

# Queueing Theory: Simulation in discrete time

EBB074A05

Nicky D. van Foreest Nicky

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 it will be quite 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.

The idea of the code is that you run it yourself as you read through the tutorial.

## 2 Deterministic queues

Here we consider a simple queueing system in which all is deterministic.

### 2.1 one period, demand, service capacity, and queue

There is one server, jobs enter a queue in front of the server, and the server serves batches of customers, every hour say.

---

Python Code

---

```
1 L = 10
2 a = 5
3 d = 8
4 L = L + a - d
5 print(L)
```

---

---

Python Code

---

```
1 L = 3
2 a = 5
3 c = 7
4 d = min(c, L)
5 L += a - d
6 print(d, L)
```

---

### 2.2 two periods

---

Python Code

---

```
1 L = 3
2 a = 5
3 c = 7
```

```

4  d = min(c, L)
5  L += a - d
6
7  a = 6
8  d = min(c, L)
9  L += a - d
10 print(d, L)

```

---

**Ex 2.1.** Add a third period, and report your result in the assignment.

## 2.3 many periods, use vectors and for loops

	Python Code	
<pre> 1  num = 5 2 3  a = 9*np.ones(num) 4  c = 10*np.ones(num) 5  L = np.zeros_like(a) 6 7  L[0] = 20 8  for i in range(1, num): 9      d = min(c[i], L[i-1]) 10     L[i] = L[i-1] + a[i] - d 11 12 print(L) </pre>		

---

**Ex 2.2.** Run the code for 10 periods and report your result.

**Ex 2.3.** What will be the queue after 100 or so periods?

## 3 Stochastic/random behavior

Real queueing systems show lots of stochasticity: in the number of jobs that arrive in a period, and the time it takes to serve these jobs. In this section we extend the previous models so that we can deal with randomness.

### 3.1 simulate demand

	Python Code	
<pre> 1  import numpy as np 2 3  a = np.random.randint(5, 8, size=5) 4  print(a) </pre>		

---

### 3.2 Set seed

	Python Code	
<pre> 1  import numpy as np 2 3  np.random.seed(3) 4 5  a = np.random.randint(5, 8, size=5) 6  print(a) </pre>		

---

**Ex 3.1.** Why do we set the seed?

### 3.3 Compute mean and std of simulated queue length for $\rho \approx 1$

We discuss the concept of load more formally in the course at a later point in time. Conceptually the load is the rate at which work arrives. For instance, if  $\lambda = 5$  jobs arrive per hour, and each requires 20 minutes of work (on average), then we say that the load is  $5 \times 20/60 = 5/3$ . Since one server can do only 1 hour of work per hour, we need at least two servers to deal with this load. We define the utilization  $\rho$  as the load divided by the number of servers present.

In discrete time, we define  $\rho$  as the average number of jobs arriving per period divided by the average number of jobs we can serve per period. Slightly more formally, for discrete time,

$$\rho \approx \frac{\sum_{k=1}^n a_k}{\sum_{k=1}^n c_k}. \quad (1)$$

And formally, we should take the limit  $n \rightarrow \infty$  (but such limits are a bit hard to obtain in simulation).

Python Code

```
1 import numpy as np
2
3 np.random.seed(3)
4 num = 5000
5
6 a = np.random.randint(5, 10, size=num)
7 c = 7 * np.ones(num)
8 L = np.zeros_like(a)
9
10 L[0] = 20
11 for i in range(1, num):
12     d = min(c[i], L[i-1])
13     L[i] = L[i-1] + a[i] - d
14
15 print(a.mean(), c.mean(), L.mean(), L.std())
```

**Ex 3.2.** Read the numpy documentation on `numpy.random.randint` to explain the range of random numbers that are given by this function. In view of that, what should `a.mean()` approximately be? Is that larger, equal or smaller than `c`?

### 3.4 plot the queue length process

Glue the next code after the other code.

Python Code

```
1 import matplotlib.pyplot as plt
2
3 plt.clf()
4 plt.plot(L)
5 plt.savefig('queue-discrete_1.pdf')
```

**Ex 3.3.** Comment on the graph you get. Did you expect such large excursions of the queue length process?

### 3.5 A trap: integers versus floats

Suppose that we change the arrival rate a bit.

