

Exactly computing the tail of the Poisson-Binomial Distribution

Noah Peres, Andrew Lee and Uri Keich

April 17, 2020

Abstract

We offer ShiftConvolvePoibin, a fast exact method to compute the tail of a Poisson-Binomial distribution (PBD). Our method employs an exponential shift to retain its accuracy when computing a tail probability, and in practice we find that it is immune to the significant relative errors that other methods, exact or approximate, can suffer from when computing very small tail probabilities of the PBD. The accompanying R package is also competitive with the fastest implementations for computing the entire PBD.

1 Introduction

Let X_i for $i = 1, \dots, N$ be independent Bernoulli random variables (RVs) with corresponding probabilities of success p_i . The distribution of $X = \sum_1^N X_i$ is called a

Poisson-binomial distribution (PBD) and it clearly generalizes the binomial distribution ($p_i \equiv p$ for all i).

Biscarri et al. highlighted the considerable interest in using the PBD in diverse areas of scientific research ranging from genetics through survey sampling to sports analysis [1]. Coinciding with this scientific interest was a continuing effort on part of the statistical community to develop efficient and accurate tools for evaluating the significance of tests based on the PBD. Most of these tools, which range from approximation-based ones, including Poisson and normal approximations, to a growing recent interest in exact methods are also reviewed in [1].

In this paper we focus on computing the tail probability of the PBD. Specifically we look at evaluating $P(X \geq x)$ for a Poisson-binomial (PB) RV X focusing on accurately computing this (right) tail probability when it is small. Note that a left tail can be readily transformed into a right tail by considering $X' = N - X = \sum_1^N (1 - X_i)$ which is again a PB RV.

Recently Madsen et al. pointed out that when using Hong's popular DFT-CF [5] to compute the right tail of the PBD one can expect extremely high *relative errors* when the actual value is smaller than about 10^{-16} [8].¹ While many statisticians would not consider that a problem, and indeed it is not a problem in the canonical 5% significance cutoff scenario, Madsen et al.'s motivation comes from bioinformatics research where one often needs to accurately evaluate even much smaller tails. In such cases relying on DFT-CF could lead to significant errors in the downstream analysis.

¹If the true value $s \neq 0$ is computed as \tilde{s} then the associated *relative error* is defined as $|(\tilde{s} - s)/s|$.

As an alternative, Madsen et al. develop a saddle-point based approximation method which, as they demonstrate, has a much better-behaved relative error when computing the p-value of large values of the statistic X [8]. Moreover, as Madsen et al. argue, the runtime complexity of their approach is $O(N)$ compared with DFT-CF’s complexity of $O(N^2)$ making their method a seemingly win-win proposition for someone interested in evaluating the right tail of the PBD.

In this paper we show that while the saddlepoint approximation method of Madsen et al. (SA) is generally quite accurate, there are cases where it also suffers from significantly large errors. Particularly problematic are cases where the actual p-value is close to 1 but SA reports much smaller values (e.g., Figure 2 below). Madsen et al. acknowledged that SA “is not suited for calculating large (not significant) p-values (> 0.1)”, however how would the user know the p-value is > 0.1 if SA reports a value that is much closer to 0 (e.g., consider $s_0 < 4000$ in panel A of Figure 2 below)?

Going back to the extremely large relative errors that DFT-CF can induce, we show below that the same applies to DC-FFT, which is the more recent method of Biscarri et al. [1]. Note that both DFT-CF and DC-FFT are exact methods, that is, they compute the probability mass function (pmf) of the PBD using the underlying distribution rather than relying on an asymptotic approximation. So how can the approximation-based SA be much more accurate than those exact methods?!

Before explaining this we would like to clarify that both DFT-CF and DC-FFT are accurate as long as you gauge their accuracy using the total absolute error (TAE) as your figure of merit. Indeed, as reported by both Hong and Biscarri et al., in that case these exact methods live up to their name with the error rarely exceeding 10^{-10} .

However, as noted above the relative error can tell a very different story.

As explained, for example in [13], the source of the extremely large relative errors is the Discrete Fourier Transform (DFT, defined below), which both DFT-CF and DC-FFT rely on. Specifically, the DFT involves summing many positive and negative numbers so the accumulated roundoff errors are amplified by intermediate cancellations and potentially create extremely large relative errors in the final result — something that we observe with DFT-CF and DC-FFT.

In contrast, if we add up only non-negative numbers then the relative error is well-controlled. In particular, the exact method of Direct Convolution (DC) — proposed by Biscarri et al. (Algorithm 1 [1]), as well as by Madsen et al., can be considered as the gold standard in terms of accuracy. The downside of DC is that its runtime complexity is $O(N^2)$ making it potentially forbiddingly slow for large values of N . Indeed, DC-FFT was specifically designed to be much faster than DC.

In this paper we introduce ShiftConvolvePoibin (or ShiftConvolve for short), a novel exact method with a close-to-linear runtime complexity of $O(N(\log N)^2)$ and which, in practice, is on par with DC in terms of its control of the relative error. ShiftConvolve accomplishes this by relying on the same exponential shift idea that was originally employed in [7] to control the numerical errors introduced by the DFT when computing a certain tail probability (see also [13]).

2 Addressing the accuracy problem

2.1 DFT and the relative error problem

The DFT operator, D , and its inverse, D^{-1} , are linear operators defined on \mathbb{C}^n as

$$\begin{aligned}(D\mathbf{x})(k) &:= \sum_{j=0}^{n-1} \mathbf{x}(j)e^{-i2\pi kj/n} & k = 0, \dots, n-1 \\ (D^{-1}\mathbf{y})(k) &:= \frac{1}{n} \sum_{j=0}^{n-1} \mathbf{y}(j)e^{i2\pi kj/n} & k = 0, \dots, n-1.\end{aligned}\tag{1}$$

It is clear from their definitions that computing the real and imaginary components of the result generally involves adding up both positive and negative numbers creating the significant relative errors we alluded to.

In practice the DFT is almost invariably calculated by the Fast Fourier Transform (FFT) which has a time complexity of $O(n \log n)$ [3]. Brisebarre et al. provide an exhaustive overview of the analysis of the error introduced by the FFT [2] and it should be stressed that in general a naive implementation of the DFT would not be any more accurate than the FFT [11]. Regardless, our goal here is not to bound the error, rather we follow up on Madsen et al. in pointing out the potentially catastrophic effect the FFT can have on the relative error and to offer a solution to this problem.

Both DFT-CF and DC-FFT rely on the DFT (and its inverse): the first to invert the characteristic function of the PBD and the second to perform its convolutions as explained below. In both cases the DFT can severely compromise the relative accuracy of the computed values as shown in panel A of Figure 1: while the DFT-

relying methods coincide with the accurate DC for pmf entries that are larger than $\approx 10^{-16}$, values that are smaller cannot be recovered by DFT-CF and DC-FFT (DFT-CF often returns 0 in this case as is evident by the gaps in the green logarithmic curve). In particular, using DFT-CF or DC-FFT to compute a right tail probability for large values of x will typically yield a result that is orders of magnitude off in terms of the relative accuracy: the correct values as per DC are orders of magnitude smaller than the values reported by DC-FFT, whereas for DFT-CF the reported values vary between 0 (100% relative error) and the same order as DC-FFT's reported values.²

Panel A of Figure 1 suggests that it is impossible to recover the smaller entries of the pmf because of the numerical errors inherent to the DFT. There is however a solution that was first suggested in [7] and that uses an exponential shift to overcome the numerical errors.

2.2 The exponential shift

Let \mathbf{q} be a PB pmf supported on $0, 1, \dots, n$ then the exponentially shifted version of \mathbf{q} is defined as

$$\mathbf{q}_\theta(k) = \mathbf{q}(k)e^{k\theta}/M_{\mathbf{q}}(\theta) \quad k = 0, 1, \dots, n, \quad (2)$$

where $M_{\mathbf{q}}(\theta)$ is the moment generating function (MGF) of \mathbf{q} . Note that dividing by the MGF guarantees that \mathbf{q}_θ is a proper pmf.

The basic idea behind the introduction of the exponential shift is that the FFT-based methods are still accurate for the larger values so all we need to do is make

²Note that the absolute error is still small: no more than $\approx 10^{-16}$ but our interest here is in accurately recovering the small values.

