# Stochastic Simulation (MIE1613H) - Homework 3 (Solutions)

- Submit your homework on Quercus as a single PDF file by the deadline. Late submissions are penalized 10% each day the homework is late.

- At the top of your homework include your name, student number, department, and program.

- You may discuss the assignment with other students in general terms, but each student must solve the problems, write the code and the solutions individually.

- **The simulation models must be programmed in Python. For the input modelling questions you may use R or Python.** You must include both the source code (including comments to make it easy to follow) and the output of the simulation in your submission.

- **It is very important to explain your answers and solution approach clearly and not just provide the code/results**. Full marks will only be given to correct solutions that are fully and clearly explained.

**Problem 1. (25 Pts.)** You are asked to make projections about cycle times for a semiconductor manufacturer who plans to open a new plant. Here "cycle time" means the time from product release until completion. The process that you will consider has a single diffusion step with sub-steps as indicated in the diagram.

Raw material for two products (C) and (D) will begin at the CLEAN step, and make multiple passes until the product completes processing. Movement within each process is handled by robots and takes very little time (treat as 0) relative to the processing steps. The movement time from release to diffusion is 15 minutes (1/4 hour).

The anticipated product mix is to have 60% of products to be of type (C) and 40% of type (D). Product (C) requires 5 passes and product D requires 3 passes.

The OXDIZE step is deterministic but differs by product type: it takes 2.7 hours for (C) and 2.0 hours for (D).

**(a)** The CLEAN and LOAD/UNLOAD steps do not differ by product type but are subject to uncertainty. Historical data is provided in (SemiconductorData.xls). Represent them using the distributions available in PythonSim, and justify your choice using the graphical and statistical methods discussed in the lecture.

Product will be released in cassettes at the rate of 1 cassette/hour, 7 days a week, 24-hours a day (this is to achieve a desired throughput of 1 cassette/hour). Product is moved and processed in single cassette loads.

The table below shows the number of machines that are planned for each fabrication step:
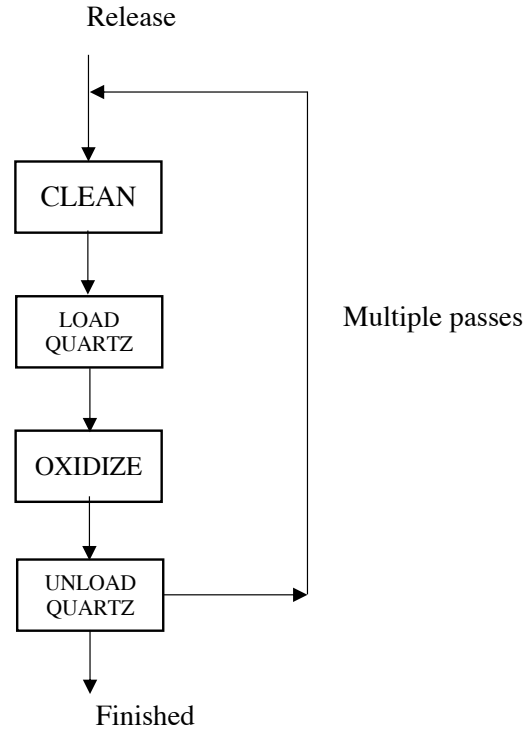
Release

CLEAN

LOAD
QUARTZ

Multiple passes

OXIDIZE

UNLOAD
QUARTZ

Finished

Figure 1: Diagram of the diffusion step.

| Step | Number of Machines |
|---|---|
| CLEAN | 9 |
| LOAD QTZ | 2 |
| OXD | 11 |
| UNLOAD QTZ | 2 |

Table 1: Number of available machines in each step.

**(b)** Provide estimates of the long-run average cycle times for each product. Determine and justify a warm-up period and run length. You may use the replication-deletion approach. Provide appropriate confidence intervals for your estimates.

(a) **(10 Points)** Since we do not have a physical justification for the choice of distributions, we can speculate a few candidate distributions (e.g., by looking at the shape of the histogram) in each case and examine the goodness of fit for them. For the CLEAN data, we pick Exponential, Weibull, and Gamma distributions. The KS test does not reject any of the distributions at $(1 - \alpha) = 95\%$ confidence level. Examining the Q-Q plot we observe that the fits are very similar, so any of them would provide an acceptable choice. We pick exponential that is available in PythonSim. The estimated parameter for exponential is $\lambda = 0.5512835$, so the estimate of mean, that is used in PythonSim, is $1/0.5512835 = 1.813949$.

