

# Decompounding discrete distributions: A nonparametric Bayesian approach

Shota Gugushvili<sup>1</sup>  | Ester Mariucci<sup>2</sup>  | Frank van der Meulen<sup>3</sup> 

<sup>1</sup>Biometris, Wageningen University & Research,

<sup>2</sup>Institut für Mathematik, Potsdam Universität,

<sup>3</sup>Delft Institute of Applied Mathematics, Delft University of Technology,

## Correspondence

Ester Mariucci, Institut für Mathematik, Potsdam Universität,  
Karl-Liebknecht-Str. 24-25, D-14476  
Potsdam OT Golm, Germany.  
Email: mariucci@ovgu.de

## Funding information

Deutsche Forschungsgemeinschaft, CRC 1294 Data Assimilation, DFG 314838170, GRK 2297 MathCoRe; European Research Council, ERC Grant Agreement 320637

## Abstract

Suppose that a compound Poisson process is observed discretely in time and assume that its jump distribution is supported on the set of natural numbers. In this paper we propose a nonparametric Bayesian approach to estimate the intensity of the underlying Poisson process and the distribution of the jumps. We provide a Markov chain Monte Carlo scheme for obtaining samples from the posterior. We apply our method on both simulated and real data examples, and compare its performance with the frequentist plug-in estimator proposed by Buchmann and Grübel. On a theoretical side, we study the posterior from the frequentist point of view and prove that as the sample size  $n \rightarrow \infty$ , it contracts around the “true,” data-generating parameters at rate  $1/\sqrt{n}$ , up to a  $\log n$  factor.

## KEYWORDS

compound Poisson process, data augmentation, diophantine equation, Gibbs sampler, Metropolis-Hastings algorithm, Nonparametric Bayesian estimation

## 1 | INTRODUCTION

### 1.1 | Problem formulation

Let  $N = (N_t : t \geq 0)$  be a Poisson process with a constant intensity  $\lambda > 0$ , and let  $Y_i$  be a sequence of independent random variables, each with distribution  $P$ , that are also independent of  $N$ . By definition, a compound Poisson process (abbreviated CPP)  $X = (X_t : t \geq 0)$  is

-----  
This is an open access article under the terms of the Creative Commons Attribution License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

© 2019 The Authors. *Scandinavian Journal of Statistics* published by John Wiley & Sons Ltd on behalf of The Board of the Foundation of the Scandinavian Journal of Statistics.

$$X_t = \sum_{j=1}^{N_t} Y_j, \quad (1)$$

where here and below the sum over an empty index set is understood to be equal to zero. In particular,  $X_0 = 0$ . CPP constitutes a classical model in, for example, risk theory, see Embrechts, Mikosch, and Klüppelberg (1997).

Assume that the process  $X$  is observed at discrete times  $0 < t_1 < t_2 < \dots < t_n = T$ , where the instants  $t_i$  are not necessarily equidistant on  $[0, T]$ . Based on the observations  $X_{t_1}, X_{t_2}, \dots, X_{t_n}$ , our goal is to estimate the jump size distribution  $P$  and the intensity  $\lambda$ . We specifically restrict our attention to the case where  $P$  is a discrete distribution,  $P(\mathbb{N}) = 1$ , and we will write  $p = (p_k)_{k \in \mathbb{N}}$  for the probability mass function corresponding to  $P$ , where  $p_k = P(\{k\})$ . A similar notation will be used for any other discrete law. The distribution  $P$  is called the base distribution. Abusing terminology, we will at times identify it with the corresponding probability mass function  $p$ . An assumption that  $P$  has no atom at zero is made for identifiability: otherwise this atom gets confounded with  $e^{-\lambda}$ , which does not allow consistent estimation of the intensity  $\lambda$ . For a discussion of applications of this CPP model in risk theory, see Zhang, Liu, and Li (2014).

Define the increments  $Z_i = X_{t_i} - X_{t_{i-1}}$ ,  $i = 1, \dots, n$ . Then  $\mathcal{Z}_n = (Z_i : i = 1, \dots, n)$  is a sequence of independent random variables. When  $\{t_i\}$  are equidistant on  $[0, T]$ , the random variables  $Z_i$  have in fact a common distribution  $Q$  satisfying  $Q(\mathbb{N}_0) = 1$ . As  $\mathcal{Z}_n$  carries as much information as  $(X_{t_i} : i = 1, \dots, n)$  does, we can base our estimation procedure directly on the increments  $\mathcal{Z}_n$ . Since summing up the jumps  $Y_j$ s amounts to compounding their distributions, the inverse problem of recovering  $P$  and  $\lambda$  from  $Z_i$  can be referred to as decompounding (Buchmann & Grübel, 2003).

There are two natural ways to parameterize the CPP model: either in terms of the pair  $(\lambda, p)$ , or in terms of the Lévy measure  $\nu = (\nu_k)_{k \in \mathbb{N}}$  of the process  $X$  (Sato, 2013). A relationship between the two is  $\lambda = \sum_{k=1}^{\infty} \nu_k$  and  $p = \nu/\lambda$ . Inferential conclusions in one parameterization can be easily translated into inferential conclusions into another parameterization. However, for our specific statistical approach the Lévy measure parameterization turns out to be more advantageous from the computational point of view.

## 1.2 | Approach and results

In this paper, we take a nonparametric Bayesian approach to estimation of the Lévy measure  $\nu$  of  $X$ . See Ghosal & van der Vaart (2017) and Müller, Quintana, Jara, and Hanson (2015) for modern expositions of Bayesian nonparametrics. A case for nonparametric Bayesian methods has already been made elsewhere in the literature, and will not be repeated here. On the practical side, we implement our procedure via the Gibbs sampler and data augmentation, and show that it performs well under various simulation setups. On the theoretical side, we establish its consistency and derive the corresponding posterior contraction rate, which can be thought of as an analogue of a convergence rate of a frequentist estimator (Ghosal and van der Vaart, (2017)). The posterior contraction rate, up to a practically insignificant  $\log n$  factor, turns out to be  $1/\sqrt{n}$ , which is an optimal rate for nonparametric estimation of cumulative distribution functions. Our contribution thus nicely bridges practical and theoretical aspects of Bayesian nonparametrics.

### 1.3 | Related literature

To provide a better motivation for our model and approach, in this subsection we briefly survey the existing literature. A Bayesian approach to nonparametric inference for Lévy processes is a very recent and emerging topic, with references limited at the moment to Belomestny, Gugushvili, Schauer, and Spreij (2019), Gugushvili, van der Meulen, and Spreij (2015), Gugushvili, van der Meulen, and Spreij (2018) and Nickl and Söhl (2017). These deal exclusively with the case when the Lévy measure is absolutely continuous with respect to the Lebesgue density. At least from the computational point of view, these works are of no help in our present context.

Related frequentist papers for CPP models with discrete base distributions are Buchmann and Grübel (2003) and Buchmann and Grübel (2004), which, after earlier contributions dating from the previous century, in fact revived interest in nonparametric techniques for Lévy processes. To estimate the base distribution  $p$ , Buchmann and Grübel (2003) employ a frequentist plug-in approach relying on the Panjer recursion (i.e., an empirical cumulative distribution estimate of  $q$  is plugged into the Panjer recursion equations to yield an estimate of  $p$ ; see below on the Panjer recursion). The drawback is that the parameter estimates are not guaranteed to be nonnegative. Buchmann and Grübel (2004) fix this problem by truncation and renormalization. This works, but looks artificial. As noted in Buchmann and Grübel (2004), in practice the latter approach breaks down if no zero values are observed among  $Z_i$ s. Buchmann and Grübel (2004) establish weak convergence of their modified estimator, but on the downside its asymptotic distribution is unwieldy to give confidence statements on  $p$ . Most importantly, the plug-in approaches in Buchmann and Grübel (2003) and Buchmann and Grübel (2004) do not allow obvious generalisations to nonequidistant observation times  $\{t_i\}$ . In Lindo, Zuyev, and Sagitov (2018), another frequentist estimator of the jump measure is introduced, that is obtained via the steepest descent technique as a solution to an optimization problem over the cone of positive measures. The emphasis in Lindo et al. (2018) is on numerical aspects; again, no obvious generalization to the case of nonequidistant  $\{t_i\}$  is available.

Finally, some important, predominantly theoretical references on inference for Lévy processes are Comte and Genon-Catalot (2011), Duval and Hoffmann (2011), Kappus (2014), Neumann and Reiß (2009), Nickl and Reiß (2012), van Es, Gugushvili, and Spreij (2007) and Trabs (2015). We refer to Belomestny, Comte, Genon-Catalot, Masuda, and Reiß (2015), Coca (2018a, 2018b), and Duval and Mariucci (2017) for more extensive literature surveys.

### 1.4 | Outline

The rest of the paper is organized as follows: in Section 2 we introduce our approach and describe an algorithm for drawing from the posterior distribution. In Sections 3 and 4 we study its performance on synthetic and real examples. Section 5 is devoted to the examination of asymptotic frequentist properties of our procedure. An outlook on our results is given in Section 6. Finally, in Appendix A technical lemmas used in the proofs of Section 5 are collected, whereas Appendix B contains some additional simulation results.

### 1.5 | Notation

For two sequences  $\{a_n\}$  and  $\{b_n\}$  of positive real numbers, the notation  $a_n \lesssim b_n$  (or  $b_n \gtrsim a_n$ ) means that there exists a constant  $C > 0$  that is independent of  $n$  and such that  $a_n \leq Cb_n$ . We write  $a_n \asymp b_n$  if both  $a_n \lesssim b_n$  and  $a_n \gtrsim b_n$  hold. We denote a prior (possibly depending on the sample

size  $n$ ) by  $\Pi_n$ . The corresponding posterior measure is denoted by  $\Pi_n(\cdot | \mathcal{Z}_n)$ . The Gamma distribution with shape parameter  $a$  and rate parameter  $b$  ( $a, b > 0$ ) is denoted by  $\text{Gamma}(a, b)$ . Its density is given by  $x \mapsto \frac{b^a}{\Gamma(a)} x^{a-1} e^{-bx}$ ,  $x > 0$ , where  $\Gamma$  is the Gamma function. The inverse Gamma distribution with shape parameter  $a$  and scale parameter  $b$  is denoted by  $\text{IG}(a, b)$ . Its density is  $x \mapsto \frac{b^a}{\Gamma(a)} x^{-a-1} e^{-b/x}$ ,  $x > 0$ . We use the notation  $\text{Exp}(a)$  for an exponential distribution with mean  $1/a$ . Finally, given a metric  $d$  on a set  $\mathcal{Q}$  and  $\epsilon > 0$ , the covering number  $N(\epsilon, \mathcal{Q}, d)$  is defined as the minimal number of balls of radius  $\epsilon$  needed to cover  $\mathcal{Q}$ .

## 2 | ALGORITHM FOR DRAWING FROM THE POSTERIOR

A Bayesian statistical approach relies on the combination of the likelihood and the prior on the parameter of interest through Bayes' formula. We start with specifying the prior. As far as the likelihood is concerned, although explicit, it is intractable from the computational point of view for nonparametric inference in CPP models. We will circumvent the latter problem by means of data augmentation, as detailed below.

### 2.1 | Prior

We define a prior  $\Pi$  on  $\nu$  through a hierarchical specification

$$\begin{aligned} \{\nu_k\}_{k=1}^{\infty} | a, m, \beta_k &\stackrel{\text{i.i.d.}}{\sim} \text{Gamma}(a, 1/\beta_k) \cdot \mathbf{1}_{\{1 \leq k \leq m\}}, \\ \beta_1, \dots, \beta_m | \gamma &\stackrel{\text{i.i.d.}}{\sim} \text{IG}(c, \gamma), \\ \gamma &\sim \text{Exp}(1). \end{aligned}$$

Note that the (fixed) hyperparameters  $m \in \mathbb{N}$ ,  $a, c > 0$  are denoted by Latin letters.

The hyperparameter  $m$  incorporates our a priori opinion on the support of the Lévy measure  $\nu$ , or equivalently, the base measure  $p$ . In applications, the support of  $p$  may be unknown, which necessitates the use of a large  $m$ , for example,  $m = \max_{i=1, \dots, n} Z_i$ ; this latter is the maximal value suggested by the data  $\mathcal{Z}_n$  at hand. Nevertheless, we may simultaneously expect that the “true,” data-generating  $\nu$  charges full mass only to a proper, perhaps even a small subset of the set  $\{1, \dots, m\}$ . In other words,  $\nu$  may form a sparse sequence, with many components equal to zero. In fact, there are at least two plausible explanations for an occurrence of a large increment  $Z_i$  in the data: either a few large jumps  $Y_j$ 's occurred, which points toward a large right endpoint of the support of  $\nu_0$ , or  $Z_i$  is predominantly formed of many small jumps, which in turn indicates that the intensity  $\lambda$  of the Poisson arrival process  $N$  may be large. To achieve accurate estimation results, a prior should take a possible sparsity of  $\nu$  into account. This is precisely the reason of our hierarchical definition of the prior  $\Pi$ : a small  $\beta_k$  encourages a priori the shrinkage of the components  $\nu_k$  of  $\nu$  toward zero.

