

Queueing theory Assignments

EBBo74Ao5

Nicky D. van Foreest

2022:02:16

CONTENTS

1	Simulation of queues in discrete time	3
2	Queueing control: Psychiatrists doing intakes	13
3	Simulation of Queueing processes in continuous time	21
4	Sakasegawa's formula and Batching rules	30
5	Control of Queues	43
6	Queueing theory assignment: Simulation with event stacks	54
7	Hints	64

mv assignments.pdf ../

INTRODUCTION

ASSIGNMENTS

These assignments are meant to show you how I build up computer programs (in python) from scratch to analyze and simulate queueing systems. I include the code with the idea that you run it yourself as you read through the tutorial. For that reason many code examples are somewhat long, so that you just have to copy-paste the code to have a fully running example.

I included many (simple) exercises to help you think about the code. As such, most of these exercises ask to explain the code. Of course, for the assignments, you only have to explain the *relevant* parts of the code, that is, part that relate to the queueing model. You should skip, so called, boiler plate code, like importing numpy. So keep your answers short; you don't have to win the Nobel prize on literature.

You should read the relevant section of my [queueing book](#) to understand what it going on in an assignment. Most of it is very easy, but without background it will be quite cryptic (I believe).

In 2020 I made some youtube movies to illustrate the material. However, this year I revised the code at several places. Hence the movies discuss most (but not all) of the material. You can find the youtube movies [here](#). Below I include the links to each of the videos organized per section.

The document `assignment-answer-template.tex` is the template that you have to use for making your L^AT_EX assignment.

RUNNING PYTHON

You can install python on your computer (preferred), or run it in a browser if you don't have to install python locally.

- [diggy](#)
- [colab](#)

There is lots of info on the web on how to setup things.

If you plan to install python on your machine, the best installation is perhaps [anaconda](#).

I prefer to work within emacs (an editor), and run the code in a terminal. This works much faster and more conveniently, but requires a bit (but not much) of intellectual effort and investment in how to configure things on your computer. If you're interested in the power tools I use, check out my [tools page](#).

SIMULATION OF QUEUES IN DISCRETE TIME

1.1 INTRO

Here I show how to set up an environment in python to simulate queueing systems in discrete time.

Perhaps you find the following youtube movies helpful.

- <https://youtu.be/DfYxayoQmjYc>
- <https://youtu.be/D8BIAoBICnw>
- https://youtu.be/_BoagRyH5c0

1.2 DETERMINISTIC QUEUES

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

1.2.1 *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)
```

1.2.2 *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 1.2.1. Add a third period, and report your result in the assignment.

1.2.3 *many periods, use vectors and for loops*

Python Code

```

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)

```

Ex 1.2.2. Run the code for 10 periods and report your result.

Ex 1.2.3. What will be the queue after 100 or so periods?

1.3 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.

1.3.1 *simulate demand*

Python Code

```

1  import numpy as np
2
3  a = np.random.randint(5, 8, size=5)
4  print(a)

```

1.3.2 *Set seed*

Python Code

```

1  import numpy as np
2
3  np.random.seed(3)
4
5  a = np.random.randint(5, 8, size=5)
6  print(a)

```

Ex 1.3.1. Why do we set the seed?

1.3.3 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 ρ as the load divided by the number of servers present.

In discrete time, we define ρ 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 1.3.2. 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`?

1.3.4 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 1.3.3. Comment on the graph you get. Did you expect such large excursions of the queue length process?

1.3.5 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 1.3.4. 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.

When you make the plot you should see that is is very different from the one before. To explain why this is so, the following somewhat cryptic question will be helpful, hopefully.

Ex 1.3.5. 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. Explain the difference between these two numbers? Explain that since `L` stores *integers*, not floats, we actually use a service capacity of 7, *not* of 6.5.

Ex 1.3.6. The code below repairs the above problem. Explain which line is different from the previous code. How did that change repair the problem? Now explain the title of this section. (BTW, I made this type of error with floats and ints 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

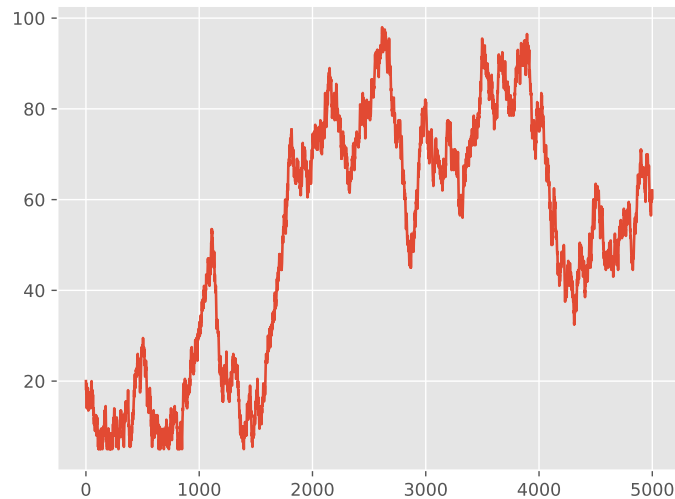
```

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')

```



1.3.6 Serve arrivals in the same period as they arrive

Ex 1.3.7. 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).

1.3.7 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 1.3.8. What is `c.mean()` - `a.mean()`? Explain the slope of the line (fitted through the points.)

1.3.8 *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 1.3.9. Explain the slope of the lines (if you would fit that through the points.)

1.3.9 *Hitting time*

Ex 1.3.10. When $\rho < 1$ and L_0 is some large number, such as here, why is it reasonable to use a normal distribution to model the time τ until the queue hits zero for the first time? (Assume that jobs can be served in the period in which they can arrive.)

Ex 1.3.11. The code below uses simulation to estimate the expected time $E[\tau]$ and variance $V[\tau]$. Explain how we use *simulation to estimate* $E[\tau]$ and $V[\tau]$. (You don't have to find an closed form expression for $E[\tau]$ or $V[\tau]$). 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 1.3.12. 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.

1.3.10 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.

1.4 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.

1.4.1 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 1.4.1. 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 $\text{loss}[i]$; why?

Ex 1.4.2. Why is, in this case, `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 1.4.3. Does the blocking rule work as it should?

1.4.2 Rejection at the start of the period

If we would assume that rejection occurs as 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 1.4.4. Explain the line in the code that changed.

Ex 1.4.5. Explain that this rule has the same effect as assuming that departures occur after the rejection.

1.4.3 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

```

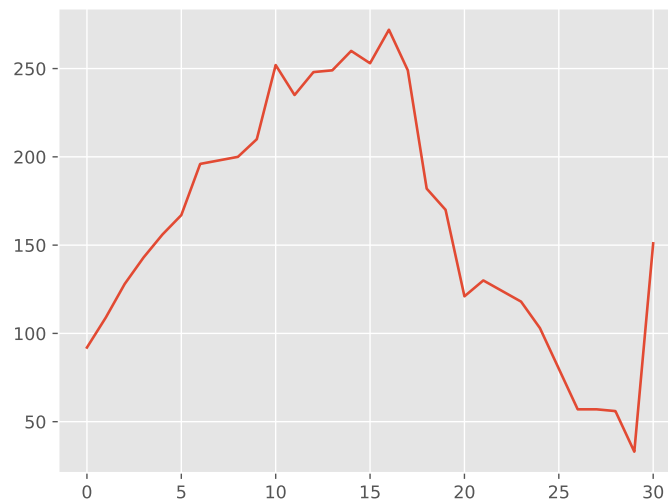
```

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

Here is the graph of p .



Ex 1.4.7. Note that p is not normalized (i.e., sums up to 1). The following code repairs that. Explain how it works.