For the LOAD/UNLOAD data, we pick Lognormal, Weibull, and Gamma. The KS test does not reject any of the distributions. Examining the Q-Q plot however suggests that Gamma provides a better fit at the tail. However, we pick Lognormal that is available in PythonSim. Let $X$ be the
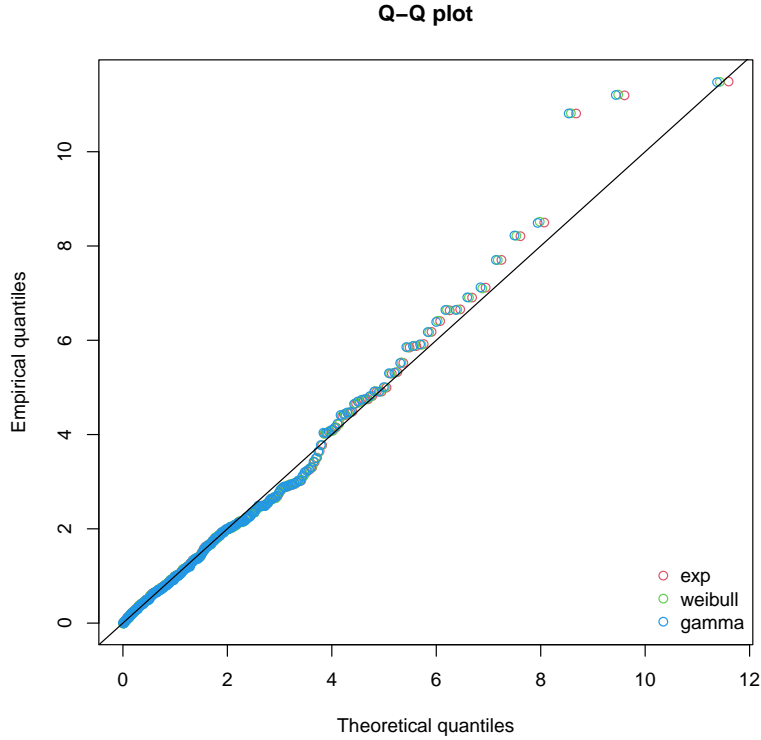
**Q–Q plot**



Figure 2: Q-Q plot for the fitted distributions to the CLEAN data.

Lognormal random variable and denote by $\mu$ and $\sigma$ the mean and standard deviation of $\log(X)$ as estimated in R. The mean and standard deviation of X, that is required in PythonSim, are given by,

$$\mathrm{E}(X) = e^{\mu + \frac{\sigma^2}{2}} = 0.2475798,$$

$$\mathrm{SD}(X) = e^{\mu + \frac{\sigma^2}{2}} \sqrt{e^{\sigma^2} - 1} = 0.2233498.$$

```
1   # This file uses the fitdistrplus package. You'll need to install it first
2   # by running the command: 'install.packages("fitdistrplus")' in R
3   require(fitdistrplus)
4
5   # import the data from .csv file
6   MyData <- read.csv("SemiconductorData.csv", header=FALSE)
7   Load <- MyData[,2]
8   Clean <- MyData[1:299,4]
9   plotdist(Load, histo = TRUE, demp = TRUE)
10  plotdist(Clean, histo = TRUE, demp = TRUE)
11
12  # fit a lognormal,gamma, and weibull to load data
13  load_fln <- fitdist(Load, "lnorm")
14  load_fw <- fitdist(Load, "weibull")
15  load_fg <- fitdist(Load, "gamma")
16
```
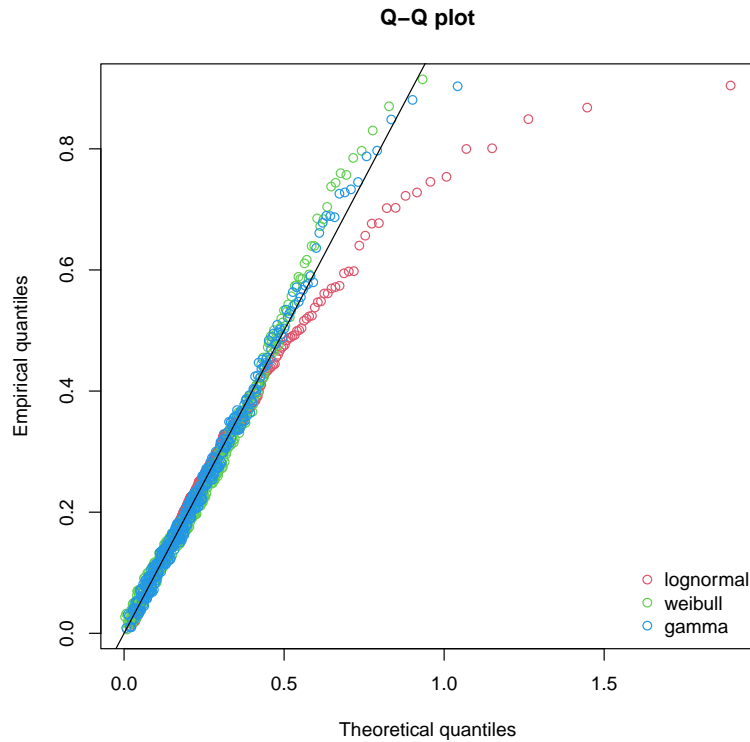
**Q–Q plot**



