

ON NON-ASYMPTOTIC OPTIMAL STOPPING CRITERIA IN MONTE CARLO SIMULATIONS*

CHRISTIAN BAYER[†], HÅKON HOEL^{‡§}, ERIK VON SCHWERIN^{‡¶}, AND RAÚL TEMPONE[‡]

Abstract. We consider the setting of estimating the mean of a random variable by a sequential stopping rule Monte Carlo (MC) method. The performance of a typical second moment based sequential stopping rule MC method is shown to be unreliable in such settings both by numerical examples and through analysis. By analysis and approximations, we construct a higher moment based stopping rule which is shown in numerical examples to perform more reliably and only slightly less efficiently than the second moment based stopping rule.

Key words. Monte Carlo methods, optimal stopping, sequential stopping rules, non-asymptotic.

AMS subject classifications. Primary 65C05; Secondary 62L12, 62L15

1. Introduction. Given i.i.d. random variables X_1, X_2, \dots the typical way of approximating their expected value $\mu = E[X]$ using M samples is the sample average

$$\bar{X}_M := \sum_{i=1}^M \frac{X_i}{M}.$$

We consider the objective of choosing M sufficiently large so that the error probability satisfies

$$P(|\bar{X}_M - \mu| > TOL) \leq \delta, \quad (1.1)$$

for some fixed small constants $TOL > 0$ and $\delta > 0$. Clearly, $P(|\bar{X}_M - \mu| > TOL)$ decreases as M increases, but at the same time the cost of computing \bar{X}_M increases. From an application and cost point of view it is therefore of interest to derive theory or algorithms that will give upper bounds on M satisfying (1.1) that are not far too large. When a-priori information about the distribution of X is available, for example if X is a bounded r.v. with an explicitly given bound, it is possible to derive good theoretical upper bounds for M using Hoeffding type inequalities, cf. Hoeffding [7]. In the general case when no or little information of the distribution is given, little theory is however known, and the typical way of estimating $E[X]$ using a sufficiently large number of samples M is through a sequential stopping rule. Below we give the general structure of the class of sequential stopping rules we have in mind.

- (I) Generate a batch of M i.i.d. samples X_1, X_2, \dots, X_M .
- (II) Infer distributive properties of \bar{X}_M from the generated batch of samples through higher order sample moments, e.g. the sample mean and the sample variance.

*This work was supported by the King Abdullah University of Science and Technology (KAUST) Strategic Research Initiative (SRI) Center for Uncertainty Quantification in Computational Science; the Center for Industrial and Applied Mathematics (CIAM) at Royal Institute of Technology, KTH; and Engineering and the VR project "Effektiva numeriska metoder för stokastiska differentialekvationer med tillämpningar". R. Tempone is a member of the KAUST SRI UQ Center.

[†]Weierstrass Institute, Mohrenstr. 39, 10117 Berlin, Germany (christian.bayer@wias-berlin.de).

[‡]Division of Mathematics, King Abdullah University of Science and Technology, Thuwal 23955-6900, Kingdom of Saudi Arabia (hakon.hoel@kaust.edu.sa, raul.tempone@kaust.edu.sa).

[§]Department of Numerical Analysis and Computer Science, KTH, SE-100 44, Stockholm, Sweden.

[¶]MATHICSE-CSQI, EPF de Lausanne, Switzerland (erik.vonschwerin@epfl.ch).

- (III) Based on the sample moments, estimate the error probability. When, based on the estimated probability, (1.1) is violated, increase the number of samples M and return to step (I).
Else, break and accept M .

Algorithm 1 Sample Variance Based Stopping Rule

Input: Number of samples M_0 , accuracy TOL , confidence δ , the cumulative distribution function of the standard normal distributed r.v. $\Phi(x)$.

Output: \bar{X}_M .

Set $n = 0$, generate M_n samples $\{X_i\}_{i=1}^{M_n}$ and compute the sample variance

$$\bar{\sigma}_{M_n}^2 := \frac{1}{M_n - 1} \sum_{i=1}^{M_n} (X_i - \bar{X}_{M_n})^2. \quad (1.2)$$

while $2\left(1 - \Phi(\sqrt{M_n}TOL/\bar{\sigma}_{M_n})\right) > \delta$ **do**

Set $n = n + 1$ and $M_n = 2M_{n-1}$.

Generate a batch of M_n i.i.d. samples $\{X_i\}_{i=1}^{M_n}$.

Compute the sample variance $\bar{\sigma}_{M_n}^2$ as given in (1.2).

end while

Set $M = M_n$, generate samples $\{X_i\}_{i=1}^M$ and compute the output sample mean \bar{X}_M . (See Section 2 for a motivation of the choice of the stopping criterion in the while loop above.)

Certainly the most natural and important representative of this class of algorithms is given in Algorithm 1. The algorithm estimates the error probability by appealing to the Central Limit Theorem (CLT). Consequently, it only relies on the sample variance in addition to the sample mean. In particular, the algorithm only requires mild additional assumptions on X , namely square integrability.

In the literature, various second moment based sequential stopping rules have been introduced to estimate the steady-state mean of stochastic processes, see for example Law and Kelton [10, 11] for comparisons of the performance of different stopping rules and Bratley, Bennet, and Fox [2] for an overview. Second moment based sequential stopping rules generally tend to perform well in the asymptotic regime when $TOL \rightarrow 0$. In fact, Chow and Robbins [3] proved that under very loose restrictions, second moment based sequential stopping rules such as Algorithm 1 are asymptotically consistent, meaning that for a fixed δ ,

$$\lim_{TOL \rightarrow 0} \mathbb{P}(|\bar{X}_M - \mu| > TOL) = \delta,$$

and in Glynn and Whitt [5] the consistency property is proven to hold for such stopping rules applied to more general stochastic processes. In the non-asymptotic regime – when TOL and δ are finite – Hickernell et al. [6] recently developed a second moment based MC method which guarantees to meet condition (1.1) under the assumption that an upper bound is given for the kurtosis of the r.v. to be sampled prior to sampling it. When no moment bounds are given prior to sampling, however, Bahadur and Savage [1] proved that it is not possible to develop an algorithm guaranteed to

meet condition (1.1). This is somewhat unsatisfactory as in applications the non-asymptotic regime with little prior information on the r.v. is a setting we believe is often encountered. While consistency is clearly a reassuring property in any case, in many situations one is in dire need of quantitative estimates of the error probability in the non-asymptotic regime, for instance when one tries to optimize the computational cost needed to meet a certain accuracy target using an adaptive algorithm. We could not find such a quantitative, non-asymptotic analysis of sequential stopping algorithms like Algorithm 1 in the literature.

In this work we demonstrate by numerical examples that second moment based stopping rules can fail convincingly in the non-asymptotic regime, especially when the underlying distribution X is heavy-tailed, see Section 2. We proceed by giving an error analysis of Algorithm 1 specifically in the non-asymptotic regime. We note a-priori that there are two obvious approximation errors in the underlying assumptions of Algorithm 1:

- (I) The algorithm appeals to the CLT to approximate the tail probabilities for \bar{X}_M even though M is finite.
- (II) In doing so, it uses the sample variance $\bar{\sigma}_M^2$ instead of the true variance σ^2 .

To get a hold on the error probability (1.1) despite the fact that the distribution of the sample mean \bar{X}_M is unknown, we again appeal to the central limit theorem, but we adjust the estimate by adding a Berry-Esseen type term, which extends the validity of the estimate to the non-asymptotic case, thereby dealing with the first approximation error. As the error probability (1.1) is a tail probability for the distribution of the sample mean and the Berry-Esseen theorem itself is rather aimed at being sharp at the center of the distribution, we appeal to non-uniform versions of the Berry-Esseen theorem, see Theorem 1.1 and Corollary 1.2 below. However, both intuition and numerical tests suggest that the approximation of the tail probabilities by the non-uniform Berry-Esseen theorem is far too pessimistic at least when the second approximation error is small, i.e., when the computed sample variance is actually close to the true variance. In this case, we adjust the normal distribution by a less stringent extra term, which is obtained from an Edgeworth expansion of the distribution function of the sample mean \bar{X}_M , cf. Feller [4].¹

Having identified possible origins of failure of Algorithm 1, we propose an improvement of Algorithm 1. However, this variant requires third and fourth sample moments, see Section 4. Finally, in Section 5, we test the new algorithm numerically. We find that the new stopping Algorithm 2 indeed satisfies the desired confidence level δ on the error probability (1.1) even when $\delta \ll TOL$.