---

Python Code

---

```

1 num = 5000
2
3 np.random.seed(3)
4 a = np.random.randint(5, 9, size=num)
5 c = (5+8)/2 * np.ones(num)
6 L = np.zeros_like(a)
7
8 L[0] = 20
9 for i in range(1, num):
10     d = min(c[i], L[i-1])
11     L[i] = L[i-1] + a[i] - d
12
13
14 plt.clf()
15 plt.plot(L)
16 plt.savefig('queue-discrete-1-1.pdf')

```

---

Don't forget that the numpy library has to be imported (as we did before).

**Ex 3.4.** Compare the definition of `a` and `c` to what we had earlier. What is `a.mean()` now approximately? Observe that the mean of `c` is now around 6.5.

Here is the figure.

You should see that it is very different from the one before. To explain why this is so, this somewhat cryptic question will be helpful, hopefully.

**Ex 3.5.** What is  $9 - 6.5$ ? What is `int(9-6.5)`, that is, run the code in python, and type in the answer in your answer sheet. What is the difference between  $9 - 6.5$  and `int(9-6.5)`? Explain that since `L` stores *integers* we actually use a service capacity of 7, *not* of 6.5.

**Ex 3.6.** The code below repairs this. Explain which line is different of the previous code. How did that change repair the problem? Now explain the title of this section. (BTW, I made this type of error many, many times. The reason to include these exercises is to make you aware of the problem, so that you can spot it when you make the same error.)

---

Python Code

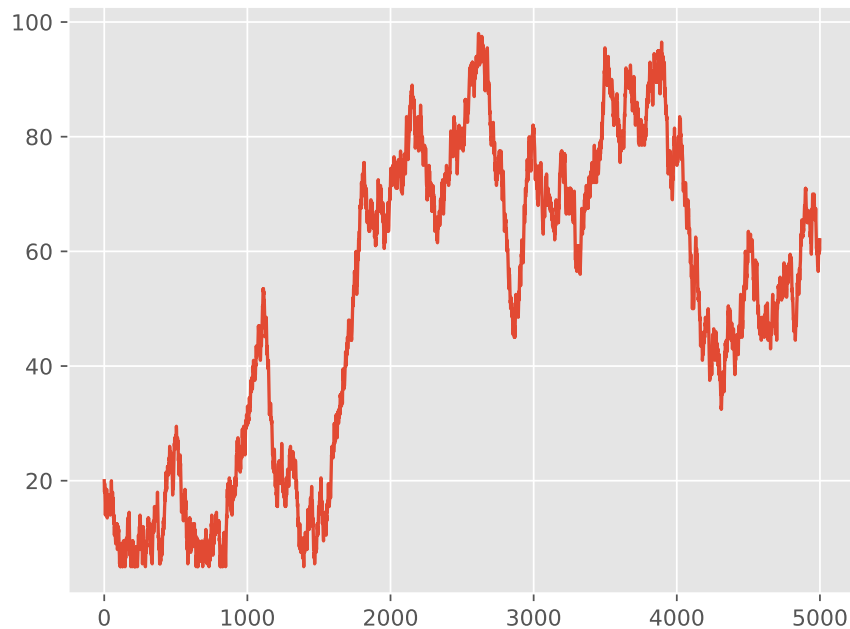
---

```

1 num = 5000
2
3 np.random.seed(3)
4 a = np.random.randint(5, 9, size=num)
5 c = (5+8)/2 * np.ones(num)
6 L = np.zeros_like(a, dtype=float)
7
8 L[0] = 20
9 for i in range(1, num):
10     d = min(c[i], L[i-1])
11     L[i] = L[i-1] + a[i] - d
12
13
14 plt.clf()
15 plt.plot(L)
16 plt.savefig('queue-discrete-1-2.pdf')

```

---



### 3.6 Serve arrivals in the same period as they arrive

**Ex 3.7.** Change the code such that the arrivals that occur in period  $i$  can also be served in period  $i$ . Include your code in your assignment. Then make a graph (include that too of course) and compare your results with the results of the simulation we do here (recall, in the simulations up to now, arrivals cannot be served in the periods in which they arrive).