### 2.2 | Data augmentation

Assume temporarily  $t_i = i$ ,  $i = 1, \dots, n$ , and write  $q = (q_k)_{k \in \mathbb{N}_0}$  for  $q_k = q(\{k\})$ . Then  $Z_i$  have the distribution

$$q = e^{-\lambda} \sum_{j=0}^{\infty} \frac{\lambda^j}{j!} p^{*j}, \quad (2)$$

with  $*$  denoting convolution. The compounding mapping  $(\lambda, p) \mapsto q$  can be expressed explicitly via the Panjer recursion (see Panjer (1981)):

$$q_0 = e^{-\lambda}, \quad q_k = \frac{\lambda}{k} \sum_{j=1}^k j p_j q_{k-j}, \quad k \in \mathbb{N}.$$

This recursion can be inverted to give the inverse mapping  $q \mapsto (\lambda, p)$  via

$$\lambda = -\log q_0, \quad p_k = -\frac{q_k}{q_0 \log q_0} - \frac{1}{k q_0} \sum_{j=1}^{k-1} j p_j q_{k-j}, \quad k \in \mathbb{N}.$$

In view of (2), the likelihood in the CPP model is explicit. Nevertheless, an attempt to directly use (2) or the Panjer recursions in posterior computations results in a numerically intractable procedure. Equally important is the fact that a Panjer recursion-based approach would not apply to nonequidistant observation times  $\{t_i\}$ . Therefore, instead of (2) and the Panjer recursion, we will employ data augmentation (Tanner & Wong, 1987). We switch back to the case when values of  $\{t_i\}$  are not necessarily uniformly spaced. The details of our procedure are as follows: when the process  $X$  is observed continuously over the time interval  $[0, T]$ , so that our observations are a full sample path  $X^{(T)} = (X_t : t \in [0, T])$  of CPP, the likelihood is tractable and is proportional to

$$e^{-T \sum_{k=1}^m \nu_k} \prod_{k=1}^m \nu_k^{\mu_k},$$

see Shreve (2004). Here

$$\mu_k = \#\{Y_j = k\}, \quad (3)$$

that is, the total number of jumps of size  $k$ . Then the prior  $\Pi$  from Subsection 2.1 leads to conjugate posterior computations. In fact, the full conditionals are

$$\begin{aligned} \nu_k | \{\mu_k\}, \{\beta_k\} &\stackrel{\text{i.i.d.}}{\sim} \text{Gamma}(a + \mu_k, 1/\beta_k + T), \quad k = 1, \dots, m, \\ \beta_k | \{\nu_k\}, \gamma &\stackrel{\text{i.i.d.}}{\sim} \text{IG}(a + c, \gamma + \nu_k), \quad k = 1, \dots, m, \\ \gamma | \{\beta_k\} &\sim \text{Gamma}\left(cm + 1, 1 + \sum_{k=1}^m \beta_k^{-1}\right). \end{aligned}$$

Therefore, the Gibbs sampler for posterior inference on  $\nu$  can be implemented. The Gibbs sampler cycles through the above conditionals a large number of times, generating approximate (dependent) samples from the posterior. See, e.g., Gelfand and Smith (1990) and section 24.5 in Wasserman (2004) on the Gibbs sampler and its use in Bayesian statistics.

As we do not observe the process  $X$  continuously, we will combine the above with the data augmentation device. First note that we have

$$Z_i = \sum_{j=1}^m j \mu_{ij},$$

where  $(\mu_{ij} : i = 1, \dots, n, j = 1, \dots, m)$  are independent, and  $\mu_{ij} \sim \text{Poisson}(\Delta_i v_j)$  for  $v_j = \lambda p_j$  and  $\Delta_i = t_i - t_{i-1}$ ; see Corollary 11.3.4 in Shreve (2004, p. 498). Furthermore, for  $\mu_k$  as in (3) we trivially have  $\mu_k = \sum_{i=1}^n \mu_{ik}$ . Data augmentation iterates the following two steps:

- (i) Draw  $(\mu_{ij})$  conditional on the data  $\mathcal{Z}_n$  and the parameter  $v$ .
- (ii) Draw  $v$  conditional on  $(\mu_{ij})$ .

Once the algorithm has been run long enough, this gives approximate (dependent) samples from the posterior of  $v$ . We already know how to deal with step (ii); now we need to handle step (i).

Thus, keeping  $v$  fixed, for each  $i$  we want to compute the conditional distribution  $(\mu_{ij} : j = 1, \dots, m) | Z_i$ , and furthermore, we want to be able to simulate from this distribution. In turn, this will immediately allow us to simulate  $\mu_k$  conditional on the data  $\mathcal{Z}_n$ . Now, with  $\text{Pr}(\cdot)$  referring to probability under the parameter  $v$ , it holds that

$$\text{Pr}(\mu_{i1} = k_1, \dots, \mu_{im} = k_m | Z_i = z_i) = \frac{1}{\text{Pr}(Z_i = z_i)} e^{-\Delta_i \sum_{j=1}^m v_j} \prod_{j=1}^m \frac{(\Delta_i v_j)^{k_j}}{k_j!} \mathbf{1} \left\{ \sum_{j=1}^m j k_j = z_i \right\}.$$

Knowledge of the normalizing constant  $\text{Pr}(Z_i = z_i)$  will not be needed in our approach.

In general, simulation from a discrete multivariate distribution is nontrivial; some general options are discussed in Devroye (1986, Chapter XI, Section 1.5), but are unlikely to work easily for a large  $m$ . We will take an alternative route and use the Metropolis-Hastings algorithm, see, for example, section 24.4 in Wasserman (2004). We start by observing that for a fixed  $i$ , the support of  $\text{Pr}(\cdot | Z_i = z_i)$  is precisely the set  $S_i$  of nonnegative solutions  $(k_1, \dots, k_m)$  of the Diophantine equation  $\sum_{j=1}^m j k_j = z_i$ . The **R** package **nilde** (see Pya Arnqvist, Voinov, & Voinov, 2018) implements an algorithm from Voinov and Nikulin (1997) that finds all such solutions for given integers  $m$  and  $z_i$ . By Markovianity of the process  $X$ , we can simulate the vectors  $(\mu_{i1}, \dots, \mu_{im})$  independently for each  $i = 1, \dots, n$ . If  $z_i = 0$  or 1, there is only one solution to the Diophantine equation: the trivial solution  $(0, \dots, 0)$  in the first case, and the solution  $(1, 0, \dots, 0)$  in the second case; for such  $z_i$ , no simulation is required, as  $(\mu_{i1}, \dots, \mu_{im})$  is known explicitly. We thus only need to consider each  $i \in \mathcal{I} = \{i : z_i \neq 0 \text{ or } 1\}$  in turn, and design a Metropolis-Hastings move on the set of the corresponding solutions  $S_i$ . Fix once and for all an ordering of elements in  $S_i$  (this could be, e.g., lexicographic, or reverse lexicographic); we use the notation  $|S_i|$  for the cardinality of  $S_i$ . Let  $\mu = (\mu_{i1}, \dots, \mu_{im})$  be the current state of the chain, corresponding to the  $\ell$ th element  $s_\ell$  of  $S_i$ . A proposal  $\mu^\circ = (\mu_{i1}^\circ, \dots, \mu_{im}^\circ)$  is obtained as follows:

- (i) If  $\ell = 1$ , draw  $\mu^\circ$  uniformly at random among the elements  $\{s_2, s_{|S_i|}\}$  of  $S_i$ .
- (ii) If  $\ell = |S_i|$ , draw  $\mu^\circ$  uniformly at random among the elements  $\{s_1, s_{|S_i|-1}\}$  of  $S_i$ .
- (iii) If  $\ell \neq 1$  or  $|S_i|$ , draw  $\mu^\circ$  uniformly at random among the elements  $\{s_{\ell-1}, s_{\ell+1}\}$  of  $S_i$ .

Occasionally, one may want to propose another type of a move too.

(iv) Draw  $\mu^\circ = (\mu_{i_1}^\circ, \dots, \mu_{i_m}^\circ)$  uniformly at random from  $S_i$ .

The two proposals lead to reversible moves, and one may also alternate them with probabilities  $\pi$  and  $1 - \pi$ , for example,  $\pi = 0.8$ . The logarithm of the acceptance probability of a move from  $(\mu_{i_1}, \dots, \mu_{i_m})$  to  $(\mu_{i_1}^\circ, \dots, \mu_{i_m}^\circ)$  is computed as

$$\log A = \sum_{k=1}^m (\mu_{i_k}^\circ - \mu_{i_k}) \log(\Delta_i \nu_k) + \sum_{k=1}^m \{\log(\mu_{i_k}!) - \log(\mu_{i_k}^\circ!)\}.$$

The move is accepted if  $\log U \leq \log A$  for  $U$  an independently generated uniform random variate on  $[0, 1]$ , and in that case the current state of the chain is reset to  $(\mu_{i_1}^\circ, \dots, \mu_{i_m}^\circ)$ . Otherwise the chain stays in  $(\mu_{i_1}, \dots, \mu_{i_m})$ .

### 3 | SIMULATION EXAMPLES

In this section, we test performance of our approach in a range of representative simulation examples. For benchmarking, we use the frequentist plug-in estimator from Buchmann and Grübel (2004). Two real data examples are given in Section 4. Unless otherwise stated, we took  $c = 2$  and  $a = 0.01$  as hyperparameters in our prior specification. As can be seen from the update formulae for the Gibbs sampler, as long as  $a$  is not taken too large, its precise value is not very influential on the posterior, given a reasonable sample size. The value  $c = 2$  ensures that the update step for  $\beta_k$  has finite variance. At each step of updating the imputed data for increment size  $z$  we have chosen with probability 0.2 to propose uniformly from all solutions to the Diophantine equation (for that particular value of  $z$ ).

We implemented our procedure in Julia, see (Bezanson, Edelman, Karpinski, & Shah, 2017). The code and datasets for replication of our examples are available on GitHub<sup>1</sup> and Zenodo (Gugushvili, Mariucci, & van der Meulen, 2019).

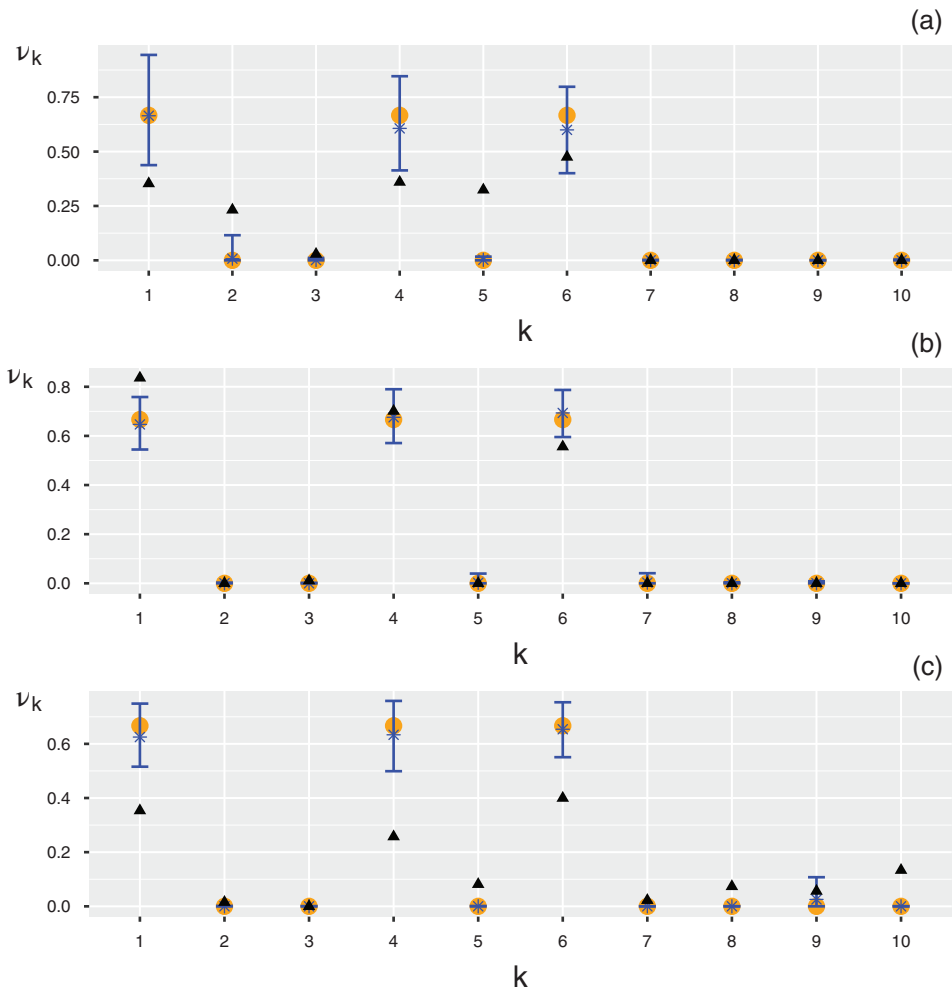
#### 3.1 | Uniform base distribution

This simulation example follows with some extensions that in Buchmann and Grübel (2004). Let  $\lambda_0 = 2$ , and let  $P_0$  be the discrete uniform distribution on  $\{1, 4, 6\}$ . We simulated data according to the following settings:

- (a)  $n = 100$ ,  $\Delta_i = 1$  for  $1 \leq i \leq n$ ;
- (b)  $n = 500$ ,  $\Delta_i = 1$  for  $1 \leq i \leq n$  (the data under (a) are augmented with 400 extra observations);
- (c)  $n = 500$ ,  $\Delta_i = \text{Unif}(0, 2)$  for  $1 \leq i \leq n$ .

