

Queueing theory Assignments

EBBo74Ao5

Nicky D. van Foreest

2022:03:14

CONTENTS

1	Simulation of queues in discrete time	3
2	Queueing control: Psychiatrists doing intakes	13
3	Simulation of Queueing processes in continuous time	20
4	Sakasegawa's formula and Batching rules	30
5	Control of Queues	38
6	Simulation with event stacks	45
7	Hints	53

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 \LaTeX 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 1.1. 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. Run the code for 10 periods and report your result.

Ex 1.3. 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 1.4. 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 1.5. 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('figures/queue-discrete_1.pdf')

```

Ex 1.6. 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('figures/queue-discrete-1-1.pdf')

```

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

Ex 1.7. 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 1.8. 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 1.9. 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.)

Include the graph in your report and explain the differences.

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('figures/queue-discrete-1-2.pdf')

```

1.3.6 Serve arrivals in the same period as they arrive

Ex 1.10. 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('figures/queue-discrete_2.pdf')
```

Ex 1.11. Include the graph in your report. What is $c.mean() - a.mean()$? Explain the slope of the line (fitted through the points.) (Check the hint.)

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('figures/queue-discrete_3.pdf')
```

Ex 1.12. Explain the slope of the lines (if you would fit that through the points.)

1.3.9 Hitting times

When $\rho < 1$ and L_0 is some large number we could be interested in estimating the time until the queue is empty. With this we can decide if extra capacity hired to remove the long queue (for instance, long waiting times in a hospital) suffices. If the time to hit zero is still too long, we should hire more capacity.

Here we study how to estimate the hitting time $\tau = \min\{k : Z_k \leq 0\}$, where $Z_k = Z_{k-1} + a_k - c_k$. (Observe that Z and L are not the same everywhere: Z can become negative, the number L in the system L is always ≥ 0 .)

Ex 1.13. Run the following code and include the figure in your report. Use the CTL (see the hint if you forgot) to explain that most sample paths of Z seem to hit 0 when

$$\tau \approx \frac{100}{(5+9)/1.9 - 7}.$$

Next, change the factor 1.9 to 1.2. Why do the graphs lie much nearer to each other?

Python Code

```

1  import numpy as np
2  import matplotlib.pyplot as plt
3
4  num = 500
5  L0 = 100
6
7
8  def hitting_time():
9      a = np.random.randint(5, 10, size=num)
10     c = (5 + 9) / 1.9 * np.ones(num)
11     a[0] = L0
12     Z = (a - c).cumsum()
13     plt.plot(Z)
14     return
15
16
17 samples = 30
18 for n in range(samples):
19     hitting_time()
20
21 plt.savefig("figures/free-random-walk.pdf")

```

We see in the figure of the previous exercise that there is considerable variation in the time the random walk hits zero—a bit more specifically, hits the set $\{\dots, -2, -1, 0\}$ —when ρ is not very small. We need some code to compute τ for a sample path of Z .

Ex 1.14. Explain how the loop over i in the function `hitting_time` computes τ for a specific sample path. You don't have to run the code, the aim is that you understand how the algorithm works.

Python Code

```

1  import numpy as np
2
3  np.random.seed(3)
4
5  num = 500
6  L0 = 100

```

```

7
8
9 def hitting_time():
10     a = np.random.randint(5, 10, size=num)
11     c = (5 + 9) / 1.9 * np.ones(num)
12     L = L0
13     for i in range(1, num):
14         L += +a[i] - c[i]
15         if L <= 0:
16             return i
17
18
19 samples = 10
20 tau = np.zeros(samples)
21 for n in range(samples):
22     tau[n] = hitting_time()
23
24 print(tau.mean(), tau.std())

```

Once again, python (and R) are rather slow in comparison to C or fortran; the factor in speed can be 100 or more in the execution of for loops. For this reason I prefer to tinker a bit with code to push as much of the computation to numpy.

Ex 1.15. Run this code, and explain the output of each print statement of this piece of code. In particular, realize that for `np.argmax`: ‘In case of multiple occurrences of the maximum values, the indices corresponding to the first occurrence are returned.’ (I found this explanation in the numpy documentation.)

<pre> 1 import numpy as np 2 3 np.random.seed(3) 4 5 num = 10 6 L0 = 10 7 samples = 3 8 9 10 dims = (samples, num) 11 a = np.random.randint(5, 10, size=dims) 12 c = 8 * np.ones(dims) 13 a[:, 0] = L0 14 Z = (a - c).cumsum(axis=1) 15 print(Z) 16 print(Z <= 0) 17 print(np.argmax(Z <= 0, axis=1)) </pre>	Python Code	
--	-------------	--

BTW, observe that I use small samples to print the output, so that it’s easy to see how everything works.

Here is the final code.

<pre> 1 import numpy as np 2 from scipy.stats import norm </pre>	Python Code	
--	-------------	--