### 3.7 Drift when $\rho > 1$

Python Code

```

1 num = 5_000
2
3 np.random.seed(3)
4 a = np.random.randint(5, 9, size=num)
5 c = (5+8)/2.3 * np.ones(num)
6 L = np.zeros_like(a, dtype=float)
7
8 L[0] = 20
9 for i in range(1, num):
10     d = min(c[i], L[i-1])
11     L[i] = L[i-1] + a[i] - d
12
13
14 plt.clf()
15 plt.plot(L)
16 plt.savefig('queue-discrete_2.pdf')

```

**Ex 3.8.** What is  $c.mean() - a.mean()$ ? Explain the slope of the line (fitted through the points.)

### 3.8 Drift when $\rho < 1$ , start with large queue.

Python Code

```
1 num = 5_000
2
3 np.random.seed(3)
4 a = np.random.randint(5, 9, size=num)
5 c = (5+8)/1.8 * np.ones(num)
6 L = np.zeros_like(a, dtype=float)
7
8 L[0] = 2_000
9 for i in range(1, num):
10     d = min(c[i], L[i-1])
11     L[i] = L[i-1] + a[i] - d
12
13
14 plt.clf()
15 plt.plot(L)
16 plt.savefig('queue-discrete_3.pdf')
```

Ex 3.9. Explain the slope of the lines (if you would fit that through the points.)

### 3.9 Hitting time

Ex 3.10. When  $\rho < 1$  and  $L_0$  is some large number, such as here, why is the normal distribution reasonable to model the time  $\tau$  until the queue hits zero for the first time?

Ex 3.11. The code below uses simulation to estimate the expected time  $E[\tau]$  and variance  $V[\tau]$ . Explain how it works; don't forget to address the following points:

1. Why don't I need to use  $d = \min(L, c)$ ?
2. Explain that running up to time 5000 is just a guess.

Python Code

```
1 import numpy as np
2
3 num = 5_000
4
5
6 def hitting_time(seed):
7     np.random.seed(seed)
8     a = np.random.randint(5, 9, size=num)
9     c = (5 + 8) / 1.8 * np.ones(num)
10
11     L = 2_000
12     for i in range(1, num):
13         L += a[i] - c[i]
14         if L <= 0:
15             return i
16
17
18 N = 10
19 tau = np.zeros(N)
20 for n in range(N):
21     tau[n] = hitting_time(n)
22
23 print(tau.mean(), tau.std())
```

**Ex 3.12.** Run the above simulation with some seeds to your liking. What are the numbers you get? Use properties of uniform random variables to explain the values you get.

### 3.10 Things to memorize

1. If the capacity is equal or less than the arrival rate, the queue length will explode.
2. If the capacity is larger than the arrival rate, the queue length will stay around 0 (between quotes).
3. If we start with a huge queue, but the service capacity is larger than the arrival rate, then the queue will drain rather fast, in fact, about linear.

## 4 Queues with blocking

Consider a queue that is subject to blocking: this means that when the queue exceeds  $K$ , say, then the excess jobs are rejected.

### 4.1 A simple rejection rule

Python Code

```

1 num = 500
2
3 np.random.seed(3)
4 a = np.random.randint(0, 20, size=num)
5 c = 10*np.ones(num)
6 L = np.zeros_like(a)
7 loss = np.zeros_like(a)
8
9 K = 30 # max people in queue, otherwise they leave
10
11 L[0] = 28
12 for i in range(1, num):
13     d = min(c[i], L[i-1])
14     loss[i] = max(L[i-1] + a[i] - d - K, 0) # this
15     L[i] = L[i-1] + a[i] - d - loss[i] # this 2
16
17 lost_fraction = sum(loss)/sum(a)
18 print(lost_fraction)

```

**Ex 4.1.** Explain how the line marked as *this* works, in other words, how does that line implement a queue with loss? In line *this 2* we subtract `loss[i]`; why?

**Ex 4.2.** Why is `dtype=float` not necessary in the definition of `L`?

Here is the graph of the (simulated) queue length process.

Python Code

```

1 plt.clf()
2 plt.plot(L)
3 plt.savefig('queue-discrete-loss.pdf')

```

**Ex 4.3.** Does the blocking rule work as it should?

## 4.2 Rejection at the start of the period

If we would assume that rejection occurs at the start of the period, the code has to be as follows:

---

Python Code

---

```
1  for i in range(1, num):
2      d = min(c[i], L[i-1])
3      loss[i] = max(L[i-1] + a[i] - K, 0)
4      L[i] = L[i-1] + a[i] - d - loss[i]
5
6  lost_fraction = sum(loss)/sum(a)
7  print(lost_fraction)
```

---

Ex 4.4. Explain the line in the code that changed.

Ex 4.5. Explain that this rule has the same effect as assuming that departures occur after the rejection.

