
Automatic Monte Carlo Algorithms Where the Integrand Size Is Unknown *

Fred J. Hickernell¹, Lan Jiang¹, Yuewei Liu², and Art Owen³

¹ Department of Applied Mathematics, Illinois Institute of Technology, Chicago, IL, USA,
hickernell@iit.edu, ljiang14@hawk.iit.edu

² School of Mathematics and Statistics, Lanzhou University, Lanzhou City, Gansu, China
730000, ???

³ Art's address here with email here

Summary. We attempt a probabilistic analysis of simple Monte Carlo, achieving probabilistic error bounds when the kurtosis is controlled. The algorithm uses a sample size that depends adaptively on the estimated variance of the integrand. Thus, the algorithm is nonlinear (depending essentially on the function). The advantage of what is done here over standard error analysis (complexity theory) is that the algorithm does not depend a priori on the scale of the problem (in this case the variance) to determine the number of samples. Our intention, if what is done here is correct, is to try to extend this to the more sophisticated sampling schemes and infinite dimensional problems.

1 Introduction

Monte Carlo algorithms for multidimensional integrals have been in use for half a century. The integral or mean, μ , of a d -variate function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, i.e.,

$$\mu = \mu(f) = \int_{\mathbb{R}^d} f(\mathbf{x}) \rho(\mathbf{x}) d\mathbf{x},$$

is approximated by the sample mean of integrand values,

$$\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n Y_i, \quad Y_i = f(\mathbf{X}_i). \quad (1)$$

Here $\rho : \mathbb{R}^d \rightarrow [0, \infty)$ is a specified probability density function, e.g., the uniform density on a cube or the normal density. For simple Monte Carlo the points $\mathbf{X}_1, \mathbf{X}_2, \dots$ independent and identically distributed (i.i.d.) random variables with the marginal probability density function ρ . Other choices, such as quasi-random points, may provide greater accuracy, but this paper deals with only simple Monte Carlo.

It is natural to that practitioners would want an automatic Monte Carlo algorithm that takes as its input the integrand, f , the probability density, ρ , and an error tolerance ε , and produces

* The first and second authors were partially supported by the National Science Foundation under DMS-0923111 and DMS-1115392

a value $\hat{\mu}_n$ that is within ε of μ while requiring only a reasonable sample size, n . While there is an extensive theory on Monte Carlo algorithms, there are some key gaps that must be filled to make it an automatic algorithm. The aim of this article is to provide such an algorithm. Next this problem of approximating μ is summarized from the perspective of a statistician and then from the perspective of an information-based complexity theorist or numerical analyst.

Art will add and edit this material in this section, especially the next subsection.

1.1 The Statistician's Perspective

The statistician's sees this a problem of constructing non-parametric, fixed-width confidence intervals for μ . Setting $Y = f(\mathbf{X})$, the goal is to approximate its mean, $\mu = E(Y)$, from an i.i.d. sample, $Y_i = f(\mathbf{X}_i)$, $i = 1, \dots, n$, to obtain an expression

$$\text{Prob}[|\hat{\mu}_n - \mu| \leq \varepsilon] \geq 1 - \alpha. \quad (2)$$

Here taking $\alpha = 5\%$ corresponds to a 95% confidence interval for μ with half-width ε . Although the probability distribution for Y is not readily obtainable, the Central Limit Theorem states that the sample mean, $\hat{\mu}_n$, is approximately normally distributed for large n . This fact may be used to determine the sample size n to give an approximate confidence interval,

$$\text{Prob}\left[|\hat{\mu}_n - \mu| \leq \frac{z_{\alpha/2}\sigma}{\sqrt{n}}\right] \approx 1 - \alpha, \quad (3)$$

where $\sigma^2 = \text{var}(Y) = E[(Y - \mu)^2]$, and z_α denotes the $(1 - \alpha)100\%$ percentile of the standard Gaussian distribution. This suggests taking a sample size of $n = \lceil (z_{\alpha/2}\sigma/\varepsilon)^2 \rceil$ to attain (2). Since σ is not known a priori, it must be approximated by the sample variance of the Y_i :

$$\hat{\sigma}^2 = \mathfrak{C}^2 \hat{v}_{n_\sigma}, \quad \text{where} \quad \hat{v}_{n_\sigma} = \frac{1}{n_\sigma - 1} \sum_{i=1}^{n_\sigma} (Y_i - \hat{\mu}_{n_\sigma})^2, \quad \hat{\mu}_{n_\sigma} = \frac{1}{n_\sigma} \sum_{i=1}^{n_\sigma} Y_i. \quad (4)$$

The variance inflation factor, $\mathfrak{C} > 1$, accounts for the fact the unbiased variance estimator, \hat{v}_{n_σ} , may be larger or smaller than the true variance, σ^2 . Then the sample size for the sample mean may be chosen as $n = \lceil (z_{\alpha/2}\hat{\sigma}/\varepsilon)^2 \rceil$. To avoid dependence between $\hat{\sigma}^2$ and $\hat{\mu}_n$, an independent sample of Y_i , should be used to compute $\hat{\mu}_n$, i.e.,

$$\hat{\mu}_n = \frac{1}{n} \sum_{i=n_\sigma+1}^{n_\sigma+n} Y_i, \quad (5)$$

for a total cost of $n_\sigma + n$ samples.

While this fixed-width confidence interval may be often correct, it is not guaranteed. There are two approximations, whose errors must be understood and controlled: approximating the distribution $\hat{\mu}_n$ as a Gaussian distribution, and approximating σ^2 by the sample variance. The key to controlling these two errors, as seen in Theorem 1, is in assuming a bound on the kurtosis, i.e.,

$$\kappa = \text{kurt}(Y) = \frac{E[(Y - \mu)^4]}{\sigma^4} \leq \kappa_{\max}. \quad (6)$$

A bounded kurtosis means that Y has or does not have ????? A bounded kurtosis allows one to reliably bound $\text{var}(Y)$ in terms of the sample variance, and also allows one to obtain a rigorous bound on the deviation of the sample mean from the true mean. It was shown by (Bahadur and Savage, 1956, Corollary 2) that nonparametric confidence intervals are impossible for

convex sets of distributions. Assuming bounded kurtosis produces a non-convex set of possible distributions for Y .

Another approach to obtaining a fixed-width confidence interval for μ would be to assume an upper bound on σ^2 and apply Chebyshev's inequality. The disadvantage of this approach is that while it might work for $Y = f(\mathbf{X})$, it would not work for $cY = cf(\mathbf{X})$ if c were large enough, since $\text{var}(cY) = c^2 \text{var}(Y)$. However, since $\text{kurt}(cY) = \text{kurt}(Y)$ if the procedure for obtaining a fixed-width confidence interval that is described in this article works for Y , it also works for cY .

1.2 The Information-Based Complexity Theorist's or Numerical Analyst's Perspective

The information-based complexity theorist or numerical analyst sees this as a multivariate integration problem. There are a myriad of cubature methods, each making certain assumptions on the domain of integration and the smoothness of the integrand, f , and deriving error bounds in terms of some semi-norm of f . Simple Monte Carlo methods work for functions with low degrees of smoothness and the natural space of integrands is \mathcal{L}_2 , where the \mathcal{L}_p norm is defined as

$$\|f\|_p := \left\{ \int_{\mathbb{R}^d} |f(\mathbf{x})|^p \rho(\mathbf{x}) d\mathbf{x} \right\}^{1/p}. \quad (7)$$

The root mean square error of the simple Monte Carlo method is $\|f - \mu(f)\|_2 / \sqrt{n}$.

