

# Queueing Theory: Simulation of Queueing processes in continuous time

EBB074A05

Nicky D. van Foreest

2022:01:19

## 1 General info

This file contains the code and the results that go with the YouTube movies mentioned in the README file in this directory.

You should read the relevant section of my queueing book to understand what it going on. Most of it is very easy, but without background a bit cryptic (I believe).

I included a number of exercises to help you think about the code. Keep your answers short; you don't have to win the Nobel prize on literature.

## 2 Computing waiting times

Here we just follow the steps of the queueing book to construct a single server FIFO queue in continuous time.

### 2.1 Load standard modules

We need the standard libraries for numerical work and plotting.

Python Code

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from matplotlib import style
4
5 style.use('ggplot')
6
7 np.random.seed(3)
```

### 2.2 Inter-arrival times

Simulate random interarrival times that are  $\sim \text{Exp}(\lambda)$ , with  $\lambda = 3$ . First I take just three jobs, so that I can print out all intermediate results and check how things work. Once I am convinced about the correctness, I run a simulation for many jobs.

Python Code

```
1 num = 3
2 labda = 3
3 X = np.random.exponential(scale=labda, size=num)
4 print(X)
```

Here is an important check (I always forget the meaning of  $\lambda$  when I provide it to the simulator)

Python Code

```
1 num = 100
2 labda = 3
3 X = np.random.exponential(scale=labda, size=num)
4 print(X.mean())
```

**Ex 2.1.** Explain that `scale=labda` sets the interarrival times to 3, but that in our queueing models,  $\lambda$  should correspond to the arrival rate. Why is the code below in line with what we want?

Python Code

```
1 num = 3
2 labda = 3
3 X = np.random.exponential(scale=1/labda, size=num)
```

## 2.3 Arrival times

Python Code

```
1 A = X.cumsum()
2 print(A)
```

**Ex 2.2.** Why do we generate first random inter-arrival times, and use these to compute the arrival times? Why not directly generate random arrival times?

Check the numbers to see that the arrival time of job 0 is  $A_0 > 0$ . But I want time to start at time  $A_0 = 0$ . Here is the trick to achieve that.

Python Code

```
1 A = np.zeros(len(X) + 1)
2 A[1:] = X.cumsum()
3 print(A)
```

This is better!

**Ex 2.3.** Why is the vector A one longer than X?

## 2.4 Service times

We have arrival times. We next need the service times of the jobs. Assume they are  $\sim \text{Exp}(\mu)$  with  $\mu$  somewhat larger than  $\lambda$ . (Recall this means that jobs can be served faster than that they arrive.)

Python Code

```
1 mu = 1.2 * labda
2 S = np.random.exponential(scale=1/mu, size=len(A))
3 S[0] = 0
4 print(S)
```

Note, `S[0]` remains unused; it should correspond to job 0, but we neglect this job 0 in the remainder.

**Ex 2.4.** Why do I use `size=len(A)` in the definition of S?

Ex 2.5. Why do we set  $\text{scale}=1/\mu$ ?

Ex 2.6. It's easy to compute the mean service time like this

```
_____ Python Code _____  
1 S.mean()
```

Explain that we need to set  $S[0]=0$  to get the correct result. If  $\text{num}$  is big number, does it matter that we set  $S[0]=0$ ?

## 2.5 Departure times

The standard recursion to compute the departure times.

```
_____ Python Code _____  
1 D = np.zeros_like(A)  
2  
3 for k in range(1, len(A)):  
4     D[k] = max(D[k - 1], A[k]) + S[k]  
5  
6 print(D)
```

Ex 2.7. Explain now why it is practical to have  $A_0 = 0$ .

## 2.6 Sojourn times

How long do you stay in the system if you arrive at some time  $A_n$  and you depart at  $D_n$ ?

