

A Novel Incremental Quantile Estimator Using the Magnitude of the Observations

Hugo Lewi Hammer¹ and Anis Yazidi¹

Abstract—Incremental quantile estimators like the deterministic multiplicative incremental quantile estimator by Yazidi and Hammer (2017) are simple and efficient algorithms to estimate and track quantiles when data are received sequentially. The estimators merely relying on the sign of the difference between the quantile estimate and the current observation which seems like a waste of information from the data stream. In this paper we suggest a novel incremental estimator that rather use the *magnitude* of the observations. The intuition behind our approach is that the magnitude is more informative than the sign of the difference. Extensive experiments show that our estimators clearly outperform legacy state-of-the-art quantile estimators.

I. INTRODUCTION

In this paper we consider the problem of estimating quantiles when data arrive sequentially (data stream). The problem has been considered for many applications like portfolio risk measurement in the stock market [10], [1], fraud detection [29], signal processing and filtering [25], climate change monitoring [30], SLA violation monitoring [23], [24] and back-bone network monitoring [8].

Suppose that we are interested in estimating the quantile related to some probability q . The most natural estimator is to use the q quantile of the sample distribution. Unfortunately, such a quantile estimator has clear disadvantages for data streams as computation time and memory requirement are linear to the number of samples received so far from the data stream. Such methods thus are infeasible for large data streams.

Several algorithms have been proposed to deal with those challenges. Most of the proposed methods fall under to the category of what can be called histogram or batch based methods. The methods are based on efficiently maintaining a histogram estimate of the data stream distribution such that only a small storage footprint is required. A thorough review of state-of-the-art histogram and batch methods is given in the related work section (Section II).

Another ally of methods are the so-called incremental update methods. The latter methods are based on performing small updates of the quantile estimate every time a new sample is received from the data stream. Generally, the current estimate is a convex combination of the estimate at the previous time step and a quantity depending on the current observation. One of the first and prominent examples of this family of methods is the algorithm attributed to

Tierney (1983) [26] which is based on the stochastic learning theory. A few modifications of the Tierney method have been suggested, see e.g. [7], [4], [5], [6].

In data stream applications, a common situation is that the distribution of the samples from the data stream varies with time. Such system or environment is referred to as a dynamical system in the literature. Unfortunately, histogram based methods usually perform poorly in estimating quantiles in such systems and we are left with incremental methods as typically the only viable lightweight alternatives [5].

Current state-of-the art incremental estimators, like [7], [4] and the more recently the deterministic multiplicative incremental quantile estimator (DUMIQE) [28], do not use the values of the received samples directly to update the estimate, but only whether the value of the sample is above or below some threshold which is in this case the value of the current quantile estimate. Intuitively, this seems like a waste of information received from the data stream. In this paper, we thus present an estimator that uses the value of the received sample directly when updating the quantile estimate. The estimator is such that the update step size is *proportional* to the distance between the current estimate and the value of the sample. If the distance is large our estimator is off track and large jumps should be initiated to rapidly get the estimator back on track.

II. RELATED WORK

In this Section we shall review some of the related work on estimating quantiles from data streams. However, as we will explain later, these related works require some memory restrictions which renders our work radically distinct from them. In fact, our approach requires storing only one sample value in order to update the estimate. The most representative work for this type of “streaming” quantile estimator is due to the seminal work of Munro and Paterson [18]. In [18], Munro and Paterson described a p -pass algorithm for selection using $O(n^{1/(2p)})$ space for any $p \geq 2$. Cormode and Muthukrishnan [9] proposed a more space-efficient data structure, called the Count-Min sketch, which is inspired by Bloom filters, where one estimates the quantiles of a stream as the quantiles of a random sample of the input. The key idea is to maintain a random sample of an appropriate size to estimate the quantile, where the premise is to select a subset of elements whose quantile approximates the true quantile. From this perspective, the latter body of research requires a certain amount of memory that increases as the required accuracy of the estimator increases [27]. Furthermore, in the case where the underlying distribution changes over time, those methods

¹Hugo Lewi Hammer and Anis Yazidi are with Department of Computer Science at Oslo and Akershus University College of Applied Sciences, Pilestredet 35, N-0166 Oslo, Norway {hugo.hammer, anis.yazidi}@hioa.no

suffer from large bias in the summary information since the stored data might be stale [7]. Examples of these works include [3], [27], [18], [11], [12]. Guha and McGregor [12] advocate the use of random-order data models in contrast to adversarial-order models. They show that computing the median requires exponential number of passes in adversarial model while requiring $O(\log \log n)$ in random order model.