Again, the practical problem again in applying this error analysis to determine the sample size, n , to guarantee the desired error, one must know the size of the integrand, $\|f - \mu(f)\|_2 = \sqrt{\text{var}(Y)} = \sigma$. Analogous to the statistical argument given above, the solution to this problem is not to look at balls of integrands, i.e.,

$$\mathcal{B}_{\sigma_{\max}} = \{f \in \mathcal{L}_2 : \|f - \mu(f)\|_2 \leq \sigma_{\max}\}, \quad (8)$$

but to look at *cones* of integrands

$$\mathcal{C}_{\kappa_{\max}} = \{f \in \mathcal{L}_4 : \|f - \mu(f)\|_4 \leq \kappa_{\max}^{1/4} \|f - \mu(f)\|_2\}. \quad (9)$$

This condition is the same as (6). Whereas f lying in the ball $\mathcal{B}_{\sigma_{\max}}$ does not guarantee cf lies in that same ball, f lying in the cone $\mathcal{C}_{\kappa_{\max}}$ does guarantee that cf lies in that same cone.

Looking for algorithms that work well for cones of integrands, $\mathcal{C}_{\kappa_{\max}}$, leads one to *adaptive* algorithms. The sample size used to estimate the integral is determined adaptively by first computing an upper bound on $\|f - \mu(f)\|_2$. In information-based complexity theory it is known that adaptive information does not help for convex sets of integrands in the worst case and probabilistic settings (Traub et al, 1988, Chapter 4, Theorem 5.2.1; Chapter 8, Corollary 5.3.1). Here, the cone, $\mathcal{C}_{\kappa_{\max}}$ is not a convex set, so adaption can help.

Again, it should be stressed that the algorithm to be presented here is automatic. It does not require information about $\|f - \mu(f)\|_2 = \sigma$, but this quantity needs to be reliably estimated by the algorithm. Thus, the sample size needed, and consequently the time required, to estimate μ to within the prescribed error tolerance depends on how large $\|f - \mu(f)\|_2 = \sigma$ is estimated to be. The algorithm is adaptive, and its cost depends on the integrand.

1.3 Illustrative Univariate Examples of Automatic Algorithms

Lan will be adding to this section.

Several commonly used software packages have automatic algorithms for integrating functions of a single variable. These include

- `quad` in MATLAB (The MathWorks, Inc., 2012), adaptive Simpson's rule based on `adaptsim` by Gander and Gautschi (2000),
- `quadgk` in MATLAB (The MathWorks, Inc., 2012), adaptive Gauss-Kronrod quadrature based on `quadva` by Shampine (2008),
- the `chebfun` (Hale et al, 2012) toolbox for MATLAB (The MathWorks, Inc., 2012), which approximates integrals by integrating interpolatory Chebyshev polynomial series for the integrands,
- `NIntegrate` in Mathematica (Wolfram Research Inc., 2011), which uses a number of adaptive rules for one dimensional integrals,

For the first three of these automatic algorithms one can easily probe where they sample the integrand, feed the algorithms zero values, and then construct fooling functions that the automatic algorithms will return a zero value for the integral. Figure 1 displays these fooling functions for the problem $\mu = \int_0^1 f(x) dx$ for the first three algorithms. Each of these algorithms is asked to provide an answer with an absolute error no greater than 10^{14} , but in fact the absolute error is 1 for these fooling functions. The algorithms `quad` and `chebfun` sample only about a dozen points before concluding that the function is zero, whereas the algorithm `quadgk` samples a much larger number of points (only those between 0 and 0.01 are shown in the plot). Algorithm `NIntegrate` is hard snoop, but it is examined for the next example.

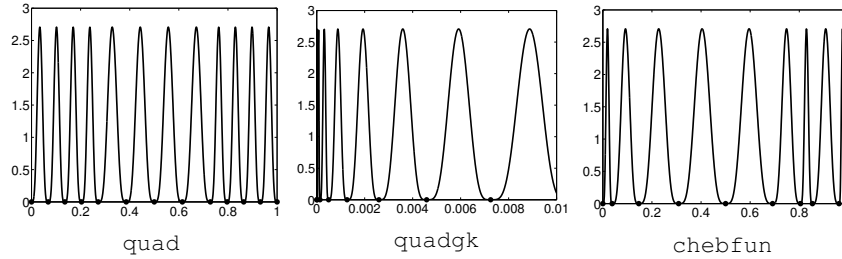


Fig. 1. Plots of fooling functions, f , with $\mu = \int_0^1 f(x) dx = 1$, but for which the corresponding algorithms return values of $\hat{\mu} = 0$.

Accuracy and timing results have been recorded for the test function

$$f(\mathbf{x}) = a_0 + b_0 \prod_{j=1}^d \left[1 + b_j \exp \left(-\frac{(x_j - h_j)^2}{c_j^2} \right) \right]. \quad (10)$$

Here \mathbf{x} is a d dimensional vector, and the b_j , c_j , and h_j are parameters. For Figure 2 shows the results of different algorithms being used to integrate 500 different instances of f being performed. For each replication, the parameters determining f are chosen as follows:

- $b_j \in [0.1, 10]$ with $\log(b_j)$ being i.i.d. uniform, $j = 1, \dots, d$,

- $c_j \in [10^{-6}, 1]$ with $\log(c_j)$ being i.i.d. uniform, $j = 1, \dots, d$,
- $h_j \in [0, 1]$ with h_j being i.i.d. uniform, $j = 1, \dots, d$,
- b_0 chosen in terms of the b_j , c_j , and h_j in order to make $\sigma^2(f) \in [10^{-2}, 10^2]$, with $\log(\sigma(f))$ being i.i.d. uniform for each instance, and
- a_0 in terms of chosen in terms of the b_j , c_j , and h_j to make $\mu(f) = 1$.
- dimension $d = 1$
- the probability density function $\rho = 1$ which is uniformly distributed between $[0, 1]$
- $n_0 = 1024$ for the simple algorithm and $n_0 = 1024 \cdot 128$ for the heavy duty algorithm.
- For the iid sampling, kurtosis max is 9.2085; for the iid heavy duty sampling, it is 1051.9.

For the iid sampling, the successful rate is only 70% whereas the heavy duty iid sampling, the successful rate is above 95%. in the graph, the diamond means it reaches the sample limit; the star means the kurtosis is less than the kurtosis maximum. the dot is the general case. For the Sobol sampling the successful rate is around 70% , the heavy duty case, the successful rate is almost 100%. Clearly, Sobol takes less time than iid sampling. For the quad, quadgk, chebfun and chebfun heavy duty, they take less time but could not reach a higher tolerance as we expected.

