the steady-state distribution. As an approximation, which we regard as reasonable as long as $p$ is not too close to $q$, we will assume that these successive first passage times are i.i.d. We estimate the expected values of these first passage times by sample averages and estimate 95% confidence intervals under the i.i.d. assumption. The rate at which these transitions occur can be defined by

$$r(p,q) \equiv \frac{1}{E[T(p,q)]} \quad \text{and} \quad r_f(p,q) \equiv \frac{1}{E[T_f(p,q)]}. \tag{12}$$

The associated $(p,q)$-*speed ratio* can be defined by

$$\omega(p,q) \equiv \frac{r(p,q)}{r_f(p,q)} = \frac{E[T_f(p,q)]}{E[T(p,q)]}. \tag{13}$$

To obtain further simplification, we assume that $q = 1 - p$ with $0 < p < 1/2$ and consider round trips, so that

$$T(p) = T(p, 1-p) + T(1-p, p) \quad \text{and} \quad T_f(p) = T_f(p, 1-p) + T_f(1-p, p), \tag{14}$$

$$r(p) \equiv \frac{1}{E[T(p)]} \quad \text{and} \quad r_f(p) \equiv \frac{1}{E[T_f(p)]} \tag{15}$$

and the *p-speed ratio* can be defined by

$$\omega(p) \equiv \frac{r(p)}{r_f(p)} = \frac{E[T_f(p)]}{E[T(p)]}. \tag{16}$$

Consistent with our previous discussion, Figure 10 shows that the speed ratios approach 1 as $\gamma$ increases, but we see very different behavior as $\gamma \downarrow 0$. The finite limit for the $M_t/M/\infty$ model confirms the limit of the steady-state distributions, whereas the divergence for the $M_t/M/40$ model shows the divergence of the 40-server models, due to the persistent overload over long time intervals.

## 2.7 Different Service Distributions

We have also conducted corresponding simulation experiments for the $M_t/GI/\infty$ model with non-exponential service-time distributions. Figure 11 shows the fitted rates for the $H_2$ service distributions with scv $c^2 = 2$ just as in §2 of [10]. The corresponding plots for the $E_2$ distribution are in the appendix; they look very similar. Figure 12 shows the associated steady-state mass functions for $H_2$ and $E_2$ service times.

Figure 13 shows that there are discernible differences among the speed ratios for the three service distributions, but the differences are not great.
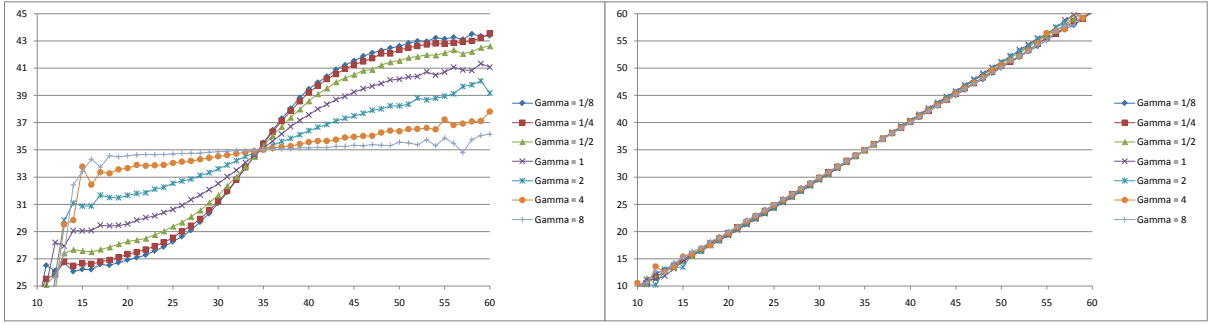
Figure 11: Fitted birth rates (left) and fitted death rates (right) for the $M_t/H_2/\infty$ model with the sinusoidal arrival rate function in (3) having parameters $\bar{\lambda} = 35$ and $\beta\bar{\lambda} = 10$ and 7 values of $\gamma$ ranging from 1/8 to 8. (The service scv is $c^2 = 2$.)
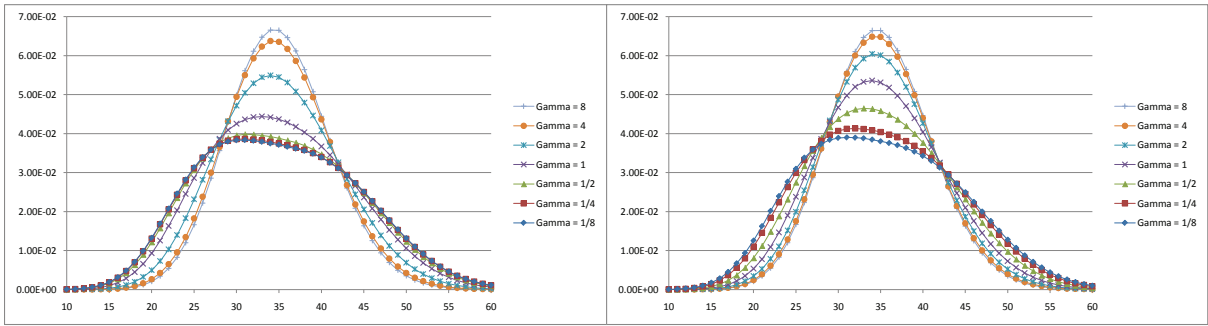


Figure 12: Fitted steady-state mass functions for the $M_t/H_2/\infty$ model (left) and the $M_t/E_2/\infty$ model (right) for with the sinusoidal arrival rate function in (3) having parameters $\bar{\lambda} = 35$ and $\beta\bar{\lambda} = 10$ and 7 values of $\gamma$ ranging from 1/8 to 8.

## 2.8 Estimating the Steady-State Distribution

In this section we investigate how we may efficiently estimate the steady-state distribution by fitting parametric functions to the estimated birth and death rates and then solve the local balance equation (2). First, for IS model we do not need to consider the death rates, because we have $\bar{\mu}_k = k$ throughout. Hence, we concentrate on the birth rates. For larger values of $\gamma$, a linear function works well, but not for smaller values of $\gamma$. As our parametric function, we choose

$$\lambda_k^p = a \arctan b(k - c) + d, \tag{17}$$

which is nondecreasing in $k$ with finite limits as $k$ increases and decreases, and has the parameter four-tuple $(a, b, c, d)$. We let $c = d = \bar{\lambda}$, so that leaves only the two parameters $a$ and $b$.