In [7], the authors proposed a modification of the stochastic approximation algorithm [26]. While Tierney [26] uses a sample mean update from previous quantile estimates, [7] propose an exponential decay in the usage of old estimates. This modification is particularly helpful in the case of non-stationary environments in order to cope with non-stationary data. Indeed, a “weighted” update scheme is applied to incrementally build local approximations of the distribution function in the neighborhood of the quantiles.

In many network monitoring applications, quantiles are key indicators for monitoring the performance of the system. For instance, system administrators are interested in monitoring the 95% quantile of the response time of a web-server so that to hold it under a certain threshold. Quantile tracking is also useful for detecting abnormal events and in intrusion detection systems in general. However, the immense traffic volume of high speed networks impose some computational challenges: little storage and the fact that the computation needs to be “one pass” on the data. It is worth mentioning that the seminal paper of Robbins and Monro [21] which established the field of research called “stochastic approximation” [15] have included an incremental quantile estimator as a proof of concept of the vast applications of the theory of stochastic approximation. An extension of the latter quantile estimator which first appeared as example in [21] was further developed in [14] in order to handle the case of “extreme quantiles”. Moreover, the estimator provided by Tierney [26] falls under the same umbrella of the example given in [21], and thus can be seen as an extension of it.

As Arandjelovic remarks [2], most quantile estimation algorithms are not single-pass algorithms and thus are not applicable for streaming data. On the other hand, the single pass algorithms are concerned with the exact computation of the quantile and thus require a storage space of the order of the size of the data which is clearly an unfeasible condition in the context of big data stream.

Thus, we submit that all work on quantile estimation using more than one pass, or storage of the same order of the size of the observations seen so far is not relevant in the context of this paper.

When it comes to memory efficient methods that require a small storage footprint, histogram based methods form an important class. A representative work in this perspective is due to Schmeiser and Deutsch [22]. In fact, they proposed to use equidistant bins where the boundaries are adjusted online. Arandjelovic et al. [2] use a different idea than equidistant bins by attempting to maintain bins in a manner that maximizes the entropy of the corresponding estimate of the historical data distribution. Thus, the bin boundaries are adjusted in an online manner. Nevertheless, histogram based

methods have problems to deal with dynamic data where the underlying distribution changes over time [5]. In addition, they are prone to outliers that might corrupt the estimates of the distribution.

In [19], the authors propose a memory efficient method for simultaneous estimation of several quantiles using interpolation methods and a grid structure where each internal grid point is updated upon receiving an observation. The application of this approach is limited for stationary data. The approximation of the quantiles relies on using linear and parabolic interpolations, while the tails of the distribution are approximated using exponential curves. It is worth mentioning that the latter algorithm is based on the P^2 algorithm [13].

In [13], Jain et al. resort to five markers so that to track the quantile, where the markers correspond to different quantiles and the min and max of the observations. Their concept is similar to the notion of histograms, where each marker has two measurements, its height and its position. By definition, each marker has some ideal position, where some adjustments are made to keep it in its ideal position by counting the number of samples exceeding the marker. In simple terms, for example, if the marker corresponds to the 80% quantile, its ideal position will be around the point corresponding to 80% of the data points below the marker. However, such approach does not handle the case of non-stationary quantile estimation as the position of the markers will be affected by stale data points. Then based on the position of the markers, quantiles are computed by supposing that the curve passing through three adjacent markers is parabolic and by using a piecewise parabolic prediction function.

In fact, the insertion can be handled easily using either sequential or batch updates, while quantile update upon deletion requires more complex forms of updates.

Finally, Lou et al. [16] perform extensive experiments to compare several of the algorithms described above.