```

_____ Python Code _____
1  p /= p.sum()
_____

```

1.4.4 *Rejection probabilities under high loads*

Ex 1.4.8. 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 1.4.9. 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 .

QUEUEING CONTROL: PSYCHIATRISTS DOING INTAKES

2.1 INTRO

There are 5 psychiatrists doing intakes. In their current organization, the queue of patients waiting for intakes is way too long, much longer than they like to see. Here I consider some strategies to deal with this controlling the queue length process, and I use simulation to evaluate how successful these are.

- <https://youtu.be/bCU3oP6r-00>

2.2 BASE SITUATION

Five psychiatrists do intakes. See the queueing book for further background.

2.2.1 Load standard modules

We need some standard libraries for numerical work and plotting.

	Python Code	
<hr/>		
<pre> 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) </pre>		
<hr/>		

2.2.2 Simulate queue length

	Python Code	
<hr/>		
<pre> 1 def computeQ(a, c, Q0=0): # initial queue length is 0 2 N = len(a) 3 Q = np.empty(N) # make a list to store the values of Q 4 Q[0] = Q0 5 for n in range(1, N): 6 d = min(Q[n - 1], c[n]) 7 Q[n] = Q[n - 1] + a[n] - d 8 return Q </pre>		
<hr/>		

Ex 2.2.1. Compute the queue lengths when

	Python Code	
<hr/>		
<pre> 1 a = [1,2,5,6,8,3,7,3] 2 c = [2,2,0,5,4,4,3,2] </pre>		
<hr/>		

and include the results in your report.

2.2.3 Arrivals

We start with run length 10 for demo purpose; later we extend to longer run times

Python Code

```

1 a = np.random.poisson(11.8, 10)
2 print(a)

```

2.2.4 Service capacity

Python Code

```

1 def unbalanced(a):
2     p = np.empty([5, len(a)])
3     p[0, :] = 1.0 * np.ones_like(a)
4     p[1, :] = 1.0 * np.ones_like(a)
5     p[2, :] = 1.0 * np.ones_like(a)
6     p[3, :] = 3.0 * np.ones_like(a)
7     p[4, :] = 9.0 * np.ones_like(a)
8     return p
9
10 p = unbalanced(a)
11 print(p)

```

Ex 2.2.2. Which psychiatrist does the most intakes per week?

2.2.5 Include holidays

Python Code

```

1 def spread_holidays(p):
2     for j in range(len(a)):
3         psych = j % 5
4         p[psych, j] = 0
5
6 spread_holidays(p)
7 print(p)

```

Ex 2.2.3. What is the long-run time average of the weekly capacity?

2.2.6 Total weekly service capacity

Python Code

```

1 s = np.sum(p, axis=0)
2 print(s)

```

Ex 2.2.4. Explain why we need to take the sum over axis=0 to compute the average weekly capacity for the intakes.

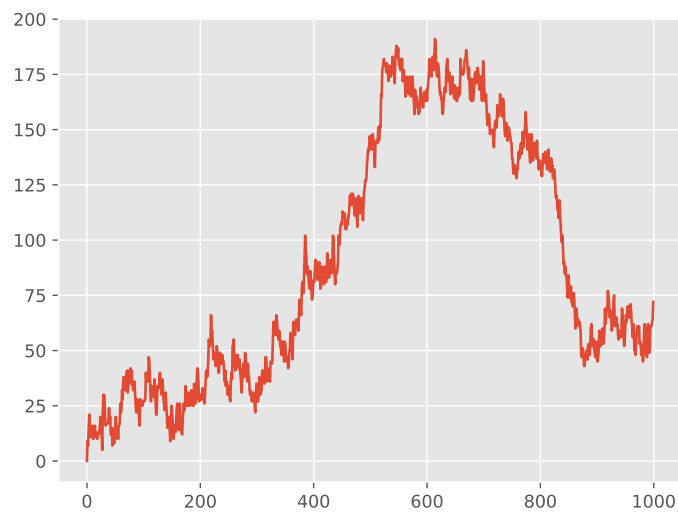
2.2.7 Simulate the queue length process

Python Code

```

1  np.random.seed(3)
2
3  a = np.random.poisson(11.8, 1000)
4  p = unbalanced(a)
5  spread_holidays(p)
6  s = np.sum(p, axis=0)
7
8  Q1 = computeQ(a, s)
9
10 plt.clf()
11 plt.plot(Q1)
12 plt.savefig("psych1.pdf")

```



2.3 EVALUATION OF BETTER (?) PLANS

2.3.1 Balance the capacity more evenly over the psychiatrists

I set the seed to enforce a start with the same arrival pattern.

Python Code

```

1  def balanced(a):
2      p = np.empty([5, len(a)])
3      p[0, :] = 2.0 * np.ones_like(a)
4      p[1, :] = 2.0 * np.ones_like(a)
5      p[2, :] = 3.0 * np.ones_like(a)
6      p[3, :] = 4.0 * np.ones_like(a)
7      p[4, :] = 4.0 * np.ones_like(a)
8      return p
9
10 np.random.seed(3)
11 a = np.random.poisson(11.8, 1000)
12
13

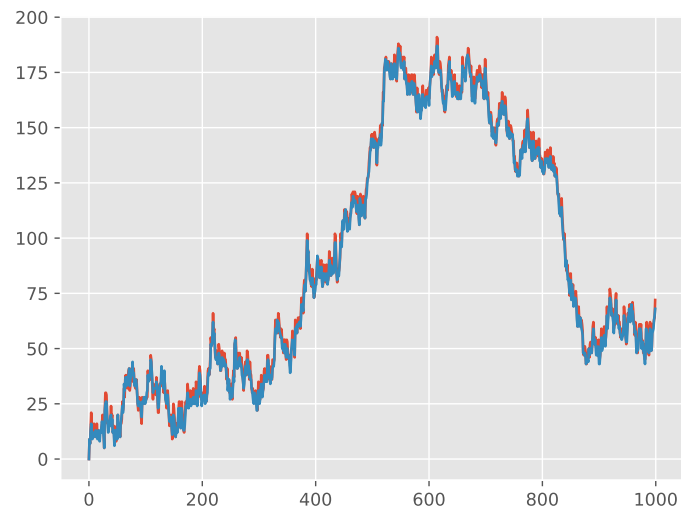
```



```

14 p = balanced(a)
15 spread_holidays(p)
16 s = np.sum(p, axis=0)
17 Q2 = computeQ(a, s)
18
19 plt.plot(Q2)
20 plt.savefig("psych2.pdf")

```

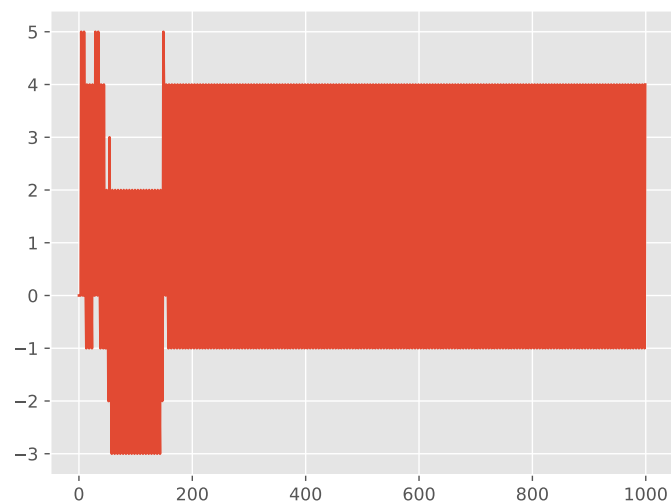


Python Code

```

1 plt.clf()
2 plt.plot(Q1-Q2)
3 plt.savefig("psych22.pdf")

```



Ex 2.3.1. How can we see that the effect of balancing capacity is totally uninteresting?

Ex 2.3.2. Change the capacities of the psychiatrists but keep the average weekly capacity the same. Include a graph of your result.

2.3.2 Synchronize holidays

What is the effect of all psychiatrists taking holidays in the same week?

Python Code

```

1  a = np.random.poisson(11.8, 10)
2
3
4  def synchronize_holidays(p):
5      for j in range(int(len(a) / 5)):
6          p[:, 5 * j] = 0 # this
7
8  p = unbalanced(a)
9  synchronize_holidays(p)
10 print(p)

```

Ex 2.3.3. Explain how the code works. Specifically, what does the line marked as this.

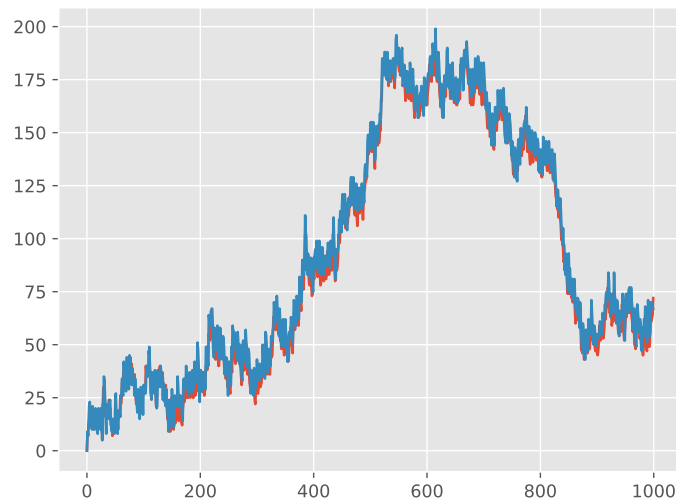
Ex 2.3.4. Change the code such that psychiatrists go on holiday every 6 weeks. However, modify the weekly capacities of the psychiatrists such that the total average weekly capacity remains the same. Include your code, and check with a sum (over an appropriate axis) that the average weekly capacity is still the same after your changes.

Python Code

