# Scaled process priors for Bayesian nonparametric estimation of the unseen genetic variation

Federico Camerlenghi [*,1], Stefano Favaro [†,2], Lorenzo Masoero [‡,3], and Tamara Broderick [§,3]

[1]Department of Economics, Management and Statistics, University of Milano - Bicocca, Piazza dell'Ateneo Nuovo 1, Milano
[2]Department of Economics and Statistics, University of Torino, Corso Unione Sovietica 218/bis, Torino
[3]Department of Electrical Engineering and Computer Science, CSAIL, Massachusetts Institute of Technology, Cambridge, Massachusetts, USA

## Abstract

There is a growing interest in the estimation of the number of unseen features, mostly driven by biological applications. A recent work brought out a peculiar property of the popular completely random measures (CRMs) as prior models in Bayesian nonparametric (BNP) inference for the unseen-features problem: for fixed prior's parameters, they all lead to a Poisson posterior distribution for the number of unseen features, which depends on the sampling information only through the sample size. CRMs are thus not a flexible prior model for the unseen-features problem and, while the Poisson posterior distribution may be appealing for analytical tractability and ease of interpretability, its independence from the sampling information makes the BNP approach a questionable oversimplification, with posterior inferences being completely determined by the estimation of unknown prior's parameters. In this paper, we introduce the stable-Beta scaled process (SB-SP) prior, and we show that it allows to enrich the posterior distribution of the number of unseen features arising under CRM priors, while maintaining its analytical tractability and interpretability. That is, the SB-SP prior leads to a negative Binomial posterior distribution, which depends on the sampling information through the sample size and the number of distinct features, with corresponding estimates being simple, linear in the sampling information and computationally efficient. We apply our BNP approach to synthetic data and to real cancer genomic data, showing that: i) it outperforms the most popular parametric and nonparametric

[*]Also affiliated to Collegio Carlo Alberto, Piazza V. Arbarello 8, Torino, and BIDSA, Bocconi University, Milano, Italy; federico.camerlenghi@unimib.it
[†]Also affiliated to Collegio Carlo Alberto, Piazza V. Arbarello 8, Torino, and IMATI-CNR "Enrico Magenes", Milan, Italy; stefano.favaro@unito.it
[‡]lom@mit.edu
[§]tbroderick@csail.mit.edu

competitors in terms of estimation accuracy; ii) it provides improved coverage for the estimation with respect to a BNP approach under CRM priors.

# 1 Introduction

The problem of estimating the number of unseen features generalizes the popular unseen-species problem [Orlitsky et al., 2016], and its importance has grown dramatically in recent years, driven by applications in biological sciences [Ionita-Laza et al., 2009, Gravel, 2014, Zou et al., 2016, Chakraborty et al., 2019]. Consider a generic population in which each individual is endowed with a finite collection of $\mathbb{W}$-valued features, with $\mathbb{W}$ possibly being an infinite space, and denote by $p_i$ the probability that an individual has feature $w_i \in \mathbb{W}$ for $i \geq 1$. The unseen-features problem assumes $N \geq 1$ observable random samples $Z_{1:N} = (Z_1, \dots, Z_N)$ from the population, such that $Z_n = (A_{n,i})_{i \geq 1}$ are independent Bernoulli random variables with unknown parameters $(p_i)_{i \geq 1}$. Then, the goal is to estimate the number of hitherto unseen features that would be observed if $M \geq 1$ additional samples were collected, i.e.

$$U = \sum_{i \geq 1} \mathbb{1}\left(\sum_{n=1}^{N} A_{n,i} = 0\right) \mathbb{1}\left(\sum_{m=1}^{M} A_{N+m,i} > 0\right),$$

with $\mathbb{1}$ being the indicator function. The unseen-species problem arises under the assumption that each individual is endowed with only one feature, i.e. a species. A wide range of approaches have been developed to estimate $U$, including Bayesian methods [Ionita-Laza et al., 2009, Masoero et al., 2021], jackknife [Gravel, 2014], linear programming [Zou et al., 2016], and variations of Good-Toulmin estimators [Orlitsky et al., 2016, Chakraborty et al., 2019].

In biological sciences, we may think of individuals as organisms and of features as groups to which organisms belong to, with each group being defined by any difference in the genome relative to a reference genome, i.e. a (genetic) variant. In human biology, the estimation of $U$ arises in the context of optimal allocation of resources between quantity and quality in genetic experiments: spending resources to sequence a greater number of genomes (quantity), which reveals more about variation across the population, or spending resources to sequence genomes with increased accuracy (quality), which reveals more about individual organisms' genomes. Accurate estimates of $U$ are critical in the experimental pipeline towards the goal of maximizing the usefulness of experiments under the trade-off between quantity and quality [Ionita-Laza and Laird, 2010, Zou et al., 2016]. While in human-biology the cost of sequencing has decreased in recent years [Schwarze et al., 2020], the expense remains non-trivial, and it is still critical in fields where scientists work with relatively budgets, e.g. non-human and non-model organisms [Souza et al., 2017]. Other applications arise in precision medicine [Momozawa and Mizukami, 2020], microbiome analysis [Sanders et al., 2019], single-cell sequencing [Zhang et al., 2020] and wildlife monitoring [Johansson et al., 2020].

## 1.1 Our contributions

We introduce a Bayesian nonparametric (BNP) approach to the unseen-features problem, which relies on a novel prior distribution for the unknown $(p_i)_{i \geq 1}$. Completely random measures (CRMs) [Kingman, 1992] provide a broad class of nonparametric priors for feature sampling problems, the most popular being the stable-Beta process prior [James, 2017, Broderick et al., 2018]. In a recent work, Masoero et al. [2021] brought out a peculiar feature of CRM priors in the unseen-features problem: they all lead to a Poisson posterior distribution of $U$, given $Z_{1:N}$ and fixed prior's parameters, which depends on $Z_{1:N}$ only through the sample size $N$. Despite the broadness of the class of CRM priors, such a common Poisson posterior structure makes CRMs not a flexible prior model for the unseen-features problem. While the Poisson posterior distribution may be appealing in principle, making posterior inferences analytically tractable and easy to interpret, its independence from $Z_{1:N}$ makes the BNP approach a questionable oversimplification, with posterior inferences being completely determined by the estimation of the unknown prior's parameters. A somehow similar scenario occurs in BNP inference for the unseen-species problem under a Dirichlet process (DP) prior [Ferguson, 1973], and led to the use of the Pitman-Yor process (PYP) prior [Pitman and Yor, 1997] for enriching the posterior distribution of the number of unseen species, while maintaining analytical tractability and interpretability of the DP prior [Lijoi et al., 2007].

We show that scaled process (SP) priors, first introduced in James et al. [2015], allow to enrich the posterior distribution of $U$ arising under CRM priors. Under SP priors, we characterize the posterior distribution of $U$ as a mixture of Poisson distributions that may include, through the mixing distribution, the whole sampling information in terms of the number of distinct features and their frequencies. While this is appealing in principle, it may be at stake with analytical tractability and interpretability, which are critical for a concrete use of SP priors. Then, we introduce the stable-Beta SP (SB-SP) prior, which provides a sensible trade-off between the amount of sampling information introduced in the posterior distribution of $U$, and analytical tractability and interpretability of the posterior inferences. In particular, we characterize the SB-SP prior as the sole SP prior for which the posterior distribution of $U$, given $Z_{1:N}$ and fixed prior's parameters, depends on $Z_{1:N}$ through the sample size $N$ and the number $K_N$ of distinct features; the SB-SP may thus be considered as the natural counterpart of the PYP for the unseen-feature problem. Under the SB-SP prior, the posterior distribution of $U$, as well as of a refinement of $U$ that deals with the number of unseen rare features, is a negative Binomial posterior distributions, whose parameters depend on $N$, $K_N$ and the prior's parameters. Corresponding Bayesian estimates of $U$, with respect to a squared loss function, are simple, linear in $K_N$ and computationally efficient.

We present an empirical validation of the effectiveness of our BNP methodology, both on synthetic and real data. As for real data, we consider cancer genomic data, where the goal is to estimate the number of new (genomic) variants to be discovered in future unobservable samples. In cancer genomics, accurate estimates of the number of new variants is of particular importance, as it might help practitioners understand the site of origin of cancers, as well as the clonal origin of metastasis, and in turn be a useful tool to develop effective clinical strategies [Chakraborty et al., 2019, Huyghe et al., 2019]. We make use of data from the cancer genome atlas (TCGA), and focus on the challenging scenario in which the sample size $N$ is particularly small, and also small with respect to the

extrapolation size $M$. Such a scenario is of interest in genomic applications, where only few samples of rare cancer might be available. We show that our BNP methodology outperforms the most popular parametric and nonparametric competitors, both classical (frequentist) and Bayesian, in terms of estimation accuracy of $U$ and a refinement of $U$ for rare features. In addition, with respect to the BNP approach under the stable-Beta process prior [Masoero et al., 2021], our approach provides improved coverage for the estimation. This is an empirical evidence of the effectiveness of replacing the Poisson posterior distribution with the negative Binomial posterior distribution, which allows to better exploit the sampling information.

## 1.2 Organization of the paper

In Section 2 we show how SP priors allow to enrich the posterior distribution of $U$ arising under CRM priors. In Section 3 we introduce and investigate the SB-SP prior in the context of the unseen-features problem: i) we characterize the SB-SP prior in the class of SP priors, providing its predictive distribution; ii) we apply the SB-SP prior to the unseen-features problem, providing the posterior distribution of $U$ and a BNP estimator. Section 4 contains illustrations of our method. In Section 5 we discuss our approach, a multivariate extension of it, and future research directions. Proofs and additional experiments are in the Appendix.

## 2 Scaled process priors for feature sampling problems

For a measurable space of features $\mathbb{W}$, we assume $N \geq 1$ observable individuals to be modeled as a random sample $Z_{1:N}$ from the $\{0,1\}$-valued stochastic process $Z(w) = \sum_{i \geq 1} A_i \delta_{w_i}(w)$, $w \in \mathbb{W}$, where $(w_i)_{i \geq 1}$ are features in $\mathbb{W}$ and $(A_i)_{i \geq 1}$ are independent Bernoulli random variables with unknown parameters $(p_i)_{i \geq 1}$, $p_i$ being the probability that an individual has feature $w_i$, for $i \geq 1$. That is, $Z$ is a Bernoulli process with parameter $\zeta = \sum_{i \geq 1} p_i \delta_{w_i}$, denoted as $\mathrm{BeP}(\zeta)$. BNP inference for feature sampling problems relies on the specification of a prior distribution on the discrete measure $\zeta$, leading to the BNP-Bernoulli model,

$$
\begin{aligned}
Z_n \,|\, \zeta &\overset{\mathrm{iid}}{\sim} \quad \mathrm{BeP}(\zeta) \qquad n = 1, \dots, N, \\
\zeta &\sim \quad \mathscr{Z},
\end{aligned}
\tag{1}
$$

namely $\zeta$ is a discrete random measure on $\mathbb{W}$ whose law $\mathscr{Z}$ takes on the interpretation of a prior distribution for the unknown feature's composition of the population. By de Finetti's theorem, the random variables $Z_n$'s in (1) are exchangeable with directing measure $\mathscr{Z}$ [Aldous, 1983]. In this section, we show how SP priors for $\zeta$ [James et al., 2015] allow to enrich the posterior distribution of the number of unseen features arising under CRM priors.

### 2.1 CRM priors for Bernoulli processes

CRMs provide a standard tool to define nonparametric prior distributions on the parameter $\zeta$ of the Bernoulli process $Z$. Consider a homogeneous CRM $\mu_0$ on $\mathbb{W}$, i.e. $\mu_0 = \sum_{i \geq 1} \rho_i \delta_{W_i}$, where the $\rho_i$'s

are $(0,1)$-valued random atoms such that $\sum_{i \geq 1} \rho_i < +\infty$, while the $W_i$'s are i.i.d. $\mathbb{W}$-valued random locations independent of the $\rho_i$'s. The law of $\mu_0$ is characterized, through Lévy-Khintchine formula, by the Lévy intensity measure $\nu_0(\mathrm{d}s, \mathrm{d}w) = \lambda_0(s)\mathrm{d}sP(\mathrm{d}w)$ on $(0,1) \times \mathbb{W}$, where: i) $\lambda_0$ is a measure on $(0,1)$, which controls the distribution of the $\rho_i$'s, and such that $\int_{(0,1)} \min\{s,1\}\lambda_0(s)\mathrm{d}s < +\infty$; ii) $P$ is a non-atomic measure on $\mathbb{W}$, which controls the distribution of the $W_i$'s. For short, $\mu_0 \sim \mathrm{CRM}(\nu_0)$. See Appendix A for an account on CRMs [Kingman, 1992, Chapter 8]. Note that, since $P$ is non-atomic, the random atoms $W_i$'s are almost surely distinct, that is to say the different features cannot coincide almost surely. The law of $\mu_0$ provides a natural prior distribution for the parameter $\zeta$ of the Bernoulli process. The Beta and the stable-Beta processes are popular examples of $\mu_0 \sim \mathrm{CRM}(\nu_0)$, for suitable specifications of $\nu_0$. A comprehensive posterior analysis of CRM priors is presented in James [2017]. In the next proposition, we recall the predictive distribution of CRM priors [James, 2017, Proposition 3.2].

**Proposition 1.** *Let $Z_{1:N}$ be a random sample from* (1) *with $\zeta \sim \mathrm{CRM}(\nu_0)$. If $Z_{1:N}$ displays $K_N = k$ distinct features $\{W_1^*, \ldots, W_{K_N}^*\}$, each feature $W_i^*$ appearing exactly $M_{N,i} = m_i$ times, then the conditional distribution of $Z_{N+1}$, given $Z_{1:N}$, coincides with the distribution of*

$$Z_{N+1} \mid Z_{1:N} \stackrel{d}{=} Z'_{N+1} + \sum_{i=1}^{K_N} A_{N+1,i}\delta_{W_i^*}, \tag{2}$$

*where: i) $Z'_{N+1} \mid \mu'_0 = \sum_{i \geq 1} A'_{N+1,i}\delta_{W'_i} \sim \mathrm{BeP}(\mu'_0)$ and $\mu'_0 \sim \mathrm{CRM}(\nu'_0)$, with $\nu'_0(\mathrm{d}s, \mathrm{d}w) = (1 - s)^N\lambda_0(s)\mathrm{d}sP(\mathrm{d}w)$; ii) the $A_{N+1,i}$'s are independent Bernoulli random variables with parameters $J_i$'s, such that $J_i$ is distributed according to the density function $f_{J_i}(s) \propto (1-s)^{N-m_i}s^{m_i}\lambda_0(s)$ for $i \geq 1$.*

According to (2), $Z_{N+1}$ displays "new" features $W'_i$'s, i.e. features not appearing in the initial sample $Z_{1:N}$, and "old" features $W_i^*$'s, i.e. features appeared in the initial sample $Z_{1:N}$. The posterior distribution of statistics of "new" features is determined by the law of $Z'_{N+1}$, which depends on $Z_{1:N}$ only through the sample size $N$; the posterior distribution of statistics of "old" features is determined by the law of $\sum_{1 \leq i \leq K_N} A_{N+1,i}\delta_{W_i^*}$, which depends on $Z_{1:N}$ through the sample size $N$, the number $K_N$ of distinct features and their frequencies $(M_{N,1}, \ldots, M_{N,K_N})$. As a corollary of Proposition 1, the posterior distribution of the number of "new" features in $(Z_{N+1}, \ldots, Z_{N+M})$, given $Z_{1:N}$ and fixed prior's parameters, is a Poisson distribution that depends on $Z_{1:N}$ only through $N$ [Masoero et al., 2021]. Such a posterior structure is peculiar to CRM priors, being inherited by the Poisson process formulation of CRMs [Kingman, 1992]. That is, despite the broadness of the class of CRM priors, all CRM priors lead to the same Poisson posterior structure for the number of unseen features, which thus makes them not a flexible prior model for the unseen-features problem. While the Poisson posterior distribution may be appealing in principle, making the posterior inferences analytically tractable and of easy interpretability, its independence from $Z_{1:N}$ makes the BNP approach under CRM priors a questionable oversimplification, with posterior inferences being completely determined by the estimation of unknown prior's parameters.

