

Simulations

Homework 4

M.S. in Stochastics and Data Science

Andrea Crusi, Lorenzo Sala

December 2, 2024

Exercise 1

1. To obtain the explicit expression for the distribution we need to substitute $m = 2$ in the probability distribution formula:

$$p(n) = \begin{cases} p(0) \left(\frac{\lambda}{\mu}\right)^n \frac{1}{n!} & \text{for } 0 \leq n \leq 2 \\ p(0) \left(\frac{\lambda}{\mu}\right)^n \frac{1}{2!2^{n-2}} = p(0) \left(\frac{\lambda}{\mu}\right)^n \frac{1}{2^{n-1}} & \text{for } n > 2. \end{cases}$$

2. We can use the fact that

$$\sum_{n=0}^N p(n) = 1$$

to find an analytical solution for $p(0)$. A simple way to do this is to check how $p(n)$ behaves for $N \rightarrow \infty$.

$$\sum_{n=0}^2 p(0) \left(\frac{\lambda}{\mu}\right)^n \frac{1}{n!} + \sum_{n=3}^{\infty} p(0) \left(\frac{\lambda}{\mu}\right)^n \frac{1}{2!2^{n-2}} = 1.$$

By explicitly computing the first sum and rearranging the second one so that it starts from

$n = 0$ we can get

$$\begin{aligned}
p(0) & \left[1 + \frac{\lambda}{\mu} + \frac{\lambda^2}{2\mu^2} + \sum_{n=3}^{\infty} \left(\frac{\lambda}{\mu} \right)^n \frac{1}{2^{n-1}} \right] = 1 \\
\Rightarrow p(0) & = \frac{1}{1 + \frac{\lambda}{\mu} + \frac{\lambda^2}{2\mu^2} + \sum_{n=3}^{\infty} \left(\frac{\lambda}{\mu} \right)^n \frac{1}{2^{n-1}}} = \\
& = \frac{1}{1 + \frac{\lambda}{\mu} + \frac{1}{2} \left(\frac{\lambda}{\mu} \right)^2 + \sum_{n=3}^{\infty} \left(\frac{\lambda}{\mu} \right)^n \frac{1}{2^{n-1}}} = \\
& = \frac{1}{1 + \frac{\lambda}{\mu} \left[1 + \frac{\frac{\lambda}{\mu}}{2} + \sum_{n=3}^{\infty} \left(\frac{\frac{\lambda}{\mu}}{2} \right)^{n-1} \right]} = \\
& = \frac{1}{1 + \frac{\lambda}{\mu} \left[\underbrace{\sum_{n=0}^{\infty} \left(\frac{\frac{\lambda}{\mu}}{2} \right)^{n-1}}_{\text{geometric series}} \right]} = \\
& = \frac{1}{1 + \frac{\lambda}{\mu} \left[\frac{1}{1 - \frac{\frac{\lambda}{\mu}}{2}} \right]} = \\
& = \frac{1}{1 + \frac{\lambda}{\mu} \left[\frac{2}{2 - \frac{\lambda}{\mu}} \right]} = \\
& = \frac{2 - \frac{\lambda}{\mu}}{2 + \frac{\lambda}{\mu}}.
\end{aligned}$$

So our new distribution becomes

$$\begin{aligned}
p(n) & = \begin{cases} \frac{2 - \frac{\lambda}{\mu}}{2 + \frac{\lambda}{\mu}} \left(\frac{\lambda}{\mu} \right)^n \frac{1}{n!} & \text{for } 0 \leq n \leq 2 \\ \frac{2 - \frac{\lambda}{\mu}}{2 + \frac{\lambda}{\mu}} \left(\frac{\lambda}{\mu} \right)^n \frac{1}{2^{n-1}} & \text{for } n > 2 \end{cases} \\
& = \begin{cases} \frac{2 - \frac{\lambda}{\mu}}{2 + \frac{\lambda}{\mu}} \left(\frac{\lambda}{\mu} \right)^n \frac{1}{n!} & \text{for } 0 \leq n \leq 2 \\ \frac{2 - \frac{\lambda}{\mu}}{2 + \frac{\lambda}{\mu}} \left(\frac{\lambda}{\mu} \right)^n \frac{1}{2^{n-1}} & \text{for } n > 2 \end{cases}
\end{aligned}$$

which can also be expressed, by denoting $\frac{\lambda}{\mu} = \rho$ as the load of the server, as

$$p(n) = \begin{cases} \frac{2-\rho}{2+\rho} (\rho)^n \frac{1}{n!} & \text{for } 0 \leq n \leq 2 \\ \frac{2-\rho}{2+\rho} (\rho)^n \frac{1}{2^{n-1}} & \text{for } n > 2. \end{cases}$$