```
_____ Python Code _____  
1 J = D - A  
2 print(J)
```

## 2.7 Waiting times

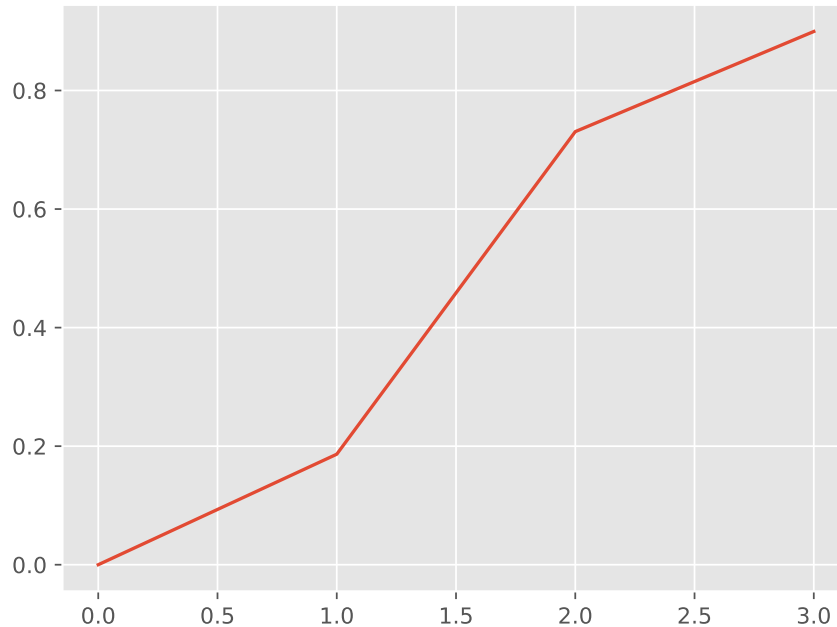
If your sojourn time is 10, say, and your service time at the server is 3 (and there is just one server and the service discipline is FIFO), then what was your time in queue?

```
_____ Python Code _____  
1 W = J - S
```

Ex 2.8. Recall that we set  $S[0] = 0$ . Suppose that we wouldn't have done this, and we would run the simulation for a small number of cases, why can  $W.\text{mean}()$  be negative?

## 2.8 KPIs and plots

```
_____ Python Code _____  
1 print(J.mean(), J.std())  
  
_____ Python Code _____  
1 plt.clf()  
2 plt.plot(J)  
3 plt.savefig("wait.pdf")
```



**Ex 2.9.** Change the simulation length to 1000 jobs. Do one run for  $\mu = 3.5$  and another for 2.8. Compute the KPIs, make a plot, and include that in your assignment. Comment on what you see.

## 2.9 Server KPI: idle time

This code computes the total time the server is idle, and then computes the fraction of time the server is idle.

Python Code

```
1 rho = S.sum() / D[-1]
2 idle = (D[-1] - S.sum()) / D[-1]
3 print(idle)
```

**Ex 2.10.** Explain the code above. Some specific points:

1. Why is `S.sum()` the total busy time of the server?
2. Why do we divide by `D[-1]` in the computation of  $\rho$ ?
3. Explain the computation of the `idle` variable.

The next code computes the separate idle times.

Python Code

```
1 idle_times = np.maximum(A[1:] - D[:-1], 0)
2 print(idle_times)
3 print(idle_times.sum())
4 print(D[-1] - S.sum())
```

**Ex 2.11.** Run this code for a simulation with 10 or so jobs (some other small number). Explain how this code works. Which line is a check on the computations?

## 2.10 Server KPI: busy time

We also like to know how long the server has to work uninterruptedly. Finding the busy times is quite a bit harder than the idle times. (A busy time starts when a job arrives at an empty system and it stops when the server becomes free again.)

**Ex 2.12.** To help you understand the code, let's first do a numerical example. Suppose jobs 1, 4, 8 find an empty system upon *arrival*. The simulation contains 10 jobs. Why do jobs 3, 7, 10 leave an empty system behind upon *departure*?

With this idea, we can compute the idle times in another way (as a check on earlier work), and then we extend the approach to the busy times.

---

Python Code

---