2 Simple Monte Carlo with Guaranteed Error Estimation

2.1 Reliably Bounding the Variance

To estimate or bound the variance of $Y = f(\mathbf{X})$ one must have a finite fourth moment. The definitions of the variance and kurtosis of Y in Section 1.1 may be extended naturally to $f \in \mathcal{L}_4$ in terms of \mathcal{L}_p norms defined in (7) as follows:

$$\sigma^2 = \text{var}(Y) = \text{var}(f) = \|f - \mu(f)\|_2^2, \quad \kappa = \text{kurt}(Y) = \text{kurt}(f) = \frac{\|f - \mu(f)\|_4^4}{\|f - \mu(f)\|_2^4}.$$

Thus, one may speak about Y or f interchangeably. If $1 \leq q \leq p$, then by Hölder's inequality,

$$\begin{aligned} \|f\|_q &= \left\{ \int_{\mathbb{R}^d} |f(\mathbf{x})|^q \rho(\mathbf{x}) d\mathbf{x} \right\}^{1/q} \\ &\leq \left\{ \int_{\mathbb{R}^d} |f(\mathbf{x})|^p \rho(\mathbf{x}) d\mathbf{x} \right\}^{1/p} \left\{ \int_{\mathbb{R}^d} 1^{p/(p-q)} \rho(\mathbf{x}) d\mathbf{x} \right\}^{(p-q)/(pq)} = \|f\|_p. \end{aligned}$$

Thus, $\kappa = \text{kurt}(f) \geq 1$, provided $\sigma^2 = \text{var}(f) > 0$. For $\sigma^2 = \text{var}(f) = 0$, one defines $\kappa = \text{kurt}(f) = 1$.

As mentioned in (4), a practical upper bound on σ^2 may be obtained in terms of the sample variance. The justification for this statement is contained in the lemma below. Two well-known probability inequalities that are needed here and later are quoted in the following lemma.

Lemma 1. *Suppose that Z is any random variable with finite second moment. For any $a > 0$,*

$$\text{Prob}[|Z - E(Z)| \geq a] \leq \frac{\text{var}(Z)}{a^2} \quad \text{Chebyshev's inequality (Lin and Bai, 2010, 6.1c),}$$

$$\text{Prob}[Z - E[Z] \geq a] \leq \frac{\text{var}(Z)}{a^2 + \text{var}(Z)} \quad \text{Cantelli's Inequality (Lin and Bai, 2010, 6.1e).}$$

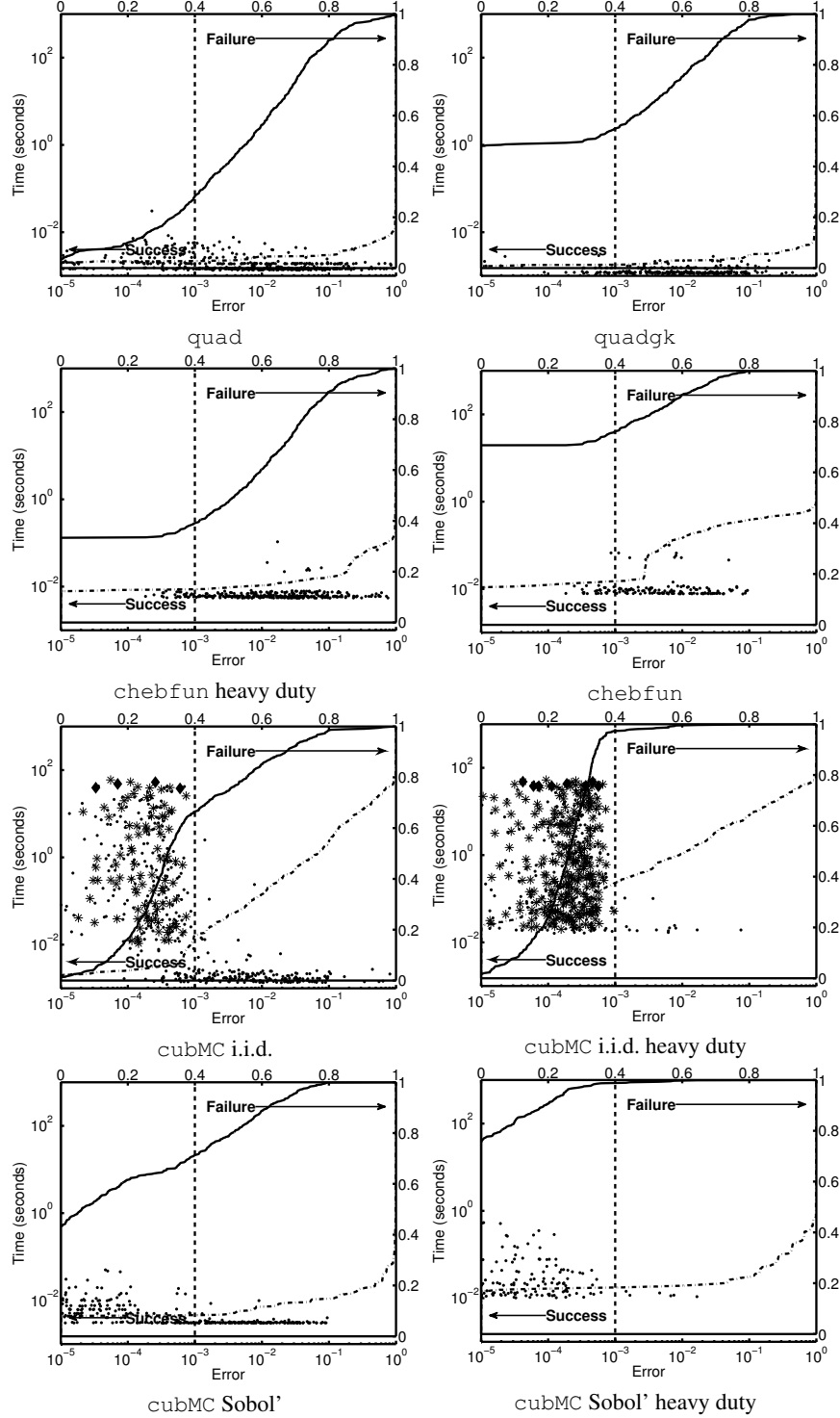


Fig. 2. Execution times and errors for test function (10) for $d = 1$ and $\varepsilon = 10^{-3}$, and a variety of parameters giving a range of $\sigma(f)$ and $\kappa(f)$. The solid line shows that cumulative distribution of actual errors, and the dot-dashed line shows the cumulative distribution of execution times. For the `cubMC i.i.d.` and `i.i.d. heavy duty` the points labeled * are those for which the Theorem 1 guarantees the error tolerance.

