

Bootstrap for quantiles estimation

Author:

Joris Limonier
University of Luxembourg
joris.limonier.001@student.uni.lu

Supervisor:

Nicolas Navet
University of Luxembourg
nicolas.navet@uni.lu

February - May 2021

Contents

1	Introduction	3
2	Bootstrap	4
3	Confidence intervals on the bootstrap	4
4	Conclusion	6
A	Supplementary tests on other datasets	7
A.1	CC1: CAM2 - ECU2	7
A.1.1	Data	7
A.1.2	Quantile 3	8
A.1.3	Quantile 4	9
A.1.4	Quantile 5	10
A.2	TFTP4 DAT	11
A.2.1	Data	11
A.2.2	Quantile 3	12
A.2.3	Quantile 4	13
A.2.4	Quantile 5	14
A.3	VD2	15
A.3.1	Data	15
A.3.2	Quantile 3	16
A.3.3	Quantile 4	17
A.3.4	Quantile 5	18

List of Figures

1	A look at the initial data	3
2	Estimation of the quantiles over bootstrap replicates	5
3	CC1: CAM2 - ECU2 data after multiple simulation times	7
4	CC1: CAM2 - ECU2 q_3 after multiple simulation times	8
5	CC1: CAM2 - ECU2 q_4 after multiple simulation times	9
6	CC1: CAM2 - ECU2 q_5 after multiple simulation times	10
7	TFTP4 DAT data after multiple simulation times	11
8	TFTP4 DAT q_3 after multiple simulation times	12
9	TFTP4 DAT q_4 after multiple simulation times	13
10	TFTP4 DAT q_5 after multiple simulation times	14
11	VD2 data after multiple simulation times	15
12	VD2 q_3 after multiple simulation times	16
13	VD2 q_4 after multiple simulation times	17
14	VD2 q_5 after multiple simulation times	18

1 Introduction

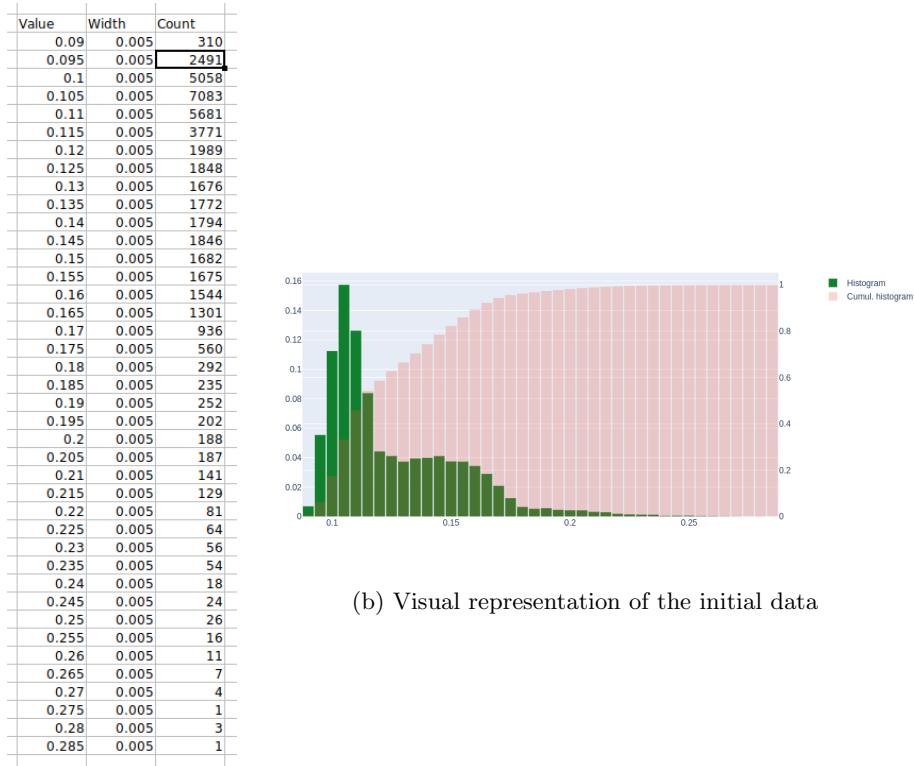
In this project we want to estimate the confidence we have when computing various quantiles. This task may be fairly straightforward when working with known distributions or arbitrary amounts of data but this seldom occurs in engineering or other real world applications. For this reason, here we try to accomplish that task on an unknown distribution and given limited amounts of data.

We define the k -th quantile as follows:

$$q_k := 1 - 10^{-k} \quad (1)$$

which gives $q_1 = 0.9$, $q_2 = 0.99$...etc. Where, simply speaking, q_k represents “0” followed by n nines. For our applications, we are mostly interested in q_3, q_4 and q_5 .

The data we work with is presented in figure 1a and plotted in the histogram in figure 1b.