Figure 3: Q-Q plot for the fitted distributions to the LOAD data.

```
17  # fit an exponential, gamma, and weibull to clean data
18  clean_fe <- fitdist(Clean, "exp")
19  clean_fw <- fitdist(Clean, "weibull")
20  clean_fg <- fitdist(Clean, "gamma")
21
22  # compare the fits for load
23  plot.legend <- c("lognormal", "weibull", "gamma")
24  qqcomp(list(load_fln, load_fw, load_fg), legendtext = plot.legend)
25
26  # compare the fits for clean
27  plot.legend <- c("exp", "weibull", "gamma")
28  qqcomp(list(clean_fe, clean_fw, clean_fg), legendtext = plot.legend)
29
30  # Perform GofFit tests
31  gof_load <- gofstat(list(load_fln, load_fw, load_fg), fitnames = c("lognormal"
        , "weibull", "gamma"))
32  gof_load
33  gof_load$kstest
34
35  gof_clean <- gofstat(list(clean_fe, clean_fw, clean_fg), fitnames = c("exp", "
        weibull", "gamma"))
36  gof_clean
37  gof_clean$kstest
```

(b) **(15 Points)** We define a special entity class in SimClasses.py as follows to represent a product.

```
1    # special entity to represent a product
2    class Product():
3        def __init__(self, type):
4            self.CreateTime = Clock
5            self.Type = type
6            if self.Type == "C":
7                self.Passes = 5
8                self.Oxd = 2.7
9            else:
10                self.Passes = 3
11                self.Oxd = 2.0
```

When a product is released, we schedule the next release for the next hour, determine the product type with the given probabilities, and schedule the end of movement from release to diffusion using the SchedulePlus function. In the EndMove function, we schedule the end of cleaning if a cleaning machine is available using the estimated exponential distribution from the data. Otherwise, we put the product in the cleaning queue. In the EndClean function, we schedule the end of loading if a loading machine is available using the estimated Lognormal distribution from the data. Otherwise, we put the product in the loading queue. We also check the cleaning queue and if there is a product waiting, we schedule the next end of cleaning for that product. Otherwise, we make a cleaning machine idle. The EndLoad and EndOxd functions are also similar to the EndClean function. In the EndUnload function, we keep track of the product passes and if a product reaches the required passes, we record the cycle time for that product type in its corresponding DTStat object. Otherwise, we check the status of the cleaning machines and repeat the same steps.