Exercise 2

We know that when a network comprises load independent stations the recursive equation for utilization in a station i is

$$\bar{n}_i(n) = U_i[1 + \bar{n}_i(n-1)].$$

We are interested in what happens when n goes to ∞ so

$$\bar{n}_i = U_i[1 + \bar{n}_i] \implies \bar{n}_i - U_i\bar{n}_i = U_i \implies \bar{n}_i = \frac{U_i}{1 - U_i}.$$

In bottleneck stations b , where the utilization U_b tends to 1 due to the saturation of the system, the queue length will tend to infinity. Moreover, we know that, by consistency laws,

$$\frac{U_b}{U_i} = \frac{V_b S_b}{V_i S_i} = \frac{D_b}{D_i} \implies U_i = \frac{D_i U_b}{D_b}$$

and therefore

$$\bar{n}_i = \frac{\frac{D_i U_b}{D_b}}{\frac{D_b - D_i U_b}{D_b}} = \frac{D_i U_b}{D_b - D_i U_b}$$

but since as load increases $U_b \rightarrow 1$ we get

$$\bar{n}_i = \frac{D_i}{D_b - D_i}$$

for non-bottleneck stations.

Let's turn to average waiting time. Little's Law tells us that

$$\bar{n}_i = \bar{w}_i X_i \implies \bar{w}_i = \frac{\bar{n}_i}{X_i}$$

but we know that as the load increases $\bar{n}_i \rightarrow U_i[1 + \bar{n}_i]$ so

$$\bar{w}_i = \frac{U_i[1 + \bar{n}_i]}{X_i}.$$

In bottleneck stations, as the queue goes to ∞ , so does the average waiting time. Remember that for load independent stations we have $U_i = X_i(n) \cdot S_i$, so

$$\bar{w}_i = \frac{\cancel{X_i(n)} S_i [1 + \bar{n}_i]}{\cancel{X_i(n)}} = S_i [1 + \bar{n}_i].$$

Now let's take into account the throughput. As we said before, by Little's Formula we know that

$$\overline{n_i} = \overline{w_i} X_i \implies X_i(n) = \frac{\overline{n_i}}{\overline{w_i}}$$

and by substituting we get

$$X_i(n) = \frac{\overline{n_i}}{S_i[1 + n_i]}.$$

In bottleneck stations, as $\overline{n_i}$ becomes increasingly larger, $[1 + \overline{n_i}]$ becomes closer and closer to $\overline{n_i}$ so that the throughput for bottleneck stations will tend to

$$\frac{1}{S_i}.$$

Exercise 3

The normalization constant is defined in two ways:

$$\underbrace{g(n, m) = \sum_{\mathbf{n} \in S(n, m)} \prod_{i=1}^m f_i(n_i)}_{\text{definition}} \quad \underbrace{\sum_{k=0}^n f_m(k) g(n - k, m - 1)}_{\text{convolution method}}$$

We know that the distribution of $p_i(k)$, which is the fraction of time the station i spends with k customers inside, can be computed as

$$\begin{aligned} p_i(k) &= \sum_{\substack{\mathbf{n} \in S(n, i) \\ n_i = k}} P(\mathbf{n}) \\ &= \sum_{\substack{\mathbf{n} \in S(n, i) \\ n_i = k}} \frac{1}{g(n, i)} \prod_{j=1}^M f_j(n_j) \\ &= \frac{1}{g(n, i)} \sum_{\substack{\mathbf{n} \in S(n, i) \\ n_i = k}} \prod_{j=1}^M f_j(n_j) \\ &= \frac{f_i(k)}{g(n, i)} \underbrace{\sum_{\mathbf{n} \in S(n-k, i-1)} \prod_{j=1}^{i-1} f_j(n_j)}_{g(n-k, i-1)} \quad \text{since } n_i = k \text{ for all states we can take it out of the sum} \\ &= \frac{f_i(k) g(n-k, i-1)}{g(n, i)} \end{aligned}$$

