

Modeling, Simulating, and Measuring the Performance of Queues

Anthony Hung

2019-03-01

Prerequisites

This vignette continues from concepts covered in the vignette: “Introduction to Queueing Theory and Queueing Models”.

Introduction

In addition to simply being able to represent waiting lines mathematically, queueing theory allows for the evaluation of the behavior and performance of queues. Being able to measure the performance of queues also allows us to determine the effects of altering components of the queue on performance.

Simulating a M/M/1 queue

In simulating a M/M/1 queue, we want to keep track of three values of the queue over time.

1. Arrival times of customers
2. Departure times of customers
3. The number of customers in the system at every moment of arrivals or departures

In simulating the queue behavior, we can take advantage of the superposition property of combined independent Poisson processes. Since arrivals and departures are independent, the number of events in the combined process can be represented as a Poisson process with parameter $\lambda_{sum} = \lambda + \mu$. The probability of an event in this combined process being an arrival is $\frac{\lambda}{(\lambda + \mu)}$, and the probability of it being a departure is $\frac{\mu}{(\lambda + \mu)}$.

The function “simulate_MM1” simulates the number of customers in a M/M/1 queue over time given values for lambda, mu, and N_0 from T_0 to T_{max} . It also keeps track of when events (arrivals or departures) occur during the time periods and what type of event occurs at each of those moments.

```
lambda <- 4
mu <- 5

simulate_MM1 <- function(lambda=lambda, mu=mu, NO=0, Tmax=1000){
  #Initialize vectors to store each of the values of interest throughout the simulation
  events <- 0 #stores the type of event (1 for arrival, -1 for departure)
  Times <- 0 #times of events
  customers <- NO #number of customers at each time in Times

  while(tail(Times,1) < Tmax){ #keep simulating until you have an event at a time greater than Tmax

    if(tail(customers,1)==0){ #separate behavior occurs if system currently has 0 customers
      tau <- rexp(1, rate=lambda) #interarrival intervals are exponentially distributed
      event <- 1 #only an arrival can occur if there are 0 customers
```

```

} else {
  tau <- rexp(1, rate=lambda+mu) #inter-event intervals are exponentially distributed
  if(runif(1,0,1) < lambda/(lambda+mu)){ #if runif is less than P(event = arrival)...
    event <- 1 #call the event an arrival
  } else{
    event <- -1 #otherwise, call the event a departure
  }
}

#now that we have simulatd one event, we need to do some accounting
customers <- c(customers, tail(customers,1)+event)
Times <- c(Times, tail(Times,1)+tau)
events <- c(events, event)
}

#we need to toss out the information from the last event (it occurred after Tmax)
events <- head(events, -1)
Times <- head(Times, -1)
customers <- head(customers, -1)

return(list(events,Times,customers))
}

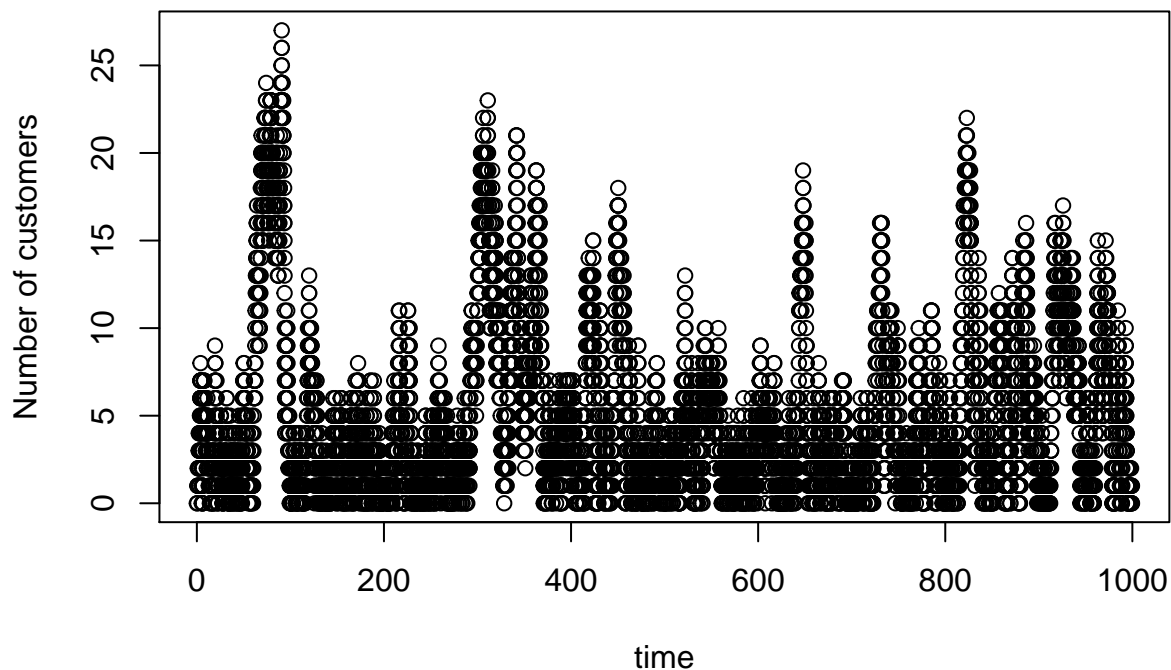
```