**Remark 1.** *For the sake of mathematical convenience, and in agreement with the work of James [2017], in the sequel we maintain the random measure formulation for both the prior model $\mu_0$ and the Bernoulli processes $Z_n$. However, we point out that each $Z_n$ is equivalently characterized by means*

of the Bernoulli variables $(A_{n,i})_{i \geq 1}$ and the random features $(W_i)_{i \geq 1}$. In other terms, there exits a one-to-one correspondence between $Z_n$ and the sequence of points $\{(A_{n,i}, W_i)\}_{i \geq 1}$. Finally, note that, although the values of features' labels $W_i$ are immaterial, the features $W_i$'s are assumed to be random. This is in line with the BNP literature on species sampling models, where the species' labels are assumed to be random [Pitman, 1996].

## 2.2   SP priors for Bernoulli processes

Consider a homogeneous CRM $\mu = \sum_{i \geq 1} \tau_i \delta_{W_i}$ on $\mathbb{W}$, where the $\tau_i$'s are non-negative and such that $\sum_{i \geq 1} \tau_i < +\infty$, and the $W_i$'s are i.i.d. and independent of the $\tau_i$'s. We denote by $\nu(\mathrm{d}s, \mathrm{d}w) = \lambda(s)\mathrm{d}sP(\mathrm{d}w)$ on $\mathbb{R}_+ \times \mathbb{W}$, with $\int_{\mathbb{R}_+} \min\{s, 1\}\lambda(s)\mathrm{d}s < +\infty$, the Lévy intensity measure of $\mu$. Let $\Delta_1 > \Delta_2 > \ldots$ be the decreasingly ordered $\tau_i$'s, and consider the discrete random measure

$$\mu_{\Delta_1} = \sum_{i \geq 1} \frac{\Delta_{i+1}}{\Delta_1} \delta_{W_{i+1}},$$

such that $\Delta_{i+1}/\Delta_1 \in (0, 1)$, for $i \geq 1$, and $\sum_{i \geq 1} \Delta_{i+1}/\Delta_1 < +\infty$. A SP on $\mathbb{W}$ is defined from $\mu_{\Delta_1}$ as follows. Let $F_{\Delta_1}(\mathrm{d}a) = \exp\left\{-\int_a^\infty \lambda(s)\mathrm{d}s\right\}\lambda(a)\mathrm{d}a$ be the distribution of $\Delta_1$ [Ferguson and Klass, 1972, pg. 1636], and let $G_a$ be the conditional distribution of $(\Delta_{i+1}/\Delta_1)_{i \geq 1}$ given $\Delta_1 = a$. Moreover, let $\Delta_{1,h}$ denote a random variable whose distribution has a density function $f_{\Delta_{1,h}}(a) = h(a)f_{\Delta_1}(a)$, where $h$ is a non-negative function and $f_{\Delta_1}$ is the density function of $F_{\Delta_1}$. If $(\rho_i)_{i \geq 1}$ are $(0, 1)$-valued random variables with distribution $G_{\Delta_{1,h}}$ then

$$\mu_{\Delta_{1,h}} = \sum_{i \geq 1} \rho_i \delta_{W_{i+1}}. \tag{3}$$

is a SP. For short, $\mu_{\Delta_{1,h}} \sim \mathrm{SP}(\nu, h)$. The law of $\mu_{\Delta_{1,h}}$ is a prior distribution for the parameter $\zeta$ of the Bernoulli process. The next proposition characterizes the predictive distribution of SP priors. See also James et al. [2015, Proposition 2.2] for a posterior analysis of SP priors.

**Proposition 2.** *Let $Z_{1:N}$ be a random sample from (1) with $\zeta \sim \mathrm{SP}(\nu, h)$. If $Z_{1:N}$ displays $K_N = k$ distinct features $\{W_1^*, \ldots, W_{K_N}^*\}$, each feature $W_i^*$ appearing exactly $M_{N,i} = m_i$ times, then the conditional distribution of $\Delta_{1,h}$, given $Z_{1:N}$, has a density function of the form*

$$g_{\Delta_{1,h} \mid Z_{1:N}}(a) \propto \frac{\prod_{i=1}^k \int_0^1 s^{m_i}(1-s)^{N-m_i}a\lambda(as)\mathrm{d}s}{\exp\left\{\sum_{n=1}^N \int_0^1 s(1-s)^{n-1}a\lambda(as)\mathrm{d}s\right\}} f_{\Delta_{1,h}}(a). \tag{4}$$

*Moreover, the conditional distribution of $Z_{N+1}$, given $(\Delta_{1,h}, Z_{1:N})$, coincides with the distribution of*

$$Z_{N+1} \mid (\Delta_{1,h}, Z_{1:N}) \overset{d}{=} Z'_{N+1} + \sum_{i=1}^{K_N} A_{N+1,i}\delta_{W_i^*}, \tag{5}$$

*where: i) $Z'_{N+1} \mid \mu'_{\Delta_{1,h}} = \sum_{i \geq 1} A'_{N+1,i}\delta_{W'_i} \sim \mathrm{BeP}(\mu'_{\Delta_{1,h}})$ and $\mu'_{\Delta_{1,h}} \mid \Delta_{1,h} \sim \mathrm{CRM}(\nu'_{\Delta_{1,h}})$, with $\nu'_{\Delta_{1,h}}(\mathrm{d}s, \mathrm{d}w) = (1-s)^N \Delta_{1,h}\lambda(s\Delta_{1,h})\mathbb{1}_{(0,1)}(s)\mathrm{d}sP(\mathrm{d}w)$; ii) the $A_{N+1,i}$'s are independent Bernoulli random variables with parameters $J_i$'s, respectively, such that $J_i \mid \Delta_{1,h}$ is distributed according to the*

*density function $f_{J_i \mid \Delta_{1,h}}(s) \propto (1-s)^{N-m_i} s^{m_i} \Delta_{1,h} \lambda(\Delta_{1,h}s) \mathbb{1}_{(0,1)}(s)\mathrm{d}s$ for $i \geq 1$.*

See Appendix B for the proof of Proposition 2. The marginalization of (5) with respect to (4) leads to the predictive distribution of SP priors: i) $Z_{N+1}$ displays "new" features $W_i'$'s, and the posterior distribution of statistics of "new" features, given $Z_{1:N}$, is determined by the law of $(\Delta_{1,h}, Z'_{N+1})$; ii) $Z_{N+1}$ displays "old" features $W_i^*$'s, and the posterior distribution of statistics of "old" features, given $Z_{N+1}$, is determined by the law of $(\Delta_{1,h}, \sum_{1 \leq i \leq K_N} A_{N+1,i}\delta_{W_i^*})$. Because of (4) and (5), the law of $(\Delta_{1,h}, Z'_{N+1})$ may include the whole sampling information, depending on the specification of $\nu$ and $h$, and hence the posterior distribution of statistics of "new" features, given $Z_{1:N}$, also includes such an information. As a corollary of Proposition 2, the posterior distribution of the number of unseen features, given $Z_{1:N}$ and fixed prior's parameters, is a mixture of Poisson distributions that may include the whole sampling information; in particular, the amount of sampling information in the posterior distribution is uniquely determined by the mixing distribution, namely by the conditional distribution of $\Delta_{1,h}$, given $Z_{1:N}$. SP priors thus allow to enrich the Poisson posterior structure arising from CRM priors, in terms of both a more flexible distribution and the inclusion of more sampling information than the sole sample size $N$, though they may lead to unwieldy posterior inferences due to the marginalization with respect to (4).

The use of the sampling information in the predictive structure of SPs somehow resembles that of Poisson-Kingman (PK) models [Pitman, 2006]. PK models form a broad class of nonparametric priors for species sampling problems. The DP prior is a PK model whose predictive distribution is such that: i) the conditional probability that the $(N+1)$-th draw is a "new" species, given $N$ observable samples, depends only on the sample size; ii) the conditional probability that the $(N+1)$-th draw is an "old" species, given $N$ observable samples, depends on the sample size, the number of distinct species and their frequencies. Such a behaviour resembles that of CRM priors, i.e. Proposition 1. PK models allow to include more sampling information in the probability of discovering a "new species" arising under the DP prior, which typically determines a loss of the analytical tractability of posterior inferences for the number of unseen species [Bacallado et al., 2017]. Such a behaviour resembles that of SP priors, i.e. Proposition 2. The PYP prior is arguably the most popular PK model. It stands out for enriching the probability of discovering a "new" species arising under the DP prior, by including the sampling information on the number of distinct species, while maintaining the analytical tractability and interpretability of the DP prior.

# 3 Stable-Beta Scaled Process (SB-SP) priors for the unseen-features problem

In Section 2 we showed how SP priors allow to enrich the Poisson posterior structure of the number of unseen features arising under CRM priors, e.g. the Beta and the stable-Beta process priors. While this is an appealing property, it may lead to a lack of analytical tractability and interpretability of posterior inferences, thus making SP priors not of practical interest in applications. In this section, we introduce and investigate a peculiar SP prior, which is referred to as the SB-SP prior, and we show that: i) it leads to a negative Binomial posterior distribution for the number of unseen features, which

generalizes the Poisson distribution while maintaining its analytical tractability and interpretability; ii) it leads to a posterior distribution for the number of unseen features, which depends on the sampling through the sample size and the number of distinct features. The SB-SP prior thus provides a sensible trade-off between the enrichment of the Poisson posterior structure of the number of unseen features arising under CRM priors and the analytical tractability and interpretability of posterior inferences. In particular, we characterize the SB-SP prior as the sole SP prior for which the posterior distribution of the number of unseen features depends on the observable sample only through the sample size and the number of distinct features. The SB-SP may thus be considered as a natural counterpart of the PYP for the unseen-feature problem.

## 3.1 SB-SP priors for Bernoulli processes

Stable scaled processes (S-SP) [James et al., 2015] form a subclass of SPs, and hence their definition follows from Section 2. In particular, for any $\sigma \in (0, 1)$, let $\mu_\sigma$ be the $\sigma$-stable CRM on $\mathbb{W}$ [Kingman, 1975], which is characterized by the Lévy intensity measure $\nu_\sigma(\mathrm{d}s, \mathrm{d}w) = \lambda_\sigma(s)\mathrm{d}sP(\mathrm{d}w)$ on $\mathbb{R}_+ \times \mathbb{W}$, with $\int_{\mathbb{R}_+} \min\{s, 1\}\lambda_\sigma(s)\mathrm{d}s < +\infty$, where $\lambda_\sigma(s) = \sigma s^{-1-\sigma}$. We recall that the largest atom $\Delta_1$ of $\mu_\sigma$ is distributed according to the density function

$$f_{\Delta_1}(a) = \sigma a^{-1-\sigma} \exp\left\{-a^{-\sigma}\right\}. \tag{6}$$

That is, $\Delta_1 = E^{-1/\sigma}$, where $E$ denotes a negative exponential random variable with parameter 1. For any non-negative function $h$, a S-SP on $\mathbb{W}$ is defined as the SP with law $\mathrm{SP}(\nu_\sigma, h)$. S-SP priors generalizes the Beta process prior, which is recovered by setting $h$ to be the identity function, and then letting $\sigma \to 0$ [James et al., 2015]. The predictive distribution of $\zeta \sim \mathrm{SP}(\nu_\sigma, h)$ is obtained from Proposition 2. In the next theorem, we characterize the S-SP priors as the sole SP priors for which the conditional distribution of $\Delta_{1,h}$, given $Z_{1:N}$, depends on $Z_{1:N}$ only through the sample size $N$ and the number $K_N$ of distinct features in $Z_{1:N}$.

**Theorem 1.** *Let $Z_{1:N}$ be a random sample from (1) with $\zeta \sim \mathrm{SP}(\nu, h)$, and let $Z_{1:N}$ displays $K_N$ distinct features with corresponding frequencies $(M_{N,1}, \ldots, M_{N,K_N})$. Moreover, let $\nu(\mathrm{d}s, \mathrm{d}w) = \lambda(s)\mathrm{d}sP(\mathrm{d}w)$, and let $f_{\Delta_{1,h}}$ be the density function of $\Delta_{1,h}$. If $f_{\Delta_{1,h}} > 0$ on $\mathbb{R}_+$ and the functions $\lambda$ and $f_{\Delta_{1,h}}$ are continuously differentiable, then the conditional distribution of $\Delta_{1,h}$, given $Z_{1:N}$, depends on $Z_{1:N}$ only through $N$ and $K_N$ if and only if $\nu = \nu_\sigma$.*

See Appendix C for the proof of Theorem 1. We recall from Section 2 that the conditional distribution of $\Delta_{1,h}$, given $Z_{1,N}$, uniquely determines the amount of sampling information included in the posterior distribution of statistics of "new" features. Then, according to Theorem 1, S-SP priors are the sole SP priors for which the posterior distribution of the number of unseen features, given $Z_{1:N}$ and fixed prior's parameters, depends on $Z_{1:N}$ only through $N$ and $K_N$. As a corollary of Theorem 1, the Beta process prior is the sole S-SP prior for which the posterior distribution of statistics of "new" features depends on $Z_{1:N}$ only through $N$. Analogous predictive characterizations are well-known in species sampling problems, and they are typically referred to as "sufficientness" postulates' [Bacallado et al., 2017]. In particular, the DP prior is characterized as the sole species sampling prior for which

8

the conditional probability that the $(N + 1)$-th draw is a "new" species, given $N$ observable samples, depends only on the sample size [Regazzini, 1978]. Moreover, the PYP prior is characterized as the sole species sampling prior for which the conditional probability that the $(N + 1)$-th draw is a "new" species, given $N$ observable samples, depends only on the sample size and the number of distinct species in the sample [Zabell, 2005]. Theorem 1 provides a "sufficientness" postulates' in the context of feature sampling problems.