We set  $m = \min(15, Z_{(n)})$ , where  $Z_{(n)} = \max_{1 \leq i \leq n} Z_i$ . In all cases this led to  $m = 15$ , as the value of  $Z_{(n)}$  was equal to 30, 35, and 40 for the simulated data under settings (a), (b), and (c), respectively. The Gibbs sampler was run for 500, 000 iterations, of which the first 250, 000 were discarded

<sup>1</sup>See <https://github.com/fmeulen/Bdd>

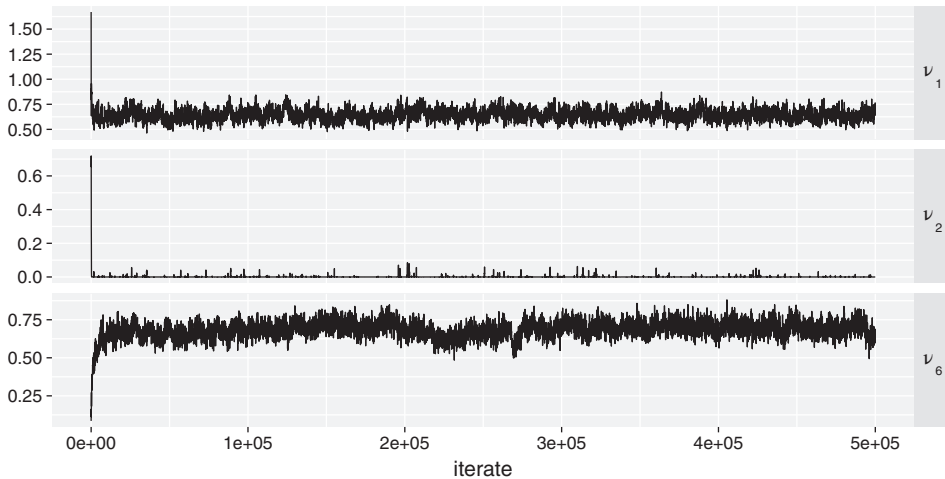


**FIGURE 1** Simulation example from Section 3.1. In each figure, the horizontal axis gives the magnitudes of  $v_k$ ,  $k \in \{1, \dots, 10\}$ . The orange balls denote the true values, the black triangles the Buchmann-Grübel estimator. The blue crosses give the posterior means, whereas the vertical blue line segments represent (pointwise) 95% credible intervals. The settings corresponding to (a), (b), and (c) are explained in the main text. Note the differences in vertical scale across the figures [Colour figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]

as burn-in. From the remaining samples, the posterior mean and 2.5% and 97.5% percentiles were computed for each coefficient  $v_k$ . The results for the first 10 coefficients are shown in Figure 1. For comparison, the estimator from Buchmann and Grübel (2004) is also included in the figure.

For setting (b), traceplots of every 50th iteration for a couple of coefficients  $v_k$  are shown in Figure 2. We measure the error of an estimate  $\{\hat{v}_k\}$  by  $\text{Err}(v, \hat{v}) = \sum_{k=1}^{\infty} |\hat{v}_k - v_k|$ . The errors are reported in Table 1. In all settings, for these particular realizations of the simulated data, the Bayesian procedure outperformed the truncated estimator from Buchmann and Grübel (2004). For setting (c), the latter produces a poor result, as was to be expected, given that it is derived under the assumption  $\Delta_i = 1$  for all  $i$ . An advantage of the Bayesian procedure is the included measure of uncertainty, namely the credible intervals for  $v_k$ . On the other hand, for the Buchmann-Grübel estimator it is hardly possible to derive confidence intervals via an asymptotic method, since the limiting distribution of the estimator is fairly complicated. Although not considered in the original





**FIGURE 2** Traceplots for the simulation example from Section 3.1 under setting (b). The posterior samples were subsampled, with every 50th iteration kept. The displayed results are for parameters  $v_1$ ,  $v_6$ , and  $v_9$

**TABLE 1** Results for scenarios (a)–(c) from Section 3.1

Simulation setting	(a)	(b)	(c)
Buchmann-Grübel estimator	1.40	0.32	1.44
Posterior mean	0.15	0.07	0.12

publications Buchmann and Grübel (2003) and Buchmann and Grübel (2004), a natural alternative is the bootstrap. A detailed examination of the performance of the latter and its comparison to that of the Bayesian method lies beyond the scope of the present paper. Indeed, any thorough study would require, on one hand, the asymptotic justification of bootstrap confidence intervals, and on another hand establishing frequentist coverage properties of our Bayesian procedure. In that respect, good performance of neither method is automatically warranted (e.g., van der Pas, Szabó, and van der Vaart (2017) and van der Vaart (1998), Chapter 23). Here instead we opt for a numerical illustration, which is reported in Appendix B.

### 3.2 | Geometric base distribution

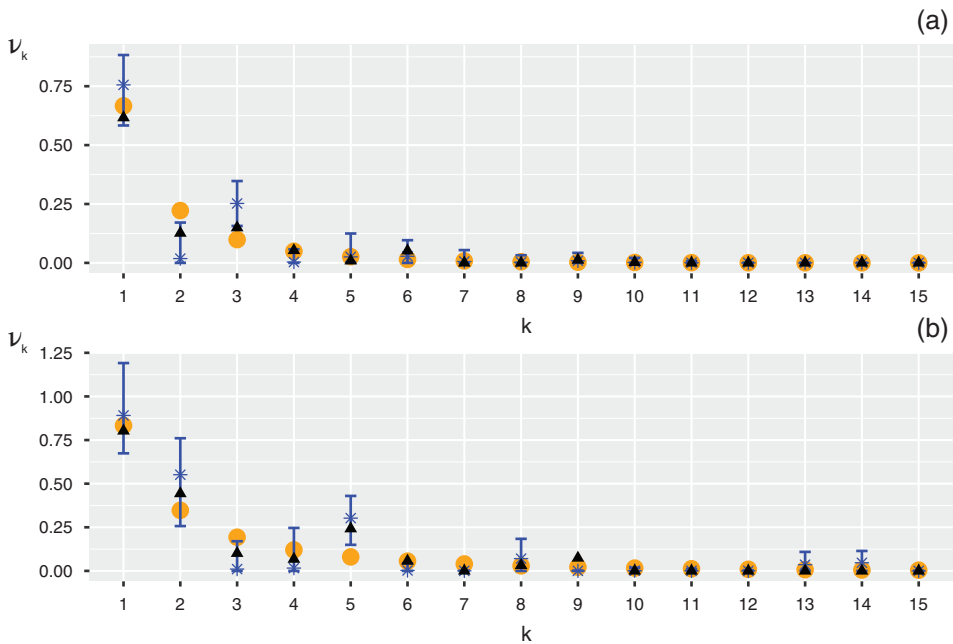
The setup of this synthetic data example likewise follows that in Buchmann and Grübel (2004). Assume  $q$  is a geometric distribution with parameter  $\alpha$ , that is,  $q_k = (1 - \alpha)^k \alpha$  for  $0 < \alpha < 1$ ,  $k \in \mathbb{N}_0$ . Then  $\lambda = -\log \alpha$ , and

$$p_k = -\frac{(1 - \alpha)^k}{k \log \alpha}, \quad k \in \mathbb{N}.$$

Hence,  $v_k = (1 - \alpha)^k / k$ .

We consider two simulation setups:

- (a)  $n = 500$ ,  $\Delta_i = 1$  for  $1 \leq i \leq n$  and  $\alpha = 1/3$ ;
- (b)  $n = 500$ ,  $\Delta_i = 1$  for  $1 \leq i \leq n$  and  $\alpha = 1/6$ .



**FIGURE 3** Simulation example from Section 3.2. Settings (a) and (b) correspond to the true jump distributions Geom(1/3) and Geom(1/6), respectively. The horizontal axis gives the magnitudes of  $\nu_k$ ,  $k \in \{1, \dots, 15\}$ . The orange balls denote the true values, the black triangles the Buchmann-Grübel estimator. The blue crosses give the posterior means, whereas the vertical blue line segments represent (pointwise) 95% credible intervals [Colour figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]

**TABLE 2** Results for scenarios (a)–(b) from Section 3.2

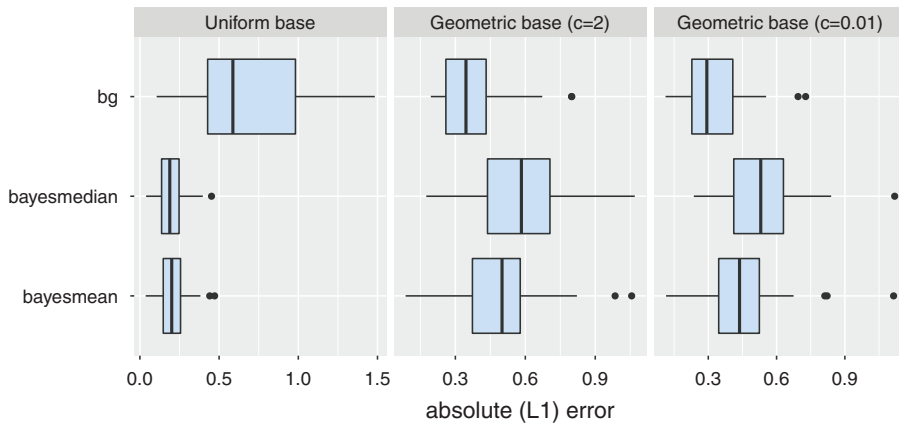
Simulation setting	(a)	(b)
Buchmann-Grübel estimator	0.28	0.60
Posterior mean	0.52	1.05

We set  $m = \min(15, Z_{(n)})$  and ran the sampler according to the settings of Section 3.1. The results for both scenarios (a) and (b) are reported in Figure 3. In Table 2 we also report estimation errors in one simulation run. For this example and these generated data, the Bayesian procedure gives less precise point estimates than the Buchmann-Grübel method. Note that estimation error for  $\alpha = 1/3$  is smaller than that for  $\alpha = 1/6$ . This appears intuitive, as a smaller value of  $\alpha$  corresponds to a larger value of  $\lambda$ . The latter implies that on average each  $Z_i$  is a superposition of a larger number of jumps, which renders the decomposing problem more difficult. However, this argument is hard to formalise.

### 3.3 | Monte Carlo study

For a more thorough comparison of the Buchmann-Grübel estimator and our Bayesian method, we performed a small Monte Carlo experiment. We considered two settings:

- (i) The setting from Section 3.1 with  $n = 250$ . We took  $m = \min(15, Z_{(n)})$ .
- (ii) The setting from Section 3.2 with  $\alpha = 1/3$ . We took  $m = \min(20, Z_{(n)})$ .



**FIGURE 4** Monte Carlo study from Subsection 3.3 comparing the Buchmann-Grübel estimator and the Bayesian method proposed in this paper. In this figure “bg” refers to the Buchmann-Grübel estimator, while “bayesmedian” and “bayesmean” refer to the Bayesian method, where either the median or mean was used as a point estimator for each  $v_i$ . The leftmost panel corresponds to the setting (i), whereas the other two panels to the setting (ii). In the latter we used both  $c = 2$  and  $c = 0.01$  in the prior specification [Colour figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]

In both cases we assumed  $\Delta_i = 1$  for all  $1 \leq i \leq n$ . The number of Monte Carlo repetitions was taken equal to 50. We took 400,000 Markov chain Monte Carlo (MCMC) iterations and discarded the first half of these as burn-in. In Figure 4 we give a graphical display of the results by means of boxplots of the errors. Here, as before, if the true values are denoted by  $v_k$  and the estimate within a particular simulation run by  $\hat{v}_k$ , the error is defined by  $\text{Err}(v, \hat{v}) = \sum_{k=1}^{\infty} |\hat{v}_k - v_k|$  (we truncated the infinite summation to 50). The results agree with our earlier findings, in that there is no clear “winner” in the comparison. Note that for the setting (ii) we considered both  $c = 2$  and  $c = 0.01$  in the prior specification. Both values give similar performance of the Bayesian method. This provides insight into sensitivity of our results with respect to the prior specification. A minor difference between the middle and rightmost panel of Figure 4 may be attributed to Monte Carlo error: the 50 simulated datasets on which these panels are based are not the same. Note that the prior promotes sparsity, and in that respect it is not surprising it does better when the true data-generating Lévy measure is sparse.

### 3.4 | Computing time

In terms of computational effort, the time it takes to evaluate the Buchmann-Grübel estimator is negligible compared to our algorithm for sampling from the posterior. This is not surprising, as that frequentist estimator relies on a plug-in approach, whereas in our case an approximation to the posterior is obtained by MCMC simulation. However, if proper uncertainty quantification is desired, then the Bayesian method is advantageous in the sense that it does not solely produce a point estimate.

