# A FAST AND ACCURATE NUMERICAL METHOD FOR THE LEFT TAIL OF SUMS OF INDEPENDENT RANDOM VARIABLES

NADHIR BEN RACHED, HÅKON HOEL, AND JOHANNES VINCENT MEO⋆

ABSTRACT. We present a flexible, deterministic numerical method for computing left-tail rare events of sums of non-negative, independent random variables. The method is based on iterative numerical integration of linear convolutions by means of Newtons–Cotes rules. The periodicity properties of convoluted densities combined with the Trapezoidal rule are exploited to produce a robust and efficient method, and the method is flexible in the sense that it can be applied to all kinds of non-negative continuous RVs. We present an error analysis and study the benefits of utilizing Newton–Cotes rules versus the fast Fourier transform (FFT) for numerical integration, showing that although there can be efficiency-benefits to using FFT, Newton–Cotes rules tend to preserve the relative error better, and indeed do so at an acceptable computational cost. Numerical studies on problems with both known and unknown rare-event probabilities showcase the method's performance and support our theoretical findings.

## 1. INTRODUCTION

For a sequence of non-negative and independent continuous random variables (RVs) $X_1, X_2, \ldots, X_n$, we seek to estimate the probability of failure

$$\alpha = \mathbb{P}\left(\sum_{i=1}^n X_i < \gamma\right)$$

with low relative error for small values of $\gamma > 0$. Left-tail rare-event problems of this kind are used for instance to estimate the outage probability in wireless communications, as is described further in the next paragraph. Our quite straightforward deterministic method approximates the density of $\sum_{i=1}^n X_i$ through iterative numerical convolution, and we show both theoretically and in numerical experiments that this approach is robust and that it performs very well in terms of accuracy and efficiency.

The main inspiration for our approach is Keich [25, 37], where a similar approach is studied for approximating the density of a sum of independent and identically distributed (i.i.d.) discrete RVs. The approaches are similar in how rounding errors propagate in computations of linear convolutions and both can be combined with the fast Fourier transform to speed up computations of linear convolutions, at the price of introducing more rounding errors. A notable difference is that for continuous RVs, the periodicity of convolutions of densities can be utilized to produce high-order convergence rates in the numerical integration. It is not clear if periodicity can be exploited to improve tractability also for discrete RVs.

Right and left tails of sums of RVs have gained significant attention in the literature due to their broad range of applications. In financial engineering, for example,

the value at risk for a portfolio based on multiple assets can be represented as the cumulative distribution function (CDF) of sums of RVs [9]. In the performance analysis of wireless communication systems, the outage probability/capacity can be expressed as the CDF of sums of RVs corresponding to either the fading channel envelopes or channel gains [17, 11]. Further applications extend to insurance risk and queueing systems. Within the Cramer-Lundberg model, the total sum of claims is modeled by a random sum of independent RVs, and the ruin probability is defined as the probability that this sum exceeds a large threshold [3].

Generally, a closed-form expression of the CDF of sums of RVs does not exist for most of the distributions. This is for instance the case for the Log-Normal distribution which has attracted a substantial interest [9, 19, 21, 17, 11, 13, 14, 41, 15, 3]. Although several closed-form approximations have been devised, their accuracy is not guaranteed and may degrade for specific parameter choices [11, 14, 41, 39, 30, 10, 27, 40, 36, 8]. Furthermore, these closed-form approximations are not generic, as they are generally tailored to specific distributions. Efficient numerical methods have been proposed in the literature to approximate the distributions of sums of RVs [29, 22, 32, 12, 19, 34, 35, 31]. For instance, in [29], Smolyak's algorithm, belonging to the family of numerical integration methods on sparse grids, has been developed for the accurate analysis of correlated Log-Normal power sums. Convenient numerical methods for Log-Normal characteristic functions have also been proposed, as seen in [12, 32, 35]. In [22], the authors used a saddle point approximation to evaluate the outage probability of wireless cellular networks. However, this method assumes the existence of the cumulant generating function, a requirement that is not met by many practical distributions, including the Log-Normal distribution. A general numerical approach, presented in [31], has also been developed for computing wireless outages. Similar to [22], this approach is general, provided that the moment generating function is known.

Monte Carlo (MC) methods are versatile tools employed to provide approximations of the CDF or complementary CDF of sums of RVs. However, it is widely recognized that naive MC simulations are computationally expensive, especially when addressing the right and left tails of sum distributions [26]. To mitigate this computational inefficiency, various efficient variance reduction techniques have been proposed. While much of the existing research has concentrated on the right tail of the sum distribution [3, 28, 18, 16, 24, 5, 7, 6, 4], the left tail region, which is the primary focus of this work, has only recently gained attention. In [9], the authors utilized the exponential twisting technique, a well-known importance sampling scheme, to efficiently estimate the left tail of sums of i.i.d. Log-Normal RVs. Additionally, in [17], two generic importance sampling schemes based on the hazard rate twisting technique of [24] were proposed to estimate the tail of the CDF of independent RVs. These algorithms have proven to be efficient for a wide range of well-known distributions within the context of wireless communication systems. An efficient importance sampling scheme has also been developed for the left tail of correlated Log-Normals [21]. This scheme was further enhanced in [2], where importance sampling and control variates were combined to achieve a further reduction in variance. Recently, state-dependent importance sampling has been proposed using a stochastic optimal control formulation [15]. This generic approach has been found to be efficient when considering rare event quantities that take the form of an expected value of some functions applied to sums of independent RVs.

The rest of this paper is organized as follows. Section 2 describes the numerical method for estimating left-tail rare events. Section 3 presents error and cost analysis of the rare-event estimation method both when using Newton–Cotes rules and FFT numerical integrators for discrete convolution. Section 4 studies the method

numerically on a collection of rare-event problems, and compares its performance to the saddle-point method for sums of Log-Normal RVs. Section 5 summarizes our findings and discusses possible future extensions.

## 2. The numerical method

In this section we construct an iterative Newton-Cotes quadrature method for approximating linear convolutions of probability densities, and ultimately the probability of failure. For the sake of streamlining the exposition, we will restrict ourselves to the setting with i.i.d. RVs, but an extension to independent (and not identically distributed) RVs is exemplified in Section 4.3. Let $f : [0, \infty) \to [0, \infty)$ denote the probability density function (PDF) of the i.i.d. RVs $X_1, \ldots, X_n$. We consider the problem of estimating the rare event

$$\alpha = \mathbb{P}\left(\sum_{i=1}^{n} X_i < \gamma\right) = \int_0^\gamma f_{S_n}(x)\, dx$$

for small values of $\gamma > 0$, where $f_{S_n}$ is the PDF of $S_n = \sum_{i=1}^n X_i$. For simplicity, we will make the assumption that $f(0) := \lim_{x \downarrow 0} f(x) = 0$ throughout this section, and describe the extension to settings when $f(0) \neq 0$ in equation (9).

2.1. **Numerical integration.** The first step in our deterministic approach for estimating $\alpha$ by discrete linear convolution is to observe that due to the independence and identical distribution of $X_1, \ldots, X_n$, the PDF of $S_n := \sum_{i=1}^n X_i$ is equal to the $n-$fold linear convolution of the density $f$:

$$f_{S_n} = f^{*n} := \text{n-fold convolution of } f.$$

In particular, we have that

$$f_{S_2}(x) = f * f(x) = \int_0^x f(y)f(x-y)\, dy,$$

and for any integers $k, \ell, n \geq 1$ such that $n = k + \ell$, it holds that

$$f_{S_n}(x) = f^{*k} * f^{*\ell}(x) = \int_0^x f^{*k}(y)f^{*\ell}(x-y)\, dy. \tag{1}$$

This means that the probability of failure can be expressed as

$$\alpha = \mathbb{P}\left(\sum_{i=1}^{n} X_i < \gamma\right) = \int_0^\gamma f^{*n}(x)\, dx,$$

and that a good approximation of $\alpha$ can be obtained from a good approximation of $f^{*n}$.

The second step is to approximate $f^{*n}$ on $[0, \gamma]$ by numerical integration. Since $f(0) = 0$, any integrand of the form $g(y; x) = f^{*k}(y)f^{*\ell}(x-y)$ is periodic on $[0, x]$ for any $x \in [0, \gamma]$, as

$$g(0; x) = f^{*k}(0)f^{*\ell}(x) = 0, \quad \text{and} \quad g(x; x) = f^{*k}(x)f^{*\ell}(0) = 0.$$

Conveniently, the trapezoidal rule yields a high order of convergence for sufficiently smooth periodic integrands, and it is therefore a suitable quadrature rule for linear convolution of (1), as described through the following steps: Consider the uniform mesh

$$x_j = jh \quad j = 0, 1, \ldots, N \quad \text{with} \quad h = \frac{\gamma}{N},$$

and the discrete function

$$\bar{f}(x_j) := f(x_j) \qquad j = 0, 1, \ldots, N.$$

Let the discrete approximation of $f^{*2}$ be given by the trapezoidal rule/discrete linear convolution

$$\bar{f}^{\circledast 2}(x_k) = h \sum_{j=0}^{k} \bar{f}(x_j)\bar{f}(x_k - x_j) \qquad k = 0, \dots, N. \tag{2}$$

The operator notation $\circledast$ represents discrete linear convolution scaled by the step-size factor $h$. It is introduced to distinguish discrete linear convolution from continuous-space linear convolution, which we denoted by $*$. To define higher-order discrete convolutions we proceed as follows: for any two $\mathbb{R}^{N+1}$-vectors $\bar{g}_1 = (\bar{g}_1(x_0), \dots, \bar{g}_1(x_N))$ and $\bar{g}_2 = (\bar{g}_2(x_0), \dots, \bar{g}_2(x_N))$, let

$$(\bar{g}_1 \circledast \bar{g}_2)(x_k) := h \sum_{j=0}^{k} \bar{g}_1(x_j)\bar{g}_2(x_k - x_j) \,.$$

It then holds that $\circledast$ is an associative operation, as for any $\bar{g}_1, \bar{g}_2, \bar{g}_3 \in \mathbb{R}^{N+1}$,

$$\left(\bar{g}_1 \circledast (\bar{g}_2 \circledast \bar{g}_3)\right)(x_k) = h^2 \sum_{j_1 + j_2 + j_3 = k} \bar{g}_1(x_{j_1})\bar{g}_2(x_{j_2})\bar{g}_3(x_{j_3}) = \left((\bar{g}_1 \circledast \bar{g}_2) \circledast \bar{g}_3\right)(x_k),$$

This shows that $\bar{f}^{\circledast n}$ is a well-defined operation for any $n \geq 2$, since it does not matter which order the convolutions are taken in. So for any $n > 2$ and mesh point $x_k$,

$$\begin{aligned}
\bar{f}^{\circledast n}(x_k) &= (\bar{f}^{\circledast(n-1)} \circledast \bar{f})(x_k) = (\bar{f}^{\circledast(n-2)} \circledast \bar{f}^{\circledast 2})(x_k) \\
&= \cdots = (\bar{f} \circledast \bar{f}^{\circledast(n-1)})(x_k) = h^{n-1} \sum_{j_1 + \cdots + j_n = k} \bar{f}(x_{j_1}) \cdots \bar{f}(x_{j_n}) \,.
\end{aligned} \tag{3}$$

2.2. **Efficient computation of $\bar{f}^{\circledast n}$.** In this section we describe the sequence of discrete convolutions used to compute $\bar{f}^{\circledast n}$ efficiently. This matters for the performance of the method when implemented on a computer.

Let $m \in \mathbb{N}$ denote the largest integer such that $m \leq \log_2(n)$ and for $\ell = 2, \dots, m$, we minimize the number of convolutions through computing

$$\bar{f}^{\circledast 2^{\ell}}(x_k) = h \sum_{j=0}^{k} \bar{f}^{\circledast 2^{\ell-1}}(x_j)\bar{f}^{\circledast 2^{\ell-1}}(x_k - x_j) \qquad k = 0, \dots, N. \tag{4}$$

If $2^m = n$, then the above computation represents discrete approximation $\bar{f}^{\circledast n}$ of the density of $\sum_{j=1}^{n} X_j$, otherwise we obtain the approximation by

$$\bar{f}^{\circledast n}(x_k) = h \sum_{j=0}^{k} \bar{f}^{\circledast 2^m}(x_j)\bar{f}^{\circledast(n-2^m)}(x_k - x_j) \qquad k = 0, \dots, N,$$

where $\bar{f}^{\circledast(n-2^m)}$ is computed in at most $m-1$ steps by a similar iterative application of the Trapezoidal rule.