```

1  np.random.seed(3)
2
3  a = np.random.poisson(11.8, 1000)
4  p = unbalanced(a)
5  spread_holidays(p)
6  s = np.sum(p, axis=0)
7  Q3 = computeQ(a, s)
8
9  plt.clf()
10 plt.plot(Q3)
11
12 p = balanced(a)
13 synchronize_holidays(p)
14 s = np.sum(p, axis=0)
15 Q4 = computeQ(a, s)
16
17 plt.plot(Q4)
18 plt.savefig("psych3.pdf")

```



Ex 2.3.5. Just to improve your coding skills (and your creativity), formulate another vacation plan. Implement this idea in code, and test its success/failure. Make a graph to show its effect on the dynamics of the queue length. (I don't mind whether your proposal works or not; as long as you 'play' and investigate, all goes.) Include your code—if you ported all this code to R, then include your R code—and comment on the difficult points.

Most probably, your proposals will also not solve the problem. We need something smarter.

2.4 CONTROL CAPACITY AS A FUNCTION OF QUEUE LENGTH

2.4.1 Simple on-off strategies

Let's steal an idea from supermarkets: dynamic control.

Python Code

```

1 lower_thres = 12
2 upper_thres = 24
3
4 def computeQExtra(a, c, e, Q0=0): # initial queue length is 0
5     N = len(a)
6     Q = [0] * N # make a list to store the values of Q
7     Q[0] = Q0
8     for n in range(1, N):
9         if Q[n-1] < lower_thres:
10             C = c - e
11         elif Q[n-1] >= upper_thres:
12             C = c + e
13         d = min(Q[n-1], C)
14         Q[n] = Q[n-1] + a[n] - d
15     return Q
16
17
18 np.random.seed(3)

```

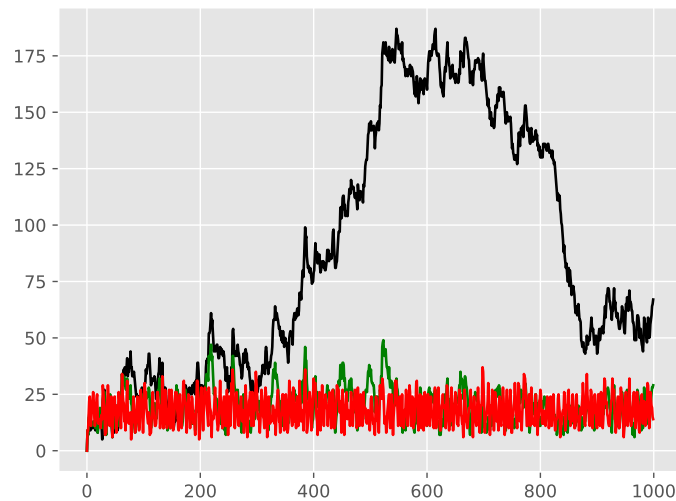
```

19 a = np.random.poisson(11.8, 1000)
20 c = 12
21 Q = computeQ(a, c * np.ones_like(a))
22 Qe1 = computeQExtra(a, c, 1)
23 Qe5 = computeQExtra(a, c, 5)
24
25 plt.clf()
26 plt.plot(Q, label="Q", color='black')
27 plt.plot(Qe1, label="Qe1", color='green')
28 plt.plot(Qe5, label="Qe5", color='red')
29 plt.savefig("psychfinal.pdf")

```

Ex 2.4.1. Explain how the if statements in the code above work.

Ex 2.4.2. Explain how this idea relates to what happens in a supermarket if there are still open service desks but queues become very long.



We see, dynamically controlling the service capacity (as a function of queue length) is a much better plan.

Ex 2.4.3. Use simulation to show that the psychiatrists don't have more work.

Ex 2.4.4. Choose some other control thresholds (something reasonable of course, but otherwise you are free to select your own values.) Run the simulation with your values, include a graph and explain what you see.

2.4.2 Hire an extra server for a fixed amount of time

In the real case the psychiatrists hired an extra person to do intakes when the queue became very long, 100 or higher, and then they hired this person for one month (you may assume that a month consists of 4 weeks). Suppose this person can do 2 intakes a day and works for 4 days a week.

The code below implements this control algorithm.

Ex 2.4.5. Explain the code below.

Python Code

```

1  import numpy as np
2  import matplotlib.pyplot as plt
3  from matplotlib import style
4
5  style.use('ggplot')
6  np.random.seed(3)
7
8  extra_capacity = 8 # extra weekly capacity
9  contract_duration = 4 # weeks
10
11
12 def compute_Q_control(a, c, Q0=0):
13     N = len(a)
14     Q = np.empty(N)
15     Q[0] = Q0
16     extra = False
17     mark_time = 0
18     for n in range(1, N):
19         if Q[n - 1] > 100:
20             extra = True
21             mark_time = n
22         if extra and n >= mark_time + contract_duration:
23             extra = False
24         d = min(Q[n - 1], c[n] + extra * extra_capacity)
25         Q[n] = Q[n - 1] + a[n] - d
26     return Q
27
28
29 a = np.random.poisson(11.8, 1000)
30 c = 12
31 Q = compute_Q_control(a, c * np.ones_like(a), Q0=110)
32 # print(Q)
33 plt.clf()
34 plt.plot(Q, label="Q", color='black')
35 plt.savefig("psych-extra.pdf")

```

Ex 2.4.6. Do a number of experiments to see the effect of the duration of the contract by making it longer (experiment 1), or shorter (experiment 2). Run the simulation, Include graphs, and discuss the effect of these changes.

Ex 2.4.7. Now change the number of intakes per day done by the extra person. (For instance, an experienced person can do more intakes in the same amount of time than a newbie. However, this comes at an additional cost of course.) Make a graph, and compare the effect of this change to the previous (changing the duration).

Ex 2.4.8. If you were a consultant, what would you advice the psychiatrists on how to control their waiting lists?

SIMULATION OF QUEUEING PROCESSES IN CONTINUOUS TIME

3.1 INTRO

I discuss two elegant algorithms to simulate the waiting time process. One is for a system with one server, and jobs are served in the order in which they arrive. The second is for a multi-server FIFO queue.

- Youtube: <https://youtu.be/h10TvdLs9ik>

3.2 COMPUTING WAITING TIMES

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

3.2.1 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 np.random.seed(3)

```

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 λ 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 3.2.1. Explain that `scale=labda` sets the interarrival times to 3, but that in our queueing models, λ 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)

```

3.2.3 Arrival times

Ex 3.2.2. Why do we generate first random inter-arrival times, and use these to compute the arrival times? Why not directly generate random arrival times?

Python Code

```

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

```

Check the output 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 X[0] = 0
2 A = X.cumsum()
3 print(A)

```

This is better!

Ex 3.2.3. Why can we remove `X[0]` without fundamentally changing the analysis?

3.2.4 Service times

We have arrival times. We next need the service times of the jobs. Assume they are $\sim \text{Exp}(\mu)$ with μ somewhat larger than λ . (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 3.2.4. Why do I use `size=len(A)` in the definition of `S`?

Ex 3.2.5. Why do we set $\text{scale}=1/\mu$?

Ex 3.2.6. It's easy to compute the mean service time like this

```
1 print(S.mean())
```

Python Code

Explain that this is not exactly equal to $E[S]$.

3.2.5 *Departure times*

The standard recursion to compute the departure times.

```
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)
```

Python Code

Ex 3.2.7. Explain now why it is practical to have $A_0 = 0$.

3.2.6 *Sojourn times*

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

```
1 J = D - A
2 print(J)
```

Python Code

3.2.7 *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?

```
1 W = J - S
2 print(W)
```

Python Code

Ex 3.2.8. 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?

3.3 KPIS AND PLOTTING

3.3.1 *Relevant averages*

```
1 print(W.mean(), W.std())
```

Python Code

Python Code

```
1 print(J.mean(), J.std())
```

Python Code

```
1 plt.clf()
2 plt.plot(J)
3 plt.savefig("sojourn.pdf")
```

Ex 3.3.1. 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.

3.3.2 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 3.3.2. 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 ρ ?
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 3.3.3. 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?

3.3.3 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 3.3.4. 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  X[0] = 0
10 A = 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 3.3.5. 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 3.3.6. Explain these lines. About the last line, explain why this acts as a check.

3.3.4 Virtual waiting time

Plotting the virtual waiting time is subtle. (The code below is short, hence may seem to be easy to find, but for me it wasn't. To get it right took me two to three hours, and I also discovered other bugs I made elsewhere. Coding is hard!)