Note that the proposed MCMC scheme requires determination of the solutions to the Diophantine equation  $\sum_{j=1}^m jk_j = z$  for all unique values  $z$  in the observation set. For moderate values of  $z$ , say  $z \leq 30$ , this is rather quick, but for large values of  $z$  the computing time increases exponentially, as does the amount of the allocated memory. The computing time of each Metropolis-Hastings step is then very small, but we potentially need a very large number of iterations to reach stationarity. The latter is caused firstly by the fact that at a particular iteration, our

proposals for  $\mu_{ij}$  do not take into account the current values of  $v_1, \dots, v_m$ ; secondly, the size of the state space that needs to be explored increases exponentially with  $m$ .

## 4 | REAL DATA EXAMPLES

### 4.1 | Horse kick data

To further illustrate our procedure, we will use the von Bortkewitsch data on the number of soldiers in the Prussian cavalry killed by horse kicks (available by year and by cavalry corps); this example was also employed in Buchmann and Grübel (2003). Each observation is an integer from 0 to 4, giving the number of deaths for a given year and a given cavalry corps, with overall counts reported in Table 3. The data are extracted from the table on p. 25 in von Bortkewitsch (1898). Note that von Bortkewitsch corrects for the fact that the Guards and I, VI, and XI cavalry corps of the Prussian army had a different organization from other units, and justifiably omits the corresponding counts from consideration.

It has been demonstrated by von Bortkewitsch that the Poisson distribution fits the horse kick data remarkably well. Assuming instead that observations follow a compound Poisson distribution is a stretch of imagination, as that would correspond to a horse running amok and killing possibly more than one soldier in one go. Nevertheless, this example provides a good sanity check for our statistical procedure.

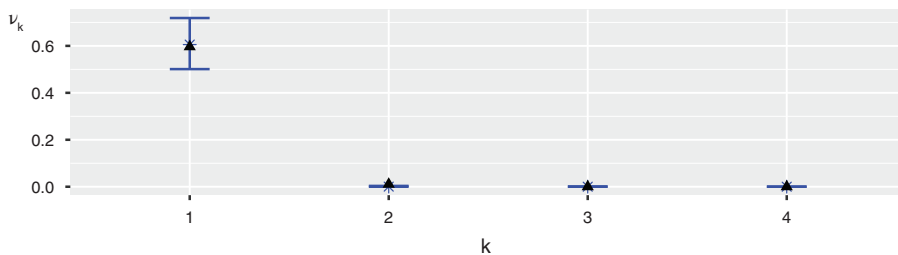
The estimation results are graphically depicted in Figure 5. Clearly, point estimates of both methods are in agreement and lend support to the Poisson model for this dataset.

### 4.2 | Plant data

Our second real example is the one used in Buchmann and Grübel (2004). Consider the data in Table 4, taken from Evans (1953). The data were collected as follows: the area was divided into plots of equal size and in each plot the number of plants was counted; the number of plants in each

**TABLE 3** Data on the number of soldiers in the Prussian cavalry killed by horse kicks. See the table on p. 25 in von Bortkewitsch (1898)

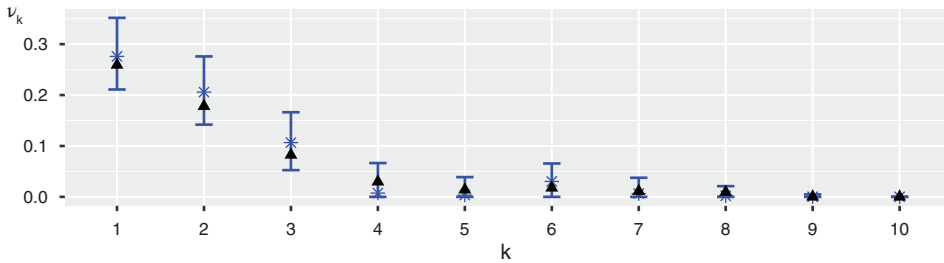
Deaths	0	1	2	3	4
Counts	109	65	22	3	1



**FIGURE 5** Estimation for the horse kick data from Subsection 4.1. The horizontal axis gives the magnitudes of  $v_k$ ,  $k \in \{1, \dots, 4\}$ . The black triangles denote the Buchmann-Grübel estimator, the blue crosses give the posterior means, whereas the vertical blue line segments represent (pointwise) 95% credible intervals [Colour figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]

**TABLE 4** Plant population data from Evans (1953)

Plants	0	1	2	3	4	5	6	7	8	9	10	11	12
Counts	274	71	58	36	20	12	10	7	6	3	0	2	1



**FIGURE 6** Estimation results for the plant data from Subsection 4.2. The horizontal axis gives the magnitudes of  $\nu_k$ ,  $k \in \{1, \dots, 10\}$ . The black triangles denote the Buchmann-Grübel estimates, the blue crosses give the posterior means, whereas the vertical blue line segments represent (pointwise) 95% credible intervals [Colour figure can be viewed at wileyonlinelibrary.com]

plot ranges from 0 to 12. The second row of Table 4 gives the counts of plots containing a given number of plants; thus, there were 274 plots that contained no plant, 71 that contained 1 plant, etc. It is customary in the ecological literature to model such count data as i.i.d. realizations from a compound Poisson distribution. Thus, for example, Neyman (1939) advocated the use of a Poisson base distribution in this context; another option here is a geometric base distribution. Given existence of several distinct modeling possibilities, performing an exploratory nonparametric analysis appears to be a sensible strategy.

The estimation results are graphically depicted in Figure 6. There are some small differences between the posterior mean and the Buchmann-Grübel estimate, but overall they are very similar.

## 5 | FREQUENTIST ASYMPTOTICS

In this section we assume that the observation times  $\{t_i\}$  are equidistant:  $t_i = i, i = 1, \dots, n$ . To evaluate our Bayesian method from a theoretical point of view, we will verify that it is consistent, and we will establish the rate at which the posterior contracts around the “true,” data-generating Lévy measure  $\nu_0$ ; see Ghosal and van der Vaart (2017) for a thorough treatment of Bayesian asymptotics from the frequentist point of view. From now on the subscript 0 in various quantities will refer to the data-generating distribution.

Our strategy consists in proving that the posterior contraction rate for  $\nu_0$ , given the sample  $\mathcal{Z}_n = (Z_1, \dots, Z_n)$ , can be derived from the posterior contraction rate for  $q_0$  given  $\mathcal{Z}_n$ , which is mathematically easier since  $Z_1, \dots, Z_n$  is a sequence of independent and identically distributed random variables with distribution  $q_0$ . We therefore effectively avoid dealing directly with the inverse nature of the problem of estimating  $p_0$ .

The prior we consider in this section is defined as follows:

- Endow the rate  $\lambda$  of the Poisson process with a prior distribution.
- Independently, endow the vector  $(p_1, \dots, p_m)$  with a Dirichlet distribution with parameter  $(\alpha_1, \dots, \alpha_m)$ .
- Set a priori  $p_k = 0$  for all  $k > m$ .

This is a somewhat simplified version of the prior we used in Section 2, which allows us to concentrate on essential features of the problem, without need to clutter the analysis with extra and unenlightening technicalities. Also remember the well-known relationship between the Gamma and Dirichlet distributions: if  $\xi_1, \dots, \xi_m$  are independent Gamma distributed random variables,  $\xi_i \sim \text{Gamma}(\alpha_i, 1)$ , then for  $\eta_i = \xi_i / \sum_{j=1}^m \xi_j$ , the vector  $(\eta_1, \dots, \eta_m)$  follows the Dirichlet distribution with parameter  $(\alpha_1, \dots, \alpha_m)$ ; furthermore, we have that  $\sum_{j=1}^m \xi_j \sim \text{Gamma}(\sum_{j=1}^m \alpha_j, 1)$ , and  $\eta_i$  are independent of  $\sum_{j=1}^m \xi_j$ .

In our asymptotic setting, we will make  $m = m_n$  dependent on  $n$  and let  $m_n \rightarrow \infty$  at a suitable rate as  $n \rightarrow \infty$ .

Recall that we write  $Q = (q_k)_{k \in \mathbb{N}_0}$  for  $q_k = Q(\{k\})$ . Let  $\mathcal{Q}$  denote the collection of all probability measures supported on  $\mathbb{N}$ .

**Theorem 1.** *Suppose there exists  $\underline{\alpha}$ , such that  $0 < \underline{\alpha} \leq \alpha_i \leq 1$  for all  $1 \leq i \leq m_n$ . Suppose  $\lambda \sim \text{Gamma}(a, b)$  with  $a \in (0, 1]$  and that  $\nu_0$  has a compact support. Then, for any  $\gamma > 1$ ,*

$$\Pi_n \left( \| \nu - \nu_0 \|_1 \geq \frac{\log^\gamma n}{\sqrt{n}} \mid \mathcal{Z}_n \right) \rightarrow 0$$

in  $Q_0^n$ -probability, as  $n \rightarrow \infty$ .

*Remark 1.* Note that since the support of  $\nu_0$  is not assumed to be known, our CPP model is still naturally nonparametric. The assumption of the compact support of  $\nu_0$  does not cover the simulation example of Section 3.2. However, its relaxation appears to pose very difficult technical challenges and is not attempted in this work.

The remainder of this section is devoted to the proof of Theorem 1.

## 5.1 | Basic posterior inequality via the stability estimate

A key step of the proof of Theorem 1 is the stability estimate in Equation (5) below, which bounds the total variation distance between the Lévy measures  $\nu, \nu'$  in terms of the total variation distance between the corresponding compound distributions  $q, q'$ .

In principle, it is conceivable that the Panjer recursion should allow one to bound probability distances between  $P$ -probabilities via distances between  $Q$ -probabilities; we call such a bound a stability estimate. Nevertheless, explicit as the equations of the Panjer recursion are, they are still somewhat unwieldy for that purpose. Hence we will use another inversion formula from Buchmann & Grübel (2003), that will lead to the stability estimate we are after.

First we introduce some notation, and also recall a few useful facts summarised in Buchmann & Grübel (2003). The space of absolutely summable sequences is defined as  $\ell_1 := \{a \in \mathbb{R}^{\mathbb{N}_0} : \sum_{j=0}^{\infty} |a_j| < \infty\}$ , with a norm given by  $\|a\|_1 = \sum_{j=0}^{\infty} |a_j|$ . For probability vectors  $a, b$ , the norm  $\|a - b\|_1$  is (twice) the total variation distance between  $a$  and  $b$ . For any  $a, b \in \ell_1$ , we have the inequality

$$\|a * b\|_1 \leq \|a\|_1 \|b\|_1, \quad (4)$$

where  $*$  denotes convolution of  $a$  and  $b$ . We define a mapping  $a \mapsto \exp(a)$  from  $\ell_1$  into  $\ell_1$  via

$$\exp(a) = \sum_{j=0}^{\infty} \frac{a^{*j}}{j!}.$$

The exponential has the following two useful properties:

$$\exp(a + b) = \exp(a) * \exp(b), \quad a, b \in \ell_1,$$

and

$$\exp(a) = \exp(b) \Rightarrow a = b, \quad a, b \in \ell_1.$$

We define a sequence  $\delta_0 = (\delta_{0,k})_{k \in \mathbb{N}_0}$ , such that  $\delta_{0,0} = 1$  and its all other entries are equal to zero. Then, using the above properties of the exponential, we can write concisely the compounding mapping in (2) in terms of convolutions of infinite sequences:  $q = \exp(\lambda(p - \delta_0))$ . Its convolution inverse, that is,  $q^{*(-1)}$  such that  $q^{*(-1)} * q = \delta_0$ , is given by  $r = q^{*(-1)} = \exp(-\lambda(p - \delta_0))$ . Note that  $r \in \ell_1$ . We have the following recursive expressions

$$r_0 = \frac{1}{q_0}, \quad r_k = -\frac{1}{q_0} \sum_{j=1}^k q_j r_{k-j}, \quad k \in \mathbb{N}.$$

**Lemma 1.** *Let  $q, q'$  correspond to two pairs  $(\lambda, p)$  and  $(\lambda', p')$ , respectively (and  $r$  correspond to  $q$ , i.e. the pair  $(\lambda, p)$ ). Then, in accordance with the notation introduced above and provided that  $\|q' - q\|_1 < \|r\|_1^{-1}$ , it holds that*

$$\|v' - v\|_1 = \|\lambda' p' - \lambda p\|_1 \leq \frac{\|r\|_1 \|q' - q\|_1}{1 - \|r\|_1 \|q' - q\|_1}. \tag{5}$$

*Proof.* The result is a direct consequence of Lemma 3 in Buchmann & Grübel (2003), which states that

$$(\lambda' - \lambda)\delta_0 + \lambda p - \lambda' p' = \sum_{j=1}^{\infty} \frac{1}{j} (r * (q - q'))^{*j},$$

whenever  $\|q' - q\|_1 < \|r\|_1^{-1}$ . Taking the  $\|\cdot\|_1$ -norm on both sides and some elementary bounding via (4) imply that

$$|\lambda' - \lambda| + \|\lambda' p' - \lambda p\|_1 \leq \frac{\|r\|_1 \|q' - q\|_1}{1 - \|r\|_1 \|q' - q\|_1},$$

and thus Equation (5) follows. ■

We will use Equation (5) to establish the key inequality for the posterior measure  $\Pi(\cdot | \mathcal{Z}_n)$ . We recall once again that the subscript 0 refers to “true,” data-generating quantities.