(b) Visual representation of the initial data

(a) Initial data

Figure 1: A look at the initial data

2 Bootstrap

We use the bootstrap to estimate the confidence in the computed quantiles, even with small to moderate sample size.

Let n be the size of the data at our disposal and let q_k be the quantile we want to compute. Here is how we proceed:

1. Do the following R times.
 - (a) Draw a random sample of size n with replacement from the initial data.
 - (b) Compute q_k on the sample which has just been drawn.
2. Compute the mean over the set of values resulting from step 1b.
3. Compute the confidence intervals (details in section 3)

The plots in figure 2 show the evolution of our estimate for the value of the quantiles as we go through replicates of the bootstraps. The grey areas represent the 95% confidence intervals during that evolution. We will see how to get the confidence intervals from our histogram in section 3.

3 Confidence intervals on the bootstrap

Let us assume that we want to compute the 95% confidence interval. First, let us note the following:

$$\frac{1 - 0.95}{2} = 0.025 \quad \text{and} \quad 1 - \frac{(1 - 0.95)}{2} = 0.975 \quad (2)$$

Now we compute the 95% confidence interval. First, we sort the set of all the q 's that have been computed for each individual sample. Then we take the 0.025-th and 0.975-th quantile respectively as our lower-bound and upper bound for the confidence intervals.

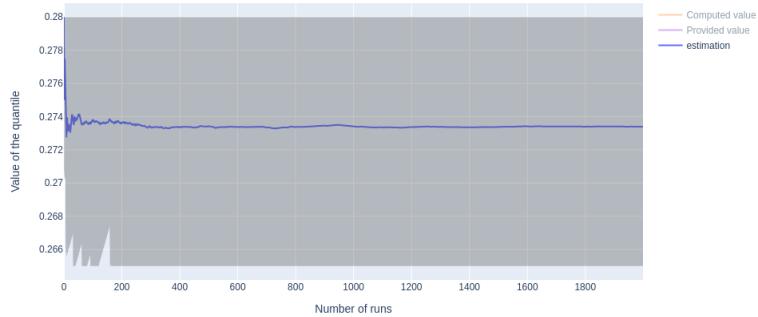
More generally, for a γ -confidence interval (instead of 95%), one has that the lower and upper bounds of the confidence interval are respectively

$$\gamma_{lo} := \frac{1 - \gamma}{2} \quad \text{and} \quad \gamma_{up} := 1 - \frac{(1 - \gamma)}{2}$$

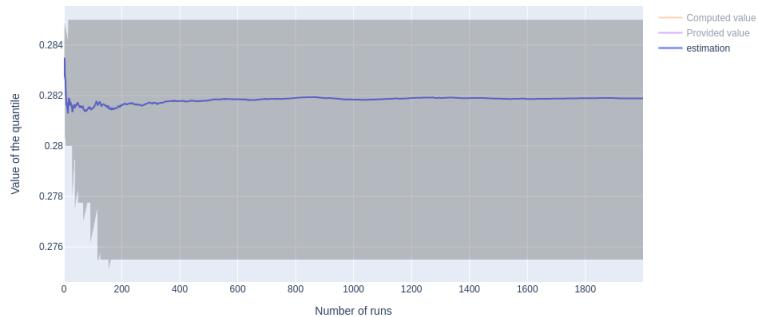
According to our data, it seems that there is close to no change in the estimation of the quantiles after 1000 replications of the bootstrap. The variations are small after 500 replicates already but for safety purposes we consider that we have our final guess after 1000 replicates. As for the confidence intervals, only in some edge cases do we have changes past the 1000 mark. The fact that our estimate seems to be stable after 1000 replicates does not matter for the number 1000 itself. However, it matters because we seem to “converge” to some value and reach a final value.



(a) Estimation of q_3



(b) Estimation of q_4



(c) Estimation of q_5

Figure 2: Estimation of the quantiles over bootstrap replicates

NB: Our estimate of 1000 replicates is based on empirical evidence. It is not a theoretical result, however, we believe that it is suitable for engineering purposes.

According to the litterature, that is according to Jean-Yves Le Boudec's *Performance Evaluation of Computer and Communication Systems*, there exists some "good value" for the number of bootstrap replicates R . With γ being the confidence level, Le Boudec advises to choose R as follows

$$R = \frac{50}{1 - \gamma} - 1 \quad (3)$$

therefore for $\gamma = 95\% = 0.95$ we deduce that the value of R is

$$\begin{aligned} R &= \frac{50}{1 - 0.95} - 1 \\ &= 1000 - 1 \\ &= 999 \end{aligned}$$

We find that our estimates verifies Le Boudec's result (1000 compared to 999).