As a noteworthy example of S-SPs, we introduce the SB-SP. The SB-SP is a S-SP obtained by a suitable specification of the non-negative function $h$. In particular, for any $c, \beta > 0$ let

$$h_{c,\beta}(a) = \frac{\beta^{c+1}}{\Gamma(c+1)} a^{-c\sigma} \exp\left\{-(\beta - 1)a^{-\sigma}\right\}, \tag{7}$$

where $\Gamma(\cdot)$ denotes the Gamma function. Then a SB-SP on $\mathbb{W}$ is defined as the SP with law $SP(\nu_\sigma, h_{c,\beta})$. For short, we denote the law of a SB-SP by SB-SP$(\sigma, c, \beta)$. The SB-SP prior generalizes the Beta process prior, which is recovered by setting $c = 0$ and $\beta = 1$, and then letting $\sigma \to 0$. According to the construction of SPs, the distribution of $\Delta_{1,h_{c,\beta}}$ has a density function obtained by combining (6) and (7); this is a polynomial-exponential tilting of the density function (6). In particular, $\Delta_{1,h_{c,\beta}}^{-\sigma}$ is distributed as a Gamma distribution with shape $(c + 1)$ and rate $\beta$. Such a straightforward distribution for $\Delta_{1,h_{c,\beta}}$ is at the core of the analytical tractability of posterior inferences under the SB-SP prior; this fact will be clear in the application of the SB-SP prior to the problem of estimating the number of unseen features. The next proposition characterizes the predictive distribution of the SB-SP prior.

**Proposition 3.** *Let $Z_{1:N}$ be a random sample from (1) with $\zeta \sim$ SB-SP$(\sigma, c, \beta)$. If $Z_{1:N}$ displays $K_N = k$ distinct features $\{W_1^*, \ldots, W_{K_N}^*\}$, each feature $W_i^*$ appearing exactly $M_{N,i} = m_i$ times, then the conditional distribution of $\Delta_{1,h_{c,\beta}}$, given $Z_{1:N}$, has a density function of the form*

$$g_{\Delta_{1,h_{c,\beta}} \mid Z_{1:N}}(a) = \sigma \frac{(\beta + \gamma_0^{(N)})^{k+c+1}}{\Gamma(k + c + 1)} a^{-k\sigma - (c+1)\sigma - 1} e^{-a^{-\sigma}(\beta + \gamma_0^{(N)})}, \tag{8}$$

*where $\gamma_0^{(N)} = \sigma \sum_{1 \leq i \leq N} B(1 - \sigma, i)$, with $B(\cdot, \cdot)$ being the (Euler) Beta function. Moreover, the conditional distribution of $Z_{N+1}$, given $(\Delta_{1,h_{c,\beta}}, Z_{1:N})$, coincides with the distribution of*

$$Z_{N+1} \mid (\Delta_{1,h_{c,\beta}}, Z_{1:N}) \stackrel{d}{=} Z_{N+1}' + \sum_{i=1}^{K_N} A_{N+1,i} \delta_{W_i^*}, \tag{9}$$

*where:*

*i)* $Z_{N+1}' \mid \mu_{\Delta_{1,h_{c,\beta}}}' = \sum_{i \geq 1} A_{N+1,i}' \delta_{W_i'} \sim \text{BeP}(\mu_{\Delta_{1,h_{c,\beta}}}')$ *such that $\mu_{\Delta_{1,h_{c,\beta}}}' \mid \Delta_{1,h_{c,\beta}} \sim \text{CRM}(\nu_{\Delta_{1,h_{c,\beta}}}')$, with*

$$\nu_{\Delta_{1,h_{c,\beta}}}'(\mathrm{d}s, \mathrm{d}w) = \Delta_{1,\Delta_{1,h_{c,\beta}}}^{-\sigma}(1 - s)^N \sigma s^{-1-\sigma} \mathbb{1}_{(0,1)}(s)\mathrm{d}s P(\mathrm{d}w);$$

*ii) the $A_{N+1,i}$'s are independent Bernoulli random variables with parameters $J_i$'s, respectively, such*

*that each $J_i \,|\, \Delta_{1,h_{c,\beta}}$ is distributed according to a density function of the form*

$$f_{J_i \,|\, \Delta_{1,h_{c,\beta}}}(s) = \frac{1}{B(m_i - \sigma, N - m_i + 1)} s^{m_i - \sigma}(1-s)^{N-m_i+1} \mathbb{1}_{(0,1)}(s).$$

See Appendix C for the proof of Proposition 3. According to Equation (8), the conditional distribution of $\Delta_{1,h_{c,\beta}}$, given $Z_{1:N}$, depends on $Z_{1:N}$ only through the sample size $N$ and the number $K_N$ of distinct features in $Z_{1:N}$. This agrees with Theorem 1, implying that the posterior distribution of the number of unseen features, given $Z_{1:N}$ and fixed prior's parameters, depends on $Z_{1:N}$ only through $N$ and $K_N$. Because of (8) and (9), the posterior distribution of statistics of "new" features stands out for analytical tractability, thus being competitive with that arising from CRMs, e.g. the Beta and the stable-Beta processes. In particular, from Equation (9), the conditional distribution of $Z'_{N+1}$, given $(\Delta_{1,h_{c,\beta}}, Z_{1:N})$ is a Poisson distribution that depends on $Z_{1:N}$ only through $N$. Then, from (8), its marginalization with respect to the conditional distribution of $\Delta_{1,h_{c,\beta}}$, given $Z_{1:N}$, leads to a negative Binomial posterior distribution. Such an appealing property arises from the peculiar form $h_{c,\beta}$ that, combined with $\nu_\sigma$, leads to a conjugacy property for the conditional distribution of $\Delta_{1,h_{c,\beta}}$, given $Z_{1:N}$. That is, the conditional distribution of $\Delta_{1,h_{c,\beta}}^{-\sigma}$, given $Z_{1:N}$, is a Gamma distribution with shape $(K_N + c + 1)$ and rate $\beta + \gamma_0^{(N)}$, which is the distribution $\Delta_{1,h_{c,\beta}}^{-\sigma}$ with shape and rate being updated through $Z_{1:N}$. The next proposition establishes the distribution of a random sample $Z_{1:N}$ from a SB-SP prior. See Appendix C for details.

**Proposition 4.** *Let $Z_{1:N}$ be a random sample from (1) with $\zeta \sim \text{SB-SP}(\sigma, c, \beta)$. The probability that $Z_{1:N}$ displays a particular feature allocation of $k$ distinct features with frequencies $(m_1, \ldots, m_k)$ is*

$$p_k^{(N)}(m_1, \ldots, m_k) = \frac{\frac{\sigma^k \beta^{c+1}}{(\beta + \gamma_0^{(N)})^{k+c+1}}}{\frac{\Gamma(c+1)}{\Gamma(k+c+1)}} \prod_{i=1}^{k} \frac{\Gamma(m_i - \sigma)\Gamma(N - m_i + 1)}{\Gamma(N - \sigma + 1)}. \tag{10}$$

## 3.2 BNP inference for the unseen-features problem

Now, we apply the SB-SP prior to the unseen-features problem. For any $N \geq 1$ let $Z_{1:N}$ be an observable sample modeled as the BNP Bernoulli model (1), with $\zeta \sim \text{SB-SP}(\sigma, c, \beta)$. Moreover, under the same model of the $Z_n$'s, for any $M \geq 1$ let $(Z_{N+1}, \ldots, Z_{N+M})$ be additional unobservable sample. Then, the unseen-feature problem calls for the estimation of

$$U_N^{(M)} = \sum_{i \geq 1} \mathbb{1}\left(\sum_{m=1}^{M} A_{N+m,i} > 0\right) \mathbb{1}\left(\sum_{n=1}^{N} A_{n,i} = 0\right), \tag{11}$$

namely the number of hitherto unseen features that would be observed in $(Z_{N+1}, \ldots, Z_{N+M})$. As generalization of the unseen-feature problem (11), for $r \geq 1$ we consider the estimation of

$$U_N^{(M,r)} = \sum_{i \geq 1} \mathbb{1}\left(\sum_{m=1}^{M} A_{N+m,i} = r\right) \mathbb{1}\left(\sum_{n=1}^{N} A_{n,i} = 0\right), \tag{12}$$

namely the number of hitherto unseen features that would be observed with prevalence $r$ in $(Z_{N+1}, \ldots, Z_{N+M})$. Of special interest is $r = 1$, which concerns rare (unique) features. The next theorem characterizes the posterior distributions of $U_N^{(M)}$ and $U_N^{(M,r)}$, given $Z_{1:N}$. We denote by $\text{NegativeBinonial}(n, p)$ the negative Binomial distribution with parameter $n$ and $p \in (0, 1)$.

**Theorem 2.** *Let $Z_{1:N}$ be a random sample from* (1) *with $\zeta \sim \text{SB-SP}(\sigma, c, \beta)$, and let $Z_{1:N}$ displays $K_N = k$ distinct features with frequencies $(M_{N,1}, \ldots, M_{N,K_N}) = (m_1, \ldots, m_k)$. Then, the posterior distributions of $U_N^{(M)}$ and of $U_N^{(M,r)}$, given $Z_{1:N}$, coincide with the distributions of*

$$U_N^{(M)} \mid Z_{1:N} \sim \text{NegativeBinonial}\left(K_N + c + 1, \frac{\gamma_N^{(M)}}{\beta + \gamma_0^{(N+M)}}\right), \tag{13}$$

*and*

$$U_N^{(M,r)} \mid Z_{1:N} \sim \text{NegativeBinonial}\left(K_N + c + 1, \frac{\rho_N^{(M,r)}}{\beta + \gamma_0^{(N)} + \rho_N^{(M,r)}}\right), \tag{14}$$

*for any index of prevalence $r \geq 1$, respectively, where $\gamma_N^{(M)} = \sigma \sum_{1 \leq i \leq M} B(1 - \sigma, N + i)$ and where $\rho_N^{(M,r)} = \binom{M}{r}\sigma B(r - \sigma, N + M - r + 1)$, with $B(\cdot, \cdot)$ denoting the (Euler) Beta function.*

See Appendix D for the proof of Theorem 2. The posterior distributions (13) and (14) depend on $Z_{1:N}$ through the sample size $N$ and the number $K_N$ of distinct features. This is in contrast with the corresponding posterior distributions obtained under the Beta and the stable-Beta process priors, which are Poisson distributions that depend on $Z_{1:N}$ only through $N$ [Masoero et al., 2021, Proposition 1]. BNP estimators of $U_N^{(M)}$ and $U_N^{(M,r)}$, with respect to a squared loss function, are obtained as the posterior expectations of (13) and (14), i.e.

$$\hat{U}_N^{(M)} = (K_N + c + 1) \frac{\gamma_N^{(M)}}{\beta + \gamma_0^{(N+M)} - \gamma_N^{(M)}} \tag{15}$$

and

$$\hat{U}_N^{(M,r)} = (K_N + c + 1) \frac{\rho_N^{(M,r)}}{\beta + \gamma_0^{(N)}} \tag{16}$$

respectively. The estimators (15) and (16) are simple, linear in the sampling information and computationally efficient. In the next theorem we establish the large $M$ asymptotic behaviour of the posterior distributions (13) and (14), showing that the number of unseen features has a power-law growth in $M$. The same growth in $M$ holds under the stable-Beta process prior [Masoero et al., 2021, Proposition 2], though the limiting distribution is degenerate.

**Theorem 3.** *Let $Z_{1:N}$ be a random sample from* (1) *with $\zeta \sim \text{SB-SP}(\sigma, c, \beta)$, and let $Z_{1:N}$ displays $K_N = k$ distinct features with frequencies $(M_{N,1}, \ldots, M_{N,K_N}) = (m_1, \ldots, m_k)$. As $M \to +\infty$*

$$\frac{U_N^{(M)}}{M^\sigma} \mid Z_{1:N} \xrightarrow{\text{a.s.}} W_N, \tag{17}$$

11

where $W_N$ is a Gamma random variable with shape $(K_N + c + 1)$ and rate $(\beta + \gamma_0^{(N)})/\Gamma(1-\sigma)$, and

$$\frac{U_N^{(M,r)}}{M^\sigma} \mid Z_{1:N} \xrightarrow{\text{a.s.}} W_{N,r}, \tag{18}$$

where $W_{N,r}$ is a Gamma random variable with shape $(K_N+c+1)$ and rate $\Gamma(r+1)(\beta+\gamma_0^{(N)})/\sigma\Gamma(r-\sigma)$.

## 4   Experiments

Over the last decade, genomics has witnessed an extraordinary improvement in the data availability due to the advent of next generation sequencing technologies. Thanks to larger and richer datasets, researchers have started uncovering the role and impact of rare genetic variants in heritability and human disease [Hernandez et al., 2019, Momozawa and Mizukami, 2020]. The development of methods for estimating the number of new genomic variants to be observed in future studies is an active research area, as it can aid the design of effective clinical procedures in precision medicine [Ionita-Laza et al., 2009, Zou et al., 2016], enhance understanding of cancer biology [Chakraborty et al., 2019], and help to optimize sequencing procedures [Rashkin et al., 2017, Masoero et al., 2021]. Here, we consider datasets of individual genomic sequences. Following common practice, we assume that an underlying fixed and idealized genomic sequence (the "reference") is given. Then, each coordinate of an individual sequence reports the presence (1) or absence (0) of variation at a given locus with respect to the reference. All variants are treated equally, namely, any expression differing from the underlying reference at a given locus counts as a variant. We make use of our methodology to estimate the number of genomic loci at which variation was not observed in the original sample, and is going to be observed in (at least one of) $M$ additional datapoints.

We find in our experiments that the estimates of the total number of new variants to be observed produced using the SB-SP-Bernoulli model, hereafter referred to as SSB, tend to be more accurate than other available methods in the literature. This phenomenon is particularly evident when the sample size $N$ of the training set is small, and when the extrapolation size $M$ is large with respect to $N$. Moreover, the SSB model is particularly effective in estimating the number of new rare variants, e.g. variants appearing only once in the additional unobservable samples. Accurate estimation of rare variants is particularly important, as these are believed to be largely responsible for heritability of human disease [Rashkin et al., 2017, Chakraborty et al., 2019]. To benchmark the quality of the SSB, we consider a number of competing methodologies for the feature prediction problem available in the literature: i) Jackknife estimators (J) [Gravel, 2014]; ii) a linear programming method (LP) [Zou et al., 2016] and variations of Good-Toulmin estimators (GT) [Chakraborty et al., 2019]. We also compare our empirical findings to a BNP estimator obtained under the stable-Beta process prior (3BB), which has been introduced in Masoero et al. [2021]. We complete our analysis with a thorough investigation on synthetic data in Appendix F and Appendix G, as well as on additional real data from the gnomAD database [Karczewski et al., 2020] in Appendix H.