Lastly, the probability of failure is also approximated by a suitable closed Newton–Cotes formula:

$$\bar{\alpha}_N := \sum_{j=0}^{N} w_j \bar{f}^{\circledast n}(x_j), \tag{5}$$

where the weights sum to the the length of the interval $[0, \gamma]$, meaning $\sum_{j=0}^{N} w_j = \gamma$, and we restrict ourselves to the class of Newton–Cotes formulas with non-negative weights, which means formulas with degree $d \leq 8$, cf. [33, Chapter 7]. See Theorem 1 for further details on how to choose the Newton–Cotes rule. Since it holds that $\bar{f}^{\circledast j}(0) = 0$ for any $j \geq 1$, all of the numerical integration above can be viewed as applications of the Trapezoidal rule.

| Fading Type | PDF |
|---|---|
| Rayleigh | $\frac{2x}{\Omega}\exp\left(-\frac{x^2}{\Omega}\right)$ |
| Nakagami-m | $\frac{2m^m}{\Omega^m\Gamma(m)}x^{2m-1}\exp\left(-\frac{m}{\Omega}x^2\right)$ |
| Rice | $\frac{2(K+1)x}{\Omega}\exp\left(-K-\frac{K+1}{\Omega}x^2\right)I_0\left(2\sqrt{\frac{K(K+1)}{\Omega}}x\right)$ |
| Weibull | $k\left(\frac{\beta}{\Omega}\right)^k x^{k-1}\exp\left(-\left(\frac{x\beta}{\Omega}\right)^k\right),\quad \beta=\Gamma\left(1+\frac{1}{k}\right)$ |
| Log-Normal | $\frac{1}{x\sigma\sqrt{2\pi}}\exp\left(-\frac{(\log(x)-\mu)^2}{2\sigma^2}\right)$ |
| generalized Gamma | $\frac{p(\frac{\beta}{\Omega})^d}{\Gamma(\frac{d}{p})}x^{d-1}\exp\left(-(\frac{\beta}{\Omega}x)^p\right),\quad \beta=\frac{\Gamma\left(\frac{d+1}{p}\right)}{\Gamma\left(\frac{d}{p}\right)}$ |
| $\kappa-\mu$ | $\frac{2\mu(K+1)^{\frac{1+\mu}{2}}x^\mu}{\Omega^{\frac{1+\mu}{2}}\kappa^{\frac{\mu-1}{2}}}\exp\left(-\mu K-\frac{K+1}{\Omega}\mu x^2\right)I_{\mu-1}\left(2\mu\sqrt{\frac{K(K+1)}{\Omega}}x\right)$ |

TABLE 1. Some common PDF satisfying $f(0)=0$. Functions $\Gamma(\cdot)$ and $I_\xi(\cdot)$ are respectively the Gamma function and the modified Bessel function of the first kind and order $\xi$ [20]

*Remark* 1. The assumption that $f(0)=0$ is not restrictive as it is satisfied by most of the common fading channel envelopes. A non exhaustive list of common fading channel having the previous assumption satisfied is in Table I.

*Remark* 2. For the simpler setting when $f$ is a probability mass function instead of a probability density function, a similar approach to approximating the density of sums $Y_1+\cdots+Y_n$ where $Y_k\overset{iid}{\sim} f$ through FFT-based convolution has been studied in [38].

2.3. **Implementations of discrete linear convolution.** There are two standard approaches to implement the above discrete linear convolution (4) in most programming languages, with different strengths and weaknesses:

(1) **Direct convolution:** For each $k=0,1,\ldots,N$ compute the RHS sliding sum of (4). In Matlab, this can be achieved by calling the **conv**() function as follows:

```
fBar_2l = conv(g,g);
fBar_2l   = fBar_2l(1:N+1)*h;
```

for an input vector $g:=\bar{f}^{2^{\ell-1}\circledast}\in\mathbb{R}^{N+1}$. The computational cost of this function call, measured in number of floating point operations is $\mathcal{O}(N^2)$, which is quite high. But Theorem 2 and our numerical experiments in Section 4.1 show that direct convolution is accurate even for very rare events and it does not appear sensitive to round-off errors.

(2) **FFT-based convolution:** The second approach is to append/pad $N$ zeros to the vector $f^{\circledast 2^{\ell-1}}\in\mathbb{R}^{N+1}$, and use the fast Fourier transform (FFT) to compute the linear convolution as follows:

$$\bar{f}^{\circledast 2^{\ell-1}}=[\bar{f}^{\circledast 2^{\ell-1}}(x_0),\ldots,\bar{f}^{\circledast 2^\ell}(x_N),\underbrace{0,\ldots,0}_{N}]$$

$$\bar{f}^{\circledast 2^\ell}=\mathbf{IFFT}(\mathbf{FFT}(\bar{f}^{\circledast 2^{\ell-1}}).\,\hat{}2)\times h \tag{6}$$

$$\bar{f}^{\circledast 2^\ell}=[\bar{f}^{\circledast 2^\ell}(x_0),\ldots,\bar{f}^{\circledast 2^\ell}(x_N)].$$

which in Matlab takes the form

```
g         = [g zeros(1,N)];
```

```
fBar_2l = ifft(fft(g).^2)*h;
fBar_2l = fBar_2l(1:N+1);
```

An advantage with this approach is that the computational cost of the three assignments (6) is $\mathcal{O}(N \log(N))$, which for large $N$ will be far lower than the cost of direct convolution. A downside is that FFT-based convolution can be very sensitive to rounding errors when the floating point precision is low, which occurs when the machine epsilon is greater or of the similar magnitude as $\alpha$, cf. Section 3.1. This is illustrated in an numerical example in Section 4.1.

*Remark* 3. The method straightforwardly extends to settings where $X_1, \ldots, X_n$ are independent but not identically distributed, cf. (9).

## 3. THEORETICAL RESULTS

In this section we first prove that $\bar{f}^{\circledast n}(x_k) \to f^{*n}(x_k)$ and that $\bar{\alpha}_N \to \alpha$ as $N \to \infty$ and obtain convergence rates for these results in Lemma 1 and Theorem 1, respectively. Thereafter, we bound the relative approximation error of $\mathrm{fl}[\bar{\alpha}_N] \approx \bar{\alpha}_N$ for direct convolution in terms of the floating-point precision machine epsilon in Lemma 2. This leads to the upper bound for the computable approximation $\mathrm{fl}[\bar{\alpha}_N] \approx \alpha$ that is described in Theorem 2. A similar results for FFT-based convolution is given in 3. Finally, we bound the computational cost of our method in Theorem 4 and combine the results on error and cost to compare the efficiency of direct convolution and FFT-based convolution in (16).

Let $C_0^{2p}([0, \gamma])$ denote the set of $2p$ times continuously differentiable functions on $[0, \gamma]$ for which $f^{(k)}(0) = 0$ for all $k \in \{0, 1, \ldots, 2p - 1\}$, and let $C_{0,\mathrm{per}}^{2p}([0, \gamma]) \subset C_0^{2p}([0, \gamma])$ denote the subset of such functions that also are periodic on $[0, \gamma]$ up to the $(2p-1)$-th derivative, meaning that $f^{(k)}(0) = f^{(k)}(\gamma)$ for all $k \in \{0, 1, \ldots, 2p - 1\}$.

The first Lemma provides a convergence rate for $\bar{f}^{\circledast m}(x_k) \to f^{*m}(x_k)$ as $N \to \infty$. It does not take rounding errors into account.

**Lemma 1.** *Let $f \in C_0^{2p}([0, \gamma])$ for some integer $p \geq 1$ and let $f^{\circledast n}(x_k)$ for $n \geq 2$ be defined by (3). Then, there exists a constant $C_1 > 0$ that depends on $p$ such that*

$$|\bar{f}^{\circledast n}(x_k) - f^{*n}(x_k)| \leq (1 + x_k)^{n-2} C_1 C_2(n, x_k) N^{-2p} \text{ for all } k = 0, 1, \ldots, N, \quad (7)$$

*where*

$$C_2(n, x_k) :=$$

$$\max_{m=2,\ldots,n} \max_{x_{j_1} + \cdots + x_{j_{m-1}} = x_k} \max_{0 \leq y \leq x_{j_1}} \left| \frac{d^{2p}}{dy^{2p}} f^{*(n+1-m)}(y) f(x_{j_1} - y) \right| \prod_{\ell=2}^{m-1} f(x_{j_\ell})$$

*with the conventions that $\prod_{\ell=2}^{m-1} f(x_{j_\ell}) \equiv 1$ when $m = 2$ and $f^{*1} = f$.*

*Proof.* Assume that (7) holds for all $2 \leq n \leq \bar{n}$ for some $\bar{n} \geq 2$. Recalling from (3) that $\bar{f}^{\circledast(\bar{n}+1)}(x_k) = (\bar{f}^{\circledast \bar{n}} \circledast \bar{f})(x_k)$ and also that $\bar{f}(x_j) = f(x_j)$ for all

$j = 0, 1, \ldots, N$, we obtain

$$|\bar{f}^{\circledast(\bar{n}+1)}(x_k) - f^{*(\bar{n}+1)}(x_k)| =$$

$$h \left| \sum_{j=0}^{k} \bar{f}^{\circledast\bar{n}}(x_j)\bar{f}(x_{k-j}) - f^{*\bar{n}}(x_j)f(x_{k-j}) \right|$$

$$+ \left| h \sum_{j=0}^{k} f^{*\bar{n}}(x_j)f(x_{k-j}) - \int_0^{x_k} f^{*\bar{n}}(y)f(x_k - y)dy \right| =: I + II.$$

For the first term,

$$I = h \left| \sum_{j=0}^{k} \bar{f}^{\circledast\bar{n}}(x_j)\bar{f}(x_{k-j}) - f^{*\bar{n}}(x_j)f(x_{k-j}) \right|$$

$$\leq h \sum_{j=0}^{k} |\bar{f}^{\circledast\bar{n}}(x_j) - f^{*\bar{n}}(x_j)|f(x_{k-j})$$

$$\leq hC_1(1 + x_k)^{\bar{n}-2} \sum_{j=0}^{k} C_2(\bar{n}, x_j)f(x_{k-j})$$

$$\leq (k-1)h \times (1 + x_k)^{\bar{n}-2}C_1 C_2(\bar{n}+1, x_k)N^{-2p}$$

$$\leq x_k(1 + x_k)^{\bar{n}-2}C_1 C_2(\bar{n}+1, x_k)N^{-2p},$$

where the penultimate inequality follows from the change of subindex in $x_{k-j} = x_{j_m}$ and

$$f(x_{k-j})C_2(\bar{n}, x_j)$$

$$= f(x_{k-j}) \max_{m=2,\ldots,\bar{n}} \max_{x_{j_1}+\cdots+x_{j_{m-1}}=x_j} \max_{0 \leq y \leq x_{j_1}} \left| \frac{d^{2p}}{dy^{2p}} f^{*(\bar{n}+1-m)}(y)f(x_{j_1} - y) \right| \prod_{\ell=2}^{m-1} f(x_{j_\ell})$$

$$\leq \max_{m=2,\ldots,\bar{n}} \max_{x_{j_1}+\cdots+x_{j_m}=x_k} \max_{0 \leq y \leq x_{j_1}} \left| \frac{d^{2p}}{dy^{2p}} f^{*(\bar{n}+1-m)}(y)f(x_{j_1} - y) \right| \prod_{\ell=2}^{m} f(x_{j_\ell})$$

$$= \max_{m=3,\ldots,\bar{n}+1} \max_{x_{j_1}+\cdots+x_{j_{m-1}}=x_k} \max_{0 \leq y \leq x_{j_1}} \left| \frac{d^{2p}}{dy^{2p}} f^{*(\bar{n}+2-m)}(y)f(x_{j_1} - y) \right| \prod_{\ell=2}^{m-1} f(x_{j_\ell})$$

$$\leq C_2(\bar{n}+1, x_k).$$

The second term is the quadrature error of the composite Trapezoidal rule applied to the integrand $g(y) = f^{*\bar{n}}(y)f(x_k - y)$. Thanks to $g \in C^{2p}_{0,\text{per}}([0, x_k])$, we obtain that

$$II \leq C_1 \max_{0 \leq y \leq x_k} \left| \frac{d^{2p}}{dy^{2p}} f^{*\bar{n}}(y)f(x_k - y) \right| N^{-2p} \leq C_1 C_2(\bar{n}+1, x_k)N^{-2p},$$

for the constant $C_1 > 0$ introduced above. This yields,

$$|\bar{f}^{\circledast(\bar{n}+1)}(x_k) - f^{*(\bar{n}+1)}(x_k)| \leq (1 + x_k)^{\bar{n}-1}C_1 C_2(\bar{n}+1, x_k).$$

We next verify that (7) holds for $n = 2$. Since $\bar{f}^{\circledast 2}(x_k)$ is the composite Trapezoidal rule approximation of $f^{*2}(x_k)$, as is apparent from

$$|\bar{f}^{\circledast 2}(x_k) - f^{*2}(x_k)| = \left| h \sum_{j=0}^{k} f(x_j)f(x_{k-j}) - \int_0^{x_k} f(y)f(x - y)dy \right|,$$

and $f \in C_0^{2p}([0, \gamma])$ implies that $g(y) := f(y)f(x_k - y)$ belongs to $C_{0,\mathrm{per}}^{2p}([0, x_k])$, it follows from [33, Chapter 7.6] that

$$|\bar{f}^{\circledast 2}(x_k) - f^{*2}(x_k)| \le C_1 \underbrace{\max_{0 \le y \le x_k} \left| \frac{d^{2p}}{dy^{2p}} f(y)f(x_k - y) \right|}_{=C_2(2, x_k)} N^{-2p},$$

for the constant $C_1 > 0$ introduced above. The proof follows by induction. $\qquad \square$

The next theorem proves a convergence rate for $\alpha_N \to \alpha$ as $N \to \infty$ in the setting of no rounding errors.