Figures 14, 15 and 16 show the fitted mass function and birth rates for the three gamma values: $\gamma = 1/8$, 1/2 and 2, respectively. These were constructed using the Matlab curve fitting toolbox, which fits by least squares. The figures show that the special arctangent function in (17) does much better than a linear fit for small $\gamma$, but a simple linear fit works well for large $\gamma$. The parameter pairs in the three cases were $(a, b) = (7.541, 0.125)$, $(6.682, 0.1253)$ and $(3.577, 0.0744)$, respectively. The main point is that a parametric fit based on only two parameters yields an accurate fit to a mass function that can be quite complicated.
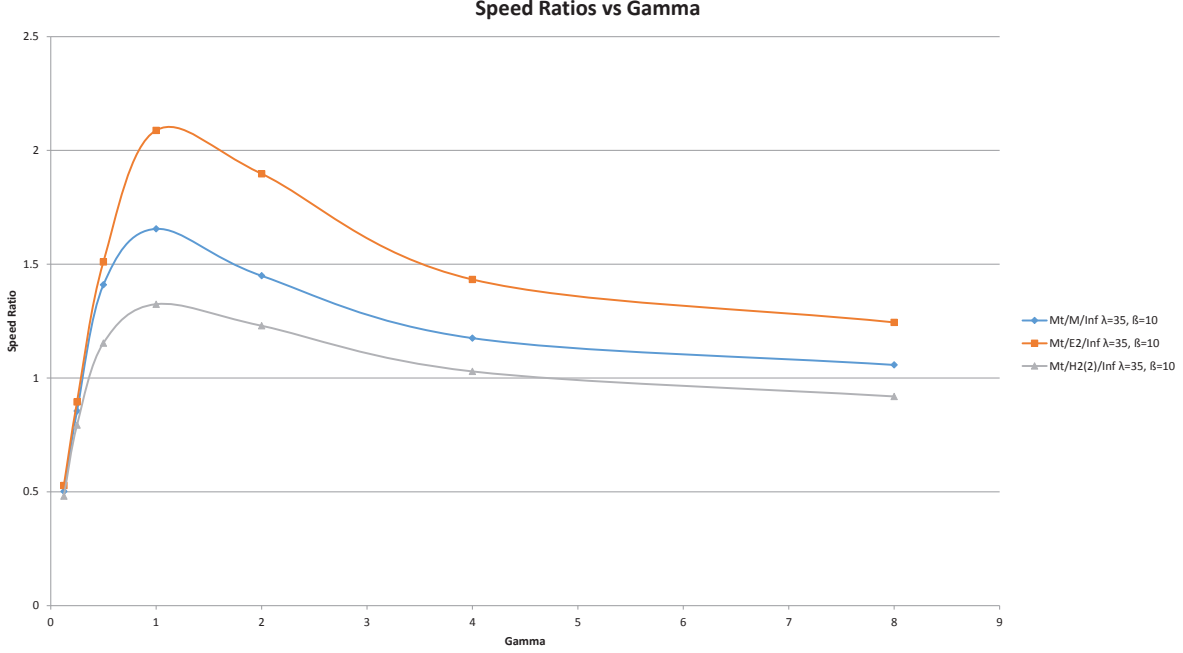
Figure 13: Speed ratios for the $M_t/GI/\infty$ model with the sinusoidal arrival rate function in (3) having parameters $\bar{\lambda} = 35$ and $\beta\bar{\lambda} = 10$ as a function of $\gamma$ for three different service distributions.

# 3   Supporting Theory: The Periodic $M_t/GI/s$ Queueing Model

We now develop supporting theory. Let $A(t)$ count the number of arrivals in the interval $[0, t]$. We assume that the arrival rate function $\lambda(t)$ is a periodic continuous function with periodic cycle of length $c$. Let $\bar{\lambda}$ be the long-run average arrival rate, with

$$\bar{\lambda} \equiv \frac{1}{c} \int_0^c \lambda(s)\,ds = \lim_{t\to\infty} \frac{A(t)}{t}. \tag{18}$$

Let the service times be distributed as a random variable $S$ with cumulative distribution function (cdf) $G$ and mean $E[S] \equiv 1/\mu < \infty$. Let the (long-run) traffic intensity be defined by $\bar{\rho} \equiv \bar{\lambda}E[S] = \lambda/\mu$.

Let $Q(t)$ denote the number of customers in the system at time $t$ and let $P(Q(t) = k)$, $k \geq 0$, be its time-dependent probability mass function. As indicated in [20], the stochastic process $\{Q(t) : t \geq 0\}$ is a regenerative processes, with the events $\{Q(nc + t) = 0\}$, $n \geq 1$, for any fixed $t$, $0 \leq t < c$, being regeneration times. As a consequence, we have a well defined periodic steady-state distribution when $\bar{\rho} < 1$.

**Theorem 3.1** (*periodic steady-state distribution*) *If $\bar{\rho} < 1$ in the periodic $M_t/GI/s$ queueing model, then $\alpha(t)$, $0 \leq t < c$ and $\alpha^c$ are well defined probability vectors with*

$$\alpha_k(t) \quad \equiv \quad \lim_{n\to\infty} P(Q(nc + t) = k) = \lim_{n\to\infty} \frac{1}{n}\sum_{j=1}^n 1_{\{Q(jc+t)=k\}}, \quad k \geq 0, \quad and$$

$$\alpha_k^c \quad \equiv \quad \frac{1}{c}\int_0^c \alpha_k(t)\,dt = \lim_{t\to\infty} \frac{1}{t}\int_0^t 1_{\{Q(s)=k\}}\,ds, \quad k \geq 0. \tag{19}$$