## 4.1 Empirics and evaluation metrics

For the SSB method to be useful, we need to estimate the underlying, unknown, parameters of the SB-SP prior. To learn these prior's parameters, we here adopt an empirical Bayes procedure, which consists in maximizing the marginal distribution (10). In particular, we maximize numerically Equation (10) with respect to the parameters $\beta > 0$, $c > 0$ and $\sigma \in (0, 1)$ of the SB-SP prior, and use the resulting values to produce our estimators. That is, we let

$$(\hat{\beta}, \hat{c}, \hat{\sigma}) = \arg \max_{(\beta, c, \sigma)} \left\{ p_k^{(N)}(m_1, \cdots, m_k) \right\},$$

and plug these values in the BNP estimator (13) and (14). The resulting values provide our BNP estimates of the number $U_N^{(M)}$ of new variants and the number $U_N^{(M,r)}$ of new variants with prevalence $r$.

To assess the accuracy of our estimates, we consider the percent deviation of the estimate from the truth to be the achieved accuracy. That is, the accuracy of the estimator $\hat{U}_N^{(M)}$ is defined as

$$v_N^{(M)} := 1 - \min \left\{ \frac{|U_N^{(M)} - \hat{U}_N^{(M)}|}{U_N^{(M)}}, 1 \right\}. \tag{19}$$

In particular, the accuracy $v_N^{(M)}$ equals 1 when the estimate is perfect (no error is incurred), and decreases to 0 as the estimate deviates from the truth. The min operator in (19) ensures that $v_N^{(M)}$ lies in $[0, 1]$: we let the accuracy to be equal to 0 whenever there is a severe overestimation, and the percentage estimation error exceeds 100%, i.e. when $\hat{U}_N^{(M)} \geq 2 \times U_N^{(M)}$. The SSB, 3BB and LP methods also offer an estimate for the number of new features observed with a given prevalence $r$. We let $v_N^{(M,r)}$ be the accuracy metric, where we replace in (19) $U_N^{(M)}$ with $U_N^{(M,r)}$, the number of new features observed with prevalence $r$, and $\hat{U}_N^{(M)}$ with $\hat{U}_N^{(M,r)}$.

## 4.2 Estimating the number of new variants in cancer genomics

Following the empirical study of Chakraborty et al. [2019], we make use of data from the Cancer Genome Atlas (TCGA), the largest publicly available cancer genomics dataset, containing somatic mutations from 10,295 patients and spanning 33 different cancer types. We partition the samples into 33 smaller datasets according to cancer-type annotation of each patient. See Chakraborty et al. [2019] and Masoero et al. [2021, Appendix F] for details on the data and the experimental setup. For each cancer type, we retain a small fraction of the data for purposes of training, and consider the task of estimating the number of new variants that will be observed in a follow-up sample given a pilot sample. We validate our estimates by comparing the estimate $\hat{U}_N^{(M)}$ of the number of distinct variants to the true value, obtained by extrapolating to the remaining data. To assess the variability and error in our estimates, we repeat for every cancer type the experiment on $S = 1,000$ subsets of the data, each obtained by randomly subsampling without replacement from the full sample.

We find that the SSB and 3BB methods perform particularly well when the training sample size $N$ is small compared to the extrapolation sample size $M$. This setting is relevant in the context of cancer genomics, as scientists are interested in understanding the "unexploited potential" of the

genetic information, especially for rare cancer subtypes [Chakraborty et al., 2019, Huyghe et al., 2019]. To compare and quantify the performance of the available methodologies in this setting, we report in Figure 1 the distribution of the estimation accuracy when retaining only $N = 10$ samples for training and extrapolating to the largest possible sample size $M$ for which we can compute the accuracy metric (Equation (19)). We report results for the 10 cancer types with the largest number of samples in the original dataset. For each cancer type and for each method, the distribution of the estimation accuracy is obtained by considering its performance across the $S = 1{,}000$ replicates. Across all cancer types, the estimates obtained from the SSB method achieve higher accuracy.

We show in Figure 2 the behavior of $\hat{U}_N^{(i)}$ for five different cancer types as $i = 1, \ldots, M$. Again, we let $N = 10$, and $M$ be the largest possible extrapolation value, as dictated by the dataset size. We report the estimates obtained from a fixed sample of size $N = 10$, as well as the variability around such estimates obtained by re-fitting each model, iteratively leaving one datapoint out from the sample. In this setting, the SSB method outperforms competing methods in terms of estimation accuracy. Moreover, the variability in the estimates arising from re-fitting the model on subsets of the data provides a useful measure of uncertainty in such estimation.



Figure 1: Estimation accuracy $v_N^{(M)}$ for the number of new genomic variants $\hat{U}_{10}^{(M)}$. For each method and each cancer type, we retain $N = 10$ random samples and use them to estimate up to $M$ total observations, where $N + M$ is the size of the original sample.

## 4.3 Estimating the number of new rare variants in cancer genomics

In recent years, the cancer genomics research community has become increasingly interested in studying and understanding the role of extremely rare variants, such as singletons, i.e. observed in only
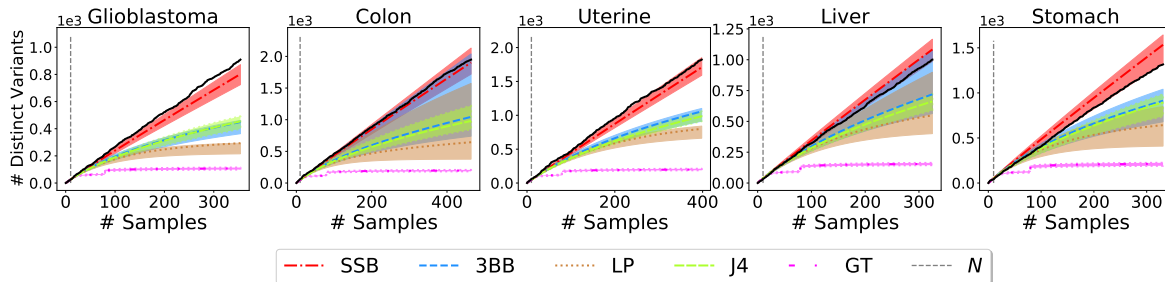
Figure 2: Estimation of the number of new genomic variants $\hat{U}_N^{(i)}$, for $i = 1, \ldots, M$. For each method and cancer type, we retain $N = 10$ random samples and use them to estimate up to the largest possible size. We fit each model on the full sample, as well as $N = 10$ additional times by iteratively leaving one datapoint out from the training sample. The solid black line is the true number of features that would have been observed (vertical axis) for any extrapolation size $N + M$ (horizontal axis), for a fixed ordering of the data. Shaded regions report the prediction range obtained from the estimates from the leave-one-out fits.

one patient. Evidence suggests that rare deleterious variants can have far stronger effect sizes than common variants [Rasnic et al., 2020] and can play an important role in the development of cancer. For example, in breast cancer, it is well accepted that the risk of a variant is inversely proportional with respect to its prevalence: the rarer the variant, the higher the risk [Wendt and Margolin, 2019]. Therefore, effective identification and discovery of rare variants is an active, is an ongoing research area [Lawrenson et al., 2016, Lee et al., 2019]. This phenomenon is not limited to breast cancer, but is progressively being studied across different cancer types. See, e.g. the recent works on ovarian [Phelan et al., 2017], skin [Goldstein et al., 2017], prostate [Nguyen-Dumont et al., 2020] and lung [Liu et al., 2021] cancers and references therein. In downstream analysis, these estimates could be useful for planning and designing future experiments, e.g. informing scientists on the number of new samples to be collected in order to observe a target number of new variants, or for power analysis considerations in rare variants association tests [Rashkin et al., 2017].

The BNP framework considered here allows us to estimate the number of new rare variants to be discovered. While Zou et al. [2016] did not consider the problem of estimating rare variants, it is straightforward to obtain an estimate for this quantity from their framework. Indeed, for every prevalence $x \in [0, 1]$, the LP estimates the histogram $h(x)$, which counts the number of variants appearing with prevalence $x$ in the population, and the number of variants appearing with prevalence $r$ follows from the binomial sampling model assumption, namely $\hat{U}_N^{(M,r)} = \sum_x h(x) \left\{ \binom{N+M}{r} x^r (1-x)^{N+M-r} - \binom{N}{r} x^r (1-x)^{N-r} \right\}$. We show in Figure 4 that the SSB method provides better estimates than the 3BB and LP methods.

## 4.4 Coverage and calibrated uncertainties

One of the benefits of the BNP approach is that it automatically yields a notion of variability of the estimate of $U$ via posterior credible intervals. We here check whether these intervals produce a useful notion of uncertainty, by investigating their calibration. For $\alpha \in (0, 1)$, we say that a $100 \times \alpha\%$ credible interval is calibrated if it contains the true value of interest, arising from hypothetical repeated draws, $100 \times \alpha\%$ of the times. We here assess the calibration of a $100 \times \alpha\%$ credible interval for $U_N^{(M)}$ conditionally given $Z_{1:N}$ as follows. Let $S$ be a large number ($S = 1,000$ in our experiments). For
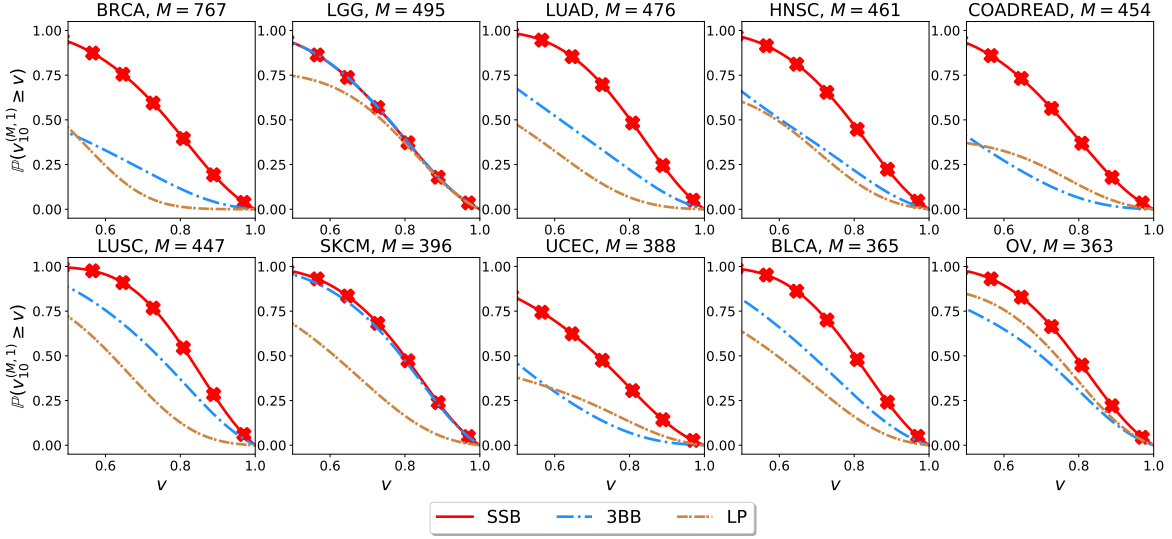
Figure 3: Estimation accuracy $v_N^{(M,1)}$ for new variants appearing with prevalence one in future unobservable samples for different cancer types. For each method and each cancer, we retain $N = 10$ random samples and use them to estimate up to the largest possible size.
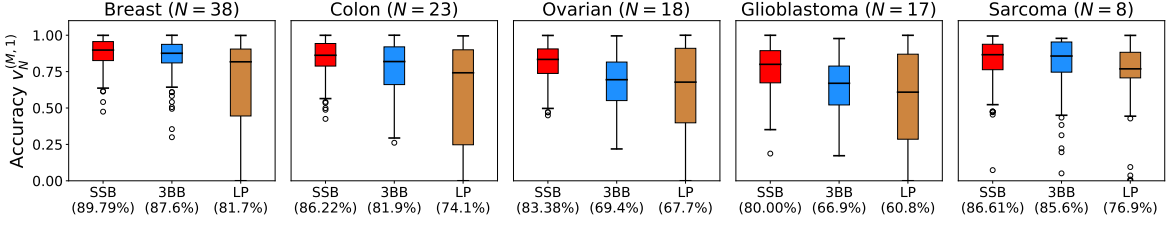


Figure 4: Estimation accuracy $v_N^{(M,1)}$ for new variants appearing with prevalence one in future samples. For each method and different cancer types, we retain a random sample of size $N = 5\%$ of the available dataset, and use it to estimate up to the largest possible size.

each $s = 1, \ldots, S$, we retain a random subset of the data of size $N$, and estimate the corresponding parameters $\hat{\beta}, \hat{c}, \hat{\sigma}$ as discussed in Section 4.1. Then, we let $\hat{W}_{N,s,low}^{(M)}(\alpha), \hat{W}_{N,s,hi}^{(M)}(\alpha)$ be the endpoints of a $100 \times \alpha\%$ credible interval for the distribution of the number of new features, as given by Equation (13), centered around the posterior predictive mean. We compute coverage calibration via

$$w_N^{(M)}(\alpha) = \frac{1}{S} \sum_{s=1}^{S} \mathbb{1} \left\{ \hat{W}_{N,s,low}^{(M)}(\alpha) \leq K_{N+M} \leq \hat{W}_{N,s,hi}^{(M)}(\alpha) \right\}.$$

This is the fraction of the $S$ experiments in which the true value was contained by an $100 \times \alpha\%$ credible interval. The closer $w_N^{(M)}(\alpha)$ to $\alpha$, the better calibrated the credible intervals. We compute the same quantity for the 3BB method using the results in Masoero et al. [2021]. Although still not perfect, we find that the posterior predictive intervals obtained from the SSB method are better calibrated than the ones under the 3BB method (see Figure 5).
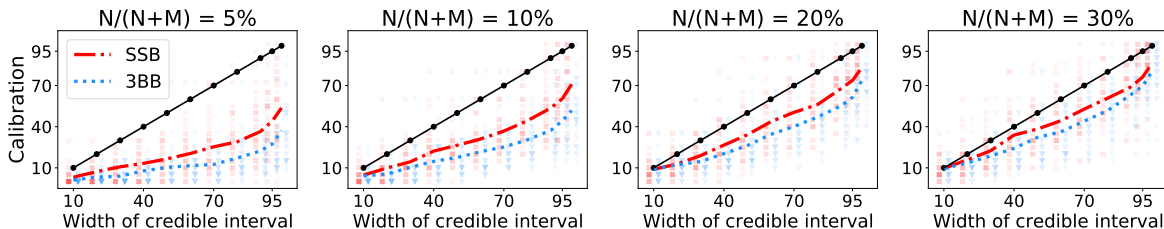
16

Figure 5: Coverage calibraiton of BNP estimators for number of new variants in future samples across all cancer types in TCGA. Different subplots refer to different ratios of the training $N$ with respect to the extrapolation $M$. For each cancer, we retain a training sample of size $N \in \{5\%, 10\%, 20\%, 30\%\}$ of the total available dataset, and extrapolate up to the largest available $M$. Colored lines report the average coverage $w_N^{(M)}(\alpha)$ across all cancer types ($y$-axis) as a function of $\alpha$ ($x$-axis). Faded dots refer to coverage for individual cancer types.