**Proposition 1.** *For any prior  $\Pi$  on  $v$ , for any  $\varepsilon \in (0, 1]$  and for any  $n \geq 1$ , the following posterior inequality holds:*

$$\Pi(\|v - v_0\|_1 \geq \varepsilon | \mathcal{Z}_n) \leq 2\Pi\left(\|q - q_0\|_1 \geq \frac{\varepsilon}{2\|r_0\|_1} \mid \mathcal{Z}_n\right).$$

*Proof.* Write  $\{v : \|v - v_0\|_1 \geq \varepsilon\}$  as a union of the sets

$$\{v : \|v - v_0\|_1 \geq \varepsilon\} \cap \{v : \|r_0\|_1 \|q - q_0\|_1 < 1/2\},$$

and

$$\{\nu : \|\nu - \nu_0\|_1 \geq \varepsilon\} \cap \{\nu : \|r_0\|_1 \|q - q_0\|_1 \geq 1/2\}.$$

Thanks to Lemma 1, the set

$$\{\nu : \|\nu - \nu_0\|_1 \geq \varepsilon\} \cap \{\nu : \|r_0\|_1 \|q - q_0\|_1 < 1/2\},$$

is a subset of  $\{\nu : \|q - q_0\|_1 \geq \varepsilon/(2\|r_0\|_1)\}$ . The proof is concluded by observing that  $\{\nu : \|\nu - \nu_0\|_1 \geq \varepsilon\} \cap \{\nu : \|r_0\|_1 \|q - q_0\|_1 \geq 1/2\}$  is a subset of  $\{\nu : \|q - q_0\|_1 \geq \varepsilon/(2\|r_0\|_1)\}$ , too, since  $\varepsilon \leq 1$ . ■

In general, stability estimates like the one in Equation (5) are unknown in the literature on Lévy processes. Consequently, studying Bayesian asymptotics for Lévy models, even in the CPP case, necessitates the use of very intricate arguments under restrictive assumptions (e.g., Nickl and Söhl (2017)).

## 5.2 | Proof of Theorem 1

The usefulness of Proposition 1 above lies in the fact that the posterior contraction rate in the inverse problem of estimating the Lévy measure  $\nu_0$  from indirect observations  $\mathcal{Z}_n$  can be now deduced from the posterior contraction rate in the direct problem of estimating the compound distribution  $q_0$ , which is easier (observe that  $r_0$  is determined by  $\nu_0$  and is therefore fixed in the proofs). The general machinery developed in Ghosal, Ghosh, and van der Vaart (2000) can be applied to handle the latter, and also several inequalities obtained in Gugushvili et al. (2015) are useful in that respect. In particular, we make use of the following inequality for the Hellinger distance,

$$h(q_{\lambda,p}, q_{\lambda',p'}) \leq \sqrt{\lambda} h(p, p') + \left| \sqrt{\lambda} - \sqrt{\lambda'} \right|, \quad (6)$$

compare Lemma 1 in Gugushvili et al. (2015). To ease our notation, in the sequel we will often write  $q$  and  $q'$  instead of  $q_{\lambda,p}$  and  $q_{\lambda',p'}$ , respectively.

Denote

$$\text{KL}(q_0, q) = \mathcal{Q}_0 \left( \log \frac{q_0}{q} \right), \quad V(q_0, q) = \mathcal{Q}_0 \left( \log \frac{q_0}{q} \right)^2.$$

Another two inequalities we will use are the following: let  $\lambda, \lambda_0 \in [\underline{\lambda}, \bar{\lambda}]$ . Then there exists a positive constant  $\bar{C}$ , such that

$$\begin{aligned} \text{KL}(q_0, q) &\leq \bar{C}(\text{KL}(p_0, p) + |\lambda_0 - \lambda|^2), \\ V(q_0, q) &\leq \bar{C}(V(p_0, p) + \text{KL}(p_0, p) + |\lambda_0 - \lambda|^2); \end{aligned} \quad (7)$$

Compare with Equations (14) and (15) in Lemma 1 in Gugushvili et al. (2015).

These three inequalities can be obtained by adjustment of the arguments used in Gugushvili et al. (2015). However, we opted to give their direct proofs in Lemma 5 from Appendix A under slightly weaker conditions than required in Gugushvili et al. (2015).

Our proof of Theorem 1 proceeds via verification of the conditions for posterior contraction in theorem 2.1 in Ghosal et al. (2000). In our setting, the latter result reads as follows.



**Theorem 2.** Assume  $\mathcal{Z}_n = (Z_1, \dots, Z_n)$ , where  $Z_1, \dots, Z_n$  are independent and identically distributed with distribution  $q_0$ . Let  $h$  denote the Hellinger metric on  $\mathcal{Q}$ , a collection of all measures with support in  $\mathbb{N}$ . Suppose that for a sequence  $\{\epsilon_n\}$  with  $\epsilon_n \rightarrow 0$  and  $n\epsilon_n^2 \rightarrow \infty$ , a constant  $C > 0$  and sets  $\mathcal{Q}_n \subset \mathcal{Q}$ , we have

$$\log N(\epsilon_n, \mathcal{Q}_n, h) \leq n\epsilon_n^2,$$

$$\Pi_n(\mathcal{Q} \setminus \mathcal{Q}_n) \leq \exp(-n\epsilon_n^2(C + 4)),$$

$$\Pi_n(q : \text{KL}(q_0, q) \leq \epsilon_n^2, V(q_0, q) \leq \epsilon_n^2) \geq \exp(-Cn\epsilon_n^2).$$

Then, for sufficiently large  $M > 0$ , we have that  $\Pi_n(\mathcal{Q} : h(q, q_0) \geq M\epsilon_n | \mathcal{Z}_n) \rightarrow 0$  in  $\mathcal{Q}_0^n$ -probability.

We will now verify the three conditions of this theorem, which we refer to as the entropy condition, the remaining mass condition, and the prior mass condition, respectively. To that end, fix strictly positive sequences  $\{\underline{\Lambda}_n\}$ ,  $\{\bar{\Lambda}_n\}$ , and define the sieves

$$\mathcal{Q}_n = \{q_{\lambda,p} : \lambda \in [\underline{\Lambda}_n, \bar{\Lambda}_n], \text{supp } p \subseteq \{1, \dots, m_n\}\}.$$

### 5.2.1 | Entropy

We start with bounding the entropy of the sieve  $\mathcal{Q}_n$  for  $h$ -balls of radius  $\epsilon_n$ .

**Lemma 2.** Assume that as  $n \rightarrow \infty$ ,

$$m_n \rightarrow \infty, \quad \epsilon_n \rightarrow 0, \quad \underline{\Lambda}_n \rightarrow 0, \quad \bar{\Lambda}_n \rightarrow \infty. \quad (8)$$

Then

$$\log N(\epsilon_n, \mathcal{Q}_n, h) \lesssim m_n \left\{ \log(m_n) + \log(\bar{\Lambda}_n) + \log\left(\frac{1}{\epsilon_n}\right) \right\} + \log\left(\frac{1}{\underline{\Lambda}_n}\right). \quad (9)$$

*Proof.* For  $\lambda, \lambda' \geq \underline{\Lambda}_n$ ,

$$|\sqrt{\lambda} - \sqrt{\lambda'}| = \frac{|\lambda - \lambda'|}{\sqrt{\lambda} + \sqrt{\lambda'}} \leq \frac{1}{2\sqrt{\underline{\Lambda}_n}} |\lambda - \lambda'|.$$

Furthermore, from section 3.3 in Pollard (2002),

$$h(p, p') \leq \sqrt{\|p - p'\|_1} \leq \sqrt{m_n \|p - p'\|_\infty}.$$

Combining the preceding two displays and Equation (6), we get

$$h(q_{\lambda,p}, q_{\lambda',p'}) \leq \sqrt{\bar{\Lambda}_n m_n \|p - p'\|_\infty} + \frac{1}{2\sqrt{\underline{\Lambda}_n}} |\lambda - \lambda'|.$$

Hence, if

$$\|p - p'\|_\infty \leq \frac{\epsilon_n^2}{4\bar{\Lambda}_n m_n}, \quad |\lambda - \lambda'| \leq \sqrt{\underline{\Lambda}_n} \epsilon_n,$$

then the Hellinger distance between  $q_{\lambda,p}$  and  $q_{\lambda',p'}$  is bounded by  $\epsilon_n$ . To cover  $[\underline{\Lambda}_n, \bar{\Lambda}_n]$ , we need at most  $\left\lfloor \frac{\bar{\Lambda}_n}{2\epsilon_n\sqrt{\underline{\Lambda}_n}} \right\rfloor + 1$  intervals of length  $2\sqrt{\underline{\Lambda}_n}\epsilon_n$ . To cover discrete distributions with support in  $\{1, \dots, m_n\}$ , we need at most

$$\left( \left\lfloor \frac{2\bar{\Lambda}_n m_n}{\epsilon_n^2} \right\rfloor + 1 \right)^{m_n}$$

$L_\infty$ -balls of radius  $\epsilon_n^2/(4\bar{\Lambda}_n m_n)$ . Under assumption (8), the summand 1 in the above display is asymptotically negligible and can be omitted. In that case, the number of  $h$ -balls that we need to cover  $\mathcal{Q}_n$  is of order

$$\left( \frac{\bar{\Lambda}_n m_n}{\epsilon_n^2} \right)^{m_n} \times \frac{\bar{\Lambda}_n}{\epsilon_n \sqrt{\underline{\Lambda}_n}}.$$

Taking the logarithm and next a straightforward rearrangement of the terms gives the statement of the lemma. ■

### 5.2.2 | Remaining prior mass

Now we will derive an inequality for the remaining prior mass.

**Lemma 3.** For  $\lambda \sim \text{Gamma}(a, b)$  with  $0 < a \leq 1$ ,

$$\Pi_n(\mathcal{Q} \setminus \mathcal{Q}_n) \lesssim \bar{\Lambda}_n^{a-1} e^{-b\bar{\Lambda}_n} + \underline{\Lambda}_n.$$

*Proof.* We have (with a slight abuse of notation)

$$\Pi_n(\mathcal{Q} \setminus \mathcal{Q}_n) = \Pi_n([\bar{\Lambda}_n, \infty)) + \Pi_n([0, \underline{\Lambda}_n)).$$

Now,

$$\Pi_n(\lambda \geq \bar{\Lambda}_n) = \frac{b^a}{\Gamma(a)} \int_{\bar{\Lambda}_n}^\infty \lambda^{a-1} e^{-b\lambda} d\lambda \lesssim \bar{\Lambda}_n^{a-1} e^{-b\bar{\Lambda}_n}.$$

Furthermore,

$$\Pi_n([0, \underline{\Lambda}_n)) = \frac{b^a}{\Gamma(a)} \int_0^{\underline{\Lambda}_n} \lambda^{a-1} e^{-b\lambda} d\lambda \lesssim \underline{\Lambda}_n^a.$$

The proof is concluded. ■

### 5.2.3 | Prior mass

Finally, we lower bound the prior mass in a small Kullback-Leibler neighbourhood of the data-generating compound distribution  $q_0$ . Define the function  $g : (0, \infty) \times (0, 1) \rightarrow (0, \infty)$  by

$$g(\epsilon, c) = C \frac{\epsilon^2}{2[\log(e/c)]^2},$$

where  $C$  is the constant appearing in the statement of Lemma 6 below.

**Lemma 4.** Assume that

- (i) there exists  $\alpha$ , such that  $0 < \underline{\alpha} \leq \alpha_i \leq 1$  for all  $1 \leq i \leq m_n$ ;
- (ii) strictly positive sequences  $\underline{p}_n \rightarrow 0$ ,  $\epsilon_n \rightarrow 0$  and  $m_n \rightarrow \infty$  satisfy the inequalities  $m_n g(\epsilon_n, \underline{p}_n) < 1$  and  $\underline{p}_n < g(\epsilon_n, \underline{p}_n)^2$ .

Define

$$B_n(\epsilon) = \{q \in \mathcal{Q}_n : \text{KL}(q_0, q) \leq \epsilon^2, V(q_0, q) \leq \epsilon^2\}.$$

Then

$$\Pi_n(B_n(\epsilon_n)) \gtrsim \Pi_n(|\lambda_0 - \lambda| \leq \tilde{\epsilon}_n) \times \Gamma \left( \sum_{i=1}^{m_n} \alpha_i \right) \exp(-m_n \log(1/(g(\tilde{\epsilon}_n, \underline{p}_n)^2 - \underline{p}_n)) - m_n \log(1/\underline{\alpha})).$$

Here  $\tilde{\epsilon}_n = \epsilon_n/\sqrt{3\bar{C}}$ , with a constant  $\bar{C} > 0$  not depending on  $n$ .