**Theorem 1.** *Let $f \in C_0^{2p}([0, \gamma])$ for some integer $p \ge 1$, and let $r \le 2p$ be the order of convergence for the Newton–Cotes rule with non-negative weights that is used to compute $\bar{\alpha}$ in (5). Then there exist a constant $C_3 > 0$ that depends on $r$ such that*

$$|\bar{\alpha}_N - \alpha| \le \gamma C_1 \overline{C}_2(n, \gamma) N^{-2p} + C_3 \max_{x \in [0, \gamma]} \left| \frac{d^r}{dx^r} f^{*n}(x) \right| N^{-r}, \qquad (8)$$

*where $\overline{C}_2(n, \gamma) := \max_{k=0,\dots,N}(1 + x_k)^{n-2} C_2(n, x_k)$, and the constant $C_1$ and the mapping $C_2(n, x_k)$ are defined in Lemma 1.*

*Proof.* For a closed Newton–Cotes formula with convergence rate $r \le 2p$, it follows by [23, Chapter 7.1.1] that

$$\left| \alpha - \sum_{j=0}^{N} w_j f^{*n}(x_j) \right| \le C_3 N^{-r}.$$

The triangle inequality and the non-negativity of the weights $w_j$ yield the final bound

$$|\alpha - \bar{\alpha}_N| = \left| \alpha - \sum_{j=0}^{N} w_j \bar{f}^{\circledast n}(x_j) \right|$$

$$\le \left| \alpha - \sum_{j=0}^{N} w_j f^{*n}(x_j) \right| + \sum_{j=0}^{N} w_j |f^{*n}(x_j) - \bar{f}^{\circledast n}(x_j)|$$

$$\le C_3 \max_{x \in [0, \gamma]} \left| \frac{d^r}{dx^r} f^{*n}(x) \right| N^{-r} + \sum_{j=0}^{N} w_j (1 + x_k)^{n-2} C_1 C_2(n, x_j) N^{-2p}$$

$$\le C_3 \max_{x \in [0, \gamma]} \left| \frac{d^r}{dx^r} f^{*n}(x) \right| N^{-r} + \gamma C_1 \overline{C}_2(n, \gamma) N^{-2p}$$

$$\qquad \square$$

*Remark 4.* If $f \notin C_0^{2p}[0, \gamma]$, but we have $f \in C^2[0, \gamma]$ (which is the case if e.g. $f(0) \ne 0$ or $f'(0) \ne 0$), then the slightly altered direct convolution (compare to (4))

$$\bar{f}^{\circledast 2^\ell}(x_k) = h \sum_{j=1}^{k-1} \bar{f}^{\circledast 2^{\ell-1}}(x_j) \bar{f}^{\circledast 2^{\ell-1}}(x_k - x_j)$$

$$+ h \frac{\bar{f}^{\circledast 2^{\ell-1}}(x_0) \bar{f}^{\circledast 2^{\ell-1}}(x_k) + \bar{f}^{\circledast 2^{\ell-1}}(x_k) \bar{f}^{\circledast 2^{\ell-1}}(x_0)}{2} \qquad (9)$$

lead to the error bounds as in Lemma 1 and Theorem 1 with $p = 1$.

3.1. **Rounding errors.** In practice, approximations of $\alpha$ are computed using floating-point arithmetic where a float is represented by $x = s \times b \times 2^e$ with sign $s$ (1-bit), the significand $b \in [1, 2)$ and exponent $e$. The standard IEEE 754 64-bit floats, for example, has $p = 53$-bit significand precision (52-bits stored) and 11-bit exponent. For estimating rounding errors in relative error, the machine epsilon $\varepsilon = 2^{-p}$, which is equal to half the distance between the number 1 and the closest floating point number to 1, is important. For an $x \in \mathbb{R}$, let $\mathrm{fl}[x]$ denote the closest number to $x$ among the floating point numbers. Then it holds that $|x - \mathrm{fl}[x]| \leq (1 + \varepsilon)|x|$.

More generally, we let $\mathrm{fl}[\bar{\alpha}_N]$ denote the value of $\bar{\alpha}_N$ that is obtained when **all underlying arithmetic operations** are computed with the given floating point precision, and thus possibly all being subject to rounding errors, and similarly also for $\mathrm{fl}\big[\bar{f}^{\circledast n}(x_k)\big]$. Observe that this notation is recursive, it assumes that a quantity is computed in a uniquely specified manner (otherwise it would not be clear how to estimate rounding errors), and it is extremely compact, as is illustrated when applying it to the formula (4):

$$\mathrm{fl}\Big[\bar{f}^{\circledast 2^{\ell}}(x_k)\Big] = \mathrm{fl}\Bigg[\mathrm{fl}[h]\mathrm{fl}\Bigg[\sum_{j=0}^{k}\mathrm{fl}\Big[\mathrm{fl}\Big[\bar{f}^{\circledast 2^{\ell-1}}(x_j)\Big]\mathrm{fl}\Big[\bar{f}^{\circledast 2^{\ell-1}}(x_k - x_j)\Big]\Big]\Bigg]\Bigg].$$

**Lemma 2** (Rounding error direct convolution)**.** *Let $\bar{\alpha}_N$ be computed by direct convolution, $\mathrm{fl}[\bar{\alpha}_N]$ be computed with floating point arithmetic with machine epsilon $\varepsilon > 0$ and let $n$ be the number of i.i.d RVs in the underlying sum. Furthermore, set $\bar{m} = \lceil \log_2(n) \rceil$. Assume that for each $x$ in the codomain of $f$ we have $|\mathrm{fl}[x] - x| \leq xc\varepsilon$ for some $c \in (0, 1]$. Moreover, assume that $N \geq 2^{10}$, $\varepsilon \leq 2^{-53}$ and that $2^{2\bar{m}+2}N\varepsilon < 1/10$. Then it holds that*

$$|\mathrm{fl}[\bar{\alpha}_N] - \bar{\alpha}_N| \leq 4\bar{\alpha}_N nN\varepsilon.$$

*Proof.* We begin the proof by looking at some results from [25] showing how rounding errors propagate when adding and multiplying already estimated values. Assume that we have a set of non-negative real numbers $A = \{a_1, a_2, \ldots, a_N\}$ that are estimated by use of floating point arithmetic with an arbitrary number of operations used for calculating the approximations. We denote the estimated values $\tilde{A} = \{\tilde{a}_1, \tilde{a}_2, \ldots, \tilde{a}_N\}$ and have that the absolute accumulated error of our estimates $\tilde{a}_i$ can be bounded by some constant $c_a > 0$, giving $|\tilde{a}_i - a_i| \leq |a_i|c_a\varepsilon$. Moreover, let $\mathrm{fl}[h]$ be our floating point estimation of the step length $h \in \mathbb{R}_{>0}$ with $|\mathrm{fl}[h] - h| \leq h\varepsilon$. Based on the proof of [25, Lemma 3], it is straightforward to check that

$$|\mathrm{fl}[\mathrm{fl}[h]\tilde{a}_i\tilde{a}_j] - ha_ia_j| \leq ha_ia_j\left[2 + 2c_a + (1 + 4c_a + c_a^2)\varepsilon + (2c_a + 2c_a^2)\varepsilon^2 + c_a^2\varepsilon^3\right]\varepsilon, \tag{10}$$

for $i, j \in \{1, 2, \ldots, N\}$. Then, by letting $S_k = \sum_{j=0}^{k} ha_ja_{k-j}, k \in \{1, 2, \ldots, N\}$ we have from [25, Lemma 2] that

$$\left|\mathrm{fl}\left[\sum_{j=0}^{k}\mathrm{fl}[\mathrm{fl}[h]\tilde{a}_i\tilde{a}_{k-j}]\right] - S_k\right| \leq$$
$$S_k\left[k + 2 + 2c_a + (1 + 4c_a + c_a^2)\varepsilon + (2c_a + 2c_a^2)\varepsilon^2 + c_a^2\varepsilon^3\right]\varepsilon(1 + k\varepsilon), \quad (11)$$

when (10) holds for each term in the sum and $(N + c_{ha_ja_{k-j}})\varepsilon < 1$, where

$$c_{ha_ja_{k-j}} = 2 + 2c_a + (1 + 4c_a + c_a^2)\varepsilon + (2c_a + 2c_a^2)\varepsilon^2 + c_a^2\varepsilon^3.$$

Note that 11 only holds as long as all terms of $S_k$ have the same sign, which in our case is non-negative. For simplicity, we will bound the rounding error of all $S_k$ by

inserting $N$ instead of $k$ in 11, yielding

$$\left| \mathrm{fl}\left[ \sum_{j=0}^{k} \mathrm{fl}[\mathrm{fl}[h]\tilde{a}_i \tilde{a}_{k-j}] \right] - S_k \right| \leq$$

$$S_k \left[ N + 2 + 2c_a + (1 + 4c_a + c_a^2)\varepsilon + (2c_a + 2c_a^2)\varepsilon^2 + c_a^2\varepsilon^3 \right]\varepsilon(1 + N\varepsilon)$$

for all $k \in \{0, 1, \dots N\}$.

We now move on to prove the following statement by induction on $l$:

$$\left| \mathrm{fl}\left[ \bar{f}^{\circledast 2^l}(x_k) \right] - \bar{f}^{\circledast 2^l}(x_k) \right| \leq \bar{f}^{\circledast 2^l}(x_k)(2^{l+1} - 2)N\varepsilon. \tag{12}$$

For $l = 1$ we first observe that from (10) and our assumptions we have

$$\begin{aligned}
\left| \mathrm{fl}\left[ h\bar{f}(x_j)\bar{f}(x_{k-j}) \right] - h\bar{f}(x_j)\bar{f}(x_{k-j}) \right| &\leq h\bar{f}(x_j)\bar{f}(x_{k-j})[2 + 2c + (1 + 4c + c^2)\varepsilon \\
&\qquad + (2c + 2c^2)\varepsilon^2 + c^2\varepsilon^3]\varepsilon \\
&\leq h\bar{f}(x_j)\bar{f}(x_{k-j})5\varepsilon.
\end{aligned}$$

Then, from (11) we have that

$$\begin{aligned}
\left| \mathrm{fl}\left[ \bar{f}^{\circledast 2}(x_k) \right] - \bar{f}^{\circledast 2}(x_k) \right| &\leq \bar{f}^{\circledast 2}(x_k)(N+5)(1+N\varepsilon)\varepsilon \\
&\leq \bar{f}^{\circledast 2}(x_k)\left( N + 5 + N^2\varepsilon + 5N\varepsilon \right)\varepsilon \leq \bar{f}^{\circledast 2}(x_k)2N\varepsilon,
\end{aligned}$$

as we needed to show.

Assume now that (12) holds for some $q \in \mathbb{N} \setminus \{0\}$, that is

$$\left| \mathrm{fl}\left[ \bar{f}^{\circledast 2^q}(x_k) \right] - \bar{f}^{\circledast 2^q}(x_k) \right| \leq \bar{f}^{\circledast 2^q}(x_k)(2^{q+1} - 2)N\varepsilon$$

Furthermore, assume that given some value for $N$ our $q$ satisfies $2^{2q+2}N\varepsilon < 1/10$. Then we have from (11) that

$$\begin{aligned}
\left| \mathrm{fl}\left[ \bar{f}^{\circledast 2^{q+1}}(x_k) \right] - \bar{f}^{\circledast 2^{q+1}}(x_k) \right| &\leq \bar{f}^{\circledast 2^{q+1}}(x_k)\Bigg[ N + 2 + 2\left( 2^{q+1} - 2 \right)N \\
&\qquad + \left( 1 + 4\left[ 2^{q+1} - 2 \right]N + \left[ 2^{q+1} - 2 \right]^2 N^2 \right)\varepsilon \\
&\qquad + \left( 2\left[ 2^{q+1} - 2 \right]N + 2\left[ 2^{q+1} - 2 \right]^2 N^2 \right)\varepsilon^2 \\
&\qquad + \left( 2^{q+1} - 2 \right)^2 N^2\varepsilon^3 \Bigg]\varepsilon(1 + N\varepsilon) \\
&\leq \bar{f}^{\circledast 2^{q+1}}(x_k)\Bigg[ 2 + \left( 2^{q+2} - 3 \right)N + \varepsilon + \frac{1}{10} + \frac{N}{10} \\
&\qquad + \frac{\varepsilon}{10} + \frac{2N\varepsilon}{10} + \frac{N\varepsilon^2}{10} \Bigg]\varepsilon(1 + N\varepsilon) \\
&\leq \bar{f}^{\circledast 2^{q+1}}(x_k)\left[ \left( 2^{q+2} - \frac{5}{2} \right)N + \frac{N}{10} \right]\varepsilon \\
&\leq \bar{f}^{\circledast 2^{q+1}}(x_k)\left( 2^{q+2} - 2 \right)N\varepsilon.
\end{aligned}$$

Proceeding with the last step we need to calculate the accumulated rounding error for $\bar{\alpha}_N = \sum_{j=0}^{N} w_j \bar{f}^{\circledast n}(x_j)$. We then set $m = \lfloor \log_2(n) \rfloor$. Next we need to consider two separate cases, one where $m = \log_2(n)$, and the alternative case $m < \log_2(n)$. In the former case we have

$$\bar{\alpha}_N = \sum_{j=0}^{N} w_j \bar{f}^{\circledast 2^{\log_2(n)}}(x_j) = \sum_{j=0}^{N} w_j \bar{f}^{\circledast 2^m}(x_j).$$

In the latter case we calculate $\bar{f}^{\circledast 2^m}(x_j)$ and $\bar{f}^{\circledast(n-2^m)}(x_j)$, which can be done in at most $m-1$ steps. We then have that $\bar{f}^{\circledast n}(x_j) = \bar{f}^{\circledast 2^m}(x_j) \circledast \bar{f}^{\circledast n-2^m}(x_j)$, which would have an error bounded by $\bar{f}^{\circledast 2^{m+1}}$. Therefore, by setting $\bar{m} = \lceil \log_2(n) \rceil$, we have that the error of $\bar{f}^{\circledast n}(x_j)$ is bounded by the error of $\bar{f}^{\circledast 2^{\bar{m}}}(x_j)$.