14

Figure 14: Fitted mass function (left) and birth rates (right) for the $M_t/M/\infty$ model with the sinusoidal arrival rate function in (3) having parameters $\bar{\lambda} = 35$, $\beta\bar{\lambda} = 10$ and $\gamma = 0.125$
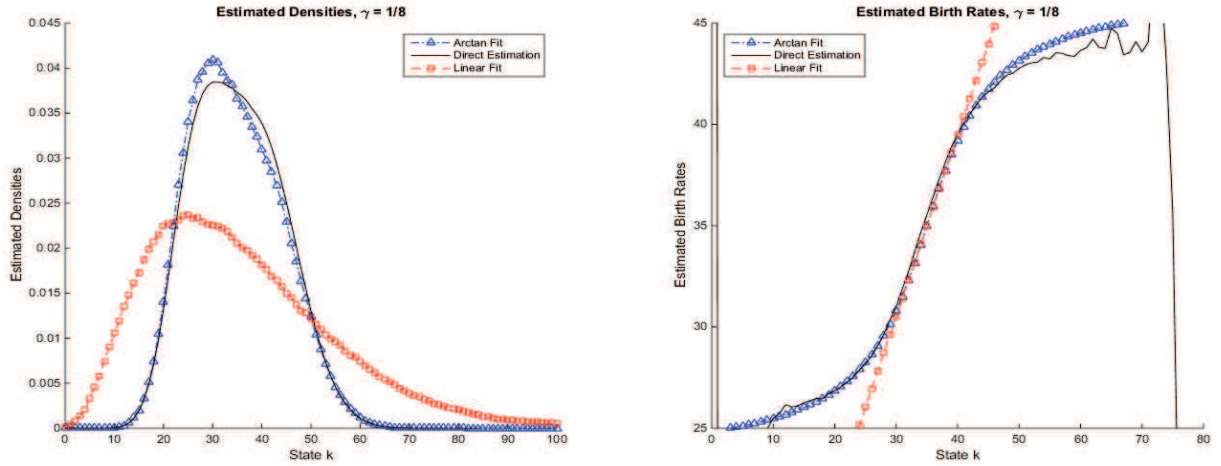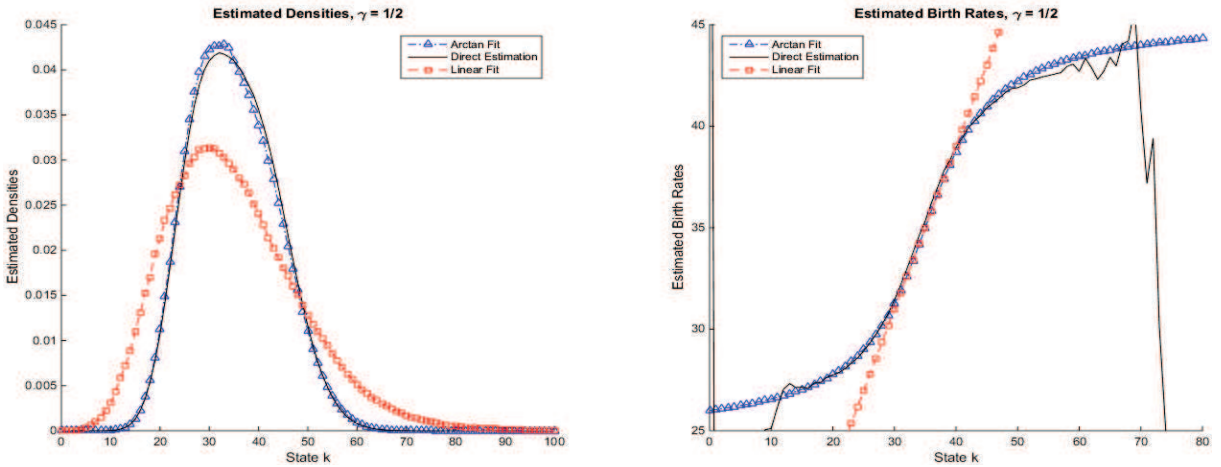


Figure 15: Fitted mass function (left) and birth rates (right) for the $M_t/M/\infty$ model with the sinusoidal arrival rate function in (3) having parameters $\bar{\lambda} = 35$, $\beta\bar{\lambda} = 10$ and $\gamma = 0.5$

Let $\bar{\lambda}_k^e(t)$ and $\bar{\mu}_k^e(t)$ be the estimated birth rate and death rate in state $k$ from data over $[0, t]$. In the $M_t/GI/s$ model, the arrival rate actually depends only on time, not the state. Hence, we can obtain the following explicit expressions for the asymptotic values as the sample size increases, $\bar{\lambda}_k(\infty)$ and $\bar{\mu}_k(\infty)$.

**Theorem 3.2** (*estimated birth and death rates with ample data*) *In the periodic $M_t/GI/s$ queueing model with $\bar{\rho} < 1$,*

$$\bar{\lambda}_k^e(\infty) = \frac{\int_0^c \alpha_k(t)\lambda(t)\,dt}{\int_0^c \alpha_k(t)\,dt} = \frac{\int_0^c \alpha_k(t)\lambda(t)\,dt}{c\alpha_k^c} \tag{20}$$

*and*

$$\bar{\mu}_{k+1}^e(\infty) = \frac{\alpha_k^c \bar{\lambda}_k^e(\infty)}{\alpha_{k+1}^c} = \frac{\int_0^c \alpha_k(t)\lambda(t)\,dt}{c\alpha_{k+1}^c}. \tag{21}$$

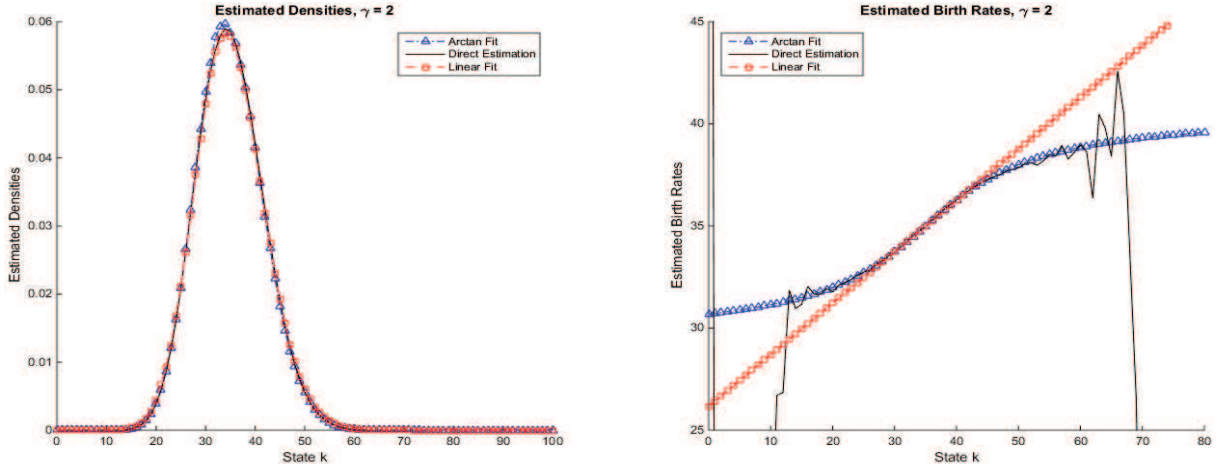*for $\alpha_k(t)$ and $\alpha_k^c$ in (19).*

Figure 16: Fitted mass function (left) and fitted birth rates (right) for the $M_t/M/\infty$ model with the sinusoidal arrival rate function in (3) having parameters $\bar{\lambda} = 35$, $\beta\bar{\lambda} = 10$ and $\gamma = 2.0$

**Proof.** Since the arrival rate depends only on time, we have (20). We then can apply the detailed balance equation in (2) to get (21). ∎

Theorems 3.1 and 3.2 can be applied in two ways. First, we can apply these theorems to learn about the fitted birth and death rates. They pose a strong constraint on the fitted birth and death rates because the detailed balance equation in (2) must hold. As a consequence, if we know either the fitted birth rates or the fitted death rates, then the others are determined as well. We will illustrate in our specific results below.