4 Conclusion

The in-depth study in the main body of this article, as well as the supplementary tests provided in appendix A show a wide variety in the amounts of data at our disposal. The simulation times are 24 seconds, 4 minutes, 40 minutes and 400 minutes; in other words, 24 seconds times powers of 10, which gives an overview over multiple orders of magnitude. The number of data points ranges from a mere 1355 to almost 24 000 000 depending on the simulation time and the test case. The distributions themselves are very different too, some have well-spread values while others are very narrowly concentrated. As a result, we get some fairly contrasted results.

We notice that simulation times of 24 seconds appear too short and the results usually get better for 4, 40 and 400 minutes. However, and this is more surprising, in some cases, 400 minutes of simulation appear to yield worse results than 4 or 40 minutes (figure 8d for example). Having few cases makes it hard to state whether this is an edge case or a general trend. However, besides some rare cases (such as figure 13d), the improvements between 40 and 400 minutes of simulation time do not seem substantial. On the other hand, performing the bootstrap on an order of magnitude more data points takes a far greater time. A vast portion of the runtime was dedicated to performing the bootstrap on the 400 minutes simulation times (about half of the total ≈ 48 hours runtime on an 8-core i7-4710HQ CPU at 2.50GHz). In conclusion, it appears from our data that 4 to 40 minutes of simulation yields the best performance-to-runtime ratio.

A Supplementary tests on other datasets

To verify our approach and make sure that it generalises well, we test it on other datasets.