III. INCREMENTAL QUANTILE ESTIMATOR USING THE MAGNITUDE OF OBSERVATIONS

Let X_n denote a stochastic variable representing the possible outcomes from a data stream at time n and let x_n denote a random sample (realization) of X_n . We assume that X_n is distributed according to some distribution $f_n(x)$ that varies dynamically with time n . We denote the cumulative distribution of X_n with $F_n(x)$, i.e. $P(X_n \leq x) = F_n(x)$. Further let $Q_n(q)$ denote the quantile associated with probability q , i.e. $P(X_n \leq Q_n(q)) = F_n(Q_n(q)) = q$.

Yazidi and Hammer [28] introduced the DUMIQE estimator given by

$$\begin{aligned} \tilde{Q}_{n+1}(q) &\leftarrow \tilde{Q}_n(q) + \lambda q \tilde{Q}_n(q) & \text{if } x_n > \tilde{Q}_n(q) \\ \tilde{Q}_{n+1}(q) &\leftarrow \tilde{Q}_n(q) - \lambda(1-q) \tilde{Q}_n(q) & \text{if } x_n \leq \tilde{Q}_n(q) \end{aligned} \quad (1)$$

The intuition behind the algorithm is simple. If the received sample gets a value below (above) the current quantile estimate, the estimate is decreased (increased). The “weights”

q and $1 - q$ is included to ensure convergence to the true quantile. Even though the estimator is really simple, it was shown to document state-of-the-art performance for tracking quantiles of dynamic data streams [28]. However, the estimator only uses the sign of the difference between the current estimate and the received sample and it may be a potential to achieve more efficient estimators by using more information from the data stream.

We now suggest a novel quantile estimator where the update step size is *proportional* to the distance between the received sample and current estimate. The intuition is that if the distance is large the estimator is off track and large jumps should be initiated to rapidly get the estimator back on track. The suggested estimator is as follows

$$\begin{aligned}\hat{Q}_{n+1}(q) &\leftarrow \hat{Q}_n(q) + \lambda c_n \frac{q}{\mu_n^+ - X_n} |x_n - \hat{Q}_n(q)| && \text{if } x_n > \hat{Q}_n(q) \\ \hat{Q}_{n+1}(q) &\leftarrow \hat{Q}_n(q) - \lambda c_n \frac{1-q}{X_n - \mu_n^-} |x_n - \hat{Q}_n(q)| && \text{if } x_n \leq \hat{Q}_n(q)\end{aligned}\quad (2)$$

where $\mu_n^+ = E(X_n | X_n > \hat{Q}_n(q))$ and $\mu_n^- = E(X_n | X_n < \hat{Q}_n(q))$. The constants c_n can be any sequence positive and bounded values. The estimator performed well, when the fractions in (2) were “normalized”, i.e. we chose

$$c_n = \left(\frac{q}{\mu_n^+ - \hat{Q}_n(q)} + \frac{1-q}{\hat{Q}_n(q) - \mu_n^-} \right)^{-1}$$

Finally note that the conditinal expectations must be on each side of the current estimate, i.e. $\mu_n^- < \hat{Q}_n(q) < \mu_n^+$.

Now we will present a theorem that catalogs the properties of the estimator $\hat{Q}_n(q)$ in (2) for a stationary data stream, i.e. $X_n = X \sim F(x)$, $n = 1, 2, \dots$

Theorem 1: Let $Q(q) = F^{-1}(q)$ be the true quantile to be estimated. Applying the updating rule (2), we obtain:

$$\lim_{n \rightarrow \infty, \lambda \rightarrow 0} \hat{Q}_n(q) = Q(q)$$

The theorem can be proved using the theory of stochastic learning due to Norman [20], but is omitted for the sake of brevity.