As already discussed above, we need to approximate the unknown distribution of a sample mean in a general, non-asymptotic regime. The uniform and non-uniform

¹Note that here we are introducing a gap in the analysis: the estimate based on the non-uniform Berry-Esseen theorem is reliable in the sense that it always leads to an upper bound of the error probability (1.1). For the Edgeworth expansion, however, there might be situations when the true error probability is underestimated, and, consequently, the accuracy target might still be missed. Numerical evidence, however, suggests that the estimate obtained from solely relying on the non-uniform Berry-Esseen theorem is usually by orders of magnitude too pessimistic. Apart from intrinsic reasons, one reason might be that the constants known in the non-uniform Berry-Esseen theorems might be far from being optimal. In the end, we think that the above compromise between Berry-Esseen type estimations and estimations based on the Edgeworth expansion might be a good compromise which retains the goal of reliably meeting the accuracy target – except maybe for very extreme situations – while keeping a certain level of efficiency. We note, however, that it is also possible to construct even more conservative stopping rules which are only based on the Berry-Esseen theorem.

Berry-Esseen theorems provide quantitative bounds for the difference between the true distribution of the sample mean and its asymptotic limit, namely the normal distribution. The following Berry-Esseen bounds, which can be found in their classical formulations in Petrov [13], is here presented with the present optimal bound constants.

THEOREM 1.1 (Uniform and Non-Uniform Berry-Esseen). *Suppose X_1, X_2, \dots are i.i.d. r.v. with $E[X] = 0$, $\sigma^2 = \text{Var}(X)$ and $\beta = \frac{E[|X|^3]}{\sigma^3} < \infty$. Then the following uniform bound*

$$\left| \mathbb{P}\left(\frac{1}{\sigma\sqrt{n}} \sum_{i=1}^n X_i < x\right) - \Phi(x) \right| \leq 0.3328 \cdot \frac{\beta + 0.429}{\sqrt{n}}$$

holds, cf. Shevtsova [14]. Furthermore, the following non-uniform bound holds

$$\left| \mathbb{P}\left(\frac{1}{\sigma\sqrt{n}} \sum_{i=1}^n X_i < x\right) - \Phi(x) \right| \leq 18.1139 \cdot \frac{\beta}{\sqrt{n}(1 + |x|^3)},$$

cf. Hickernell et al. [6], and the references therein. For the purpose of this work, it will be useful to combine the uniform and non-uniform Berry-Esseen bound as follows.

COROLLARY 1.2 (Berry-Esseen). *Suppose X_1, X_2, \dots are i.i.d. r.v. with $E[X] = 0$. Then*

$$\left| \mathbb{P}\left(\frac{1}{\sigma\sqrt{n}} \sum_{i=1}^n X_i < x\right) - \Phi(x) \right| \leq \frac{C_{\text{BE}}(x, \beta)}{\sqrt{n}}$$

where the bound function $C_{\text{BE}} : \mathbb{R} \rightarrow [0, C_0]$ is defined by

$$C_{\text{BE}}(x, \beta) := \min\left(0.3328 \cdot (\beta + 0.429), 18.1139 \cdot \frac{\beta}{(1 + |x|^3)}\right).$$

In the asymptotic regime, the distribution of $\mathbb{P}\left(\frac{1}{\sigma\sqrt{n}} \sum_{i=1}^n X_i < x\right)$ can be expressed by so called Edgeworth expansions. Here we present the one-term Edgeworth expansion.

THEOREM 1.3 (Edgeworth expansion, cf. Feller [4] p. 541). *Suppose X_1, X_2, \dots are i.i.d. r.v. with a distribution which is not a lattice distribution and $E[X] = 0$, $\sigma^2 = \text{Var}(X)$ and $E[X^3] < \infty$. Then*

$$\mathbb{P}\left(\frac{1}{\sigma\sqrt{n}} \sum_{i=1}^n X_i < x\right) = \Phi(x) + \frac{(x^2 - 1)e^{-x^2/2}}{6\sqrt{2\pi n}} \frac{E[X^3]}{\sigma^3} + o(n^{-1/2}),$$

uniformly for $x \in \mathbb{R}$.

2. Stopping rule failures. Suppose we seek to estimate $\mu = E[X]$ using Monte Carlo simulation and we actually *do know* the variance $\sigma^2 = \text{Var}(X)$. As before, our objective is to achieve $\mathbb{P}(|\bar{X}_M - \mu| > \text{TOL}) \leq \delta$, for some fixed, small constants $\text{TOL}, \delta > 0$. The CLT motivates the stopping rule

$$M = \frac{C_{\text{CLT}}^2 \sigma^2}{\text{TOL}^2}, \quad C_{\text{CLT}} := \Phi^{-1}\left(\frac{2 - \delta}{2}\right), \quad (2.1)$$

which would exactly fulfill our objective (1.1) in the asymptotic regime $M \rightarrow \infty$. Of course, this conflicts with our choice (2.1) for M , since we treat δ and TOL as finite constants. However, we can still estimate the probability in (1.1) using Corollary 1.2 and obtain

$$\begin{aligned}
\mathbb{P}(|\bar{X}_M - \mu| > TOL) &= \mathbb{P}\left(\sqrt{M} \frac{|\bar{X}_M - \mu|}{\sigma} > \frac{\sqrt{MTOL}}{\sigma}\right) \\
&\leq 2 \left(1 - \Phi\left(\frac{\sqrt{MTOL}}{\sigma}\right)\right) + 2C_{\text{BE}}\left(\frac{\sqrt{MTOL}}{\sigma}, \beta\right) \frac{1}{\sqrt{M}} \\
&= 2(1 - \Phi(C_{\text{CLT}})) + 2C_{\text{BE}}(C_{\text{CLT}}, \beta) \frac{1}{\sqrt{M}} \\
&= \delta + 2 \frac{C_{\text{BE}}(C_{\text{CLT}}, \beta)}{\sigma C_{\text{CLT}}} TOL,
\end{aligned} \tag{2.2}$$

where $\beta = \frac{\mathbb{E}[|X - \mu|^3]}{\sigma^3}$. This means that in the worst case, the actual error probability could be $\delta + \mathcal{O}(TOL)$ instead of δ .² For instance, in situations where the statistical confidence in the result is more stringent than the accuracy so that $\delta \ll TOL$, the asymptotically motivated choice of M in (2.1) could, granted the bound (2.2) is sharp, fail to deliver the expected level of confidence. For most r.v. however, the bound (2.2) is far too conservative, and one might ask whether it is reasonable to fear underestimating the error probability in the fashion we have described. The following numerical example shows the existence of r.v. for which the stopping rule (2.1) fails in the non-asymptotic regime

EXAMPLE 2.1. *The heavy-tailed Pareto-distribution has the probability distribution function*

$$f(x) = \begin{cases} \alpha x_m^\alpha x^{-(\alpha+1)} & \text{if } x \geq x_m \\ 0 & \text{else,} \end{cases} \tag{2.3}$$

where $\alpha, x_m \in \mathbb{R}_+$ are respectively the shape and the scale parameter. The moments of $\mathbb{E}[X^n]$ for the Pareto r.v. only exists for $n < \alpha$ and, supposing $\alpha > 2$, its mean and variance are given by

$$\mu = \frac{\alpha x_m}{\alpha - 1} \text{ and } \sigma^2 = \frac{x_m^2 \alpha}{(\alpha - 1)^2 (\alpha - 2)}.$$

It is further easy to derive that for a Pareto r.v. with $\alpha = 3 + \gamma$ and $0 < \gamma < 1$,

$$\beta = \frac{\mathbb{E}[|X - \mu|^3]}{\sigma^3} = \mathcal{O}(\gamma^{-1}).$$

This implies that there exists r.v. for which the second summand of the bound (2.2) can become arbitrary large. So for such r.v. the stopping rule (2.1) might fail. Let us investigate by numerical approximations. Considering the distribution with $\alpha = 3.1$ (and $x_m = 1$), yields a heavy-tailed r.v. with known mean, variance and third moment.