To get $p_i(k, N)$ we must sum the probabilities for station i to have k customers over all possible distribution in the other station with a maximum of N customers. This happens over the reduced state space

$$S^{[-i]}(N, M)$$

(which is basically the whole state space of N maximum clients and M stations but with the number of customers in i -th station fixed to 0) and with the reduced normalization constant with the i -th station missing over the reduced state space

$$g^{[-i]}(N, M) = \sum_{\mathbf{n} \in S^{[-i]}(N, M)} \prod_{j \neq i} f_j(n_j).$$

So, in a similar manner as before, we get

$$\begin{aligned} p_i(k, N) &= \sum_{\mathbf{n} \in S^{[-i]}(N, i)} P(\mathbf{n}) \\ &= \frac{f_i(k)}{g(N, i)} \underbrace{\sum_{\mathbf{n} \in S^{[-i]}(N-k, i-1)} \prod_{j=1}^{i-1} f_j(n_j)}_{g^{[-i]}(N-k, i-1)} \\ &= \frac{f_i(k) g^{[-i]}(N-k, i-1)}{g(N, i)}. \end{aligned}$$

The service function $f_i(k)$ is defined as

$$f_i(k) = \begin{cases} 1 & k = 0 \\ V_i S_i f_i(k-1) & k > 0. \end{cases}$$

This allows us to write

$$\begin{aligned} p_i(k, N) &= \frac{f_i(k) g^{[-i]}(N-k, i-1)}{g(N, i)} \\ &= V_i S_i f_i(k-1) \cdot \frac{g^{[-i]}(N-k, i-1)}{g(N, i)} \\ &= V_i S_i f_i(k-1) \cdot \frac{g^{[-i]}(N-k, i-1)}{g(N, i)} \cdot \frac{g(N-1, i)}{g(N-1, i)} \\ &= \underbrace{V_i \frac{g(N-1, i)}{g(N, i)}}_{X_i(N)} \cdot S_i \cdot \underbrace{\frac{f_i(k-1) g^{[-i]}(\overbrace{(N-1) - (k-1)}^{N-k}, i-1)}{g(N-1, i)}}_{p_i(k-1, N-1)} \\ &= X_i(N) S_i p_i(k-1, N-1). \end{aligned}$$

Exercise 4

This code computes both the MVA computations and the bottleneck analysis computations.

```

1  /* -----
2  Università di Torion
3  M.S. in STOCHASTICS AND DATA SCIENCE
4  Course in Simulation
5  Homework 4
6
7  By Andrea Crusi and Lorenzo Sala
8  * -----
9  */
10 #include <stdio.h>
11
12 // Constants
13 #define M 4      // Number of stations
14 #define N 80     // Maximum number of customers
15
16 // Input parameters
17 double Z = 10.0; // Delay time of the station
18 double S[M] = {0, 0.04, 0.06, 0.04}; // Service times
19 char ST[M] = {'D', 'L', 'L', 'L',}; // Station types ('L' for load independent, 'D' for delay.
20 ↪ Our reference station is the delay)
21 double Q[M][M] = {
22     {0, 1, 0, 0},
23     {0.1, 0, 0.55, 0.35},
24     {0, 1, 0, 0},
25     {0, 1, 0, 0}
26 };
27 double V[M];
28
29 // We need bi-dimensional arrays for results
30 double X[M][N+1];
31 double U[M][N+1];
32 double n[M][N+1];
33 double w[M][N+1];
34 double R[M][N+1];
35
36 int main() {
37     // Initialization
38     for (int i = 0; i < M; i++) {
39         n[i][0] = 0.0; // put 0 as the value for all stations
40     }
41
42     // Compute visit count
43     V[0]=1.0;
44     V[1]=V[0]/Q[1][0];
45     V[2]=V[1]*Q[1][2];
46     V[3]=V[1]*Q[1][3];
47
48     // Compute performance measures for each population size k (from 1 to N)
49     for (int k = 1; k <= N; k++) {
50
51         // Compute the waiting time w_i[k] for each station i
52         for (int i = 0; i < M; i++) {
53             if (ST[i] == 'D') {
54                 w[i][k] = Z; // Delay station
55             } else {
56                 w[i][k] = S[i] * (1 + n[i][k - 1]); // Queue station
57             }
58         }
59         double sum = 0.0; // initialize sum
60
61         // Compute the sum of V_i * w_i[k] across all stations
62         for (int i = 0; i < M; i++) {
63