A. Quantile Estimation Algorithms

The update rules (2) and Theorem 1 constitute some intriguing results. However, the estimator cannot be used directly since the conditional expectations μ_n^+ and μ_n^- are unknown. The natural solution is to track the conditinal expectations as well from the data stream observations. The quantile estimation algorithm is then as follows:

We start with some initial values $\hat{Q}_0(q)$, μ_0^+ and μ_0^- satisfying $\mu_0^- < \hat{Q}_0(q) < \mu_0^+$. Every time we receive a new observation from the data stream, we do the following

updates

$$\hat{Q}_{n+1}(q) \leftarrow \hat{Q}_n(q) + \lambda \hat{c}_n \frac{q}{\hat{\mu}_n^+ - X_n} |x_n - \hat{Q}_n(q)| \quad \text{if } x_n > \hat{Q}_n(q) \quad (3)$$

$$\hat{Q}_{n+1}(q) \leftarrow \hat{Q}_n(q) - \lambda \hat{c}_n \frac{1-q}{X_n - \hat{\mu}_n^-} |x_n - \hat{Q}_n(q)| \quad \text{if } x_n \leq \hat{Q}_n(q)$$

$$\hat{\mu}_{n+1}^+ \leftarrow \hat{Q}_{n+1}(q) - \hat{Q}_n(q) + (1-\gamma)\hat{\mu}_n^+ + \gamma x_n \quad \text{if } x_n > \hat{Q}_n(q) \quad (4)$$

$$\hat{\mu}_{n+1}^+ \leftarrow \hat{\mu}_n^+ \quad \text{if } x_n \leq \hat{Q}_n(q)$$

$$\hat{\mu}_{n+1}^- \leftarrow \hat{\mu}_n^- \quad \text{if } x_n > \hat{Q}_n(q)$$

$$\hat{\mu}_{n+1}^- \leftarrow \hat{Q}_{n+1}(q) - \hat{Q}_n(q) + (1-\gamma)\hat{\mu}_n^- + \gamma x_n \quad \text{if } x_n \leq \hat{Q}_n(q) \quad (5)$$

$$\hat{c}_{n+1} = \left(\frac{q}{\hat{\mu}_{n+1}^+ - \hat{Q}_{n+1}(q)} + \frac{1-q}{\hat{Q}_{n+1}(q) - \hat{\mu}_{n+1}^-} \right)^{-1}$$

We see that (4) and (5) are updated using a exponentially weighted moving average, but with the additional part $\hat{Q}_{n+1}(q) - \hat{Q}_n(q)$. The part $\hat{Q}_{n+1}(q) - \hat{Q}_n(q)$ is included to ensure that $\hat{\mu}_n^- < \hat{Q}_n(q) < \hat{\mu}_n^+$ in every iteration. The tuning parameter, γ , makes sure that the conditional expectations $\hat{\mu}_{n+1}^+$ and $\hat{\mu}_{n+1}^-$ are estimated correctly relative to the current quantile estimate $\hat{Q}_{n+1}(q)$. In other words, (3) tracks the overall trends of the dynamical data stream while (4) and (5) are responsible of estimating the conditional expectations relative to the quantile estimate. We thus expect that for most dynamical data streams, it is reasonable to use a value of γ that is on a smaller scale than λ . This is verified in our experiments.

IV. SIMULATION EXPERIMENTS

In this section we evaluate the performance of the suggested algorithm against four state-of-the-art quantile estimators namely the DUMIQE and RUMIQE by Yazidi and Hammer [28], the estimator due to Cao et al. [4] and the Frugal approach by Ma et al. [17]. The estimator in this paper is designed to perform well for dynamically changing data streams and the experiments will focus on such streams. It would have be interesting to evaluate the performance of the different methods for real life data, but this is challenging to do in a systematic way for dynamical data streams as the ground truth generally is missing.

We assume a data stream where the outcomes are from a normal distribution and where the expectation jumps between values a and $-a$.

$$\mu_n = \begin{cases} a & \text{if } n \bmod T \leq T/2 \\ -a & \text{else} \end{cases} \quad (6)$$

We assume that the standard deviation of the normal distribution does not vary with time and is equal to one.

We estimated quantiles two different periods, namely $T = 100$ (rapid variation) and $T = 500$ (slow variation). For each data stream we estimated the 50, 70 and 90% quantiles ending up with a total of six different estimation tasks.