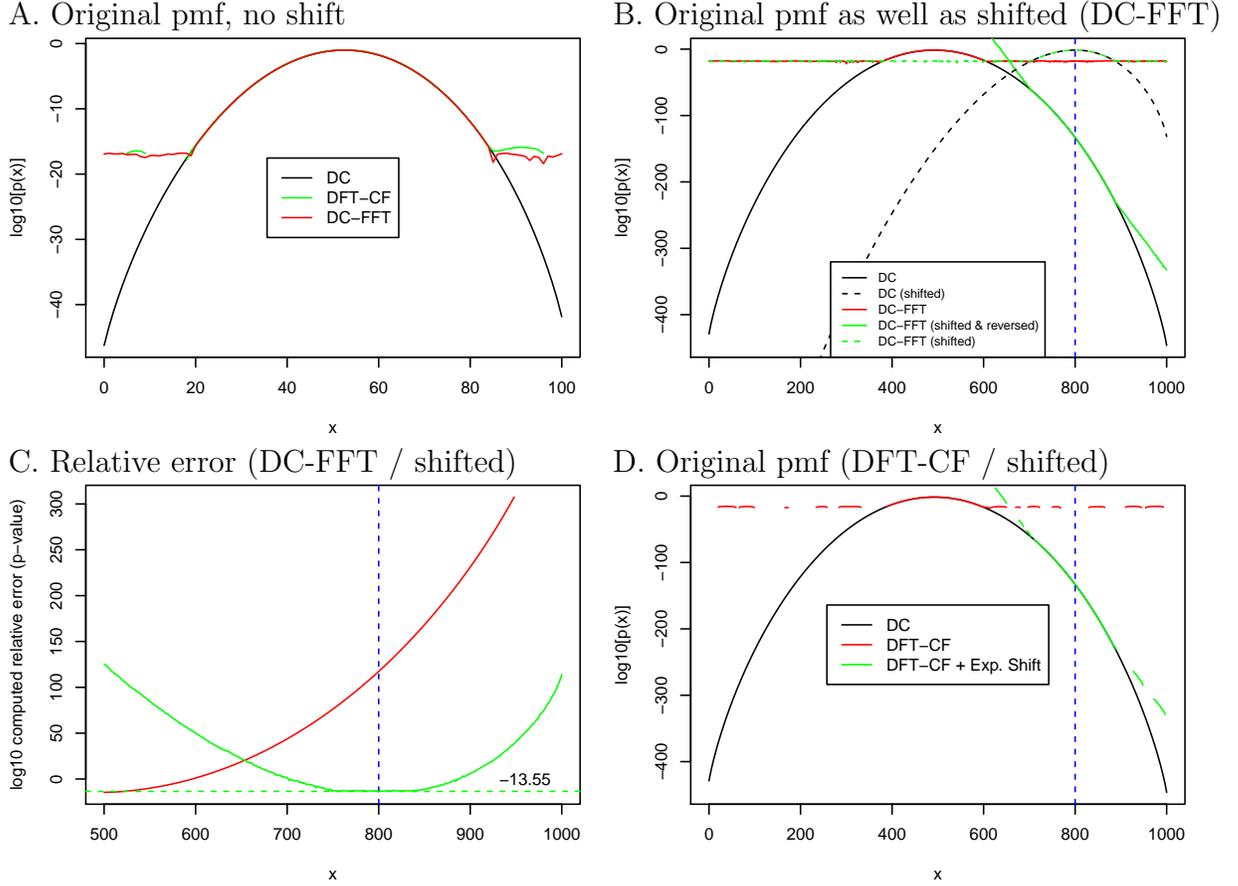


Figure 1: **The accuracy problem.** (A) The (log of the) PB pmf computed by the accurate DC is contrasted with the values computed by DFT-CF and DC-FFT both of which rely on the DFT: values below $\approx 10^{-16}$ cannot be recovered. Missing values of DFT-CF correspond to reported 0 values. The PBD here is defined using $N = 100$ values of p_i that were independently and uniformly sampled. (B) The log of the PB pmf computed by DC-FFT compared with a shifted variant. (C) The (log of the) relative error in computing the right tail using DC-FFT and its shifted variant. (D) Same as panel B but using DFT-CF and a shifted variant of DFT-CF. Gaps correspond to reported 0 values.

sure the values we are actually interested in are “large enough”.

The approach is visualized in panel B of Figure 1, which focuses on accurately recovering $P(X \geq 800)$. As seen by the red curves in panels B and C of the figure, DC-FFT does not allow us to accurately gauge this probability which is many order of magnitude smaller than the values reported by DC-FFT.

However, when the proper exponential shift is applied to the pmf (panel B, dashed black curve) the values that we are interested in are sufficiently inflated so that DC-FFT accurately recovers the section of the pmf about $x = 800$ (panel B, dashed green curve).

Finally, reversing the shift we note that the combination of DC-FFT with the shift and its reversal allows us to accurately recover the tail probability for any x in the neighborhood of 800: the relative error is miniscule ($\approx 10^{-13}$, panels B and C, solid green curves). Panel D suggests that the same principle would work for DFT-CF.

The reason we use an exponential shift rather than some other arbitrary way to inflate the values we are interested in is that the exponential shift can be readily applied, as well as peeled off, or reversed in our context. Specifically:

Claim 1. The exponential shift commutes with the convolution operation $(*)$, that is, if \mathbf{p} and \mathbf{q} are pmfs defined on $0, 1, \dots, n$ then $(\mathbf{p} * \mathbf{q})_\theta \equiv \mathbf{p}_\theta * \mathbf{q}_\theta$.

Proof. Recall that if X and Y are independent \mathbb{N} -valued RVs with corresponding pmfs \mathbf{p} and \mathbf{q} then

$$M_{\mathbf{p} * \mathbf{q}} \equiv M_{X+Y} \equiv M_X \cdot M_Y \equiv M_{\mathbf{p}} \cdot M_{\mathbf{q}}.$$

With this in mind for $k = 0, \dots, 2n$ we have

$$\begin{aligned}
(\mathbf{p} * \mathbf{q})_\theta(k) &= (\mathbf{p} * \mathbf{q})(k) \cdot e^{k\theta} / M_{\mathbf{p} * \mathbf{q}}(\theta) \\
&= \sum_{i=0}^k (\mathbf{p}(i) \cdot e^{i\theta}) (\mathbf{q}(k-i) \cdot e^{(k-i)\theta}) / (M_{\mathbf{p}}(\theta) \cdot M_{\mathbf{q}}(\theta)) \\
&= (\mathbf{p}_\theta * \mathbf{q}_\theta)(k).
\end{aligned}$$

□

It follows that the same holds for convolutions of any order and hence that we can compute \mathbf{q}_θ , the θ -shifted version of the PB pmf \mathbf{q} , by convolving the θ -shifted versions of each of the Bernoulli pmfs. The latter of course can be trivially found because a θ -shifted Bernoulli(p) pmf is again a Bernoulli pmf only with probability of success $p_\theta = pe^\theta / (1 - p + pe^\theta)$. Note that, in particular, an exponentially shifted PBD is also a PBD.

It follows that we can use DFT-CF and DC-FFT to compute the exponentially shifted pmf \mathbf{p}_θ of the PBD by applying them to θ -shifted Bernoullis. Of course, recovering \mathbf{q} from \mathbf{q}_θ is a trivial exercise of inverting (2), or equivalently applying a shift of $-\theta$ to \mathbf{q}_θ . This is indeed the procedure we applied when obtaining the green curves of panels B and D of Figure 1.

What is missing at this point is the protocol for determining θ given x . Here again we follow [7], and we define the shift as the value of θ that minimizes the expression $\log M_{\mathbf{q}}(\theta) - \theta \cdot x$. The critical point of the latter function is the value θ for which $M'_{\mathbf{q}}(\theta) / M_{\mathbf{q}}(\theta) = x$ but $M'_{\mathbf{q}}(\theta) / M_{\mathbf{q}}(\theta)$ is simply the expected value of the

shifted \mathbf{q}_θ . Hence, finding θ amounts to solving

$$x = E(X_\theta) = \sum_{i=1}^N E(X_{i,\theta}) = \sum_{i=1}^N \frac{p_i e^\theta}{1 - p_i + p_i e^\theta}, \quad (3)$$

where $X_\theta = \sum_1^N X_{i,\theta}$ is a \mathbf{q}_θ distributed RV and $X_{i,\theta}$ are independent Bernoulli($\mathbf{p}_\theta(i) = p_i e^\theta / (1 - p_i + p_i e^\theta)$) RVs. In practice we solve (3) numerically using the `uniroot` function in R.

3 ShiftConvolvePoibin

Combining our exponential shift protocol with either DFT-CF or DC-FFT produces an exact method that is able to practically recover any right tail probability (p-value) with a negligible relative error. The combined general procedure is outlined in Algorithm 1.