## 5   Discussion

Masoero et al. [2021] first applied CRM priors to the unseen-features problem, showing that: i) despite the broadness of the class of CRM priors, all CRM priors lead to the same Poisson posterior structure for the number of unseen features, which thus makes them not a flexible prior model for the unseen-features problem; ii) while the Poisson posterior distribution may be appealing in principle, making the posterior inferences analytically tractable and of easy interpretability, its independence from $Z_{1:N}$ makes the BNP approach a questionable oversimplification, with posterior inferences being completely determined by the estimation of unknown prior's parameters. In this paper, we introduced the SB-SP prior, and showed that: i) it enriches the posterior distribution of the number of unseen features arising under CRM priors, which results in a negative Binomial distribution whose parameters depend on the sample size and the number of distinct features; ii) it maintains the same analytical tractability and interpretability as CRM priors, which results in BNP estimators that are simple, linear in the sampling information and computationally efficient. The effectiveness of the SB-SP prior is showcased through an empirical analysis on synthetic and real data. Under the SB-SP prior, we found that estimates of the unseen number of features are accurate, and they outperform the most popular competitors in the challenging scenario where the sample size $N$ is particularly small, and also small with respect to the extrapolation size $M$.

Our approach admits an extension to the multiple-feature setting, which takes into account of the many forms of variation, e.g. single nucleotide changes, tandem repeats, insertions and deletions, copy number variations [Zou et al., 2016]. We briefly describe the multiple-feature setting, and defer to Appendix E for details. It is assumed that a feature $w_i$ comes with a characteristic, i.e. the form of variation, chosen among $q > 1$ characteristics. For $N \geq 1$, the observable sample $\boldsymbol{Z}_{1:N} = (\boldsymbol{Z}_1, \ldots, \boldsymbol{Z}_N)$ is modeled as a $\{0,1\}^q$-valued stochastic process $\boldsymbol{Z} = \sum_{i \geq 1} \boldsymbol{A}_i \delta_{w_i}$, where $\boldsymbol{A}_i := (A_{i,1}, \ldots, A_{i,q})$ is a Multinomial random variable with parameter $\boldsymbol{p}_i = (p_{i,1}, \ldots, p_{i,q})$ such that $|\boldsymbol{p}_i| = \sum_{1 \leq j \leq q} p_{i,j} < 1$, and the $\boldsymbol{A}_i$'s are i.i.d. That is, for any $i \geq 1$ all the $A_{i,j}$'s are equal to 0 with probability $(1 - |\boldsymbol{p}_i|)$, i.e. $w_i$ does not display variation, or only one $A_{i,j}$'s is equal to 1 with probability $p_{i,j}$, i.e. $w_i$ displays variation with characteristic $j$. $\boldsymbol{Z}$ is a multivariate Bernoulli process with parameter $\boldsymbol{\zeta} = \sum_{i \geq 1} \boldsymbol{p}_i \delta_{w_i}$. The stable-Beta-Dirichlet process prior for $\boldsymbol{\zeta}$ is a multivariate generalization of the stable-Beta process prior [James, 2017], and it leads to a Poisson posterior distribution for the number of unseen features,

given $\boldsymbol{Z}_{1:N}$, which depends on $\boldsymbol{Z}_{1:N}$ only through $N$. In Appendix E we introduce a scaled version of the stable-Beta-Dirichlet process, and show that it leads to a negative Binomial posterior distribution for the number of unseen features, which depends on $\boldsymbol{Z}_{1:N}$ through $N$ and the number of distinct features in $\boldsymbol{Z}_{1:N}$.

SP priors have been introduced in James et al. [2015] and, to the best of our knowledge, since then no other works have further investigated such a class of priors. To date, the peculiar predictive properties of SP priors appear to be unknown in the BNP literature. Our work on the unseen-features problem is the first to highlight the great potential of SP priors in BNPs, showing that they provide a critical tool for enriching the predictive structure of the popular CRM priors [James, 2017, Broderick et al., 2018]. We believe that SPs may be of interest beyond the unseen-features problem, and more generally beyond the broad class of feature sampling problems. CRM priors, and in particular the Beta and stable-Beta process priors, have been widely used in several contexts, with a broad range of applications in topic modeling, analysis of social networks, binary matrix factorization for dyadic data, analysis of choice behaviour arising from psychology and marketing surveys, graphical models, and analysis of similarity judgement matrices. See Griffiths and Ghahramani [2011] and references therein for details. In all these contexts, SP priors may be more effective than CRM priors, as they allow to better exploit the sampling information in posterior inferences.

Among applications of SP priors beyond features sampling problems, it is worth mentioning the use of SP priors as hierarchical (or latent) priors in models of unsupervised learning [Griffiths and Ghahramani, 2011, Section 5], the most popular being Gaussian latent feature modeling. Differently from features sampling problems, where the values of features' labels $W_i$s are immaterial, in Gaussian latent feature modeling the values the $W_i$'s become material. That is, under the Gaussian latent feature model with a SP prior, observations are assumed to modeled as a multivariate Gaussian distribution, whose mean depends on latent features that are modeled with a SP prior, thus making the values of features' labels $W_i$'s of critical importance for the analysis. Bayesian factor analysis [Knowles and Ghahramani, 2011] provides another context where SP priors may be usefully applied as hierarchical priors. Within the context of factor analysis, we also mention the work of Ayed and Caron [2021] with applications to network analysis. There, the authors exploit CRM priors to recover the latent community structure in a network between individuals, and the features' labels describe the level of affiliation of a certain individual to a latent community. In such a context, we believe that SP priors may be used in place of CRM priors, with the advantage of introducing richer predictive structure. In this respect, our work paves the way to promising directions of future research, in terms of both methods and applications.

## Acknowledgement

# A   A brief account on completely random measures

In this section we provide a short account on completely random measures (CRMs). For a more exhaustive treatment refer to Daley and Vere-Jones [2008], Kingman [1992]. Let us denote by $\mathbb{W}$ a Polish space equipped with its Borel $\sigma$-field $\mathcal{W}$, and we also indicate by $\mathcal{B}_{\mathbb{R}_+}$ the Borel $\sigma$-field of the positive real line $\mathbb{R}_+$. Denote by $\mathsf{M}_{\mathbb{W}}$ the space of all bounded and finite measures on $(\mathbb{W}, \mathcal{W})$, in other words $\mu \in \mathsf{M}_{\mathbb{W}}$ iff $\mu(A) < +\infty$ for any bounded set $A \in \mathcal{W}$. The space $\mathsf{M}_{\mathbb{W}}$ is usually assumed to be equipped with a proper Borel $\sigma$-algebra, which is induced by the so called *weak-hash convergence* and denoted here as $\mathcal{M}_{\mathbb{W}}$ (see Daley and Vere-Jones [2008] for details).

**Definition 1.** *A Completely Random Measure (CRM) $\mu$ on $(\mathbb{W}, \mathcal{W})$ is a random element defined on a suitable probability space and taking values in $(\mathsf{M}_{\mathbb{W}}, \mathcal{M}_{\mathbb{W}})$ such that the random variables $\mu(A_1), \ldots, \mu(A_n)$ are independent for any choice of bounded and disjoint sets $A_1, \ldots, A_n \in \mathcal{W}$ and for any $n \geq 1$.*

Kingman [1967] proved that a CRM may be decomposed as the sum of three main components: i) a deterministic drift $u$, namely a deterministic measure on $(\mathbb{W}, \mathcal{W})$; ii) a part with random jumps $(\tau_i)_{i \geq 1}$ at random locations $(W_i)_{i \geq 1}$, denoted here as $\mu_c = \sum_{i \geq 1} \tau_i \delta_{W_i}$; iii) a component with random jumps $(\eta_i)_{i \geq 1}$ at fixed locations $w_1, w_2, \ldots \in \mathbb{W}$. That is to say

$$\mu(\,\cdot\,) = u(\,\cdot\,) + \mu_c(\,\cdot\,) + \sum_{i \geq 1} \eta_i \delta_{w_i}(\,\cdot\,). \tag{20}$$

See Daley and Vere-Jones [2008] for a proof.

Following standard practice in the nonparametric literature, in this paper we deal with CRMs without deterministic drift and without fixed atoms, namely we assume that $\mu \equiv \mu_c$. In this case $\mu = \mu_c$ is characterized through the Lévy-Khintchine representation of its Laplace functional:

$$\mathbb{E}\left[e^{-\int_{\mathbb{W}} f(w)\mu_c(\mathrm{d}w)}\right] = \exp\left\{-\int_{\mathbb{R}_+ \times \mathbb{W}} (1 - e^{-sf(w)})\nu(\mathrm{d}s, \mathrm{d}w)\right\}, \tag{21}$$

for any measurable function $f : \mathbb{W} \to \mathbb{R}_+$, where $\nu$ is a measure on $\mathbb{R}_+ \times \mathbb{W}$ and it is referred to as the Lévy intensity of the CRM $\mu_c$. The measure $\nu$ is also required to satisfy the following conditions

$$\nu(\mathbb{R}_+ \times \{w\}) = 0 \quad \forall w \in \mathbb{W}, \quad \text{and} \quad \int_{\mathbb{R}_+ \times A} \min\{s, 1\}\nu(\mathrm{d}s, \mathrm{d}w) < \infty$$

for any bounded $A \in \mathcal{W}$. The representation (21) is of paramount importance to prove all our posterior results, and it clarifies the pivotal role of $\nu$ in the determination of the distributional properties of $\mu_c$. Kallenberg [2010] provides a very general decomposition for such a measure $\nu$ as follows: $\nu(\mathrm{d}s, \mathrm{d}w) = \lambda_w(\mathrm{d}s)\Lambda(\mathrm{d}w)$, where $\Lambda$ is a $\sigma$-finite measure on $(\mathbb{W}, \mathcal{W})$ and $\lambda_w$ is a transition kernel, i.e., $w \to \lambda_w(A)$ is $\mathcal{W}$-measurable for all Borel sets $A \in \mathcal{B}_{\mathbb{R}_+}$ and $A \to \lambda_w(A)$ is a measure on $(\mathbb{R}_+, \mathcal{B}_{\mathbb{R}_+})$. When $\lambda_w(\mathrm{d}s) \equiv \lambda(\mathrm{d}s)$ does not depend on $w \in \mathbb{W}$, we say that the CRM is homogeneous, which is

tantamount to saying that the atoms $W_i$'s and the jumps $\tau_i$'s are independent random variables. In BNP problems, it is common to suppose that $\Lambda(\mathrm{d}w) = \alpha P(\mathrm{d}w)$, where $P$ is a probability measure on $(\mathbb{W}, \mathcal{W})$ and $\alpha > 0$. Two remarkable examples of CRMs are the $\sigma$-stable process, which can be recovered by choosing $\lambda(\mathrm{d}s) = \sigma s^{-1-\sigma}\mathrm{d}s$, and the gamma process, which corresponds to the choice $\lambda(\mathrm{d}s) = e^{-s}/s\,\mathrm{d}s$. See also [Lijoi and Prünster, 2010] for additional details and connections with the BNP literature.

In Section E, we will make use of multivariate CRMs to define a multivariate extension of the Bernoulli process model, called the Bernoulli process model with a condiment. For this reason we now specify what we mean for a multivariate CRM. A vector $\boldsymbol{\mu} = (\mu_1, \ldots, \mu_q)$ of completely random measures is said to be a multivariate CRM if the random variables

$$(\mu_1(A_1), \ldots, \mu_q(A_1)), \ldots, (\mu_1(A_n), \ldots, \mu_q(A_n))$$

are independent for any choice of bounded and disjoint Borel sets $A_1, \ldots, A_n \in \mathcal{W}$ and for any $n \geq 1$. A decomposition similar to the one stated in Equation (20) holds true for multivariate CRMs as well [Kallenberg, 2010]. In the present paper we focus on multivariate CRMs which are functionals of marked Poisson point processes on $\mathbb{R}_+^q \times \mathbb{W}$, i.e.,

$$\boldsymbol{\mu} = \sum_{i \geq 1} \boldsymbol{\tau}_i \delta_{W_i},$$

where $(\boldsymbol{\tau}_i)_{i \geq 1}$ are random jumps in $\mathbb{R}_+^q$ and $(W_i)_{i \geq 1}$ is a sequence of random atoms in $\mathbb{W}$. Such a multivariate CRM has the following Lévy-Khintchine representation which generalizes Equation (21):

$$\mathbb{E}[e^{-\int_{\mathbb{W}} f_1(w)\mu_1(\mathrm{d}w) - \cdots - \int_{\mathbb{W}} f_q(w)\mu_q(\mathrm{d}w)}]$$

$$= \exp\left\{ -\int_{\mathbb{W}} \int_{\mathbb{R}_+^q} (1 - e^{-s_1 f_1(w) - \cdots - s_q f_q(w)}) \nu_{(q)}(\mathrm{d}s_1, \ldots, \mathrm{d}s_q, \mathrm{d}w) \right\} \tag{22}$$

for arbitrary measurable functions $f_1, \ldots, f_d : \mathbb{W} \to \mathbb{R}_+$. The intensity measure $\nu_{(q)}$ in (22) is required to simultaneously satisfy

$$\nu_{(q)}(\mathbb{R}_+^q \times \{w\}) = 0 \quad \forall w \in \mathbb{W}$$

and

$$\int_{\mathbb{R}_+ \times A} \min\{||\boldsymbol{s}||, 1\} \nu_{(q)}(\mathrm{d}s_1, \ldots, \mathrm{d}s_q, \mathrm{d}w) < \infty,$$

for any bounded $A \in \mathcal{W}$, and having denoted by $||\boldsymbol{s}||$ the Euclidean norm of the vector $\boldsymbol{s} := (s_1, \ldots, s_q)$. In the present paper, we will work with a *homogeneous* Lévy intensity measure of the following form $\nu_{(q)}(\mathrm{d}s_1, \ldots, \mathrm{d}s_q, \mathrm{d}w) = \lambda_{(q)}(s_1, \ldots, s_q)\mathrm{d}s_1 \cdots \mathrm{d}s_q P(\mathrm{d}w)$, where $P$ is a diffuse probability measure on $(\mathbb{W}, \mathcal{W})$ and $\lambda_{(q)} : \mathbb{R}_+^q \to \mathbb{R}_+$ is measurable. See, e.g., [Kallenberg, 2017] for further details.

# B  Posterior analysis for SP priors: proofs and details

In the present section we derive the marginal, posterior and predictive distributions for the Bernoulli process model under a scaled process prior. Specifically we focus on the following statistical model throughout the section:

$$
\begin{aligned}
Z_n \mid \mu &\overset{\text{iid}}{\sim} \text{BeP}(\mu_{\Delta_{1,h}}), \quad \text{for} \quad n = 1, \dots, N \\
\mu_{\Delta_{1,h}} &\sim \text{SP}(\nu, h),
\end{aligned}
\tag{23}
$$

where $\mu_{\Delta_{1,h}}$ has been defined at the beginning of Section 2.2. In Subsection B.1 we provide some lemmas regarding SP priors, then Subsection B.2 is concerned with the Bayesian posterior analysis of the model in (23).

