

Queueing theory assignment: Controlled queues

EBB074A05

Nicky D. van Foreest

2022:01:19

1 Blocking on waiting time

In this assignment we investigate queues under certain control rules. First we focus on blocking, then on switching on and off the server depending on the queue length.

Suppose we don't allow jobs to enter the system when the waiting time becomes too long. A simple rule is the block job k when $W_k \geq T$, for some threshold T .

How to simulate that?

1.1 A start with no blocking

Before doing something difficult, I tend to start from a situation that I do understand, which, in this case, is the single server queue without blocking.

Here is our standard code. With this we can find parameters that are suitable to see that blocking will have an impact. (Suppose $\lambda = 1$, $\mu = 1000$, $T = 100$, we will see no job being blocked during any simulation, at least not in this universe.)

Python Code

```
1 import numpy as np
2
3
4 np.random.seed(3)
5
6 num = 10000
7 labda = 1
8 mu = 1.1 * labda
9 T = 5
10 X = np.random.exponential(scale=1 / labda, size=num)
11 S = np.random.exponential(scale=1 / mu, size=len(X))
12 S[0] = 0
13
14 W = np.zeros_like(S)
15 for k in range(1, len(S)):
16     W[k] = max(W[k - 1] + S[k - 1] - X[k], 0)
17
18 print(W.mean(), W.max())
19 print((W >= T).mean()) # this
```

Ex 1.1. Run the code and check that indeed long waiting times do occur. What does the `this` line do?

Ex 1.2. Give the Kendall notation for the queueing model that we simulate here.

1.2 A hack to implement blocking

Here is a dirty hack to implement blocking. When $W_k \geq T$, job k should not enter. That means that its service time should not be added to the waiting time. But not adding the service time can be achieved by setting $S_k = 0$. To realize this, change the for loop to

Python Code

```
1 for k in range(1, len(S)):
2     W[k] = max(W[k - 1] + S[k - 1] - X[k], 0)
3     if W[k] >= T:
4         S[k] = 0
```

Ex 1.3. Run this code and check its effect on $E[W]$, $V[W]$, and $\max\{W\}$. Explain how it can occur that the maximum can still be larger than T .

1.3 A better method

As a matter of principle, I don't like the code of the previous section. In my opinion such hacks are a guarantee on bugs that can be very hard to find later. Mind, with this trick I am changing my primary data, in this case the service times. Reuse these service times at a later point in the code, for instance for a comparison with other models or for testing, has become impossible. And if I forget this (when I use this code maybe half a year later), then finding the bug will be very hard. Hence, as a golden rule: don't touch the primary data.

However, modifying this code so that I can block seems not straightforward.

Ex 1.4. The recursion we use, i.e.,

Python Code

```
1 for k in range(1, len(S)):
2     W[k] = max(W[k - 1] + S[k - 1] - X[k], 0)
```

assumes that job $k - 1$ is accepted. Suppose that $W[k] \geq T$ so that job k is blocked. Explain that with this recursion, we now have a problem to compute the waiting time for job $k + 1$.

Here is better code.

Python Code

```
1 W = np.zeros_like(S)
2 I = np.ones_like(S)
3 for k in range(1, len(S)):
4     W[k] = max(W[k - 1] + S[k - 1] * I[k - 1] - X[k], 0)
5     if W[k] >= T:
6         I[k] = 0
7
8 print(W.mean(), I.mean())
```

Ex 1.5. What does the vector I represent?

1.4 Some other blocking rules

There are other rules to block jobs.

Ex 1.6. In this code,

Python Code

```

1 W = np.zeros_like(S)
2 I = np.ones_like(S)
3 for k in range(1, len(S)):
4     W[k] = max(W[k - 1] + S[k - 1] * I[k - 1] - X[k], 0)
5     if W[k] + S[k] >= T:
6         I[k] = 0
7
8 print(W.mean(), W.max(), I.mean())

```

how does the rule below block jobs?

Ex 1.7. Likewise,

Python Code

```

1 W = np.zeros_like(S)
2 V = np.ones_like(S)
3 for k in range(1, len(S)):
4     W[k] = max(W[k - 1] + V[k - 1] - X[k], 0)
5     V[k] = min(T - W[k], S[k])
6
7 print(W.mean(), W.max(), S.mean() - V.mean())

```

how does the rule below block jobs? What is the meaning of V?

2 Batch queues and blocking on waiting time

Let us now set up a simulation to see the combined effect of batch arrivals and blocking on waiting time.

Recall, in the queueing book we discuss some methods to block jobs in the $M^X/M/1$ queue when the queue length (not the waiting time) is too long. We tackle blocking on queue length in a separate section below.

2.1 Again start without blocking

We need a slightly different way to generate service times. When a batch of B_k jobs arrives at time A_k , then the service time added to the waiting is the sum of the service times of all B_k jobs in the batch.