Second, we can apply the estimated birth and death rates to estimate the steady-state probability vector $\alpha^c$ in Theorem 3.1. Let $\bar{\alpha}^e(\infty)$ be the steady-state probability vector of the fitted BD process obtained from (2). Since $\bar{\alpha}^e$ coincides with $\alpha^c$ in (19), we can use the fitted BD model to calculate the steady-state distribution $\alpha^c$ in (19). To do so, we estimate the birth and death rates and then apply the detailed balance equation in (2). Moreover, by developing analytical approximations for the fitted birth and death rates, we succeed in developing an analytical approximation for $\alpha^c$.

## 3.1 The Periodic $M_t/M/s$ Model

For the special case of an exponential service-time distribution, i.e., for the $M_t/M/s$ model, the stochastic process $\{Q(t) : t \geq 0\}$ is Markov and more convenient explicit formulas are available.

We first observe that an analog of Theorem 3.1 of [10] also holds for the fitted death rates in the present time-varying case.

**Theorem 3.3** (*asymptotically correct death rates*) *For the periodic $M_t/M/s$ model with $\bar{\rho} < 1$, the fitted death rates are asymptotically correct as the sample size increases, i.e.,*

$$\bar{\mu}_k(\infty) = \min\{k, s\}\mu, \quad k \geq 0. \tag{22}$$

*Hence, the fitted birth rates can be expressed as*

$$\bar{\lambda}_k(\infty) = \frac{\alpha_{k+1}^c \min\{k+1, s\}\mu}{\alpha_k^c}, \quad k \geq 0, \tag{23}$$

*for $\alpha_k^c$ in (19).*

16

**Proof.** As for Theorem 3.1 of [10], (22) follows from the lack of memory property of the exponential distribution. However, we show that it is possible to directly apply Theorem 3.1 of [10] here. We use the fact that the $M_t/M/s$ model has a proper dynamic periodic steady-state distribution with a period equal to the period of the arrival process, cf. [20]. For that model we can convert the arrival process to a stationary point process by simply randomizing where we start in the first cycle. If the period is of length $d$, then we start the arrival process at time $t$, where $t$ is uniformly distributed over the interval $[0, d]$. That randomization converts the arrival process to a stationary point process, so that we can apply Theorem 3.1 of [10] (a). But then we observe that the randomization does not alter the limit (22). We then apply (2) to get (23). ∎

For the $M_t/M/s$ model, we are primarily interested in the fitted arrival rates $\bar{\lambda}_k^e(t)$, where the run length $t$ is sufficiently long that we can regard them as essentially the limiting values $\bar{\lambda}_k^e(\infty)$. We want to compare the fitted arrival rates in the $M_t/M/s$ model to the associated fitted arrival rates in the corresponding $M/M/s$ and $H_2/M/s$ models, where the average arrival rates and other parameters are hold fixed. We want to see if the fitted birth rates allow us to distinguish between extra stochastic variability, as illustrated by having an $H_2$ renewal arrival process instead of an $M$ Poisson arrival process, and extra time-variability, as illustrated by having an $M_t$ NHPP arrival process instead of an $M$ Poisson arrival process. Both of these can be contrasted with the constant arrival rate $\lambda$ with an $M$ arrival process.

We next observe that a geometric tail holds for the $M_t/M/s$ model with the same decay rate as for the associated stationary $M/M/s$ model with arrival rate $\bar{\lambda}$. Recall that a probability vector $\alpha$ has a geometric tail with decay rate $\sigma$ if

$$\alpha_k \sim \beta\sigma^k \quad \text{as} \quad k \to \infty, \tag{24}$$

i.e., if the ratio of the two sides converges to 1 as $k \to \infty$; see §3.2 of [10].

**Theorem 3.4** (*geometric tail*) *For the $M_t/M/s$ model with $s < \infty$ and $\bar{\lambda} < s\mu$, the periodic steady-state distribution has a geometric tail as in* (24) *with the same decay rate as in the associated stationary $M/M/s$ model with arrival rate $\bar{\lambda}$; i.e.,*

$$\bar{\alpha}_k(\infty) \sim \beta_t\sigma_t^k \quad as \quad k \to \infty, \tag{25}$$

*where*

$$\sigma_t = \sigma = \rho \equiv \frac{\bar{\lambda}}{s\mu} \quad and \quad \beta_t \geq \beta \geq (1 - \rho) \tag{26}$$

*with $(\beta, \sigma)$ and $(\beta_t, \sigma_t)$ denoting the asymptotic parameter pairs for the $M/M/s$ and $M_t/M/s$ models, respectively. As a consequence,*

$$\bar{\lambda}_k(\infty) \to \bar{\lambda} \quad as \quad k \to \infty. \tag{27}$$

**Proof.** The tail behavior can be deduced by considering bounding discrete-time processes, looking at the system at times $t_0 + kc$. Both systems are bounded below by the discrete-time model that has all arrivals in each interval at the end of the interval and all departures at the beginning of the interval, while both systems are bounded above by the discrete-time model that has all arrivals in each interval at the beginning of the interval and all departures at the end of the interval. These two-discrete time systems are random walks with steady-state distributions satisfying (24) with common decay factor $\sigma = \rho$. A step in the random walk is the difference of two Poisson random variables $U - D$, where $EU = \bar{\lambda}c$ and $ED = s\mu c$, which have ratio $EU/ED = \bar{\lambda}/s\mu$, which in turn determines the decay rate. A stochastic comparison [6] then implies that $\beta_t \geq \beta$. For the final

inequality in (26), we can compare the $M/M/s$ system to the corresponding $M/M/1$ model with a fast server, working at rate $s\mu$. The two systems have the same birth rate, while the $M/M/1$ system has death rates that are greater than or equal to those in the $M/M/s$ model. Hence, the steady-state distributions are ordered stochastically. Finally, the final limit in (27) follows from Theorem 3.3 and (25), where here $s\mu\sigma = s\mu\rho = \bar{\lambda}$. ∎

We remark in closing this section that the periodic $M_t/M/\infty$ has different tail behavior; hence the assumption that $s < \infty$. We next start considering the infinite-server model.