## B.1  Preparatory lemmas

Some preparatory lemmas are required before the posterior analysis. The first lemma provides the reader with the conditional distribution of $\mu_{\Delta_{1,h}}$ given $\Delta_{1,h}$.

**Lemma 1.** *Let* $\mu_{\Delta_{1,h}} \sim \text{SP}(\nu, h)$, *governed by the Lévy intensity measure* $\nu(\mathrm{d}s, \mathrm{d}w) = \lambda(s)\mathrm{d}sP(\mathrm{d}w)$ *on* $\mathbb{R}_+ \times \mathbb{W}$. *The conditional distribution of* $\mu_{\Delta_{1,h}}$, *given* $\Delta_{1,h}$, *equals the one of a CRM on* $(\mathbb{W}, \mathcal{W})$ *with Lévy intensity*

$$
\Delta_{1,h}\lambda(\Delta_{1,h}s)\mathbb{1}_{(0,1)}(s)\mathrm{d}sP(\mathrm{d}w).
$$

*Proof.* Recall the construction of a SP prior, as detailed in Section 2.2. It starts from an underlying CRM $\mu = \sum_{i \geq 1} \tau_i \delta_{W_i}$ with intensity $\nu$ on $\mathbb{R}_+ \times \mathbb{W}$. Moreover, having denoted by $\Delta_1 > \Delta_2 > \dots$ the decreasingly ordered jumps $\tau_i$'s of $\mu$, one considers:

$$
\mu_{\Delta_1} = \sum_{i \geq 1} \frac{\Delta_{i+1}}{\Delta_1} \delta_{W_{i+1}},
$$

and the SP process is defined by a change of measure of the largest jump $\Delta_1$, replaced with the distribution of $\Delta_{1,h}$. As a consequence it is sufficient to prove that $\mu_{\Delta_1} \mid \Delta_1$ is a CRM with Lévy intensity

$$
\Delta_1\lambda(\Delta_1 s)\mathbb{1}_{(0,1)}(s)\mathrm{d}sP(\mathrm{d}w).
\tag{24}
$$

In order to prove this remind that $(\Delta_i)_{i \geq 2}\mid\Delta_1$ are the points of a Poisson process with Lévy intensity $\lambda(s)\mathbb{1}_{(0,\Delta_1)}(s)\mathrm{d}s$, thanks to the representation by Ferguson and Klass [1972]. Therefore, the conditional distribution of $\mu_{\Delta_1}$, given $\Delta_1$, may be found by a simple evaluation of the Laplace functional. To this end, consider a measurable function $f : \mathbb{W} \to \mathbb{R}_+$ and compute

$$
\begin{aligned}
\mathbb{E}[e^{-\int_{\mathbb{W}} f(w)\mu_{\Delta_1}(\mathrm{d}w)}\mid\Delta_1] &= \mathbb{E}\left[e^{-\sum_{i \geq 1} f(W_{i+1})\Delta_{i+1}/\Delta_1}\mid\Delta_1\right] \\
&= \exp\left\{-\int_{\mathbb{W}}\int_0^{+\infty}(1 - e^{-f(w)s/\Delta_1})\mathbb{1}_{(0,\Delta_1)}(s)\lambda(s)\mathrm{d}s\ P(\mathrm{d}w)\right\} \\
&= \exp\left\{-\int_{\mathbb{W}}\int_0^{+\infty}(1 - e^{-f(w)s})\mathbb{1}_{(0,1)}(s)\lambda(s\Delta_1)\Delta_1\mathrm{d}s\ P(\mathrm{d}w)\right\}
\end{aligned}
$$

which is exactly the Laplace functional of a CRM having Lévy intensity (24).

□                                                                              □

We now provide the reader with a sufficient condition to ensure that each $Z_n$ in (23) is almost surely finite, for any $n \geq 1$.

**Lemma 2.** *Consider the model in Equation* (23). *If*

$$\mathbb{E}\left[\int_0^1 \Delta_{1,h}\lambda(s\Delta_{1,h})\mathrm{d}s\right] < \infty, \tag{25}$$

*then each $Z_n$ displays almost surely finitely many features — i.e. $\sum_{i\geq 1} A_{n,i} < \infty$, almost surely, for every $n \geq 1$.*

*Proof.* For a fixed $n \geq 1$, it is sufficient to show that condition (25) entails

$$\mathbb{E}\left[\sum_{i=1}^{\infty} A_{n,i}\right] < \infty.$$

The expected value in the previous formula may be computed as follows

$$\mathbb{E}\left[\sum_{i=1}^{\infty} A_{n,i}\right] = \mathbb{E}\left[\mathbb{E}\left[\sum_{i=1}^{\infty} A_{n,i}\Big|\Delta_{1,h}\right]\right] = \mathbb{E}\left[\mathbb{E}[\mu_{\Delta_{1,h}}(\mathrm{W})|\Delta_{1,h}]\right]$$

$$= \mathbb{E}\left[\int_{\mathrm{W}}\int_0^1 s\Delta_{1,h}\lambda(\Delta_{1,h}s)\mathrm{d}sP(\mathrm{d}w)\right] = \mathbb{E}\left[\int_0^1 s\Delta_{1,h}\lambda(\Delta_{1,h}s)\mathrm{d}s\right]$$

where we have applied the Campbell theorem [Kingman, 1992] and Lemma 1 to evaluate the total mass $\mu_{\Delta_{1,h}}(\mathrm{W})$ of $\mu_{\Delta_{1,h}}$. As a consequence, condition (25) is sufficient for the finiteness of the Bernoulli process $Z_n$.

□                                                                              □

## B.2  Posterior analysis

We start with the marginal distribution of the observations $Z_{1:N}$ induced by the model. Our derivation closely follows the proof in James [2017]. The marginal distribution is the counterpart of the "exchangeable feature probability function" (EFPF) for the Indian Buffet Process (IBP; see, e.g., Broderick et al. [2013]).

**Proposition 5** (Joint marginal distribution). *For any $N \geq 1$, let $Z_{1:N}$ be a random sample modeled as the BNP-Bernoulli model* (23), *where $\mu_{\Delta_{1,h}} \sim \mathrm{SP}(\nu, h)$. The probability that the observations $Z_{1:N}$ display $K_N = k$ distinct features, labelled by $\{W_1^*, \ldots, W_{K_N}^*\}$, with corresponding frequencies $(M_{N,1}, \ldots, M_{N,K_N}) = (m_1, \ldots, m_k)$, equals*

$$p_k^{(N)}(m_1, \ldots, m_k) = \int_0^{+\infty} e^{-\sum_{n=1}^N \phi_n(a)} \prod_{i=1}^k \int_0^1 s^{m_i}(1-s)^{N-m_i} a\lambda(as)\mathrm{d}s \; f_{\Delta_{1,h}}(a)\mathrm{d}a,$$

*where $\phi_n(a) = \int_0^1 s(1-s)^{n-1}a\lambda(as)\mathrm{d}s$.*

*Proof.* From the result showed in Lemma 1, we know that conditionally on a known value of $\Delta_{1,h} = a$, the random measure $\mu_{\Delta_{1,h}}$ is completely random. Therefore, we can exploit the result in James [2017, Proposition 3.1] to characterize the marginal distribution of the feature counts $m_{N,1}, \ldots, m_{N,K_N}$. This is given by

$$
\begin{aligned}
p_k^{(N)}(m_1, &\ldots, m_k \mid \Delta_{1,h} = a) \\
&= \exp\left\{ -\sum_{n=1}^N \phi_n(a) \right\} \prod_{i=1}^k \left\{ \int_0^1 s^{m_i}(1-s)^{N-m_i} a\lambda(as)\mathrm{d}s \right\},
\end{aligned}
\tag{26}
$$

with $\phi_n(a) = \int_0^1 s(1-s)^{n-1} a\lambda(as)\mathrm{d}s$. Integrating with respect to $f_{\Delta_{1,h}}$ — the mixing distribution of $\Delta_{1,h}$ — yields the desired result. □ □

Next, we characterize the posterior distribution of the random measure $\mu_{\Delta_{1,h}} \sim \mathrm{SP}(\nu, h)$. The posterior distribution of the law of $\Delta_{1,h}$ is an important ingredient in the study of the predictive properties of the model. We mention that the posterior characterization of Proposition 6 is a consequence of [James et al., 2015, Propositions 2.2] and the results developed by James [2017].

**Proposition 6** (Posterior distribution)**.** *For any $N \geq 1$, let $Z_{1:N}$ be a random sample modeled as the BNP-Bernoulli model* (23)*, where $\mu_{\Delta_{1,h}} \sim \mathrm{SP}(\nu, h)$. Suppose that the observations $Z_{1:N}$ display $K_N = k$ distinct features, labelled by $W_1^*, \ldots, W_{K_N}^*$, with corresponding frequencies $(M_{N,1}, \ldots, M_{N,K_N}) = (m_1, \ldots, m_k)$, then the conditional distribution of $\Delta_{1,h}$, given $Z_{1:N}$, has density function*

$$
g_{\Delta_{1,h}|Z_{1:N}}(a) \propto \exp\left\{ -\sum_{n=1}^N \phi_n(a) \right\} \prod_{i=1}^k \left\{ \int_0^1 s^{m_i}(1-s)^{N-m_i} a\lambda(as)\mathrm{d}s \right\} f_{\Delta_{1,h}}(a),
\tag{27}
$$

*with $\phi_n(a) = \int_0^1 s(1-s)^{n-1} a\lambda(as)\mathrm{d}s$. Moreover, the posterior distribution of the random measure $\mu_{\Delta_{1,h}}$, conditionally given $Z_{1:N}$ and $\Delta_{1,h}$, equals*

$$
\mu_{\Delta_{1,h}} \mid (\Delta_{1,h}, Z_{1:N}) \overset{d}{=} \mu'_{\Delta_{1,h}} + \sum_{i=1}^{K_N} J_i \delta_{W_i^*},
\tag{28}
$$

*where*

*i. $\mu'_{\Delta_{1,h}}|\Delta_{1,h} \sim \mathrm{CRM}(\nu'_{\Delta_{1,h}})$ with*

$$
\nu'_{\Delta_{1,h}}(\mathrm{d}s, \mathrm{d}w) = (1-s)^N \Delta_{1,h}\lambda(s\Delta_{1,h})\mathbb{1}_{(0,1)}(s)\mathrm{d}s\, P(\mathrm{d}w);
\tag{29}
$$

*ii. $J_{1:K_N}$ are $K_N$ independent random jumps and independent of $\mu'_{\Delta_{1,h}}$, with density on $[0,1]$ proportional to*

$$
f_{J_i|\Delta_{1,h}}(s) \propto (1-s)^{N-m_i} s^{m_i} \Delta_{1,h}\lambda(\Delta_{1,h}s).
\tag{30}
$$

*Proof.* Again, leveraging the result showed in Lemma 1, we know that conditionally on a known

value of $\Delta_{1,h}$, the measure $\mu_{\Delta_{1,h}}$ is completely random. Therefore, we can simply apply James [2017, Theorem 3.1] to obtain the posterior distribution of $\mu_{\Delta_{1,h}}|(\Delta_{1,h}, Z_{1:N})$ as described in Equation (28). Finally, the posterior distribution of the largest jump $\Delta_{1,h}$ conditionally on the observations $Z_{1:N}$ derived in Equation (27) follows by direct application of Bayes' theorem, recognizing that $f_{\Delta_{1,h}}$ is the prior distribution for $\Delta_{1,h}$, and the distribution in (26) as the likelihood of the observations $Z_{1:N}|\Delta_{1,h}$.

$\square$          $\square$

Last, we prove the predictive characterization provided in Proposition 2, which has a pivotal role in our analysis, as it is the conceptual starting point in order to study the predictive behavior of the model, and it again follows form [James, 2017].

*Proof of Proposition 2.* We consider $\zeta \overset{d}{=} \mu_{\Delta_{1,h}}$, thus we are dealing with the model (23). The posterior distribution of $\Delta_{1,h}$ in (8) follows from (27), by the argument used in Proposition 6. In order to prove the characterization in Equation (9), we use once again the fact that conditionally on a known value of $\Delta_{1,h}$, $\mu_{\Delta_{1,h}}$ is a completely random measure (see Lemma 1). Thus, we can exploit the results in [James, 2017] to characterize the predictive distribution of $Z_{N+1}$ given the sample $Z_{1:N}$ and the jump $\Delta_{1,h}$. More specifically the form of the predictive distribution in (9) follows by a plain application of James [2017, Proposition 3.2].

$\square$          $\square$

# C   Posterior analysis for SB-SP priors: proofs and details

Here we provide details and proofs of the results in Section 3.1, i.e. a full Bayesian analysis for the SB-SP prior. More specifically we prove Theorem 1, then we move to characterize the posterior distribution of $\Delta_{1,h_{c,\beta}}$, marginal, predictive and posterior distributions of the SB-SP model.