Consider the variant of Algorithm 1 where the pmf of the (shifted) PBD is computed with DC-FFT with its default setting of $M = N$. In this case, starting with the N (shifted) Bernoulli pmfs at each step DC-FFT pairs all intermediate (shifted) PB pmfs and convolves each pair while passing the resulting PB pmf to the next step until at the last step it ends up with a single PB pmf (Algorithm (2)).

Each such pairwise convolution is executed using the FFT based on the following identity. Suppose that the pmfs \mathbf{p} and \mathbf{q} are supported on $\{0, 1, \dots, m\}$ and $\{0, \dots, n\}$ respectively and for $Q \geq m + n + 1$ embed \mathbf{p} and \mathbf{q} in \mathbb{R}^Q by extending them with $Q - m$, respectively $Q - n$, zeros. Then using the Q -dimensional DFT

Algorithm 1: shift-compute-PB-pmf-unshift

Input: $\mathbf{p} = (p_0, p_1, \dots, p_{N-1})$ vector of success probabilities, observed value s_0
/* Lines 1-5: apply the exponential shift to \mathbf{p} */
1 **Function:** $\mu(\theta, \mathbf{p}) \leftarrow \text{sum}([p_i \exp(\theta)/(1 - p_i + p_i \exp(\theta)) \text{ for } p_i \in \mathbf{p}])$
2 $\theta^* \leftarrow$ find root θ such that $\mu(\theta, \mathbf{p}) - s_0 = 0$
3 $\mathbf{p}_\theta \leftarrow$ empty list
4 **for** $p_i \in \mathbf{p}$ **do**
5 | $\mathbf{p}_\theta(i)^+ \leftarrow p_i \exp(\theta^*)/(1 - p_i + p_i \exp(\theta^*))$
6 $\mathbf{q}_\theta \leftarrow$ apply DC-FFT or DFT-CF to \mathbf{p}_θ
/* Lines 7-10: reverse the exponential shift */
7 $\mathbf{q} \leftarrow$ empty list
8 $M \leftarrow \text{prod}([(1 - p_i + p_i \exp(\theta^*)) \text{ for } p_i \in \mathbf{p}])$
9 **for** $j \in 0, 1, \dots, N - 1$ **do**
10 | $\mathbf{q}(j) \leftarrow \mathbf{q}_\theta(j) \exp(-j \cdot \theta) \cdot M$
Output: \mathbf{q} : the convolved pmf

operator and its inverse we have [12]

$$\mathbf{p} * \mathbf{q} = D^{-1}(D\mathbf{p} \odot D\mathbf{q}), \quad (4)$$

where for Q -dimensional vectors \mathbf{u} and \mathbf{v} , $(\mathbf{u} \odot \mathbf{v})(i) = \mathbf{u}(i)\mathbf{v}(i)$.

At first glance DC-FFT ($M = N$) seems wasteful because each intermediate pmf is computed by applying D^{-1} while at the next step D is applied to the same pmf (lines 13 and 15). Of course, these two operators have different dimensions (the latter's twice the former's) so this is not as wasteful as it might initially look. Still, we next show how some work can be saved by making a more efficient use of what has already been computed in the previous step.

For $\mathbf{u} \in \mathbb{C}^n$ let $\mathbf{u}^* \in \mathbb{C}^{2n}$ be the zero-padded $2n$ -dimensional version of \mathbf{u} : $\mathbf{u}^*(i) =$

Algorithm 2: Pair-aggregated-FFT-convolution (DC-FFT with $M = N$)

Input: $p = (p_0, p_1, \dots, p_{N-1})$ vector of success probabilities

- 1 $V \leftarrow$ empty list
- 2 **for** $p_i \in p$ **do**
- 3 | add FFT $((1 - p_i, p_i, 0, 0))$ to list V
- 4 $L \leftarrow$ length(V)
- 5 **while** $L > 2$ **do**
- 6 | $V^* \leftarrow$ empty list
- 7 | **if** L is odd **then**
- 8 | | $n \leftarrow$ length(\mathbf{v} in V)
- 9 | | add the n -dimensional vector $(1, 1, \dots, 1)$ to list V
- 10 | split V into pairs $(\mathbf{v}_i, \mathbf{v}_j)$
- 11 | **for** each pair $(\mathbf{v}_i, \mathbf{v}_j)$ in V **do**
- 12 | | $\mathbf{u} \leftarrow$ FFTInverse(PointwiseMultiply($\mathbf{v}_i, \mathbf{v}_j$))
- 13 | | $\mathbf{u} \leftarrow$ pad \mathbf{u} with length(\mathbf{u}) 0s
- 14 | | $\mathbf{v}^* \leftarrow$ FFT(\mathbf{u})
- 15 | | add \mathbf{v}^* to list V^*
- 16 | $V \leftarrow V^*$
- 17 | $L \leftarrow$ length(V)
- 18 $\mathbf{u} \leftarrow$ FFTInverse(PointwiseMultiply($V[0], V[1]$))

Output: \mathbf{u}

$\mathbf{u}(i)$ for $i = 0, \dots, n-1$ and $\mathbf{u}^*(i) = 0$ for $i = n, \dots, 2n-1$. Then for $k = 0, \dots, n-1$

$$(D_{2n}\mathbf{u}^*)(2k) = \sum_{j=0}^{2n-1} \mathbf{u}^*(j) e^{-i2\pi(2k)j/(2n)} = \sum_{j=0}^{n-1} \mathbf{u}(j) e^{-i2\pi kj/n} = (D_n\mathbf{u})(k). \quad (5)$$

That is, the even entries of $D_{2n}\mathbf{u}^*$ coincide with $D_n\mathbf{u}$ and therefore we do not need to recompute them.

As for the odd entries of $D_{2n}\mathbf{u}^*$, for $k = 0, \dots, n-1$

$$(D_{2n}\mathbf{u}^*)(2k+1) = \sum_{j=0}^{2n-1} \mathbf{u}^*(j) e^{-i2\pi(2k+1)j/(2n)} = \sum_{j=0}^{n-1} (\mathbf{u}(j) e^{-i\pi j/n}) e^{-i2\pi kj/n} = (D_n(\mathbf{u} \odot \boldsymbol{\omega}))(k), \quad (6)$$

where $\boldsymbol{\omega}(j) = e^{-i\pi j/n}$ for $j = 0, \dots, n-1$ and \odot again is the coordinate-wise product. So while these odd entries do not come for free they can be computed using an n -dimensional DFT rather than a $2n$ -dimensional one.

Algorithm 3 takes advantage of the last two identities to speed up DC-FFT ($M = N$), or Algorithm (2), by about 50% for large N (more on that in Section 4.2 below).

Our ShiftConvolvePoibin algorithm (Algorithm 4) combines Algorithm 1 with the latter, more efficient, Algorithm 3. Note that due to its built-in exponential shift ShiftConvolve can return the accurate logarithm of the right tail probability even when the actual number can create an underflow or 0 when not using logs.

We have two implementations of ShiftConvolve where in both cases the critical part of the code is written in C and is wrapped in an R package. The two versions differ in which code they use to execute the FFT: one version relies on FFTW [4] and

Algorithm 3: Frugal-pair-aggregated-FFT-convolution (FPA-FFTC)

Input: $\mathbf{p} = (p_0, p_1, \dots, p_{N-1})$ vector of success probabilities

- 1 $V \leftarrow$ empty list
- 2 **for** $p_i \in \mathbf{p}$ **do**
- 3 | add FFT($(1 - p_i, p_i, 0, 0)$) to list V
- 4 $L \leftarrow$ length(V)
- 5 **while** $L > 2$ **do**
- 6 | $V^* \leftarrow$ empty list
- 7 | $n \leftarrow$ length(\mathbf{v} in V)
- 8 | $\boldsymbol{\omega} \leftarrow [\exp(ij\pi/n)$ **for** $j \in [0, 1, \dots, n - 1]$
- 9 | **if** L is odd **then**
- 10 | | add the n -dimensional vector $(1, 1, \dots, 1)$ to list V
- 11 | split V into pairs $(\mathbf{v}_i, \mathbf{v}_j)$
- 12 | **for** each pair $(\mathbf{v}_i, \mathbf{v}_j)$ in V **do**
- 13 | | $\mathbf{u} \leftarrow$ FFTInverse(PointwiseMultiply($\mathbf{v}_i, \mathbf{v}_j$))
- 14 | | $\mathbf{w} \leftarrow$ FFT(PointwiseMultiply($\mathbf{u}, \boldsymbol{\omega}$))
- 15 | | | /* assign the components of \mathbf{v} and \mathbf{w} to even and odd
- 16 | | | components of \mathbf{v}^* respectively: */
- 15 | | $\mathbf{v}^* \leftarrow (\mathbf{v}[0], \mathbf{w}[0], \mathbf{v}[1], \mathbf{w}[1], \dots, \mathbf{v}[n - 1], \mathbf{w}[n - 1])$
- 16 | | add \mathbf{v}^* to list V^*
- 17 | $V \leftarrow V^*$
- 18 | $L \leftarrow$ length(V)
- 19 $\mathbf{u} \leftarrow$ FFTInverse(PointwiseMultiply($V[0], V[1]$))

Output: \mathbf{u}

requires the user to install the FFTW package whereas the other uses minFFT [9] and is self-contained. In both cases ShiftConvolve saves some runtime by taking advantage of the fact that the FFTInverse operation on line 13 of Algorithm 3 should produce a real-valued vector.