Fig. 4 shows the mean cycle time by product number for the first $m = 1000$ products of each type averaged across $n = 30$ replications. At somewhere around 100 products the average cycle times seem to vary around a central value, so we may take $d = 100$ and choose the run length at least $m = 1100$ of each product. The estimate of the average cycle times for product C and D is 31.296 and 16.700 hours with the 95% CI of [31.28, 31.31] and [16.69, 16.71], respectively.

```
1    # Cycle times for a semiconductor manufacturer
2    import SimFunctions
3    import SimRNG
4    import SimClasses
5    import scipy.stats as sp
6    import numpy as np
7    import pandas as pd
8    import matplotlib.pyplot as plt
9
10   def mean_confidence_interval(data, confidence=0.95):
11       a = 1.0*np.array(data)
12       n = len(a)
13       m, se = np.mean(a), sp.sem(a)
14       h = se * sp.t._ppf((1+confidence)/2., n-1)
15       return m, m-h, m+h
16
17   ZSimRNG = SimRNG.InitializeRNSeed()
18   Calendar = SimClasses.EventCalendar()
19
20   # set up queues
```
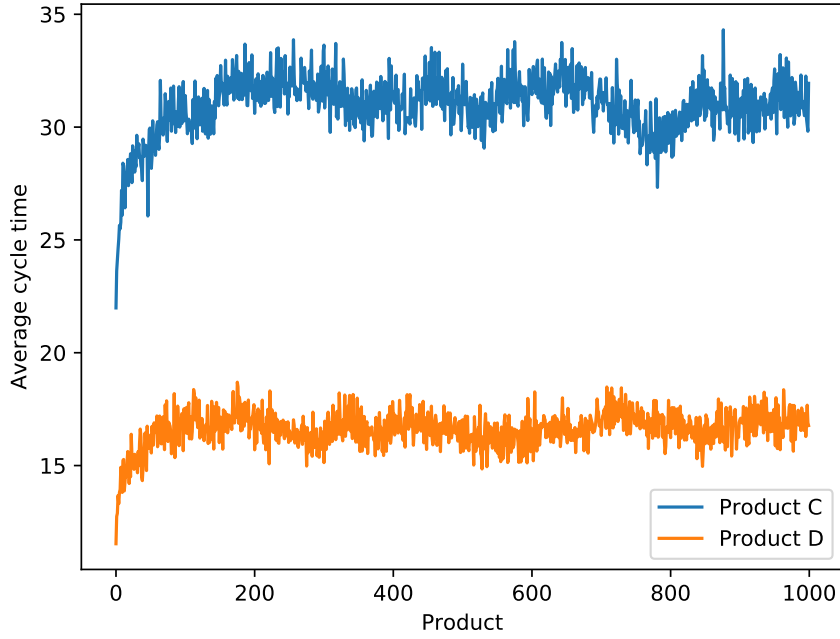
Figure 4: Mean plot, by product number, averaged across $n = 30$ replications.

```
21  CleanQ = SimClasses.FIFOQueue()
22  LoadQ = SimClasses.FIFOQueue()
23  OxdQ = SimClasses.FIFOQueue()
24  UnloadQ = SimClasses.FIFOQueue()
25  # set up resources
26  Clean = SimClasses.Resource()
27  Load = SimClasses.Resource()
28  Oxd = SimClasses.Resource()
29  Unload = SimClasses.Resource()
30
31  Clean.SetUnits(9)
32  Load.SetUnits(2)
33  Oxd.SetUnits(11)
34  Unload.SetUnits(2)
35
36  # set up statistics
37  CycleTimeC = SimClasses.DTStat()
38  CycleTimeD = SimClasses.DTStat()
39
40  TheCTStats = []
41  TheDTStats = [CycleTimeC, CycleTimeD]
42  TheQueues = [CleanQ, LoadQ, OxdQ, UnloadQ]
43  TheResources = [Clean, Load, Oxd, Unload]
44
45
46  def Release():
47      SimFunctions.Schedule(Calendar, "Release", 1.0)
```

```python
48          if SimRNG.Uniform(0, 1, 1) < 0.6:
49              NewProduct = SimClasses.Product("C")
50          else:
51              NewProduct = SimClasses.Product("D")
52
53          SimFunctions.SchedulePlus(Calendar, "EndMove", 0.25, NewProduct)
54
55  def EndMove(TheProduct):
56      if Clean.Busy < Clean.NumberOfUnits:
57          Clean.Seize(1)
58          SimFunctions.SchedulePlus(Calendar, "EndClean", SimRNG.Expon(1.813949,
                    2), TheProduct)
59      else:
60          CleanQ.Add(TheProduct)
61
62
63  def EndClean(TheProduct):
64      if Load.Busy < Load.NumberOfUnits:
65          Load.Seize(1)
66          SimFunctions.SchedulePlus(Calendar, "EndLoad", SimRNG.Lognormal
                    (0.2475798, 0.2233498 ** 2, 3), TheProduct)
67      else:
68          LoadQ.Add(TheProduct)
69
70      if CleanQ.NumQueue() > 0:
71          NextProduct = CleanQ.Remove()
72          SimFunctions.SchedulePlus(Calendar, "EndClean", SimRNG.Expon(1.813949,
                    2), NextProduct)
73      else:
74          Clean.Free(1)
75
76
77  def EndLoad(TheProduct):
78      if Oxd.Busy < Oxd.NumberOfUnits:
79          Oxd.Seize(1)
80          SimFunctions.SchedulePlus(Calendar, "EndOxd", TheProduct.Oxd,
                    TheProduct)
81      else:
82          OxdQ.Add(TheProduct)
83
84      if LoadQ.NumQueue() > 0:
85          NextProduct = LoadQ.Remove()
86          SimFunctions.SchedulePlus(Calendar, "EndLoad", SimRNG.Lognormal
                    (0.2475798, 0.2233498 ** 2, 3), NextProduct)
87      else:
88          Load.Free(1)
89
90
91  def EndOxd(TheProduct):
92      if Unload.Busy < Unload.NumberOfUnits:
93          Unload.Seize(1)
94          SimFunctions.SchedulePlus(Calendar, "EndUnload", SimRNG.Lognormal
                    (0.2475798, 0.2233498 ** 2, 3), TheProduct)
95      else:
```

```python
96              UnloadQ.Add(TheProduct)
97
98          if OxdQ.NumQueue() > 0:
99              NextProduct = OxdQ.Remove()
100             SimFunctions.SchedulePlus(Calendar, "EndOxd", NextProduct.Oxd,
                    NextProduct)
101         else:
102             Oxd.Free(1)
103
104
105   def EndUnload(TheProduct):
106       global CTC, CTD, TC, TD
107       TheProduct.Passes = TheProduct.Passes - 1
108       if TheProduct.Passes > 0:
109           if Clean.Busy < Clean.NumberOfUnits:
110               Clean.Seize(1)
111               SimFunctions.SchedulePlus(Calendar, "EndClean", SimRNG.Expon
                      (1.813949, 2), TheProduct)
112           else:
113               CleanQ.Add(TheProduct)
114       else:
115           if TheProduct.Type == "C":
116               CycleTimeC.Record(SimClasses.Clock - TheProduct.CreateTime)
117               CTC.append(SimClasses.Clock - TheProduct.CreateTime)
118               TC += 1
119           else:
120               CycleTimeD.Record(SimClasses.Clock - TheProduct.CreateTime)
121               CTD.append(SimClasses.Clock - TheProduct.CreateTime)
122               TD += 1
123
124       if UnloadQ.NumQueue() > 0:
125           NextProduct = UnloadQ.Remove()
126           SimFunctions.SchedulePlus(Calendar, "EndUnload", SimRNG.Lognormal
                  (0.2475798, 0.2233498 ** 2, 3), NextProduct)
127       else:
128           Unload.Free(1)
129
130   AcrossC = []
131   AcrossD = []
132   AllCTC = []  # these will be a list of lists, each
133   AllCTD = []  # list corresponding to one replication's cycle times
134
135   m = 1100  #
136   d = 100  # Deletion point
137   for reps in range(0, 30, 1):
138       CTC = []  # cycle times from current replication
139       CTD = []  # cycle times from current replication
140       SimFunctions.SimFunctionsInit(Calendar, TheQueues, TheCTStats, TheDTStats,
                  TheResources)
141       SimFunctions.Schedule(Calendar, "Release", 1.0)
142
143       NextEvent = Calendar.Remove()
144       SimClasses.Clock = NextEvent.EventTime
145       Release()
```

```python
146
147        TC = 0  # Number of completed products of type C
148        TD = 0  # Number of completed products of type D
149        while min(TC, TD) < m:
150            if min(TC, TD) == d:
151                SimFunctions.ClearStats(TheCTStats, TheDTStats)
152
153            NextEvent = Calendar.Remove()
154            SimClasses.Clock = NextEvent.EventTime
155            if NextEvent.EventType == "Release":
156                Release()
157            elif NextEvent.EventType == "EndMove":
158                EndMove(NextEvent.WhichObject)
159            elif NextEvent.EventType == "EndClean":
160                EndClean(NextEvent.WhichObject)
161            elif NextEvent.EventType == "EndLoad":
162                EndLoad(NextEvent.WhichObject)
163            elif NextEvent.EventType == "EndOxd":
164                EndOxd(NextEvent.WhichObject)
165            elif NextEvent.EventType == "EndUnload":
166                EndUnload(NextEvent.WhichObject)
167
168
169        AcrossC.append(CycleTimeC.Mean())
170        AcrossD.append(CycleTimeD.Mean())
171        AllCTC.append(CTC)
172        AllCTD.append(CTD)
173
174    # output results
175    print("Estimated_expected_cycle_time_for_product_C:", mean_confidence_interval
         (AcrossC, 0.05))
176    print("Estimated_expected_cycle_time_for_product_D:", mean_confidence_interval
         (AcrossD, 0.05))
177
178    # Create mean plot to determine the warm-up period
179    MeanC = []
180    MeanD = []
181    for i in range(1000):  # Average cycle time for the first 1000 products
182        MeanC.append(np.mean([rep[i] for rep in AllCTC]))
183        MeanD.append(np.mean([rep[i] for rep in AllCTD]))
184    plt.plot(MeanC, label = "Product_C")
185    plt.plot(MeanD, label = "Product_D")
186    plt.xlabel("Product")
187    plt.ylabel("Average_cycle_time")
188    plt.legend()
189    plt.show()
```