## C.1   Proof of Theorem 1

The posterior density of $\Delta_{1,h}$, given $Z_{1:N}$, has density proportional to

$$\prod_{n=1}^{N} e^{-\phi_n(a)} \prod_{i=1}^{K_N} \int_0^1 s^{m_{N,i}}(1-s)^{N-m_{N,i}} a\lambda(as) \mathrm{d}s\, f_{\Delta_{1,h}}(a),$$

where we used the notation $\phi_n(a) = \int_0^1 s(1-s)^{n-1} a\lambda(as)\mathrm{d}s$. Hence, there exists a normalizing factor $c(m_{N,1}, \ldots, m_{N,k}, N, K_N)$, depending on the sample size $N$, the distinct number of features $K_N$ and the frequency counts, such that

$$g_{\Delta_{1,h}|Z_{1:N}}(a) = \frac{\prod_{n=1}^{N} e^{-\phi_n(a)} \prod_{i=1}^{K_N} \int_0^1 s^{m_{N,i}}(1-s)^{N-m_{N,i}} a\lambda(as)\mathrm{d}s\, f_{\Delta_{1,h}}(a)}{c(m_{N,1}, \ldots, m_{N,K_N}, N, K_N)},$$

or equivalently we can write

$$g_{\Delta_{1,h}|Z_{1:N}}^{-1}(a) \prod_{n=1}^{N} e^{-\phi_n(a)} \prod_{i=1}^{K_N} \int_0^1 s^{m_{N,i}}(1-s)^{N-m_{N,i}} a\lambda(as)\mathrm{d}s \, f_{\Delta_{1,h}}(a)$$
$$= c(m_{N,1}, \ldots, m_{N,K_N}, N, K_N). \tag{31}$$

If the posterior density $g_{\Delta_{1,h}|Z_{1:N}}(a)$ does not depend on $m_{N,1}, \ldots, m_{N,K_N}$, then the function

$$g_{\Delta_{1,h}|Z_{1:N}}^{-1}(a) \prod_{n=1}^{N} e^{-\phi_n(a)} \, g(a) = f_1(a, K_N, N)$$

depends only on $K_N, N$ and $a$, but not on the frequency counts. Therefore, (31) boils down to

$$f_1(a, K_N, N) \cdot \prod_{i=1}^{K_N} \int_0^1 s^{m_{N,i}}(1-s)^{N-m_{N,i}} a\lambda(as)\mathrm{d}s = c(m_{N,1}, \ldots, m_{N,K_N}, N, K_N). \tag{32}$$

As a consequence, the function on the right hand side of (32) is independent of $a$, for any choice of the vector $(m_{N,1}, \ldots, m_{N,K_N}, N, K_N)$. Now we consider $m_{N,1} = \cdots = m_{N,K_N} = m > 0$, and we can say that the function

$$\left[ w(a, K_N, N) \int_0^1 s^m (1-s)^{N-m} a\lambda(as)\mathrm{d}s \right]^{K_N} \tag{33}$$

does not depend on $a \in \mathbb{R}_+$, where $w(a, K_N, N) = \sqrt[K_N]{f_1(a, K_N, N)}$. We now select $m = N$, thus the function

$$w(a, K_N, N) \int_0^1 s^N a\lambda(as)\mathrm{d}s \tag{34}$$

does not depend on $a \in \mathbb{R}_+$. Note that, since $f_{\Delta_{1,h}}$ and $\lambda$ are functions of class $C^1(\mathbb{R}_+)$, i.e., derivable with continuous derivative, also $w$ is in class $C^1(\mathbb{R}_+)$ with respect to the variable $a$. Thus, we can take the derivative of (34), and this is equal to 0:

$$\frac{\mathrm{d}}{\mathrm{d}a} w(a, K_N, N) \int_0^a s^N \lambda(s)\mathrm{d}s a^{-N} - N a^{-N-1} w(a, K_N, N) \int_0^a s^N \lambda(s)\mathrm{d}s + w(a, K_N, N)\lambda(a) = 0$$

which is an ordinary differential equation in $w$, and it can be easily solved by separation of variables, thus obtaining

$$w(a, K_N, N) = a^N \cdot \frac{R}{\int_0^a s^N \lambda(s)\mathrm{d}s}$$

where $R > 0$ is a suitable constant independent of $a$. As a consequence, the function in (33) equals

$$\left[ \frac{R}{\int_0^1 s^N \lambda(as)\mathrm{d}s} \cdot \int_0^1 s^m (1-s)^{N-m} \lambda(as)\mathrm{d}s \right]^{K_N}$$

25

and this is independent of $a \in \mathbb{R}_+$. It is possible to choose $m = N - 1$ in the previous function, and we can state that

$$\int_0^1 s^{N-1}\lambda(as)\mathrm{d}s - \int_0^1 s^N\lambda(as)\mathrm{d}s = C\int_0^1 s^N\lambda(as)\mathrm{d}s$$

where $C$ is constant with respect to $a$. If one takes the derivative of the previous equation two times with respect to $a$, then she obtains

$$\lambda(a)(1 - NC) = a\lambda'(a)C,$$

which is an ordinary differential equation in $\lambda$ that can be solved by separation of variables. In particular we get the following result

$$\lambda(a) = \alpha a^{(1-NC)/C}, \quad \text{for } \alpha > 0. \tag{35}$$

The exponent of $a$ in (35) should satisfy

$$\int_0^{+\infty} \min\{1, a\}\lambda(a)\mathrm{d}a < +\infty,$$

from which it is easy to realize that $-2 < (1 - NC)/C < -1$, hence

$$\lambda(a) = \alpha\frac{1}{a^{1+\sigma}}$$

where $\alpha > 0$ and $\sigma \in (0, 1)$. The reverse implication of the theorem is trivially true, hence the proof is completed.

$\square$

## C.2 Detailed derivation of the distribution of $\Delta_{1,h_{c,\beta}}$

We first derive explicitly the distribution of the largest jump given in Equation (6). This follows from direct application of the law of the largest jump,

$$F_{\Delta_1}(\mathrm{d}a) = \exp\left\{-\int_a^\infty \lambda_\sigma(s)\mathrm{d}s\right\}\lambda_\sigma(a)\mathrm{d}a$$

when the Lévy measure is

$$\lambda_\sigma(s)\mathrm{d}s = \sigma s^{-\sigma-1}\mathbb{1}_{\mathbb{R}_+}(s)\mathrm{d}s.$$

Having denoted by $f_{\Delta_1}$ the density function of $F_{\Delta_1}$, we get

$$f_{\Delta_1}(a) = \lambda_\sigma(a)e^{-\Lambda(a)}\mathbb{1}_{\mathbb{R}_+}(a) = \sigma a^{-\sigma-1}\exp\left\{-\int_a^\infty \sigma u^{-1-\sigma}\mathrm{d}u\right\}\mathbb{1}_{\mathbb{R}_+}(a)$$
$$= \sigma a^{-\sigma-1}e^{-a^{-\sigma}}\mathbb{1}_{\mathbb{R}_+}(a).$$

From direct inspection, we recognize that this is the density function of $\Delta_1 = T^{-1/\sigma}$, where $T$ is a Gamma with parameters $(1, 1)$. The mixing measure is then obtained by tilting the density $f_{\Delta_1}$ as

26

follows:

$$f_{\Delta_{1,h_{c,\beta}}}(a) \propto f_{\Delta_1}(a)h_{c,\beta}(a) = \sigma a^{-\sigma(c+1)-1}\exp\left\{-\beta a^{-\sigma}\right\}\mathbb{1}_{\mathbb{R}_+}(a),$$

i.e. letting

$$h_{c,\beta}(a) \propto a^{-\sigma c}\exp\left\{-(\beta-1)a^{-\sigma}\right\}.$$

By integration, we get the normalizing constant:

$$\int_0^\infty a^{-\sigma(c+1)-1}\exp\left\{-\beta a^{-\sigma}\right\}\mathrm{d}a = \frac{\Gamma(c+1)}{\sigma\beta^{c+1}}.$$

from which

$$f_{\Delta_{1,h_{c,\beta}}}(a) = \frac{\sigma\beta^{c+1}}{\Gamma(c+1)}a^{-\sigma(c+1)-1}\exp\left\{-\beta a^{-\sigma}\right\}\mathbb{1}_{\mathbb{R}_+}(a). \tag{36}$$

## C.3   Posterior distribution of SB-SP priors

Here we characterize the posterior distribution of SB-SP priors: the result is not included in the paper, but we think it is useful to have a full picture on SB-SP priors from a Bayesian viewpoint.

**Proposition 7.** *For $N \geq 1$ let $Z_{1:N}$ be a random sample modeled as the BNP-Bernoulli model* (1), *with $\zeta \sim \mathrm{SB\text{-}SP}(\sigma,c,\beta)$. If $Z_{1:N}$ displays $K_N = k$ distinct features $\{W_1^*,\ldots,W_{K_N}^*\}$, each feature $W_i^*$ appearing exactly $M_{N,i} = m_i$ times in the samples, then the conditional distribution of $\Delta_{1,h_{c,\beta}}$, given $Z_{1:N}$, has a density function of the form*

$$g_{\Delta_{1,h_{c,\beta}}\,|\,Z_{1:N}}(a) = \sigma\frac{(\beta+\gamma_0^{(N)})^{k+c+1}}{\Gamma(k+c+1)}a^{-k\sigma-(c+1)\sigma-1}\exp\{-a^{-\sigma}(\beta+\gamma_0^{(N)})\}, \tag{37}$$

*where $\gamma_0^{(n)} = \sigma\sum_{1\leq i\leq n}B(1-\sigma,i)$, with $B(\cdot,\cdot)$ denoting the (standard) Beta function. Moreover, the conditional distribution of $\zeta$, given $(\Delta_{1,h_{c,\beta}}, Z_{1:N})$, coincides with the distribution of*

$$\zeta\,|\,(\Delta_{1,h_{c,\beta}}, Z_{1:N}) \overset{d}{=} \mu'_{\Delta_{1,h_{c,\beta}}} + \sum_{i=1}^{K_N}J_i\delta_{W_i^*}, \tag{38}$$

*where:*

i) *$\mu'_{\Delta_{1,h_{c,\beta}}}$ is a discrete random measure such that $\mu'_{\Delta_{1,h_{c,\beta}}}\,|\,\Delta_{1,h_{c,\beta}} \sim \mathrm{CRM}(\nu'_{\Delta_{1,h_{c,\beta}}})$, with $\nu'_{\Delta_{1,h_{c,\beta}}}$ being*

$$\nu'_{\Delta_{1,h_{c,\beta}}}(\mathrm{d}s,\mathrm{d}w) = \Delta_{1,\Delta_{1,h_{c,\beta}}}^{-\sigma}(1-s)^N\sigma s^{-1-\sigma}\mathbb{1}_{(0,1)}(s)\mathrm{d}sP(\mathrm{d}w); \tag{39}$$

ii)

$$J_i|\Delta_{1,h_{c,\beta}} \sim \mathrm{Beta}(m_i-\sigma, N-m_i+1), \tag{40}$$

*where* Beta *denotes the beta distribution.*

*Proof.* We apply Proposition 6, which describes the general posterior distribution of a SP process. We first compute the posterior distribution (8) of the largest jump conditionally on observations $Z_{1:N}$.

To do so we specify (27) in our case, and we first compute the exponent $\phi_n(a)$. In our case the Lévy density equals $\lambda_\sigma(s) = \sigma s^{-\sigma-1}$ and the mixing density of $\Delta_{1,h_{c,\beta}}$ is provided in Equation (36), thus the exponent $\phi_n$ takes the form

$$\phi_n(a) = \sigma \int_0^1 s(1-s)^{n-1} a^{-\sigma} s^{-\sigma-1} \mathrm{d}s = \sigma a^{-\sigma} B(1-\sigma, n). \tag{41}$$

Recalling the shorthand notation $\gamma_0^{(N)} = \sigma \sum_{1 \le n \le N} B(1-\sigma, n)$, the posterior distribution of $\Delta_{1,h_{c,\beta}}$ is then proportional to

$$a^k \exp\left\{ -a^{-\sigma} \gamma_0^{(N)} \right\} \prod_{i=1}^k \int_0^1 t^{m_i} (1-t)^{N-m_i} \lambda_\sigma(at) \mathrm{d}t \, f_{\Delta_{1,h_{c,\beta}}}(a)$$
$$\propto a^{-\sigma(k+c+1)-1} \exp\left\{ -a^{-\sigma} \left[ \beta + \gamma_0^{(N)} \right] \right\},$$

where $f_{\Delta_{1,h_{c,\beta}}}$ has been specified in (36). As a consequence we get

$$\Delta_{1,h_{c,\beta}}^{-\sigma} \mid Z_{1:N} \sim \mathrm{Gamma}\left( k + c + 1, \beta + \gamma_0^{(N)} \right),$$

which corresponds to the posterior density in (37). The characterization of the posterior distribution in (38) is an easy consequence of Proposition 6, by a specialization of this result with the choice $\lambda(s) = \lambda_\sigma(s) = \sigma s^{-\sigma-1}$ for the underlying Lévy intensity. □ □

## C.4  Proof of Proposition 3

The predictive characterization is a simple consequence of the general characterization in Proposition 2 with the SB-SP specifications $\lambda(s) = \lambda_\sigma(s) = \sigma s^{-\sigma-1}$. □

## C.5  Proof of Proposition 4

We apply Proposition 5 to obtain the marginal distribution for the SB-SP prior. Conditionally on $\Delta_{1,h_{c,\beta}} = a$, using the form $\phi_n(a)$ derived in (41), the marginal distribution is given by

$$p_k^{(N)}(m_1, \ldots, m_k \mid \Delta_{1,h_{c,\beta}} = a) = (\sigma a^{-\sigma})^k \exp\left\{ -a^{-\sigma} \gamma_0^{(N)} \right\} \prod_{i=1}^k \int_0^1 s^{m_i-\sigma-1}(1-s)^{N-m_i} \mathrm{d}s,$$

that may be written in terms of the Beta function as follows

$$p_k^{(N)}(m_1, \ldots, m_k \mid \Delta_{1,h_{c,\beta}} = a) = (\sigma a^{-\sigma})^k \exp\left\{ -a^{-\sigma} \gamma_0^{(N)} \right\} \prod_{i=1}^k B(m_i - \sigma, N - m_i + 1).$$

Last, we obtain the marginal distribution in Equation (10) by randomizing with respect to the mixing distribution of the largest jump given in Equation (36). We need to compute

$$
\begin{aligned}
p_k^{(N)}(m_1, \ldots, m_k) &= \int_0^\infty p_k^{(N)}(m_1, \ldots, m_k \mid \Delta_{1,h_{c,\beta}} = a) f_{\Delta_{1,h_{c,\beta}}}(a) \mathrm{d}a \\
&= \frac{\sigma^{k+1}\beta^{c+1}}{\Gamma(c+1)} \prod_{i=1}^k B(m_i - \sigma, N - m_i + 1) \\
&\quad \times \int_0^\infty a^{-\sigma(k+c+1)-1} \exp\left\{ -a^{-\sigma}\left[\beta + \gamma_0^{(N)}\right]\right\} \mathrm{d}a \\
&= \frac{\sigma^k \beta^{c+1}}{(\beta + \gamma_0^{(N)})^{k+c+1}} \frac{\Gamma(k+c+1)}{\Gamma(c+1)} \prod_{i=1}^k B(m_i - \sigma, N - m_i + 1),
\end{aligned}
$$

and the thesis now follows. $\qquad\square$

# D   Estimation of the unseen features via SB-SP priors: proofs

Here we detail the proofs of Section 3.2, which is devoted to the unseen-features problem under the SB-SP prior.