```

3  import matplotlib.pyplot as plt
4
5  np.random.seed(3)
6
7  num = 400
8  L0 = 100
9  samples = 3000
10
11
12  dims = (samples, num)
13  a = np.random.randint(5, 10, size=dims)
14  c = (5 + 9) / 1.9 * np.ones(dims)
15  a[:, 0] = L0
16  Z = (a - c).cumsum(axis=1)
17  tau = np.argmax(Z <= 0, axis=1)
18  print(tau.mean(), tau.std())
19
20  tau_scaled = (tau - tau.mean()) / tau.std()
21  print(tau_scaled.mean(), tau_scaled.std())
22  bins = np.linspace(-3, 3, 40)
23
24  B = (bins[1:] + bins[:-1]) / 2
25
26  plt.hist(tau_scaled, bins=bins, density=True)
27  plt.plot(B, norm.pdf(B))
28  plt.savefig("figures/tau.pdf")

```

Ex 1.16. Use the CLT to provide some intuition to see why `tau_scaled` is approximately normally distributed. Include the plot in your report.

Here are some specific points of interest in the code:

1. `num=400` just to guess to ensure that $Z[400] < 0$ for all sample paths. This trick allowed me to use numpy. Otherwise I would have needed a for loop in python, which I wanted to avoid.
2. `bins` contains the edges of the bins. To plot the pdf of the standard normal distribution, I need the mid points of the bins. This explains `B`.

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’, roughly speaking.
3. If we start with a huge queue, but the service capacity is larger than the arrival rate, then the queue will drain like a straight line (roughly).

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

Ex 1.18. Why is, in this case, dtype=float not necessary in the definition of L?

Ex 1.19. Add the following code to make the graph of the (simulated) queue length process.

Python Code

```

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

```

Include the graph in your report. 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 1.20. Explain the line in the code that changed.

Ex 1.21. 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('figures/queue-discrete-loss-p.pdf')
```

Ex 1.22. Why should p have a length of $K+1$? Then explain how the code estimates p_i .

Ex 1.23. 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 1.24. 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 1.25. Now change the parameters such $\rho \approx 10$, use a simple argument to show that π_{K-2} should be really small. 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 then 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

```

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

```

2.2.2 Simulate queue length

Python Code

```

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

```

Ex 2.1. Compute the queue lengths when

Python Code

```

1 a = [1,2,5,6,8,3,7,3]
2 c = [2,2,0,5,4,4,3,2]

```

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, 12)
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. 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 2.3. 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 2.4. 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("figures/psych1.pdf")

```

Ex 2.5. Choose your own seed, run the code, include the figure in your report and comment on what you see.

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.

```

#+begin_src python
def balanced(a):
    p = np.empty([5, len(a)])
    p[0, :] = 2.0 * np.ones_like(a)
    p[1, :] = 2.0 * np.ones_like(a)
    p[2, :] = 3.0 * np.ones_like(a)
    p[3, :] = 4.0 * np.ones_like(a)
    p[4, :] = 4.0 * np.ones_like(a)
    return p

np.random.seed(3)
a = np.random.poisson(11.8, 1000)
p = balanced(a)
spread_holidays(p)
s = np.sum(p, axis=0)
Q2 = computeQ(a, s)
plt.plot(Q2)
plt.savefig("figures/psych2.pdf")
#+end_exercise

```

Ex 2.6. How can we see that the effect of balancing capacity is totally uninteresting?

Ex 2.7. Change the capacities of the psychiatrists but keep the average weekly capacity the same. Include a graph of your result, and comment on the effect of your changes.

2.3.2 Synchronize holidays

Supposwe we would synchronize the holidays so that all psychiatrists take holiday in the same week. Would that have an effect on the queue process?

Python Code

```

1 a = np.random.poisson(11.8, 12)
2
3
4 def synchronize_holidays(p):
5     for j in range(len(a) // 5 + 1):
6         p[:, 5 * j] = 0 # this
7     return p
8
9 p = unbalanced(a)
10 p = synchronize_holidays(p)
11 print(p)

```

Ex 2.8. Explain how the code works. Specifically, what does the line marked as this?

Let's do a longer run to see the effect.

Ex 2.9. In the code below, choose your own seed, run it, include the figure in your report and comment on what you see.

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("figures/psych3.pdf")

```

Ex 2.10. 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.