Ex 3.3.7. Make a plot of the virtual waiting time by hand. Figure out which points are the most important ones to characterize the virtual waiting times, and explain why this is so.

Here is the code to plot the virtual waiting time.

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 W = np.zeros_like(A)
19 for k in range(1, len(X)):
20     W[k] = max(W[k - 1] + S[k - 1] - X[k], 0)
21     D[k] = A[k] + W[k] + S[k]
22
23 idx = np.where(W <= 0)[0] # this
24 empty = D[idx[1:] - 1]
25
26 E = np.zeros((2 * len(A) + len(empty), 2)) # epochs
27 E[: len(A), 0] = A
28 E[: len(A), 1] = W
29 E[len(A) : 2 * len(A), 0] = A
30 E[len(A) : 2 * len(A), 1] = W + S
31 E[2 * len(A) : 2 * len(A) + len(empty), 0] = empty
32 E[2 * len(A) : 2 * len(A) + len(empty), 1] = 0
33 E = E[np.lexsort((E[:, 1], E[:, 0]))]
34
35 plt.plot(E[:, 0], E[:, 1])
36 plt.savefig("virtual-waiting-time.pdf")

```

The this line is perhaps somewhat strange. By printing the result, we can find out that `np.where(W <= 0)` returns a tuple of which the first element is an array of the indices where `W` is zero. To get that first element we add the extra `[0]`.

Ex 3.3.8. Explain how we store the relevant epochs in `E`. Why do we use `idx[1:]` (What is `W[0]`?). Why do we subtract 1 from `idx[1:]`? Why do we use `np.lexsort`? (Check the documentation to see how lexical sorting works. It is important to know what lexical sorting is.)

3.4 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.

Python Code

```

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

Ex 3.4.2. The above procedure to compute the number of jobs in the system is pretty inefficient. Why is that so?

Ex 3.4.3. 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.

3.5 MULTI-SERVER QUEUE

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

3.5.1 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 3.5.1. Explain that we can implement a fast server by specifying the number of servers c , and change the service times as follows:

Python Code

```

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

```

3.5.2 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(1, 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 3.5.2. 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 3.5.3. 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.argmax() # 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 3.5.4. 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 3.5.5. 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 3.5.6. 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.

SAKASEGAWA'S FORMULA AND BATCHING RULES

4.1 INTRO

Here we study how to apply Sakasegawa's formula in various cases with and without batching, and check the quality of this approximation.

4.2 VARIOUS MODELS COMPARED

We have (at least) three different types of models for a queueing system: simulation in discrete time, simulation in continuous, and Sakasegawa's formula to compute the expected waiting time in queue. Let's see how these models compare.

4.2.1 Discrete time simulation

We have a machine that can serve c_k jobs on day k . When the period length is T , then $c_k \sim \text{Pois}(\mu T)$. The number jobs that arrive on k is $a_k \sim \text{Pois}(\lambda T)$. The arrivals in period k cannot be served on day k . The code to simulate this should be familiar to you by now.

Python Code

```

1  import numpy as np
2
3  np.random.seed(3)
4
5  labda = 3
6  mu = 1.2 * labda
7  T = 8 # period length, 8 hours in a day
8
9  num = 1000
10
11 a = np.random.poisson(labda * T, size=num)
12 c = np.random.poisson(mu * T, size=num)
13 L = np.zeros_like(a, dtype=int)
14
15 L[0] = 5
16 for i in range(1, num):
17     d = min(c[i], L[i - 1])
18     L[i] = L[i - 1] + a[i] - d
19
20
21 print(L.mean(), L.std())

```

Ex 4.2.1. Run this code (and write down the results). Then change the code such that the arrivals can be served on the day they arrive. Rerun and compare the results.

Ex 4.2.2. Now change the period time, which was 8 hours in a day, to $T = 1$ (so that we concentrate on what happens during an hour instead of a day).

Rerun the code with and without the arrivals being served on the period of arrival, and compare to the previous exercise. Explain the difference.

Ex 4.2.3. What are the advantages and disadvantages of using small values for T ?

4.2.2 Continuous time simulation

Here is the code.

Python Code

```

1  import numpy as np
2
3  np.random.seed(3)
4
5  labda = 3
6  mu = 1.2 * labda
7  num = 1000
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
15 for k in range(1, len(A)):
16     D[k] = max(D[k - 1], A[k]) + S[k]
17
18 J = D - A
19 EL = labda * J.mean()
20 print(EL)

```

Ex 4.2.4. Which result did we use to compute $E[L]$? Can we use that to estimate $V[L]$?

Ex 4.2.5. Run the code and compare with the discrete time simulation.

4.2.3 Sakasegawa's formula

Here we use Sakasegawa's formula to compute $E[L]$.

Python Code

```

1  import numpy as np
2
3
4  def sakasegawa(labda, ES, Ca2=1, Cs2=1, c=1):
5      rho = labda * ES
6      V = (Ca2 + Cs2) / 2
7      U = rho ** (np.sqrt(2 * (c + 1)) - 1) / (1 - rho)
8      T = ES / c
9      return V * U * T
10

```



```

11
12 labda = 3
13 mu = 1.2 * labda
14 ES = 1 / mu
15 W = sakasegawa(labda, 1 / mu, 1, 1, 1)
16 L = labda * (W + ES)
17 print(L)

```

Ex 4.2.6. Explain the code, in particular the values of the parameters.

Ex 4.2.7. Run the code, and compare the value with the discrete and continuous time simulations.

Ex 4.2.8. What are advantages and disadvantages of using Sakasegawa's formula?

4.3 SERVER SETUP AND BATCHING

We can now setup a model in which jobs arrive in batches of size B and in between batches the server needs a constant setup time R . Check the queueing book for further background; we are going to build the model of the related section.

4.3.1 Sakasegawa's formula

I build up the code in small blocks. You should put the code blocks underneath each other as you progress.

Since we add setup times, it's easy to have too large loads.

Python Code

```

1 import numpy as np
2
3 np.random.seed(3)
4
5 B = 13
6 labda = 3
7 mu = 2 * labda
8 R = 2
9
10 rho = labda * (1 / mu + R / B)
11 assert rho < 1, f"rho={rho} >= 1"

```

Ex 4.3.1. What does the assert statement? Why do I put it before doing any other work? What happens if you would set $\text{labda}=10$ and $\text{mu}=3$.

Next, we need to get the parameters correct for the batches. I just follow the book.

Python Code

```

1 W1 = (B - 1) / 2 / labda

```

Ex 4.3.2. What is $W1$ conceptually:?

For the queueing time, we have this:

Python Code

```

1 labdaB = labda / B
2 VR = 0 # constant R
3 ca2 = 1 / B
4 ES0 = 1 / mu
5 VS0 = ES0 * ES0
6 VSB = VR + B * VS0
7 ES = ES0 + R / B
8 ESB = R + B * ES0
9 cs2 = VSB / ESB ** 2

```

Ex 4.3.3. We set $ca2 = 1$. What is the assumption about the distribution of the interarrival times X ?

Ex 4.3.4. Print $cs2$. Why is that smaller than 1? `#+end_src`

For the average queueing time:

Python Code

```

1 def sakasegawa(labda, ES, Ca2, Cs2, c=1):
2     rho = labda * ES
3     V = (Ca2 + Cs2) / 2
4     U = rho ** (np.sqrt(2 * (c + 1)) - 1) / (1 - rho)
5     T = ES / c
6     return V * U * T
7
8 W2 = sakasegawa(labdaB, ESB, ca2, cs2)

```

Ex 4.3.5. In view of the previous exercise, explain that

Python Code

```

1 sakasegawa(labda, ES, 1, 1)

```

gives the wrong result, but

Python Code

```

1 sakasegawa(labda, ES, 1, cs2)

```

agrees with

Python Code

```

1 sakasegawa(labdaB, ESB, ca2, cs2)

```

The last step:

Python Code

```

1 W3 = R + (B + 1) / 2 * ES0

```

Ex 4.3.6. What is the meaning of W_3 ?

The sojourn time:

Python Code

```

1 J = W1 + W2 + W3
2 print(J)

```

4.3.2 Simulation

To set up the simulation requires a bit fiddling with slicing. It took me a bit of time, and print statements, to get the details right. Here is the code, with the print statements so that you can figure out how it works.

Python Code

```

1 # Don't forget the copy the parameters like B so that you work with
2 # the same numbers.
3
4 num = 10
5 num = B * (num // B) # get multiple of B
6 X = np.ones(num)
7 X[0] = 0
8 A = X.cumsum()
9 A = np.zeros_like(A0)
10 for i in range(num // B):
11     st = i * B + 1 # start
12     fi = (i + 1) * B # finish
13     print(st, fi, A0[fi])
14     A[st : fi + 1] = A0[fi]
15     print(A)
16
17 S0 = np.ones_like(A0)
18 S0[0] = 0
19 S = S0.copy()
20 S[1::B] += R
21 print(S)

```

Ex 4.3.7. Use the print statements to explain how the slicing, i.e., notation like $A[\text{st} : \text{fi} + 1]$, works. Explain how A and S correspond to batch arrivals and services with setup times.