To measure estimation error, we used the root mean squares error (RMSE) for each quantile given as:

$$\text{RMSE} = \sqrt{\frac{1}{N} \sum_{n=1}^N \left(Q_n(q) - \hat{Q}_n(q) \right)^2} \quad (7)$$

where N is the total number of samples in the data stream. In the experiments, we used $N = 10^6$ which efficiently removes any Monte Carlo errors in the experimental results. In order to get a good overview of the performance of the algorithms, we measure the estimation error for a large set of different values of the tuning parameters of the algorithms.

Figure 1 illustrate the results of our experiments. We see that the suggested algorithm (blue color) outperforms the state-of-the-art algorithms with a clear margin. Interestingly, the suggested algorithm performs the best with a small value of the ratio γ/λ as small as $1/100$. This is in accordance with what we expected in Section III-A. The Cao et al. algorithm struggled with numerical problems for some choices of the tuning parameters and therefore some of the curves are short.

Figure 2 illustrates why the suggested algorithm performs so well. The expectation changes as given in (6) with period $T = 500$ and the estimators track the quantile related to the probability $q = 0.7$. The true quantile is given as the dashed black line. We compare the suggested algorithm against DUMIQE since it performed well in the experiments (Figure 1). The tuning parameters are adjusted such that the estimation error in the stationary parts after convergence is the same for the two algorithms. Still we see that the suggested algorithm track the true quantile more efficiently after a switch than the DUMIQE. The explanation is the step size of the suggested algorithm is proportional to the difference between the observations and the quantile estimate (recall (2)). After a switch, these differences are large, and the suggested algorithm makes large steps to get back on track. The DUMIQE, and the other state of the incremental algorithms, use the same step size independent of the these difference, resulting in poorer tracking.

V. CLOSING REMARKS

In this paper we suggest a new algorithm where the update step size is proportional to the difference between the received observation and the current estimate. The current state-of-the-art quantile estimators merely relying on the sign of the difference between the quantile estimate and the current observation which seems like a waste of information from the data stream.

In the future we will evaluate the performance of the algorithm in more experiments and apply the estimator to real-life problems.