## 3.2 The Periodic Infinite-Server Model

We now consider the special case of the periodic $M_t/GI/\infty$ infinite-server (IS) model, because it admits many explicit formulas, as shown in [11, 12, 25]. We let the model start in the indefinite past, so that it can be regarded as in periodic steady-state at time 0. This is achieved by assuming an explicit form for the arrival rate function, as in (3), and then assuming that the system started empty in the indefinite past.

By Theorem 1 of [12], the number in system has a Poisson distribution for each $t$ with periodic mean function $m(t)$, with the same period $c$, where

$$m(t) = E[\lambda(t - S_e)]E[S] = E[S]\int_0^\infty \lambda(t - s)dG_e(s), \quad t \geq 0, \tag{28}$$

and $S_e$ is a random variables with the stationary-excess cdf $G_e$ associated with the service-time cdf $G$, i.e.,

$$G_e(t) \equiv P(S_e \leq t) \equiv \frac{1}{E[S]}\int_0^t (1 - G(s))\, ds, \quad t \geq 0. \tag{29}$$

Moreover, the departure process in the $M_t/GI/\infty$ model is a Poisson process with periodic rate function $\delta(t)$, with the same period $c$, where

$$\delta(t) = E[\lambda(t - S)] = \int_0^\infty \lambda(t - s)dG(s), \quad t \geq 0. \tag{30}$$

For the special case of a sinusoidal arrival rate function, an explicit expression for $m(t)$ is given in Theorem 4.1 of [11].

As a consequence, we have the following corollary to Theorem 3.1.

**Corollary 3.1** (*periodic steady-state distribution in the IS model*) *In the periodic $M_t/GI/\infty$ queueing model starting empty in the distant past, $\alpha(t)$, $0 \leq t < c$ and $\alpha^c$ are well defined probability vectors with*

$$\alpha_k(t) = \pi_k(m(t)), \quad 0 \leq t < c, \quad and \quad \alpha_k^c = \frac{1}{c}\int_0^c \pi_k(m(t))\, dt, \tag{31}$$

*for $m(t)$ in (28), where $\pi_k(m)$ be the Poisson distribution with mean $m$, i.e.,*

$$\pi_k(m) \equiv \frac{e^{-m}m^k}{k!}, \quad k \geq 0. \tag{32}$$

We now consider the estimated death rates with ample data, i.e., $\bar{\mu}_k^e(\infty)$. To obtain the departure rate conditional on the number of busy servers, we use use the following consequence of Theorem 2.1 of [17].

18

**Theorem 3.5** (*remaining service times conditional on the number*) *Consider the periodic $M_t/GI/\infty$ queueing model starting empty in the distant past, where the service-time cdf $G$ has pdf $g$. Conditional on $Q(t) = k$, the remaining service times at time $t$ are distributed as $k$ i.i.d. random variables with pdf*

$$g_{k,t}(x) = \frac{\int_0^\infty \lambda(t-u)g(x+u)\,du}{\int_0^\infty \lambda(t-u)G^c(u)\,du}, \quad x \geq 0, \tag{33}$$

*which is independent of $k$. Hence, conditional on $Q(t) = k$, the departure rate at time $t$ is*

$$\delta_k(t) = k\delta_1(t) = \frac{k\mu g_{k,t}(0)}{E[\lambda(t-S_e)]} = \frac{k\mu E[\lambda(t-S)]}{E[\lambda(t-S_e)]} = \frac{k\delta(t)}{m(t)}. \tag{34}$$

From Theorem 3.5, we can recover the result that $\delta_k(t) = k\mu$ for the $M_t/M/\infty$ model, because $S_e$ is distributed the same as $S$ if and only if $S$ is exponential. That in turn implies that $\mu_k^e(\infty) = k\mu$ as well, as implied by Theorem 3.3. We now apply Theorem 3.5 to deduce a rate conservation property for this $M_t/GI/\infty$ model in each state over a periodic cycle. We also deduce alternative expressions for the estimated death rates.

**Theorem 3.6** (*arrival and departure rates over a cycle*) *For the periodic $M_t/GI/\infty$ queueing model starting empty in the distant past,*

$$\int_0^c \alpha_k(t)\lambda(t)\,dt = \int_0^c \alpha_k(t)\delta(t)\,dt \quad \text{for each} \quad k \geq 0 \tag{35}$$

*for $\alpha_k(t)$ in (31), so that*

$$\int_0^c \lambda(t)\,dt = \int_0^c \delta(t)\,dt. \tag{36}$$

*In addition, for each $k \geq 0$,*

$$\bar{\mu}_{k+1}^e(\infty) = \frac{\int_0^c \alpha_{k+1}(t)\delta_{k+1}(t)\,dt}{\int_0^c \alpha_{k+1}(t)\,dt} = \frac{\int_0^c \alpha_k(t)\delta(t)\,dt}{c\alpha_{k+1}^c} = \frac{\int_0^c \alpha_k(t)\lambda(t)\,dt}{c\alpha_{k+1}^c} = \frac{\bar{\lambda}_k^e(\infty)\alpha_k^c}{\alpha_{k+1}^c} \tag{37}$$

*for $\delta_k(t)$ in (34), $\alpha_k(t)$ and $\alpha_k^c$ in (31) and $\delta(t)$ in (30).*

**Proof.** Since $\lambda_k(t) = \lambda(t)$, independent of $k$, we can apply first (2) and then (34) to obtain

$$\int_0^c \alpha_k(t)\lambda(t)\,dt = c\alpha_k^c\bar{\lambda}_k^e(\infty) = c\alpha_{k+1}^c\bar{\mu}_{k+1}^e(\infty) = \int_0^c \alpha_{k+1}(t)\delta_{k+1}(t)\,dt$$

$$= \int_0^c \alpha_{k+1}(t)(k+1)[\delta(t)/m(t)]\,dt = \int_0^c \alpha_k(t)\delta(t)\,dt, \tag{38}$$

as in (35). We add over $k$ to get (36). The first expression in (37) is the direct rate expression for $\bar{\mu}_k^e(\infty)$. Then we apply (34), (35) and (2). ∎

## 4 The IS Model with a Sinusoidal Arrival-Rate Function

We now consider the special case of the periodic $M_t/GI/\infty$ model with a sinusoidal arrival rate function, as in [11, 25]; i.e., now we consider arrival rate functions of the form (3). By Theorem 4.1 of [11], the mean function is as in (6). We first exploit bounds on the mean $m(t)$ in (6) from §4 of [11] and (37) to obtain upper and lower bounds on the ratio $(k+1)\bar{\lambda}_k^e(\infty)/\bar{\mu}_{k+1}^e(\infty)$.