Now that we know how to construct batch arrivals and job service times that include regular setups, the rest of the simulation is standard.

Python Code

```

1 # put the parameters here, or glue this code after the code for
2 # Sakasegawa's formula
3
4 num = 1000
5 num = B * (num // B) # get multiple of B
6 X = np.random.exponential(scale=1 / labda, size=num)
7 A0 = np.zeros(len(X) + 1)
8 A0[1:] = X.cumsum()
9 A = np.zeros_like(A0)
10 for i in range(num // B):
11     st = i * B + 1 # start

```

```

12     fi = (i + 1) * B # finish
13     A[st : fi + 1] = A0[fi]
14
15
16     S0 = np.random.exponential(scale=1 / mu, size=len(A))
17     S0[0] = 0
18     S = S0.copy()
19     S[1::B] += R
20
21     D = np.zeros_like(A)
22     for k in range(1, len(A)):
23         D[k] = max(D[k - 1], A[k]) + S[k]
24
25     J = D - A0
26     print(J.mean())

```

Ex 4.3.8. Why is $D - A$ not the sojourn time?

Ex 4.3.9. Run the code for $\text{num} = 1000$ and compare the results of the formulas and the simulation. (Ensure that both models use the same data.) Then extend to $\text{num} = 1_000_000$ and check again. What do you see, and conclude?

4.3.3 Getting things right

While making the above code, I made several (tens of) errors, so that the simulation and the formulas gave different results. Here are the steps that I followed to get things right. Only after one step was correct, I moved on to the next.

1. Check with $B = 1$ and $R = 0$, since $B = 1$ is the single job case.
2. Keep $B = 1$, set $R = 0.2$. I had to change μ so that still $\rho < 1$.
3. Set $B = 2$, $R = 0$. Compare ES (input for Sakasegawa's formula) to $S.\text{mean}()$ (input for simulation).
4. In the previous step I did not get corresponding results for $\text{num} = 10000$. Changing it to 1 million helped.

After these four steps, the simulation and the model gave similar results. However, from a higher level of abstraction, I am not quite happy about this. It is not realistic to wait until we have seen a million or so arrivals in any practical setting. My personal way to deal with this situation is like this (but not all people agree on what to do though):

- Simple formulas are tremendously useful to get *insight* into the main drivers of the behavior of a system. In other words, there is not better way to get *qualitative* understanding than with simple formulas.
- The quantitative quality of a formula need to not be too good.
- Building a simulator is intellectually rewarding as it helps understand the *dynamics* of a system.

- Building a simulator is difficult; it's easy to make mistakes, in the code, in the model, in the data...
- Simulation depends on large quantities of data. It's very hard (next to impossible) to *understand* the output.
- The simple formulas can be used to check the output of a simulator when applied to simple cases.

All in all, I think that simulation and theoretical models should go hand in hand, as they offer different type of insight, and have different strengths and weaknesses.

4.3.4 Random setup times

4.4 SERVER ADJUSTMENTS

Consider now a queueing system in which the server needs an adjustments with probability p (see the section on server adjustments in the book). The repair times are assumed constant, at first, with mean 2. Here is the simulator.

4.4.1 Check work

First we should check the formulas for $E[S]$ and $V[S]$.

Python Code

```

1  import numpy as np
2
3  np.random.seed(3)
4
5  labda = 3
6  mu = 2 * labda
7  ES0 = 1 / mu
8  p = 1 / 20
9  num = 10000
10
11
12  S0 = np.random.exponential(scale=ES0, size=num)
13  R = 2 * np.ones_like(S0)
14  I = np.random.uniform(size=len(R)) <= p
15  S = S0 + R * I
16  print(S.mean(), S.var())
17
18  ER = R.mean()
19  ES = ES0 + p * ER
20  VS0 = ES0 * ES0
21  VR = R.var()
22  VS = VS0 + p * VR + p * (1 - p) * ER * ER
23  print(ES, VS)

```

Ex 4.4.1. Explain what is I . Then explain the line $S0 + R * I$ works.

Ex 4.4.2. Run the code; include the numbers in your assignment. Are the numbers nearly the same?

4.4.2 The simulations

Now that we checked the formulas to compute $E[S]$ and $V[S]$, we can see how well Sakasegawa's formula applies for a queueing system in which a server requires regular, but random, adjustments.

Python Code

```

1  import numpy as np
2
3  np.random.seed(3)
4
5  labda = 3
6  mu = 2 * labda
7  p = 1 / 20
8  num = 10000
9
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 R = 2 * np.ones(len(S)) # this
15 I = np.random.uniform(size=len(S)) <= p
16 D = np.zeros_like(S)
17
18 for k in range(1, len(A)):
19     D[k] = max(D[k - 1], A[k]) + S[k] + R[k] * I[k]
20
21 W = D - A - S
22 print(W.mean())

```

Ex 4.4.3. Explain how D is computed.

To see how the approximation works, glue the next code below the previous code.

Python Code

```

1  def sakasegawa(labda, ES, Ca2=1, Cs2=1, c=1):
2      rho = labda * ES
3      V = (Ca2 + Cs2) / 2
4      U = rho ** (np.sqrt(2 * (c + 1)) - 1) / (1 - rho)
5      T = ES / c
6      return V * U * T
7
8
9  ES0 = 1 / mu
10 VS0 = ES0 * ES0
11 ER = R.mean()
12 ES = ES0 + p * ER
13 rho = labda * ES
14 assert rho < 1, "rho >= 1"
15 VR = R.var()
16 VS = VS0 + p * VR + p * (1 - p) * ER * ER
17 Cs2 = VS / ES / ES
18 W = sakasegawa(labda, ES, 1, Cs2, 1)
19 print(W)

```

Ex 4.4.4. To test the code I set at first $R = 0 * \text{np.ones}(\text{len}(A))$ in the line marked as this. Why is this a good test?

Ex 4.4.5. Now run the code, with R as in the code (not set as 0 such as in the previous exercise). Compare the answers. Then set $\text{num} = 100000$, i.e., 10 times as large. What is the effect?

Ex 4.4.6. Now set $R = \text{np.random.exponential}(\text{scale}=2, \text{size}=\text{len}(A))$. What is the effect on $E[W]$? In general, do you see that indeed $E[W]$ increases with the variability of the adjustments?

Ex 4.4.7. What is the model behind this code?

Python Code

```

1  import numpy as np
2
3  np.random.seed(3)
4
5  labda = 3
6  mu = 4
7  N = 1000
8
9  X = np.random.exponential(scale=1 / labda, size=N)
10 A = np.zeros(len(X) + 1)
11 A[1:] = X.cumsum()
12 S = np.random.exponential(scale=1 / mu, size=len(A))
13 R = np.random.uniform(0, 0.1, size=len(A))
14
15 D = np.zeros_like(A)
16 for k in range(1, len(A)):
17     D[k] = max(D[k - 1], A[k]) + S[k] + R[k]
18
19 W = D - A - S
20 print(W.mean(), W.std())

```

Ex 4.4.8. In stead of

Python Code

```

1  for k in range(1, len(A)):
2      D[k] = max(D[k - 1], A[k]) + S[k] + R[k] * I[k]

```

we could write

Python Code

```

1  for k in range(1, len(A)):
2      D[k] = max(D[k - 1] + R[k] * I[k], A[k]) + S[k]

```

What modeling choice would this change reflect? Which of these two models makes the sojourn smaller?

4.5 SERVER FAILURES

This time we focus on a server that can fail; again check the queueing book for the formulas. Here we just implement them.

4.5.1 *Check work*

Again, first we need to check that our (implementation of the) formulas for $E[S]$ and $V[S]$ are correct.

Python Code

```

1  import numpy as np
2  from scipy.stats import expon
3
4  np.random.seed(3)
5
6  labda = 3
7  mu = 2 * labda
8  ES0 = 1 / mu
9  labda_f = 2
10 ER = 0.5
11 num = 10000
12
13 S0 = np.random.exponential(scale=ES0, size=num + 1)
14 N = np.random.poisson(labda_f * ES0, len(S0))
15 R = expon(scale=ER)
16 S = np.zeros_like(S0)
17 for i in range(len(S0)):
18     S[i] = S0[i] + R.rvs(N[i]).sum()
19
20 A = 1 / (1 + labda_f * ER)
21 ES = ES0 / A
22 print(ES, S.mean(), S0.mean() + N.mean() * R.mean())
23
24 C02 = 1
25 Cs2 = C02 + 2 * A * (1 - A) * ER / ES
26 print(Cs2, S.var() / (S.mean() ** 2))

```

Ex 4.5.1. Run this code, and check the result. Then change num to 100000 to see that the estimate improves.