*Proof.* Define

$$\tilde{B}_n(\epsilon) = \left\{ (\lambda, p) : \lambda \in [\underline{\Lambda}_n, \bar{\Lambda}_n], \min_{1 \leq i \leq m_n} p_i \geq \underline{p}_n, \text{supp } p \subseteq \{1, \dots, m_n\}, \right. \\ \left. \text{KL}(p_0, p) \leq \epsilon^2, V(p_0, p) \leq \epsilon^2, |\lambda_0 - \lambda| \leq \epsilon \right\}.$$

For all  $n$  large enough and  $\epsilon$  small, we have  $\{\lambda : |\lambda_0 - \lambda| \leq \epsilon\} \subseteq [\underline{\Lambda}_n, \bar{\Lambda}_n]$ . Then by inequalities in Lemma 5,  $\tilde{B}_n(\epsilon) \subset B_n(\sqrt{3\bar{C}}\epsilon)$ , with a constant  $\bar{C}$  that can be taken the same for all large enough  $n$ ; see the arguments in Section 4.2 in Gugushvili et al. (2015). Hence, using the a priori independence of  $p$  and  $\lambda$ ,

$$\Pi_n(B_n(\epsilon_n)) \geq \Pi_n(\tilde{B}_n(\tilde{\epsilon}_n)) = \Pi_n(|\lambda_0 - \lambda| \leq \tilde{\epsilon}_n) \times U_n,$$

where

$$U_n = \Pi_n \left( \left\{ p : \text{KL}(p_0, p) \leq \tilde{\epsilon}_n^2, V(p_0, p) \leq \tilde{\epsilon}_n^2, \min_{1 \leq i \leq m_n} p_i \geq \underline{p}_n \right\} \right).$$

Furthermore, by Lemma 6 from Appendix A, we have

$$U_n \geq \Pi_n \left( \left\{ p : \sum_{i=1}^{m_n} |p_{0i} - p_i| \leq 2g(\tilde{\epsilon}_n, \underline{p}_n), \min_{1 \leq i \leq m_n} p_i \geq \underline{p}_n \right\} \right).$$

The statement of the lemma now follows upon applying Lemma 7 from Appendix A with  $\eta = \underline{p}_n$  and  $\epsilon = g(\tilde{\epsilon}_n, \underline{p}_n)$ . ■

## 5.2.4 | Using bounds in Theorem 2

We take

$$m_n \asymp \log n, \quad \epsilon_n \asymp \frac{\log^\gamma n}{\sqrt{n}}, \quad \underline{p}_n \asymp \frac{1}{n^2}, \\ \bar{\Lambda}_n \asymp \log^{2\gamma} n, \quad \underline{\Lambda}_n \asymp \exp(-\text{const} \cdot \log^{2\gamma} n)$$

with appropriately selected proportionality constants, and verify the conditions in Theorem 2.

Firstly, condition (8) is trivially satisfied. Therefore, we can invoke Lemma 2 and conclude that the entropy is upper bounded by a multiple of  $\log^{2\gamma} n$ , since  $\gamma > 1$ . Now  $\log^{2\gamma} n \lesssim n\epsilon_n^2$ , and this verifies the entropy condition in Theorem 2.

Be Lemma 3, for a suitable choice of the constant  $C$  the remaining prior mass condition is likewise satisfied.

Finally, for the prior mass condition in a small Kullback-Leibler neighbourhood to hold, by Lemma 4 we need that the term

$$\Pi_n(|\lambda_0 - \lambda| \leq \tilde{\epsilon}_n) \exp(-m_n \log(1/(g(\tilde{\epsilon}_n, \underline{p}_n)^2 - \underline{p}_n))) - m_n \log(1/\underline{\alpha})$$

is lower bounded by  $\exp(-Cn\epsilon_n^2)$  for some large enough  $C > 0$ . Now,  $\Pi_n(|\lambda_0 - \lambda| \leq \tilde{\epsilon}_n) \asymp \tilde{\epsilon}_n$ . Take the logarithm on both sides of the above display and note that by our conditions

$$\log(\Pi_n(|\lambda_0 - \lambda| \leq \tilde{\epsilon}_n)) \gtrsim \log(\tilde{\epsilon}_n) \gtrsim -n\epsilon_n^2.$$

Likewise,

$$m_n \log(1/(g(\tilde{\epsilon}_n, \underline{p}_n)^2 - \underline{p}_n)) + m_n \log(1/\underline{\alpha}) \lesssim n\epsilon_n^2,$$

so that the prior mass condition holds.

Thus we have verified all the conditions of Theorem 2. The resulting posterior contraction rate is  $\epsilon_n \asymp \log^\gamma n / \sqrt{n}$ .

## 6 | OUTLOOK

In this paper we introduced a nonparametric Bayesian approach to estimation of the Lévy measure  $\nu$  of a discretely observed CPP, when the support of  $\nu$  is a subset of  $\mathbb{N}$ . We constructed an algorithm for sampling from the posterior distribution of  $\nu$ , and showed that in practice our procedure performs well and measures up to a benchmark frequentist plug-in approach from (Buchmann & Grübel, 2004). Although computationally more demanding and slower than the latter, our method has an added benefit of providing uncertainty quantification in parameter estimates through the spread of the posterior distribution. On the theoretical side we show that our procedure is consistent, in that asymptotically, as the sample size  $n \rightarrow \infty$ , the posterior concentrates around the “true,” data-generating distribution. The corresponding posterior contraction rate is the (nearly) optimal rate  $\log^\gamma n / \sqrt{n}$  for an arbitrary  $\gamma > 1$ , if we are to ignore a practically insignificant  $\log n$  factor.

Among several generalizations of our results, the one that looks the most promising is extension of our methodology to CPP processes with jump size distributions supported on the set of integers  $\mathbb{Z}$ . The corresponding model has garnered substantial interest in financial applications (see Barndorff-Nielsen, Pollard, & Shephard, 2012). We leave this extension as a topic of possible future research.

## ACKNOWLEDGEMENTS


The research leading to the results in this paper has received funding from the European Research Council under ERC Grant Agreement 320637, from the Deutsche Forschungsgemeinschaft (DFG,

German Research Foundation) – 314838170, GRK 2297 MathCoRe, and from the Deutsche Forschungsgemeinschaft (DFG) through the grant CRC 1294 "Data Assimilation." The authors would like to thank the Associate Editor and the referee for their detailed and constructive comments on the paper.

## ORCID

Shota Gugushvili  <https://orcid.org/0000-0002-6963-295X>

Ester Mariucci  <https://orcid.org/0000-0003-0409-4131>

Frank van der Meulen  <https://orcid.org/0000-0001-7246-8612>

## REFERENCES

- Barndorff-Nielsen, O. E., Pollard, D. G., & Shephard, N. (2012). Integer-valued Lévy processes and low latency financial econometrics. *Quantitative Finance*, 12(4), 587–605.
- Belomestny, D., Comte, F., Genon-Catalot, V., Masuda, H., & Reiß, M. (2015). *Lévy matters IV. Estimation for discretely observed Lévy processes* (Vol. 2128). Cham, Switzerland: Springer.
- Belomestny, D., Gugushvili, S., Schauer, M., & Spreij, P. (2019). Nonparametric Bayesian inference for Gamma-type Lévy subordinators. *Communications in Mathematical Sciences*, 17(3), 781–816.
- Bezanson, J., Edelman, A., Karpinski, S., & Shah, V. (2017). Julia: A fresh approach to numerical computing. *SIAM Review*, 59(1), 65–98.
- Buchmann, B., & Grübel, R. (2003). Decompounding: An estimation problem for Poisson random sums. *The Annals of Statistics*, 31(4), 1054–1074.
- Buchmann, B., & Grübel, R. (2004). Decompounding Poisson random sums: Recursively truncated estimates in the discrete case. *Annals of the Institute of Statistical Mathematics*, 56(4), 743–756.
- Coca, A. J. (2018a). *Adaptive nonparametric estimation for compound Poisson processes robust to the discrete-observation scheme*. arXiv e-prints. Retrieved from <https://arxiv.org/abs/1803.09849>
- Coca, A. J. (2018b). Efficient nonparametric inference for discretely observed compound Poisson processes. *Probability Theory and Related Fields*, 170(1), 475–523.
- Comte, F., & Genon-Catalot, V. (2011). Estimation for Lévy processes from high frequency data within a long time interval. *The Annals of Statistics*, 39(2), 803–837.
- Devroye, L. (1986). *Non-uniform random variate generation*. New York, NY: Springer-Verlag.
- Duval, C., & Hoffmann, M. (2011). Statistical inference across time scales. *Electronic Journal of Statistics*, 5, 2004–2030.
- Duval, C., & Mariucci, E. (2017). *Compound Poisson approximation to estimate the Lévy density*. ArXiv e-prints. Retrieved from <https://arxiv.org/abs/1702.08787>
- Embrechts, P., Mikosch, T., & Klüppelberg, C. (1997). *Modelling extremal events: For insurance and finance*. Berlin, Germany: Springer-Verlag.
- Evans, D. A. (1953). Experimental evidence concerning contagious distributions in ecology. *Biometrika*, 40, 186–211.
- Gelfand, A. E., & Smith, A. F. M. (1990). Sampling-based approaches to calculating marginal densities. *Journal of the American Statistical Association*, 85(410), 398–409.
- Ghosal, S., Ghosh, J. K., & van der Vaart, A. W. (2000). Convergence rates of posterior distributions. *The Annals of Statistics*, 28(2), 500–531.
- Ghosal, S., & van der Vaart, A. (2007). Posterior convergence rates of Dirichlet mixtures at smooth densities. *The Annals of Statistics*, 35(2), 697–723.
- Ghosal, S., & van der Vaart, A. (2017). *Fundamentals of nonparametric Bayesian inference Cambridge Series in Statistical and Probabilistic Mathematics* (Vol. 44). Cambridge, UK: Cambridge University Press.
- Gugushvili, S., Mariucci, E., & van der Meulen, F. (2019). *Bdd: Julia code for Bayesian decompounding of discrete distributions*. doi:<https://doi.org/10.5281/zenodo.2598802>

- Gugushvili, S., van der Meulen, F., & Spreij, P. (2015). Nonparametric Bayesian inference for multidimensional compound Poisson processes. *Modern Stochastics: Theory & Applications*, 2(1), 1–15.
- Gugushvili, S., van der Meulen, F., & Spreij, P. (2018). A non-parametric Bayesian approach to decomposing from high frequency data. *Statistical Inference for Stochastic Processes*, 21(1), 53–79.
- Kappus, J. (2014). Adaptive nonparametric estimation for Lévy processes observed at low frequency. *Stochastic Processes and their Applications*, 124(1), 730–758.
- Lindo, A., Zuyev, S., & Sagitov, S. (2018). Nonparametric estimation for compound Poisson process via variational analysis on measures. *Statistics and Computing*, 28(3), 563–577.
- Müller, P., Quintana, F. A., Jara, A., & Hanson, T. (2015). *Bayesian nonparametric data analysis Springer Series in Statistics*. Cham, Switzerland: Springer.
- Neumann, M. H., & Reiß, M. (2009). Nonparametric estimation for Lévy processes from low-frequency observations. *Bernoulli*, 15(1), 223–248.
- Neyman, J. (1939). On a new class of “contagious” distributions, applicable in entomology and bacteriology. *The Annals of Mathematical Statistics*, 10, 35–57.
- Nickl, R., & Reiß, M. (2012). A Donsker theorem for Lévy measures. *Journal of Functional Analysis*, 263(10), 3306–3332.
- Nickl, R., & Söhl, J. (2017). *Bernstein-von Mises theorems for statistical inverse problems II: Compound Poisson processes*. ArXiv e-prints. Retrieved from <https://arxiv.org/abs/1709.07752>
- Panjer, H. H. (1981). Recursive evaluation of a family of compound distributions. *ASTIN Bulletin*, 12(1), 22–26.
- Pollard, D. (2002). *A user's guide to measure theoretic probability* (Vol. 8). Cambridge, UK: Cambridge University Press.
- Pya Arnqvist, N., Voinov, V., & Voinov, Y. (2018). *nilde: Nonnegative integer solutions of linear Diophantine equations with applications*. R package version 1.1-2. Retrieved from <https://CRAN.R-project.org/package=nilde>
- Sato, K.-I. (2013). *Lévy processes and infinitely divisible distributions* (Vol. 68, 2nd revised ed.). Cambridge, UK: Cambridge University Press.
- Shreve, S. E. (2004). *Stochastic calculus for finance. II: Continuous-time models*. New York, NY: Springer.
- Tanner, M. A., & Wong, W. H. (1987). The calculation of posterior distributions by data augmentation. *Journal of the American Statistical Association*, 82(398), 528–550 With discussion and with a reply by the authors.
- Trabs, M. (2015). Information bounds for inverse problems with application to deconvolution and Lévy models. *Annales de l'Institut Henri Poincaré, Probabilités et Statistiques*, 51(4), 1620–1650.
- van der Pas, S., Szabó, B., & van der Vaart, A. (2017). Uncertainty quantification for the horseshoe (with discussion). *Bayesian Analysis*, 12(4), 1221–1274.
- van der Vaart, A. W. (1998). *Asymptotic statistics Cambridge series in statistical and probabilistic mathematics* (Vol. 3). Cambridge, UK: Cambridge University Press.
- van Es, B., Gugushvili, S., & Spreij, P. (2007). A kernel type nonparametric density estimator for decomposing. *Bernoulli*, 13(3), 672–694.
- Voinov, V. G., & Nikulin, M. S. (1997). *On a subset sum algorithm and its probabilistic and other applications*. In N. Balakrishnan (Ed.), *Advances in combinatorial methods and applications to probability and statistics* (pp. 153–163). Boston, MA: Birkhäuser Boston.
- von Bortkewitsch, L. (1898). *Das Gesetz der kleinen Zahlen*. Leipzig, Germany: B.G. Teubner.
- Wasserman, L. (2004). *All of statistics. A concise course in statistical inference*. New York, NY: Springer.
- Zhang, H., Liu, Y., & Li, B. (2014). Notes on discrete compound Poisson model with applications to risk theory. *Insurance: Mathematics & Economics*, 59, 325–336.