Moving on we have from [25, Lemma 3], the bound (12) and from our assumptions that

$$
\begin{aligned}
\left| \mathrm{fl}\left[w_j \bar{f}^{\circledast n}(x_j)\right] - w_j \bar{f}^{\circledast n}(x_j) \right| \leq & w_j \bar{f}^{\circledast 2^{\bar{m}}}(x_j)[2 + \left(2^{\bar{m}+1} - 2\right) N \\
& + \left(1 + \left(2^{\bar{m}+2} - 4\right) N\right) \varepsilon \\
& + \left(2^{\bar{m}+1} - 2\right) N \varepsilon^2]\varepsilon \\
\leq & w_j \bar{f}^{\circledast 2^{\bar{m}}}(x_j) \left[2 + \left(2^{\bar{m}+1} - 2\right) N + \varepsilon + \frac{1}{10} + \frac{\varepsilon}{10}\right] \varepsilon \\
\leq & w_j \bar{f}^{\circledast 2^{\bar{m}}}(x_j) \left[\left(2^{\bar{m}+1} - 2\right) N + 3\right] \varepsilon.
\end{aligned}
$$

Then, by [25, Lemma 2] we get the following bound for the rounding error of $\bar{\alpha}_N$:

$$
\begin{aligned}
|\mathrm{fl}[\bar{\alpha}_N] - \bar{\alpha}_N| \leq & \bar{\alpha}_N \left[\left(2^{\bar{m}+1} - 2\right) N + N + 3\right] (1 + N\varepsilon)\varepsilon \\
\leq & \bar{\alpha}_N \left[\left(2^{\bar{m}+1} - 1\right) N + 3 + \frac{N}{10} + 3N\varepsilon\right] \varepsilon \\
\leq & \bar{\alpha}_N 2^{\bar{m}+1} N\varepsilon \\
\leq & 4\bar{\alpha}_N n N\varepsilon
\end{aligned}
$$

which is what we set out to prove. $\qquad \square$

This leads to our main convergence result.

**Theorem 2** (Approximation error direct convolution). *Let the assumptions in Theorem 1 and Lemma 2 hold. Then it holds that*

$$
|\mathrm{fl}[\bar{\alpha}_N] - \alpha| \leq (1 + 4nN\epsilon) \left(\gamma C_1 \overline{C}_2(n, \gamma) N^{-2p} + C_3 \max_{x \in [0, \gamma]} \left|\frac{d^r}{dx^r} f^{*n}(x)\right| N^{-r}\right) + 4\alpha n N\varepsilon,
$$

*for all integers $n$ and $N$ such that $N \leq \varepsilon/C_5$ and $(4n)^2 N\varepsilon < 1/10$.*

The result follows from Lemma 2 and Theorem 1 and using the triangle inequality.

We continue with a lemma needed for the proceeding result. The lemma is a version of [25, Lemma 5] and the proof of our lemma is based on the one given in the cited paper. Note first that the discrete version $\bar{f}$ of $f$ can be associated with a vector $q \in \mathbb{R}^N_{\geq 0}$ by letting $q_i = \bar{f}(x_i)$. Then we let

$$
\|q\|_1 := \sum_{i=1}^N |q_i| \text{ and } \|q\|_\infty := \max_{1 \leq i \leq N} |q_i|.
$$

We also need to define the DFT and IDFT operators, denoted $D_N$ and $D_N^{-1}$ respectively. Let

$$
D_{N,k,j} = e^{\frac{ikj2\pi}{N}}, \text{ and } D_{N,k,j}^{-1} = \frac{1}{N} e^{\frac{-ikj2\pi}{N}}.
$$

We will denote the operators by $D$ and $D^{-1}$ when the dimension is clear from the context. We are then ready to proceed with the lemma.

**Lemma 3.** *Let $q \in \mathbb{R}^N_{\geq 0}$ where the numerical approximation $\mathrm{fl}[q] \in \mathbb{R}^N$ satisfies $|q_i - \mathrm{fl}[q_i]| \leq q_i c_\delta \varepsilon$ for some $c_\delta \in (0,1]$ and let $k = \lceil \log_2(N) \rceil + 1$. Assume that $\mathrm{fl}[q^{\circledast 2}]$ is computed by the method described in (6), that $13k\varepsilon \leq 1$ as well as*

$|\text{fl}[h] - h| = 0$, *i.e. that we are able to accurately represent the constant $h = \frac{\gamma}{N}$ numerically. Furthermore, we also assume that $\text{fl}[D]\text{fl}[q]$ and its square can be calculated exactly with floating point arithmetic, that is $\text{fl}[D]\text{fl}[q] = \text{fl}[\text{fl}[D]\text{fl}[q]]$ and $(\text{fl}[D]\text{fl}[q])^2 = \text{fl}\big[(\text{fl}[D]\text{fl}[q])^2\big]$. Then*

$$\left\| \text{fl}\big[q^{\circledast 2}\big] - q^{\circledast 2} \right\|_\infty \leq 2h(c_\delta + 9k)\varepsilon \|q\|_1^2 + ch\varepsilon^2,$$

*where $c > 0$ is a constant depending on $k$ and $\|q\|_1^2$ capturing higher-order terms of $\varepsilon$.*

*Proof.* Let $q \in \mathbb{R}^{2^k}_{\geq 0}$ be a zero-padded version of our original $q \in \mathbb{R}^N_{\geq 0}$. Note that we for simplicity collect all higher-order terms of $\varepsilon$ in constants $c_j$ throughout this proof. We then have that

$$
\begin{aligned}
\|\text{fl}[D]\text{fl}[q] - Dq\|_\infty &\leq \|D(\text{fl}[q] - q)\|_\infty + \|(D - \text{fl}[D])\text{fl}[q]\|_\infty \\
&\leq \|\text{fl}[q] - q\|_1 + 6k\varepsilon \|\text{fl}[q]\|_1 \\
&\leq c_\delta \varepsilon \|q\|_1 + 6k\varepsilon(1 + c_\delta \varepsilon) \|q\|_1 \\
&\leq (c_\delta + 6k)\varepsilon \|q\|_1 + c_1 \varepsilon^2
\end{aligned}
$$

where, for the first inequality, we have used the triangle inequality and the transition from the first to the second line holds due to the fact that $\|Dx\|_\infty \leq \|x\|_1$, $x \in \mathbb{R}^{2^k}$ and [25, Lemma 4] together with our assumption on $13k\varepsilon$. The jump from the second to the third line comes from our assumptions on the differences $|q_i - \text{fl}[q_i]| \leq q_i c_\delta \varepsilon$ as this implies

$$\|\text{fl}[q] - q\|_1 = \sum_{i=0}^{2^k} |\text{fl}[q(x_i)] - q(x_i)| \leq \sum_{i=0}^{2^k} q(x_i) c_\delta \varepsilon \leq c_\delta \varepsilon \|q\|_1,$$

and further

$$
\begin{aligned}
\|\text{fl}[q]\|_1 = \sum_{i=0}^{2^k} |\text{fl}[q(x_i)]| &\leq \sum_{i=0}^{2^k} (|\text{fl}[q(x_i)] - q(x_i)| + q(x_i)) \\
&\leq \|\text{fl}[q] - q\|_1 + \|q\|_1 = (1 + c_\delta \varepsilon) \|q\|_1.
\end{aligned}
$$

The final transition follows by choosing an appropriate constant $c_1$. Let now $r(x) = [(Dq)(x)]^2$ and similarly $\text{fl}[r(x)] = [(\text{fl}[D]\text{fl}[q])(x)]^2$, where we have used our assumption stating that the multiplication of $\text{fl}[D]$, $\text{fl}[q]$ and the square of their product can be represented exactly in floating point arithmetic. Then, we have

$$
\begin{aligned}
\|r - \text{fl}[r]\|_1 &\leq 2^k \left\| (Dq)^2 - (\text{fl}[D]\text{fl}[q])^2 \right\|_\infty \\
&\leq 2^k \left( \left\| (Dq)^2 - (\text{fl}[D]\text{fl}[q])(Dq) \right\|_\infty + \left\| (\text{fl}[D]\text{fl}[q])(Dq) - (\text{fl}[D]\text{fl}[q])^2 \right\|_\infty \right) \\
&\leq 2^k \left( \|Dq\|_\infty \|Dq - \text{fl}[D]\text{fl}[q]\|_\infty + \|\text{fl}[D]\text{fl}[q]\|_\infty \|Dq - \text{fl}[D]\text{fl}[q]\|_\infty \right) \\
&\leq 2^k \left( \|q\|_1 + \|\text{fl}[D]\text{fl}[q]\|_\infty \right) \|Dq - \text{fl}[D]\text{fl}[q]\|_\infty \\
&\leq 2^k \left( 2\|q\|_1 + \|Dq - \text{fl}[D]\text{fl}[q]\|_\infty \right) \|Dq - \text{fl}[D]\text{fl}[q]\|_\infty \\
&\leq 2^k \left( 2\|q\|_1 + (c_\delta + 6k)\varepsilon \|q\|_1 + c_1 \varepsilon^2 \right) \left( (c_\delta + 6k)\varepsilon \|q\|_1 + c_1 \varepsilon^2 \right) \\
&\leq 2^k \left( 2(c_\delta + 6k)\varepsilon \|q\|_1^2 + (c_\delta + 6k)^2 \varepsilon^2 \|q\|_1^2 \right) + c_2 \varepsilon^2 + c_3 \varepsilon^3 + c_4 \varepsilon^4 \\
&\leq 2^k 2(c_\delta + 6k)\varepsilon \|q\|_1^2 + c_5 \varepsilon^2
\end{aligned}
$$

where we have used the triangle inequality, the bound found above and absorbed the higher-order term in $\varepsilon$ by an appropriate constant $c_5$. We also have the following

inequality

$$\|r\|_1 = \sum_{i=0}^{2^k} [Dq(x_i)]^2 \leq \sum_{i=0}^{2^k} \|Dq\|_\infty^2 = 2^k \|Dq\|_\infty^2 \leq 2^k \|q\|_1^2,$$

which we use in order to show that

$$\|\mathrm{fl}[r]\|_1 \leq \sum_{i=0}^{2^k} \left( |\mathrm{fl}[r(x_i)] - r(x_i)| + r(x_i) \right) \leq \|\mathrm{fl}[r] - r\|_1 + \|r\|_1$$

$$\leq 2^k 2(c_\delta + 6k)\varepsilon \|q\|_1^2 + c_5\varepsilon^2 + 2^k \|q\|_1^2 = 2^k [1 + 2(c_\delta + 6k)\varepsilon] \|q\|_1^2 + c_5\varepsilon^2.$$

Moving on, we have

$$\left\| D^{-1}r - \mathrm{fl}[D^{-1}]\mathrm{fl}[r] \right\|_\infty \leq \left\| D^{-1}(r - \mathrm{fl}[r]) \right\|_\infty + \left\| (D^{-1} - \mathrm{fl}[D^{-1}])\mathrm{fl}[r] \right\|_\infty$$

$$\leq \frac{\|r - \mathrm{fl}[r]\|_1}{2^k} + \frac{6k\varepsilon \|\mathrm{fl}[r]\|_1}{2^k}$$

$$\leq 2(c_\delta + 6k)\varepsilon \|q\|_1^2 + c_6\varepsilon^2$$

$$+ 6k\varepsilon[1 + 2(c_\delta + 6k)\varepsilon] \|q\|_1^2 + c_7\varepsilon^3$$

$$\leq 2(c_\delta + 9k)\varepsilon \|q\|_1^2 + c_8\varepsilon^2$$

where again, in the first inequality, we have used the triangle inequality followed by the fact that $\left\| D^{-1}x \right\|_\infty \leq \frac{\|x\|_1}{N}$, $x \in \mathbb{R}^{2^k}$ and once again we use [25, Lemma 4] to proceed from the first to the second line. The transition to the last line is done by choosing an appropriate constant $c_8$. Finally, we have that

$$\left\| \mathrm{fl}[q^{\circledast 2}] - q^{\circledast 2} \right\|_\infty = \left\| \mathrm{fl}[h] \left( \mathrm{fl}[D^{-1}] \left[ (\mathrm{fl}[D]\mathrm{fl}[q])^2 \right] \right) - h \left( D^{-1} \left[ (Dq)^2 \right] \right) \right\|_\infty$$

$$= \left\| \mathrm{fl}[h] \left( \mathrm{fl}[D^{-1}]\mathrm{fl}[r] \right) - h \left( D^{-1}r \right) \right\|_\infty$$

$$\leq 2h(c_\delta + 9k)\varepsilon \|q\|_1^2 + c_8 h\varepsilon^2,$$

for a suitable constant $c_8$ that depends on $c_\delta, k$ and $\|q\|_1^2$. $\qquad\square$

The next Lemma shows that FFT-based convolution may be more sensitive to rounding errors, since we can only bound its absolute error.