Ex 4.5.2. Explain how we compute $S[i]$.

4.5.2 *The simulations*

Here is all the code to compare the results of a simulation to Sakasagawa's formula.

Python Code

```

1  import numpy as np
2  from scipy.stats import expon
3
4  np.random.seed(3)
5

```



```

6  labda = 2
7  mu = 6
8  ES0 = 1 / mu
9  labda_f = 2
10 ER = 0.5
11 num = 10000
12
13 S0 = np.random.exponential(scale=ES0, size=num + 1)
14 N = np.random.poisson(labda_f * ES0, len(S0))
15 R = expon(scale=ER)
16 S = np.zeros_like(S0)
17 for i in range(len(S0)):
18     S[i] = S0[i] + R.rvs(N[i]).sum()
19
20
21 X = np.random.exponential(scale=1 / labda, size=num)
22 A = np.zeros(len(X) + 1)
23 A[1:] = X.cumsum()
24 D = np.zeros_like(A)
25 for k in range(1, len(A)):
26     D[k] = max(D[k - 1], A[k]) + S[k]
27
28 W = D - A - S
29 print(W.mean())
30
31
32 def sakasegawa(labda, ES, Ca2=1, Cs2=1, c=1):
33     rho = labda * ES
34     V = (Ca2 + Cs2) / 2
35     U = rho ** (np.sqrt(2 * (c + 1)) - 1) / (1 - rho)
36     T = ES / c
37     return V * U * T
38
39
40 A = 1 / (1 + labda_f * ER)
41 ES = ES0 / A
42 C02 = 1
43 Cs2 = C02 + 2 * A * (1 - A) * ER / ES
44 rho = labda * ES
45 assert rho < 1, "rho >= 1"
46 W = sakasegawa(labda, ES, 1, Cs2, 1)
47 print(W)

```

Ex 4.5.3. Run the code and include the results in your assignment.

Ex 4.5.4. Suppose you can choose between two alternative ways to improve the system. Increase λ_f , and decrease $E[R]$, but such that $\lambda_f E[R]$ remains constant; or the other way around, decrease λ_f and increase $E[R]$. Which alternative has better influence on $E[W]$? (Use Sakasegawa's formula for this; you don't have to do the simulations. In general, computing functions is much faster than simulation.)

4.6 A SIMPLE TANDEM NETWORK

We have two queues in tandem. Again we compare the approximations with simulation.

4.6.1 Simulation in continuous time

This code simulates two queues in tandem in continuous time.

Python Code

```

1  import numpy as np
2
3  np.random.seed(4)
4
5  labda = 3
6  mu1 = 1.2 * labda
7  num = 100000
8  X = np.random.exponential(scale=1 / labda, size=num)
9  A1 = np.zeros(len(X) + 1)
10 A1[1:] = X.cumsum()
11 S1 = np.random.exponential(scale=1 / mu1, size=len(A1))
12 D1 = np.zeros_like(A1)
13
14 for k in range(1, len(A1)):
15     D1[k] = max(D1[k - 1], A1[k]) + S1[k]
16
17 W1 = D1 - A1 - S1
18
19 # queue two
20 A2 = D1
21 mu2 = 1.1 * labda
22 S2 = np.random.exponential(scale=1 / mu2, size=len(A2))
23 D2 = np.zeros_like(A2)
24
25 for k in range(1, len(A2)):
26     D2[k] = max(D2[k - 1], A2[k]) + S2[k]
27
28 W2 = D2 - A2 - S2
29
30 J = D2 - A1
31 print(W1.mean(), S1.mean(), W2.mean(), S2.mean(), J.mean())

```

Ex 4.6.1. Explain how the code works. Why did I choose these parameters?

4.6.2 Sakasegawa's formula plus tandem formula

Now we can check the quality of the approximations.

Python Code

```

1  import numpy as np
2
3
4  def sakasegawa(labda, ES, Ca2=1, Cs2=1, c=1):
5      rho = labda * ES
6      V = (Ca2 + Cs2) / 2

```

```

7      U = rho ** (np.sqrt(2 * (c + 1)) - 1) / (1 - rho)
8      T = ES / c
9      return V * U * T
10
11
12     labda = 3
13     mu1 = 1.2 * labda
14     ES1 = 1 / mu1
15     rho1 = labda * ES1
16     Ca2 = 1
17     Cs2 = 1
18     W1 = sakasegawa(labda, ES1, Ca2, Cs2, 1)
19     Cd2 = rho1 ** 2 * Ca2 + (1 - rho1 ** 2) * Cs2
20
21     Ca2 = Cd2
22     mu2 = 1.1 * labda
23     ES2 = 1 / mu2
24     Cs2 = 1
25     W2 = sakasegawa(labda, ES2, Ca2, Cs2, 1)
26
27
28     J = W1 + ES1 + W2 + ES2
29     print(W1, ES1, W2, ES2, J)

```

Ex 4.6.2. How does the code work? Compare the results obtained from simulation and by formula.

Ex 4.6.3. Assume that the service times at the first queue are constant and equal to $1/\mu_1$. What should be the parameter values for Sakasegawa's formula (explain). Then run both codes and comment on the results.

CONTROL OF QUEUES

5.1 INTRO

We simulate queues that are controlled by some policy that uses information on waiting time or queue length. We also develop algorithms to compute the state probabilities of the $M/G/1$ queue under the policies.

5.2 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?

5.2.1 *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 5.2.1. Run the code and check that indeed long waiting times do occur. What does the this line do?

Ex 5.2.2. Give the Kendall notation for the queueing model that we simulate here.

5.2.2 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 5.2.3. 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 .

5.2.3 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 5.2.4. 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 5.2.5. What does the vector I represent?5.2.4 *Some other blocking rules*

There are other rules to block jobs.

Ex 5.2.6. 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 5.2.7. 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?

5.3 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.

5.3.1 *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 5.3.1. Explain how this code works.

Ex 5.3.2. Run the code. Why do I take B as it is here? Why should $W.mean()$ and $\rho^2/(1 - \rho)$ be approximately equal

5.3.2 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 5.3.3. 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!)

5.4 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.)

5.4.1 *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 $\{A_k\}$.

Python Code

```

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 5.4.1. 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 5.4.2. Just to check that you really understand: why is there no problem with the statement $(L == 0)$?

Ex 5.4.3. Why do I compare $L.mean()$ to $\rho/(1 - \rho)$ and not to $\rho^2/1 - \rho$?

Ex 5.4.4. 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.

5.4.2 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.

Python 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 5.4.5. Read the documentation of how a deque works, then explain the code.

Ex 5.4.6. At first I had these lines:

Python Code

```

1  elif len(Q) < T:
2      Q.append(k)
3      D[k] = D[Q[-1]] + S[k]
```

Why is that wrong?

Ex 5.4.7. What goes wrong with this code:

Python Code

```

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

```

Ex 5.4.8. What queueing discipline would result if we would use the `pop()` and `appendleft()` methods of a deque?

Ex 5.4.9. What queueing discipline would result if we would use the `pop()` and `append()` methods of a deque?

Ex 5.4.10. 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 5.4.11. 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.)

5.5 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.

5.5.1 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(labda 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 5.5.1. Which formulas (give the numbers) of the queueing book have we implemented?

Ex 5.5.2. 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 5.5.3. 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 5.5.4. When the service times are contant, 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(lambda / 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 5.5.5. 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.

5.5.2 Effect of blocking on performance

Ex 5.5.6. 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 5.5.7. Do the same computations for $\mu = 0.5\lambda$. Why is the loss probability not so sensitive to T ?

Ex 5.5.8. 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 ?

5.6 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 5.6.1. 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 5.6.2. Explain how the code works. What is dx , what is dy ? Why do we divide by 2?

Ex 5.6.3. Explain the procedure to find the best T .

QUEUEING THEORY ASSIGNMENT: SIMULATION WITH EVENT STACKS

6.1 GENERAL INFO

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

We will discuss *heap queues* and *classes*. With these we will build a very powerful simulation environment. You should memorize in particular that we implement an *event stack* with a heap queue.

6.1.1 *TODO* 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 240)
```

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

6.1.2 Load standard modules

I need a new modules heapq.