Algorithm 4: ShiftConvolvePoibin

Input: $\mathbf{p} = (p_0, p_1, \dots, p_{N-1})$ vector of success probabilities, observed value s_0
1 Apply steps 1-5 of Algorithm 1 to get the shifted \mathbf{p}_θ
2 $\mathbf{q}_\theta \leftarrow \text{FPA-FFTC}(\mathbf{p}_\theta)$ /* apply Algorithm 3 to \mathbf{p}_θ */
3 Apply steps 7-10 of Algorithm 1 to get the unshifted \mathbf{q}
Output: \mathbf{q} : the convolved pmf

To conserve runtime and accuracy in practice ShiftConvolve treats the special degenerate cases $p_i \in \{0, 1\}$ differently. That is, before any convolutions are done on the input vector, each $p_i = 0$ is discarded (these correspond to guaranteed failures, or adding the constant random variable 0) and each $p_i = 1$ is also discarded, but represents a (+1) shift in index for the final pmf (guaranteed successes, or adding the constant random variable 1).

Finally, the ShiftConvolvePoibin package offers the user the option of forgoing any exponential shift if one wants to compute the entire pmf rather than a tail.

4 Comparative analysis

In this section we look at the performance of the exact methods DFT-CF, DC-FFT and ShiftConvolvePoibin, as well as the approximation method SA. We start with the analysis of the accuracy of the computed right tail probability.

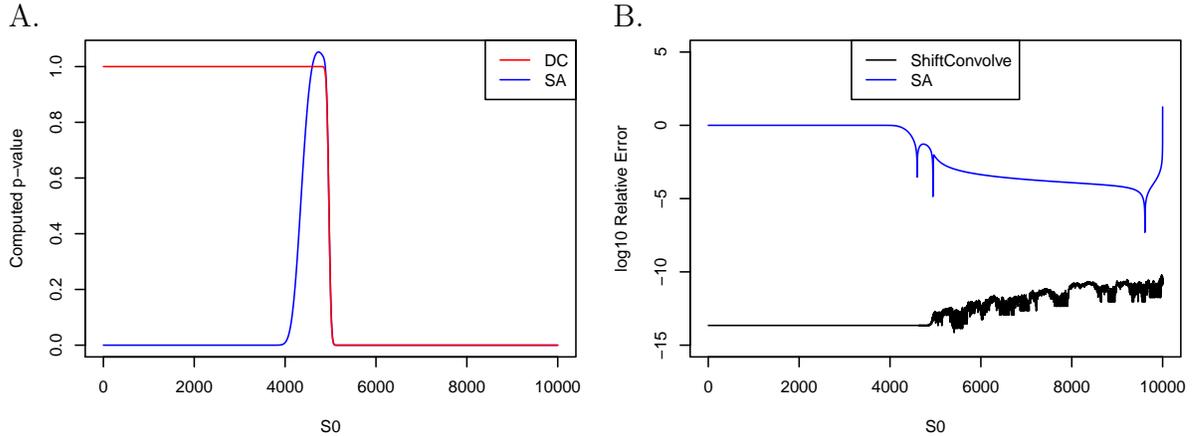


Figure 2: **Failures of SA.** The panels highlight regions where SA fails to accurately compute the right tail probability. Here, $N = 10000$ and the p_i were sampled uniformly from $[0, 1]$. (A) The computed right tail probability (p-value) reported by SA and the accurate DC. (B) The corresponding (log base 10 of) the relative error where we added for reference the significantly smaller relative error of ShiftConvolve. The performance of SA is of particular concern for $s_0 \leq 4000$ as it reports p-values that are smaller than 0.1 giving the user no indication that it might be off from the correct tail probability which is close to 1.

4.1 Accuracy

We first look only at SA and ShiftConvolve. Panel B of Figure 2 confirms that SA offers a good relative accuracy for most of the range of values s_0 for which the tail probability is small (panel A), however, for extremely large values of s_0 the relative accuracy is compromised. More alarming is the fact that for values of $s_0 < 4000$ SA consistently and inaccurately reports very small tail probabilities instead of values close to 1. In addition, SA occasionally reports probabilities that are larger than 1.

Figures 3 and 4 as well as Supplementary Figures 6 and 7 confirm what we noted above: DFT-CF and DC-FFT cannot recover entries that are smaller than approximately 10^{-16} . SA's performance in those figures is consistent with Figure 2 analyzed above: its accuracy is compromised as we are close to the maximal possible value and it is significantly off when the p-value is close to 1. In contrast, ShiftConvolve

retains very good accuracy (10 accurate digits or more) throughout the entire range of values.

We also note that for the more skewed Beta(3, 0.1) distribution with higher proportions of 1s present (due to roundoff errors), the SA algorithm failed to run, returning an error (as in panel C of Supplementary Figure 6, where the blue curve is absent). Similarly, DC-FFT breaks when N is very large ($N \geq 5.5 \cdot 10^5$).

4.2 Complexity / runtime

The memory and runtime complexities of an algorithm are important practical considerations. In our case the memory complexity of all the procedures is linear, or $O(N)$, but their runtimes differ considerably.

Starting with the approximation methods, they all require computing at least a couple of moments of the distribution hence they are typically $O(N)$, which is clearly the case for the normal approximation, as well as the refined normal approximation (RNA) mentioned in [5], but it also applies to the more computationally-intensive SA. The computationally-demanding part of SA is finding the desired exponential shift θ but in practice this requires evaluating the MGF at only a small number of candidate values of θ so the number of operations of this step is still $O(N)$. This is borne out empirically in panel B of Figure 5 where the execution time of SA follows a line of slope 1 in log-log scale.

As for the exact methods, the runtime complexity of DC is $O(N^2)$: there are N steps (the outer loop of Algorithm 1 in [1]) and the number of operations in the k th step is $O(k)$. Moving on to DFT-CF, as noted by Madsen et al., its runtime