Python Code

```

1 import numpy as np
2 from scipy.stats import expon
3
4 np.random.seed(3)
5
6 num = 100000
7 labda = 1
8 mu = 1.1 * labda
9 X = np.random.exponential(scale=1 / labda, size=num)
10 B = np.random.randint(1, 2, size=num)
11 S = expon(scale=1 / mu)
12
13 W = np.zeros_like(X)
14 for k in range(1, len(W)):
15     W[k] = max(W[k - 1] + S.rvs(B[k]).sum() - X[k], 0)
16

```

```

17 print(S.mean(), W.mean(), W.max())
18 rho = labda / mu
19 print(rho**2 / (1 - rho))

```

Ex 2.1. Explain how this code works.

Ex 2.2. Run the code. Why do I take B as it is here? Why should $W.mean()$ and $\rho^2/(1 - \rho)$ be approximately equal

2.2 Include blocking

Here is the code with a blocking rule.

	Python Code
1	import numpy as np
2	from scipy.stats import expon
3	
4	np.random.seed(3)
5	
6	num = 1000
7	labda = 1
8	mu = 3.1 * labda
9	T = 5
10	X = np.random.exponential(scale=1 / labda, size=num)
11	B = np.random.randint(1, 5, size=num)
12	S = expon(scale=1 / mu)
13	
14	W = np.zeros_like(X)
15	V = np.zeros_like(W)
16	for k in range(1, len(W)):
17	W[k] = max(W[k - 1] + V[k - 1] - X[k], 0)
18	V[k] = S.rvs(B[k]).sum() if W[k] < T else 0
19	
20	print(S.mean() * B.mean() - V.mean())
21	print(W.mean(), W.max())
22	print(np.isclose(V, 0).mean())
23	print((V <= 0).mean()) # this

Ex 2.3. Explain how the code works. What do the printed KPIs denote? Finally, why is the `this` line very dangerous to use, hence to be avoided? (Use `np.isclose` instead!)

3 Blocking on queue length

Blocking on queue length is quite a bit harder with a simulation in continuous time because we need to keep track of the number of jobs in the system. (Recall in discrete time the recursions to compute $\{L\}$ are easy, while in continuous time the recursions for $\{W\}$ or $\{J\}$ are easy.)

3.1 Start without blocking

As before, I start from a code that I really understand, and then I extend it to a situation that I find more difficult. So, here is code to find the system length L at *arrival* epochs $\{A_k\}$.

```

1  import numpy as np
2
3  np.random.seed(3)
4
5  num = 10000
6  labda = 1
7  mu = 1.5 * labda
8  X = np.random.exponential(scale=1 / labda, size=num)
9  A = np.zeros(len(X) + 1)
10 A[1:] = X.cumsum()
11 S = np.random.exponential(scale=1 / mu, size=len(A))
12 S[0] = 0
13 D = np.zeros_like(A)
14 L = np.zeros_like(A, dtype=int)
15
16 idx = 0
17 for k in range(1, len(A)):
18     D[k] = max(D[k - 1], A[k]) + S[k]
19     while D[idx] < A[k]:
20         idx += 1
21     L[k] = k - idx
22
23 rho = labda / mu
24 print(L.mean(), rho/(1-rho), L.max())
25 print((L == 0).mean(), 1 - rho)
26 print((L == 1).mean(), (1 - rho)*rho)

```

Ex 3.1. Explain how this computes $L[k]$. Do we count the system length as seen upon arrival, or does $L[k]$ include job k , i.e., the job that just arrived?

Ex 3.2. Just to check that you really understand: why is there no problem with the statement $(L == 0)$?

Ex 3.3. Why do I compare $L.mean()$ to $\rho/(1 - \rho)$ and not to $\rho^2/1 - \rho$?

Ex 3.4. Change μ to 1.05λ . Now the results of the simulation are not very good if $num=1000$ or so. Making num much larger does the job, though.

3.2 Include blocking

It might seem that we are now ready to implement a continuous time queueing system with blocking on the queue length. Why not merge the ideas we developed above? Well, because this does not work.

(If you like a challenge, stop reading here, and try to see how far you can get with developing a simulation for this situation.)

Only after having worked for 3 hours I finally saw ‘the light’. As a matter of fact, I needed a new data structure, a deque from which we can pop and append jobs at either end of a list. Here is the code.

```

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

```

```

5
6 num = 10000
7 labda = 1
8 mu = 1.2 * labda
9 T = 5
10 X = np.random.exponential(scale=1 / labda, size=num)
11 A = np.zeros(len(X) + 1)
12 A[1:] = X.cumsum()
13 S = np.random.exponential(scale=1 / mu, size=len(A))
14 S[0] = 0
15 D = np.zeros_like(A)
16 L = np.zeros_like(A, dtype=int)
17
18 Q = deque(maxlen=T + 1)
19 for k in range(1, len(A)):
20     while Q and D[Q[0]] < A[k]:
21         Q.popleft()
22     L[k] = len(Q)
23     if len(Q) == 0:
24         D[k] = A[k] + S[k]
25         Q.append(k)
26     elif len(Q) < T:
27         D[k] = D[Q[-1]] + S[k]
28         Q.append(k)
29     else:
30         D[k] = A[k]

```

Ex 3.5. Read the documentation of how a deque works, then explain the code.