Python Code

```
1 from heapq import heappop, heappush
2 import numpy as np
3 from scipy.stats import expon, uniform
4
5 import matplotlib.pyplot as plt
6 import seaborn as sns; sns.set()
7
8 np.random.seed(3)
```

6.2 SORTING WITH HEAP QUEUES

Here is a simple example of using heaps to sort data

We are going to sort a bunch of students by age. The heap uses the first element of the *tuple* of a student's age and name to sort the students.

Ex 6.2.1. Look up on the web: what is the difference between a tuple and a list? When to use one or the other? (As always, keep your answer brief.)

Let's put a number of students on the stack.

Python Code

```
1 stack = []
2
3 heappush(stack, (25, "Cynthia"))
4 heappush(stack, (21, "James"))
5 heappush(stack, (21, "James"))
6 heappush(stack, (25, "Cynthia"))
```

```

7
8 print(stack)

```

Python Code

```

1 age, name = heappop(stack)
2 print(stack)

```

Python Code

```

1 age, name = heappop(stack)
2 print(stack)

```

Python Code

```

1 heappush(stack, (20, "Pete"))
2 heappush(stack, (18, "Clair"))
3 heappush(stack, (14, "Jim"))
4
5 print(stack)

```

With popping, we take an entry from the top of the stack.

Python Code

```

1 age, name = heappop(stack)
2 print(stack)

```

Python Code

```

1 heappush(stack, (15, "John"))
2 print(stack)

```

Ex 6.2.2. Add an extra attribute to the tuple, for instance a student's height, and make a stack by which you can sort by height. Include your code.

You must have seen during the youtube video that by putting the same element more than once on the stack resulted in a longer, and unsorted, stack. (This came to a surprise to me!) Here I made a simple mistake. A heap queue is not necessarily sorted. However, when *popping* from the stack, it is guaranteed that the elements are popped in sorted order. In other words, it is irrelevant how the stack is sorted, as long as the elements come out in a sorted way. Let's check this.

Python Code

```

1 stack = []
2
3 heappush(stack, (25, "Cynthia"))
4 heappush(stack, (21, "James"))
5 heappush(stack, (21, "James"))
6 heappush(stack, (25, "Cynthia"))
7
8 print(stack)

```

Apparently, the stack does not appear to be sorted. In the code below I pop the elements of the stack, with `heappop`, and append the popped items to the list `res`. You should know that anytime we pop from the stack, the stack becomes one element shorter. So the code below moves elements from

stack to res. Any pass through the while loop removes an element of the stack, and the loop keeps running until the stack is empty. (To help you think about this, this is an ideal exam question.)

Python Code

```

1 res = []
2
3 while stack:
4     res.append(heapop(stack))
5
6 print(res)

```

Indeed, res *is* sorted.

6.3 CLASSES

In a class we organize information.

Python Code

```

1 class Student:
2     def __init__(self, name, age, phone):
3         self.name = name
4         self.age = age
5         self.phone = phone
6
7     def __repr__(self):
8         return f"{self.name}, {self.age}, {self.phone}"
9

```

Ex 6.3.1. After you have read this section on classes, extend the class such that we can give the student also a surname. Include your code.

Ex 6.3.2. Look up on the web: what does the `repr` method do?

Making an *object* of a *class* is called *instantiation*.

Python Code

```

1 hank = Student("Hank", "21", "Huawei")
2 print(hank)

```

Let's add some more students and put them in a list.

Python Code

```

1 students = [
2     Student("Joseph", "18", "Motorola"),
3     Student("Maria", "21", "Huawei"),
4     Student("Natasha", "20", "Apple"),
5     Student("Chris", "25", "Nexus"),
6 ]
7 print(students)

```

Ex 6.3.3. Make two more students, e.g., take your phones, ages and names. (And if you don't like to spread such details, just lie about your age :-)) Then add these students to the stack. Show your code and the output.

With heaps we can sort the students in any sequence we like. Let's sort them by phone brand.

Python Code

```

1  stack = []
2
3  for s in students:
4      heappush(stack, (s.phone, s))
5
6  res = [] # get a sorted list students in order of their name.
7  while stack:
8      res.append(heapop(stack))
9
10 print(res)
11

```

Note once more that we can provide a *key* as a first element in a tuple, and then insert the entire tuple into a heap. In the rest of the tuple we can put anything we like. Like this, we associate things (here student objects) to keys.

Ex 6.3.4. Change the code so that you can sort the students by age. Show the line(s) of your code to achieve this. Then sort by name.

6.4 SORTING JOBS WITH A HEAP QUEUE

How can use heap queues in the simulation of queueing systems? To see, think of time as a sequence of events in which things happen. Then, in a queueing system, two things can happen: a job can arrive or a job leaves. So, by specifying the arrival times of jobs, and their departure times, and storing these times in a heap queue, we use the heap queue to sort all these times. Then we jump from event to event, and this is nothing but the queueing simulation!. In other words, with a heap queue, we can let the heap queue do all the work of tracking time. You'll see below we can nearly forget about it. Once you get the idea (which takes some time), you'll see how neat is all is.

6.4.1 A decent job class

We store the arrival, service, and departure time as job *attributes*. We also store the queue length at arrival times to gather statistics at the end. We can then compute the job sojourn time by means of *class methods*.

Let's start from scratch again with the code so that you have a simple starting point.

Python Code

```

1  class Job:
2      def __init__(self):
3          self.arrival_time = 0
4          self.service_time = 0
5          self.departure_time = 0
6          self.queue_length_at_arrival = 0
7
8      def sojourn_time(self):

```

```

9         return self.departure_time - self.arrival_time
10
11     def waiting_time(self):
12         return self.sojourn_time() - self.service_time
13
14     def service_start(self):
15         return self.departure_time - self.service_time
16
17     def __repr__(self):
18         return f"{self.arrival_time}, {self.service_time}, {self.service_start()}, \
19                 {self.departure_time}\n"
20
21     def __le__(self, other):
22         # this is necessary to sort jobs when they have the same arrival times.
23         return self.id <= other.id

```

Ex 6.4.1. Explain the waiting time and sojourn functions.

Let's make a few jobs and store them in a heap queue. For later purposes, we have to add labels to indicate what type of event we are dealing with, an arrival or a departure.

Python Code

```

1  ARRIVAL, DEPARTURE = 0, 1
2
3  events = [] # event stack, global
4  num_jobs = 5
5
6  time = 0
7  for i in range(num_jobs):
8      time += 3
9      job = Job()
10     job.arrival_time = time
11     job.service_time = 5
12     heappush(events, (job.arrival_time, job, ARRIVAL))
13
14  while events:
15     time, job, typ = heappop(events)
16     print(job)

```

6.5 A VERY POWERFUL GG1 SERVER SIMULATOR

6.5.1 The class definition

We need two states to indicate whether the server is busy or idle.

Python Code

```

1  IDLE, BUSY = 0, 1

```

Python Code

```

1  class GG1:
2      def __init__(self, F, G, num_jobs):
3          self.F = F # interarrival time distribution
4          self.G = G # service time distribution
5          self.num_jobs = num_jobs

```

```

6         self.queue = []
7         self.served_jobs = [] # assemble statistics
8         self.state = IDLE
9
10        def make_jobs(self):
11            time = 0
12            for i in range(num_jobs):
13                time += self.F.rvs()
14                job = Job()
15                job.arrival_time = time
16                job.service_time = self.G.rvs()
17                heappush(events, (job.arrival_time, job, ARRIVAL))
18
19        def run(self):
20            while events: # not empty
21                time, job, typ = heappop(events)
22
23                if typ == ARRIVAL:
24                    self.handle_arrival(time, job)
25                else:
26                    self.handle_departure(time, job)
27
28        def handle_arrival(self, time, job):
29            job.queue_length_at_arrival = len(self.queue)
30            if self.state == IDLE:
31                self.state = BUSY
32                self.start_service(time, job)
33            else:
34                self.put_job_in_queue(job)
35
36        def start_service(self, time, job):
37            job.departure_time = time + job.service_time
38            heappush(events, (job.departure_time, job, DEPARTURE))
39
40        def put_job_in_queue(self, job):
41            heappush(self.queue, (job.arrival_time, job))
42
43        def handle_departure(self, time, job):
44            if self.queue: # not empty
45                _, next_job = heappop(self.queue)
46                self.start_service(time, next_job)
47            else:
48                self.state = IDLE
49            self.served_jobs.append(job)
50

```

Ex 6.5.1. Explain how the `handle_departure` method works.

Ex 6.5.2. Explain how the `run` method works.