**Lemma 4** (Rounding error FFT-based convolution). *Let $\bar\alpha_N$ be computed by FFT-based convolution with $n = 2^m, m \in \mathbb{N}$, let $\mathrm{fl}[\alpha_N]$ be computed with floating point arithmetic with machine epsilon $\varepsilon > 0$. Assume further that $\mathrm{fl}[f(x_j)] = f(x_j)$ for $j \in \{0, 1, \ldots, N\}$ with $N = 2^r, r \in \mathbb{N}$, and that $2m \leq r$. Then, when disregarding higher-order epsilon terms, we have*

$$\left\| \mathrm{fl}\left[ f^{\circledast 2^m} \right] - f^{\circledast 2^m} \right\|_\infty \leq 18hc \log_2(nN) \log_2(n)\varepsilon \|f\|_1^n$$

*where $c = \max\{1, \gamma\}$.*

*Proof.* Note that by applying recurrently the triangle inequality and by our assumption that $\mathrm{fl}[f] = f$ we have

$$\left\| \mathrm{fl}\left[ f^{\circledast 2^m} \right] - f^{\circledast 2^m} \right\|_\infty \leq \sum_{i=0}^{m-1} \left\| \mathrm{fl}\left[ \mathrm{fl}\left[ f^{\circledast 2^{m-i-1}} \right]^{\circledast 2} \right]^{\circledast 2^i} - \left( \mathrm{fl}\left[ f^{\circledast 2^{m-i-1}} \right]^{\circledast 2} \right)^{\circledast 2^i} \right\|_\infty$$

For the sake of lighter notation we let $g_i = \mathrm{fl}\left[ f^{\circledast 2^i} \right]$, as we then can rewrite the right-hand side of the equation above as

$$\sum_{i=0}^{m-1} \left\| \mathrm{fl}\left[ g_{m-i-1}^{\circledast 2} \right]^{\circledast 2^i} - \left( g_{m-i-1}^{\circledast 2} \right)^{\circledast 2^i} \right\|_\infty. \tag{13}$$

We now need to find an expression for each term in the sum above. First, for readability we introduce the shorthand notation $\bar{g} := g_{m-i-1}^{\circledast 2}$ for some arbitrary value of $i \in \{1, \ldots, m-1\}$ and let $\bar{g}$ be zero-padded such that $\bar{g} \in \mathbb{R}^{2^{r+i}}$. Note then that

$$\|D\mathrm{fl}[\bar{g}]\|_\infty = \|D\mathrm{fl}[\bar{g}] - D\bar{g} + D\bar{g}\|_\infty \leq \|D\mathrm{fl}[\bar{g}] - D\bar{g}\|_\infty + \|D\bar{g}\|_\infty$$
$$\leq \|\mathrm{fl}[\bar{g}] - \bar{g}\|_1 + \|\bar{g}\|_1 . \quad (14)$$

Furthermore, we have that

$$\left\|(D\mathrm{fl}[\bar{g}])^{2^i} - (D\bar{g})^{2^i}\right\|_\infty \leq \left\|(D\mathrm{fl}[\bar{g}])^{2^i} - (D\mathrm{fl}[\bar{g}])^{2^i-1}(D\bar{g})\right\|_\infty$$
$$+ \left\|(D\mathrm{fl}[\bar{g}])^{2^i-1}(D\bar{g}) - (D\mathrm{fl}[\bar{g}])^{2^i-2}(D\bar{g})^2\right\|_\infty$$
$$+ \cdots + \left\|(D\mathrm{fl}[\bar{g}])(D\bar{g})^{2^i-1} - (D\bar{g})^{2^i}\right\|_\infty$$
$$\leq \left( \left\|(D\mathrm{fl}[\bar{g}])^{2^i-1}\right\|_\infty + \left\|(D\mathrm{fl}[\bar{g}])^{2^i-2}(D\bar{g})\right\|_\infty \right.$$
$$+ \cdots + \left\|(D\mathrm{fl}[\bar{g}])(D\bar{g})^{2^i-2}\right\|_\infty$$
$$\left. + \left\|(D\bar{g})^{2^i-1}\right\|_\infty \right) \|D(\mathrm{fl}[\bar{g}] - \bar{g})\|_\infty$$
$$\leq \left( \|D\mathrm{fl}[\bar{g}]\|_\infty^{2^i-1} + \|D\mathrm{fl}[\bar{g}]\|_\infty^{2^i-2} \|D\bar{g}\|_\infty \right.$$
$$\left. + \cdots + \|D\mathrm{fl}[\bar{g}]\|_\infty \|D\bar{g}\|_\infty^{2^i-2} + \|D\bar{g}\|_\infty^{2^i-1} \right) \|\mathrm{fl}[\bar{g}] - \bar{g}\|_1 .$$

Then from (14) and Lemma 3 we have that

$$\|D\mathrm{fl}[\bar{g}]\|_\infty^s \leq (\|\mathrm{fl}[\bar{g}] - \bar{g}\|_1 + \|\bar{g}\|_1)^s \leq \left( 18h(r+i+1)\varepsilon \|g_{m-i-1}\|_1^2 + ch\varepsilon^2 + \|\bar{g}\|_1 \right)^s .$$

Thus, we end up with $\|D\mathrm{fl}[\bar{g}]\|_\infty^s \leq \|\bar{g}\|_1^s + c\varepsilon$, where $c$ captures all terms multiplied with $\varepsilon$. Then we can write

$$\left\|(D\mathrm{fl}[\bar{g}])^{2^i} - (D\bar{g})^{2^i}\right\|_\infty \leq \left( \|\bar{g}\|_1^{2^i-1} + \|\bar{g}\|_1^{2^i-1} + \cdots \right.$$
$$\left. + \|\bar{g}\|_1^{2^i-1} + \|\bar{g}\|_1^{2^i-1} + c\varepsilon \right) 2^i \|\mathrm{fl}[\bar{g}] - \bar{g}\|_\infty$$
$$\leq 2^i \|\bar{g}\|_1^{2^i-1} 2^i 18h(r+i+1)\varepsilon \|g_{m-i-1}\|_1^2 + c\varepsilon^2$$
$$\leq 18h2^{2i} \|g_{m-i-1}\|_1^{2^{i+1}-2} (r+i+1)\varepsilon \|g_{m-i-1}\|_1^2 + c\varepsilon^2$$
$$\leq 18\gamma(r+i+1)\varepsilon \|g_{m-i-1}\|_1^{2^{i+1}} + c\varepsilon^2,$$

where $c$ is a constant capturing the higher-order terms in $\varepsilon$ that we adjust appropriately from line to line and from our assumption that $2m \leq r$ we have $2^{2i}h \leq \gamma$ as well as the fact that $\left\|g_{m-i-1}^{\circledast 2}\right\|_1 \leq \|g_{m-i-1}\|_1^2$. We can then show that

$$\left\|\mathrm{fl}[\bar{g}]^{\circledast 2^i} - \bar{g}^{\circledast 2^i}\right\|_\infty = \left\|hD^{-1}\left[(D\mathrm{fl}[\bar{g}])^{2^i} - (D\bar{g})^{2^i}\right]\right\|_\infty$$
$$\leq h \frac{\left\|(D\mathrm{fl}[\bar{g}])^{2^i} - (D\bar{g})^{2^i}\right\|_1}{2^{i+r}}$$
$$\leq h \left\|(D\mathrm{fl}[\bar{g}])^{2^i} - (D\bar{g})^{2^i}\right\|_\infty$$
$$\leq 18h\gamma(r+i+1)\varepsilon \|g_{m-i-1}\|_1^{2^{i+1}} + c\varepsilon^2 .$$

We can now apply this inequality in order to get bounds on the terms in the sum in (13). Consider first the term in (13) where $i = 0$, by applying Lemma 3 we have

$$\left\| \mathrm{fl}\big[g_{m-1}^{\circledast 2}\big] - \big(g_{m-1}^{\circledast 2}\big) \right\|_\infty \leq 18h(r+1)\varepsilon \left\| g_{m-1} \right\|_1^2.$$

Then, moving on to the case $i \geq 1$ we have from the bound above that

$$\left\| \mathrm{fl}\big[g_{m-i-1}^{\circledast 2}\big]^{\circledast 2^i} - \big(g_{m-i-1}^{\circledast 2}\big)^{\circledast 2^i} \right\|_\infty \leq 18h\gamma(r+i+1)\varepsilon \left\| g_{m-i-1} \right\|_1^{2^{i+1}} + c\varepsilon^2.$$

By inserting the above estimates in (13) we achieve the estimate

$$
\begin{aligned}
\left\| \mathrm{fl}\big[f^{\circledast 2^m}\big] - f^{\circledast 2^m} \right\|_\infty \leq{}& 18h(r+1)\varepsilon \left\| g_{m-1} \right\|_1^2 \\
&+ \sum_{i=1}^{m-1} 18h\gamma(r+i+1)\varepsilon \left\| g_{m-i-1} \right\|_1^{2^{i+1}} + c\varepsilon^2 \\
\leq{}& 18h(r+1)\varepsilon \left\| g_{m-1} \right\|_1^2 \\
&+ 18h\gamma(r+m)\varepsilon \sum_{i=1}^{m-1} \left\| g_{m-i-1} \right\|_1^{2^{i+1}} + c\varepsilon^2.
\end{aligned}
$$

Thus, when only considering the leading term in $\varepsilon$ we get

$$
\begin{aligned}
\left\| \mathrm{fl}\big[f^{\circledast 2^m}\big] - f^{\circledast 2^m} \right\|_\infty \leq 18h\log_2(nN)\varepsilon \Bigg(& \left\| \mathrm{fl}\big[f^{\circledast 2^{m-1}}\big] \right\|_1^2 \\
&+ \gamma \sum_{i=1}^{m-1} \left\| \mathrm{fl}\big[f^{\circledast 2^{m-i-1}}\big] \right\|_1^{2^{i+1}} \Bigg).
\end{aligned}
$$

By now recursively applying this relation on the norms on the left hand side and ignoring higher-order terms in $\varepsilon$ we are able to rewrite the equation above as

$$
\begin{aligned}
\left\| \mathrm{fl}\big[f^{\circledast 2^m}\big] - f^{\circledast 2^m} \right\|_\infty \leq{}& 18h\log_2(nN)\varepsilon \left( \left\| f \right\|_1^{2^m} + \gamma \sum_{i=1}^{m-1} \left\| f \right\|_1^{2^m} \right) \\
\leq{}& 18hc\log_2(nN)\log_2(n)\varepsilon \left\| f \right\|_1^n,
\end{aligned}
$$

where $c = \max\{\gamma, 1\}$. $\qquad \square$

We are then ready to prove the following theorem, giving a bound on the error of performing convolution using FFT.

**Theorem 3** (Approximation error FFT-based convolution)**.** *Let the assumptions in Theorem 1 and Lemma 4 hold. Then it holds that*

$$
\begin{aligned}
|\mathrm{fl}[\bar{\alpha}_N] - \alpha| \leq{}& \gamma C_1 \overline{C}_2(n, \gamma) N^{-2p} + C_3 \max_{x \in [0, \gamma]} \left| \frac{d^r}{dx^r} f^{*n}(x) \right| N^{-r} \\
&+ 18hc \log_2(nN) \log_2(n)\varepsilon \left\| f \right\|_1^n).
\end{aligned}
$$

*Proof.* The result follows directly from Theorem 1 and Lemma 4. $\qquad \square$

3.2. **Computational cost.** In this section we compare the computational cost and accuracy of direct-based convolution against FFT-based convolution as a function of the numerical resolution $N$ and the number of RVs $n$. We restrict ourselves to settings where the lemmas and theorems in Section 3 apply.

**Theorem 4.** *The computational cost of computing $\bar{\alpha}_N$, counted in the number of floating point operations, is*

$$COST(\bar{\alpha}_N) = \begin{cases} \mathcal{O}(\log_2(n)N^2) & \textit{when using direct convolution} \\ \mathcal{O}(\log_2(n)N\log_2(N)) & \textit{when using FFT-based convolution.} \end{cases}$$

*Proof.* Recall that $m$ denotes the largest integer such that $m \leq \log_2(n)$. For each $\ell = 1, \ldots, m$ and $k = 0, 1, \ldots, N$, the computation $\bar{f}^{\circledast 2^{\ell}}(x_k) = \bar{f}^{\circledast 2^{\ell-1}} \circledast \bar{f}^{\circledast 2^{\ell-1}}(x_k)$ costs $\mathcal{O}(N)$. The cost of computing $\bar{f}^{\circledast n}$ thus becomes $\mathcal{O}(mN^2)$ and computing $\bar{\alpha}$ adds an additional (relatively speaking, negligible) cost of $\mathcal{O}(N)$. The upper bound in cost for FFT-based convolution follows by a similar argument. $\square$

When disregarding factors of $\log_2(N)$ in the cost estimate and also disregarding rounding errors, we obtain the following relation between cost and absolute approximation error:

$$|\bar{\alpha}_N - \alpha| \leq \begin{cases} \frac{\widehat{C}}{(\text{COST})^{r/2}} & \text{when using direct convolution} \\ \frac{\widehat{C}}{(\text{COST})^r} & \text{when using FFT-based convolution.} \end{cases} \tag{15}$$

Supposing further that there exists a constant $C > 0$ such that

$$\frac{(1 + 4nN\epsilon)\left(\gamma C_1 \overline{C}_2(n,\gamma)N^{r-2p} + C_3 \max_{x \in [0,\gamma]}\left|\frac{d^r}{dx^r}f^{*n}(x)\right|\right)}{\alpha} \leq C$$

holds for all relevant $\gamma, n, N$ and $\varepsilon$, we obtain the following error estimate for the relative error of approximating $\alpha$:

$$\frac{|\bar{\alpha}_N - \alpha|}{\alpha} \leq \begin{cases} CN^{-r} + nN\varepsilon & \text{for direct convolution} \\ CN^{-r} + C_6\frac{hc\log_2(nN)\log_2(n)\|f\|_1^n\varepsilon}{\alpha} & \text{for FFT-based convolution.} \end{cases} \tag{16}$$

We note that when $\alpha \ll 1$, the result indicates that for a given resolution $N$, the relative error may be substantially smaller for direct-based convolution than for FFT-based convolution, precisely as we observe in the numerical examples in Section 4.