A.1 CC1: CAM2 - ECU2

A.1.1 Data

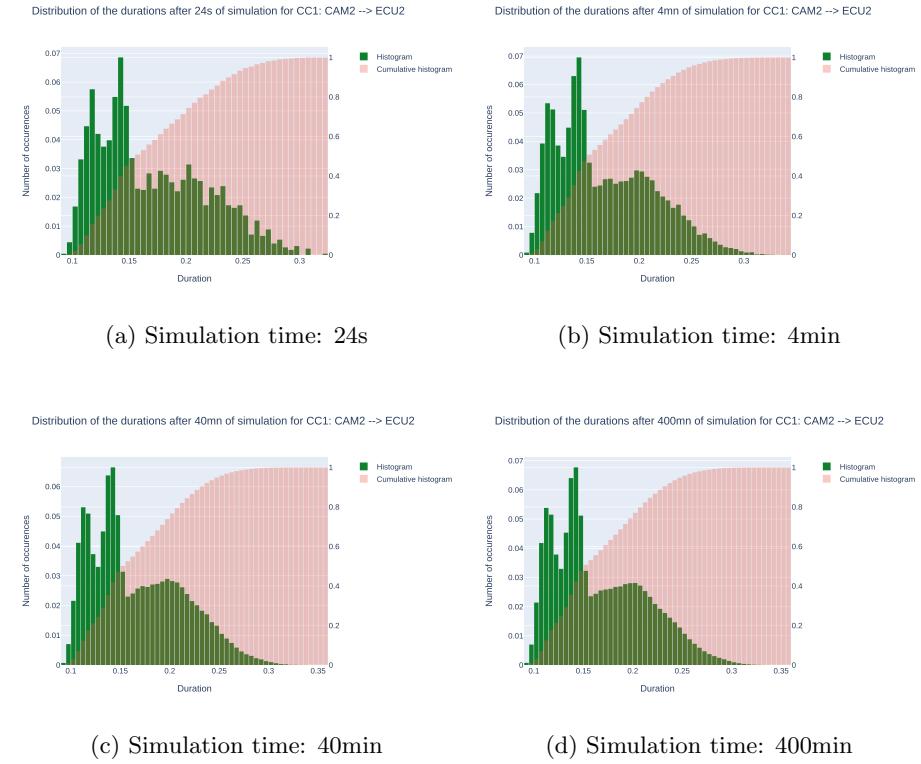


Figure 3: CC1: CAM2 - ECU2 data after multiple simulation times

A.1.2 Quantile 3

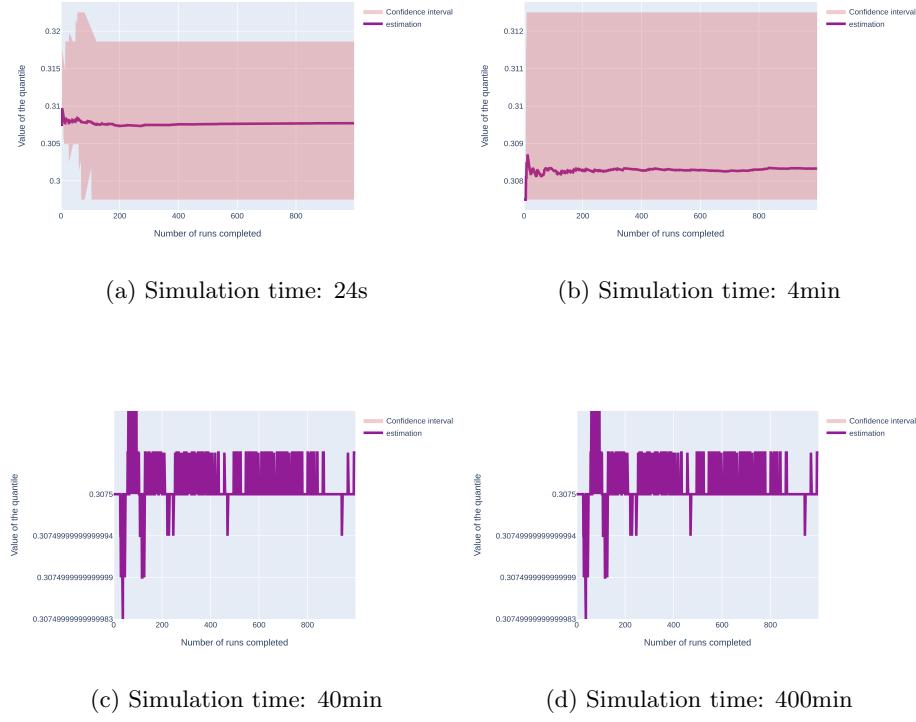


Figure 4: CC1: CAM2 - ECU2 q_3 after multiple simulation times

A.1.3 Quantile 4