REFERENCES

- [1] Babak Abbasi and Montserrat Guillen. Bootstrap control charts in monitoring value at risk in insurance. *Expert Systems with Applications*, 40(15):6125–6135, 2013.
- [2] Ognjen Arandjelovic, Duc-Son Pham, and Svetha Venkatesh. Two maximum entropy-based algorithms for running quantile estimation in nonstationary data streams. *Circuits and Systems for Video Technology, IEEE Transactions on*, 25(9):1469–1479, 2015.
- [3] Arvind Arasu and Gurmeet Singh Manku. Approximate counts and quantiles over sliding windows. In *Proceedings of the twenty-third ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pages 286–296. ACM, 2004.
- [4] Jin Cao, Li Li, Aiyu Chen, and Tian Bu. Tracking quantiles of network data streams with dynamic operations. In *INFOCOM, 2010 Proceedings IEEE*, pages 1–5. IEEE, 2010.
- [5] Jin Cao, Li Erran Li, Aiyu Chen, and Tian Bu. Incremental tracking of multiple quantiles for network monitoring in cellular networks. In *Proceedings of the 1st ACM workshop on Mobile internet through cellular networks*, pages 7–12. ACM, 2009.
- [6] John M Chambers, David A James, Diane Lambert, and Scott Vander Wiel. Monitoring networked applications with incremental quantile estimation. *Statistical Science*, pages 463–475, 2006.
- [7] Fei Chen, Diane Lambert, and José C Pinheiro. Incremental quantile estimation for massive tracking. In *Proceedings of the sixth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 516–522. ACM, 2000.
- [8] Baek-Young Choi, Sue Moon, Rene Cruz, Zhi-Li Zhang, and Christophe Diot. Quantile sampling for practical delay monitoring in internet backbone networks. *Computer Networks*, 51(10):2701–2716, 2007.
- [9] Graham Cormode and S Muthukrishnan. An improved data stream summary: the count-min sketch and its applications. *Journal of Algorithms*, 55(1):58–75, 2005.
- [10] Manfred Gilli et al. An application of extreme value theory for measuring financial risk. *Computational Economics*, 27(2-3):207–228, 2006.
- [11] Michael Greenwald and Sanjeev Khanna. Space-efficient online computation of quantile summaries. In *ACM SIGMOD Record*, volume 30, pages 58–66. ACM, 2001.
- [12] Sudipto Guha and Andrew McGregor. Stream order and order statistics: Quantile estimation in random-order streams. *SIAM Journal on Computing*, 38(5):2044–2059, 2009.
- [13] Raj Jain and Imrich Chlamtac. The p 2 algorithm for dynamic calculation of quantiles and histograms without storing observations. *Communications of the ACM*, 28(10):1076–1085, 1985.
- [14] V Roshan Joseph. Efficient robbins-monro procedure for binary data. *Biometrika*, 91(2):461–470, 2004.
- [15] Harold Kushner and G George Yin. *Stochastic approximation and recursive algorithms and applications*, volume 35. Springer Science & Business Media, 2003.
- [16] Ge Luo, Lu Wang, Ke Yi, and Graham Cormode. Quantiles over data streams: experimental comparisons, new analyses, and further improvements. *The VLDB Journal*, pages 1–24, 2016.
- [17] Qiang Ma, S Muthukrishnan, and Mark Sandler. Frugal streaming for estimating quantiles. In *Space-Efficient Data Structures, Streams, and Algorithms*, pages 77–96. Springer, 2013.
- [18] J Ian Munro and Mike S Paterson. Selection and sorting with limited storage. *Theoretical computer science*, 12(3):315–323, 1980.
- [19] Valeriy Naumov and Olli Martikainen. Exponentially weighted simultaneous estimation of several quantiles. *World Academy of Science, Engineering and Technology*, 8:563–568, 2007.
- [20] M Frank Norman. *Markov processes and learning models*, volume 84. Academic Press New York, 1972.
- [21] Herbert Robbins and Sutton Monro. A stochastic approximation method. *The annals of mathematical statistics*, pages 400–407, 1951.
- [22] Bruce W Schweiser and Stuart Jay Deutsch. Quantile estimation from grouped data: The cell midpoint. *Communications in Statistics-Simulation and Computation*, 6(3):221–234, 1977.
- [23] Joel Sommers, Paul Barford, Nick Duffield, and Amos Ron. Accurate and efficient sla compliance monitoring. In *ACM SIGCOMM Computer Communication Review*, volume 37, pages 109–120. ACM, 2007.
- [24] Joel Sommers, Paul Barford, Nick Duffield, and Amos Ron. Multi-objective monitoring for sla compliance. *IEEE/ACM Transactions on Networking (TON)*, 18(2):652–665, 2010.