## 4. Numerical experiments

To verify numerically that the proposed method produces satisfactory results and to confirm that the theoretical error rate identified in the previous section holds in practice, we conducted a series of experiments. First, in Section 4.1 we compare the FFT implementation of the convolution method with the direct method in terms of how well they are able to approximate the rare-event probability of a sum of RVs. As the results from the first experiment shows that the direct method gives low rounding errors, we run the rest of the experiments using the direct method only. In Section 4.2, we look at how well the convolution method estimate the CDF for the sum of RVs for which the distribution of the sum is indeed known. Then, in Section 4.3, we consider the Log-Normal distribution with two goals in mind: 1) We want to explore the convergence properties of the convolution method and check if we empirically are able to observe the theoretical convergence rate as given by Theorem 1, 2) We compare the calculated estimates of the CDF with approximations calculated using an alternative method, in this instance a saddlepoint method presented in [9]. Then, in the last Section 4.4 we look at how the convolution method performs when approximating the CDF for the sum of RVs for other distributions where the distribution of the sum is not known.
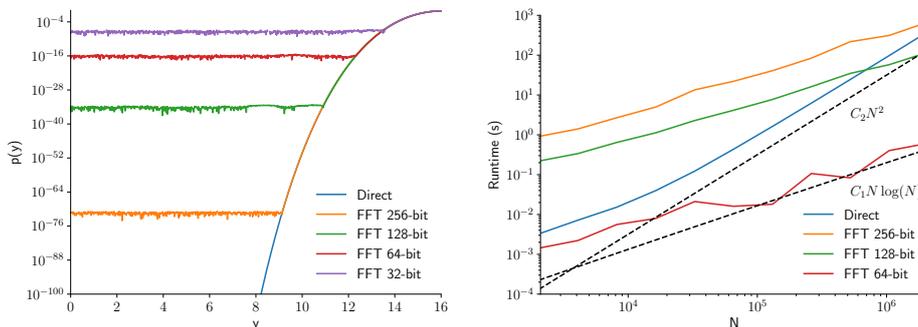
FIGURE 1. Left: The probability density function $p(y) = \bar{f}^{\circledast 16}(y)$ for direct convolution and FFT-based convolution for the rare-event problem studied in Section 4.1. Right: The runtime for direct convolution and FFT-based convolution for the rare-event problem studied in Section 4.1.

### 4.1. Comparison of direct- and FFT-based convolution.

In this section, we compare the performance of direct convolution and FFT-based convolution for left-tail rare-event estimation. In agreement with the theoretical results in Section 3.1, we show that FFT-based convolution is more sensitive to rounding errors than direct convolution in two problem settings where $\alpha \ll 1$.

4.1.1. *Log-Normal distribution.* We estimate the probability of $Y = \sum_{i=1}^{16} X_i \leq \gamma$, where $X_i$ for $i = 1, 2, \cdots, 16$ are i.i.d Log-Normal$(0, 1/64)$ with density denoted by $f$. The large variation in magnitude for the density of $Y$ is illustrated in the left plot of Figure 1, where we numerically have computed $p(y) := \bar{f}^{\circledast 16}(y)$ over the interval $y \in [8, 16]$ using $N = 10^6$ quadrature points. The density is computed by direct convolution with Matlab's **conv()** function and 64-bit floating point precision, and by the FFT-based method for a range of different floating point precisions, using the multiple precision toolbox [1]. We observe that the higher the precision, the better the FFT-based convolution approximates the direct convolution's density, and that FFT introduces an approximation error that is proportional to the machine epsilon. This is consistent with the observations in [38]. For reference, we note that the machine epsilon is approximately $1.19 \times 10^{-7}$ for 32-bit floats, $2.22 \times 10^{-16}$ for 64-bit floats, $1.93 \times 10^{-34}$ for 128-bit floats, and $1.81 \times 10^{-71}$ for 256-bit floats. Table 2 presents the relative error $|\alpha - \text{fl}[\bar{\alpha}_N]|/\alpha$ for different values of $\gamma$ using 64-bit precision direct convolution and FFT-based convolution for a range of precisions. All methods use $N = 10^6$ quadrature points. The pseudo-reference solution is computed using $N = 2^{21}$ quadrature points with 512-bit precision FFT-based convolution. The FFT-based convolution only approximates the rare event well when $\alpha$ is much larger than the machine epsilon for the floating point precision.

Our numerical studies indicate that FFT-based convolution is much more sensitve to rounding errors than direct convolution, but it may still be useful in computations when the number of quadrature points $N$ is large since the method has a lower asymptotic computational cost than direct convolution, cf. Theorem 4. The right plot in Figure 1 measures the computational cost of the two convolution methods in runtime, displays cost rates that are consistent with Theorem 4, and shows that for all considered floating point precisions, FFT-based convolution will eventually, for sufficiently large $N$, outperform direct convolution in terms of runtime.

| $\gamma$ | Ref. sol. CDF | Dir. conv. | Saddlp. | FFT 32-bit | FFT 64-bit | FFT 128-bit | FFT 256-bit |
|---|---|---|---|---|---|---|---|
| 8.8 | $2.05\times10^{-83}$ | $4.99\times10^{-13}$ | $4.90\times10^{-06}$ | $4.92\times10^{+75}$ | $1.67\times10^{+67}$ | $3.58\times10^{+49}$ | $1.48\times10^{+11}$ |
| 9.6 | $1.02\times10^{-61}$ | $5.77\times10^{-13}$ | $5.46\times10^{-06}$ | $5.34\times10^{+53}$ | $1.27\times10^{+45}$ | $1.25\times10^{+28}$ | $1.01\times10^{-10}$ |
| 10.4 | $1.04\times10^{-44}$ | $5.47\times10^{-13}$ | $5.61\times10^{-06}$ | $1.66\times10^{+37}$ | $2.23\times10^{+28}$ | $1.02\times10^{+11}$ | $2.24\times10^{-28}$ |
| 11.2 | $1.76\times10^{-31}$ | $6.00\times10^{-13}$ | $5.24\times10^{-06}$ | $1.74\times10^{+24}$ | $2.94\times10^{+15}$ | $1.56\times10^{-03}$ | $1.46\times10^{-40}$ |
| 12 | $2.45\times10^{-21}$ | $5.81\times10^{-13}$ | $4.27\times10^{-06}$ | $4.83\times10^{+13}$ | $3.20\times10^{+04}$ | $1.91\times10^{-13}$ | $1.16\times10^{-50}$ |
| 12.8 | $9.81\times10^{-14}$ | $5.74\times10^{-13}$ | $2.72\times10^{-06}$ | $1.96\times10^{+05}$ | $2.37\times10^{-03}$ | $1.87\times10^{-20}$ | $3.80\times10^{-59}$ |
| 13.6 | $3.03\times10^{-08}$ | $5.97\times10^{-13}$ | $6.88\times10^{-07}$ | $5.38\times10^{+00}$ | $1.23\times10^{-08}$ | $3.93\times10^{-26}$ | $1.79\times10^{-63}$ |
| 14.4 | $1.63\times10^{-04}$ | $6.02\times10^{-13}$ | $1.66\times10^{-06}$ | $3.01\times10^{-03}$ | $4.37\times10^{-13}$ | $1.09\times10^{-29}$ | $5.45\times10^{-68}$ |

TABLE 2. Comparison of the relative error $|\mathrm{fl}[\bar{\alpha}_N] - \alpha|/\alpha$ for sums of Log-Normal-distributed RVs using 64-bit direct convolution, the saddlepoint method (computed with 512-bit floating bit precision) and FFT-based convolution computed with four different floating point precisions.

4.1.2. *Lévy distribution.* We next estimate the probability of $Y = \sum_{i=1}^{16} X_i \leq \gamma$, with $X_i$, $i = 1, 2, \cdots, 16$ are i.i.d. Lévy$(0, 0.1)$ whose density is denoted by $f$ (given in Table 4). This is a stable distribution, which is very suitable for validation of the numerical methods since its density and CDF are known: $Y \sim$ Lévy$(0, 25.6)$ and $\alpha = \mathbb{P}(Y \leq \gamma) = \mathrm{erfc}(\sqrt{12.8/\gamma})$. The large variation in magnitude for the density of $Y$ is illustrated in Figure 2, where we compare the exact density of $Y$ to numerical approximations using $N = 10^6$ quadrature points. The density is
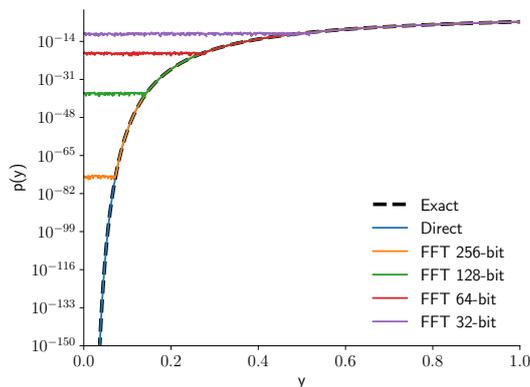


FIGURE 2. The probability density function to $Y = \sum_{i=1}^{16} X_i$ where $X_i$ are i.i.d. Lévy distributed RVs. The density is computed by the exact formula, by direct convolution and FFT-based convolution.

computed numerically by direct convolution using 64-bit precision and FFT-based convolution for a range of floating point precisions. Remarkably, direct convolution agrees well with the exact density over the full range of values displayed, while FFT-based convolution agrees well only when the precision is sufficiently high.

Table 3 presents the relative error $|\alpha - \mathrm{fl}[\bar{\alpha}_N]|/\alpha$ for different values of $\alpha$ using 64-bit precision direct convolution and FFT-based convolution for a range of precisions. All numerical methods use $N = 10^6$ quadrature points. The reference solution is computed by the exact CDF, $\alpha = \mathrm{erfc}(\sqrt{12.8/\gamma})$ with 512-bit floating point precision. We again observe that FFT-based convolution only approximates the rare event well when $\alpha$ is much larger than the floating point precision machine epsilon.

| $\gamma$ | $\alpha = \mathbb{P}(Y \leq \gamma)$ | Dir. conv. | FFT 32-bit | FFT 64-bit | FFT 128-bit | FFT 256-bit |
|------|-------------------------|-------------------------|-------------------------|-------------------------|-------------------------|-------------------------|
| 0.05 | $2.33 \times 10^{-113}$ | $6.74 \times 10^{-13}$ | $7.82 \times 10^{+68}$ | $6.12 \times 10^{+50}$ | $2.16 \times 10^{+33}$ | $1.58 \times 10^{-04}$ |
| 0.1  | $1.28 \times 10^{-57}$  | $6.78 \times 10^{-13}$ | $7.45 \times 10^{+25}$ | $6.05 \times 10^{+15}$ | $3.58 \times 10^{-02}$ | $7.47 \times 10^{-27}$ |
| 0.2  | $1.12 \times 10^{-29}$  | $6.05 \times 10^{-13}$ | $1.15 \times 10^{+09}$ | $2.61 \times 10^{-02}$ | $3.45 \times 10^{-18}$ | $9.20 \times 10^{-29}$ |
| 0.5  | $8.34 \times 10^{-13}$  | $4.24 \times 10^{-13}$ | $2.68 \times 10^{-02}$ | $1.51 \times 10^{-10}$ | $6.18 \times 10^{-28}$ | $1.54 \times 10^{-31}$ |
| 1    | $4.20 \times 10^{-07}$  | $2.80 \times 10^{-13}$ | $4.95 \times 10^{-05}$ | $1.24 \times 10^{-13}$ | $2.22 \times 10^{-32}$ | $1.11 \times 10^{-34}$ |

TABLE 3. Comparison of the relative error $|\text{fl}[\bar{\alpha}_N] - \alpha|/\alpha$ for sums of Lévy-distributed RVs using 64-bit direct convolution and FFT-based convolution computed with four different floating point precisions.

| Distribution | Parameters | PDF |
|--------------|------------|-----|
| Chi-squared ($\chi^2$) | $df \in \mathbb{N}$ | $\frac{1}{2^{df/2}\Gamma(df/2)} x^{df/2-1} e^{-x/2}$ |
| Lévy | $c > 0$ | $\sqrt{\frac{c}{2\pi}} \frac{e^{-\frac{c}{2x}}}{x^{3/2}}$ |

TABLE 4. Probability density function for the Chi-squared and Lévy distributions

4.2. **Estimating known distributions.** There exist several probability distributions for which the distribution of the sum $Y = \sum_{i=1}^{n} X_i$ is known, given that the RVs $X_i$ are all independent. If for example $X_i$ is Chi-squared distributed with $r_i$ degrees of freedom $(X_i \sim \chi^2(r_i))$ for $i \in \{1, 2, \ldots, n\}$ we have that

$$Y \sim \chi^2 \left( \sum_{i=1}^{n} r_i \right).$$

Thus, to check empirically that the algorithm presented in this paper indeed is able to accurately calculate the PDF and the CDF of the sum of RVs we check against distributions where the resulting distribution is known. We chose to do numerical experiments with the Chi-squared ($\chi^2$) and Lévy distributions. The PDFs are given in Table 4. In the case of the Lévy distribution if we let $X_i \sim \text{Lévy}(\mu_i, c_i)$ with $\mu_i \in (-\infty, \infty), c_i > 0$ for $i \in \{1, 2, \ldots, n\}$ we have that

$$Y \sim \text{Lévy} \left( \sum_{i=1}^{n} \mu_i, \left( \sum_{i=1}^{n} \sqrt{c_i} \right)^2 \right).$$

For these experiments we are estimating the value $\alpha = F_Y(\gamma), \gamma = xn$ with $x = 0.05$ and $n = 16$, i.e.