Lemma 2. Suppose that Y is a random variable with finite fourth moment, and with mean μ , variance σ^2 , and kurtosis κ . Let \hat{v}_n denote the unbiased sample variance based on n i.i.d. samples, Y_1, \dots, Y_n as defined in (4). Then for any $\alpha \in (0, 1]$ it follows that

$$\text{Prob} \left[\frac{\hat{v}_n}{1 - \sqrt{\left(\kappa - \frac{n-3}{n-1}\right) \left(\frac{1-\alpha}{\alpha n}\right)}} > \sigma^2 \right] \geq 1 - \alpha, \quad (11a)$$

$$\text{Prob} \left[\frac{\hat{v}_n}{1 + \sqrt{\left(\kappa - \frac{n-3}{n-1}\right) \left(\frac{1-\alpha}{\alpha n}\right)}} < \sigma^2 \right] \geq 1 - \alpha. \quad (11b)$$

Proof. It is known that the sample variance is an unbiased estimator of the variance, i.e., $E(\hat{v}_n) = \sigma^2$. It is also known from ? that the variance of the sample variance can be expressed in terms of σ^2 and κ as

$$\text{var}(\hat{v}_n) = \frac{\sigma^4}{n} \left(\kappa - \frac{n-3}{n-1} \right).$$

Choosing

$$a = \sqrt{\text{var}(\hat{v}_n) \frac{1-\alpha}{\alpha}} = \sigma^2 \sqrt{\left(\kappa - \frac{n-3}{n-1}\right) \left(\frac{1-\alpha}{\alpha n}\right)} > 0,$$

it follows from Cantelli's inequality (Lemma 1) that

$$\begin{aligned} \text{Prob} \left[\hat{v}_n - \sigma^2 \geq \sigma^2 \sqrt{\left(\kappa - \frac{n-3}{n-1}\right) \left(\frac{1-\alpha}{\alpha n}\right)} \right] &= \text{Prob} \left[\hat{v}_n - \sigma^2 \geq a \right] \\ &\leq \frac{\text{var}(\hat{v}_n)}{a^2 + \text{var}(\hat{v}_n)} = \frac{\text{var}(\hat{v}_n)}{\text{var}(\hat{v}_n) \frac{1-\alpha}{\alpha} + \text{var}(\hat{v}_n)} = \frac{1}{\left(\frac{1-\alpha}{\alpha}\right) + 1} = \alpha. \end{aligned}$$

Then (11a) follows directly. By a similar argument, applying Cantelli's inequality to the expression $\text{Prob} [-\hat{v}_n + \sigma^2 \geq a]$ implies (11b). \square

Lemma 2 provides probabilistic justification to use $\hat{\sigma}^2 = \mathfrak{C}^2 \hat{v}_{n_\sigma}$ as a reliable upper bound for $\sigma^2 = \text{var}(Y) = \text{var}(f)$. One can claim that $\text{Prob}(\hat{\sigma}^2 \geq \sigma^2) \geq 1 - \alpha$ assuming that

$$\begin{aligned} \frac{1}{1 - \sqrt{\left(\kappa - \frac{n_\sigma-3}{n_\sigma-1}\right) \left(\frac{1-\alpha}{\alpha n_\sigma}\right)}} &\leq \mathfrak{C}^2 \\ \iff \kappa &\leq \frac{n_\sigma-3}{n_\sigma-1} + \left(\frac{\alpha n_\sigma}{1-\alpha}\right) \left(1 - \frac{1}{\mathfrak{C}^2}\right)^2 =: \kappa_{\max}(\alpha, n_\sigma, \mathfrak{C}). \end{aligned} \quad (12)$$

Figure 3a shows how large a kurtosis can be accommodated for a given n_σ , α , and variance inflation factor $\mathfrak{C} = 1.5$. Note that for $n = 30$, a common rule of thumb for applying the central limit theorem, even $\alpha = 0.1$ gives κ_{\max} of only about 2, which is rather restrictive.

2.2 Determining the Sample Size

The other issue that needs to be addressed is a tight probabilistic error bound to replace the asymptotic or approximate error bound given by the Central Limit Theorem, (3). Chebyshev's

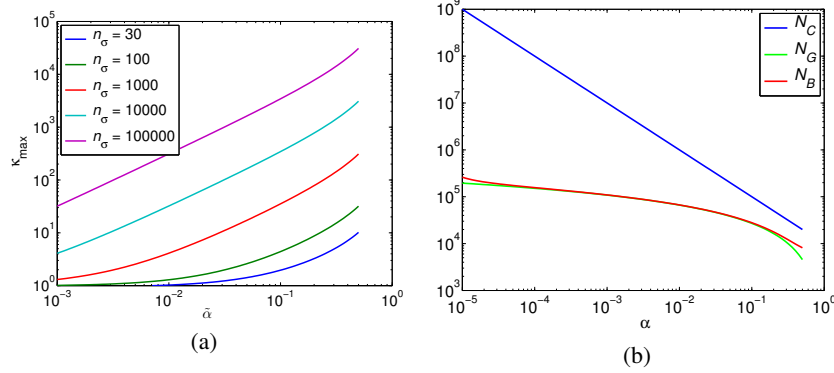


Fig. 3. (a) The maximum kurtosis, $\kappa_{\max}(\alpha, n_\sigma, 1.5)$, as defined in (12); (b) comparison of sample sizes $N_G(0.01, \alpha)$, $N_C(0.01, \alpha)$, and $N_B(0.01, \alpha, \kappa_{\max}^{3/4}(\alpha, 1000, 1.5))$.

inequality implies that the number of function evaluations needed to ensure that $\hat{\mu}_n$ satisfies the error tolerance with high probability is

$$\text{Prob}[|\hat{\mu}_n - \mu| \leq \varepsilon] \geq 1 - \alpha \quad \text{for } n \geq N_C(\varepsilon/\sigma, \alpha), \quad f \in \mathcal{L}_2, \quad (13a)$$

where

$$N_C(\varepsilon, \alpha) := \left\lceil \frac{1}{\alpha \varepsilon^2} \right\rceil. \quad (13b)$$

However, this sample size is much larger than that given by the Central Limit Theorem, as shown in Figure 3b.

Since higher order moments of the integrand are required to guarantee an upper bound on the true variance in terms of the sample variance, it is sensible to use these higher order moments to obtain smaller sample sizes. A smaller sample size than (13b) with a rigorous probabilistic bound can be found by invoking the following inequality.

Lemma 3 (Non-uniform Berry-Esseen Inequality). (Petrov, 1995, Theorem 5.16, p. 168) *Let Y_1, \dots, Y_n be i.i.d. random variables. Suppose that $\mu = E(Y_i)$, $\text{var}(Y_i) = \sigma^2 > 0$, and $M = E|Y_i - \mu|^3 / \sigma^3 < \infty$. Then*

$$\left| \text{Prob} \left[\frac{1}{\sigma \sqrt{n}} \sum_{i=1}^n (Y_i - \mu) < x \right] - \Phi(x) \right| \leq \frac{AM}{\sqrt{n}(1 + |x|)^3}.$$

for all x , where Φ is the cumulative distribution function of a standard Gaussian (normal) random variable, and A is some number satisfying $0.4097 \leq A \leq 0.5600$.