²Note that C_{CLT} as defined in (2.1) grows only very slowly as δ decreases, since we have $C_{\text{CLT}} < \sqrt{2 \log(\delta^{-1})}$. So as a function of δ , the factor in front of TOL is approximately constant.

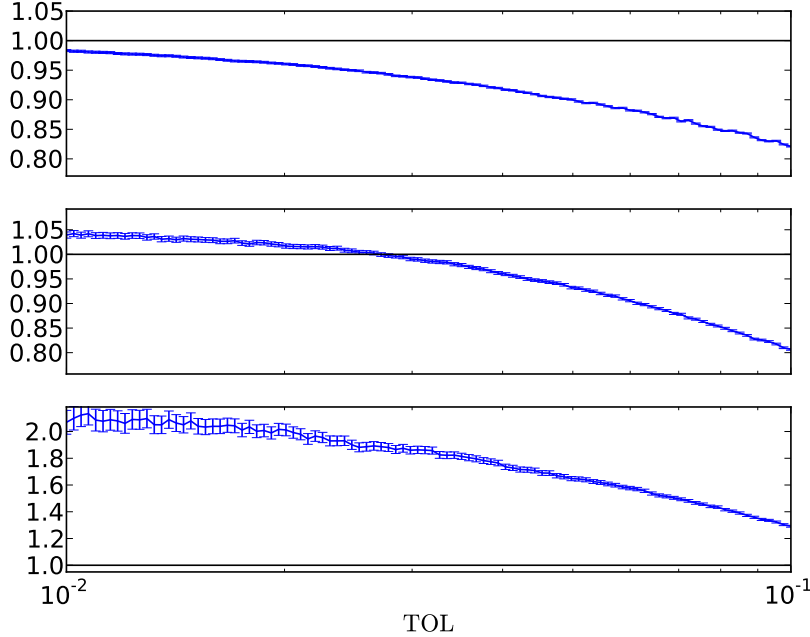


FIG. 2.1. MC estimate using the stopping rule (2.1) for *i.i.d.* Pareto r.v. with parameters $\alpha = 3.1$ and $x_m = 1$. The quantity $\mathbb{P}(|\bar{X}_M - \mu| > TOL) / \delta$ (blue line) is plotted for the settings $\delta(TOL) = TOL^\ell$ with $\ell = 0.5$ (top), $\ell = 1$ (middle), and $\ell = 2$ (bottom). The probability of failure is estimated by $\bar{p}_N(TOL) = N^{-1} \sum_{i=1}^N \mathbf{1}_{|\bar{X}_M(\omega_i) - \mu| > TOL}$ with $N = 10^7$, and the error bars for the estimate of \bar{p}_N by $1.96 \sqrt{\bar{p}_N(1 - \bar{p}_N)/N}$. The stopping rule fails when $\mathbb{P}(|\bar{X}_M - \mu| > TOL) / \delta > 1$. We observe that the smaller δ is relative to TOL , the larger is the probability of failure for the stopping rule.

For a set of accuracies $TOL \in [0.05, 0.2]$ and confidences $\delta = TOL^\ell$, $\ell = 1$ and 2 we have numerically approximated $\mathbb{P}(|\bar{X}_M - \mu| > TOL) \leq \delta$ using, in accordance with (2.1), the stopping rule

$$M = \left\lceil \frac{C_{\text{CLT}}^2 \sigma^2}{TOL^2} \right\rceil$$

The results, illustrated in Figure 2.1, show that when $\delta \ll TOL$, the sought confidence is far from met.

The demonstrated stopping rule failure motivated us to study and develop ways of constructing more reliable stopping rules. In Section 3, we first analyze the stopping rule of Algorithm 1, and derive an approximate upper bound for the failure probability expressed in terms of M, TOL and δ . In Section 4, we develop a more reliable stopping rule algorithm, which in addition to sampled second moments of the r.v. in question also depends on sampled third and fourth moments. The paper is concluded with numerical examples comparing the reliability and computational cost of Algorithm 1 with the stopping rule developed in Section 4.

3. Error analysis for Algorithm 1. Example 2.1 illustrate that for some r.v. the stopping rule in Algorithm 1 does not meet the accuracy-confidence con-

straint (1.1). To construct a more reliable stopping rule, penalty terms have to be added to the stopping criterion in Algorithm 1. Some care should be taken to make the penalty terms of the right size: if too large penalties are added, the new stopping rule will be reliable but very inefficient, while if too small penalty terms are added, the algorithm will of course be efficient but unreliable.

In this section, we first derive an approximate upper bound for the failure probability

$$\mathbb{P}\left(|\bar{X}_M - \mu| > TOL \mid M\right) \quad (3.1)$$

corresponding to the stopping rule of Algorithm 1 conditional on the (random) final number of samples M . Clearly, the bound for (3.1) will also be a r.v. Using the bound for (3.1), we thereafter construct reasonable penalty terms to be added to the stopping criterion of our new stopping rule.

As a general idea, we are going to use a weighted average of one very reliable, but typically overly pessimistic error bound based on the Berry-Esseen theorem, and another error bound based on the Edgeworth expansion, which is typically too optimistic. The (critical) choice of the weights is based on the following consideration: If we have successfully estimated the variance of X to high accuracy, then we give ourselves some leeway for the bias and choose the more optimistic bound. If, on the other hand, we have already mis-estimated the variance σ^2 , then we want to be highly conservative in our conditional bias estimate. Thus, our first step towards an upper bound for (3.1) is partitioning the probability (3.1) into two parts as follows

$$\begin{aligned} \mathbb{P}\left(|\bar{X}_M - \mu| > TOL \mid M\right) &= \mathbb{P}\left(|\bar{X}_M - \mu| > TOL \mid M\right) \mathbb{P}\left(|\bar{\sigma}_M^2 - \sigma^2| \geq \sigma^2/2 \mid M\right) \\ &\quad + \mathbb{P}\left(|\bar{X}_M - \mu| > TOL \mid M\right) \mathbb{P}\left(|\bar{\sigma}_M^2 - \sigma^2| < \sigma^2/2 \mid M\right). \end{aligned} \quad (3.2)$$

In the next step, we will use two different (asymptotic) error bounds for $\mathbb{P}\left(|\bar{X}_M - \mu| > TOL \mid M\right)$: a conservative one weighted by $\mathbb{P}\left(|\bar{\sigma}_M^2 - \sigma^2| \geq \sigma^2/2 \mid M\right)$ and an efficient one weighted by $\mathbb{P}\left(|\bar{\sigma}_M^2 - \sigma^2| < \sigma^2/2 \mid M\right)$. To derive the conservative error bound, we recall that in Algorithm 1 the samples used in the output estimate \bar{X}_M and for $\bar{\sigma}_M$ are independent of the samples used to determine M . Keeping this in mind, we derive the following approximate upper bound

$$\begin{aligned} \mathbb{P}\left(|\bar{X}_M - \mu| > TOL \mid M\right) &= \mathbb{P}\left(|\bar{X}_n - \mu| > TOL\right) \Big|_{n=M} \\ &\lesssim 2 \left(1 - \Phi\left(\frac{\sqrt{MTOL}}{\sigma}\right) + C_{\text{BE}}\left(\frac{\sqrt{MTOL}}{\sigma}, \beta\right) \frac{1}{\sqrt{M}} \right). \end{aligned} \quad (3.3)$$

Here the Berry-Esseen bound of Corollary 1.2 was used to derive the approximate bound of the last line.

