

A multi-server queue with changing interarrival rates

EBB074A05

Nicky D. van Foreest

2020:12:17

1 General info

This file contains the code and the results that go with this youtube movie:

1.1 DONE Set theme and font size

Set the theme and font size so that it is easier to read on youbute

```
(load-theme 'material-light t)
(set-face-attribute 'default nil :height 200)
```

By the way, this is emacs lisp; you cannot run it in python.

1.2 Load standard modules

```
1 import numpy as np
2 from numpy.random import poisson, uniform, exponential
3 import matplotlib.pyplot as plt
4 import seaborn as sns; sns.set()
5
6 np.random.seed(3)
```

2 Supermarket with one fast server.

In the simulations models we considered up to now the arrival rate was constant. For supermarkets we know that this is not true. In this simulation we learn how to simulate a (simple) queueing system with time changing arrival rates.

2.1 Make data

We take the average number of arrivals per minute constant during one hour. Thus, $\lambda = 1$ means that during the first hour 1 customer arrives per minute, on average.

```
1 labdas = np.array([1, 2, 3, 3, 2.5, 2, 2, 2, 2, 3, 4, 3, 2.5, 2, 1])
2 # labdas = np.array([1, 2]) # Short case for testing
```

For ease, the shop opens at hour 1, and is open for this amount of hours.

```
1 len(labdas)
```

15

We now generate the number of arrivals in an hour. Let the number N_i of customers in hour i be $\sim \text{Pois}(60\lambda)$; recall λ is number of arrival per minute.

```
1 N = np.zeros(len(labdas), dtype=int)
2
3 for i, labda in enumerate(labdas):
4     N[i] = poisson(labda * 60)
5 N
```

61 113 162 150 155 126 112 113 118 182 231 176 169 116 63

The total number of arrivals is

```
1 N.sum()
```

2047

The average service time per customer is $E[S]$, in minutes.

```
1 ES = 2.5 # expected service time in minutes.
```

Finally, we need the service times. Let's take it as exponential.

```
1 S = exponential(ES, size=N.sum())
2 S.mean(), S.std() # check
```

2.506528761470696 2.4599902225851897

Here is an exercise. Do it after you have read the entire document.

Ex 2.1. Change the service time distribution for another distribution (take any distribution to your liking.). Analyze the impact on the KPIs.

We need some extra capacity, for otherwise the queue is critically loaded. Take the slack s ; initially I set it to 1. The number of servers planned for the i th hour is then $c_i = \lfloor \rho + s \rfloor$. Hence, c is at least as large as the loads.

```
1 loads = labdas * ES
2 loads
```

2.5 5 7.5 7.5 6.25 5 5 5 5 7.5 10 7.5 6.25 5 2.5

```
1 slack = 1
2 num_servers = np.array.loads + slack, dtype=int)
3 num_servers
```

3 6 8 8 7 6 6 6 6 8 11 8 7 6 3

2.2 First set of KPIs

The amount of planned capacity (hence payed fees to personnel) is $\sum_{i=1}^n c_i$.

```
1 total_offered_service_capacity = num_servers.sum() # in hours
2 total_offered_service_capacity
```

99

For the total amount of required service, add up all services and divide by 60 to get the time in hours.

```
1 total_offered_service = S.sum() / 60 # in hours
2 total_offered_service
```

85.51440624550858

```
1 idle = 1 - total_offered_service / total_offered_service_capacity
2 idle
```

0.13621811873223655

2.3 Make arrival times.

There is an interesting theorem that relates the Poisson and exponential distribution in the following way. Suppose that the inter-arrival times or jobs are $\sim \text{Exp}(\lambda)$. Then, given N arrivals during the interval $[0, T]$, the arrival times of the N jobs are uniform on $[0, T]$.

```
1 A = np.zeros(N.sum() + 1)
2 start = 1
3 for n in N:
4     A[start : start + n] = A[start - 1] + sorted(uniform(0, 60, size=n))
5     start += n
6
7 A[:10]
```

0 1.04562269 1.13263927 1.28932131 2.17167152 2.90243072 3.98007067 4.24975945 4.96445375 6.

Ex 2.2. While building the model the code was initially like this.

```
1 A = np.zeros(N.sum())
2 start = 1
3 for n in N:
4     A[start : start + n] = A[start - 1] + uniform(0, 60, size=n)
5     start = n
```

Explain which lines contain errors (there are 3) and why these are wrong.

2.4 A 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 just model the system as if there is one fast server.

```
1 W = np.zeros_like(S)
2 hour = 0
3 for k in range(1, len(W)):
4     X = A[k] - A[k - 1]
5     c = num_servers[int(A[k] / 60)]
6     W[k] = max(W[k - 1] + S[k - 1] / c - X, 0)
```