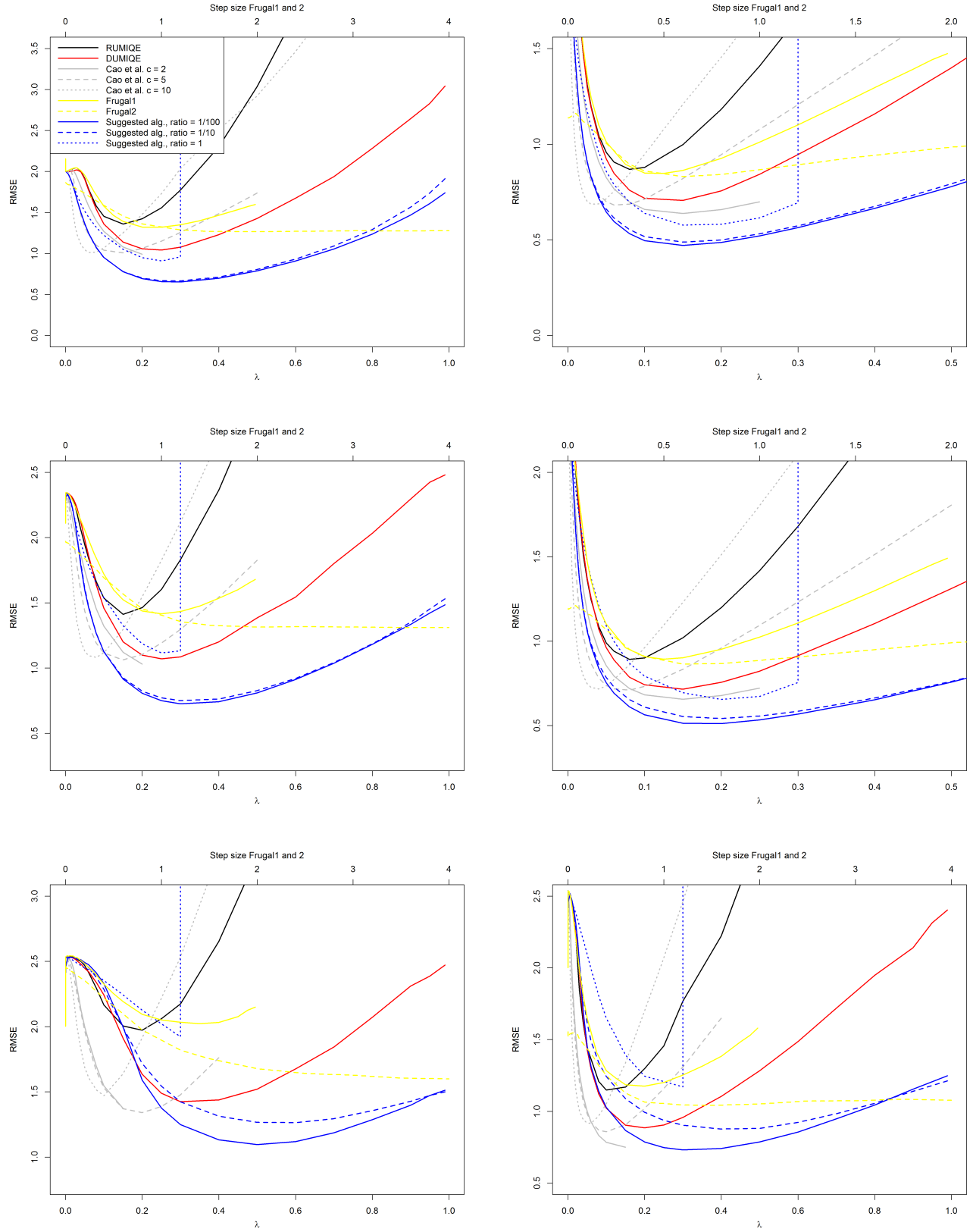


Fig. 1. The left and right columns show results for $T = 100$ and $T = 500$, respectively. The rows from top to bottom show results when estimating quantile $Q_n(q = 0.5)$, $Q_n(q = 0.7)$ and $Q_n(q = 0.9)$, respectively. Ratio refers to the ratio between the tuning parameters, i.e. ratio = γ/λ . The upper x axis refers to the step size in the Frugal algorithms.