After simulating the number of customers in the queue over one run of the simulation, we can plot it.

```

sim <- simulate_MM1(lambda = 4, mu=5, NO=0, Tmax=1000)
plot(x=sim[[2]], y=sim[[3]], xlab="time", ylab="Number of customers") #number of customers vs time

```

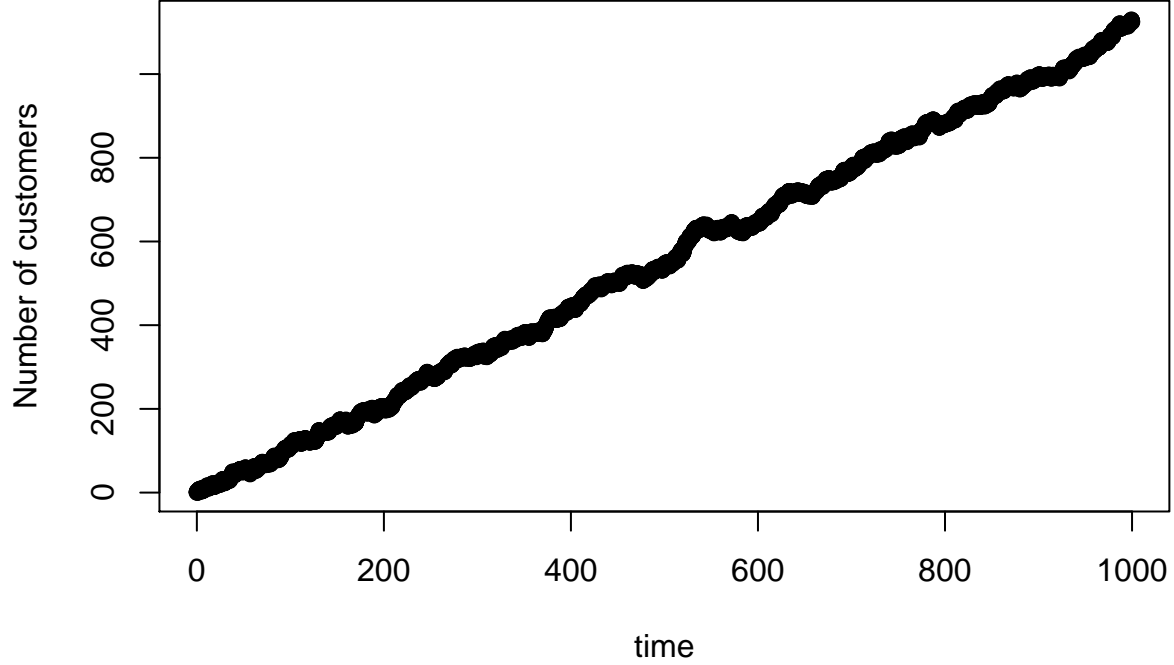


Notice that if $\lambda > \mu$, the customer number explodes and will never reach a steady state.

```

sim2 <- simulate_MM1(lambda = 6, mu=5, NO=0, Tmax=1000)
plot(x=sim2[[2]], y=sim2[[3]], xlab="time", ylab="Number of customers") #number of customers vs time

```



Measuring the performance of queues

There are several formal quantities used to measure the performance of a queueing system (with c servers).