For the weighting factor, we obtain the following equality

$$\mathbb{P}\left(|\bar{\sigma}_M^2 - \sigma^2| \geq \sigma^2/2 \mid M\right) = \mathbb{P}\left(|\bar{\sigma}_n^2 - \sigma^2| \geq \sigma^2/2\right) \Big|_{n=M}.$$

Furthermore, using Chebycheff's inequality for k-statistics to bound the variance of the sample variance, cf. Keeping [9], we derive that

$$\begin{aligned} \mathbb{P}\left(|\bar{\sigma}_M^2 - \sigma^2| \geq \sigma^2/2 \mid M\right) &= \mathbb{P}\left(|\bar{\sigma}_n^2 - \sigma^2| \geq \sigma^2/2\right) \Big|_{n=M} \\ &\leq 4 \mathbb{E} \left[\frac{|\sigma^2 - \bar{\sigma}_n^2|^2}{\sigma^4} \right] \Big|_{n=M} \leq 4 \frac{\sigma^4 \left(\frac{2}{M-1} + \frac{\kappa}{M} \right)}{\sigma^4} = 4 \left(\frac{2}{M-1} + \frac{\kappa}{M} \right). \end{aligned}$$

Here κ denotes the *kurtosis*, i.e.,

$$\kappa = \frac{\mathbb{E}[|X - \mu|^4]}{\sigma^4} - 3.$$

More generally, we might consider Markov's inequality for a different p than two – for instance, because the kurtosis fails to be finite. Then we can use

$$\mathbb{P}\left(|\bar{\sigma}_M^2 - \sigma^2| \geq \sigma^2/2 \mid M\right) \leq 2^p \frac{\mathbb{E}[|\bar{\sigma}_M^2 - \sigma^2|^p \mid M]}{\sigma^{2p}}.$$

In particular, we are free to choose $p < 2$, if the fourth moment of the distribution might not exist. We conclude that

$$\begin{aligned} \mathbb{P}\left(|\bar{\sigma}_M^2 - \sigma^2| \geq \sigma^2/2 \mid M\right) &\leq \min \left\{ 1, 4 \left(\frac{2}{M-1} + \frac{\kappa}{M} \right), \frac{2^{p_1}}{\sigma^{2p_1}} \mathbb{E}\left[|\bar{\sigma}_M^2 - \sigma^2|^{p_1} \mid M\right], \right. \\ &\quad \left. \dots, \frac{2^{p_n}}{\sigma^{2p_n}} \mathbb{E}\left[|\bar{\sigma}_M^2 - \sigma^2|^{p_n} \mid M\right] \right\} \\ &=: C_P(2, p_1, \dots, p_n; M), \end{aligned} \tag{3.4}$$

for some finite sequence $p_1, \dots, p_n > 0$, with the natural interpretation of $\kappa = \infty$, if fourth moments do not exist.

Next, we come to the optimistic bound of the error probability

$$\mathbb{P}\left(|\bar{X}_M - \mu| > TOL \mid M\right)$$

based on the Edgeworth expansion given in Theorem 1.3. The mild penalty obtained in this way is of the form

$$\begin{aligned} &\mathbb{P}\left(|\bar{X}_M - \mu| > TOL \mid M\right) \\ &= \mathbb{P}\left(\sqrt{n} \frac{|\bar{X}_n - \mu|}{\sigma} > \frac{\sqrt{n} TOL}{\sigma}\right) \Big|_{n=M} \\ &\lesssim 2 \left(1 - \Phi\left(\frac{\sqrt{M} TOL}{\sigma}\right) + \frac{\left| \frac{M TOL^2}{\sigma^2} - 1 \right| \exp\left(-\frac{M TOL^2}{\sigma^2}\right) |\mathbb{E}[(X - \mu)^3]|}{6\sqrt{2\pi} M \sigma^3} \right). \end{aligned} \tag{3.5}$$

Combining (3.3), (3.4) and (3.5), and noting that for all $x \in \mathbb{R}_+$ and $n \in \mathbb{N}$,

$$\frac{|x^2 - 1| e^{-x^2/2} |\mathbb{E}[(X - \mu)^3]|}{6\sqrt{2\pi} n \sigma^3} \leq C_{BE}(x, \beta) \frac{1}{\sqrt{n}},$$

we obtain the following approximate bound for failing to meet the accuracy of Algorithm 1 conditioned on the stopped number of samples M :

$$\begin{aligned} \mathbb{P}\left(|\bar{X}_M - \mu| > TOL \mid M\right) &\lesssim 2 \left(1 - \Phi\left(\frac{\sqrt{MTOL}}{\sigma}\right)\right) + \\ &+ 2C_{BE}\left(\frac{\sqrt{MTOL}}{\sigma}, \beta\right) \frac{1}{\sqrt{M}} C_P(2, p_1, \dots, p_n; M) \\ &+ \frac{\left|\frac{MTOL^2}{\sigma^2} - 1\right| |\mathbb{E}[(X - \mu)^3]|}{\exp\left(\frac{MTOL^2}{2\sigma^2}\right) \times 3\sqrt{2\pi M}\sigma^3} (1 - C_P(2, p_1, \dots, p_n; M)). \end{aligned} \quad (3.6)$$

REMARK 3.1. *The bound from equation (3.6) is exact up to the higher order terms coming from the Edgeworth expansion. Note that $1 - C_P(2, p_1, \dots, p_n; M)$ is not a lower bound for the complimentary error probability $\mathbb{P}\left(|\bar{\sigma}_M^2 - \sigma^2| < \sigma^2/2 \mid M\right)$, so we are even more conservative in our approach.*

4. A higher moments based stopping rule. From the approximate stochastic error bound (3.6) we will in this section construct a new, more reliable stopping rule with a stopping criterion based on second, third, and fourth moments of the r.v. that is sampled. The sampled moments our new algorithm will depend on are (here represented in biased form)

$$\begin{aligned} \bar{\sigma}_M &:= \sqrt{\sum_{i=1}^M \frac{(X_i - \bar{X}_M)^2}{M}}, & \bar{\beta}_M &:= \sum_{i=1}^M \frac{|X_i - \bar{X}_M|^3}{M\bar{\sigma}_M^3}, \\ \hat{\beta}_M &:= \sum_{i=1}^M \frac{(X_i - \bar{X}_M)^3}{M\bar{\sigma}_M^3}, & \text{and } \bar{\kappa}_M &:= \sum_{i=1}^M \frac{(X_i - \bar{X}_M)^4}{M\bar{\sigma}_M^4} - 3. \end{aligned} \quad (4.1)$$

Replacing moments with sample moments in (3.6), we obtain a computable approximate stochastic error bound

$$\begin{aligned} &\mathbb{P}\left(|\bar{X}_M - \mu| > TOL \mid M\right) \\ &\lesssim 2 \left(1 - \Phi\left(\frac{\sqrt{MTOL}}{\bar{\sigma}_M}\right)\right) + 2C_{BE}\left(\frac{\sqrt{MTOL}}{\bar{\sigma}_M}, \bar{\beta}_M\right) \frac{1}{\sqrt{M}} \bar{C}_P(2, p_1, \dots, p_n; M) \\ &+ \frac{\left|\frac{MTOL^2}{\bar{\sigma}_M^2} - 1\right| |\hat{\beta}_M|}{\exp\left(\frac{MTOL^2}{2\bar{\sigma}_M^2}\right) \times 3\sqrt{2\pi M}\bar{\sigma}_M^3} (1 - \bar{C}_P(2, p_1, \dots, p_n; M)). \end{aligned} \quad (4.2)$$

The resulting approximate stochastic error bound will be implemented as the following stopping criterion in Algorithm 2:

$$\begin{aligned} &2 \left(1 - \Phi\left(\frac{\sqrt{MTOL}}{\bar{\sigma}_M}\right)\right) + 2C_{BE}\left(\frac{\sqrt{MTOL}}{\bar{\sigma}_M}, \bar{\beta}_M\right) \frac{1}{\sqrt{M}} \bar{C}_P(2, p_1, \dots, p_n) \\ &+ \frac{\left|\frac{MTOL^2}{\bar{\sigma}_M^2} - 1\right| |\hat{\beta}_M|}{\exp\left(\frac{MTOL^2}{2\bar{\sigma}_M^2}\right) \times 3\sqrt{2\pi M}\bar{\sigma}_M^3} (1 - \bar{C}_P(2, p_1, \dots, p_n)) < \delta, \end{aligned} \quad (4.3)$$

where $\overline{C}_P(2, p_1, \dots, p_n)$ is obtained from $C_P(2, p_1, \dots, p_n)$ by replacing all moments of $\overline{\sigma}$ by their empirical counterparts.