```

```

64         sum += V[i] * w[i][k];
65     }
66
67     // Compute throughput for reference job. Remember that the reference station is
68     ↪ the station 0 (delay station)
69     double Xref = k / sum;
70
71     // Compute performance metrics for each station i using MVA equations seen in
72     ↪ class
73     for (int i = 0; i < M; i++) {
74         X[i][k] = V[i] * Xref; // Throughput for station i
75
76         if (ST[i] == 'D') {
77             // Delay station
78             n[i][k] = Z * X[i][k];
79             U[i][k] = n[i][k] / k;
80         } else {
81             // Computational station
82             U[i][k] = S[i] * X[i][k]; // Utilization
83             n[i][k] = U[i][k] * (1 + n[i][k - 1]); // Average queue length
84         }
85     }
86     // Compute the response time
87     double Y[M]={0.0,0.0,0.0,0.0};
88     for (int i=0; i < M; i++) {
89         for (int i = 0; i < M; i++) {
90             Y[i] += V[i] * w[i][k];
91         }
92         if (ST[i] == 'D') {
93             R[i][k] = (k/X[i][k])-Z; //formula for delay station
94         } else {
95             R[i][k] = w[i][N]+S[i];
96         }
97     }
98 }
99 int n_customers[2] = {1,N};
100 // Print results for N = 1 and N = 80;
101 for (int j=0; j<=1; j++){
102     printf("*****\n\n");
103     printf("Simulation with %d customers\n\n", n_customers[j]);
104     for (int i = 0; i < M; i++) {
105         printf("Station %d results:\n", i);
106         printf("Throughput (X[%d])\t\t= %f\n", n_customers[j],
107             ↪ X[i][n_customers[j]]);
108         printf("Utilization (U[%d])\t\t= %f\n", n_customers[j],
109             ↪ U[i][n_customers[j]]);
110         printf("Mean queue length (n[%d])\t= %f\n", n_customers[j],
111             ↪ n[i][n_customers[j]]);
112         printf("Mean waiting time (w[%d])\t= %f\n", n_customers[j],
113             ↪ w[i][n_customers[j]]);
114         //printf("Mean response time (R_0)\t\t= %f\n",
115             ↪ (n_customers[j]/X[i][n_customers[j]])-Z);
116         printf("Mean response time (R_0)\t\t= %f\n", (R[i][n_customers
117             ↪ [j]]));
118         printf("\n");
119     }
120 }
121
122 double X_0_1=X[0][1]; // We will need this later to compute the saturation point
123
124 // We now make a pair of csv files that we will plot with pgfplots in the homework
125 ↪ discussion
126 FILE *file = fopen("x0.csv", "w");

```

```

120     if (file == NULL) {
121         fprintf(stderr, "Error opening file for writing.\n");
122         return 1;
123     }
124
125     for (int k = 1; k <= N; k++) {
126         fprintf(file, "%d %f\n", k, X[0][k]);
127     }
128
129     fclose(file);
130     printf("Data written to x0.csv successfully.\n");
131
132     // We now save the results of the response time for station 0
133     FILE *file2 = fopen("R0.csv", "w");
134     if (file2 == NULL) {
135         fprintf(stderr, "Error opening file for writing.\n");
136         return 1;
137     }
138
139     for (int k = 1; k <= N; k++) {
140         fprintf(file2, "%d %f\n", k, R[0][k]);
141     }
142
143     fclose(file2);
144     printf("Data written to R0.csv successfully.\n");
145
146     // We now save the results of average queue length
147     FILE *file3 = fopen("queue.csv", "w");
148     if (file3 == NULL) {
149         fprintf(stderr, "Error opening file for writing.\n");
150         return 1;
151     }
152
153     fprintf(file3, "n station0 station1 station2 station3\n");
154
155     for (int k = 1; k <= N; k++) {
156         fprintf(file3, "%d %f %f %f %f\n", k, n[0][k], n[1][k], n[2][k], n[3][k]);
157     }
158
159     fclose(file3);
160     printf("Data written to queue.csv successfully.\n");
161
162     // We now save the results of utilization
163     FILE *file4 = fopen("utilization.csv", "w");
164     if (file4 == NULL) {
165         fprintf(stderr, "Error opening file for writing.\n");
166         return 1;
167     }
168
169     fprintf(file4, "n station0 station1 station2 station3\n");
170
171     for (int k = 1; k <= N; k++) {
172         fprintf(file4, "%d %f %f %f %f\n", k, U[0][k], U[1][k], U[2][k], U[3][k]);
173     }
174
175     fclose(file4);
176     printf("Data written to utilization.csv successfully.\n");
177
178     // Now we perform the bottleneck analysis using the equations derived in the previous
179     ↪ exercise.
180     printf("\n");
181
182     double D[M];           // Service demands for each station
183     double Xb;             // Variable to store the max throughput of the bottleneck station