Ex 2.11. 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

```

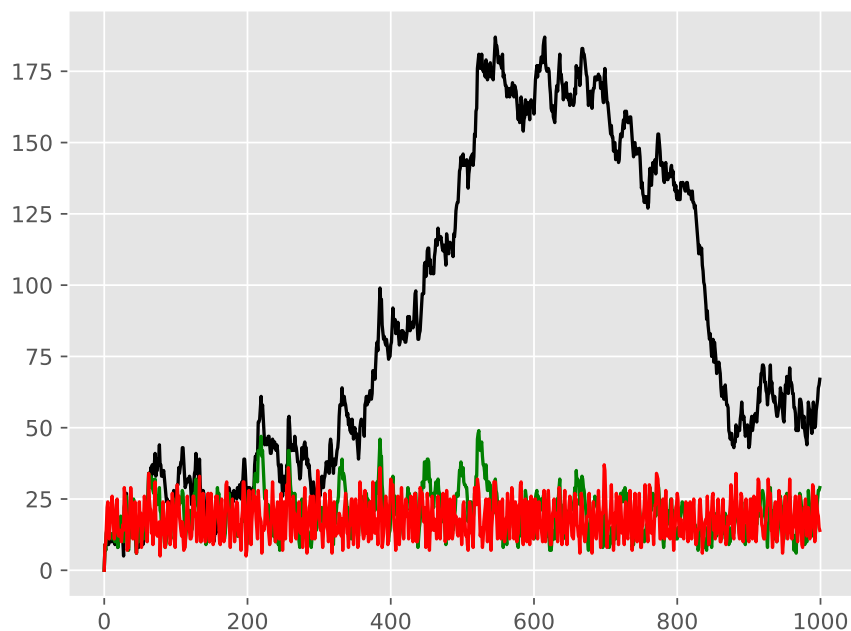
```

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("figures/psychfinal.pdf")

```

Ex 2.12. Explain how the if statements in the code above work.

Ex 2.13. 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 2.14. Use simulation to show that the psychiatrists don't have more work.

Ex 2.15. 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 2.16. 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("figures/psych_extra.pdf")

```

Ex 2.17. 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 2.18. 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 2.19. If you were a consultant, what would you advise 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	<code>import numpy as np</code>
2	<code>import matplotlib.pyplot as plt</code>
3	<code>from matplotlib import style</code>
4	
5	<code>np.random.seed(3)</code>

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	<code>num = 3</code>
2	<code>labda = 3</code>
3	<code>X = np.random.exponential(scale=labda, size=num)</code>
4	<code>print(X)</code>

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

	Python Code
1	<code>num = 100</code>
2	<code>labda = 3</code>
3	<code>X = np.random.exponential(scale=labda, size=num)</code>
4	<code>print(X.mean())</code>

Ex 3.1. 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. 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)

```

- Ex 3.3.**
1. Why is this better?
 2. 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 3.4. Why do I use `size=len(A)` in the definition of S ?

Ex 3.5. Why do we set `scale=1/mu`?

Ex 3.6. It's easy to compute the mean service time like this

Python Code

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

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

3.2.5 Departure times

The standard recursion to compute the departure times.

Python Code

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

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 ?

Python Code

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

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?

Python Code

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

Ex 3.7. 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 $W.mean()$ be negative?

3.3 KPIS AND PLOTTING

The next step is to see how to compute the most important indicators to assess the performance of the system. We can use these, so-called, Key Performance Indicators (KPIs) to see whether we should add, for instance, service capacity.

3.3.1 Relevant averages

Python Code

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

Python Code

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

Ex 3.8. The mean of W is not entirely correct in the way we compute it here. What is (just a bit) wrong?

Python Code

```

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

```

Ex 3.9. Change the simulation length to 1000 jobs. Do one run for $\mu = 3.5$ and another for 2.8. Compute the KPIS, make a plot, and include that in your assignment. Comment on what you see.

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 3.10. Explain the code above. Some specific points:

1. Why is `S.sum()` the total busy time of the server?
2. Why do we divide by `D[-1]` in the computation of 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 3.11. Run this code for a simulation with 10 or so jobs (some other small number). Explain how this code works. Which line is a check on the computations?

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

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

Python Code

```

1 import numpy as np
2
3 np.set_printoptions(suppress=True)
4 np.random.seed(3)
5
6 num = 10
7 labda = 3
8 X = np.random.exponential(scale=1 / labda, size=num)
9 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 3.13. What is stored in `idx`? Why do we strip `A[0]`? Why do we subtract `D[idx-1]` and not `D[idx]`? (Print out the variables to understand what they mean, e.g., `print(idx)`.)