Ex 3.6. At first I had these lines:

	Python Code	
<pre> 1 elif len(Q) < T: 2 Q.append(k) 3 D[k] = D[Q[-1]] + S[k] </pre>		

Why is that wrong?

Ex 3.7. What goes wrong with this code:

	Python Code	
<pre> 1 if len(Q) < T: 2 D[k] = max(D[Q[-1]], A[k]) + S[k] 3 Q.append(k) </pre>		

Ex 3.8. What queueing discipline would result if we would use the `pop()` and `appendleft()` methods of a deque?

Ex 3.9. What queueing discipline would result if we would use the `pop()` and `append()` methods of a deque?

Ex 3.10. Run this code with $T=100$ and compare this with the queueing system without blocking. Why should you get the same results? (Realize that this is a check on the correctness of our code.)

Glue the next code (for the theoretical model) at the end of the previous code.

Python Code

```
1 rho = labda / mu
2 p = np.ones(T + 1)
3 for i in range(1, T + 1):
4     p[i] = rho * p[i - 1]
5 p /= p.sum()
6 for i in range(T + 1):
7     print((L == i).mean(), p[i])
```

Ex 3.11. Now set $T=5$ and $\text{num} = 10000$ or so. Run the code. Why do the result agree with the theoretical model? Why is this the $M/M/1$ queue?

In fact, I used the above theoretical model to check whether the simulation was correct. (My first 20 or so attempts weren't.)

4 An algorithm for the $M/G/1$ queue with blocking

In the queueing book we develop an algorithm to compute $\pi(n)$. Here we implement this, use this as another test on the simulator, and improve our understanding of queueing systems.

4.1 The algorithm

This is the code.

Python Code

```
1 import numpy as np
2 from scipy.integrate import quad
3 from scipy.stats import expon
4
5 np.random.seed(3)
6
7 labda = 1
8 mu = 1.2 * labda
9 T = 5
10 S = expon(scale=1 / mu)
11
12
13 def g(j, x):
14     res = np.exp(-labda * x) * (labda * x) ** j * S.pdf(x)
15     return res / np.math.factorial(j)
16
17
18 f = np.zeros(T + 1)
19 for j in range(T + 1):
20     f[j] = quad(lambda x: g(j, x), 0, np.inf)[0]
21
22 F = f.cumsum()
23 G = 1 - F
24
25 pi = np.ones(T + 1)
26 for n in range(T):
```

```

27     pi[n + 1] = pi[0] * G[n]
28     pi[n + 1] += sum(pi[m] * G[n + 1 - m] for m in range(1, n + 1))
29     pi[n + 1] /= f[0]
30
31 pi /= pi.sum()
32 print(pi)

```

Ex 4.1. Which formulas (give the numbers) of the queueing book have we implemented?

Ex 4.2. Run this code after the computation of f.

Python Code

```

1 j = 2
2 print(mu / (mu + labda) * (labda / (labda + mu)) ** j, f[j])

```

Why should these numbers be the same?

Ex 4.3. Run the code for $\mu = 0.3$ and compare the numerical results to what you get from:

Python Code

```

1 rho = labda / mu
2 p = np.ones(T + 1)
3 for i in range(1, T + 1):
4     p[i] = rho * p[i - 1]
5 p /= p.sum()
6 print(pi)

```

Explain why you should get the same numbers.

Ex 4.4. When the service times are constant, explain that this code computes f correctly:

Python Code

```

1 from scipy.stats import expon, uniform, poisson
2 # include useful code here
3 f = poisson(labda / mu).pmf(range(T + 1))

```

Then change the S in the simulation part to

Python Code

```

1 S = np.ones(len(A)) / mu

```

Run the code and include your results; of course the simulation and the algorithm should give more or less the same results.

Ex 4.5. As another good example, take $S \sim U(0, 2/\mu)$. The relevant code changes are this:

Python Code

```

1 from scipy.stats import expon, uniform, poisson
2 # other stuff for the model
3 S = uniform(0, 2 / mu)

```

and for the simulator:

Python Code

```

1 S = np.random.uniform(0, 2 / mu, size=len(A))

```

Run the code, and include your output.