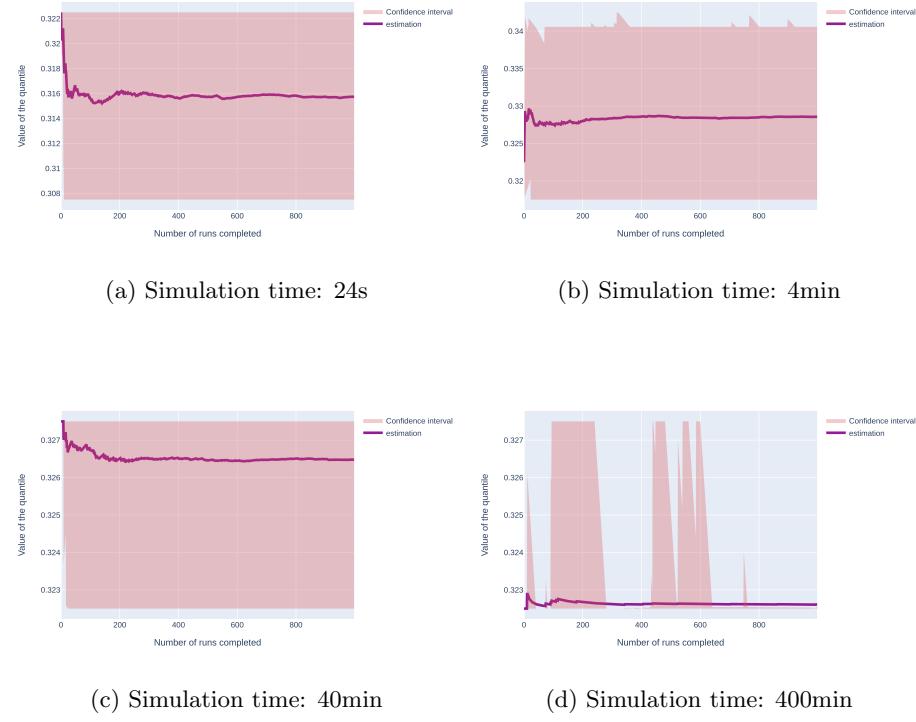


Figure 5: CC1: CAM2 - ECU2 q_4 after multiple simulation times

A.1.4 Quantile 5

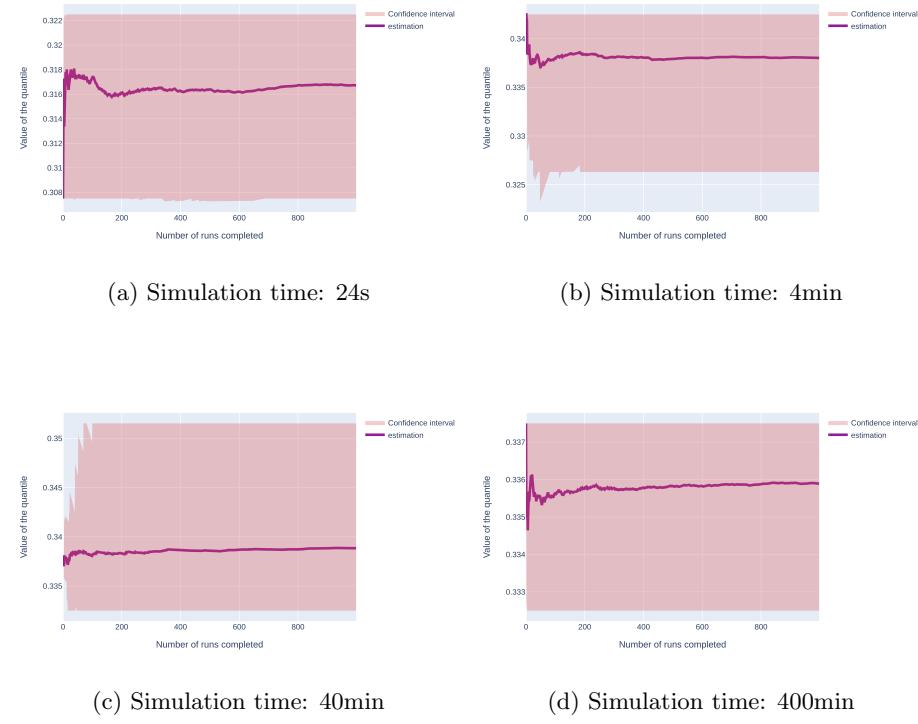


Figure 6: CC1: CAM2 - ECU2 q_5 after multiple simulation times

A.2 TFTP4 DAT

A.2.1 Data

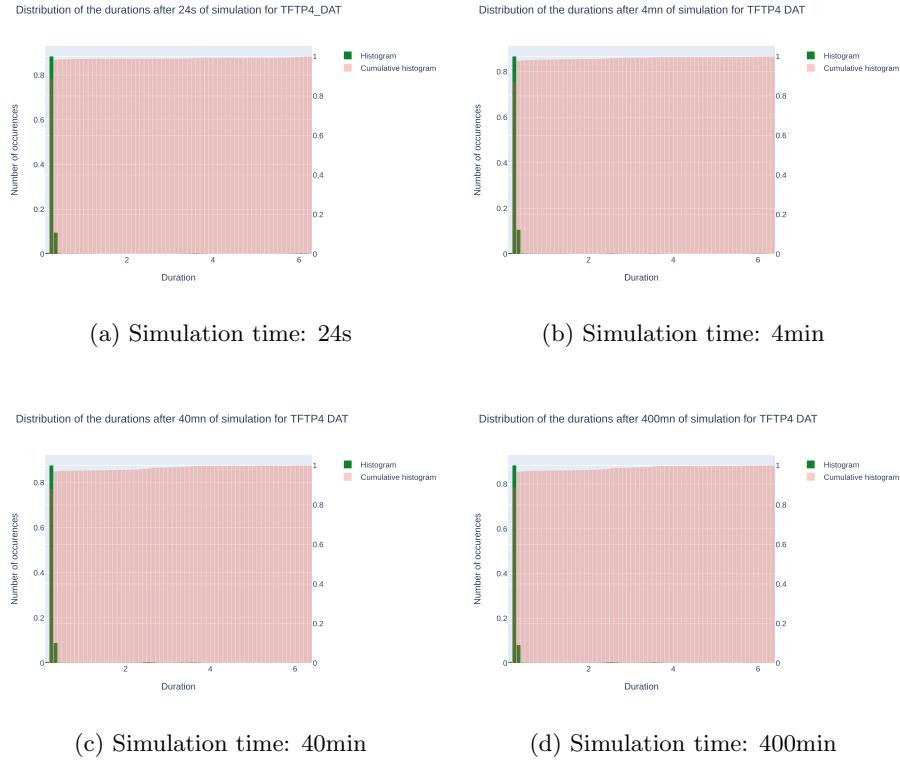