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

Python Code

```

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

```

Ex 3.14. 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 3.15. Make a plot of the virtual waiting time by hand (you don't have a make a large plot, just show that you understand what the virtual waiting process looks like). Find 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("figures/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 3.16.**
1. Explain how we store the relevant epochs in `E`.
 2. Why do we use `idx[1:]` (What is `W[0]`)?
 3. Why do we subtract 1 from `idx[1:]`?
 4. 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 THE NUMBER OF JOBS IN THE SYSTEM

We have the waiting times, but not the number of jobs in the system (queue). Here we show how to plot L , i.e., the number of jobs in the system as seen by a job *upon* arrival.

A simple, but *inefficient*, algorithm to construct L is the following.

Python Code

```

1  import numpy as np
2

```

```

3  np.random.seed(3)
4
5  num = 10
6  labda = 1
7  mu = 1.1 * labda
8
9  X = np.random.exponential(scale=1 / labda, size=num)
10 X[0] = 0
11 A = np.zeros_like(X)
12 A = X.cumsum()
13 S = np.random.exponential(scale=1 / mu, size=len(A))
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_like(A)
20 for k in range(1, len(A)):
21     idx = 0
22     while D[idx] < A[k]:
23         idx += 1
24     L[k] = k - idx
25
26 print(L)

```

- Ex 3.17.**
1. Explain how this algorithm works.
 2. Why does this algorithm find the number of jobs as seen *just before the arrival* of a job?
 3. What line should be changed so that we count the number of jobs in the system *just after* the arrival? What should it become?
 4. Why is this algorithm (very) inefficient?

Making the algorithm much more efficient is not hard.

Python Code

```

1  L = np.zeros_like(A)
2  idx = 0
3  for k in range(1, len(A)):
4      while D[idx] < A[k]:
5          idx += 1
6      L[k] = k - idx

```

Ex 3.18. Why is this much better?

Ex 3.19. Here is another algorithm that plots $\{L(t)\}$. Explain how it works, and its difference with the previous algorithm.

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("figures/wait4.pdf")

```

3.5 MULTI-SERVER QUEUE

Let us now generalize the simulation to a $G/G/c$ queue, i.e., there are $c \geq 1$ servers present to serve jobs. However, before doing that, we approximate the $G/G/c$ queueing process by a $G/G/1$ queue in which the server works c times as fast. Like this we can reuse the code above to approximate a $G/G/c$ queue.

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 . As a first simple case, we model the multi-server queueing system as if there is one fast server working at rate c .

Ex 3.20. Explain that we can implement a fast server by changing the service times as follows:

Python Code

```

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

```

Here is the same idea, but the implementation is slightly different. This is useful because we will see that we can generalize this to a multi-server queue in which the servers have different speeds.

Python Code

```

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

```

3.5.2 A real multi-server queue

Here is the code to implement the $G/G/2$ queue; see the queueing book to understand the algorithm. (I include the print statements so that you can see step by step how it works.)

There is one very convenient feature of numpy, which is used in the this line of the $G/G/2$ queue code below. The feature is called *broadcasting* and is best explained with an example.

Ex 3.21. Run this code and explain what it does.

Python Code

```

1 import numpy as np
2
3 a = np.array([2, 3, 4])
4 a = a - 1 # or a -= 1
5 print(a)

```

This is the code for the $G/G/2$ queue.

Python Code

```

1 import numpy as np
2
3 num = 5
4
5 X = np.ones(num + 1)
6 S = 5 * np.ones(num)
7
8 c = np.array([1.0, 1.0])
9 W = np.zeros_like(S) # store the waiting time as seen by the kth job
10 w = np.zeros_like(c) # waiting times at each of the servers
11 for k in range(1, len(S)):
12     s = w.argmax() # server with smallest waiting time
13     W[k] = w[s]
14     print(f"{k=}", w, W[k])
15     w[s] += S[k]
16     print(w)
17     print(w - X[k + 1])
18     w = np.maximum(w - X[k + 1], 0) # this
19     print(w)
20
21 print(W)

```

Ex 3.22. Run the code, and explain the results.

As as a test we could set the vector of server capacities like $c=[1]$ so that we reduce our multi-server queue to a single-server queue, and compare the output with our earlier simulators. 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.

Now that we have tested the implementation (in part), here is the code for a queue served by three servers, but they can work at different speeds.

Python Code