We now present the new stopping rule algorithm.

Algorithm 2 Higher Moments Based Stopping Rule

Input: Accuracy TOL , confidence δ , and initial number of samples M_0 .

Output: \overline{X}_M .

Set $n = 0$, generate i.i.d. samples $\{X_i\}_{i=1}^{M_n}$ and compute the sample moments $\overline{\sigma}_{M_n}$, $\overline{\beta}_{M_n}$, $\widehat{\beta}_{M_n}$ and $\overline{\kappa}_{M_n}$ and all (other) moments needed for \overline{C}_P according to (4.1).

while Inequality (4.3) is not fulfilled. **do**

Set $n = n + 1$ and $M_n = 2M_{n-1}$.

Generate M_n i.i.d. samples $\{X_i\}_{i=1}^{M_n}$ and compute the sample moments $\overline{\sigma}_{M_n}$, $\overline{\beta}_{M_n}$, and $\overline{\kappa}_{M_n}$ and all (other) moments needed for \overline{C}_P .

end while

Set $M = M_n$, generate i.i.d. samples $\{X_i\}_{i=1}^M$ and return the sample mean \overline{X}_M .

REMARK 4.1. *The reasoning in Section 3 could also be used to construct alternative stopping rule algorithms. For instance, instead of using the respective probabilities for misestimation of σ , $\mathbb{P}\left(|\overline{\sigma}_M^2 - \sigma^2| \geq \sigma^2/2 \mid M\right)$ and $\mathbb{P}\left(|\overline{\sigma}_M^2 - \sigma^2| < \sigma^2/2 \mid M\right)$ as respective weights for the conservative and the optimistic bound for the misestimation error probability for the mean μ itself, we could also take the argument more literally and use the pessimistic error bound if we have misestimated the variance and the optimistic one otherwise, i.e., using random indicator functions instead of their expectations as weights. We leave this as a remark.*

5. Numerical experiments. In the numerical experiments we will estimate the mean of four differently distributed r.v. by using three different stopping rules: the sample variance based stopping rule in Algorithm 1, the new higher moments based stopping rule in Algorithm 2 and the stopping rule recently introduced by Hickernell et al. [6] which they prove performs reliably, with guaranteed conservative fixed width confidence intervals, provided an upper bound for the kurtosis of the sampled r.v. is known prior to sampling.

5.1. Experimental setup and implementational details. The distributions of the r.v. we will consider will be the light-tailed Uniform, the Exponential, the heavy-tailed Normal-inverse Gaussian and the heavy-tailed Pareto distribution. The initial number of samples used in Algorithm 1 and 2 is set to $M_0 = 32$ for all experiments. The version of Algorithm 2 we will implement for these experiments uses the only sample moments defined in (4.1), i.e., $C_P = C_P(2; M)$. For the implementation of Hickernell et al.’s algorithm an upper bound on the kurtosis is required, so we will assume the exact kurtosis of the r.v. is known prior to sampling and set the “variance inflation factor” $\mathfrak{C} = 1.1$, cf. [6]. Let us stress that the assumption of knowing the exact value of kurtosis prior to sampling is only made in the runs for Hickernell et al.’s algorithm, and due to this, its performance is presented in a slightly more favorable light than Algorithm 1 and 2’s.

For each of the three algorithms we will numerically estimate the probability of

failure $P(|\bar{X}_M - \mu| > TOL)$ by the MC outer loop

$$\bar{p}_N(TOL, \delta) := N^{-1} \sum_{i=1}^N \mathbf{1}_{|\bar{X}_M(\omega_i) - \mu| > TOL} \quad (5.1)$$

on a 100×100 grid of $(TOL_i, \delta_j) \in [10^{-2}, 10^{-1}] \times [10^{-3}, 10^{-1}]$ values where $TOL_i = 10^{-(1+i/99)}$ and $\delta_i = 10^{-(1+2i/99)}$ for $i = 0, 1, \dots, 99$. The number of outer loop samples N we are able to use in an experiment will depend on both the cost of sampling the given r.v. and the cost of the used stopping algorithm. The error in estimating the probability of failure by $\bar{p}_N(TOL, \delta)$ is estimated by the standard deviation approximation of a sampled Bernoulli r.v.

$$\frac{\sqrt{\bar{p}_N(TOL, \delta)(1 - \bar{p}_N(TOL, \delta))}}{\sqrt{N}}.$$

At (TOL_i, δ_j) points where for a given algorithm $\bar{p}_N(TOL_i, \delta_j)/\delta_j > 1$, that algorithm is considered unreliable. So to visualize the domain of reliability for an algorithm we plot the function \bar{p}_N/δ . Furthermore, Table 5.1 provides the following values

$$\max_{i,j=0,1,\dots,99} \frac{\bar{p}_N(TOL_i, \delta_j)}{\delta_j}, \quad \text{and} \quad \max_{i,j=0,1,\dots,99} \frac{\sqrt{\bar{p}_N(TOL_i, \delta_j)(1 - \bar{p}_N(TOL_i, \delta_j))}}{\sqrt{N}\delta_j} \quad (5.2)$$

as estimates of an algorithm's maximum unreliability and the uncertainty in the unreliability estimate, respectively.

The computer code for the algorithms is written in the Java programming language and uses the "Stochastic Simulation in Java" library to sample the r.v. in parallel with the LFSR258 pseudo random number generator, cf. [12]. The experiments were run in parallel using 10 threads on an Intel Xeon(R) CPU X5650, 2.66 GHz, and the plots were made in Python using the open source plotting library Matplotlib, cf. [8].

5.2. Results. Figures 5.1, 5.2, 5.3, and 5.4 visualize the results of the numerical experiments and Table 5.1 contains numerical and parameter values used in the experiments.

For the Pareto distributed r.v. considered in Figure 5.1, we see that while Algorithm 2 is reliable everywhere, Algorithm 1 is not reliable in the region where TOL is large and $\delta \ll TOL$. Furthermore, the complexity of Algorithm 2 is only slightly higher than Algorithm 1's. (Hickernell et al.'s algorithm requires an upper bound on the kurtosis prior to sampling to be applicable, so that algorithm is not applicable in this experiment.)

For the Normal-inverse Gaussian distributed r.v. studied in Figure 5.2, we see that Algorithm 1 is unreliable in the region where TOL is large. Algorithm 2 and Hickernell et al.'s algorithm are on the other hand both very reliable. Hickernell et al.'s algorithm does however seem to be more reliable than required, and due to this overkill the complexity of Hickernell et al.'s algorithm is much higher than Algorithm 2's. The Normal-inverse Gaussian distributed r.v. is quite expensive to generate computationally, therefore we have had to reduce the number of outer samples N substantially in this experiment.

For the uniformly and exponentially distributed r.v. studied in Figure 5.3 and 5.4, respectively, all three algorithms are reliable and Algorithm 1 and 2 perform very

similarly in terms of reliability. We also see that for both of these experiments the complexity measured in terms of the average number of r.v. used to generate \bar{X}_M is very similar for Algorithm 1 and 2. This observation is at odds with the corresponding measured average computer run times given in Table 5.1 which show Algorithm 2's average run time is considerably longer than Algorithm 1's. The reason for the disagreement between the complexity measurements and the run time measurements is that for uniformly distributed r.v. (and to some weaker degree also for exponentially distributed r.v.) the generation of r.v. is not so costly for the computer program, and therefore do also other arithmetical operations contribute considerably to the run time of the computer program in these experiments.