$$\alpha = F_Y(0.8) = \mathbb{P}(Y \leq 0.8) = P \left( \sum_{i=1}^{16} X_i \leq 0.8 \right)$$

where $F_Y$ is the CDF of $Y = \sum_{i=1}^{16} X_i$ with the RVs $X_i, i \in \{1, 2, \ldots, 16\}$ all independent. The estimates $\bar{\alpha}$ are calculated using equation (5) with Boole's rule as the closed Newton-Cotes formula in the last step (see Table 5 for an overview of the weights). In the left plot of Figure 3 we display the relative error $\delta = \frac{|\alpha - \bar{\alpha}|}{\alpha}$ as a function of the mesh size $N$ when estimating $\alpha$ with $X_i \sim \chi^2(df)$ for a number of different parameter values $df$. Note that the legend also display the value of $\alpha$, showing that we indeed are estimating rare events. It is apparent from the figure that there is a large difference in the convergence rate depending on the value of $df$, with the convolution method converging faster to the correct value of $\alpha$ when the value of $df$ increase. It is straight forward to check that $f_3'(x) \xrightarrow{x \to 0} \infty, f_2'(0) = c_1$, $f_4'(0) = c_2$ and $f_6'(0) = 0$ for some constants $c_1, c_2 \in \mathbb{R}$ and where $f_{df}$ is the PDF
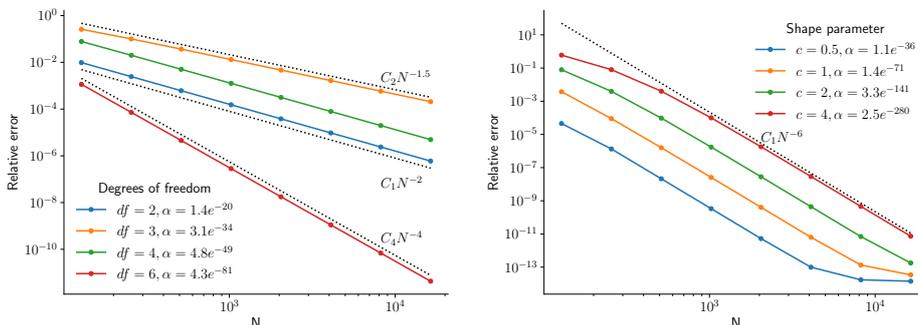
FIGURE 3. Left: Relative error as a function of the mesh-size when estimating $F_Y(0.8), Y \sim \chi^2(16df)$. Right: Relative error as a function of the mesh-size when we are estimating $F_Y(0.8), Y \sim$ Lévy $\left(0, \left(\sum_{i=1}^{n} \sqrt{c_i}\right)^2\right)$

for the $\chi^2$-distribution with parameter value $df$. Note also that $f_6''(0) = c_3$ for some $c_3 \in \mathbb{R}$. The observed convergence rates are therefore coherent with the theory, with the notable exception of the case $df = 3$. We suspect that the bad convergence rate is a result of $f_3'(x) \xrightarrow{x \to 0} \infty$. Note also that $f_2(0) \neq 0$ thus our implementation of the convolution method used for this experiment utilize the formula given in 9.

The results of a similar experiment where $X_i \sim$ Lévy$(0, c)$ for a number of different values for the shape parameter $c$ is shown in the right plot in Figure 3. From the figure we see that the convolution method performs really well when estimating the Lévy distribution, with the relative error $\delta$ being less than $10^{-9}$ for all values of $c$ that we tested when the mesh-size $N > 10^4$. Furthermore, we observe from the figure that a convergence rate of about $N^{-6}$ were attained for all choices of $c$. In contrast with the $\chi^2$-distribution we have $f^{(k)}(0) = 0$ for all $k \in \mathbb{N}$ for our choices of $c$. The result in Theorem 1 therefore implies that we should observe a convergence rate of at least $N^{-6}$ as we use Boole's rule for the integration in the last step. This fits well with the observed convergence rate for all choices of $c$. Furthermore, we observe that the relative error flattens out when we reach errors bellow $10^{-13}$, which is probably due to round-off errors and in agreement with with the errors observed in Section 4.1.

4.3. **Convergence properties for the convolution method.** In this subsection we again consider the problem of estimating the value

$$\alpha = F_Y(\gamma)$$

where $\gamma = xn$ with $x = 0.7$ and $n = 16$, i.e.

$$\alpha = F_Y(11.2) = \mathbb{P}(Y \leq 11.2) = \mathbb{P}\left(\sum_{i=1}^{16} X_i \leq 11.2\right),$$

where $F_Y$ is the CDF of $Y = \sum_{i=1}^{16} X_i$ and $X_i \sim$ Log-Normal$(0, \sigma^2), i \in \{1, \dots, 16\}$ are i.i.d. Here, we let $\sigma = 0.125$. The PDF is given in Table 1. The approximations $\bar{\alpha}$ are calculated using equation (5). In this experiment we aim to observe how the convergence rate vary when utilize three different closed Newton-Cotes formulas in the last step, Trapezoid, Simpson and Boole. The weights $w_j$ in equation (5) depend on the used Newton-Cotes formula. The resulting formula for each of the Newton-Cotes formulas utilized in this experiment is given in Table 5. We also compare our estimates of the CDF value $\alpha$ and the PDF value $f(11.2)$ with approximations

| Newton-Cotes formula | Formula |
|---|---|
| Trapezoid | $h\left(\frac{1}{2}(\bar{f}^{\circledast n}(x_0) + \bar{f}^{\circledast n}(x_N)) + \sum_{j\in\{1,2,\dots,N-1\}} \bar{f}^{\circledast n}(x_j)\right)$ |
| Simpson | $\frac{h}{3}\left(\bar{f}^{\circledast n}(x_0) + \bar{f}^{\circledast n}(x_N) + 4\sum_{j\in\{1,3,\dots,N-1\}} \bar{f}^{\circledast n}(x_j)\right.$ |
|  | $\left.+ 2\sum_{j\in\{2,4,\dots,N-2\}} \bar{f}^{\circledast n}(x_j)\right)$ |
| Boole's | $\frac{2h}{45}\left(7\big(\bar{f}^{\circledast n}(x_0) + \bar{f}^{\circledast n}(x_N)\big)\right.$ |
|  | $+ 32\sum_{j\in\{1,3,\dots,N-1\}} \bar{f}^{\circledast n}(x_j)$ |
|  | $+ 12\sum_{j\in\{2,6,\dots,N-2\}} \bar{f}^{\circledast n}(x_j)$ |
|  | $\left.+ 14\sum_{j\in\{4,8,\dots,N-4\}} \bar{f}^{\circledast n}(x_j)\right)$ |

TABLE 5. Explicit Newton-Cotes formulas

generated by a saddlepoint method presented in [9]. For details on the saddlepoint method we refer the reader to the cited paper. Furthermore, the estimate of $f(11.2)$ is numerically found by the value $f^{\circledast n}(x_N)$.

In order to test the convergence rate using the different rules in the last step we apply an iterative scheme where we first calculate a pseudo-reference solutions $\bar{\alpha}_{N_M}$ using Boole's rule on a mesh of size $N_M$ for some large $N_M$. We then calculate estimates $\bar{\alpha}_{N_m}$ using the three different rules on a mesh of size $N_m$, where $N_m \ll N_M$, and calculate the relative error

$$\bar{\delta} = \frac{|\bar{\alpha}_{N_M} - \bar{\alpha}_{N_m}|}{\alpha_{N_M}}$$

for each of the estimates. We then check if all of the calculated relative errors are smaller than some threshold $\epsilon > 0$. If this is not the case we repeat the calculation on a mesh of size $N'_m = 2N_m$. This is repeated until all calculated errors are bellow the given threshold. Here we set $N_M = 2^{17}, N_m = 2^{10}$ and $\epsilon = 10^{-8}$.

The results are shown in the left plot of Figure 4 together with reference slopes showing the theoretical convergence rates and the relative error of the above mentioned saddlepoint method. Given the fact that for the PDF $f$ of the Log-Normal distribution we have $f^{(k)}(0) = 0$ for all $k \in \mathbb{N}$, we would from Theorem 1 expect convergence rates similar to the convergence rate of the chosen Newton-Cotes formula. Our results are in accordance with this expectation as the convergence rates are $2, 4$ and $6$ for the Trapezoid, Simpson and Boole's respectively. We also see that we quickly achieve relative errors smaller than the saddlepoint method with all three rules.

In order to further validate the correctness of the convolution method when applied to the Log-Normal distribution, we compared estimates of the CDF and PDF for different values of $x$ generated with our method with approximations calculated by the saddlepoint method from [9]. For the convolution method we utilized a mesh-size of $10^4$ and Boole's rule in the last step, while the saddelpoint approximations were calculated using our implementation of the saddlepoint method presented in [9]. The results are given in 6. Note that the values calculated with our implementation of the saddlepoint method gives slightly different values compared with the values listed in the paper [9] ($\pm 1$ in the third digit of the significand).

The results show that the two methods gives similar approximations of $\alpha$ for all chosen values of $x$, with only the fourth non-zero decimal being different between the two estimates in some cases. This indicates that the convolution method is indeed able to accurately estimate the desired probabilities. Both our implementation of the convolution method and the saddlepoint method has for this experiment a negligible CPU-time. The advantage with the convolution method as opposed to
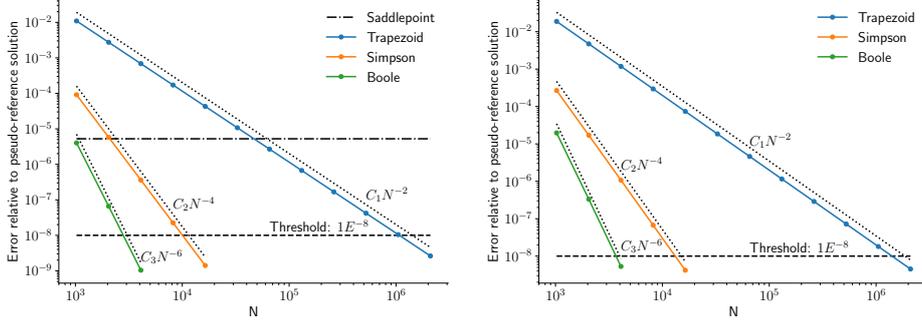
FIGURE 4. Left: Relative error as a function of the mesh-size when approximating $\alpha_{N_M}$ in the case when $X_i \sim$ Log-Normal$(0, 0.125)$ where ($\alpha_{N_M}$ is a pseudo-reference solution calculated using the convolution method with $N_M = 1e6$ using Boole's rule in the last step. Right: Relative error as a function of the mesh-size when estimating pseudo-reference solution $\alpha_{N_M}$ of the CDF of a sum of Log-Normals with varying $\sigma$ calculated using the convolution method with $N_M = 1e6$ using Boole's rule in the last step

| x | Convolution CDF | Saddle CDF | Convolution PDF | Saddle PDF |
|------|------------------|------------------|------------------|------------------|
| 0.70 | $1.761 \times 10^{-31}$ | $1.761 \times 10^{-31}$ | $5.873 \times 10^{-30}$ | $5.873 \times 10^{-30}$ |
| 0.80 | $9.806 \times 10^{-14}$ | $9.806 \times 10^{-14}$ | $1.829 \times 10^{-12}$ | $1.829 \times 10^{-12}$ |
| 0.85 | $3.031 \times 10^{-8}$ | $3.031 \times 10^{-8}$ | $3.975 \times 10^{-7}$ | $3.975 \times 10^{-7}$ |
| 0.90 | $1.631 \times 10^{-4}$ | $1.631 \times 10^{-4}$ | $1.388 \times 10^{-3}$ | $1.388 \times 10^{-3}$ |
| 0.91 | $5.955 \times 10^{-4}$ | $5.955 \times 10^{-4}$ | $4.577 \times 10^{-3}$ | $4.577 \times 10^{-3}$ |
| 0.92 | $1.911 \times 10^{-3}$ | $1.911 \times 10^{-3}$ | $1.318 \times 10^{-2}$ | $1.318 \times 10^{-2}$ |
| 0.93 | $5.423 \times 10^{-3}$ | $5.423 \times 10^{-3}$ | $3.332 \times 10^{-2}$ | $3.332 \times 10^{-2}$ |
| 0.94 | $1.368 \times 10^{-2}$ | $1.368 \times 10^{-2}$ | $7.416 \times 10^{-2}$ | $7.416 \times 10^{-2}$ |
| 0.95 | $3.081 \times 10^{-2}$ | $3.081 \times 10^{-2}$ | $1.460 \times 10^{-1}$ | $1.460 \times 10^{-1}$ |
| 0.98 | $1.901 \times 10^{-1}$ | $1.901 \times 10^{-1}$ | $5.520 \times 10^{-1}$ | $5.520 \times 10^{-1}$ |

TABLE 6. Approximations of the CDF and PDF of Y

the more intricate saddlepoint method is that the convolution method is generic in the way that it can handle multiple distributions and that the RVs does not need to be identically distributed. However, for the convolution method when a larger mesh size $N$ is needed to get accurate estimates the cost increases by $\mathcal{O}(N^2)$ and when the number $n$ of RVs in the sum increases the complexity increase by $\mathcal{O}(\log_2(n))$. On the other hand, the saddlepoint method has a negligible computational cost.