```



```

183     double X[M];           // Throughput
184     double U[M];           // Utilization for each station
185     double n[M];           // Mean queue length for each station
186     double w[M];           // Mean waiting time for each station
187
188     double Db = 0.0;
189     int bottleneck_station = 0;
190
191     // Calculate service demand for each station and identify the bottleneck as the station
192     // ↳ with higher service demand
193     for (int i = 0; i < M; i++) {
194         D[i] = V[i] * S[i];
195         if (D[i] > Db) {
196             Db = D[i];
197             bottleneck_station = i;
198         }
199     }
200
201     // Maximum throughput as N -> infinity: we can use this formula because we are in a
202     // ↳ closed system
203     Xb = 1.0 / Db;
204
205     // Calculate utilization, mean queue length, and mean waiting time for each station
206     for (int i = 0; i < M; i++) {
207         X[i] = V[i]*Xb; //our reference station is the bottleneck station
208         U[i] = X[i] * S[i];
209
210         if (ST[i] == 'D') {
211             // For the delay station
212             n[i] = Z * Xb;
213             w[i] = Z;
214         } else {
215             // For load independent station stations
216             n[i] = U[i] / (1 - U[i]); // Mean queue length
217             w[i] = n[i] / Xb; // Mean waiting time
218         }
219     }
220
221     // Print the bottleneck analysis results
222     printf("=====\n");
223     printf("Bottleneck Analysis (using formulas of the previous exercise)\n");
224     printf("=====\n");
225     printf("The bottleneck station is %d, with the highest service demand that is %f\n",
226     // ↳ bottleneck_station, Db);
227     printf("Maximum Throughput (Xb) as N -> infinity: %f\n\n", Xb);
228
229     // Print results for each station
230     for (int i = 0; i < M; i++) {
231         printf("Station %d:\n", i);
232         printf("  Service Demand (D_%d(N))      = %f\n", i, D[i]);
233         printf("  Throughput (X_%d(N))                    = %f\n", i, X[i]);
234         printf("  Utilization (U_%d(N))                    = %f\n", i, U[i]);
235         printf("  Mean Queue Length (n_%d(N))              = %f\n", i, n[i]);
236         printf("  Mean Waiting Time (w_%d(N))              = %f\n\n", i, w[i]);
237     }
238     printf(".....\n");
239     printf("For plotting purposes:\n");
240     printf("Vb*Sb=Db=%f\n", Db);
241     printf("Point of saturation N_star: %f clients", (1/(X_0_1*Db)));
242     return 0;
243 }

```

This prints the following results for the simulation with $N = 1$:

	Station 0	Station 1	Station 2	Station 3
Throughput ($X_i(1)$)	0.091996	0.919963	0.505980	0.321987
Utilization ($U_i(1)$)	0.919963	0.036799	0.030359	0.012879
Mean queue length ($\bar{n}_i(1)$)	0.919963	0.036799	0.030359	0.012879
Mean waiting time ($\bar{w}_i(1)$)	10.000000	0.040000	0.060000	0.040000
Mean response time (R_0)	0.870000	0.870000	0.870000	0.870000

For $N = 80$ we have:

	Station 0	Station 1	Station 2	Station 3
Throughput ($X_i(80)$)	2.499979	24.999794	13.749887	8.749928
Utilization ($U_i(80)$)	0.312497	0.999992	0.824993	0.349997
Mean queue length ($\bar{n}_i(80)$)	24.999794	49.749403	4.712349	0.538454
Mean waiting time ($\bar{w}_i(80)$)	10.000000	1.989992	0.342719	0.061538
Mean response time (R_0)	22.000263	22.000263	22.000263	22.000263