**Theorem 4.1** (*bounds on the ratio*) *For the $M_t/GI/\infty$ model with sinusoidal arrival rate function in* (3),

$$\frac{(k+1)\bar{\lambda}^e_k(\infty)}{\bar{\mu}^e_{k+1}(\infty)} = \frac{(k+1)\alpha^c_{k+1}}{\alpha^c_k}, \tag{39}$$

*where*

$$\left|\frac{(k+1)\alpha^c_{k+1}}{\alpha^c_k} - \frac{\bar{\lambda}}{\mu}\right| \le \left(\frac{\beta\bar{\lambda}}{\mu}\right)\left(E[\cos(\gamma S_e)]^2 + E[\sin(\gamma S_e)]^2\right)^{1/2} \le \frac{\beta\bar{\lambda}}{\mu}. \tag{40}$$

**Proof.** We obtain (39) directly from (37). We bound the term $(k+1)m(t)$ above and below by exploiting (12) of [11]. After removing this term from the $\alpha^c_{k+1}$, that term coincides with $\alpha^c_k$. ∎

We now establish asymptotic results for the extreme cases in which the cycles are very long ($\gamma \downarrow 0$) or are very short ($\gamma \uparrow \infty$). We directly show the dependence on $\gamma$; e.g., by writing $\bar{\lambda}_k(\infty; \gamma)$.

**Theorem 4.2** (*short cycles*) *For the $M_t/GI/\infty$ model with sinusoidal arrival rate function in* (3),

$$\bar{\lambda}_k(\infty; \gamma) \to \bar{\lambda} \quad and \quad \bar{\mu}_{k+1}(\infty; \gamma) \to (k+1)\mu \quad as \quad \gamma \uparrow \infty \quad for\ all \quad k \ge 0. \tag{41}$$

**Proof.** First, it is helpful to rewrite (20) so that the integrals are over a fixed interval, independent of $\gamma$. By making a change of variables $s = \gamma t$, we obtain

$$\bar{\lambda}^e_k(\infty) = \frac{\int_0^{2\pi/\gamma} \alpha_k(t)\lambda(t)\,dt}{\int_0^{2\pi/\gamma} \alpha_k(t)\,dt} = \frac{\int_0^{2\pi} \alpha_k(s/\gamma)\lambda(s/\gamma)\,ds}{\int_0^{2\pi} \alpha_k(s/\gamma)\,ds} \tag{42}$$

First, $\lambda(t; \gamma) \to \bar{\lambda}$ as $\gamma \uparrow \infty$, uniformly in $t$. By Theorem 4.5 of [11], $m(t; \gamma) \to \bar{\lambda}/\mu$ as $\gamma \uparrow \infty$, uniformly in $t$. Hence, $\alpha_k(t; \gamma) \to \alpha_k(t; \infty)$ as $\gamma \uparrow \infty$, uniformly in $t$, where $\alpha_k(t; \infty)$ is the Poisson pmf with mean $\bar{\lambda}/\mu$, independent of $t$. The bounded convergence theorem then implies the convergence of the integrals in (18). We then can apply (37) to deduce that

$$\bar{\mu}^e_{k+1}(\infty; \gamma) = \frac{\bar{\lambda}_k(\infty; \gamma)\alpha^c_{k;\gamma}}{\alpha^c_{k+1;\gamma}} \to \frac{\bar{\lambda}\alpha^c_{k;\infty}}{\alpha^c_{k+1;\infty}} = (k+1)\mu \quad as \quad \gamma \uparrow \infty, \tag{43}$$

because $\alpha_k(t; \infty)$ is the Poisson pmf with mean independent of $t$. ∎

**Theorem 4.3** (*long cycles*) *For the $M_t/GI/\infty$ model with sinusoidal arrival rate function in* (3),

$$\bar{\lambda}_k(\infty; \gamma) \to \frac{(k+1)\mu\alpha^c_{k+1;0}}{\alpha^c_{k;0}} \quad and \quad \bar{\mu}_{k+1}(\infty; \gamma) \to (k+1)\mu \quad as \quad \gamma \downarrow 0 \tag{44}$$

*for all $k \ge 0$, where $\alpha^c_{k;0}$ is the time average of $\alpha^c_k(t; 0)$ which is the Poisson pmf with mean $\bar{\lambda}\lambda_1(t)/\mu$, where $\lambda_1(t) = 1 + \beta\sin(t)$, $0 \le t \le 2\pi$.*

**Proof.** By Theorem 4.4 of [11], $m(t/\gamma) \to \lambda(t)/\mu$ as $\gamma \downarrow 0$ uniformly in $t$. Hence, $\alpha_k(t; \gamma) \to \alpha_k(t; 0)$ uniformly in $t$. We then apply this starting from (42), getting

$$
\begin{aligned}
\bar{\lambda}^e_k(\infty) &= \frac{\int_0^{2\pi/\gamma} \alpha_k(t)\lambda(t)\,dt}{\int_0^{2\pi/\gamma} \alpha_k(t)\,dt} = \frac{\int_0^{2\pi} \alpha_k(s/\gamma)\lambda(s/\gamma)\,ds}{\int_0^{2\pi} \alpha_k(s\gamma)\,ds} \\
&\to \frac{\int_0^{2\pi} \alpha_k(s; 0)\lambda(s; 0)\,ds}{\int_0^{2\pi} \alpha_k(s; 0)\,ds} = \frac{\int_0^{2\pi}(k+1)\mu\alpha_{k+1}(s; 0)\,ds}{\int_0^{2\pi} \alpha_k(s; 0)\,ds} = \frac{(k+1)\mu\alpha^c_{k+1}}{\alpha^c_k},
\end{aligned} \tag{45}
$$

20

because $\alpha_k(s; 0)$ is the Poisson pmf with mean $\lambda(s; 0)/\mu$ at time $s$. Finally,

$$\bar{\mu}_{k+1}(\infty; \gamma) = \frac{\bar{\lambda}_k(\infty; \gamma)\alpha_{k;\gamma}^c}{\alpha_{k+1;\gamma}^c} \rightarrow \frac{\bar{\lambda}_k(\infty; 0)\alpha_{k;0}^c}{\alpha_{k+1;0}^c} = (k+1)\mu \quad as \quad \gamma \downarrow 0. \quad \blacksquare$$

## 4.1 The $M_t/M/\infty$ Model with Sinusoidal Arrival Rate

For the $M_t/M/\infty$ model with sinusoidal arrival rate function in (3), the mean has an especially tractable from. From (15) of [11], the number in system, $Q(t)$, has a Poisson distribution for each $t$ with mean in (6).

In addition to the regularity in the estimated death rates exposed in Theorem 3.3, we have the following result for the estimated birth rates, which includes an explicit expression, upper and lower bounds, and an asymptotic result for short cycles (large $\gamma$).