Estimated expected cycle time for product C: (31.29, 31.28, 31.31)
Estimated expected cycle time for product D: (16.70, 16.69, 16.71)

**Problem 2. (20 Pts.)** An Automated external defibrillator (AED) can save a person's life in event of a cardiac arrest. To accelerate delivery, start ups have been developing drone technology to quickly deliver AEDs to the scene of the cardiac arrest in case of an emergency call. The AEDs will be maintained at various bases across the city to respond to calls in the area designated to each base.

You have been tasked to estimate the minimum number of drones at a certain base to ensure that with 95% probability there will be a drone available at the event of a cardiac arrest in the area covered by that base.

Assume that calls reporting cardiac arrests arrive according to a non-homogeneous Poisson process with (per minute) rate function,

$$\lambda(t) = \begin{cases} 4, & \text{7AM-12PM,} \\ 2, & \text{12PM-12AM,} \\ 1, & \text{12AM-7AM.} \end{cases}$$

throughout a day. Further, the total "drone busy time" i.e. from the time the call is received until the drone is back to the base and ready to dispatch again is exponentially distributed with mean 45 minutes. Assume at time 12AM, there are no calls (requests for AEDs) in the system.

**(a)** Explain, in words, your approach in determining the minimum number of drones at the base required to satisfy the provided service level. Specify the queueing model, corresponding parameters, and the performance measure you are estimating.

**(b)** Determine the minimum number of drones at the base using the inputs provided. Provide a 95% CI for your estimate.

(a) **(5 Points)** The approach is the same as the parking lot example and the $M(t)/M/\infty$ queue. The piecewise-constant arrival rate function $\lambda(t)$ can be used together with the thinning method to generate the NSPP. Note that the inter-arrival times are no longer exponentially distributed and hence generating exponential inter-arrival times with time-dependent rate is NOT a valid approach. The service times are i.i.d exponential with mean 45 minutes. The performance measure of interest is the 0.95th quantile of the maximum number of calls reporting cardiac arrests during the day.