We further note that the complexity contours of Algorithm 1 and 2 in the figures seem proportional to

$$\widetilde{M}(TOL, \delta) = \left(\frac{\Phi^{-1}(1 - \delta/2)}{TOL} \right)^2,$$

up to rounding. This property, which by the CLT is expected in the asymptotic regime, can be verified by comparing the contour of the function $\widetilde{M}(TOL, \delta)$ to the complexity contour of Algorithm 1 and 2, respectively. The sharp changes in complexity that we observe in the contour plots for Algorithm 1 and 2, especially visible in Figure 5.3 and 5.4, are due to the doubling procedure ($M_n = 2M_{n-1}$) which is implemented in these two algorithms.

Considered r.v.	Algorithm	Algorithm performance			
		N	$\max \bar{p}_N / \delta$	$\max(\bar{p}_N(1 - \bar{p}_N)/N)^{1/2} / \delta$	runtime/ N
Pareto, cf. Figure 5.1 ($\sigma = 1$ and $\kappa = \infty$)	Alg 1	$5 \cdot 10^6$	3.529600	0.026522	0.124712 s
	Alg 2	$5 \cdot 10^6$	0.686309	0.009506	0.189077 s
Normal-inv. Gaussian cf. Figure 5.2 ($\sigma = 1, \kappa = 123$)	Alg 1	$5 \cdot 10^5$	12.014000	0.154076	3.019442 s
	Alg 2	$5 \cdot 10^5$	0.755712	0.028699	3.048026 s
	Hick. et al.'s	$5 \cdot 10^4$	0.000419	0.000296	77.322941 s
Uniform $U(-\sqrt{3}, \sqrt{3})$, cf. Figure 5.3 ($\sigma = 1$ and $\kappa = -6/5$)	Alg 1	$5 \cdot 10^6$	0.970623	0.013598	0.008996 s
	Alg 2	$5 \cdot 10^6$	0.970600	0.013597	0.048990 s
	Hick. et al.'s	$5 \cdot 10^6$	0.278880	0.001089	0.097067 s
Exponential $\lambda = 1$, cf. Figure 5.4 ($\sigma = 1$ and $\kappa = 6$)	Alg 1	$5 \cdot 10^6$	0.919659	0.012642	0.033680 s
	Alg 2	$5 \cdot 10^6$	0.912216	0.012040	0.076942 s
	Hick. et al.'s	$5 \cdot 10^6$	0.206530	0.000676	0.229495 s

TABLE 5.1

We compare the reliability and complexity of three MC algorithms approximating the mean of given r.v. The function $\bar{p}_N(TOL, \delta)$ is an estimate of the probability of failure and N is the number samples used to compute \bar{p}_N , cf. (5.1). The values $\max \bar{p}_N / \delta$ and $\max(\bar{p}_N(1 - \bar{p}_N)/N)^{1/2} / \delta$ are estimates of the maximum unreliability and its uncertainty, cf. (5.2). The expression runtime/ N is an estimate of the average time it takes for one thread of the computer program to compute a single output \bar{X}_M at all points (TOL_i, δ_j) , $i, j = 0, 1, \dots, 99$.

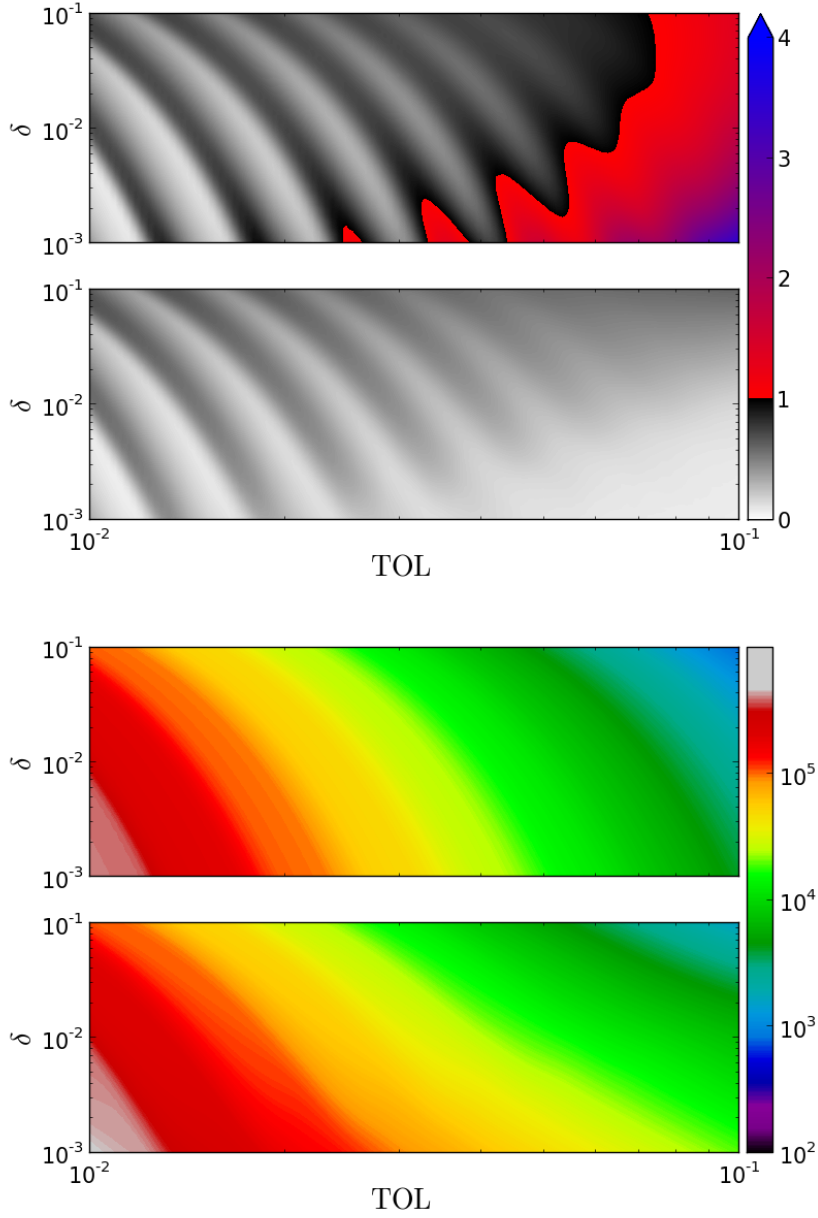


FIG. 5.1. (**Pareto Distribution**) We sample a Pareto distributed r.v. with parameters $\alpha = 3.1$ and $x_m = (\alpha - 1)\sqrt{(\alpha - 2)/\alpha}$. Top two plots: Plots of $\bar{p}_N(TOL, \delta)/\delta$, cf. (5.1), where \bar{X}_M respectively is computed by Algorithm 1 in the top plot, Algorithm 2 in the bottom plot. We consider the algorithm unreliable at (TOL, δ) points where $\bar{p}_N/\delta > 1$. The bottom two plots are corresponding plots of the respective algorithms' complexity in terms of the average number of r.v. samples required to generate \bar{X}_M at the given (TOL, δ) .