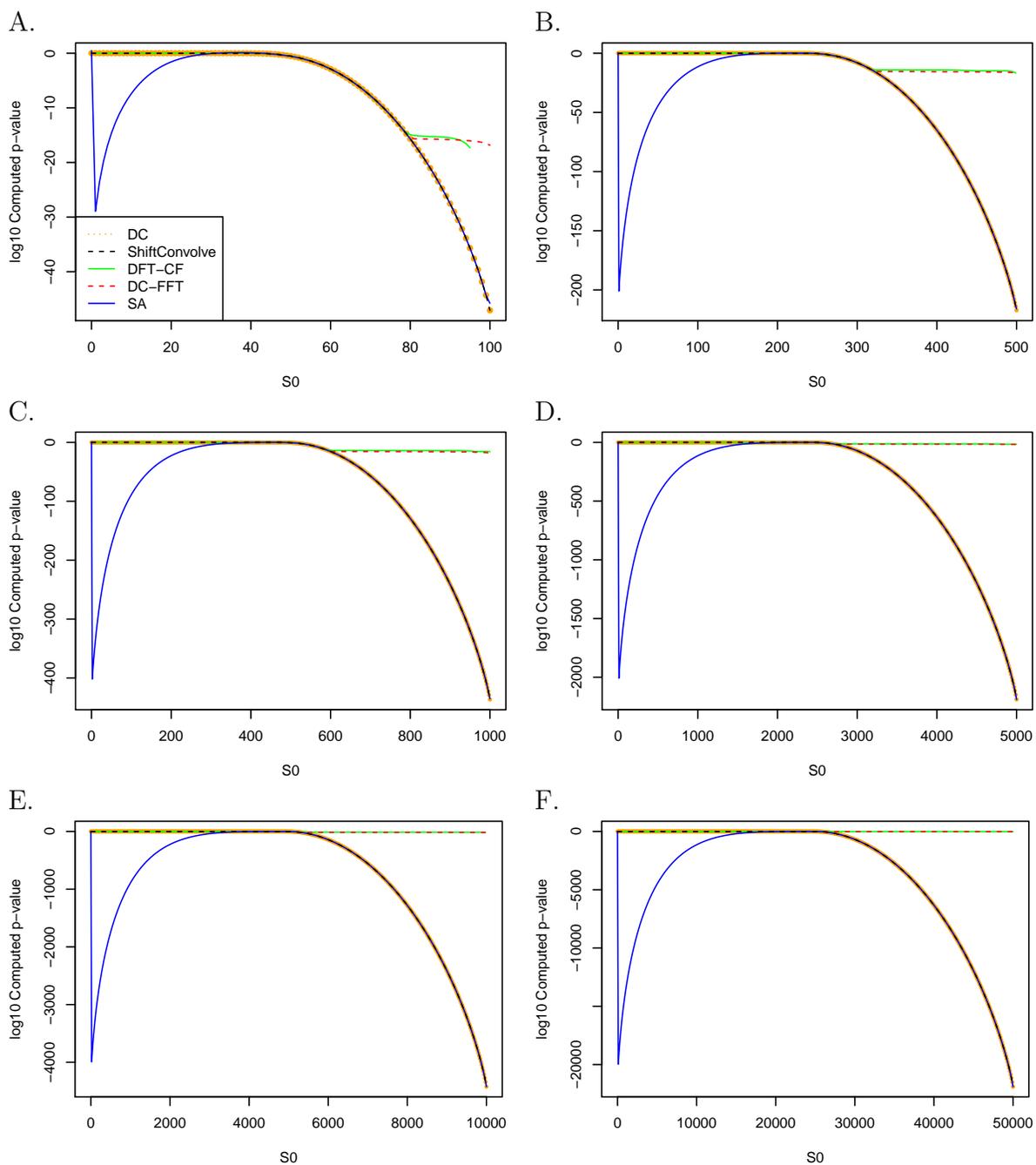


Figure 3: **Right tail probability.** The panels show the (log of the) computed right tail probability (p-value) for varying settings of the parameters. For the corresponding relative errors see Figure 4. (A) $N = 100$, p_i sampled uniformly. DFT-CF does not always extend across the entire range of possible values of the statistic because it reports 0s (log is $-\infty$ in some cases). (B) $N = 500$, p_i sampled uniformly. (C) $N = 1000$, p_i sampled uniformly. (D) $N = 5000$, p_i sampled uniformly. (E) $N = 10000$, p_i sampled uniformly. (F) $N = 50000$, p_i sampled uniformly.

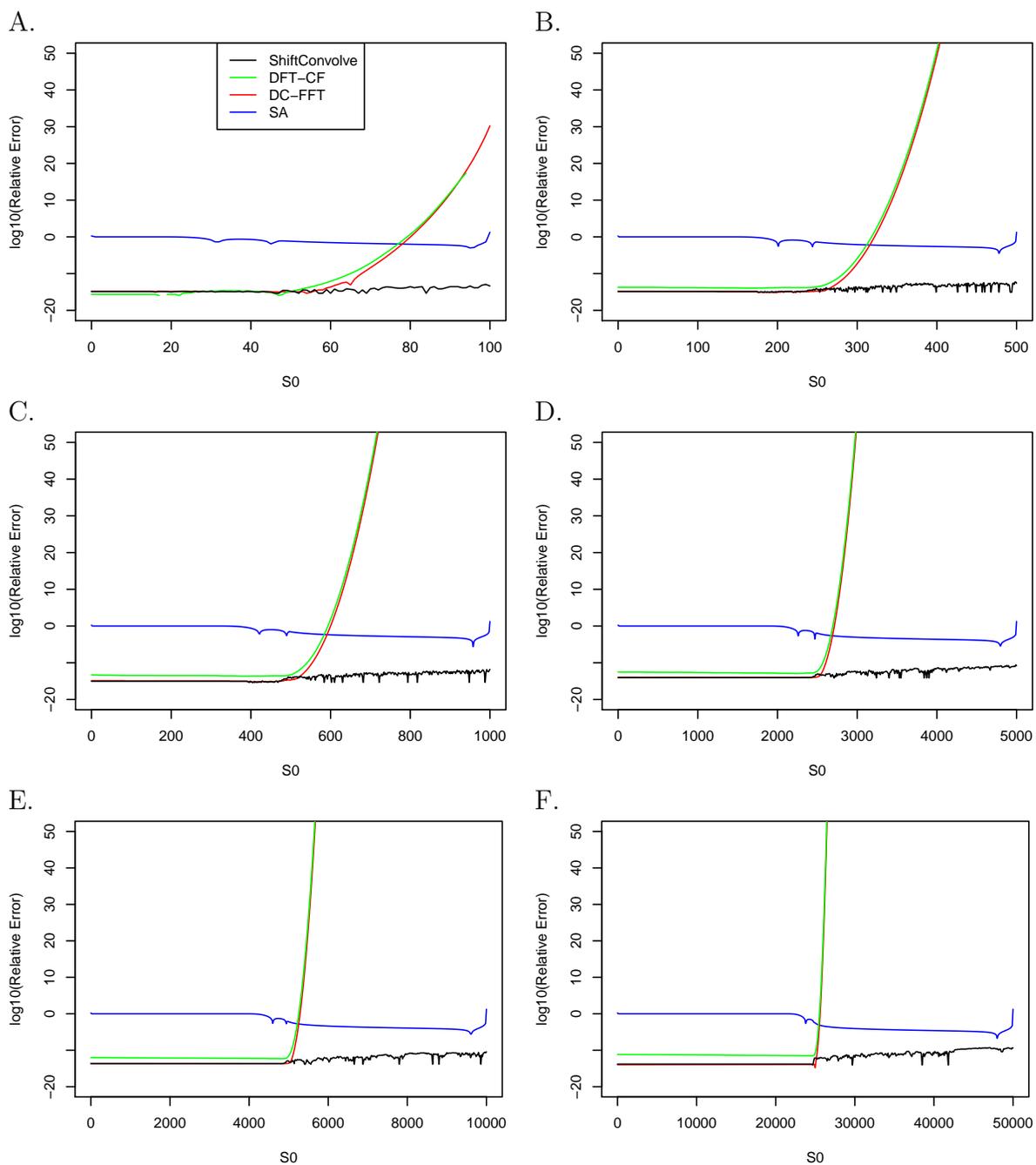


Figure 4: **Gauging the accuracy.** The panels show the (log of the) relative error (with DC taken as the gold standard) in computing the right tail probability for varying settings of the parameters. See Figure 3 for the right tail probability itself. (A) $N = 100$, p_i sampled uniformly. (B) $N = 500$, p_i sampled uniformly. (C) $N = 1000$, p_i sampled uniformly. (D) $N = 5000$, p_i sampled uniformly. (E) $N = 10000$, p_i sampled uniformly. (F) $N = 50000$, p_i sampled uniformly.

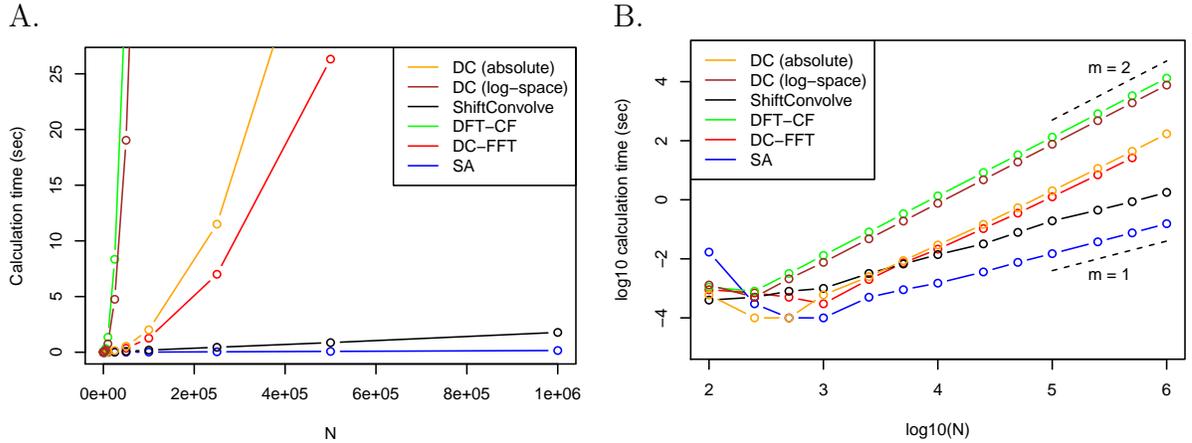


Figure 5: **Runtime comparison.** The panels yield the runtime of the various procedures we consider (panel A is in natural scale whereas panel B is in log-log scale). DC-FFT was run in its default setting of $M = 2$ (the missing point for $N = 10^6$ is due to its failure on that input size). Each data point represents the average runtime over 10 applications of the procedure each using an independently and uniformly drawn N -dimensional vector of Bernoulli success probabilities (different vector for each of the 10 applications). DC (absolute) is Algorithm 1 of Biscarri et al. and DC (log-space) is a version of DC that is part of the ShiftConvolvePoibin package and that works with logs to extend DC's dynamic range. In panel B we added short line segments of slopes $m = 1$ and $m = 2$ for reference.