Figure 7: TFTP4 DAT data after multiple simulation times

A.2.2 Quantile 3

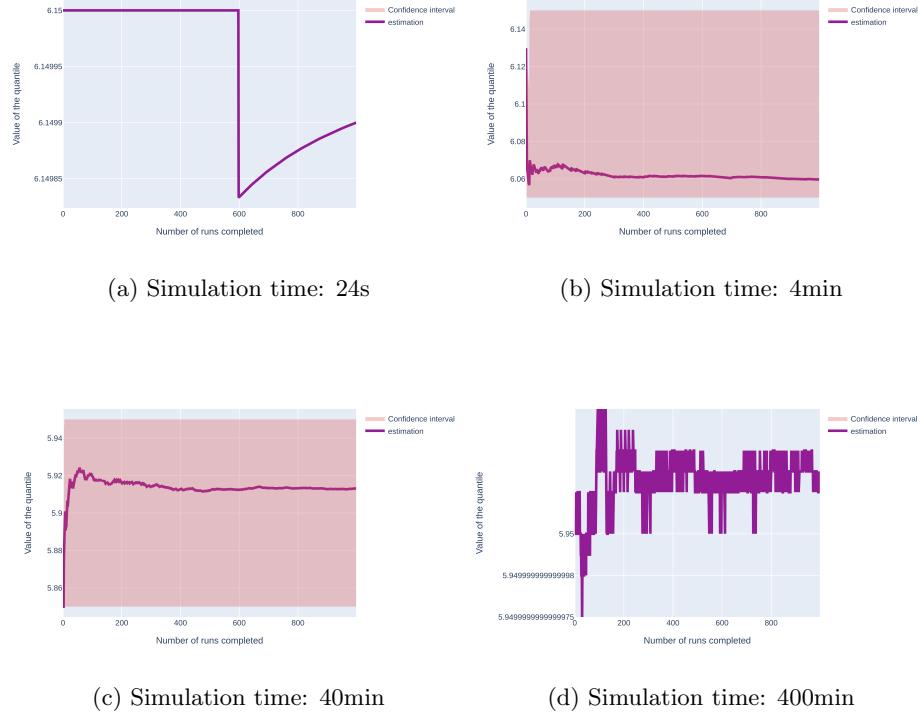


Figure 8: TFTP4 DAT q_3 after multiple simulation times

A.2.3 Quantile 4

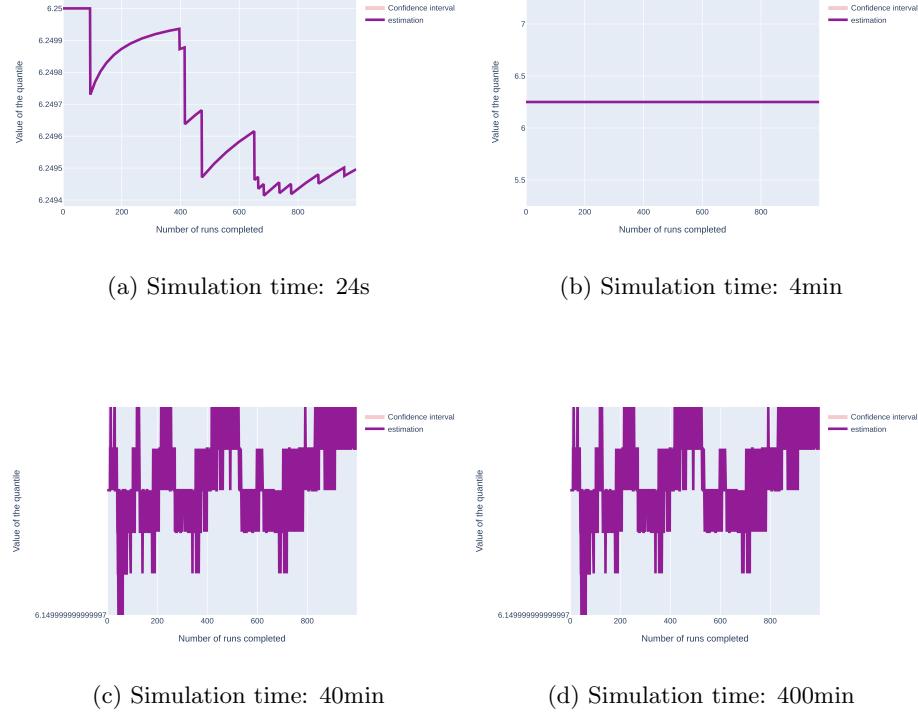


Figure 9: TFTP4 DAT q_4 after multiple simulation times