```
1 import numpy as np
2
3 np.set_printoptions(suppress=True)
4 np.random.seed(3)
5
6 num = 10
7 labda = 3
8 X = np.random.exponential(scale=1 / labda, size=num)
9 A = np.zeros(len(X) + 1)
10 A[1:] = X.cumsum()
11 mu = 1.2 * labda
12 S = np.random.exponential(scale=1 / mu, size=len(A))
13 S[0] = 0
14 D = np.zeros_like(A)
15
16 for k in range(1, len(A)):
17     D[k] = max(D[k - 1], A[k]) + S[k]
18
19
20 W = D - S - A # waiting times
21 idx = np.argwhere(np.isclose(W, 0))
22 idx = idx[1:] # strip A[0]
23 idle_times = np.maximum(A[idx] - D[idx - 1], 0)
24 print(idle_times.sum())
```

---

**Ex 2.13.** What is stored in `idx`? Why do we strip `A[0]`? Why do we subtract `D[idx-1]` and not `D[idx]`? (Print out the variables to understand what they mean, e.g., `print(idx)`.)

Now put the next piece of code behind the previous code so that we can compute the busy times.

---

Python Code

---

```
1 busy_times = D[idx - 1][1:] - A[idx][:-1]
2 last_busy = D[-1] - A[idx[-1]]
3 print(busy_times.sum() + last_busy, S.sum())
```

---

**Ex 2.14.** Explain these lines. About the last line, explain why this acts as a check.

## 3 Computing Queue length

We have the waiting times, but not the number of jobs in queue. What if we would like to plot the queue length process?

A simple, but inefficient, algorithm to construct the queue length process is to walk backwards in time.

---

```

1 import numpy as np
2 np.random.seed(3)
3
4 num = 10
5 X = np.random.exponential(scale=labda, size=num)
6 A = np.zeros(len(X) + 1)
7 A[1:] = X.cumsum()
8 mu = 0.8 * labda
9 S = np.random.exponential(scale=mu, size=len(A))
10 D = np.zeros_like(A)
11
12 for k in range(1, len(A)):
13     D[k] = max(D[k-1], A[k]) + S[k]
14
15 L = np.zeros_like(A)
16 for k in range(1, len(A)):
17     l = k - 1
18     while D[l] > A[k]:
19         l -= 1
20     L[k] = k - l
21
22 print(L)

```

---

**Ex 3.1.** Explain how this code works. At what points in time do we sample the queue length?

**Ex 3.2.** The above procedure to compute the number of jobs in the system is pretty inefficient. Why is that so?

**Ex 3.3.** Try to find a (more efficient) algorithm to compute  $L$ . If you cannot solve this yourself, explain my code that is provided in the hint.

## 4 Multi-server queue

Let us now generalize the simulation to a queue that is served by multiple servers.

### 4.1 A single fast server

While you are *still in queue* of a multi-server queue, the rate at which jobs are served is the same whether there are  $c$  servers or just one server working at a rate of  $c$ , i.e.,  $c$  times as fast as a server in the multi-server. As a first simple case, we model the multi-server queueing system as if there is one fast server working at rate  $c$ .

**Ex 4.1.** Explain that we can implement a fast server by specifying the number of servers  $c$ , and change the service times as follows:

---

```

1 c = 3
2 S = np.random.exponential(scale=1 / (c * mu), size=len(A))

```

---

### 4.2 A real multi-server queue

Here is the code to implement a real multi-server queue; see the queueing book to see how it works.

---

Python Code

---

```

1  import numpy as np
2
3  np.random.seed(3)
4
5  labda = 3
6  mu = 4
7  num = 3
8
9  X = np.random.exponential(scale=1 / labda, size=num + 1)
10 S = np.random.exponential(scale=1 / mu, size=num)
11
12 # single server queue
13 W = np.zeros_like(S)
14 for k in range(len(S)):
15     W[k] = max(W[k - 1] + S[k - 1] - X[k], 0)
16
17 print(W)
18 print(W.mean(), W.std())
19
20 # code for multi server queue
21 c = np.array([1.0])
22 W = np.zeros_like(S)
23 w = np.zeros_like(c)
24 for k in range(len(S)):
25     s = w.argmax() # server with smallest waiting time
26     W[k] = w[s]
27     w[s] += S[k] # assign arrival to this server
28     w = np.maximum(w - X[k + 1] * c, 0)
29
30 print(W)
31 print(W.mean(), W.std())

```

---

**Ex 4.2.** First a test, we set the vector of server capacities  $c=[1]$  so that we reduce our multi-server queue to a single-server queue. Run the code and check that both simulations give the same result.

BTW: such ‘dumb’ corner cases are necessary to test code. In fact, it has happened many times that I tested code of which I was convinced it was correct, but I still managed to make bugs. A bit of paranoia is a good state of mind when it comes to coding.

**Ex 4.3.** We can modify the code for the single server queue such that it represents a single fast server working at rate  $c$ , rather than at 1.

---

Python Code

---