(b) **(15 Points)** We simulate 1000 samples of the maximum number of calls during the day, sort them, and pick the $950^{\text{th}}$ value as the minimum number of drones required at the base for the desired service level. The estimated required drones is 220 with the 95% CI of $[Y_{(936)}, Y_{(964)}] = [218, 221]$.

```
1   # Drones
2   import SimClasses
3   import SimFunctions
4   import SimRNG
5   import math
6   import pandas
7   import numpy as np
8
9   ZSimRNG = SimRNG.InitializeRNSeed()
10  Calendar = SimClasses.EventCalendar()
11  TheCTStats = []
12  TheDTStats = []
13  TheQueues = []
14  TheResources = []
```

```
15  AllMaxQueue = []
16
17  MeanBusyTime = 45.0
18  RepNum = 1000
19
20  ArrivalRates = [1, 4, 2]
21  MaxRate = 4
22
23  # Picewise-constant arrival rate function
24  def PW_ArrRate(t):
25      if t < 7*60:        # 12AM–7AM
26          return ArrivalRates[0]
27      elif t < 12*60:        # 7AM–12PM
28          return ArrivalRates[1]
29      else:        # 12PM–12AM
30          return ArrivalRates[2]
31
32  # Thinning method used to generate NSPP with the piecewise-constant arrival
        rate
33  def NSPP(Stream):
34      PossibleArrival = SimClasses.Clock + SimRNG.Expon(1/MaxRate, Stream)
35      while SimRNG.Uniform(0, 1, Stream) >= PW_ArrRate(PossibleArrival)/MaxRate:
36          PossibleArrival = PossibleArrival + SimRNG.Expon(1/MaxRate, Stream)
37      nspp = PossibleArrival - SimClasses.Clock
38      return nspp
39
40  def Arrival():
41      global MaxQueue
42      global N
43      interarrival = NSPP(1)
44      SimFunctions.Schedule(Calendar, "Arrival", interarrival)
45      N = N + 1
46      if N > MaxQueue:
47          MaxQueue = N
48      SimFunctions.Schedule(Calendar, "EndOfService", SimRNG.Expon(MeanBusyTime,
          2))
49
50  def EndOfService():
51      global N
52      N = N - 1
53
54  for Reps in range(0, RepNum, 1):
55      N = 0
56      MaxQueue = 0
57
58      SimFunctions.SimFunctionsInit(Calendar, TheQueues, TheCTStats, TheDTStats,
            TheResources)
59      interarrival = NSPP(1)
60      SimFunctions.Schedule(Calendar, "Arrival", interarrival)
61      SimFunctions.Schedule(Calendar, "EndSimulation", 24*60)
62
63      NextEvent = Calendar.Remove()
64      SimClasses.Clock = NextEvent.EventTime
65      if NextEvent.EventType == "Arrival":
```

11

```
66                 Arrival()
67          elif NextEvent.EventType == "EndOfService":
68                 EndOfService()
69
70          while NextEvent.EventType != "EndSimulation":
71              NextEvent = Calendar.Remove()
72              SimClasses.Clock = NextEvent.EventTime
73              if NextEvent.EventType == "Arrival":
74                  Arrival()
75              elif NextEvent.EventType == "EndOfService":
76                  EndOfService()
77
78          AllMaxQueue.append(MaxQueue)
79
80  # estimating the 0.95 quantile
81  print("Estimated_required_drones_is:", np.sort(AllMaxQueue)[950])
82  l = int(np.floor(950 - 1.96*np.sqrt(1000*0.95*0.05)))
83  u = int(np.ceil(950 + 1.96*np.sqrt(1000*0.95*0.05)))
84  print("The_95%_CI:_[{},_{}]".format(np.sort(AllMaxQueue)[l], np.sort(
        AllMaxQueue)[u]))
```

---

```
Estimated required drones is:   220
The 95% CI: [218, 221]
```

**Problem 3. (20 Pts.)** The attached file (TSLA.csv) provides daily prices for Tesla's stock from March 2022 - March 2023.

**(a)** Fit a Geometric Brownian Motion (GBM) to the **Close** price by estimating the drift $\mu$ and volatility terms $\sigma$. Recall that for GBM the log-returns are independent Normal random variables, i.e.,

$$\log\left(\frac{S(t_i)}{S(t_{i-1})}\right) \sim N(\mu(t_i - t_{i-1}), \sigma(t_i - t_{i-1})),$$

for all $t_i, i = 1, \ldots, k$.

**(b)** Does the GBM fit the data well? Use your estimates from part (a) together with a Q-Q plot and statistical goodness of fit tests discussed in class to support your answer.

(a) **(12 Points)** We first calculate the log-returns using the **Close** price data. We then fit a normal distribution to the log-returns data and estimate $\mu = -0.001518507$ and $\sigma = 0.042617534$.