1. π_j := The stationary probability that there are j customers in the system
2. a := Offered load. The mean number of requests per service time.
3. ρ := Server utilization or traffic intensity. Offered load per server (a/c).
4. L := Mean number of customers in the system, including those in the buffer and at servers.
5. L_s := Mean number of customers being served.
6. L_q := Mean number of customers waiting in the buffer.
7. W := Mean length of time between a customer's arrival and the customer's departure from the system.
8. W_s := Mean length of time a customer spends at the server.
9. W_q := Mean length of time between a customer's arrival and when the customer's service starts.

Here are the equations for the performance metrics for a M/M/1 queue

π_j

For the M/M/1 queue, we previously calculated the stationary probabilities:

$$\pi_j = (1 - \frac{\lambda}{\mu})(\frac{\lambda}{\mu})^j$$

a, ρ

The offered load of a queue is given by the ratio of the arrival rate to the departure rate (i.e. the mean number of arrivals that occur during the mean service time). Server utilization or traffic intensity is the offered load

per server. One can think about the server utilization as, on average, what proportion of the time each server in the system is occupied. For a single server queue, the offered load equals the server utilization.

$$a = \frac{\lambda}{\mu}$$

$$\rho = \frac{a}{c} = \frac{\frac{\lambda}{\mu}}{1} = \frac{\lambda}{\mu}$$

With our solution for ρ above, the stationary probabilities of the M/M/1 queue are commonly represented as:

$$\pi_j = (1 - \frac{\lambda}{\mu})(\frac{\lambda}{\mu})^j = (1 - \rho)\rho^j$$

A brief interlude: Little's Law

Little's law states that the long-term average of the number of customers in any queue at stationarity is equal to the long-term average arrival rate λ multiplied by the average time that a customer spends in the system. Expressed algebraically using our defined variables:

$$L = \lambda W$$

Little's law also holds for the number of customers being served and for the number of customers waiting in the queue buffer:

$$L_s = \lambda W_s$$

$$L_q = \lambda W_q$$

Little's Law was originally presented by John Little in 1954 without proof, but multiple proofs of the relationship have been published since (<https://pubsonline.informs.org/doi/abs/10.1287/opre.20.6.1115>). We can make use of the relationships stated in Little's Law in working with the last four performance measures.

L, W

The mean number of customers in the system at stationarity can be computed knowing the stationary probabilities of having j customers in the system:

$$L = \sum_{j=0}^{\infty} j\pi_j = \sum_{j=0}^{\infty} j(1 - \rho)\rho^j$$

When $\rho < 1$, this expression represents the expectation of a geometric distribution with parameter $p = 1 - \rho$. Therefore,

$$L = \sum_{j=0}^{\infty} j\pi_j = \sum_{j=0}^{\infty} j(1 - \rho)\rho^j = \frac{\rho}{1 - \rho}$$

Using Little's Law, we can find W :

$$W = \frac{L}{\lambda} = \frac{\rho}{(1 - \rho)\lambda}$$

L_s, W_s

The mean length of time a customer spends at the server is simply the mean service duration:

$$W_s = \frac{1}{\mu}$$

Using Little's Law, we can find L_s :

$$L_s = \lambda W_s = \lambda \frac{1}{\mu} = \frac{\lambda}{\mu} = \rho$$

L_q, W_q

As the number of individuals in a queue system is equal to the sum of the number of individuals in the buffer and the number of individuals currently being served:

$$L_q = L - L_s = \frac{\rho}{1 - \rho} - \rho = \frac{\rho - \rho(1 - \rho)}{1 - \rho} = \frac{\rho^2}{1 - \rho}$$

Using Little's Law, we can find W_q :

$$W_q = \frac{L_q}{\lambda} = \frac{\frac{\rho^2}{1 - \rho}}{\lambda} = \frac{\rho^2}{(1 - \rho)\lambda} = \frac{\rho \frac{\lambda}{\mu}}{(1 - \rho)\lambda} = \frac{\rho}{(1 - \rho)\mu}$$

Measuring performance from simulated queue data

In addition to knowing the equations to compute these performance metrics for the M/M/1 queue, we can also compute them for the data we simulated above.

Multiple servers: the M/M/c queue