complexity is also $O(N^2)$ because this is the complexity of the part of the procedure that computes the characteristic function (CF) and inverting it is done using FFT in $O(N \log N)$. Again, these quadratic complexities can be observed in panel B of Figure 5 where the execution times of DFT-CF and DC (we have both a version that works in log-space as well as the default version) follow a line of slope 2 in log-log scale.

Biscarri et al. empirically demonstrated that DC-FFT is significantly faster than DFT-CF but did not offer a complexity analysis of DC-FFT and in fact its runtime complexity varies with the parameter M : with $M = N$ the runtime is the same $O(N(\log N)^2)$ as ShiftConvolve's (see our analysis below) but with its default setting of $M = 2$, or any other constant, it is actually $O(N^2)$ (as can be observed in Figure 5). Biscarri et al. also suggest a heuristic of choosing $M = \max\{2, 2^{\lceil \log_2(N/750) \rceil}\}$, where

$\lceil x \rceil$ is the round-to-nearest-integer function which again has the same asymptotic complexity of $O(N(\log N)^2)$ as when $M = N$, although in practice it can be much faster for any given N (Supplementary Figure 8).

The runtime complexity of our ShiftConvolve is $O(N(\log N)^2)$. Indeed, assuming that $N = 2^K$ there are K steps each of which consists of using FFT to convolve pairs of equal-length vectors: at the k th step there are $N/2^k$ pairwise convolutions of 2^k -dimensional vectors. Hence the number of operations in the k th step is $O(N/2^k \cdot 2^k \log 2^k) = O(N \cdot k)$ and summing over $k = 1, \dots, K$ we get the stated $O(N \cdot K^2)$.

Note that the modifications that we made to speed up Algorithm 2 do not change its runtime complexity. Rather, the improvement is in the constant. Specifically, in Algorithm 2 we compute the n -dimensional inverse DFT $\mathbf{u} = D^{-1}(\mathbf{v}_i \odot \mathbf{v}_j)$ (line 12), as well as the $2n$ -dimensional DFT of the padded \mathbf{u} (line 14). Compare that with applying the same n -dimensional inverse DFT $\mathbf{u} = D^{-1}(\mathbf{v}_i \odot \mathbf{v}_j)$ in Algorithm 3 (line 13), as well as the n -dimensional (rather than $2n$ -dimensional) DFT of $\mathbf{u} \odot \boldsymbol{\omega}$ (line 14). The former requires about 50% more work, hence ShiftConvolve’s improvement offers a reduction of about 50% in the runtime for large N .

ShiftConvolve differs from DC-FFT with its default setting of $M = N$ by employing the exponential shift and its different approach to FFT-based convolution outlined in Algorithm 3. While the two algorithms share the same runtime complexity we found that in practice ShiftConvolve is typically significantly faster (Supplementary Figure 8). In fact, we observed that ShiftConvolve speed is comparable to DC-FFT’s when using the above formula of $M = \max\{2, 2^{\lceil \log_2(N/750) \rceil}\}$ (Supplemen-

tary Figure 8).

Generally, the version of ShiftConvolve that uses minFFT is slightly faster than the version the uses FFTW, hence we used the former version in our runtime benchmarks. All timing runs were executed on a 3.2GHz Intel Core i7 MacMini with 32GB of RAM.

5 Discussion

Madsen et al. pointed out that DFT-CF fails to accurately recover small tail probabilities of the PBD resulting in extremely large relative errors. As we show here, the same applies to DC-FFT, which is also an FFT-based exact algorithm for computing the PBD.

As an alternative Madsen et al. offer a saddle point approximation method (SA), which does a significantly better job at recovering these small tail probabilities. However, SA has its own accuracy issues as we approach the maximal possible value of N , as well as when the right tail probability is rather large. The latter is particularly troubling because the typical user will not be aware that the small p-value that SA reports is in fact very close to 1 (Figure 2).

Following [7, 6, 10] our proposed solution to this problem combines the same exponential shift (sometime referred to as an exponential tilt) that the saddlepoint approximation is based on with the FFT-based exact method. Specifically, ShiftConvolvePoibin uses an exponential shift combined with a souped-up version of the FFT-based convolution by aggregated pairs approach as implemented in DC-FFT

($M = N$). The result is a relatively fast exact algorithm that accurately computes tail probabilities across the entire range of the PBD.

In terms of runtime complexity at $O(N(\log N)^2)$ ShiftConvolve is equivalent to DC-FFT ($M = N$), however both DFT-CF as well as DC-FFT with its default setting of $M = 2$ have a runtime complexity of $O(N^2)$ which makes a significant difference for large N .

It is worth noting that Biscarri et al. recommend using the refined normal approximation (RNA) for $N \geq 10^5$, however Supplementary Figure 9 shows that as far as computing small tail probabilities SA does a significantly better job than RNA and both are inferior when compared with the accuracy of ShiftConvolve.

In terms of future research, while in practice ShiftConvolve is accurate throughout the entire range of possible values it would be useful to obtain upper bounds on its cumulative numerical error which can then be compared with the computed result to guarantee its accuracy (cf. [14]). A related point is that ShiftConvolve is designed to accurately recover the relevant section of the pmf (panel B of Figure 1) however it is worth noting that a few well-selected shifts θ_0 should allow us to accurately recover the entire pmf. Regardless, the ShiftConvolvePoibin package allows the user to compute the entire PBD while bypassing the exponential shift. As such it is competitive with the fastest exact methods in this category.

The ShiftConvolvePoibin package is available to download from <https://github.com/andrew12678/ShiftConvolve>

6 Acknowledgements

UK would like to thank Jakob Pedersen for initial discussions on this problem as well as to Denzel Florez, Tingyue Liu, Xuanchi Liu, and Buqing Yang for initial work on this project.