**Theorem 4.4** (*estimated rates for the $M_t/M/\infty$ model with sinusoidal arrival rate function*) *In the $M_t/M/\infty$ IS queueing model with periodic arrival rate function, starting empty in the distant past,*

$$\bar{\mu}_{k+1}(\infty) = (k+1)\mu \quad and \quad \bar{\lambda}_k(\infty) = \frac{\alpha_{k+1}^c(k+1)\mu}{\alpha_k^c} = \frac{\mu \int_0^c e^{-m(t)}m(t)^{k+1}\,dt}{\int_0^c e^{-m(t)}m(t)^k\,dt} \tag{46}$$

*for $k \geq 0$, so that*

$$\bar{\lambda}\left(1 - \frac{\beta}{\sqrt{1+\gamma^2}}\right) \leq \bar{\lambda}_k(\infty) \leq \bar{\lambda}\left(1 + \frac{\beta}{\sqrt{1+\gamma^2}}\right) \quad for\ all \quad k \geq 0. \tag{47}$$

*and*

$$\bar{\lambda}_k(\infty) \rightarrow \bar{\lambda} \quad as \quad \gamma \rightarrow \infty \quad for\ all \quad k \geq 0. \tag{48}$$

**Proof.** The death rate expression and the first birth rate expression in (46) are immediate consequences of Theorem 3.3 and (2). The bounds then follow from Theorem 4.1, using the explicit expression from (18) of [11]. $\blacksquare$

## 4.2 Heavy-Traffic Limits for the Fitted Birth Rates

We conclude by deriving a heavy-traffic limit for the constant vales at large and small arguments.

**Theorem 4.5** (*heavy-traffic limits*) *In the $M_t/M/\infty$ IS queueing model with periodic arrival rate function, starting empty in the distant past,*

$$\frac{\bar{\lambda}_k(\infty)}{\bar{\lambda}} \rightarrow 1 - \frac{\beta}{\sqrt{1+\gamma^2}} \quad as \quad \bar{\lambda} \rightarrow \infty \quad and$$

$$\frac{\bar{\lambda}_{\lfloor m\bar{\lambda}\rfloor + k}(\infty)}{\bar{\lambda}} \rightarrow 1 + \frac{\beta}{\sqrt{1+\gamma^2}} \quad as \quad \bar{\lambda} \rightarrow \infty \quad for \quad m > 1/\log_e 2 \approx 1.44. \tag{49}$$

**Proof.** In each case, we apply Laplace's method to the numerator and denominator of (46), after pre-multiplying both by the same appropriate term (so this term cancels). Let $x \equiv \bar{\lambda}/\mu$ and consider the first expression. In particular, After multiplying the numerator and denominator by $e^x/x^k$, we can express the denominator as

$$\int_0^c e^{-xs(t)}(1+s(t))^k \, dt \sim \sqrt{\frac{2\pi}{x|s''(x_0)|}}(1+s(x_0))^k e^{xs(x_0)} \quad \text{as} \quad x \to \infty, \quad (50)$$

where $\sim$ means that the ratio of the two sides converges to 1, $s(t) \equiv m_1(t)-1$ for $m_1(t)$ in (6), where $c = 2\pi/\gamma$ and $x_0 = c - cot^{-1}(1/\gamma))/\gamma$ and $m(x_0) = (\bar{\lambda}/\mu)(1 - \beta/(\sqrt{1+\gamma^2}))$, by virtue of (16) and (18) in [11]. (The minus sign in the exponent of $e^{-xs(t)}$ means that we look for the most negative value of $s(t)$.) We have used the fact that the integral is dominated by an appropriate modification of the integrand at a single point when $x$ becomes large. The ratio in (46) thus approaches $1+s(x_0)$.

For the second expression, after multiplying the numerator and denominator by $e^x/x^{x+k}$, we can express the denominator as

$$\int_0^c e^{-xs(t)}(1+s(t))^{mx+k} \, dt = \int_0^c e^{+x[m\log_e\{1+s(t)\}-s(t)]}(1+s(t))^k \, dt$$

$$\sim \sqrt{\frac{2\pi}{x|f''(x_0)|}}(1+s(x_0))^k e^{xf(x_0)} \quad \text{as} \quad x \to \infty, \quad (51)$$

where $f(t) \equiv m\log_e\{1+s(t)\} - s(t)$, so that $x_0 = (c/4) + cot^{-1}(1/\gamma))/\gamma$ and $m(x_0) = (\bar{\lambda}/\mu)(1 + \beta/(\sqrt{1+\gamma^2}))$, again by (16) and (18) in [11]. (The plus sign in the exponent of $e^{+x[m\log_e\{1+s(t)\}-s(t)]}$ with $m > 1/\log_2 2$ means that we look for the most positive value of $s(t)$.) The ratio in (46) again approaches $1 + s(x_0)$. ∎

## 5 Conclusions

We have conducted extensive simulation experiments to study the potential of fitting general state-dependent birth-and-death (BD) processes to queueing system data. As indicated in §1, this can be an effective way to simultaneously fit and test the classical Erlang A model. Here we have studied the consequence of fitting a BD process to other models. In particular, here we focused on the BD processes fit to multi-server queues with NHPP arrival processes having periodic arrival rate functions. We fit BD processes to the sample path of the number in system in an $M_t/M/\infty$ model and related $M_t/GI/s$ models with $s = 40$ and non-exponential service times. These models have the sinusoidal arrival rate function in (3) with relative amplitude $\beta = 10/35$. In the experiments we considered arrival rates $\bar{\lambda} = 35$ and 100 (moderately large scale) for a range of scaling factors $\gamma$, yielding a range of sine cycles of length $2\pi/\gamma$.

From these experiments, we see that the death rates have the same linear structure as for the many-server $GI/GI/s$ models studied in [10], but we see significantly different fitted birth rates, as can be seen by comparing Figures 1 and 2. Theorems 4.4 and 4.5 establish finite bounds and heavy-traffic limits for the fitted birth rates, consistent with these figures. The simulation results in §§2.3-2.8 indicate that (i) for larger $\gamma$ (shorter cycles) such as $\gamma \geq 1$, the fitted BD process may serve as a useful direct approximation for the original queue-length process, but (ii) for smaller $\gamma$ (longer cycles) such as $\gamma \leq 0.1$, the transient behavior of the fitted BD process is very different. However, consistent with the theory in [33], we see that the fitted BD process consistently describes the steady-state distribution. In §2.8 we showed that a relatively simple two-parameter parametric function can be fit to the estimated birth rates in order to efficiently estimate first the fitted birth

rate and then the steady-state distribution of the original system. The results here for known stochastic models should help interpret similar fitting to data from complicated service systems, as in work in progress [36].