---

```
1  # This file uses the fitdistrplus package. You'll need to install it first
2  # by running the command: 'install.packages("fitdistrplus")' in R
3  require(fitdistrplus)
4
5  # import the data from .csv file
6  Data <- read.csv("TSLA.csv", header=TRUE,
7                      sep=",", na.strings=c("NA", ""), stringsAsFactors=FALSE, as
                       .is=TRUE)
8  Data$LogReturn = 0
9  for (i in 2:nrow(Data)){
10   Data$LogReturn[i] = log(Data$Close[i]/Data$Close[i-1])
11 }
12
```

```
13   MyData <- Data$LogReturn
14
15   # fit a normal distribution to the log-returns
16   fln <- fitdist(MyData, "norm")
17   fln$estimate
18
19
20   # Graphical method
21   par(mfrow = c(1, 2))
22   plot.legend <- c("normal")
23   denscomp(fln, legendtext = plot.legend)
24   qqcomp(fln, legendtext = plot.legend)
25
26   # Perform GofFit tests
27   gof_results <- gofstat(fln)
28   gof_results
29   gof_results$kstest
30
31   gof_results$chisq
32   gof_results$chisqpvalue
33   gof_results$chisqtable
```

|    mean     |    sd      |
|-------------|------------|
| -0.001518507 | 0.042617534 |

(b) **(8 Points)** Yes, for this data set the Normal distribution seems appropriate. As shown in Fig. 5, both the Q-Q plot and histogram suggest that the Normal distribution is a good fit for the log-returns. Note that having larger deviation around the tails on the Q-Q plot is expected since the variance is higher at extreme quantiles. Also, the KS test does not reject the Normal distribution at $(1 - \alpha) = 95\%$ confidence level.

**Problem 4. (15 Pts.)** Normal distribution is sometimes not a good fit for log-returns. Cauchy distribution is another (heavy-tailed) distribution used to model the log-returns. The CDF of the Cauchy distribution is given by:

$$F(x) = \frac{1}{\pi} \arctan\left(\frac{x - x_0}{\gamma}\right) + \frac{1}{2},$$

where $x_0$ and $\gamma$ are parameters.

**(a)** Propose an inversion algorithm to generate samples of a Cauchy distribution with $x_0 = 0$ and $\gamma = 2$.

**(b)** Estimating the parameters of the Cauchy distribution using MLE is challenging. One approach to estimated the parameters using matching is to match median (1/2th quantile) and half of the inter-quartile range (difference between 1/4th and 3/4th quantiles) with their estimates from the data. Use this method to estimate the parameters using 1000 samples generated from part (a). How do they compare with the original parameters, i.e., $x_0 = 0$ and $\gamma = 2$?

(a) **(7 Points)** We denote the inverse cdf by $F^{-1}(.)$. So,

$$F(x) = u \iff x = F^{-1}(u).$$

13

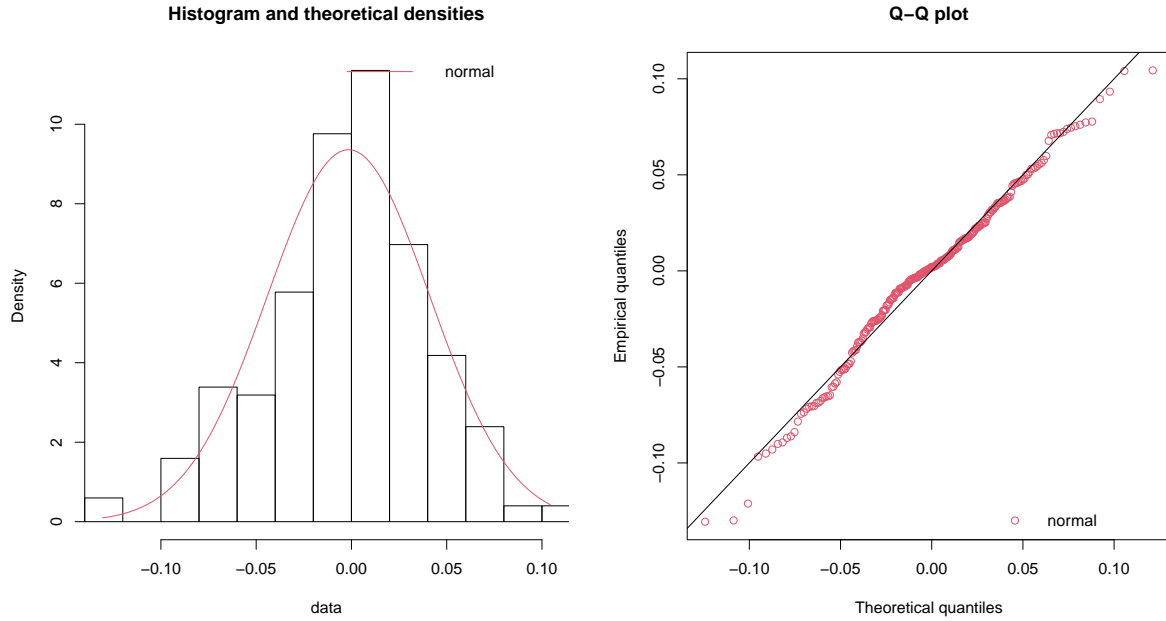Figure 5: Graphical support for the goodness of fit of GBM.