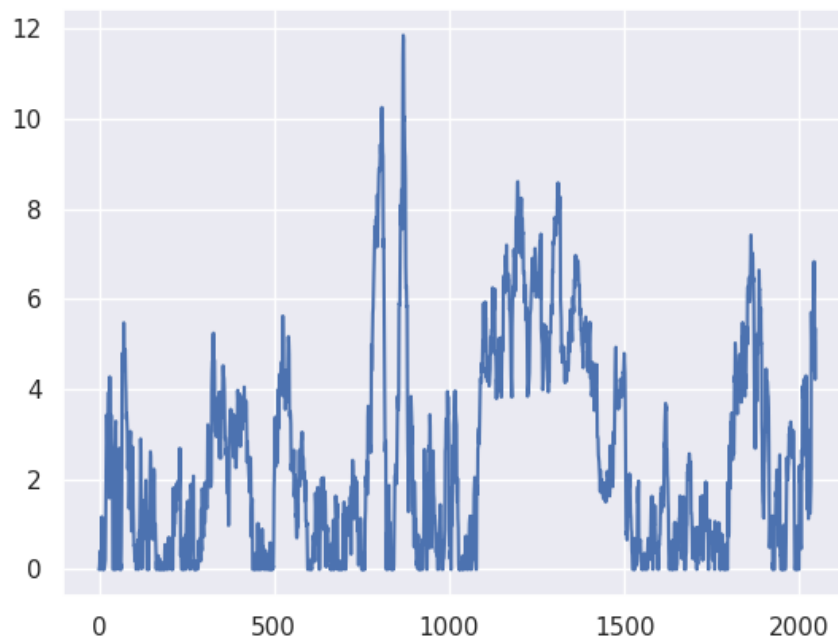
This is a bit tricky: $A[k]$ is the arrival time of job k . Then $A[k]/60$ is the hour in which job k arrives. Then we look up the number of servers planned for that hour.

Here are some KPIs.

```
1 W.mean(), W.std(), np.max(W)
```

2.4888178640316547 2.2978630075196977 11.854734178715747

```
1 plt.clf()
2 plt.plot(W)
3 plt.savefig('multi0.png')
4 'multi0.png'
```



Ex 2.3. Change the slack and see what happens to the queue length process and the KPIs.

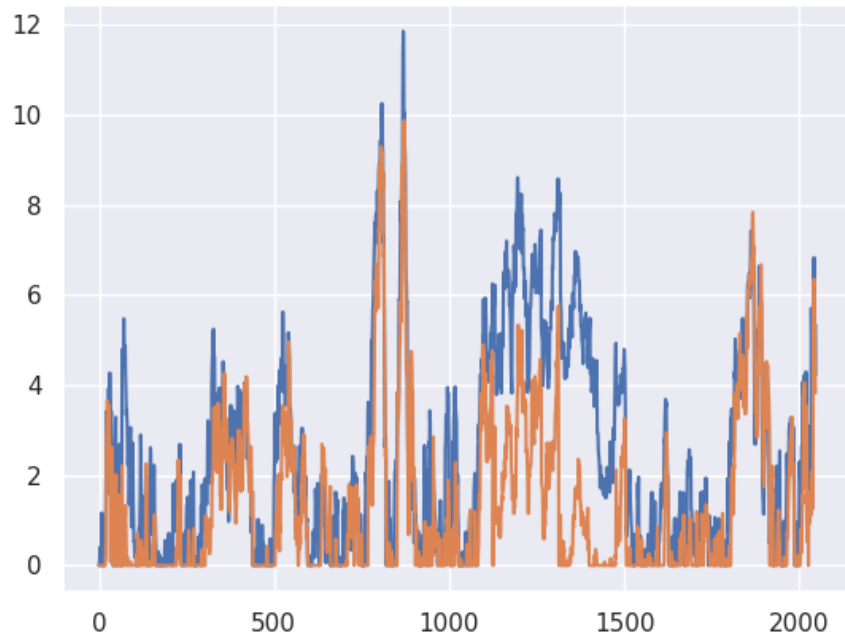
3 A real multiserver queue

```
1 m = max(num_servers)
2
3 one = np.ones(m, dtype=int) # vector with ones
4
5 w = np.zeros(m)
6 W = np.zeros_like(S)
7 server = np.zeros_like(S)
8
9 # explain why W[k] has to be above w[s]. And how I found this error.
10
11 for k in range(1, len(W)):
12     X = A[k] - A[k - 1]
13     c = num_servers[int(A[k] / 60)]
14     s = w[:c].argmin()
15     server[k] = s
16     W[k] = w[s]
17     w[s] += S[k] # store waiting time
18     w = np.maximum(0, w - X * one)
19
20
21 W.mean(), W.std(), np.max(W)
```

1.3141839485857727 1.7606254417452725 9.866951032209062

Let's plot it.

```
1 plt.plot(W)
2 plt.savefig('multi1.png')
3 'multi1.png'
```



Here I don't use `plt.clf()` because I want to compare the graphs of both models.

Ex 3.1. Explain line 4. Hint, compare the code of multi server we discussed earlier.

Ex 3.2. While building the model the code was initially like this.

```

1 for k in range(1, len(W)):
2     X = A[k] - A[k - 1]
3     c = num_servers[int(A[k] / 60)]
4     s = w[:c].argmin()
5     server[k] = s
6     w[s] += S[k]
7     W[k] = w[s]
8     w = np.maximum(0, w - X * one)

```

Explain which line contain an errors. What is wrong? How do you think I discovered the mistake (i.e., what tests did I do so that I found the problem)?

Ex 3.3. Make some graphs, 4 or so, in which you analyze different aspects of the setting. Recall, the aim is that you play with the code, you don't have to do anything special, or bright. Just changing a few of the parameters and making graphs is fine. (Do put labels, and explain what you see, and whether this is in line with what you expected before you did the experiments.)

Ex 3.4. Explain how a supermarket could use this simulator. What is the data they have to provide?

Ex 3.5. Finally, point out one aspect in which the simulator can be improved.

4 Restore my emacs settings

```
(load-theme 'material t)
(set-face-attribute 'default nil :height 100)
```