The right hand side of this inequality bounds how far the probability distribution of sample mean deviates from the Central Limit Theorem approximation. This right hand side vanishes as either the sample size or $|x|$ tends to infinity. This inequality requires that $Y = f(\mathbf{X})$, have a finite third moment, i.e., $f \in \mathcal{L}_3$, which is a weaker assumption than is needed to bound the variance in terms of the sample variance in Lemma 2. Recalling that $Y_i = f(\mathbf{X}_i)$, $\mu = E(Y_i)$, and $\hat{\mu}_n = (Y_1 + \dots + Y_n)/n$, this Berry-Esseen inequality then implies that for positive x ,

$$\begin{aligned}
\text{Prob} \left[|\hat{\mu}_n - \mu| < \frac{\sigma x}{\sqrt{n}} \right] &= \text{Prob} \left[\hat{\mu}_n - \mu < \frac{\sigma x}{\sqrt{n}} \right] - \text{Prob} \left[\hat{\mu}_n - \mu < -\frac{\sigma x}{\sqrt{n}} \right] \\
&\geq \left[\Phi(x) - \frac{0.56M}{\sqrt{n}(1+|x|)^3} \right] - \left[\Phi(-x) + \frac{0.56M}{\sqrt{n}(1+|x|)^3} \right] \\
&= 1 - 2 \left(\Phi(-x) - \frac{0.56M}{\sqrt{n}(1+|x|)^3} \right). \tag{14}
\end{aligned}$$

Letting $x = \varepsilon \sqrt{n}/\sigma$, the probability of making an error less than ε is bounded below by $1 - \alpha$, i.e.,

$$\text{Prob}[|\hat{\mu}_n - \mu| < \varepsilon] \geq 1 - \alpha, \quad \text{provided } n \geq N_B(\varepsilon/\sigma, \alpha, M), \quad f \in \mathcal{L}_3, \tag{15a}$$

where

$$N_B(b, \alpha, M) := \min \left\{ n \in \mathbb{N} : \Phi(-b\sqrt{n}) + \frac{0.56M}{\sqrt{n}(1+b\sqrt{n})^3} \leq \frac{\alpha}{2} \right\}. \tag{15b}$$

As shown in Figure 3b, for a range of α , the sample size guaranteeing coverage of the confidence interval, N_B , is quite close to the sample size for the approximate Central Limit Theorem confidence interval, N_G , however, N_B may be somewhat larger for very small or rather large α . In general N_B is significantly smaller than N_G , but not always. A disadvantage of (15) is that class of integrands, \mathcal{L}_3 , is smaller than that in (13), but this typically a small price to pay given the much smaller cost of computation.

The theorem below combines the results on estimating the variance with the sample sizes arising from Chebyshev's inequality and the Berry-Esseen inequality. These lead to an adaptive Monte Carlo algorithm with a probabilistic error guarantee.

Theorem 1. *Specify the following parameters defining the algorithm:*

- sample size for variance estimation, $n_\sigma \in \mathbb{N}$,
- inflation factor for variance estimation, $\mathfrak{C} \in (1, \infty)$,
- uncertainty tolerance, $\alpha \in (0, 1)$, and $\tilde{\alpha} = 1 - \sqrt{1 - \alpha}$, and
- absolute error tolerance, $\varepsilon \in (0, \infty)$.

For $\kappa_{\max} = \kappa_{\max}(n_\sigma, \tilde{\alpha}, \mathfrak{C})$, as defined in (12), define the cone of integrands functions with bounded kurtosis, $\mathcal{C}_{\kappa_{\max}}$, according to (9). For any $f \in \mathcal{C}_{\kappa_{\max}}$, compute the sample variance, \hat{v}_{n_σ} using a simple random sample of size n_σ . Use this to approximate the variance of f by $\hat{\sigma}^2 = \mathfrak{C}^2 \hat{v}_{n_\sigma}$ as in (4). Next define a sample size $n = N_{CB}(\varepsilon/\hat{\sigma}, \tilde{\alpha}, \kappa_{\max}^{3/4})$, where

$$N_{CB}(b, \alpha, M) := \min(N_C(b, \alpha), N_B(b, \alpha, M)), \tag{16}$$

N_C is defined in (13b) and N_B is defined in (15b). Compute $\hat{\mu}_n$, the simple Monte Carlo estimator of μ based on n samples, as in (5). A probabilistic absolute error bound is given by

$$\text{Prob}[|\hat{\mu}_n - \mu| \leq \varepsilon] \geq 1 - \alpha.$$

Proof. Note that (1) implies that the third moment can be bounded in terms of the fourth moment, namely M in Lemma 3 is no greater than $\kappa^{3/4}$. There are three primary random variables in this algorithm: the estimated upper bound on the standard deviation, $\hat{\sigma}$, the sample size to estimate the mean, n , which is an explicit function of $\hat{\sigma}$, and the estimated mean, $\hat{\mu}_n$. By (13) and (15) it then follows that $\text{Prob}[|\hat{\mu}_n - \mu| \leq \varepsilon] \geq 1 - \tilde{\alpha}$, provided that $\hat{\sigma} \geq \sigma$. Thus,

$$\begin{aligned}
\text{Prob}[|\hat{\mu}_n - \mu| \leq \varepsilon] &= E_{\hat{\sigma}} \{ \text{Prob}[|\hat{\mu}_n - \mu| \leq \varepsilon \mid \hat{\sigma}] \} \\
&\geq E_{\hat{\sigma}} \left\{ (1 - \tilde{\alpha}) 1_{[\sigma, \infty)}(\hat{\sigma}) \right\} \\
&\geq (1 - \tilde{\alpha})(1 - \tilde{\alpha}) = 1 - \alpha,
\end{aligned}$$

since $\hat{\sigma} \geq \sigma$ with probability $1 - \tilde{\alpha}$ by (12). \square

Remark 1. If one is willing to invest n_σ samples to estimate σ , it makes practical sense to choose the sample size for the sample mean at least that large, i.e.,

$$n = \max(n_\sigma, N_{CB}(\varepsilon/\hat{\sigma}, \tilde{\alpha}, \kappa_{\max}^{3/4})).$$

By the error bound following from Chebyshev's inequality, (13), this means that the probabilistic absolute error bound in Theorem 1 also holds for integrands, f , lying in the ball $\mathcal{B}_{\sigma_{\max}}$, defined in (8), where $\sigma_{\max} = \varepsilon\sqrt{\alpha n_\sigma}$.

2.3 Cost of the Algorithm

The sample size of the adaptive algorithm defined in Theorem 1 is a random variable, and so the cost of this algorithm might best be defined probabilistically. Moreover, the cost depends strongly on σ as well as the ε , and its definition should reflect this dependence.

Let A be any random algorithm defined for a set of integrands \mathcal{F} that takes as its input an error tolerance, ε , an uncertainty level, α , and a procedure for computing values of $f \in \mathcal{F}$. The algorithm then computes an approximation to the integral, $A(f, \varepsilon, \alpha)$. This approximation is based solely on $f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)$, where the choice of each \mathbf{x}_i may depend iteratively on $(\mathbf{x}_1, f(\mathbf{x}_1)), \dots, (\mathbf{x}_{i-1}, f(\mathbf{x}_{i-1}))$, and the decision to stop with n data depends on all function data collected, up to and including the n^{th} . The cost of this algorithm for this set of inputs is denoted $\text{cost}(A, \varepsilon, \alpha, f)$, is the number of function values, n , which as noted is a random variable. The probabilistic cost of the algorithm, with uncertainty β , for integrands of variance no greater than σ_{\max}^2 is defined as

$$\text{cost}(A, \varepsilon, \alpha, \beta, \mathcal{F}, \sigma_{\max}) := \sup_{\substack{f \in \mathcal{F} \\ \text{var}(f) \leq \sigma^2}} \min \{N : \text{Prob}[\text{cost}(A, \varepsilon, \alpha, f) \leq N] \geq 1 - \beta\}. \quad (17)$$