We have

$$\frac{1}{\pi}\arctan\left(\frac{x}{2}\right)+\frac{1}{2}=u \implies \arctan\left(\frac{x}{2}\right)=\pi(u-\frac{1}{2}) \implies x = 2\tan\left(\pi\left(u-\frac{1}{2}\right)\right)$$

Therefore,

$$F^{-1}(u)=2\tan\left(\pi\left(u-\frac{1}{2}\right)\right)$$

So, the inversion algorithm works as follows:

1. Generate $U \sim U[0,1]$.

2. Set $X = 2\tan\left(\pi\left(U-\frac{1}{2}\right)\right)$ and return $X$.

(b) **(8 Points)** We generate 1000 samples of the Cauchy distribution with $x_0 = 0$ and $\gamma = 2$ and estimate median and half of the inter-quartile range to be 0.05 and 3.98, respectively.

For the Cauchy distribution, we can calculate the q-th quantile using $x_0 + \gamma\tan\left(\pi(q-1/2)\right)$. So, matching median and half of the inter-quartile range with their estimates from the data, we have

$$x_0 + \gamma\tan(0) = 0.05 \implies x_0 = 0.05,$$

$$\gamma\tan\left(\frac{\pi}{4}\right) - \gamma\tan\left(-\frac{\pi}{4}\right) = 3.98 \implies \gamma = 1.99.$$

We observe that the estimated parameters are close enough to the original parameters.

```
1  #  Cauchy distribution
2  import numpy as np
3
4  n = 1000
5  X_list = []
6  np.random.seed(1)
7  for i in range(n):
8      U = np.random.random()
9      X = 2*np.tan(np.pi*(U-1/2))
10     X_list.append(X)
11
12 print("Estimate_of_median:", np.median(X_list))
13 print("Estimate_of_difference_between_1/4th_and_3/4th_quantiles:", np.quantile
       (X_list, 3/4)-np.quantile(X_list, 1/4))
```

```
Estimate of median:  0.04713779707427136
Estimate of difference between 1/4th and 3/4th quantiles:  3.9842903109340155
```

**Problem 5**. **(20 Pts.)** Data on call counts to a call centre, by hour, for 1 month (31 days) are provided in CallCounts.xls. Let $N(t)$ represent the cumulative number of arrivals by time $t$. If the process is nonstationary Poisson, then $Var(N(t))/E(N(t)) = 1$ for all $t$, or stated differently $Var(N(t)) = \beta E(N(t))$ with $\beta = 1$. Since you have arrival count data, you can estimate $Var(N(t))$ and $E(N(t))$ at $t = 1, 2, \ldots, 8h$. Use these data to fit the regression model $Var(N(t_i)) = \beta E(N(t_i))$ and see if the estimated value of $\beta$ supports the choice of a nonstationary Poisson arrival process. **Hints**: This is regression through the origin. Also, remember that $N(t_i)$ represents the total number of arrivals by time $t_i$.

**(20 Points)** We follow these steps to estimate $\beta$:

1. Transform hourly count data into cumulative count data by adding the number of calls in each hour to that of the previous hour.

2. For each hour $t_i, i \in \{1, 2, \ldots, 8\}$, estimate the expected cumulative arrivals up to that hour by computing sample average $\bar{N}(t_i)$ across the 31 days, i.e.,

$$\bar{N}(t_i) = \frac{1}{31} \sum_{j=1}^{31} N_j(t_i), \quad \forall i \in \{1, 2, \ldots, 8\},$$

where $N_j(t_i)$ is the cumulative arrivals by hour $t_i$ on day $j$.

3. Similarly, compute the sample variance $V(t_i)$ for each hour $t_i$, i.e.,

$$V(t_i) = \frac{1}{30} \sum_{j=1}^{31} (N_j(t_i) - \bar{N}(t_i))^2, \quad \forall i \in \{1, 2, \ldots, 8\}.$$

4. Perform regression of $V(t_i)$ on $\bar{N}(t_i)$ with zero intercept (8 data points) to estimate $\beta$.

```r
1  Data <- read.csv("CallCounts.csv", header=TRUE,
2                   sep=",", na.strings=c("NA", ""), stringsAsFactors=FALSE, as.
                        is=TRUE)
3
4  # Create a new data frame with the same columns as Data
5  df <- data.frame(matrix(0, nrow=nrow(Data), ncol=ncol(Data)))
6
7  # Calculate cumulative sums for each column
8  for (i in 1:ncol(Data)) {
9    df[,i] <- Data[,i]
10   if (i > 1) {
11     df[,i] <- df[,i] + df[,i-1]
12   }
13 }
14
15 # Compute sample averages
16 sa <- apply(df, MARGIN = 2, mean)
17
18 # Compute sample variances
19 sv <- apply(df, MARGIN = 2, var)
20
21 # Perform linear regression and output 95% CI
22 df2 <- data.frame(sample_average = sa, sample_variance = sv)
23 df2.lm <- lm(sample_variance ~ sample_average - 1, data = df2) # zero
        intercept
24 summary(df2.lm)
25 confint(df2.lm, level = 0.95)
```

```
    2.5%       97.5%
 1.069191    1.368783
```

We find that the estimated coefficient is 1.22 with the 95% CI of $[1.07, 1.37]$. This is reasonably close to 1 which supports the use of nonstationary Poisson process. (You may also argue that based on the %95 CI the variability is larger than that of Poisson and hence a more variable arrival process would be more appropriate.)