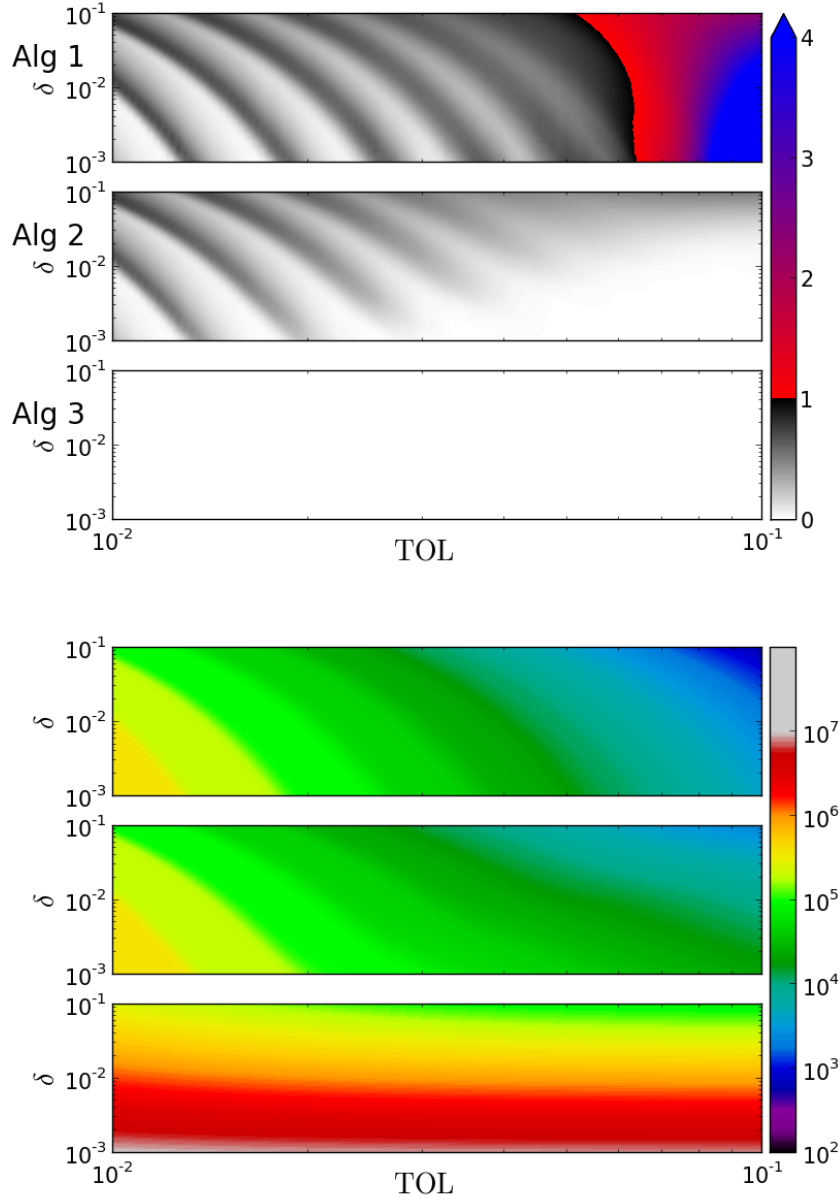


FIG. 5.2. (*Normal-inverse Gaussian Distribution*) We sample a Normal-inverse Gaussian distributed r.v. with parameters $\alpha = 3$, $\beta = \sqrt{\alpha^2 - 1}$, $\gamma = 1$, $\delta = \alpha^{-2}$, and $\mu = -\beta/\gamma$. This yields a r.v. with standard deviation $\sigma = 1$ and $\kappa = 123$. Top three plots: Plots of $\bar{p}_N(TOL, \delta)/\delta$, cf. (5.1), where samples of \bar{X}_M respectively are computed by Algorithm 1 in the top plot, Algorithm 2 in the middle plot, and Hickernell et al.'s algorithm in the bottom plot. We consider the algorithm unreliable at (TOL, δ) points where $\bar{p}_N/\delta > 1$. The bottom three plots are corresponding plots of the respective algorithms' complexity in terms of the average number of r.v. samples required to generate \bar{X}_M at the given (TOL, δ) point.

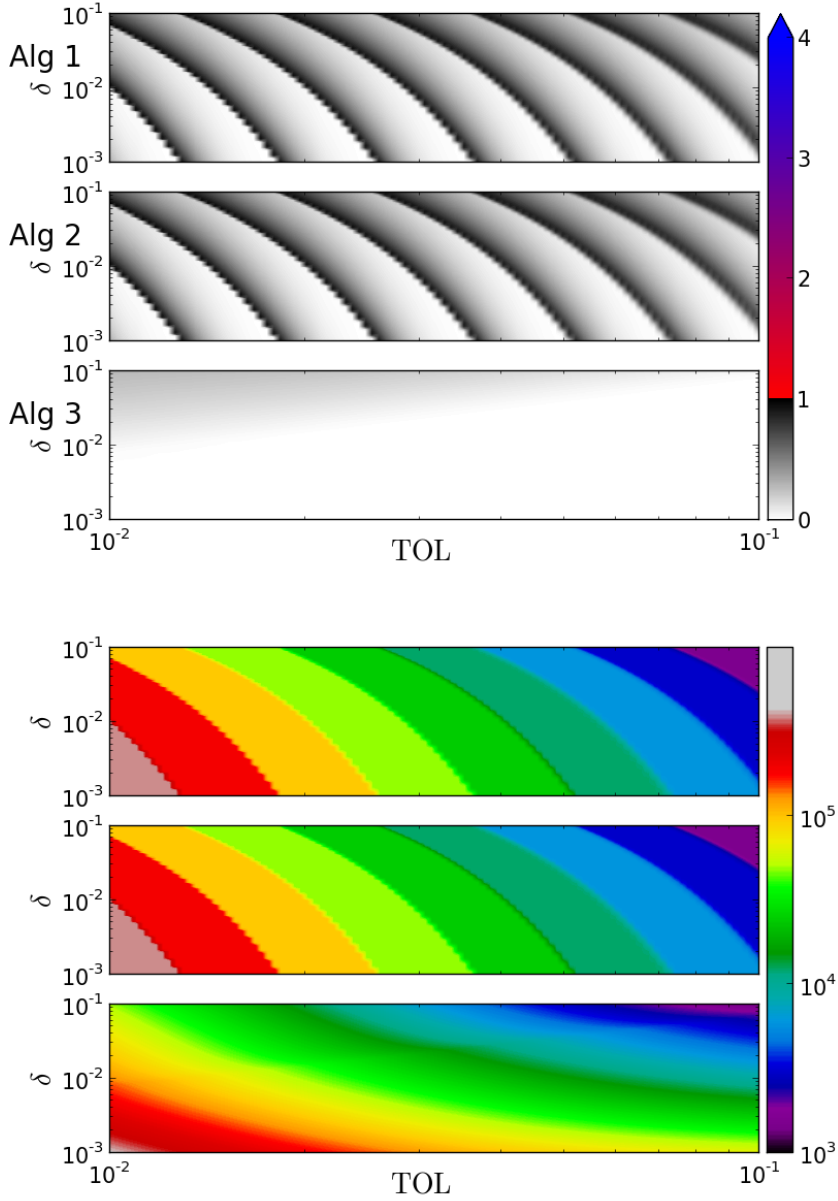


FIG. 5.3. (*Uniform Distribution*) We sample a $U(-\sqrt{3}, \sqrt{3})$ uniformly distributed r.v. Top three plots: Plots of $\bar{p}_N(TOL, \delta)/\delta$, cf. (5.1), where samples of \bar{X}_M respectively are computed by Algorithm 1 in the top plot, Algorithm 2 in the middle plot, and Hickernell et al.'s algorithm in the bottom plot. We consider the algorithm unreliable at (TOL, δ) points where $\bar{p}_N/\delta > 1$. The bottom three plots are corresponding plots of the respective algorithms' complexity in terms of the average number of r.v. samples required to generate \bar{X}_M at the given (TOL, δ) point.