The cost of the particular adaptive Monte Carlo algorithm defined in Theorem 1, denoted aMC, for the cone of integrands $\mathcal{C}_{\kappa_{\max}}$ is

$$\text{cost}(\text{aMC}, \varepsilon, \alpha, \beta, \mathcal{C}_{\kappa_{\max}}, \sigma_{\max}) := \sup_{\substack{f \in \mathcal{C}_{\kappa_{\max}} \\ \text{var}(f) \leq \sigma_{\max}^2}} \min \{N : \text{Prob}(n_\sigma + n \leq N) \geq 1 - \beta\}. \quad (18)$$

Since n_σ is fixed, bounding this cost depends on bounding n , which depends on $\hat{\sigma}$ as given by Theorem 1. Moreover, $\hat{\sigma}$ can be bounded above using (11b) in Lemma 2. For $f \in \mathcal{C}_{\kappa_{\max}}$,

$$\begin{aligned} 1 - \beta &\leq \text{Prob} \left[\hat{v}_{n_\sigma} < \sigma^2 \left\{ 1 + \sqrt{\left(\kappa - \frac{n_\sigma - 3}{n_\sigma - 1} \right) \left(\frac{1 - \beta}{\beta n_\sigma} \right)} \right\} \right] \\ &\leq \text{Prob} \left[\hat{\sigma}^2 = \mathfrak{C}^2 \hat{v}_{n_\sigma} < \mathfrak{C}^2 \sigma^2 \left\{ 1 + \sqrt{\left(\kappa_{\max}(n_\sigma, \tilde{\alpha}, \mathfrak{C}) - \frac{n_\sigma - 3}{n_\sigma - 1} \right) \left(\frac{1 - \beta}{\beta n_\sigma} \right)} \right\} \right] \\ &= \text{Prob} \left[\hat{\sigma}^2 < \sigma^2 \gamma^2(\tilde{\alpha}, \beta, \mathfrak{C}) \right], \end{aligned}$$

where

$$\gamma^2(\tilde{\alpha}, \beta, \mathfrak{C}) := \mathfrak{C}^2 \left\{ 1 + \sqrt{\left(\frac{\tilde{\alpha}}{1 - \tilde{\alpha}} \right) \left(\frac{1 - \beta}{\beta} \right) \left(1 - \frac{1}{\mathfrak{C}^2} \right)^2} \right\} > 1.$$

Noting that $N_{CB}(\cdot, \alpha, M)$ is a non-increasing function allows one to derive the following upper bound on the cost of the adaptive Monte Carlo algorithm.

Theorem 2. *The adaptive Monte Carlo algorithm described in Theorem 1, denoted aMC, has a probabilistic cost bounded above by*

$$\text{cost}(\text{aMC}, \varepsilon, \alpha, \beta, \mathcal{C}_{\kappa_{\max}}, \sigma_{\max}) \leq n_{\sigma} + N_{CB}(\varepsilon / (\sigma_{\max} \gamma(\tilde{\alpha}, \beta, \mathfrak{C})), \tilde{\alpha}, \kappa_{\max}^{3/4}).$$

The cost of the adaptive Monte Carlo algorithm aMC is roughly proportional to $\sigma_{\max}^2 \varepsilon^{-2}$. The set $\mathcal{C}_{\kappa_{\max}}$ contains integrands with arbitrarily large variance, $\sigma^2 = \text{var}(f)$, and thus with potentially arbitrarily large algorithmic cost. On the other hand, since the algorithm is adaptive, the cost may be small if σ^2 is small. The upper bound in Theorem 2 certainly scales with the σ_{\max}^2 as one might hope if σ_{\max}^2 were known. The variable cost of the algorithm for integrands in $\mathcal{C}_{\kappa_{\max}}$ is actually an advantage, rather than a drawback, of this analysis. One need not make any a priori assumptions about the size of the integrand, σ , only about its kurtosis, κ , which is unchanged when the integrand is multiplied by an arbitrary nonzero constant.

3 More Numerical Examples

We are using the function 10 and the parameters on page 5 to test the higher dimension case, set the dimension d be random between 2 to 8, then we have the graphs below:

4 A General Error Criterion

In many practical situations, one needs to approximate the integral with a certain relative accuracy. For example, one wants an answer that is correct to three significant digits. In this case, given a tolerance, ε , and a significance level α , with $\varepsilon, \alpha \in (0, 1)$, one seeks a random $\tilde{\mu}$ such that

$$\text{Prob} \left[\left| \frac{\tilde{\mu} - \mu}{\mu} \right| \leq \varepsilon \right] \geq 1 - \alpha.$$

A more general form of this criterion would be

$$\text{Prob} \left[\frac{|\tilde{\mu} - \mu|}{1 - \theta + \theta |\mu|} \leq \varepsilon \right] \geq 1 - \alpha. \quad (19)$$

for some fixed $\theta \in [0, 1]$, where $\theta = 0$ corresponds to absolute error, and $\theta = 1$ corresponds to relative error. Clearly, one must have $(1 - \theta) + |\mu| \neq 0$ for such a statement to be possible.

If $\varepsilon_A \geq 0$ is an absolute error tolerance, and $\varepsilon_R \geq 0$ is a relative error tolerance, then letting

$$\varepsilon = \frac{\varepsilon_A \varepsilon_R}{\theta \varepsilon_A + (1 - \theta) \varepsilon_R},$$

it follows that for all $\theta \in [0, 1]$,

$$[1 - \theta + \theta |\mu|] \varepsilon = (1 - \gamma) \varepsilon_A + \gamma \varepsilon_R |\mu| \leq \max(\varepsilon_A, \varepsilon_R |\mu|),$$

where

$$\gamma = \frac{\theta \varepsilon_A}{\theta \varepsilon_A + (1 - \theta) \varepsilon_R} \in [0, 1].$$

Thus, error criterion (19) implies that one has satisfied either an absolute or a relative error criterion,