Furthermore, we wanted to empirically test the performance of the convolution method when employing it on a sum of independent Log-Normals that are not identically distributed. We therefore perform a experiment similar to the ones described in above, but instead estimate

$$\alpha = F_Y(11.2), \text{ with } Y = \sum_{i=1}^{16} X_i,$$

where $X_i \sim$ Log-Normal$(0, \sigma_i)$ and $\sigma_i = \frac{1}{2^{2+j}}$ with $j = i \mod 4$ for $i \in \{1, 2, \ldots, 16\}$. The result is shown in the right plot of Figure 4. From the graphs it is apparent that the convolution method performs well in this case as well.
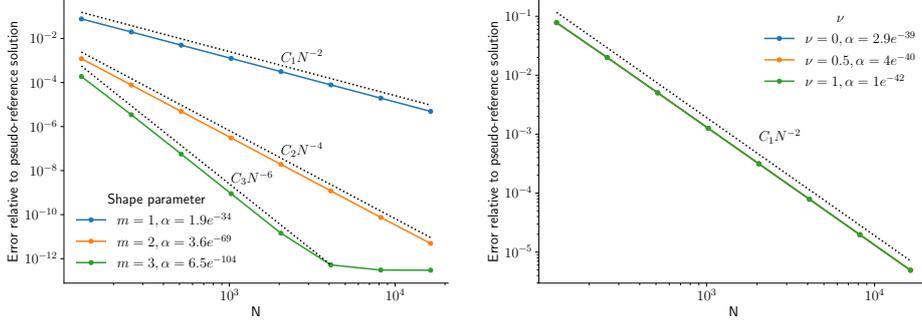
FIGURE 5. Left: Relative error as a function of the mesh-size when estimating $F_Y(0.8)$ where $Y = \sum_{i=1}^{16} X_i$ and $X_i \sim$ Nakagami-m($m$). Right: Relative error as a function of the mesh-size when estimating $F_Y(0.8)$ where $Y = \sum_{i=1}^{16} X_i$ and $X_i \sim$ Rice($\nu$)

4.4. **Estimating unknown distributions.** In this subsection we explore the convergence properties of the convolution method for Nakagami-m and the Rice distribution. The PDFs are given in Table 1. Note that we let $\Omega = 1$ for the Nakagami-m, while we for the Rice distribution use an alternative parameterization where we let $K = \frac{\nu^2}{2}$ and $\Omega = \nu^2 + 2$. Similarly to the Log-Normal distribution, which was the topic of the former subsection 4.3, we do not know the exact distribution of a sum of i.i.d Nakagami-m or Rice RVs. We run the same experiment as before, utilizing the convolution method to estimate the value

$$\alpha = F_Y(0.8) = \mathbb{P}(Y \leq 0.8) = \mathbb{P}\left(\sum_{i=1}^{16} X_i \leq 0.8\right),$$

with $X_i$, $i = 1, 2, \cdots, 16$, are i.i.d RVs drawn from either the Nakagami-m distribution or the Rice distribution. We first calculate a pseudo-reference solution $\bar{\alpha}$ using a mesh consisting of $N = 2^{20}$ intervals. We then calculate estimates $\bar{\alpha}_{m_i}, i = 7, 8, \ldots, 15$ where we utilize a mesh of size $N = 2^i$ in order to calculate $\bar{\alpha}_{m_i}$. The relative error of the approximation relative to the pseudo-reference solution is then calculated by

$$\delta_{m_i} = \frac{|\bar{\alpha} - \bar{\alpha}_{m_i}|}{\bar{\alpha}}$$

For the Nakagami-m case we see from the left plot in Figure 5 that we end up with convergence rates of $N^{-6}, N^{-4}$ and $N^{-2}$ when choosing parameter values $m = 3, m = 2$ and $m = 1$ respectively. These observations are consistent with Theorem 1 as the pdf of the Nakagami-m distribution is zero at zero for all derivatives up to and including the forth derivative when $m = 3$, while the same is true up to the second derivative for $m = 2$. For the case $m = 1$ the first derivative is not zero at zero.

The last distribution we will consider is the Rice distribution. The result from the numerical experiment is shown in the right plot of Figure 5. Here the relative error more or less coincide for all tested parameter values. We also note that the empirically observed convergence rate is $N^{-2}$. This also agree with the result from Theorem 1 (see Remark 4) as the first derivative of the pdf of the Rice distribution is not zero at zero.

## 5. Conclusion

We have presented a deterministic numerical method for estimating left-tail rare events of sums of non-negative independent RVs. The method is shown to be efficient, flexible, and accurate – even when measured in relative error. This is due to the fact that numerical integration of convoluted densities only involves sums and products of non-negative floating point values, which are operations that are insensitive to rounding errors, cf. Theorem 2. We further compare direct-based convolution to FFT-based convolution, and show by formal theoretical arguments and in numerical experiments that FFT-based convolution is more sensitive to rounding errors and a less reliable method when the magnitude of the probability of failure is sufficiently small. In the future, it would be interesting to explore whether ideas involving numerical integration of linear convolutions to could be extended to estimations of "left-tail" rare events for random vectors, and to right-tail rare events.

## Code availability

The code used to run the numerical examples presented in this paper can be found at: https://github.com/johannesvm/convolution-left-side-rare-events

## References

1. LLC Advanpix, *Multiprecision computing toolbox for matlab*, 2006.
2. Mohamed-Slim Alouini, Nadhir Ben Rached, Abla Kammoun, and Raul Tempone, *On the efficient simulation of the left-tail of the sum of correlated log-normal variates*, Monte Carlo Methods and Applications **24** (2018), no. 2, 101–115.
3. S. Asmussen and P.W. Glynn, *Stochastic simulation: Algorithms and analysis*, Stochastic Modelling and Applied Probability, Springer New York, 2007.
4. Søren Asmussen and Klemens Binswanger, *Simulation of ruin probabilities for subexponential claims*, ASTIN Bulletin: The Journal of the IAA **27** (1997), no. 2, 297–318.
5. Søren Asmussen, José Blanchet, Sandeep Juneja, and Leonardo Rojas-Nandayapa, *Efficient simulation of tail probabilities of sums of correlated lognormals*, Annals of Operations Research **189** (2011), 5–23.
6. Søren Asmussen and Dominik Kortschak, *On error rates in rare event simulation with heavy tails*, Proceedings of the 2012 Winter Simulation Conference (WSC), IEEE, 2012, pp. 1–11.
7. Søren Asmussen and Dirk P Kroese, *Improved algorithms for rare event simulation with heavy tails*, Advances in Applied Probability **38** (2006), no. 2, 545–558.
8. Søren Asmussen, Pierre-Olivier Goffard, and Patrick J. Laub, *Orthonormal polynomial expansions and lognormal sum densities*, ch. Chapter 6, pp. 127–150.
9. Søren Asmussen, Jens Ledet Jensen, and Leonardo Rojas-Nandayapa, *Exponential family techniques for the lognormal left tail*, Scandinavian Journal of Statistics **43** (2016), no. 3, 774–787.
10. N.C. Beaulieu and F. Rajwani, *Highly accurate simple closed-form approximations to lognormal sum distributions and densities*, IEEE Communications Letters **8** (2004), no. 12, 709–711.
11. N.C. Beaulieu and Qiong Xie, *An optimal lognormal approximation to lognormal sum distributions*, IEEE Transactions on Vehicular Technology **53** (2004), no. 2, 479–489.
12. Norman C. Beaulieu, *Fast convenient numerical computation of lognormal characteristic functions*, IEEE Transactions on Communications **56** (2008), no. 3, 331–333.
13. Norman C. Beaulieu and Gan Luan, *Improving simulation of lognormal sum distributions with hyperspace replication*, 2019 IEEE Global Communications Conference (GLOBECOM), 2019, pp. 1–7.
14. _____, *On the marcum q-function behavior of the left tail probability of the lognormal sum distribution*, ICC 2020 - 2020 IEEE International Conference on Communications (ICC), 2020, pp. 1–6.
15. Eya Ben Amar, Nadhir Ben Rached, Abdul-Lateef Haji-Ali, and Raúl Tempone, *State-dependent importance sampling for estimating expectations of functionals of sums of independent random variables*, Statistics and Computing **33** (2023), no. 2, 40.
16. Nadhir Ben Rached, Fatma Benkhelifa, Abla Kammoun, Mohamed-Slim Alouini, and Raul Tempone, *On the generalization of the hazard rate twisting-based simulation approach*, Statistics and Computing **28** (2018), 61–75.

17. Nadhir Ben Rached, Abla Kammoun, Mohamed-Slim Alouini, and Raul Tempone, *Unified importance sampling schemes for efficient simulation of outage capacity over generalized fading channels*, IEEE Journal of Selected Topics in Signal Processing **10** (2016), no. 2, 376–388.

18. ———, *On the efficient simulation of outage probability in a log-normal fading environment*, IEEE Transactions on Communications **65** (2017), no. 6, 2583–2593.

19. Edward Furman, Daniel Hackmann, and Alexey Kuznetsov, *On log-normal convolutions: An analytical–numerical method with applications to economic capital determination*, Insurance: Mathematics and Economics **90** (2020), 120–134.

20. I. S. Gradshteyn and I. M. Ryzhik, *Table of integrals, series, and products*, seventh ed., Elsevier/Academic Press, Amsterdam, 2007.

21. Archil Gulisashvili and Peter Tankov, *Tail behavior of sums and differences of log-normal random variables*, Bernoulli **22** (2016), no. 1, 444 – 493.

22. Sudarshan Guruacharya, Hina Tabassum, and Ekram Hossain, *Saddle point approximation for outage probability using cumulant generating functions*, IEEE Wireless Communications Letters **5** (2016), no. 2, 192–195.

23. Eugene Isaacson and Herbert Bishop Keller, *Analysis of numerical methods*, Courier Corporation, 2012.

24. Sandeep Juneja and Perwez Shahabuddin, *Simulating heavy tailed processes using delayed hazard rate twisting*, ACM Transactions on Modeling and Computer Simulation (TOMACS) **12** (2002), no. 2, 94–118.

25. Uri Keich, *sfft: a faster accurate computation of the p-value of the entropy score*, Journal of Computational Biology **12** (2005), no. 4, 416–430.

26. Dirk P Kroese, Thomas Taimre, and Zdravko I Botev, *Handbook of monte carlo methods*, John Wiley & Sons, 2011.

27. JosÉ A. Lopez-Salcedo, *Simple closed-form approximation to ricean sum distributions*, IEEE Signal Processing Letters **16** (2009), no. 3, 153–155.

28. Karthyek Rajhaa A. M. and Sandeep Juneja, *State-independent importance sampling for estimating large deviation probabilities in heavy-tailed random walks*, 6th International ICST Conference on Performance Evaluation Methodologies and Tools, 2012, pp. 127–135.

29. Marco Di Renzo, Laura Imbriglio, Fabio Graziosi, and Fortunato Santucci, *Smolyak's algorithm: a simple and accurate framework for the analysis of correlated log-normal power-sums*, IEEE Communications Letters **13** (2009), no. 9, 673–675.

30. Seyed Ali Saberali and Norman C. Beaulieu, *New approximations to the lognormal characteristic function*, 2012 IEEE Global Communications Conference (GLOBECOM), 2012, pp. 2168–2172.

31. D. Senaratne and C. Tellambura, *A general numerical method for computing the probability of outage*, 2009 IEEE Wireless Communications and Networking Conference, 2009, pp. 1–6.

32. Damith Senaratne and Chintha Tellambura, *Numerical computation of the lognormal sum distribution*, GLOBECOM 2009 - 2009 IEEE Global Telecommunications Conference, 2009, pp. 1–6.

33. Endre Süli and David F Mayers, *An introduction to numerical analysis*, Cambridge university press, 2003.

34. C. Tellambura and A. Annamalai, *An unified numerical approach for computing the outage probability for mobile radio systems*, IEEE Communications Letters **3** (1999), no. 4, 97–99.

35. C. Tellambura and D. Senaratne, *Accurate computation of the mgf of the lognormal distribution and its application to sum of lognormals*, IEEE Transactions on Communications **58** (2010), no. 5, 1568–1577.

36. V. Pham Thanh, Cong-Thang Truong, and T. Pham Anh, *On the mgf-based approximation of the sum of independent gamma-gamma random variables*, 2015 IEEE 81st Vehicular Technology Conference (VTC Spring), 2015, pp. 1–5.

37. Huon Wilson and Uri Keich, *Accurate pairwise convolutions of non-negative vectors via fft*, Computational Statistics & Data Analysis **101** (2016), 300–315.

38. ———, *Accurate small tail probabilities of sums of iid lattice-valued random variables via fft*, Journal of Computational and Graphical Statistics **26** (2017), no. 1, 223–229.

39. Jingxian Wu, N.B. Mehta, and Jin Zhang, *Flexible lognormal sum approximation method*, GLOBECOM '05. IEEE Global Telecommunications Conference, 2005., vol. 6, 2005, pp. 3413–3417.

40. Zhiqiang Xiao, Bingcheng Zhu, Julian Cheng, and Yongjin Wang, *Outage probability bounds of egc over dual-branch non-identically distributed independent lognormal fading channels with optimized parameters*, IEEE Transactions on Vehicular Technology **68** (2019), no. 8, 8232–8237.

41. Bingcheng Zhu, Zaichen Zhang, Lei Wang, Jian Dang, Liang Wu, Julian Cheng, and Geoffrey Ye Li, *Right tail approximation for the distribution of lognormal sum and its applications*, 2020 IEEE Globecom Workshops (GC Wkshps, 2020, pp. 1–6.

School of Mathematics, University of Leeds, United Kingdom
*Email address*: `n.benrached@leeds.ac.uk`

Department of Mathematics, University of Oslo, Norway
*Email address*: `haakonah@math.uio.no`

Department of Mathematics, University of Oslo, Norway
*Email address*: `johannvm@math.uio.no`