## 4.3 Estimating rejection probabilities

With the code below we can estimate the distribution  $p_i = P[L = i]$ .

---

Python Code

---

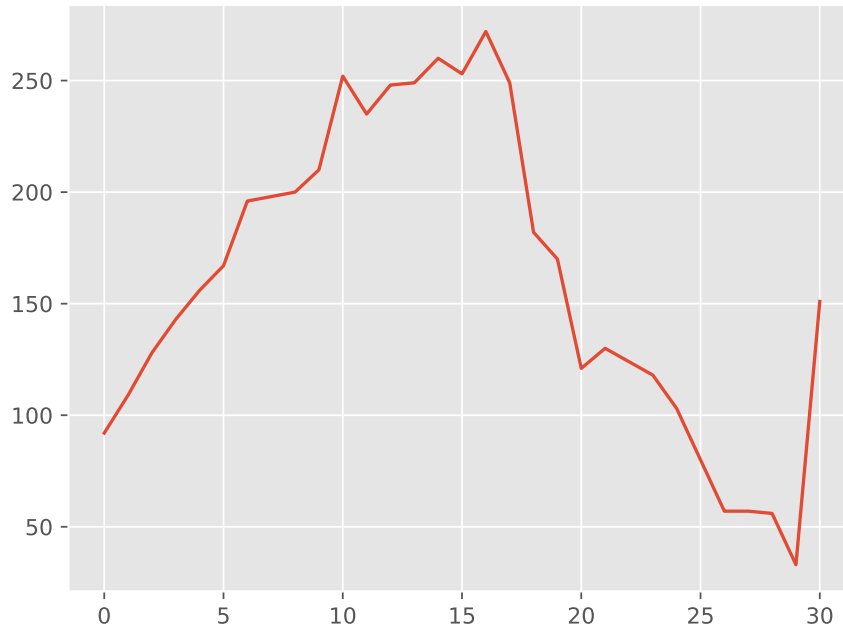
```
1  import numpy as np
2  import matplotlib.pyplot as plt
3
4  np.random.seed(3)
5
6  num = 5000
7
8  np.random.seed(3)
9  a = np.random.randint(0, 18, size=num)
10 c = 10 * np.ones(num)
11 L = np.zeros_like(a)
12
13 K = 30
14 p = np.zeros(K + 1)
15
16 L[0] = 28
17 for i in range(1, num):
18     d = min(c[i], L[i - 1])
19     L[i] = min(L[i - 1] + a[i] - d, K)
20     p[L[i]] += 1
21
22 plt.clf()
23 plt.plot(p)
24 plt.savefig('queue-discrete-loss-p.pdf')
```

---

Ex 4.6. Why should p have a length of K+1? Then explain how the code estimates  $p_i$ .

Here is the graph of p.





**Ex 4.7.** Note that  $p$  is not normalized (i.e., sums up to 1). The following code repairs that. Explain how it works.

---

```
1 p /= p.sum()
```

Python Code

---

## 4.4 Rejection probabilities under high loads

**Ex 4.8.** Change the parameters such  $\rho = 1.51$ . Then make again a plot of  $p$ . (Just copy my code, but change the parameters or the distribution of  $a$ , or  $c$ .)

**Ex 4.9.** Now change the parameters such  $\rho \approx 10$ , use a simple argument to show that  $\pi_{K-2} \approx 0$ . Check this intuition by doing a simulation with the appropriate parameters. Include your code and the graph of  $p$ .

## 5 Hints

**h.3.10.** When  $L_0 \gg 1$ , then  $a_k$  jobs arrive in period  $k$  and  $c_k$  jobs leave. For ease, write  $h_k = c_k - a_k$ . Then the time  $\tau$  to hit 0, i.e., the time until the queue is empty, is the smallest  $\tau$  such that  $\sum_{k=1}^{\tau} h_k \geq L_0$ . But then, by Wald's theorem:  $E[\tau] E[h_k] = L_0$ . What is  $E[\tau]$ ?

Moreover, consider any sum of random variables, then with  $\mu = E[X]$ ,  $\sigma$  the std of  $X$ , and  $N$  the normal distribution, and by the central limit theorem,

$$\frac{1}{n} \sum_{i=1}^n X_i \sim N\left(\mu, \frac{\sigma^2}{n}\right) \implies \sum_{i=1}^n X_i \sim N(n\mu, n\sigma^2). \quad (2)$$