**How to cite this article:** Gugushvili S, Mariucci E, der Meulen F. Decomposing discrete distributions: A nonparametric Bayesian approach. *Scand J Statist.* 2020;47:464–492. <https://doi.org/10.1111/sjos.12413>

APPENDIX A. TECHNICAL RESULTS

**Lemma 5.** Let  $q$  (resp.  $q'$ ) be the law at time 1 of a CPP with intensity  $\lambda$  (resp.  $\lambda'$ ) and jump distribution  $p$  (resp.  $p'$ ). Suppose that  $p$  and  $p'$  are distributions concentrated on  $\mathbb{C}$ . Then,

$$\begin{aligned} \text{KL}(q, q') &\leq \lambda \text{KL}(p, p') + \lambda' - \lambda + \lambda \log \frac{\lambda}{\lambda'}, \\ V(q, q') &\leq 2\lambda(V(p, p') + 2\text{KL}(p, p')) + 2\text{KL}(p, p')^2 \lambda^2 \\ &\quad + 2\lambda \log \left( \frac{\lambda}{\lambda'} \right) \left( 2(\lambda' - \lambda) + (\lambda - 1) \log \left( \frac{\lambda}{\lambda'} \right) \right) + 2(\lambda' - \lambda)^2, \\ h(q, q') &\leq \sqrt{\lambda} h(p, p') + \sqrt{1 - e^{-\frac{1}{2}(\sqrt{\lambda} - \sqrt{\lambda'})^2}} \leq \sqrt{\lambda} h(p, p') + |\sqrt{\lambda} - \sqrt{\lambda'}|. \end{aligned}$$

In particular, if  $\lambda, \lambda' \in [\underline{\Lambda}, \bar{\Lambda}]$  with  $0 < \underline{\Lambda} \leq \bar{\Lambda} < \infty$ , then there exists a positive constant  $\bar{C}$ , that depends on  $\underline{\Lambda}, \bar{\Lambda}$ , such that

$$\begin{aligned} \text{KL}(q, q') &\leq \bar{C}(\text{KL}(p, p') + |\lambda - \lambda'|), \\ V(q, q') &\leq \bar{C}(V(p, p') + \text{KL}(p, p') + (\lambda - \lambda')^2), \\ h(q, q') &\leq \bar{C}h(p, p') + |\sqrt{\lambda} - \sqrt{\lambda'}|. \end{aligned}$$

*Proof.* If  $\text{KL}(p, p')$  and  $V(p, p')$  are infinite, then the above inequalities are trivially satisfied. Therefore, we can assume these two divergences are finite. With this in mind, the proof of the lemma is divided into three steps.

**Step 1:** We begin by proving that for any  $n \geq 1$ ,

$$\begin{aligned} \text{KL}(p^{*n}, p'^{*n}) &\leq n\text{KL}(p, p'), \\ V(p^{*n}, p'^{*n}) &\leq nV(p, p') + 4n\text{KL}(p, p') + n(n - 1)\text{KL}(p, p')^2, \\ h^2(p^{*n}, p'^{*n}) &\leq nh^2(p, p'). \end{aligned}$$

The assertions are trivial for  $n = 1$ . Assuming that the first one holds for  $n - 1$  with  $n \geq 2$ , we will now show that it holds for  $n$  as well. Using the notation  $p^{*n}(i)$  for the  $i$ th element of  $p^{*n}$  and similarly in the case of  $p'^{*n}$ , we have

$$\begin{aligned} \text{KL}(p^{*n}, p'^{*n}) &= \sum_{i \in \mathbb{N}} p^{*n}(i) \log \left( \frac{p^{*n}(i)}{p'^{*n}(i)} \right) \\ &= \sum_{i \in \mathbb{N}} \sum_{k \in \mathbb{N}} p^{*(n-1)}(k) p(i - k) \log \left( \frac{\sum_{k \in \mathbb{N}} p^{*(n-1)}(k) p(i - k)}{\sum_{k \in \mathbb{N}} p'^{*(n-1)}(k) p'(i - k)} \right) \\ &\leq \sum_{i \in \mathbb{N}} \sum_{k \in \mathbb{N}} p^{*(n-1)}(k) p(i - k) \log \left( \frac{p^{*(n-1)}(k) p(i - k)}{p'^{*(n-1)}(k) p'(i - k)} \right) \\ &= \sum_{i \in \mathbb{N}} \sum_{k \in \mathbb{N}} p^{*(n-1)}(k) p(i - k) \left( \log \left( \frac{p^{*(n-1)}(k)}{p'^{*(n-1)}(k)} \right) + \log \left( \frac{p(i - k)}{p'(i - k)} \right) \right) \\ &= \text{KL}(p^{*(n-1)}, p'^{*(n-1)}) + \text{KL}(p, p'), \end{aligned}$$

where the inequality follows from the log-sum inequality, and the last equality is obtained by means of Fubini's theorem combined with the facts that

$$\sum_{k \in \mathbb{N}} p^{*(n-1)}(k) = 1, \quad \sum_{i \in \mathbb{N}} p(i-k) = 1, \quad \forall k \in \mathbb{N}.$$

By induction, we deduce that  $\text{KL}(p^{*n}, p'^{*n}) \leq n\text{KL}(p, p')$ .

The proof of the inequality for  $V$  is similar: we assume the inequality is true for  $n-1$  with  $n \geq 2$ , and will show it holds for  $n$  as well. Write  $V(p^{*n}, p'^{*n}) = R_1 + R_2$  for

$$R_1 = \sum_{i \in \mathbb{N}} p^{*n}(i) \log^2 \left( \frac{p^{*n}(i)}{p'^{*n}(i)} \right) \mathbf{1}_{\left\{ \frac{p^{*n}(i)}{p'^{*n}(i)} \geq 1 \right\}},$$

$$R_2 = \sum_{i \in \mathbb{N}} p^{*n}(i) \log^2 \left( \frac{p^{*n}(i)}{p'^{*n}(i)} \right) \mathbf{1}_{\left\{ \frac{p^{*n}(i)}{p'^{*n}(i)} < 1 \right\}}.$$

Observe that the function  $x \mapsto (x \log^2 x) \mathbf{1}_{\{x \geq 1\}}$  is convex. By Jensen's inequality we have for positive  $a_k, b_k$  that

$$\left( \sum_k a_k \right) \log^2 \left( \frac{\sum_k a_k}{\sum_k b_k} \right) \mathbf{1}_{\{\sum_k a_k \geq \sum_k b_k\}} \leq \sum_k a_k \log^2 \left( \frac{a_k}{b_k} \right) \mathbf{1}_{\{a_k/b_k \geq 1\}}.$$

Using this inequality and

$$p^{*n}(i) = \sum_{k \in \mathbb{N}} p^{*(n-1)}(k) p(i-k), \quad p'^{*n}(i) = \sum_{k \in \mathbb{N}} p'^{*(n-1)}(k) p'(i-k),$$

we get that

$$\begin{aligned} R_1 &\leq \sum_{i \in \mathbb{N}} \sum_{k \in \mathbb{N}} p(i-k) p^{*(n-1)}(k) \log^2 \left( \frac{p(i-k) p^{*(n-1)}(k)}{p'(i-k) p'^{*(n-1)}(k)} \right) \\ &= \sum_{i \in \mathbb{N}} \sum_{k \in \mathbb{N}} p(i-k) p^{*(n-1)}(k) \left( \log \left( \frac{p(i-k)}{p'(i-k)} \right) + \log \left( \frac{p^{*(n-1)}(k)}{p'^{*(n-1)}(k)} \right) \right)^2 \\ &= V(p^{*(n-1)}, p'^{*(n-1)}) + V(p, p') + 2\text{KL}(p, p') \text{KL}(p^{*(n-1)}, p'^{*(n-1)}) \\ &\leq V(p^{*(n-1)}, p'^{*(n-1)}) + V(p, p') + 2(n-1)\text{KL}(p, p')^2 \\ &\leq nV(p, p') + 4(n-1)\text{KL}(p, p') + n(n-1)\text{KL}(p, p')^2, \end{aligned}$$

where in the last inequality we used the induction hypothesis. Now recall an elementary inequality

$$e^{-x} x^2 \leq 4(e^{-x/2} - 1)^2$$

valid for  $x \geq 0$ ; see Gugushvili et al. (2015, p. 12). Applying this inequality to

$$x = -\log \left( \frac{p^{*n}(i)}{p'^{*n}(i)} \right)$$



such that  $p^{*n}(i)/p'^{*n}(i) < 1$ , we get

$$\frac{p^{*n}(i)}{p'^{*n}(i)} \log^2 \left( \frac{p^{*n}(i)}{p'^{*n}(i)} \right) \leq 4 \left( \sqrt{\frac{p^{*n}(i)}{p'^{*n}(i)}}} - 1 \right)^2.$$

By multiplying both sides of the above inequality with  $p'^{*n}(i)$ , summing the result through  $i$  and recalling the definition of the Hellinger distance, we get that

$$\begin{aligned} R_2 &\leq 4 \sum_{i \in \mathbb{N}} \left( \sqrt{p^{*n}(i)} - \sqrt{p'^{*n}(i)} \right)^2 = 4h^2(p^{*n}, p'^{*n}) \\ &\leq 4\text{KL}(p^{*n}, p'^{*n}) \leq 4n\text{KL}(p, p'). \end{aligned}$$

To conclude the proof of the inequality for  $V$ , we combine the bounds derived for  $R_1$  and  $R_2$ .

As far as the inequality for the Hellinger distance is concerned, we observe that

$$h^2(p^{*n}, p'^{*n}) = \sum_{i \in \mathbb{N}} p^{*n}(i) g \left( \frac{p'^{*n}(i)}{p^{*n}(i)} \right)$$

for a convex function  $g(x) = (1 - \sqrt{x})^2 \mathbf{1}_{[0, \infty)}(x)$ . Then, by the same reasoning as above, we have

$$\begin{aligned} h^2(p^{*n}, p'^{*n}) &\leq \sum_{i \in \mathbb{N}} \sum_{k \in \mathbb{N}} p^{*(n-1)}(k) p(i-k) g \left( \frac{p'^{*n}(i-k) p'(i-k)}{p^{*(n-1)}(k) p(i-k)} \right) \\ &= \sum_{i \in \mathbb{N}} \sum_{k \in \mathbb{N}} \left( \left( \sqrt{p^{*(n-1)}(k)} - \sqrt{p'^{*n}(i-k)} \right) \sqrt{p(i-k)} \right. \\ &\quad \left. + \sqrt{p'^{*n}(i-k)} \left( \sqrt{p(i-k)} - \sqrt{p'(i-k)} \right) \right)^2 \\ &= h^2(p^{*(n-1)}, p'^{*n}(i-k)) + h^2(p, p') \\ &\quad + 2 \sum_{i \in \mathbb{N}} \sum_{k \in \mathbb{N}} \left( \sqrt{p^{*(n-1)}(k)} - \sqrt{p'^{*n}(i-k)} \right) \\ &\quad \times \left( \sqrt{p(i-k)} - \sqrt{p'(i-k)} \right) \sqrt{p(i-k) p'^{*n}(i-k)}. \end{aligned}$$

Now note that the last summand satisfies

$$2 \left( \sum_{k \in \mathbb{N}} \sqrt{p'^{*n}(i-k) p^{*(n-1)}(k)} - 1 \right) \left( 1 - \sum_{k \in \mathbb{N}} \sqrt{p(k) p'(k)} \right) = -\frac{1}{2} h^2(p^{*(n-1)}, p'^{*n}(i-k)) h^2(p, p') \leq 0.$$

We therefore conclude that

$$h^2(p^{*n}, p'^{*n}) \leq h^2(p^{*(n-1)}, p'^{*n}(i-k)) + h^2(p, p'),$$

which leads to the desired inequality, by an induction argument.

**Step 2:** Now we prove the inequalities

$$\begin{aligned}
 \text{KL}(q, q') &\leq \sum_{n=0}^{\infty} \mathbb{P}(N = n) \text{KL}(p^{*n}, p'^{*n}) + \text{KL}(N, N'), \\
 V(q, q') &\leq 2 \sum_{n=0}^{\infty} \mathbb{P}(N = n) (V(p^{*n}, p'^{*n}) + 2\text{KL}(p^{*n}, p'^{*n})) + 2V(N, N'), \\
 h(q, q') &\leq \sqrt{\sum_{n=0}^{\infty} \mathbb{P}(N = n) h^2(p^{*n}, p'^{*n})} + h(N, N').
 \end{aligned}$$

Here  $N$  and  $N'$  are Poisson random variables with means  $\lambda$  and  $\lambda'$ , respectively, and with a slight abuse of notation,  $\text{KL}(N, N')$ ,  $V(N, N')$  and  $h(N, N')$  are the KL and  $V$  divergences and the Hellinger distance between the corresponding laws.