Finally, for the bottleneck analysis (where we have $N \rightarrow \infty$) we obtain:

	Station 0	Station 1	Station 2	Station 3
Service demand ($D_i(N)$)	0	0.4	0.33	0.14
Throughput ($X_i(N)$)	2.5	25	13.75	8.75
Utilization ($U_i(N)$)	0	1	0.825	0.35
Mean queue length ($\bar{n}_i(N)$)	25	∞	4.714286	0.538462
Mean waiting time ($\bar{w}_i(N)$)	10	∞	1.885714	0.215385

We know that the bottleneck of the system is station 1, since it has the highest service demand. We can see how results for $N = 80$ for the MVA computations are consistent with the bottleneck analysis: the mean queue length \bar{n}_1 tends to ∞ as N grows while the other values of performance tend to their respective value in the bottleneck analysis.

We now plot the data gained from the simulation along with the expected asymptotes. The value $N^* = 27.175000$ is obtained from the simulation, being simply the point of intersection of the two asymptotes in either case. In other words, this is what we would expect if our calculations are correct.

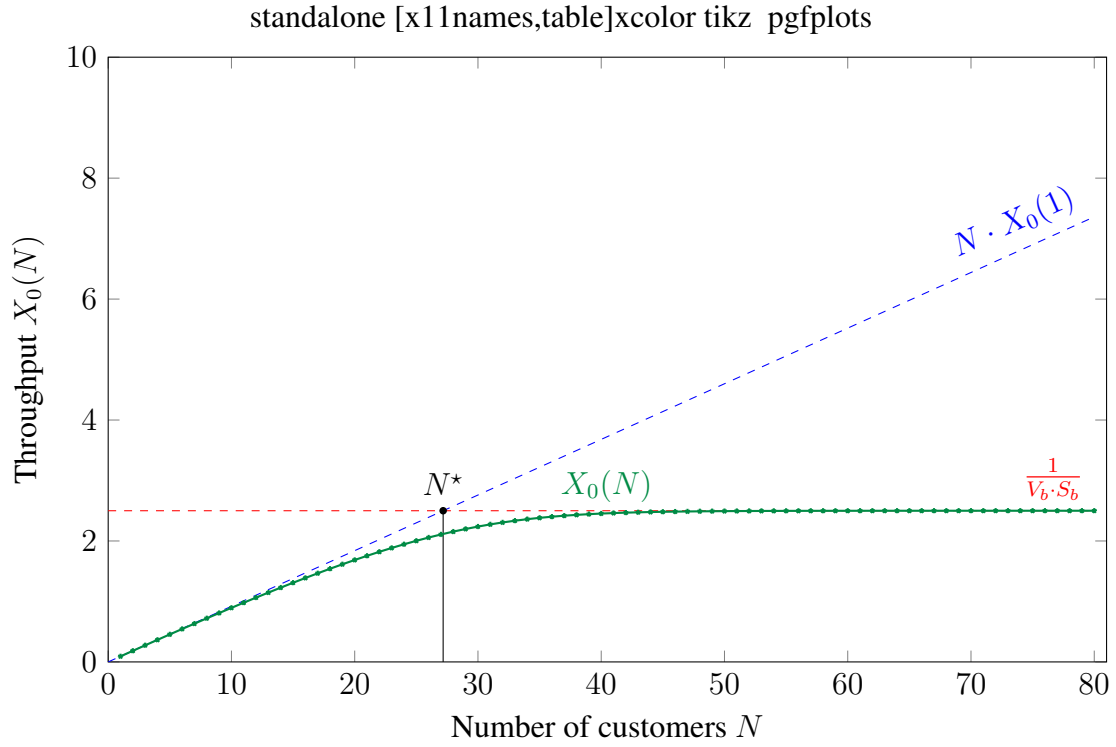


Figure 1: Evolution of throughput $X_0(N)$ as N increases. The point of saturation is $N^* = 27.175000$. We can see how the throughput reaches a plateau after the saturation point.