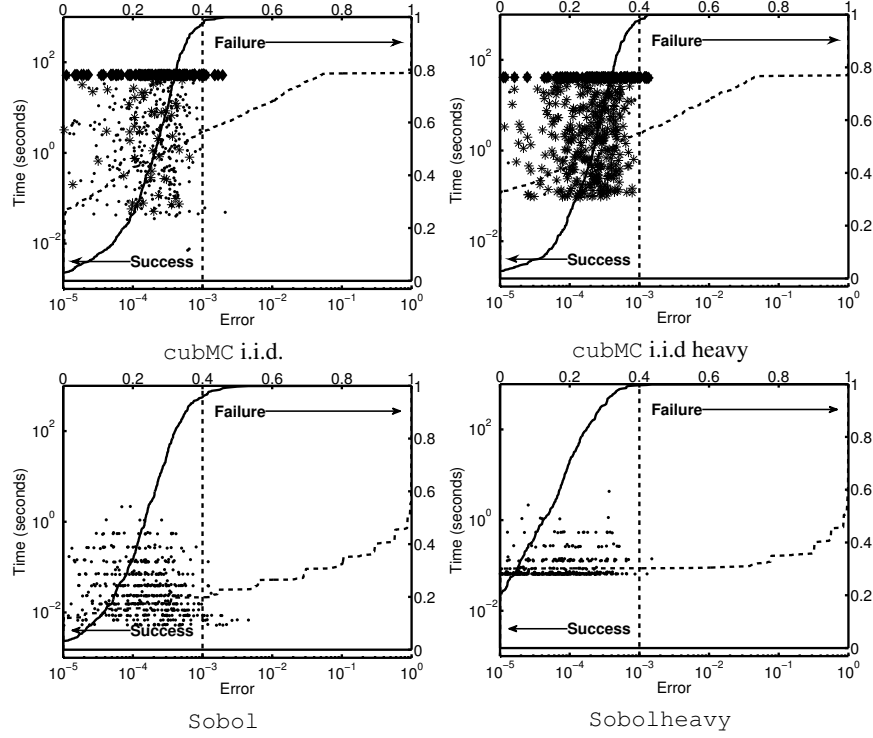


Fig. 4. Execution times and errors for test function (10) for $d = 7$ and $\varepsilon = 10^{-3}$, and a variety of parameters giving a range of $\sigma(f)$ and $\kappa(f)$. The solid line shows that cumulative distribution of actual errors, and the dot-dashed line shows the cumulative distribution of execution times. For the cubMC i.i.d. and i.i.d. heavy duty the points labeled * are those for which the Theorem 1 guarantees the error tolerance.

$$\text{Prob} \left[\frac{|\tilde{\mu} - \mu|}{1 - \theta + \theta |\mu|} \leq \varepsilon \right] \geq 1 - \alpha \implies \text{Prob} [|\tilde{\mu} - \mu| \leq \varepsilon_A \text{ or } |\tilde{\mu} - \mu| \leq |\mu| \varepsilon_R] \geq 1 - \alpha.$$

A value of γ close to zero implies a preference to fulfill the absolute error criterion while a value of γ close to one implies a preference to fulfill the relative error criterion.

Obtaining a confidence interval of the form (19), proceeds in three stages: i) obtaining an upper bound on σ^2 , ii) obtaining a lower bound on $1 - \theta + \theta |\mu|$, and iii) then using these to obtain (19). What differs from the absolute error case is step ii). For this step it is noted that

$$\begin{aligned} \text{Prob} [|\hat{\mu} - \mu| \leq \hat{\varepsilon}] &\geq 1 - \alpha \implies \text{Prob} [\max(|\hat{\mu}| - \hat{\varepsilon}, 0) \leq |\mu| \leq |\hat{\mu}| + \hat{\varepsilon}] \geq 1 - \alpha \\ \implies \text{Prob} [1 - \theta + \theta \max(|\hat{\mu}| - \hat{\varepsilon}, 0) &\leq 1 - \theta + \theta |\mu| \leq 1 - \theta + \theta (|\hat{\mu}| + \hat{\varepsilon})] \geq 1 - \alpha. \end{aligned} \quad (20)$$

Although one might be happy the left side of this inequality being positive, if it is too much smaller than the right side, then one might be eventually expending too much extra work in step iii). Thus, it makes sense to require

$$1 - \theta + \theta \max(|\hat{\mu}| - \hat{\varepsilon}, 0) \geq \hat{\delta}[1 - \theta + \theta(|\hat{\mu}| + \hat{\varepsilon})]$$

$$\iff \hat{\varepsilon} = \begin{cases} \frac{(1 - \hat{\delta})(1 - \theta)}{\hat{\delta}\theta} - |\hat{\mu}|, & 0 \leq |\hat{\mu}| < \frac{(1 - \hat{\delta})(1 - \theta)}{2\hat{\delta}\theta}, \\ \frac{1 - \hat{\delta}}{1 + \hat{\delta}} \left[\frac{1 - \theta}{\theta} + |\hat{\mu}| \right], & \frac{(1 - \hat{\delta})(1 - \theta)}{2\hat{\delta}\theta} \leq |\hat{\mu}| < \infty. \end{cases}$$

This is done iteratively in the algorithm described in Theorem 3 below. One needs to prevent $\hat{\varepsilon}$ from becoming too small. This means that $\hat{\delta}$ should be kept away from 1, which means that the lower bound on $1 - \theta + \theta|\mu|$ is allowed to be somewhat smaller than the upper bound. Preventing $\hat{\varepsilon}$ from becoming too small also means that $1 - \theta + |\hat{\mu}|$ cannot be too small. This may be unavoidable if one is interested in relative error $\theta = 1$, and the true answer, μ , is small.

Some notation is needed for this theorem. For any fixed $\alpha \in (0, 1)$, and $M > 0$, define the inverse of the functions $N_C(\cdot, \alpha)$, $N_B(\cdot, \alpha, M)$, and $N_{CB}(\cdot, \alpha, M)$,

$$N_C^{-1}(n, \alpha) := \frac{1}{\sqrt{n\alpha}},$$

$$N_B^{-1}(n, \alpha, M) := \min \left\{ b > 0 : \Phi(-b\sqrt{n}) + \frac{0.56M}{\sqrt{n}(1+b\sqrt{n})^3} \leq \frac{\alpha}{2} \right\},$$

$$N_{CB}^{-1}(n, \alpha, M) := \min(N_C^{-1}(n, \alpha), N_B^{-1}(n, \alpha, M)).$$

It then follows then by Chebyshev's inequality and the Berry-Esseen Inequality (see (14)) that

$$\text{Prob}[|\hat{\mu}_n - \mu| < \hat{\varepsilon}] \geq 1 - \alpha, \quad \text{provided } f \in \mathcal{C}_{\kappa_{\max}}, \text{ where } \hat{\varepsilon} = \sigma(f)N_{CB}^{-1}(n, \alpha, \kappa_{\max}^{3/4}),$$

and $\sigma(f) = \sqrt{\text{var}(f)}$ is the standard deviation of the integrand. Given a significance level, $\alpha \in (0, 1)$, let $\alpha_\sigma, \alpha_\mu, \alpha_1, \alpha_2, \dots$ be an infinite sequence of positive numbers all less than one, such that

$$(1 - \alpha_\sigma)(1 - \alpha_\mu)(1 - \alpha_1)(1 - \alpha_2) \cdots = 1 - \alpha. \quad (21)$$

For example, one might choose $\alpha_\sigma, \alpha_\mu$, and $\hat{\alpha}$ such that $(1 - \alpha_\sigma)(1 - \alpha_\mu)(1 - \hat{\alpha}) = 1 - \alpha$, and then

$$\alpha_i = 1 - (1 - \hat{\alpha})^{(a-1)a^{-i}}, \quad i \in \mathbb{N}, \quad \text{where } a \in (1, \infty). \quad (22)$$

Theorem 3. *Specify the following parameters defining the algorithm:*

- sample size for variance estimation, $n_\sigma \in \mathbb{N}$,
- initial sample size for mean estimation, $n_1 \in \mathbb{N}$,
- variance inflation factor for variance estimation, $\mathfrak{C} \in (1, \infty)$,
- factors for Step 2, $\hat{\delta}, \delta, \tilde{\delta} \in (0, 1)$, with $\delta < \tilde{\delta}$.
- uncertainty tolerance, $\alpha \in (0, 1)$, and a sequence $\alpha_\sigma, \alpha_\mu, \alpha_1, \alpha_2, \dots$ satisfying (21),
- the parameter $\theta \in [0, 1]$, used to define the general error criterion (19), and
- the error tolerance, $\varepsilon > 0$.