Note that

$$q(i) = \sum_{n=0}^{\infty} p^{*n}(i) P(N = n), \quad q'(i) = \sum_{n=0}^{\infty} p'^{*n}(i) P(N' = n).$$

Using this and the log-sum inequality,

$$\begin{aligned}
 \text{KL}(q, q') &= \sum_{i \in \mathbb{N}} q(i) \log \left( \frac{q(i)}{q'(i)} \right) \leq \sum_{i \in \mathbb{N}} \sum_{n \in \mathbb{N}} p^{*n}(i) P(N = n) \log \left( \frac{p^{*n}(i) P(N = n)}{p'^{*n}(i) P(N' = n)} \right) \\
 &= \sum_{i \in \mathbb{N}} \sum_{n \in \mathbb{N}} p^{*n}(i) P(N = n) \left( \log \left( \frac{p^{*n}(i)}{p'^{*n}(i)} \right) + \log \left( \frac{P(N = n)}{P(N' = n)} \right) \right) \\
 &= \sum_{n=0}^{\infty} P(N = n) \text{KL}(p^{*n}, p'^{*n}) + \text{KL}(N, N').
 \end{aligned}$$

For the divergence  $V$ , write  $V(q, q') = B_1 + B_2$  for

$$\begin{aligned}
 B_1 &= \sum_{i \in \mathbb{N}} q(i) \log^2 \left( \frac{q(i)}{q'(i)} \right) \mathbf{1}_{\left\{ \frac{q(i)}{q'(i)} \geq 1 \right\}} \\
 &\leq \sum_{i \in \mathbb{N}} \sum_{n \in \mathbb{N}} p^{*n}(i) \mathbb{P}(N = n) \log^2 \left( \frac{p^{*n}(i) \mathbb{P}(N = n)}{p'^{*n}(i) \mathbb{P}(N' = n)} \right) \\
 &\leq 2 \sum_{i \in \mathbb{N}} \sum_{n \in \mathbb{N}} p^{*n}(i) \mathbb{P}(N = n) \left( \log^2 \left( \frac{p^{*n}(i)}{p'^{*n}(i)} \right) + \log^2 \left( \frac{\mathbb{P}(N = n)}{\mathbb{P}(N' = n)} \right) \right) \\
 &= 2 \sum_{n=0}^{\infty} V(p^{*n}, p'^{*n}) \mathbb{P}(N = n) + 2V(N, N'), \\
 B_2 &= \sum_{i \in \mathbb{N}} q(i) \log^2 \left( \frac{q(i)}{q'(i)} \right) \mathbf{1}_{\left\{ \frac{q(i)}{q'(i)} < 1 \right\}}.
 \end{aligned}$$

To control  $B_2$ , we use the same arguments as in the proof of inequalities (12) and (15) in Gugushvili et al. (2015), getting

$$B_2 \leq 4\text{KL}(q, q') \leq 4 \sum_{n=0}^{\infty} \text{KL}(p^{*n}, p'^{*n}) \mathbb{P}(N = n).$$

This gives the required inequality for the  $V$  divergence.

Finally, we prove the inequality for the Hellinger distance. Denoting the law of  $\sum_{j=1}^N Y_j$  by  $\tilde{q}$ , it holds by the triangle inequality that  $h(q, q') \leq h(q, \tilde{q}) + h(\tilde{q}, q')$ . Since  $g(x) = (1 - \sqrt{x})^2 \mathbf{1}_{[0, \infty)}(x)$  is a convex function,

$$h^2(q, \tilde{q}) \leq \sum_{i \in \mathbb{N}} \sum_{n \in \mathbb{N}} p^{*n}(i) \mathbb{P}(N = n) g\left(\frac{p'^{*n}(i)}{p^{*n}(i)}\right) = \sum_{n \in \mathbb{N}} \mathbb{P}(N = n) h^2(p^{*n}, p'^{*n}).$$

It remains to prove that  $h(\tilde{q}, q') \leq h(N, N')$ . This again follows by convexity of  $g$ , since

$$\begin{aligned} h^2(\tilde{q}, q') &= \sum_{i \in \mathbb{N}} \sum_{n \in \mathbb{N}} p'^{*n}(i) \mathbb{P}(N = n) g\left(\frac{\sum_{n=0}^{\infty} p'^{*n}(i) \mathbb{P}(N' = n)}{\sum_{n=0}^{\infty} p'^{*n}(i) \mathbb{P}(N = n)}\right) \\ &\leq \sum_{i \in \mathbb{N}} \sum_{n \in \mathbb{N}} p'^{*n}(i) \mathbb{P}(N = n) g\left(\frac{\mathbb{P}(N' = n)}{\mathbb{P}(N = n)}\right) = h^2(N, N'). \end{aligned}$$

**Step 3:** From Steps 1 and 2 we derive that

$$\begin{aligned} \text{KL}(q, q') &\leq \text{KL}(p, p') \mathbb{E}[N] + \text{KL}(N, N'), \\ V(q, q') &\leq 2\mathbb{E}[N](V(p, p') + 2\text{KL}(p, p')) + 2V(N, N') \\ &\quad + 2(\mathbb{E}[N^2] - \mathbb{E}[N])\text{KL}(p, p')^2, \\ h(q, q') &\leq h(p, p') \sqrt{\mathbb{E}[N]} + h(N, N'). \end{aligned}$$

Now the proof of the lemma follows from these three inequalities upon noticing that  $\text{KL}(p, p')^2 \leq V(p, p')$ , and recalling that

$$\begin{aligned} \mathbb{E}[N] &= \lambda, \\ \mathbb{E}[N^2] &= \lambda + \lambda^2, \\ \text{KL}(N, N') &= \lambda' - \lambda + \lambda \log \frac{\lambda}{\lambda'} \lesssim (\lambda - \lambda')^2, \\ h^2(N, N') &= 1 - e^{-\frac{1}{2}(\sqrt{\lambda} - \sqrt{\lambda'})^2} \leq |\sqrt{\lambda} - \sqrt{\lambda'}|, \\ V(N, N') &= \lambda \log\left(\frac{\lambda}{\lambda'}\right) \left(2(\lambda' - \lambda) + (\lambda - 1) \log\left(\frac{\lambda}{\lambda'}\right)\right) + (\lambda' - \lambda)^2 \lesssim (\lambda - \lambda')^2. \end{aligned}$$

■

**Lemma 6.** Let  $\epsilon > 0$ . Suppose  $p = (p_1, \dots, p_m)$  and  $p' = (p'_1, \dots, p'_m)$  are points in the  $m$ -dimensional unit simplex, and let  $\min_{1 \leq i \leq m} p_i \geq c$  for some  $c \in (0, 1)$ . Then there exists a universal constant  $C > 0$ , such that the inequality

$$\sum_{i=1}^m |p'_i - p_i| \leq C \frac{\epsilon^2}{[\log(e/c)]^2},$$

implies that  $\text{KL}(p', p) \leq \epsilon^2$  and  $V(p', p) \leq \epsilon^2$  hold.

*Proof.* Lemma 8 in Ghosal and van der Vaart (2007) assures that there exists a constant  $\bar{C}$  (not depending on either  $p$  or  $p'$ ), such that

$$\text{KL}(p', p) \leq \bar{C} h^2(p', p) \left[ 1 + \log \left( \left\| \frac{p'}{p} \right\|_{\infty} \right) \right]$$

and

$$V(p', p) \leq \bar{C} h^2(p', p) \left[ 1 + \log \left( \left\| \frac{p'}{p} \right\|_{\infty} \right) \right]^2.$$

From section 3.3 in Pollard (2002) we have  $h^2(p', p) \leq \sum_{i=1}^m |p'_i - p_i|$ . Furthermore, since  $0 < c < 1$  and  $\min_{1 \leq i \leq m} p_i \geq c$ ,

$$1 \leq 1 + \log \left( \left\| \frac{p'}{p} \right\|_{\infty} \right) \leq 1 + \log(1/c) = \log(e/c).$$

Therefore,

$$\max(\text{KL}(p', p), V(p', p)) \leq \bar{C} [\log(e/c)]^2 \sum_{i=1}^m |p'_i - p_i|,$$

from which the assertion of the lemma follows trivially.  $\blacksquare$

**Lemma 7.** Let  $m \geq 2$  be an integer. Suppose  $(p_1, \dots, p_m) \sim \text{Dir}(\alpha_1, \dots, \alpha_m)$ . Let  $p_0 = (p_{01}, \dots, p_{0m})$  be an arbitrary point in the  $m$ -dimensional unit simplex. Assume there exists  $\underline{\alpha}$ , such that  $0 < \underline{\alpha} \leq \alpha_i \leq 1$  for all  $1 \leq i \leq m$ . Let  $\epsilon > 0$ , and let  $\eta$  be such that  $\eta < \epsilon^2$ . Then if  $m\epsilon < 1$ ,

$$\Pi_n \left( \sum_{i=1}^m |p_i - p_{0i}| \leq 2\epsilon, \min_{1 \leq i \leq m} p_i \geq \eta \right) \geq \Gamma \left( \sum_{i=1}^m \alpha_i \right) \exp(-m \log(1/(\epsilon^2 - \eta)) - m \log(1/\underline{\alpha})). \quad (\text{A1})$$

*Proof.* By arguments analogous to those in the proofs of lemma 6.1 in Ghosal et al. (2000) and lemma 10 in Ghosal and van der Vaart (2007), we obtain that the left-hand side of (A1) can be lower bounded by

$$\frac{\Gamma(\sum_{i=1}^m \alpha_i)}{\prod_{i=1}^m \Gamma(\alpha_i)} \prod_{i=1}^{m-1} \int_{\max(p_{0i} - \epsilon^2, \eta^2)}^{\min(p_{0i} + \epsilon^2, 1)} x_i^{\alpha_i - 1} dx_i.$$

The length of the integration interval in each of the integrals in the above product is lower bounded by  $\epsilon^2 - \eta$ . Using that  $\underline{\alpha} \leq \alpha_i \leq 1$ , we deduce that the preceding display is lower bounded by

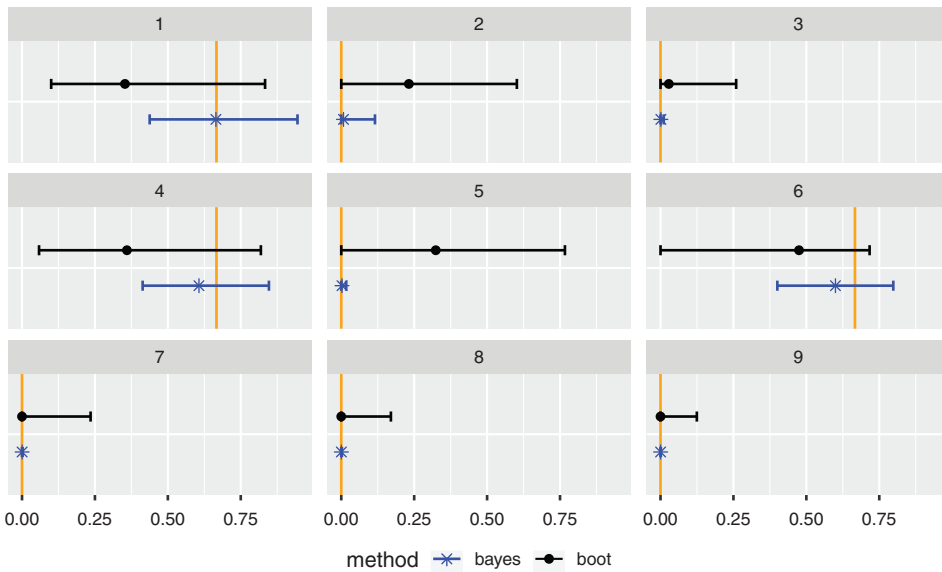
$$\Gamma \left( \sum_{i=1}^m \alpha_i \right) \underline{\alpha}^m \exp \left( -(m-1) \log \left( \frac{1}{\epsilon^2 - \eta} \right) \right).$$

This entails the result.  $\blacksquare$

## APPENDIX B. BOOTSTRAP CONFIDENCE INTERVALS

Here we report a small comparison between the Bayesian credible intervals and the bootstrap confidence intervals for the Buchmann-Grübel estimator. The setup and the simulated

dataset that we used are the same as in Subsection 3.1. The bootstrap confidence intervals were computed as follows:  $B = 9,999$  bootstrap samples were generated from the compound Poisson model under the Buchmann-Grübel estimates computed from the observed data. These bootstrap samples were then fed back to the Buchmann-Grübel procedure to yield  $B$  bootstrap estimates of the Lévy measure  $\nu$ . Finally, for each  $k$ , the  $\alpha/2$ th and  $(1 - \alpha/2)$ -th sample quantiles were obtained to yield  $1 - \alpha$  bootstrap confidence intervals for  $\nu_k$ . The results with  $\alpha = .05$  are displayed in Figure B1. One observes that while both methods result in a good coverage for this specific dataset, the bootstrap appears to be noticeably more conservative than the Bayesian approach, as evidenced by the width of the intervals.



**FIGURE B1** 95% bootstrap confidence intervals and Bayesian credible intervals for the simulation example from Section 3.1 under setting (a). The displayed results are for parameters  $\nu_1$  through  $\nu_9$ , with panels labeled sequentially from 1 to 9. The true parameter values are visualised with orange vertical lines. The coloring scheme is the same as in Figure 1. Note that some of the narrow Bayesian credible intervals are overshadowed by the symbol (star) used to visualise the posterior mean [Colour figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]