A.2.4 Quantile 5

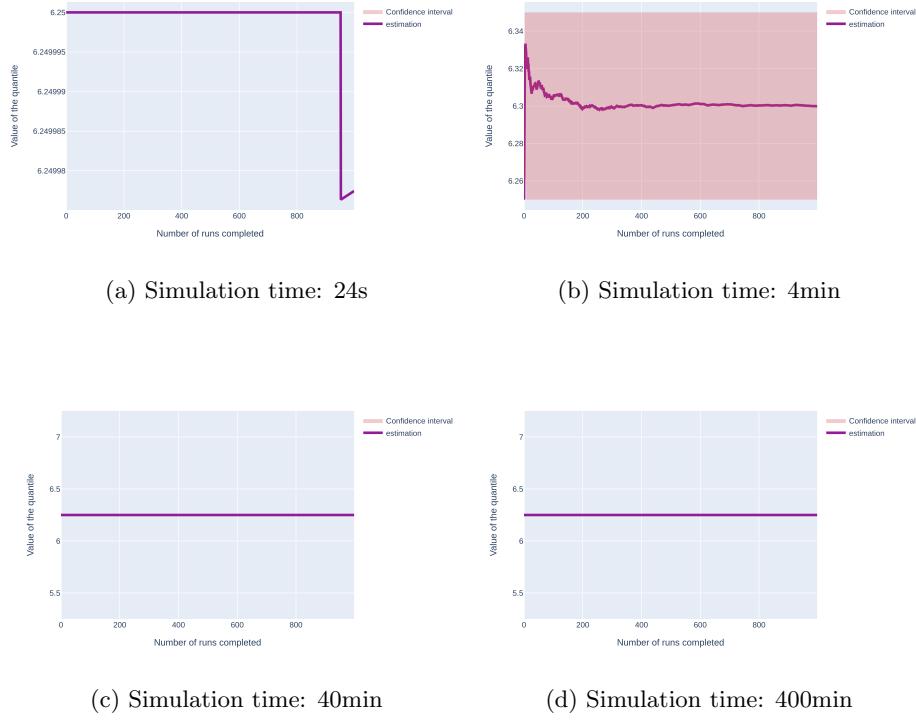


Figure 10: TFTP4 DAT q_5 after multiple simulation times

A.3 VD2

A.3.1 Data

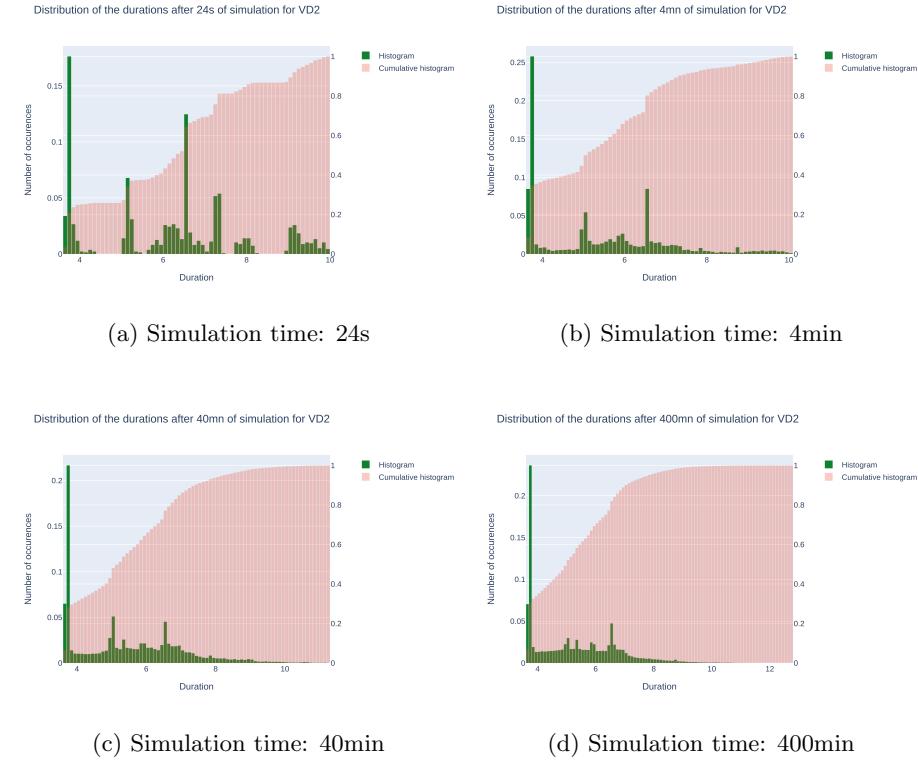


Figure 11: VD2 data after multiple simulation times

A.3.2 Quantile 3

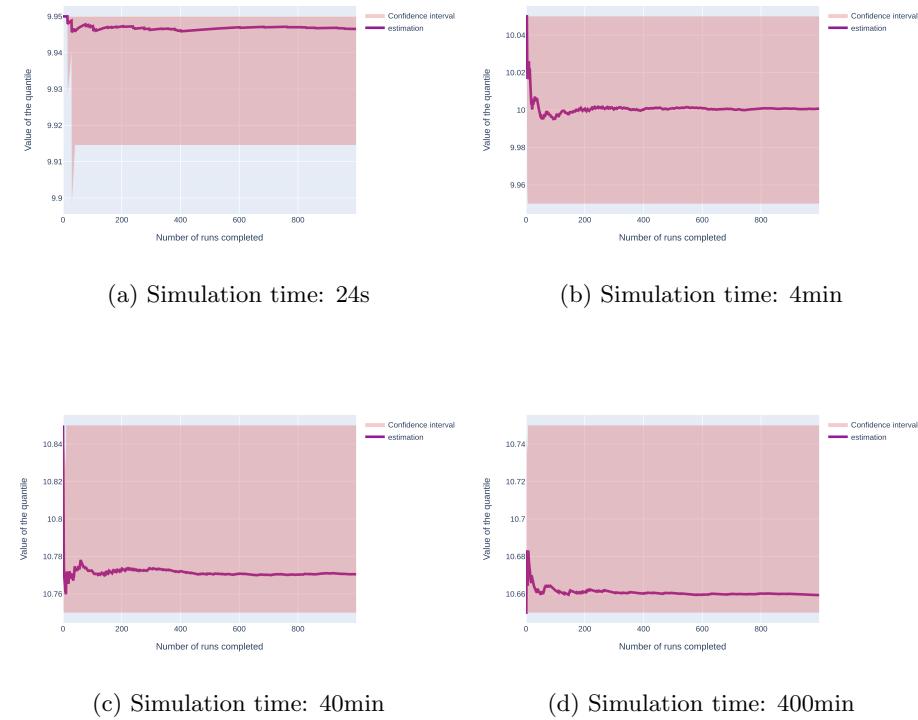