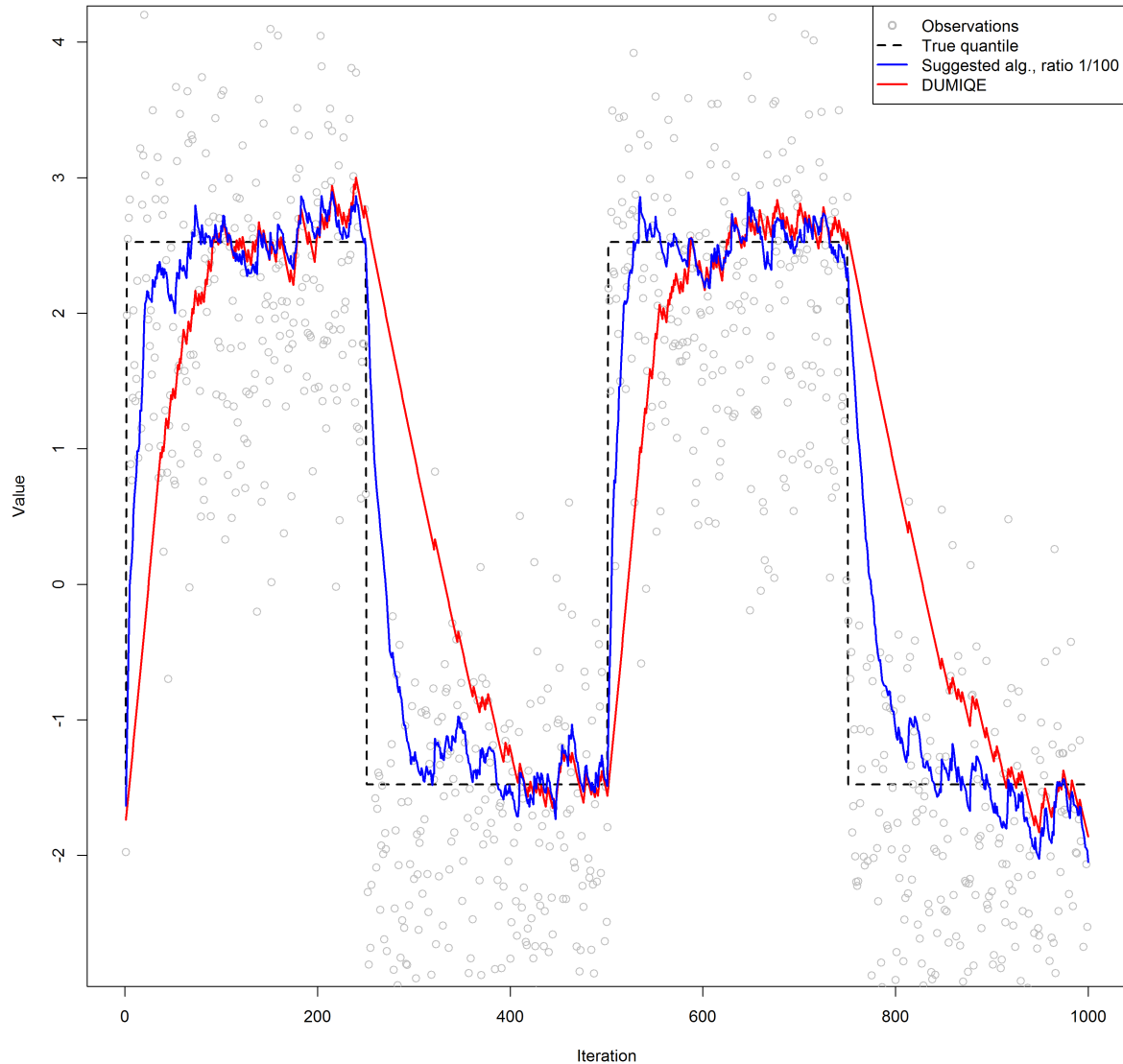


Fig. 2. Quantile estimates in every iteration using the DUMIQE and the suggested algorithms.

- [25] Volker Stahl, Alexander Fischer, and Rolf Bippus. Quantile based noise estimation for spectral subtraction and wiener filtering. In *Acoustics, Speech, and Signal Processing, 2000. ICASSP'00. Proceedings. 2000 IEEE International Conference on*, volume 3, pages 1875–1878. IEEE, 2000.
- [26] Luke Tierney. A space-efficient recursive procedure for estimating a quantile of an unknown distribution. *SIAM Journal on Scientific and Statistical Computing*, 4(4):706–711, 1983.
- [27] B Weide. Space-efficient on-line selection algorithms. In *Computer Science and Statistics: Proceedings of the Eleventh Annual Symposium on the Interface*, pages 308–311, 1978.
- [28] Anis Yazidi and Hugo Lewi Hammer. Multiplicative Update Methods for Incremental Quantile Estimation. *IEEE Transactions on Cybernetics (accepted)*, 2017.
- [29] Linfeng Zhang and Yong Guan. Detecting click fraud in pay-per-click streams of online advertising networks. In *Distributed Computing Systems, 2008. ICDCS'08. The 28th International Conference on*, pages 77–84. IEEE, 2008.
- [30] Xuebin Zhang, Lisa Alexander, Gabriele C Hegerl, Philip Jones, Albert Klein Tank, Thomas C Peterson, Blair Trewin, and Francis W Zwiers. Indices for monitoring changes in extremes based on daily temperature and precipitation data. *Wiley Interdisciplinary Reviews: Climate Change*, 2(6):851–870, 2011.