7 Supplementary Figures

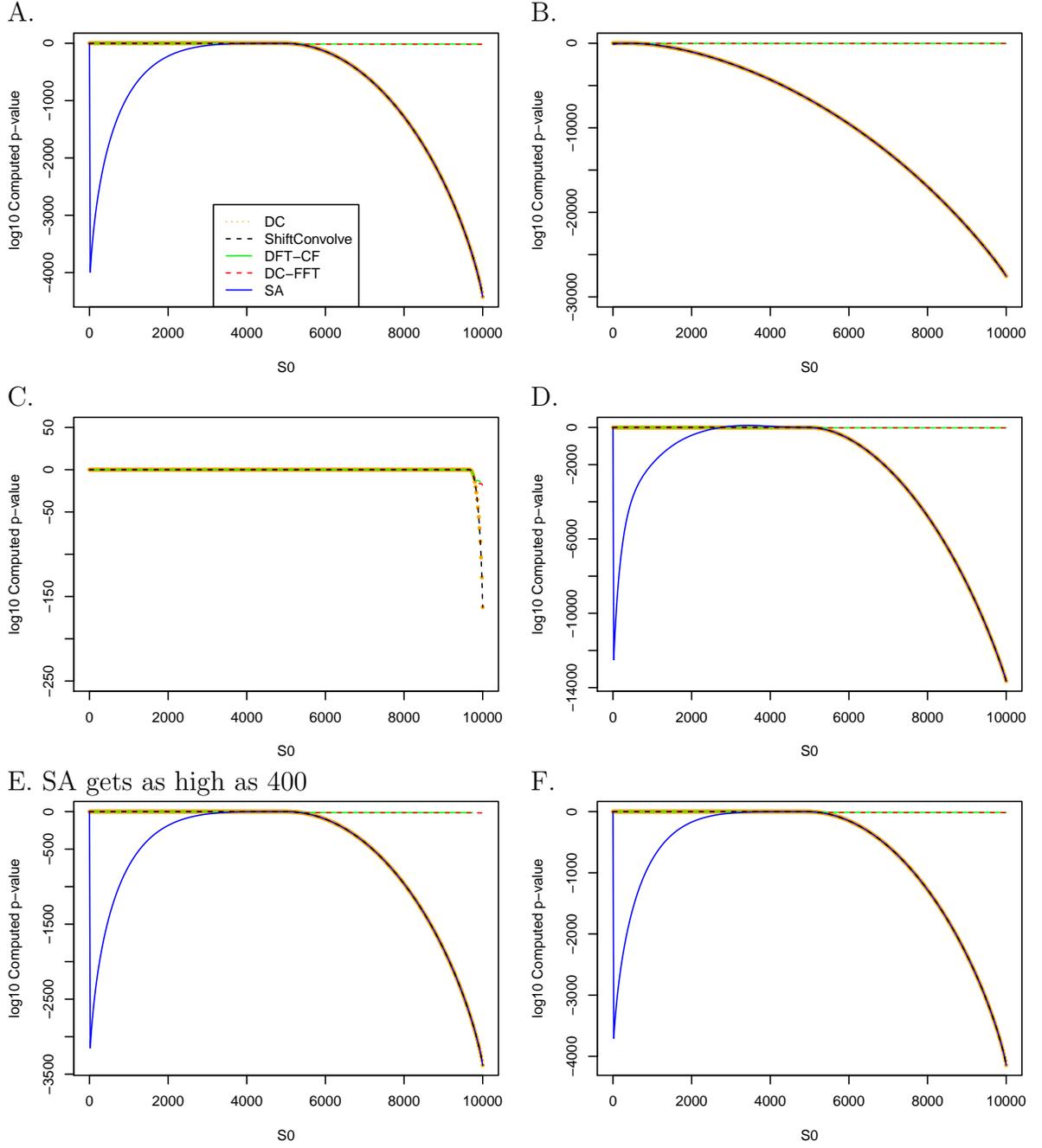


Figure 6: **Right tail probability (II)**. Similarly to Figure 3 the panels show the (log of the) computed right tail probability (p-value) for additional settings of the parameters. For the corresponding relative errors see Supplementary Figure 7. (A) $N = 10^4$, $p_i \sim U(0, 1)$. (B) $N = 10^4$, $p_i \sim \text{Beta}(0.1, 3)$. (C) $N = 10^4$, $p_i \sim \text{Beta}(3, 0.1)$. Note that in this example with higher proportions of 1s present (due to roundoff errors) the SA algorithm failed to run, returning an error where the blue curve is absent. (D) $N = 10^4$, $p_i \sim [0.5 \cdot \text{Beta}(3, 0.1) + 0.5 \cdot \text{Beta}(0.1, 3)]$ (E) $N = 10^4$, $p_i \sim \text{Beta}(3, 3)$. (F) $N = 10^4$, $p_i \sim [0.5 \cdot \text{Beta}(3, 10) + 0.5 \cdot \text{Beta}(10, 3)]$.

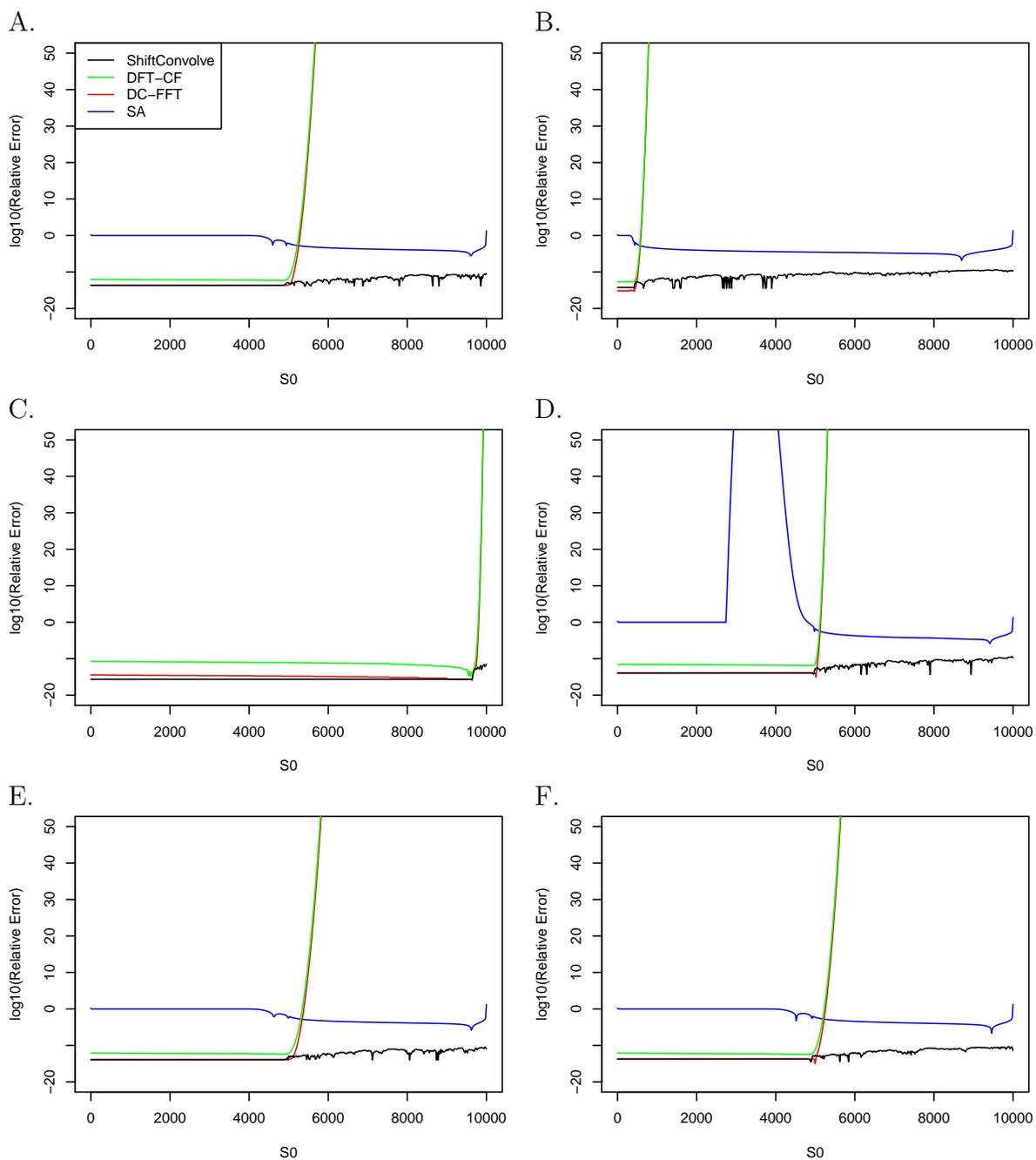


Figure 7: **Gauging the accuracy.** Similarly to Figure 4 the panels show the (log of the) relative error in computing the right tail for varying settings of the parameters. See Figure 6 for the right tail probability itself. (A) $N = 10^4$, $p_i \sim U(0,1)$. (B) $N = 10^4$, $p_i \sim \text{Beta}(0.1,3)$. (C) $N = 10^4$, $p_i \sim \text{Beta}(3,0.1)$. (D) $N = 10^4$, $p_i \sim [0.5 \cdot \text{Beta}(3,0.1) + 0.5 \cdot \text{Beta}(0.1,3)]$ (E) $N = 10^4$, $p_i \sim \text{Beta}(3,3)$. (F) $N = 10^4$, $p_i \sim [0.5 \cdot \text{Beta}(3,10) + 0.5 \cdot \text{Beta}(10,3)]$.