Figure 12: VD2 q_3 after multiple simulation times

A.3.3 Quantile 4

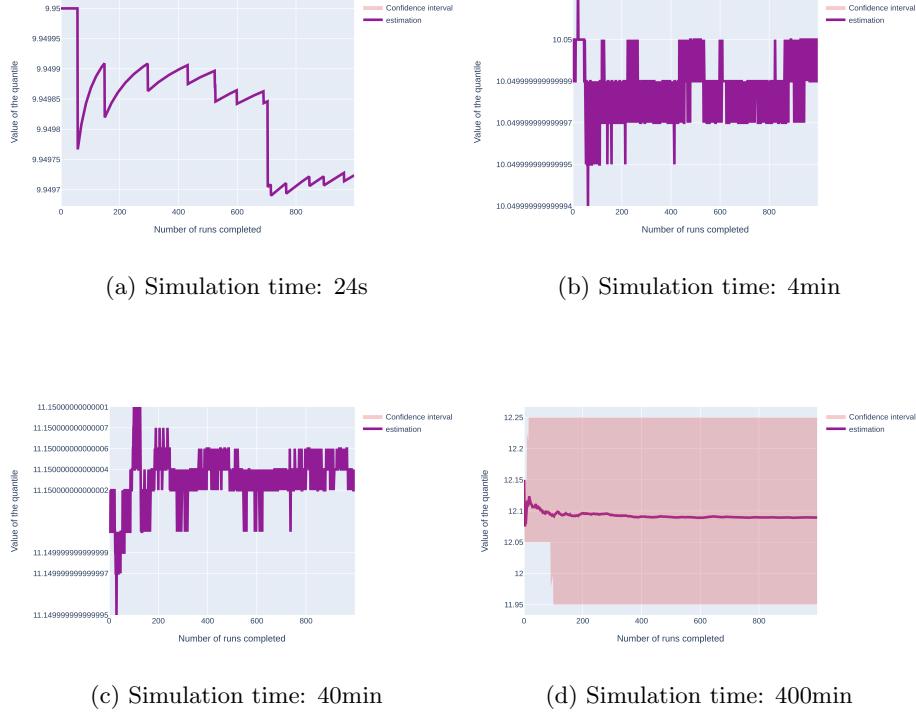


Figure 13: VD2 q_4 after multiple simulation times

A.3.4 Quantile 5

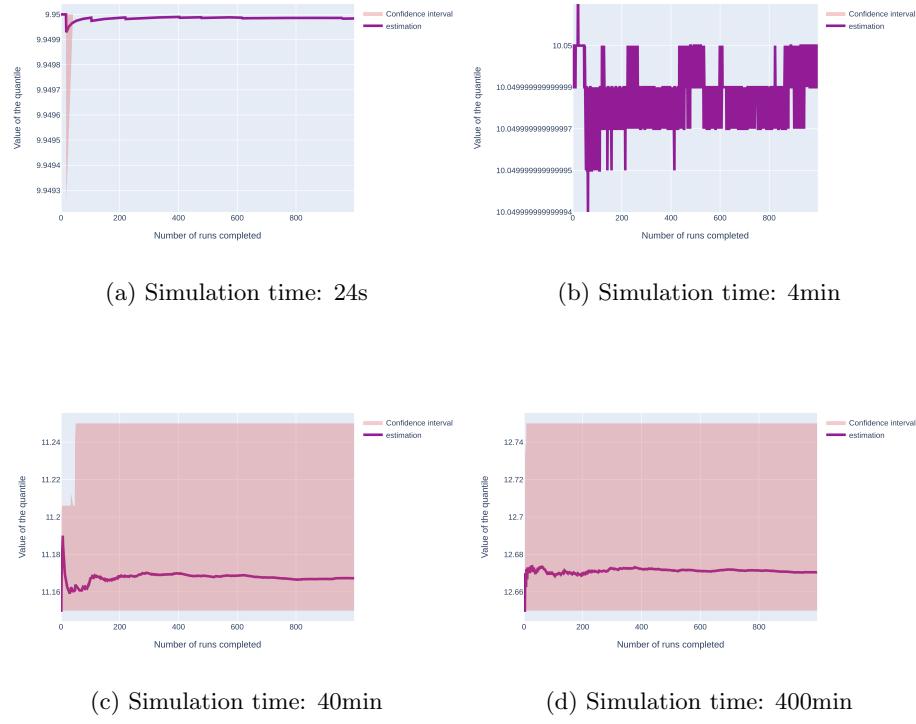


Figure 14: VD2 q_5 after multiple simulation times