Let $\kappa_{\max} = \kappa_{\max}(n_\sigma, \alpha_\sigma, \mathfrak{C})$ as defined in (12). For any f lying in the cone of functions with bounded kurtosis, $\mathcal{C}_{\kappa_{\max}}$, do the following:

1. **Bounding the variance of the integrand from above.** Compute the sample variance, \hat{v}_{n_σ} using a simple random sample of size n_σ . Use this to approximate the variance of f by $\hat{\sigma}^2 = \mathfrak{C}^2 \hat{v}_{n_\sigma}$, as in (4). Compute the width of initial the confidence interval for the mean, $\hat{\varepsilon}_1 = \hat{\sigma} N_{CB}^{-1}(n_1, \alpha_1, \kappa_{\max}^{3/4})$.

2. Bounding the denominator in the error criterion from below. For $i = 1, 2, \dots$, do the following:

- a) Compute the sample average $\hat{\mu}_{n_i}$ using a simple random sample that is independent of those used to compute \hat{v}_{n_σ} and $\hat{\mu}_{n_1}, \dots, \hat{\mu}_{n_{i-1}}$.
- b) Compute $c = 1 - \theta + \theta \max(|\hat{\mu}_{n_i}| - \hat{\epsilon}_i, 0)$, a confident lower bound on $1 - \theta + \theta |\mu|$, according to (20). If $c \geq \delta[1 - \theta + \theta(|\hat{\mu}_{n_i}| + \hat{\epsilon})]$, then c is large enough. Set $\tau = i$ and go to Step 3.
- c) Else, compute the next tolerance for the sample mean

$$\hat{\epsilon}_0 = \begin{cases} \frac{(1 - \hat{\delta})(1 - \theta)}{\hat{\delta}\theta} - |\hat{\mu}_{n_i}|, & 0 \leq |\hat{\mu}_{n_i}| < \frac{(1 - \hat{\delta})(1 - \theta)}{2\hat{\delta}\theta}, \\ \frac{1 - \hat{\delta}}{1 + \hat{\delta}} \left[\frac{1 - \theta}{\theta} + |\hat{\mu}_{n_i}| \right], & \frac{(1 - \hat{\delta})(1 - \theta)}{2\hat{\delta}\theta} \leq |\hat{\mu}_{n_i}| < \infty, \end{cases}$$

$$\hat{\epsilon}_{i+1} = \max(\min(\hat{\epsilon}_0, \hat{\delta}\hat{\epsilon}_i), \delta\hat{\epsilon}_i).$$

- d) Define the next sample size, $n_{i+1} = N_{CB}(\hat{\epsilon}_{i+1}/\hat{\sigma}, \alpha_{i+1}, \kappa_{\max}^{3/4})$, increase i by one, and go to step a).

3. Computing the sample mean to sufficient accuracy. Compute the sample size $n = N_{CB}(c\epsilon/\hat{\sigma}, \alpha_\mu, \kappa_{\max}^{3/4})$. Compute $\bar{\mu} = \hat{\mu}_n$ using a simple random sample that is independent of those used to compute \hat{v}_{n_σ} and $\hat{\mu}_{n_1}, \dots, \hat{\mu}_{n_\tau}$. Terminate the algorithm.

If this algorithm terminates, then the general error criterion, (19), is satisfied.

Proof. In this algorithm there are a number of important random variables: the estimated upper bound on the standard deviation, $\hat{\sigma}$, the sample sizes n_1, \dots, n_τ, n , the number of iterations, τ , required to get a good lower bound c , and the final estimate of the mean $\bar{\mu} = \hat{\mu}_n$. These sample means are conditionally independent given the sequence of sample sizes. The probability that the final confidence interval is correct, is then no less than the probability that all of the confidence intervals are correct, conditioned on the sample sizes. Specifically,

$$\begin{aligned} \text{Prob} \left[\frac{|\bar{\mu} - \mu|}{1 - \theta + \theta |\mu|} \leq \epsilon \right] &\geq \text{Prob} [|\hat{\mu}_n - \mu| \leq c\epsilon \text{ \& } c \leq 1 - \theta + \theta |\mu|] \\ &= E \{ \text{Prob} [|\hat{\mu}_n - \mu| \leq c\epsilon \text{ \& } |\hat{\mu}_{n_\tau} - \mu| \leq \hat{\epsilon}_\tau \mid \hat{\sigma}, \tau, n_1, \dots, n_\tau, n] \} \\ &\geq E \{ \text{Prob} [|\hat{\mu}_n - \mu| \leq c\epsilon \text{ \& } |\hat{\mu}_{n_i} - \mu| \leq \hat{\epsilon}_i \forall i \mid \hat{\sigma}, \tau, n_1, \dots, n_\tau] \} \\ &\geq E_{\hat{\sigma}} \left\{ [(1 - \alpha_\mu)(1 - \alpha_1))(1 - \alpha_2) \cdots] 1_{[\sigma, \infty)}(\hat{\sigma}) \right\} \\ &\geq (1 - \alpha_\sigma)(1 - \alpha_\mu)(1 - \alpha_1)(1 - \alpha_2) \cdots = 1 - \alpha. \quad \square \end{aligned}$$

Remark 2. Step 2 in this algorithm is not needed for the case of pure absolute error $\theta = 0$, because $c = 1$ automatically, which is large enough. As suggested earlier, a difficulty may arise if $\mu \approx 0$ and $\theta \approx 1$, in which the algorithm may fail to converge in a reasonable number of steps and overall sample size. Step 2c has safeguards against making $\hat{\epsilon}_{i+1}$ too small compared to $\hat{\epsilon}_i$, but this may also increase the number of iterations, τ , necessary for completion of Step 2. Because it is difficult to know how large τ is for a given integrand, there is no rigorous bound on the cost of this algorithm yet.

5 Discussion

Put something here.

Acknowledgements

The authors gratefully acknowledge discussions with Erich Novak and Henryk Woźniakowski. The plots of the univariate fooling functions were prepared with the help of Nicholas Clancy and Caleb Hamilton.

References

- Bahadur RR, Savage LJ (1956) The nonexistence of certain statistical procedures in nonparametric problems. *Ann Math Stat* 27:1115–1122
- Gander W, Gautschi W (2000) Adaptive quadrature — revisited. *BIT* 40:84–101
- Hale N, Trefethen LN, Driscoll TA (2012) *Chebfun Version 4*
- Lin Z, Bai Z (2010) *Probability Inequalities*. Science Press and Springer-Verlag, Beijing and Berlin
- Petrov VV (1995) *Limit Theorems of Probability Theory: Sequences of Independent Random Variables*. Clarendon Press, Oxford
- Shampine LF (2008) Vectorized adaptive quadrature in matlab. *J Comput Appl Math* 211:131–140
- The MathWorks, Inc (2012) *MATLAB 7.12*. The MathWorks, Inc., Natick, MA
- Traub JF, Wasilkowski GW, Woźniakowski H (1988) *Information-Based Complexity*. Academic Press, Boston
- Wolfram Research Inc (2011) *Mathematica 8*