```

1  W = np.zeros_like(S)
2  for k in range(len(S)):
3      W[k] = max(W[k - 1] + S[k - 1]/c - X[k], 0)

```

---

How does this work?

Now that we have tested the implementation (in part), here is the code for a queue served by three servers, all working at the same speed.

---

Python Code

---

```

1  import numpy as np
2
3  np.random.seed(3)
4

```

```

5  labda = 3
6  mu = 1.1
7  N = 1000
8
9  X = np.random.exponential(scale=1 / labda, size=N + 1)
10 S = np.random.exponential(scale=1 / mu, size=N)
11
12 c = np.array([1.0, 1.0, 1.0])
13 W = np.zeros_like(S)
14 w = np.zeros_like(c)
15 for k in range(len(S)):
16     s = w.argmin() # server with smallest waiting time
17     W[k] = w[s]
18     w[s] += S[k] # assign arrival to this server
19     w = np.maximum(w - X[k + 1] * c, 0)
20
21 print(W.mean(), W.std())

```

---

**Ex 4.4.** Compare  $E[W]$  for two cases. The first is a model with single fast server working at rate  $c = 3$ . The second is a model with three servers each working at rate 1. Include your numerical results and discuss the differences.

**Ex 4.5.** Change the code for the multi-server such that the individual servers have different speeds but such that the total service capacity remains the same. What is the impact on  $E[W]$  and  $V[W]$  as compared to the case with equal servers? Include your numerical results.

**Ex 4.6.** Once you researched the previous exercise, provide some consultancy advice. Is it better to have one fast server and several slow ones, or is it better to have 3 equal servers? What gives the least queueing times and variance? If the variance is affected by changing the server rates, explain the effects based on the intuition you can obtain from Sakasegawa's formula.

## 5 Hints

**h.2.3.** When we have  $n$  arriving jobs, how many *interarrival* times do we have?

**h.2.4.** If I would not do this, and I would want to change the simulation length (the number of jobs), at how many places should I change this number?

**h.2.6.** Did we really serve job 0?

**h.2.7.** Observe that for  $D_1$  we need  $D_0$ . If  $A_0$  would be the arrival time of the first job, then what would we take for  $D_{-1}$ ?

**h.2.8.** We subtract  $S[0]$  as if we served the corresponding job, but did we actually serve it?

**h.3.3.** Here is the code.

	Python Code	
--	-------------	--

```

1  import numpy as np
2  import matplotlib.pyplot as plt
3
4  np.random.seed(3)
5
6  num = 4
7  labda = 3
8  X = np.random.exponential(scale=1 / labda, size=num)
9  A = np.zeros(len(X) + 1)

```



```

10 A[1:] = X.cumsum()
11 mu = 0.3 * labda
12 S = np.random.exponential(scale=1 / mu, size=len(A))
13 S[0] = 0
14 D = np.zeros_like(A)
15
16 for k in range(1, len(A)):
17     D[k] = max(D[k - 1], A[k]) + S[k]
18
19 L = np.zeros((len(A) + len(D), 2))
20 L[: len(A), 0] = A
21 L[1 : len(A), 1] = 1
22 L[len(D) :, 0] = D
23 L[len(D) + 1 :, 1] = -1
24 N = np.argsort(L[:, 0], axis=0)
25 L = L[N]
26 L[:, 1] = L[:, 1].cumsum()
27 print(L)
28
29 plt.clf()
30 plt.step(L[:, 0], L[:, 1], where='post', color='k')
31 plt.plot(A[1:], np.full_like(A[1:], -0.3), '^b', markeredgewidth=1)
32 plt.plot(D[1:], np.full_like(D[1:], -0.3), 'vr', markeredgewidth=1)
33 plt.savefig("wait4.pdf")

```

---

**h.4.5.** For instance, set `c=np.array([2, 0.5, 0.5])`.