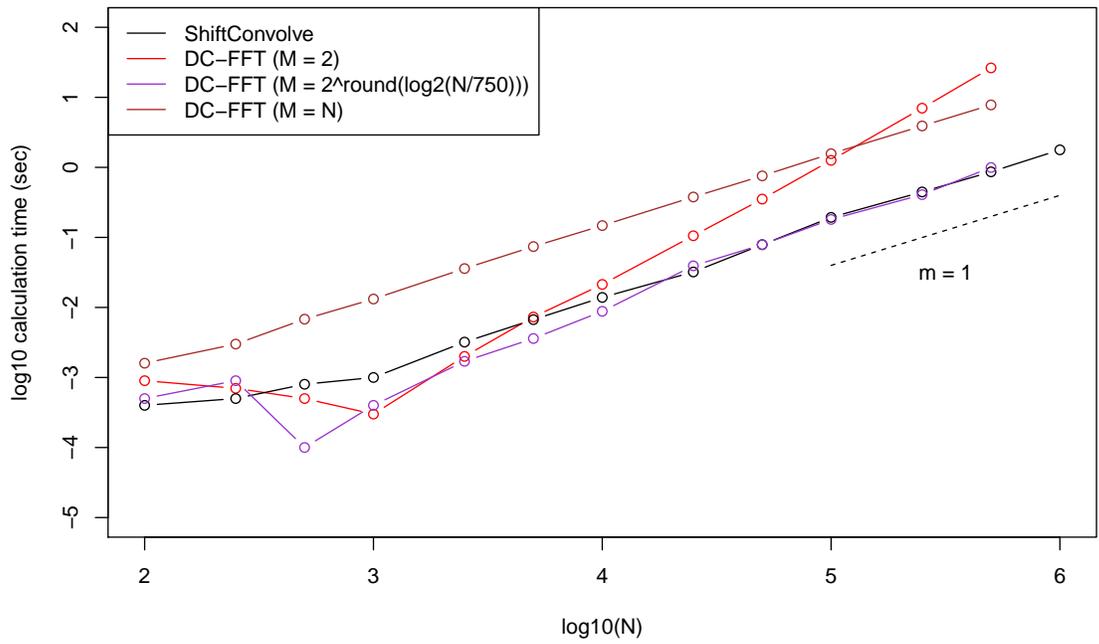


Figure 8: **Runtime: DC-FFT (varying M) and ShiftConvolve.** The runtime of DC-FFT using different settings of the parameter M and ShiftConvolve. Unfortunately we could not get DC-FFT to work for $N \geq 525 \cdot 10^3$. A line segment of slope 1 was added for reference.

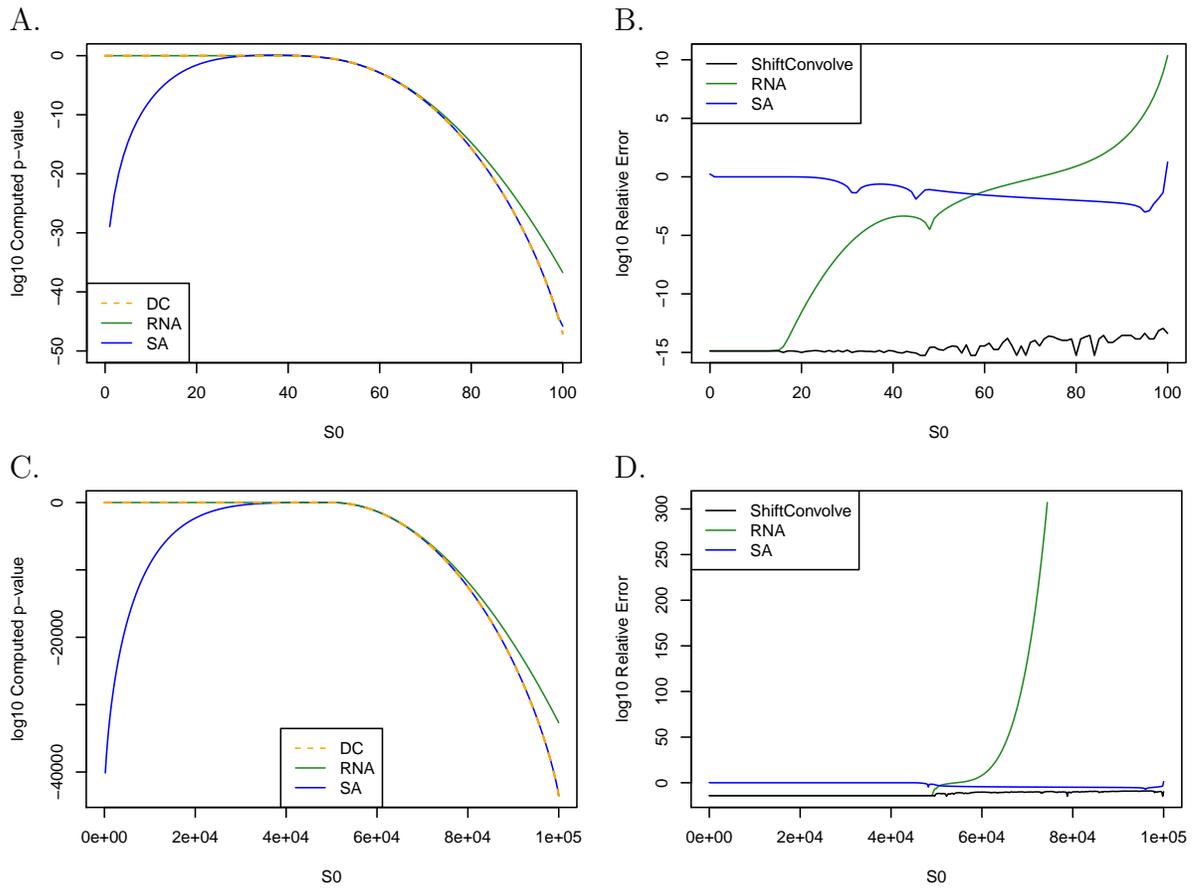


Figure 9: **Inaccuracy of RNA.** The panels show examples where RNA fails to accurately compute the right tail probability. (A) $N = 100$, p_i sampled uniformly. The computed right tail probability (p-value) reported by RNA compared with the accurate DC. (B) The corresponding (log base 10 of) the relative error where we added for reference the significantly smaller relative error of ShiftConvolve. (C) $N = 10^5$, p_i sampled uniformly. (D) The corresponding relative error.

References

- [1] William Biscarri, Sihai Dave Zhao, and Robert J Brunner. A simple and fast method for computing the poisson binomial distribution function. *Computational Statistics & Data Analysis*, 122:92–100, 2018.
- [2] Nicolas Brisebarre, Mioara Joldes, Jean-Michel Muller, Ana-Maria Naneş, and Joris Picot. Error analysis of some operations involved in the cooley-tukey fast fourier transform. *ACM Transactions on Mathematical Software*, pages 1–34, 2019.
- [3] J.W. Cooley and J.W. Tukey. An algorithm for the machine calculation of complex Fourier series. *Mathematics of Computation*, 19(90):297–301, 1965.
- [4] Matteo Frigo and Steven G. Johnson. The design and implementation of FFTW3. *Proceedings of the IEEE*, 93(2):216–231, 2005. Special issue on “Program Generation, Optimization, and Platform Adaptation”.
- [5] Yili Hong. On computing the distribution function for the poisson binomial distribution. *Computational Statistics & Data Analysis*, 59:41–51, 2013.
- [6] U. Keich and N. Nagarajan. A fast and numerically robust method for exact multinomial goodness-of-fit test. *Journal of Computational and Graphical Statistics*, 15(4):779–802, 2006.
- [7] Uri Keich. sFFT: a faster accurate computation of the p -value of the entropy score. *J Comput Biol*, 12(4):416–30, May 2005.

- [8] Tobias Madsen, Asger Hobolth, Jens Ledet Jensen, and Jakob Skou Pedersen. Significance evaluation in factor graphs. *BMC bioinformatics*, 18:199, March 2017.
- [9] Alexander Mukhin. *minFFT: A minimalist Fast Fourier Transform library.*, 2019.
- [10] Niranjan Nagarajan and Uri Keich. Reliability and efficiency of algorithms for computing the significance of the mann–whitney test. *Computational Statistics*, 24(4):605, 2009.
- [11] JC. Schatzman. Accuracy of the discrete Fourier transform and the fast Fourier transform. *SIAM J. Sci. Comput.*, 17(5):1150–1166, 1996.
- [12] T.G. Stockham, Jr. High-speed convolution and correlation. In *Proceedings of the April 26-28, 1966, Spring Joint Computer Conference, AFIPS '66 (Spring)*, pages 229–233, New York, NY, 1966. ACM.
- [13] Huon Wilson and Uri Keich. Accurate pairwise convolutions of non-negative vectors via fft. *Computational Statistics & Data Analysis*, 101:300–315, 2016.
- [14] Huon Wilson and Uri Keich. Accurate small tail probabilities of sums of iid lattice-valued random variables via fft. *Journal of Computational and Graphical Statistics*, 26(1):223–229, 2017.