```

1  import numpy as np
2
3  np.random.seed(3)
4
5  labda = 3
6  mu = 1.1
7  N = 1000
8
9  X = np.random.exponential(scale=1 / labda, size=N + 1)
10 S = np.random.exponential(scale=1 / mu, size=N)
11
12 c = np.array([1.0, 1.0, 1.0])
13 W = np.zeros_like(S)
14 w = np.zeros_like(c)
15 for k in range(len(S)):
16     s = w.argmin()
17     W[k] = w[s]
18     w[s] += S[k] / c[s]
19     w = np.maximum(w - X[k + 1], 0)
20
21 print(W.mean(), W.std())

```

Ex 3.23. Where in the code do we handle the fact that servers can work at different speeds?

Ex 3.24. Compare $E[W]$ for two cases. The first is a $G/G/1$ queue with a 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 3.25. Change the code for the multi-server such that the individual servers have different speeds but 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 3.26. 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 smallest average queueing time and variance? If the variance is larger when the service rates in the multi-server queue are different 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 4.1.**
1. Why is the expected number of arrivals in a period equal to 24?
 2. Run this code (and write down the results).
 3. Then change the code such that the arrivals can be served on the day they arrive. Rerun and compare the results to the case in which arrivals cannot be served on the day in which they arrive.

Ex 4.2. 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).

1. Why are now the expected number of arrivals during a period equal to 3?
2. Rerun the code with and without the arrivals being served on the period of arrival.
3. compare to the previous exercise. You should notice that there is less difference in $E[L]$ between the models in which you serve or don't serve jobs in the period they arrive. But why is that so?

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 4.3. Explain that we use Little's law to compute $E[L]$. Can we use this law to estimate $V[L]$?

Ex 4.4. Run the code and compare with the discrete time simulation.

4.2.3 Sakasegawa's formula

Now we use Sakasegawa's formula, rather than simulation, 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 4.5. Run the code, and compare the value with the discrete and continuous time simulations.

You should notice that there is quite a difference between the estimates for $E[L]$ we obtain from simulation, on the one hand, and from Sakasegawa's formula, on the other hand. We conclude that it takes quite a bit of time for a simulation to reach 'steady state'.

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, we should be carefull to avoid a situation which the load is too large (recall, setup times add to the service times).

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=} >= 1"

```

- Ex 4.6.**
1. What does the assert statement do?
 2. Why do I put it here, i.e., before doing any other work?
 3. What happens if you would set labda=10 and 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 4.7. 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 ES = ES0 + R / B
6 ESB = B * ES
7 VS0 = ES0 * ES0
8 VSB = VR + B * VS0
9 cs2 = VSB / ESB ** 2
```

Ex 4.8. If the interarrival times $X_k \sim \text{Exp}(\lambda)$, why should we set $ca2 = 1/B$?

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 4.9. Why should we call the function sakasegawa with the following parameters sakasegawa(labdaB, ESB, ca2, cs2)?

The last step of the queueing process can be coded like this:

Python Code

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

Ex 4.10. What is the meaning of $W3$?

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 num = 1000
2 num = B * (num // B) # get multiple of B
3 X = np.random.exponential(scale=1 / labda, size=num)
4 A0 = np.zeros(len(X) + 1)
```

```

5  A0[1:] = X.cumsum()
6  A = np.zeros_like(A0)
7  for i in range(num // B):
8      st = i * B + 1 # start
9      fi = (i + 1) * B # finish
10     A[st : fi + 1] = A0[fi]
11
12
13  S0 = np.random.exponential(scale=1 / mu, size=len(A))
14  S0[0] = 0
15  S = S0.copy()
16  S[1::B] += R
17
18  D = np.zeros_like(A)
19  for k in range(1, len(A)):
20      D[k] = max(D[k - 1], A[k]) + S[k]
21
22  J = D - A0
23  print(J.mean())

```

Ex 4.11. Clearly, A_0 are the arrival times of the jobs at the system. Explain that A corresponds to the arrival times of the *batches* at the queue, and that S are the service times including setup times. (Use print statements to understand how the slicing, i.e., notation like $A[\text{st} : \text{fi} + 1]$, works.)

Ex 4.12. Why is $D - A_0$ the sojourn time, and not $D - A$?

Ex 4.13. 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.

4.4 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. I expect you now to be able to fill in Sakasegawa's formula. It remains to show you how to simulate this system.

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 4.14. Explain how D is computed.

Ex 4.15. To test the code I set at first $R = 0 * np.ones(len(A))$ in the line marked as this. Why would I do this (to what simpler queueing system can I compare the results of this program?)

Ex 4.16. 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 $num = 100000$, i.e., 10 times as large. What is the effect?

Ex 4.17. 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?

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.