**Acknowledgement**

# References

[1] Abate, J. and Whitt, W. (1999). Computing laplace transforms for numerical inversion via continued fractions. *INFORMS Journal on Computing* 11(4):394–405.

[2] Brown, L., Gans, N., Mandelbaum, A., Sakov, A., Shen, H., Zeltyn, S. and Zhao, L. (2005). Statistical analysis of a telephone call center: a queueing-science perspective. *J Amer Stat Assoc* 100:36–50.

[3] Browne, S. and Whitt, W. (1995). Piecewise-linear diffusion processes. In Dshalalow, J. (ed.), *Advances in Queueing.* Boca Raton, FL: CRC Press, pp. 463–480.

[4] Buzen, J. (1976). Fundamental operational laws of computer system performance. *Acta Informatika* 14:167–182.

[5] Buzen, J. (1978). Operational analysis: an alternative to stochastic modeling. In Ferarri, D. (ed.), *Performance of Computer Installations.* Amsterdam: North Holland, pp. 175–194.

[6] Chang, C. S., Chao, X. L. and Pinedo, M. (1991). Monotonicity results for queues with doubly stochastic Poisson arrivals: Ross's conjecture. *Advances in Applied Probability* 12(41):210–228.

[7] Choudhury, G. L., Mandelbaum, A., Reiman, M. I. and Whitt, W. (1997). Fluid and diffusion limits for queues in slowly changing random environments. *Stochastic Models* 13(1):121–146.

[8] Crescenzo, A. D. and Nobile, A. G. (1995). Diffusion approximation to a queueing system with time-dependent arrival and service rates. *Queueing Systems* 19:41–62.

[9] Denning, P. J. and Buzen, P. J. (1978). The operational analysis of queueing network models. *Computing Surveys* 10:225–261.

[10] Dong, J. and Whitt, W. (2014). Stochastic grey-box modeling of queueing systems: fitting birth-and-death processes to data. Queueing Systems, published on line on December 2, 2014.

[11] Eick, S. G., Massey, W. A. and Whitt, W. (1993). $M_t/G/\infty$ queues with sinusoidal arrival rates. *Management Sci* 39:241–252.

[12] Eick, S. G., Massey, W. A. and Whitt, W. (1993). The physics of the $M_t/G/\infty$ queue. *Oper Res* 41:731–742.

[13] El-Taha, M. and Stidham, S. (1999). *Sample-Path Analysis of Queueing Systems.* Boston: Kluwer.

[14] Falin, G. I. (1989). Periodic queues in heavy traffic. *Advances in Applied Probability* 21:485–487.

[15] Gans, N., Liu, N., Mandelbaum, A., Shen, H. and Ye, H. (2010). Service times in call centers: Agent heterogeneity and learning with some operational consequences. *IMS Collections, Borrowing Strength: Theory Powering Applications  A Festschrift for Lawrence D Brown* 6:99–123.

[16] Garnett, O., Mandelbaum, A. and Reiman, M. I. (2002). Designing a call center with impatient customers. *Manufacturing and Service Oper Management* 4(3):208–227.

[17] Goldberg, D. and Whitt, W. (2008). The last departure time from an $M_t/G/\infty$ queue with a terminating arrival process. *Queueing Systems* 58:77–104.

[18] Green, L. V. and Kolesar, P. J. (1991). The pointwise stationary approximation for queues with nonstationary arrivals. *Management Sci* 37:84–97.

[19] Halfin, S. and Whitt, W. (1981). Heavy-traffic limits for queues with many exponential servers. *Operations Research* 29(3):567–588.

[20] Heyman, D. P. and Whitt, W. (1984). The asymptoic behavior of queues with time-varying arrival. *Journal of Applied Probability* 21(1):143–156.

[21] Ibrahim, R., L'Ecuyer, P., Regnard, N. and Shen, H. (2012). On the modeling and forecasting of call center arrivals. *Proceedings of the 2012 Winter Simulation Conference* 2012:256–267.

[22] Kim, S. and Whitt, W. (2014). Are call center and hospital arrivals well modeled by nonhomogeneous Poisson processes? *Manufacturing and Service Oper Management* 16(3):464–480.

[23] Kim, S. and Whitt, W. (2014). Choosing arrival process models for service systems: Tests of a nonhomogeneous Poisson process. *Naval Research Logistics* 17:307–318.

[24] Mandelbaum, A., Massey, W. A. and Reiman, M. I. (1998). Strong approximations for Markovian service networks. *Queueing Systems* 30:149–201.

[25] Massey, W. A. and Whitt, W. (1993). Networks of infinite-server queues with nonstationary Poisson input. *Queueing Systems* 13(1):183–250.

[26] Massey, W. A. and Whitt, W. (1996). Stationary-process approximations for the nonstationary Erlang loss model. *Oper Res* 44(6):976–983.

[27] Puhalskii, A. A. (2013). On the $M_t/M_t/K_t + M_t$ queue in heavy traffic. *Math Methods Oper Res* 78:119–148.

[28] Rolski, T. (1989). Queues with nonstationary inputs. *Queueing Systems* 5:113–130.

[29] Rothkopf, M. H. and Oren, S. S. (1979). A closure approximation for the nonstationary $M/M/s$ queue. *Management Science* 25(6):522–534.

[30] Whitt, W. (1984). Departures from a queue with many busy servers. *Mathematics of Operations Research* 9(4):534–544.

[31] Whitt, W. (1991). The pointwise stationary approximation for $M_t/M_t/s$ queues is asymptotically correct as the rates increase. *Management Science* 37(3):307–314.

[32] Whitt, W. (2004). Efficiency-driven heavy-traffic approximations for many-server queues with abandonments. *Management Sci* 50(10):1449–1461.

[33] Whitt, W. (2012). Fitting birth-and-death queueing models to data. *Statistics and Probability Letters* 82:998–1004.

[34] Whitt, W. (2014). Heavy-traffic limits for queues with periodic arrival processes. *Operations Research Letters* 42:458–461.

[35] Whitt, W. (2014). The steady-state distribution of the $M_t/M/\infty$ queue with a sinusoidal arrival rate function. *Operations Research Letters* 42:311–318.

[36] Whitt, W. and Zhang, X. (2015). Stochastic grey-box modeling of hospital occupancy levels. In preparation, Columbia University, http://www.columbia.edu/~ww2040/allpapers.html.

[37] Wolff, R. W. (1965). Problems for statistical inference for birth and death queueing models. *Operations Research* 13:343–357.