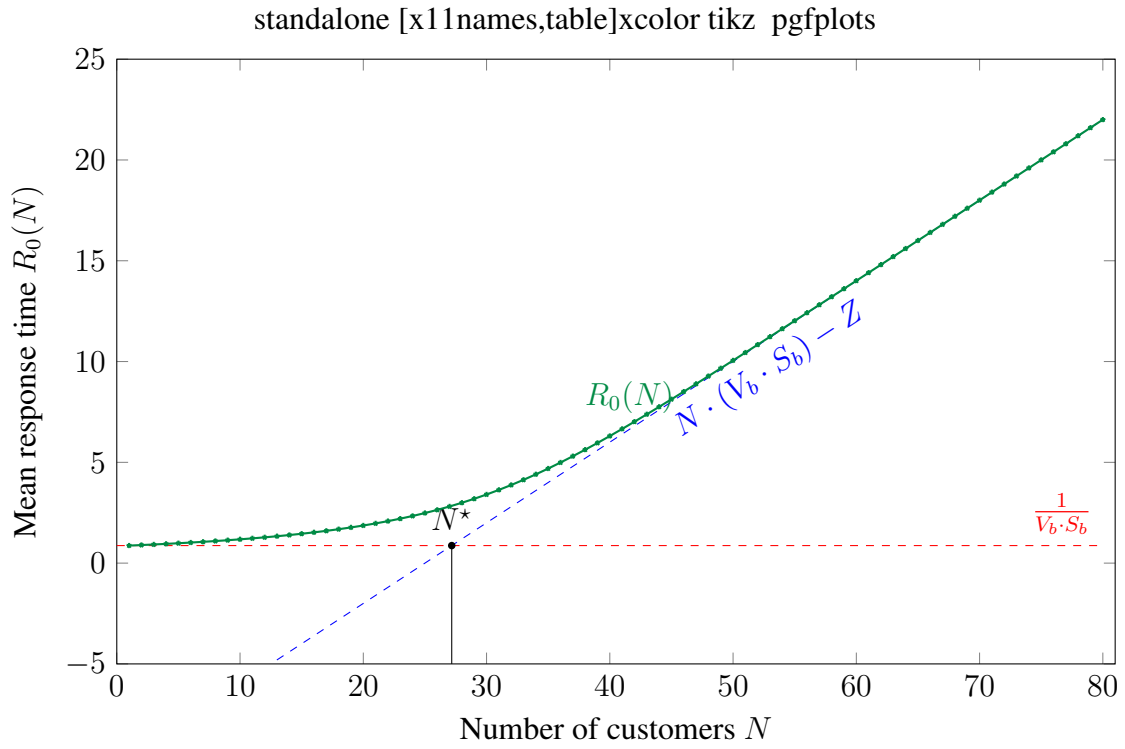


Figure 2: Evolution of mean response time $R_0(N)$ as N increases. The point of saturation is $N^* = 27.175000$. On the contrary here the response time explodes after the saturation point