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 X = np.random.exponential(scale=1 / labda, size=num)
14 X[0] = 0
15 A = X.cumsum()
16
17 S0 = np.random.exponential(scale=ES0, size=num)
18 N = np.random.poisson(labda_f * ES0, len(S0))
19 R = expon(scale=ER)
20 S = np.zeros_like(S0)
21 for i in range(len(S0)):
22     S[i] = S0[i] + R.rvs(N[i]).sum()
23
24
25 D = np.zeros_like(A)
26 for k in range(1, len(A)):
27     D[k] = max(D[k - 1], A[k]) + S[k]
28
29 W = D - A - S
30 print(W.mean())

```

Ex 4.18. Run the code and include the results in your assignment. Where do we include the break downs?

Ex 4.19. 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 $E[W]$?

4.6 A SIMPLE TANDEM NETWORK

We have two queues in tandem. Here is the code to simulate this.

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 4.20. Why is $J = D2 - A1$?

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 5.1. Run the code. What is the purpose of the this line?

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 achieve this, we can change the for loop as follows.

Python Code

```

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

```

Ex 5.2. Run this code and check its effect on $E[W]$, $V[W]$, and $\max\{W_k\}$.

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.

Here is better code.

Python Code

```

1 import numpy as np
2 from scipy.stats import expon
3
4 np.random.seed(3)
5
6 num = 100
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
15 W = np.zeros_like(S)
16 I = np.ones_like(S, dtype=int)
17 for k in range(1, len(S)):
18     W[k] = max(W[k - 1] + S[k - 1] * I[k - 1] - X[k], 0)
19     if W[k] >= T:
20         I[k] = 0
21
22 print(W.mean(), I.mean(), W[I == 1].mean())

```

Ex 5.3. What does the vector I represent? With this, explain why we need to filter only the jobs k with $I_k = 1$ in the computation of $E[W]$.

5.2.4 Some other blocking rules

There are other rules to block jobs.

Ex 5.4. In this code, what is the rule to block jobs?

Python Code

```

1 W = np.zeros_like(S)
2 I = np.ones_like(S, dtype=int)
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(), W[I == 1].mean())

```

Ex 5.5. Likewise, how does this rule work? What is the meaning of V ?

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())

```

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 = 10000
7 labda = 1
8 mu = 2 * 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-1]).sum() - X[k], 0)
16
17 ES = 1/mu
18 rho = labda * ES
19 print(S.mean(), 1 / mu)
20 print(rho / (1 - rho) * ES, W.mean(), W.max())

```

Ex 5.6. Explain how this code works.

Ex 5.7. Run the code. Why do I take B as it is here (recall, I like to test)? Why should $W.mean()$ and $\rho E[S] / (1 - \rho)$ be approximately equal?

The next remark applies only if you're interested, otherwise skip. `numpy.random` provides functionality to generate random numbers, but not more, while `scipy.stats` provides much more useful functions. For the above, observe that `expon` loads from `scipy.stats`, while the other `rvs` come from `numpy.random`. Check the manuals on the web for further information.

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 5.8. 1. Explain how the code works.

2. What do the printed KPIs mean?

3. Finally, in the this line, why is it better to 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 $\{L_k\}$ are easy, while in continuous time the recursions for $\{W_k\}$ or $\{J_k\}$ 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 5.9. 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 5.10. Just to check that you really understand: why is it ok here to use $(L == 0)$ rather than np.close ?

Ex 5.11. Why do I compare $L.\text{mean}()$ to $\rho/(1 - \rho)$ and not to $\rho^2/1 - \rho$?

Ex 5.12. <2022-03-16 wo> Please skip this exercise. Change μ to 1.05λ . Now the results of the simulation are not very good if $\text{num}=1000$ 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 5.13. Read the documentation of how a deque works, then explain the code.

Ex 5.14. Do we get a sensible queueing discipline if we would use the `pop()` instead of `popleft()`?

Ex 5.15. Run this code with $T=100$ (or some other big number) 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 5.16. 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 $M/M/1/T$ queue?

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

SIMULATION WITH EVENT STACKS

6.1 GENERAL INFO

The tools we developed up to now do not suffice to simulate more difficult (queueing) systems. For harder cases we need a concept called *event stack*. Here we show how to implement an event stack as a *heap queue*. Finally, we will use python *classes* as these are very natural objects to use in complicated simulations. At the end you will see how general and clean our final simulation environment is.

6.2 SORTING WITH HEAP QUEUES

We start with a simple problem: how to sort a bunch of students by age? For this we will use a *heap queue*. This is a stack on which we can push elements, then, when popping from the stack, the elements come out in a sorted way.

Ex 6.1. To understand the code below, first look up on the web:

1. what is the difference between a *tuple* and a *list*?
2. When to use one or the other?

As always, keep your answer brief, just show that you understand; that suffices.

The heap in our code below uses the first element of the *tuple* of a student's age and name to sort the students. Let's put a number of students on the stack. We need a new module `heapq`; this provides the algorithms for pushing to and popping from a stack.

Python Code