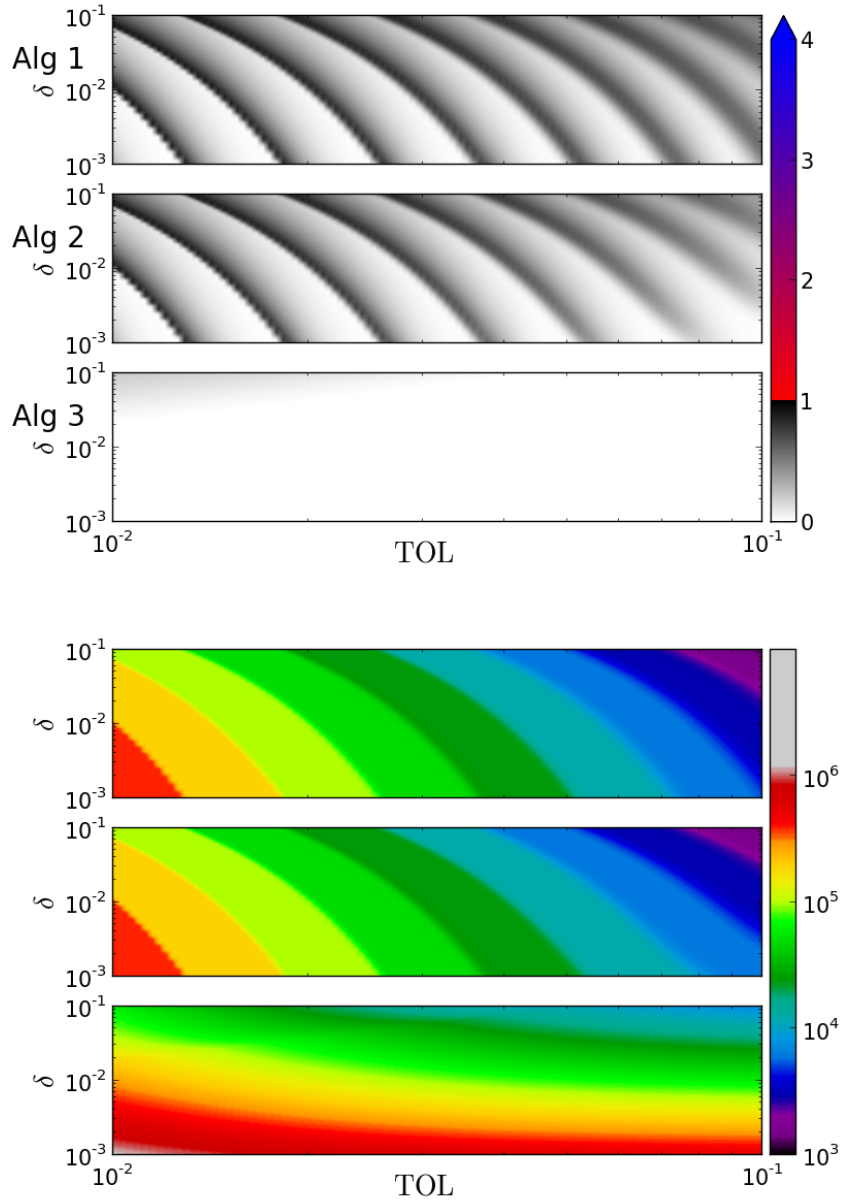


FIG. 5.4. (*Exponential Distribution*) We sample an exponentially distributed r.v. with $\mu = 1$. Top three plots: Plots of $\bar{p}_N(TOL, \delta)/\delta$, cf. (5.1), where samples of \bar{X}_M respectively are computed by Algorithm 1 in the top plot, Algorithm 2 in the middle plot, and Hickernell et al.'s algorithm in the bottom plot. We consider the algorithm unreliable at (TOL, δ) points where $\bar{p}_N/\delta > 1$. The bottom three plots are corresponding plots of the respective algorithms' complexity in terms of the average number of r.v. samples required to generate \bar{X}_M at the given (TOL, δ) point.

6. Conclusion. We have shown that second moment based sequential stopping rules such as Algorithm 1 run the risk of using too few samples in MC estimates, especially when sampling heavy-tailed r.v. in settings with very stringent confidence

requirements, i.e., $\delta \ll TOL$. Algorithm 2, a higher moment based stopping rule algorithm is proposed in this work, and, according to the numerical examples of Section 5, our new stopping rule performs much more reliable than Algorithm 1 while only moderately increasing the computational cost. In short, we believe that our new stopping rule presented in Algorithm 2 is well worth considering in settings with heavy tailed r.v. and/or $\delta \ll TOL$.

Note that our analysis of the original Algorithm 1 critically depends on three main ingredients:

- (I) a general, non-asymptotic estimate of the tail probabilities for the sample mean \bar{X}_M , for which we used either the non-uniform Berry-Esseen theorem given in Corollary 1.2 or the Edgeworth expansion given in Theorem 1.3,
- (II) a choice between the more conservative Berry-Esseen bound and the approximate Edgeworth bound made depending on whether the sample variance of the samples used to generate the output MC estimate is close to, or far from the true variance,
- (III) an estimate of the conditional distribution function of the sample variance given the output M of the stopping algorithm given in (3.4).

There is clearly room for improvement in all these steps. First of all, the second ingredient above is dangerous as we do not know how to estimate the correlation between \bar{X}_M and the events $|\bar{\sigma}_M^2 - \sigma^2| > \sigma^2/2$ and $|\bar{\sigma}_M^2 - \sigma^2| \leq \sigma^2/2$. This is problematic, as these approximations can potentially have the wrong sign, i.e., it is possible that the right-hand sides of (3.3) and (3.5) are smaller than their respective left-hand sides even though we actually seek upper bounds. It is however our hope that these approximation errors are compensated by the the overly pessimistic non-uniform Berry-Esseen estimate and by using Chebychev's inequality to bound the conditional distribution function of the sample variance. Even though the numerical evidence obtained in Section 5 seems to confirm that the compensations work well, we would prefer an analysis in which each estimation step can be controlled, at least in the sense that we indeed obtain an upper bound for the error probability.

To a lesser extent, it is not clear that the truncation of the $o(n^{-1/2})$ of the Edgeworth expansion will lead to an upper bound for the error probability, either. In this case, the approximation error is however of higher order, so a stronger case can be made on why the effect will finally be negligible. In fact, when we used the truncated Edgeworth expansion also for the estimation of (3.3) – instead of the non-uniform Berry-Esseen theorem – then the corresponding stopping rule turned out to be not much more reliable than Algorithm 1, indicating that there is a delicate balance between reliability in meeting the accuracy target (1.1) and maintaining an acceptable efficiency.

Acknowledgments. The authors would like to thank Prof. Thomas Müller-Gronbach for pointing out the existence of the preprint [6].

REFERENCES

- [1] R. R. Bahadur and Leonard J. Savage. The nonexistence of certain statistical procedures in nonparametric problems. *Ann. Math. Statist.*, 27:1115–1122, 1956.
- [2] Paul Bratley, Bennett L. Fox, and Linus E. Schrage. *A guide to simulation (2nd ed.)*. Springer-Verlag New York, Inc., New York, NY, USA, 1987.
- [3] Y. S. Chow and Herbert Robbins. On the asymptotic theory of fixed-width sequential confidence intervals for the mean. *The Annals of Mathematical Statistics*, 36(2):pp. 457–462, 1965.

- [4] W. Feller. *An introduction to probability theory and its applications. Vol II. 2nd ed.* Wiley Series in Probability and Mathematical Statistics. New York etc.: John Wiley and Sons, Inc. XXIV, 669 p., 1971.
- [5] Peter W. Glynn and Ward Whitt. The asymptotic validity of sequential stopping rules for stochastic simulations. *The Annals of Applied Probability*, 2(1):pp. 180–198, 1992.
- [6] F. J. Hickernell, L. Jiang, Y. Liu, and A. Owen. Guaranteed conservative fixed width confidence intervals via monte carlo sampling. *ArXiv e-prints*, aug 2012.
- [7] Wassily Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):pp. 13–30, 1963.
- [8] J.D. Hunter. Matplotlib: A 2d graphics environment. *Computing in Science Engineering*, 9(3):90–95, may-june 2007.
- [9] E. S. Keeping. *Introduction to statistical inference*. D. Van Nostrand Co., Inc., Princeton, N.J.-Toronto-London-New York, 1962.
- [10] Averill M. Law and W. David Kelton. Confidence intervals for steady-state simulations, ii: A survey of sequential procedures. *Management Science*, 28(5):pp. 550–562, 1982.
- [11] Averill M. Law and W. David Kelton. Confidence intervals for steady-state simulations: I. a survey of fixed sample size procedures. *Operations Research*, 32(6):pp. 1221–1239, 1984.
- [12] Pierre L’Ecuyer and Eric Buist. Simulation in java with ssj. In *Proceedings of the 37th conference on Winter simulation*, WSC ’05, pages 611–620. Winter Simulation Conference, 2005.
- [13] Valentin V. Petrov. *Limit theorems of probability theory*, volume 4 of *Oxford Studies in Probability*. The Clarendon Press Oxford University Press, New York, 1995. Sequences of independent random variables, Oxford Science Publications.
- [14] I. Shevtsova. On the absolute constants in the Berry-Esseen type inequalities for identically distributed summands. *ArXiv e-prints*, November 2011.