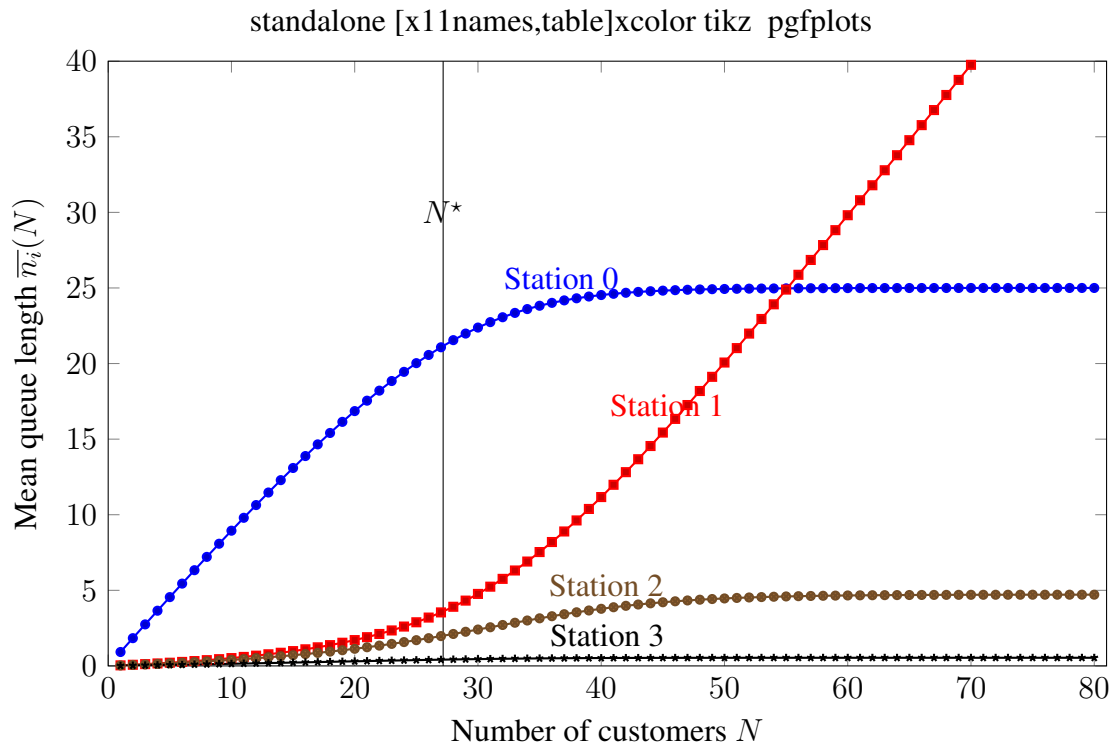


Figure 3: Behaviour of the average queue length of the four stations.

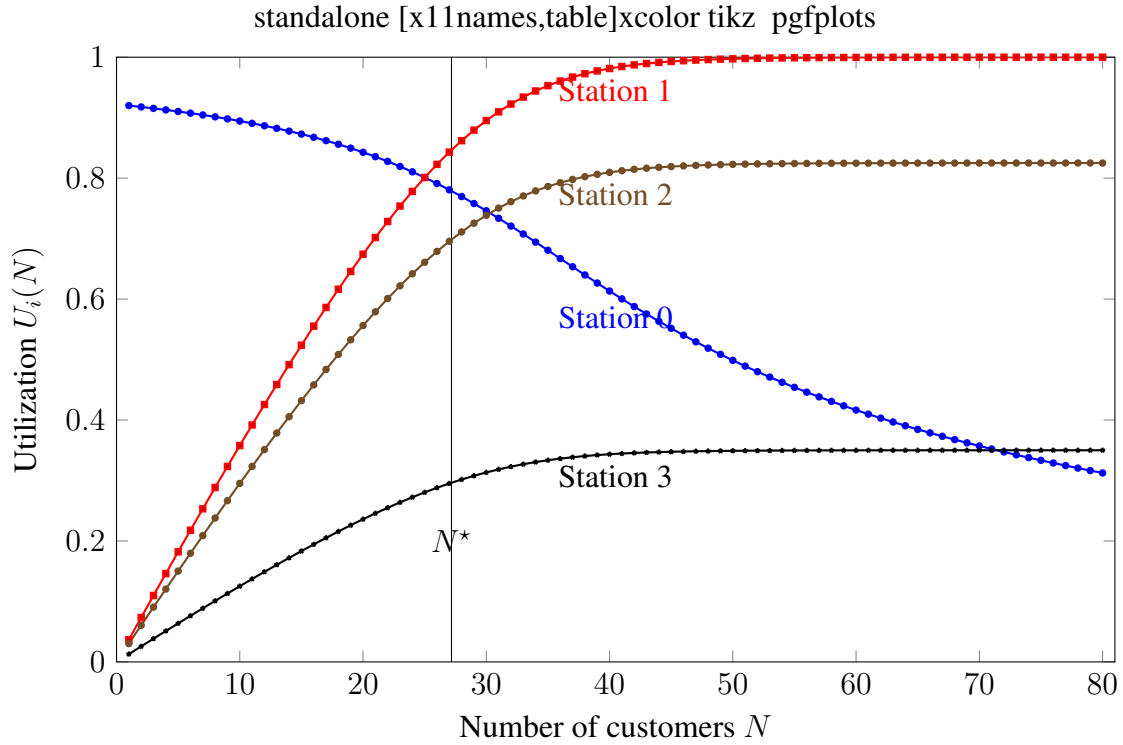


Figure 4: Behaviour of the utilization of the four stations.

Station 1 (the CPU) is the bottleneck station: we can see how, in comparison to the other stations, the queue becomes infinitely large in average while the other queues converge to a plateau. As expected, in the fourth graph the bottleneck station reaches utilization 1 after the saturation point. It is interesting to note that the utilization of the delay station keeps getting smaller even after the bottleneck stations is saturated and the others have converged to a steady value. This is not a surprise, since in our bottleneck analysis we got that the utilization should have been equal to 0 as $N \rightarrow \infty$.