```

1  from heapq import heappop, heappush
2
3  stack = []
4
5  heappush(stack, (25, "Cynthia"))
6  heappush(stack, (21, "James"))
7  heappush(stack, (21, "James"))
8  heappush(stack, (20, "Pete"))
9  heappush(stack, (18, "Clair"))
10 heappush(stack, (14, "Jim"))
11 print(stack)

```

Observe that with popping we take the element from the stack with the smallest age. You should know that anytime we pop from the stack, the stack becomes one element shorter. (And with pushing it becomes longer.)

Python Code

```

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

```

Ex 6.2. Add the following lines to the code above, and include the output. You should see that the students are printed in order of age.

Python Code

```

1 while stack:
2     e = heappop(stack)
3     print(e)

```

6.3 CLASSES

In a class we can 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. 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 ]

```

Ex 6.4. Add two more students, e.g., you and your group mate, to the list. (And if you don't like to spread such details, just lie about your age :-)) Show your code.

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

Python Code

```

1 from heapq import heappop, heappush
2
3
4 class Student:
5     def __init__(self, name, age, phone):
6         self.name = name
7         self.age = age
8         self.phone = phone
9

```

```

10     def __repr__(self):
11         return f"{self.name}, {self.age}, {self.phone}"
12
13
14     students = [
15         Student("Joseph", "18", "Motorola"),
16         Student("Maria", "21", "Huawei"),
17         Student("Natasha", "20", "Apple"),
18         Student("Chris", "25", "Nexus"),
19     ]
20
21     stack = []
22
23     for s in students:
24         heappush(stack, (s.phone, s))
25
26     while stack:
27         print(heapop(stack))

```

Ex 6.5. What line makes that the students are sorted by phone brand?

Ex 6.6. Change the code so that you can sort the students by age. Show the line(s) of your code to achieve this.

Ex 6.7. Extend the Job class such that we can give the student also a surname. Then invent some surnames, and sort the students by surname. Include your code to show how you did this.

6.4 A JOB CLASS USEFUL FOR SIMULATION

How can use heap queues in the simulation of queueing systems? To see this, think of time as a sequence of events in which things happen. In a queueing system, three things can happen: a job arrives, a job service starts, or a job leaves. For a job we should specify its arrival and departure times. We store the jobs in a heap queue and add the arrival time as a key to pop jobs in order of arrival.

We store the arrival, service, and departure time as the *attributes* of a job class. We also store the queue length at arrival times to gather statistics once the simulation is over.

Python Code

```

1  from heapq import heappop, heappush
2
3  ARRIVAL, DEPARTURE = 0, 1
4
5
6  class Job:
7      def __init__(self):
8          self.arrival_time = 0
9          self.service_time = 0
10         self.departure_time = 0
11         self.queue_length_at_arrival = 0
12
13     def sojourn_time(self):
14         return self.departure_time - self.arrival_time
15
16     def waiting_time(self):

```



```

17         return self.sojourn_time() - self.service_time
18
19     def service_start(self):
20         return self.departure_time - self.service_time
21
22     def __repr__(self):
23         return f"{self.arrival_time}, {self.service_time}, {self.service_start()}, \
24                 {self.departure_time}\n"
25
26     def __le__(self, other):
27         # this is necessary to sort jobs when they have the same arrival times.
28         return self.id <= other.id
29
30
31 events = [] # event stack, global
32 num_jobs = 5
33 interarrival_time = 3
34 service_time = 5
35
36 arrival_time = 0
37 for i in range(num_jobs):
38     arrival_time += interarrival_time
39     job = Job()
40     job.arrival_time = arrival_time
41     job.service_time = service_time
42     heappush(events, (job.arrival_time, job, ARRIVAL))
43
44 while events:
45     time, job, typ = heappop(events)
46     print(job)

```

Ex 6.8. Explain the waiting time and sojourn functions of the Job class.

For later purposes, we have to add labels to the events in the heap queue to indicate what type of event we are dealing with, an arrival or a departure.

Ex 6.9. Change the interarrival_time to some number you like. Run the code, include your results, and explain the output of the print statement.

6.5 A VERY POWERFUL GG1 SERVER SIMULATOR

The code here below is a MWE (minimal working example). You should copy it and run it. Below we'll ask questions in the exercises about how the code works.