6.5.2 Generating jobs

Python Code

```

1  lambda = 2.0
2  mu = 3.0

```

```

3 rho = labda / mu
4 F = expon(scale=1.0 / labda) # interarrival time distribution
5 G = expon(scale=1.0 / mu) # service time distribution
6 num_jobs = 100
7
8 events = []
9
10 gg1 = GG1(F, G, num_jobs)
11 gg1.make_jobs()

```

Ex 6.5.3. Suppose we would write

```

Python Code
1 gg1 = GG1(G, F, num_jobs)

```

What are then the arrival rate and service rate?

To view the contents of the first couple of events I tried this: `print(events[:5])`, but that failed. After a bit of searching on the web I found the following.

```

Python Code
1 from itertools import islice
2
3 print(list(islice(events, 0, 5)))

```

Ex 6.5.4. Read on the web what `islice` does, and explain it in your own words.

Ex 6.5.5. Explain also why we need to turn the output of `islice` into a `list`. (This requires that you read and think about what a generator is.)

Ex 6.5.6. Explain the contents of the above table; what do the rows represent? Why is the content of column 5 negative?

6.5.3 Run the simulation

Now run the simulation.

```

Python Code
1 gg1.run()

```

Print the first 5 jobs (not all, because when we simulate a 1000 or so jobs, the document explodes).

```

Python Code
1 print(gg1.served_jobs[:5])

```

Ex 6.5.7. Explain the first two lines. Are the results in line with what you expect how a $G/G/1$ FIFO queue should behave? If so, why (not)?

6.5.4 Analyze the results, statistics

Here is a plot.

Python Code

```
1 plt.clf()
2 plt.plot(sojourn)
3 plt.savefig('sojourn0.png')
4 'sojourn0.png'
```

Ex 6.5.8. Change the seed of the random number generator, choose your favorite number of jobs (something positive, reasonably small). make your own plot and include it in your document.

And some statistics. For instance the sojourn times.

Python Code

```
1 sojourn = np.zeros(len(gg1.served_jobs))
2 for i, job in enumerate(gg1.served_jobs):
3     sojourn[i] = job.sojourn_time()
4
5 print(sojourn.mean(), sojourn.std(), sojourn.max())
```

Python Code

```
1 wait = sojourn.mean() - G.mean()
2 print(wait)
```

6.5.5 Comparison with Sakasegawa's formula

We can compare the results of our simulation with the theoretical values of the stationary state.

Python Code

```
1 def sakasegawa(F, G, c):
2     labda = 1.0 / F.mean()
3     ES = G.mean()
4     rho = labda * ES / c
5     EWQ_1 = rho ** (np.sqrt(2 * (c + 1)) - 1) / (c * (1 - rho)) * ES
6     ca2 = F.var() * labda * labda
7     ce2 = G.var() / ES / ES
8     return (ca2 + ce2) / 2 * EWQ_1
9
10
11 average_wait = sakasegawa(F, G, c=1)
12 print(average_wait, average_wait + G.mean())
```

Ex 6.5.9. Compare the result of Sakasegawa's formula and the results of your simulation (with about 100 jobs). You should remark that the difference in the waiting time is quite large. Do we simulate the $M/M/1$ queue here? Is Sakasegawa's result is exact for the $M/M/1$ queue?

Ex 6.5.10. Run an example with more jobs, 1000 or so. Is now the average sojourn time in the simulation nearer to the theoretical result?

Ex 6.5.11. Another idea is this. We start with an empty system, and that is a very particular state. Suppose we would start with 10 jobs in queue. Do a simulation for 100 jobs, but with 10 jobs in queue at the start. What is then the average sojourn time? Hint, below I include some code to put 10 jobs in queue right at the start. There is subtlety. At first I set the arrival time to 0 for all these jobs, but then the heapq algorithm doesn't know how to sort the jobs. To repair this problem, I just take very small arrival times.

Python Code

```

1 eps = 0.0001
2 for i in range(10):
3     job = Job()
4     job.arrival_time = i*eps
5     job.service_time = G.rvs()
6     heappush(events, (job.arrival_time, job, ARRIVAL))

```

6.5.6 Further experiments

Ex 6.5.12. Change the service distribution to the uniform distribution $U[a, b]$. Take a and b reasonable, so that the load remains below 1. Run an example with 100 jobs. Make a graph of the waiting times and the sojourn times.

Ex 6.5.13. Change the arrival distribution to a uniform distribution, and take also the service distribution to be uniformly distributed. Take some sensible values for λ and $E[S]$, e.g., $\lambda = 1$ and $E[S] = 0.4$. Run an example with 100 jobs. Make a graph of the waiting times and the sojourn times.

6.6 SPTF SCHEDULING

Suppose we prefer to serve the shortest job in queue; it can be proven that this scheduling rule minimizes the queue length. We can inherit all of our GG1 queue, except the scheduling rule. For this we need to change just one line!

Python Code

```

1 class SPTF_queue(GG1):
2     def put_job_in_queue(self, job):
3         heappush(self.queue, (job.service_time, job))

```

And that's it! Now run it.

Python Code

```

1 sptf = SPTF_queue(F, G, num_jobs)
2 sptf.make_jobs()
3 sptf.run()
4 print(sptf.served_jobs[:5])

```

6.7 LIFO SCHEDULING

Last-in-First-Out is also trivial.

Python Code

```

1 class LIFO_queue(GG1):
2     def put_job_in_queue(self, job):
3         heappush(self.queue, (-job.arrival_time, job))

```

Ex 6.7.1. Run an example with 100 jobs. Make a graph of the waiting times and the sojourn times. Comment on your findings.

6.8 SERVE LONGEST JOB FIRST

Ex 6.8.1. Update the code of the SPTF queue such that the longest job is selected from the queue, rather than the shortest. (Hint, put a minus sign at the right position when adding a job to the queue heap.) Simulate it. Make a graph of the waiting times and the sojourn times. Comment on your findings.

HINTS

h.1.3.3. What is the average number of arrivals per period? What is the average number of jobs that can be served, i.e., the average service capacity? Are they close or not?

h.1.3.8. As a simple analogous problem: imagine you have bucket containing 10 liters. Water flows in from a hose at rate 3 liters per minute, but it flows out via another hose at rate 5 l/m. What is the net outflow? Why does it take 5 minutes before the bucket is empty?

h.1.3.10. See the queueing book exercises 2.1.8 and 3.2.3 for further explanations.

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 τ to hit 0, i.e., the time until the queue is empty, is the smallest τ 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 $\mu = E[X]$, σ the std of X , and N the normal distribution, and by the central limit theorem,

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

h.3.2.3. The X_k are iid.

h.3.2.4. 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.3.2.6. Did we really serve job 0? If num is big number, does it matter that we set $S[0]=0$?

h.3.2.7. 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.3.2.8. We subtract $S[0]$ as if we served the corresponding job, but did we actually serve it?

h.3.3.7. The crucial points are (A_k, W_k) , $(A_k, W_k + S_k)$, and $(D_k, 0)$ when $W_{k-1} = 0$. Then connect these points with straight lines.

h.3.4.3. 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.3.5.5. For instance, set $c = \text{np.array}([2, 0.5, 0.5])$.

h.4.2.1. Replace the relevant line by $d = \min(c[i], L[i - 1] + a[i])$.

h.4.2.4. Recall $L = \lambda W$.

h.4.3.4. When $S = R/B + S_0$, then part of S is constant. Hence, can it be (relatively) as variable as S_0 ?

h.4.3.7. Why do I chose np.ones to fill A_0 and S_0 ? What is the difference between A_0 and A , and S_0 and S ?

h.4.3.8. What do you think I used first?

h.4.4.8. What is the influence on the setup? Do we still require that the setup has to be done immediately before a service starts?

h.4.5.4. For instance, set $\text{labda_f} = 4$ and $\text{ER} = 0.25$.

h.4.6.1. For what situations are the formulas exact?

h.4.6.3. In the simulation, replace the line with $S1$ by $S1 = \text{np.ones}(\text{len}(A1)) / \text{mu1}$. In the formula's, don't forget to update C_s^2 for the relevant queue.

h.5.2.2. It's not the $D/D/1$ queue.

h.5.2.3. 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.5.2.4. 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.5.2.5.** If $I[k] == 1$, then what happens to job k ?
- h.5.3.2.** 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.5.4.1.** 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.5.4.2.** Is L a float?
- h.5.4.3.** What is $\rho^2/(1 - \rho)$?
- h.5.4.6.** To what job does $Q[-1]$ point to? It's one to far!
- h.5.4.7.** When Q is empty (i.e., $\text{len}(Q)$ is 0), then what is $Q[-1]$?
- h.5.4.8.** Does it matter whether we push jobs from right to left through a queue, rather than from left to right?
- h.5.4.9.** It's not FIFO.
- h.5.4.10.** Is $L.\text{max}()$ larger than 100 for this simulation?
- h.5.5.6.** Why is the loss probability equal to π_T ?