4.2 Effect of blocking on performance

Ex 4.6. Take $\lambda = 1$ and $\mu = 1.1$. Use the algorithm to compute the loss probability and $E[L]$ for $T = 5$, $T = 10$ and $T = 15$. Include the numbers.

Ex 4.7. Do the same computations for $\mu = 0.5\lambda$. Why is the loss probability not so sensitive to T ?

Ex 4.8. Set $\mu = 1.2\lambda$ again. Then compare the loss probability for $T = 5, 10, 15$ for $S \sim \text{Exp}(\mu)$ and $S \sim U(0, 2/\mu)$. What is the influence of service time variability on the loss when $T = 5$, $T = 10$, $T = 15$? Why is this influence relatively more important for larger T ?

5 Server control

With blocking we control whether jobs are allowed to enter the system. We can also focus on another type of control, namely that of the server. Here we show how to simulate a system in which the server switches on when the waiting time exceeds a level D . When the server is empty again, it switches off.

Let us first plot the virtual waiting time. Earlier we discussed how to plot the virtual waiting time for a given array of waiting times, arrival times and departure times. Thus, the only relevant code is how to find the waiting time under a D policy.

Note that we here use the letter T to refer to the threshold since the letter D is already given to the departure times.

Python Code

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 np.random.seed(3)
5
6 num = 40
7 labda = 1
8 mu = 1.1 * labda
9 T = 10 # this acts as the threshold
10 X = np.random.exponential(scale=1 / labda, size=num)
11 X[0] = 0
12 A = np.zeros_like(X)
13 A = X.cumsum()
14 S = np.ones(len(A)) / mu
15 S[0] = 0
16 D = np.zeros_like(A)
17
18
19 W = np.zeros_like(S)
20 On = False
21 for k in range(1, len(S)):
22     if On:
23         W[k] = max(W[k - 1] + S[k - 1] - X[k], 0)
24         On = False if W[k] <= 0 else True
25     else:
26         W[k] = W[k - 1] + S[k - 1]
27         On = False if W[k] < T else True
28     D[k] = A[k] + W[k] + S[k]
29
30 idx = np.where(W <= 0)[0]
31
```

```

32 empty = D[idx[1:] - 1]
33
34 E = np.zeros((2 * len(A) + len(empty), 2)) # epochs
35 E[: len(A), 0] = A
36 E[: len(A), 1] = W
37 E[len(A) : 2 * len(A), 0] = A
38 E[len(A) : 2 * len(A), 1] = W + S
39 E[2 * len(A) : 2 * len(A) + len(empty), 0] = empty
40 E[2 * len(A) : 2 * len(A) + len(empty), 1] = 0
41 E = E[np.lexsort((E[:, 1], E[:, 0]))]
42
43 plt.plot(E[:, 0], E[:, 1])
44 plt.savefig("D-policy.pdf")

```

Ex 5.1. Explain how the waiting times are computed.

Given cost K to switch on the server and holding h (per unit waiting time per unit time) we want to find the threshold T that minimizes the time-average cost. The code below shows how to compute the cost for a given T .

Python Code

```

1 h = 1.0
2 K = 3
3 cost = 0
4 epoch, height = E[:, 0], E[:, 1]
5 for i in range(1, len(epoch)):
6     dx = epoch[i] - epoch[i - 1]
7     dy = (height[i] + height[i - 1]) / 2.0
8     cost += h * dx * dy
9     if dy == 0:
10         cost += K
11 print(cost / D[-1])

```

Ex 5.2. Explain how the code works. What is dx , what is dy ? Why do we divide by 2?

Ex 5.3. Explain the procedure to find the best T .

6 Hints

h.1.2. It's not the $D/D/1$ queue.

h.1.3. Check the section in the queueing book on the $M^X/M/1$ queue, in particular the section on blocking. Which blocking method do we include here?

h.1.4. We want to block job k when $W_k \geq T$. That means that job $k+1$ should not see S_k , but (as far as I see) we cannot convert that into an elegant recursion for W_{k+1}

h.1.5. If $I[k] == 1$, then what happens to job k ?

h.2.2. Why is this the $M/M/1$ queue when the batches $B = \text{np.random.randint}(1, 2, \text{size=num})$? (Recall my obsession with testing code.)

h.3.1. When the while loop terminates, is idx the index of the last departure, or does it point to the job that is the first to leave?

h.3.2. Is L a float?

h.3.3. What is $\rho^2/(1 - \rho)$?

h.3.6. To what job does $Q[-1]$ point to? It's one to far!

h.3.7. When Q is empty (i.e., $\text{len}(Q)$ is 0), then what is $Q[-1]$?

h.3.8. Does it matter whether we push jobs from right to left through a queue, rather than from left to right?

h.3.9. It's not FIFO.

h.3.10. Is $L.\text{max}()$ larger than 100 for this simulation?

h.4.6. Why is the loss probability equal to π_T ?