6.5.1 Start

Copy this first.

```

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
7 np.random.seed(3)

```

Python Code

```
8  
9  
10 ARRIVAL, DEPARTURE = 0, 1  
11 IDLE, BUSY = 0, 1
```

Then put the Job class after this.

6.5.2 GG1 simulator.

Copy GG1 class and put this after the Job class. We put the code on one page so that the copying is easy.

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 6.10. Explain how the following methods work:

1. arrivals: `GG1.handle_arrival()`;
2. departures: `GG1.handle_departure`;
3. running the simulation: `GG1.run`.

You should observe that, with a heap queue, we can let the heap queue do all the work of tracking time. It is not easy to understand how we use a queue (that is, the heap queue) to simulate queueing systems. However, once you understand how event stacks work, you know all there is to discrete event simulation.

6.5.3 Running the simulator

Put this code the GG1 class.

Python Code

```

1  labda = 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()
12
13 print(events[:5]) # this

```

Ex 6.11. Explain the results of the 'this' line. Why is the content of column 5 negative?

6.5.4 Analyze the results, statistics

Put this at the end of the code.

Python Code

```

1  gg1.run()
2
3  sojourn = np.zeros(len(gg1.served_jobs))
4  for i, job in enumerate(gg1.served_jobs):
5      sojourn[i] = job.sojourn_time()
6
7  print(sojourn.mean(), sojourn.std(), sojourn.max())
8
9  plt.clf()
10 plt.plot(sojourn)
11 plt.savefig('figures/sojourn0.png')

```

Ex 6.12. 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. Include also the statistics such as the mean.

If we would use this code for real world problems, we of course should compare the results of this simulator to our earlier simulators; for the assignment we skip this.

6.6 OTHER SCHEDULING RULES

Suppose we prefer to serve the shortest job in queue. We can inherit all of our GG1 queue, except the scheduling rule. For this we need to change just one line of the class. The implementation of other scheduling is also extremely easy now

6.6.1 SPTF scheduling

Here is a class to simulate the SPTF (shortest processing time first). We can derive all from GG1 except that we have to change the rule on how to put the jobs on the queue.

Python Code

```

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

```

Ex 6.13. Explain how this code implements SPTF.

Ex 6.14. To run it, use this code. Include the results.

Python Code

```

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

```

Ex 6.15. If you're interested (but you can skip this) print at the end the average number of jobs in the system and compare that to $E[L]$ for the FIFO queue. What difference do you see?

6.6.2 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 6.16. Run an example with 100 jobs. Make a graph of the waiting times and the sojourn times. Comment on your findings.

6.6.3 Serve longest job first

Ex 6.17. Update the code of the SPTF queue such that the longest job is selected from the queue, rather than the shortest. Include the relevant lines of the code, and explain why it implements the SLJF rule.

Hopefully you see how easy it is (now we have done all the hard work) to compare different job scheduling rules. For instance, for the earliest due data first rule, we add a due date attribute to each job, and when selecting a job for service we choose the earliest due date.

HINTS

h.1.6. 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.1.11. 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.1.13. 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.3.3. The X_k are iid.

h.3.4. 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.3.6. Did we really serve job 0? If num is big number, does it matter that we set $S[0]=0$?

h.3.7. We subtract $S[0]$ as if we served the corresponding job, but did we actually serve it?

h.3.8. We include $W[0]$, but what is that?

h.3.15. 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.3.18. The improved algorithm in $O(n)$, while the previous was $O(n^2)$.

h.3.24. Take $c=[3.]$ for the first case, and $c=[1., 1., 1.]$ for the second. With this, the change in code is minimal.

h.3.25. For instance, set $c=\text{np.array}([2, 0.5, 0.5])$.

h.4.1. Replace the relevant line by $d = \min(c[i], L[i - 1] + a[i])$.

h.4.12. Recall that A_0 are the arrival times at the system before batching while A are the times the jobs move as a batch to the queue.

h.5.3. If $I[k] == 1$, then what happens to job k ?

h.5.7. When the batches $B = \text{np.random.randint}(1, 2, \text{size=num})$. Hence, we deal with the $M/M/1$ queue for this choice of B .

h.5.9. 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.5.10. Is L a float?

h.5.11. What is $\rho^2/(1 - \rho)$?

h.5.14. It's not FIFO.

h.5.15. Is `L.max()` larger than 100 for this simulation?

h.6.11. Realize that we have not yet called `gg1.run()`.

h.6.17. Put a minus sign at the right position when adding a job to the queue heap.