## D.1   Proof of Theorem 2

We first focus on the proof of (13), i.e. the posterior distribution of $U_N^{(M)}$. In order to do this we exploit the predictive characterization provided in Proposition 3 to evaluate the probability generating function (PGF) of the random variable $U_N^{(M)}$ *a posteriori*, conditionally on the sample $Z_{1:N}$. We denote the PGF as $\mathcal{G}_{U_N^{(M)}}(\cdot)$. If $t$ belongs to a neighborhood of the origin, then one has

$$
\mathcal{G}_{U_N^{(M)}}(t) = \mathbb{E}\left[t^{U_N^{(M)}} \mid Z_{1:N}\right] = \mathbb{E}\left[\mathbb{E}\left[t^{U_N^{(M)}} \mid Z_{1:N}, \Delta_{1,h_{c,\beta}}\right] \mid Z_{1:N}\right] \tag{42}
$$

where we have applied the tower property of the conditional expectation. We now observe that, conditionally on $Z_{1:N}$ and $\Delta_{1,h_{c,\beta}}$, the random variable $U_N^{(M)}$ may be represented as

$$
U_N^{(M)} \mid (Z_{1:N}, \Delta_{1,h_{c,\beta}}) \stackrel{\mathrm{d}}{=} \sum_{i\geq 1} \mathbb{1}\left(\sum_{m=1}^M A'_{N+m,i} > 0\right),
$$

where we used the representation given in Proposition 3. Here, independently across $i$, $A'_{N+m,i}$ is a Bernoulli random variable with parameter $\rho'_i$, conditionally on the random measure $\mu'_{\Delta_{1,h_{c,\beta}}} = \sum_{i\geq 1} \rho'_i \delta_{W'_i}$ with Lévy intensity $\sigma \Delta_{1,h_{c,\beta}}^{-\sigma}(1-s)^N s^{-1-\sigma} \mathbb{1}_{(0,1)}(s)\mathrm{d}s P(\mathrm{d}w)$. We now focus on the eval-

uation of the expected value in Equation (42):

$$\mathbb{E}\left[t^{U_N^{(M)}} \mid Z_1, \ldots, Z_N, \Delta_{1,h_{c,\beta}}\right]$$

$$= \mathbb{E}\left[\mathbb{E}\left[\prod_{i\geq 1}\left((t-1)\mathbb{1}\left\{\sum_{m=1}^{M} A'_{N+m,i} > 0\right\} + 1\right)\bigg|\mu'_{\Delta_{1,h_{c,\beta}}}\right]\right]$$

$$= \mathbb{E}\left[\prod_{i\geq 1}\left[(t-1)\mathbb{P}\left(\sum_{m=1}^{M} A'_{N+m,i} > 0 \mid \mu'_{\Delta_{1,h_{c,\beta}}}\right) + 1\right]\right]$$

$$= \mathbb{E}\left[\prod_{i\geq 1}\left[(t-1)\left\{1 - \prod_{m=1}^{M}\mathbb{P}(A'_{N+m,i} = 0 \mid \mu'_{\Delta_{1,h_{c,\beta}}})\right\} + 1\right]\right],$$

where we applied the independence of the Bernoulli random variables $A'_{N+m,i}$s, conditionally on $\mu'_{\Delta_{1,h_{c,\beta}}}$. We now recall that $\mu'_{\Delta_{1,h_{c,\beta}}}$ is a CRM with a known Lévy measure and that the $A'_{N+m,i}$s are Bernoulli with parameter $\rho'_i$ to obtain

$$\mathbb{E}\left[t^{U_N^{(M)}} \mid Z_1, \ldots, Z_N, \Delta_{1,h_{c,\beta}}\right]$$

$$= \mathbb{E}\left[\prod_{i\geq 1}((t-1)(1 - (1-\rho'_i)^M) + 1)\right]$$

$$= \mathbb{E}\left[\exp\left\{\sum_{i\geq 1}\log\left[(t-1)(1 - (1-\rho'_i)^M) + 1\right]\right\}\right]$$

$$= \exp\left\{-(1-t)\int_0^1 (1 - (1-s)^M)(1-s)^N \Delta_{1,h_{c,\beta}}^{-\sigma}\sigma s^{-1-\sigma}\mathrm{d}s\right\}.$$

$$= \exp\left\{-(1-t)\Delta_{1,h_{c,\beta}}^{-\sigma}\gamma_N^{(M)}\right\},$$

where we used the identity

$$\int_0^1 \left[1 - (1-s)^M\right](1-s)^N s^{-1-\sigma}\mathrm{d}s = \sum_{m=1}^{M} B(1-\sigma, N+m).$$

We replace this expression in Equation (42) to obtain

$$\mathcal{G}_{U_N^{(M)}}(t) = \mathbb{E}[\exp\{-(1-t)\Delta_{1,h_{c,\beta}}^{-\sigma}\gamma_N^{(M)}\} \mid Z_{1:N}]. \tag{43}$$

The results now follows by integrating with respect to the posterior distribution of $\Delta_{1,h_{c,\beta}}^{-\sigma}$, given in Equation (37):

$$\mathcal{G}_{U_N^{(M)}}(t) = \frac{(\beta + \gamma_0^{(N)})^{K_N+c+1}}{\Gamma(K_N+c+1)} \int_0^\infty \exp\left\{-(1-t)\gamma_N^{(M)}x\right\} x^{K_N+c} e^{-(\beta+\gamma_0^{(N)})x} \mathrm{d}x$$

$$= \frac{(\beta + \gamma_0^{(N)})^{K_N+c+1}}{\Gamma(K_N+c+1)} \frac{\Gamma(K_N+c+1)}{(\beta + \gamma_0^{(N)} + (1-t)\gamma_N^{(M)})^{K_N+c+1}}$$

$$= \left(\frac{\beta + \gamma_0^{(N)}}{\beta + \gamma_0^{(N+M)} - t\gamma_N^{(M)}}\right)^{K_N+c+1} = \left(\frac{1 - p_N^{(M)}}{1 - t p_N^{(M)}}\right)^{K_N+c+1},$$

for any $|t| < 1/p_N^{(M)}$, where $p_N^{(M)} := \gamma_N^{(M)}/(\beta + \gamma_0^{(N+M)}) \le 1$. This is the probability generating function of a negative binomial distribution where $K_N + c + 1$ is the number of failures, and $p_N^{(M)}$ is the success probability in each experiment.

We now apply similar arguments to derive the posterior distribution of $U_N^{(M,r)}$, provided in (14). Again, we calculate the probability generating function of $U_N^{(M,r)}$ *a posteriori*, denoted here as $\mathcal{G}_{U_N^{(M,r)}}(\cdot)$. If $t$ belongs to a neighborhood of the origin, then one has

$$\mathcal{G}_{U_N^{(M,r)}}(t) = \mathbb{E}\left[t^{U_N^{(M,r)}} \mid Z_{1:N}\right] = \mathbb{E}\left[\mathbb{E}\left[t^{U_N^{(M,r)}} \mid Z_{1:N}, \Delta_{1,h_{c,\beta}}\right] \mid Z_{1:N}\right]. \tag{44}$$

It is now easy to see that, conditionally on $Z_1, \ldots, Z_N, \Delta_{1,h_{c,\beta}}$, the random variable $U_N^{(M,r)}$ may be written as

$$U_N^{(M,r)} \mid Z_1, \ldots, Z_N, \Delta_{1,h_{c,\beta}} \stackrel{\mathrm{d}}{=} \sum_{i\ge 1} \mathbb{1}\left\{\sum_{m=1}^M A'_{N+m,i} = r\right\}$$

by applying Proposition 3. With the same notation used in the first part of the proof, we recall that the $A'_{N+m,i}$s are independent Bernoulli variables with parameters $\rho'_i$, conditionally on the CRM $\mu'_{\Delta_{1,h_{c,\beta}}} = \sum_{i\ge 1} \rho'_i \delta_{W'_i}$ with Lévy intensity $\sigma \Delta_{1,h_{c,\beta}}^{-\sigma} (1-s)^N s^{-1-\sigma} \mathbb{1}_{(0,1)}(s) \mathrm{d}s P(\mathrm{d}w)$. Across $i$, the random variables

$$S_{M,i} := \sum_{m=1}^M A'_{N+m,i}$$

are independent, each one distributed as a binomial with parameters $M$ and success probability $\rho'_i$. We then evaluate the expected value appearing in (44) as follows:

$$\mathbb{E}\left[t^{U_N^{(M,r)}} \mid Z_{1:N}, \Delta_{1,h_{c,\beta}}\right] = \mathbb{E}\left[\mathbb{E}\left[\prod_{i\ge 1}\left((t-1)\mathbb{1}\left\{\sum_{m=1}^M A'_{N+m,i} = r\right\} + 1\right) \bigg| \mu'_{\Delta_{1,h_{c,\beta}}}\right]\right]$$

$$= \mathbb{E}\left[\prod_{i\ge 1}\left[(t-1)\mathbb{P}(S_{M,i} = r \mid \mu'_{\Delta_{1,h_{c,\beta}}}) + 1\right]\right]$$

$$= \mathbb{E}\left[\prod_{i\ge 1}\left[(t-1)\binom{M}{r}(\rho'_i)^r(1-\rho'_i)^{M-r} + 1\right]\right].$$

Since $\mu'_{\Delta_{1,h_{c,\beta}}}$ is a CRM with a known Lévy measure, we can evaluate the previous expected value:

$$\mathbb{E}\left[t^{U_N^{(M,r)}} \mid Z_1, \ldots, Z_N, \Delta_{1,h_{c,\beta}}\right]$$

$$= \mathbb{E}\left[\exp\left\{\sum_{i\geq 1} \log\left((t-1)\binom{M}{r}(\rho_i')^r(1-\rho_i')^{M-r}+1\right)\right\}\right]$$

$$= \exp\left\{-(1-t)\binom{M}{r}\int_0^1 s^{r-\sigma-1}(1-s)^{M+N-r}\mathrm{d}s\sigma\Delta_{1,h_{c,\beta}}^{-\sigma}\right\}$$

$$= \exp\left\{-(1-t)\Delta_{1,h_{c,\beta}}^{-\sigma}\sigma\binom{M}{r}B(r-\sigma, M+N-r+1)\right\}$$

$$= \exp\left\{-(1-t)\Delta_{1,h_{c,\beta}}^{-\sigma}\rho_N^{(M,r)}\right\},$$

where we used the notation introduced in the statement of the theorem, i.e. $\rho_N^{(M,r)} = \sigma\binom{M}{r}B(r-\sigma, M+N-r+1)$. Then the probability generating function in Equation (44) is obtained by integrating with respect to the posterior distribution of the largest jump provided in (37):

$$\mathcal{G}_{U_N^{(M,r)}}(t) = \int_0^\infty \exp\left\{-(1-t)x\rho_N^{(M,r)}\right\} \cdot \frac{(\beta+\gamma_0^{(N)})^{K_N+c+1}}{\Gamma(K_N+c+1)}x^{K_N+c}e^{-(\beta+\gamma_0^{(N)})x}\mathrm{d}x$$

$$= \frac{\Gamma(K_N+c+1)}{(\beta+\gamma_0^{(N)}+(1-t)\rho_N^{(M,r)})^{K_N+c+1}} \cdot \frac{(\beta+\gamma_0^{(N)})^{K_N+c+1}}{\Gamma(K_N+c+1)}$$

$$= \left(\frac{\beta+\gamma_0^{(N)}}{\beta+\gamma_0^{(N)}+\rho_N^{(M,r)}-t\rho_N^{(M,r)}}\right)^{K+c+1} = \left(\frac{1-p_N^{(M,r)}}{1-tp_N^{(M,r)}}\right)^{K_N+c+1}$$

for any $|t| < 1/p_N^{(M,r)}$, where we have set

$$p_N^{(M,r)} := \frac{\rho_N^{(M,r)}}{\beta+\rho_N^{(M,r)}+\gamma_0^{(N)}}.$$

Then we conclude that the posterior distribution of $U_N^{(M,r)}$ is a negative binomial distribution where $K_N+c+1$ is the number of failures, and $p_N^{(M,r)}$ is the success probability in each experiment.

$\square$

## D.2  Proof of Theorem 3

In order to prove this result, we first exploit the Lévy continuity theorem, to obtain a convergence in distribution, and later strengthen this result to show that the convergence holds true also in the almost-sure sense. For the convergence in distribution, thanks to Theorem 2, the characteristic function of $U_N^{(M)}/M^\sigma \mid Z_{1:N}$ is given by

$$\Phi_{U_N^{(M)}/M^\sigma}(t) = \left(\frac{1-p_N^{(M)}}{1-p_N^{(M)}e^{it/M^\sigma}}\right)^{K_N+c+1}$$

where $t \in \mathbb{R}$, $K_N$ is the number of distinct features in $Z_{1:N}$ and $p_N^{(M)} = \gamma_N^{(M)}/(\gamma_0^{(N)} + \gamma_N^{(M)} + \beta)$. The quantity above can be rewritten as

$$\Phi_{U_N^{(M)}/M^\sigma}(t) = \left( \frac{\beta + \gamma_0^{(N)}}{\beta + \gamma_0^{(N)} + \gamma_N^{(M)} - \gamma_N^{(M)} e^{it/M^\sigma}} \right)^{K_N+c+1}.$$

We can exploit Masoero et al. [2021, Lemma 1] to determine the asymptotic expansion $\gamma_N^{(M)} = M^\sigma \Gamma(1-\sigma)(1+O(M^{-\sigma}))$ as $M \to +\infty$, having used the big-$O$ notation. Thus, using the asymptotic expansion of the exponential function, one has

$$\begin{aligned}
\Phi_{U_N^{(M)}/M^\sigma}(t) &= \left( \frac{\beta + \gamma_0^{(N)}}{\beta + \gamma_0^{(N)} + \gamma_N^{(M)} - \gamma_N^{(M)}(1 + itM^{-\sigma} + O(M^{-2\sigma}))} \right)^{K_N+c+1} \\
&= \left( \frac{\beta + \gamma_0^{(N)}}{\beta + \gamma_0^{(N)} - \gamma_N^{(M)} itM^{-\sigma} + O(M^{-\sigma})} \right)^{K_N+c+1} \\
&= \left( \frac{\beta + \gamma_0^{(N)}}{\beta + \gamma_0^{(N)} - it\Gamma(1-\sigma) + O(M^{-\sigma})} \right)^{K_N+c+1}
\end{aligned}$$

which converges to the characteristic function of a gamma random variable with parameters $(K_N + c + 1, (\gamma_0^{(N)} + \beta)/\Gamma(1-\sigma))$ as $M \to +\infty$. This proves that

$$U_N^{(M)}/M^\sigma \mid Z_{1:N} \xrightarrow{\mathrm{d}} W_N, \quad \text{where } W_N \sim \mathrm{Gamma}\left( K_N + c + 1, \frac{\beta + \gamma_0^{(N)}}{\Gamma(1-\sigma)} \right).$$

In order to prove convergence in the almost sure sense, we exploit the corresponding results proved for the stable beta-Bernoulli process in Masoero et al. [2021, Theorem 2] for the statistic $U_N^{(M)}$. We first notice that if we condition on the value of the largest jump $\Delta_{1,h_{c,\beta}}$, then the SB-SP-Bernoulli is a completely random measure whose asymptotic behavior is analogous to the stable beta-Bernoulli process. Thus, specializing the almost sure convergence results given in Masoero et al. [2021, Theorem 2], a posteriori, we have

$$\mathbb{P}\left( \lim_{M \to +\infty} \frac{U_N^{(M)}}{M^\sigma} = a^{-\sigma}\Gamma(1-\sigma) \Big| Z_{1:N}, \Delta_{1,h_{c,\beta}} = a \right) = 1. \tag{45}$$

The probability limit for the model in which the largest jump is random is obtained by observing that

$$\begin{aligned}
\mathbb{P}&\left( \lim_{M \to +\infty} \frac{U_N^{(M)}}{M^\sigma} = \Delta_{1,h_{c,\beta}}^{-\sigma}\Gamma(1-\sigma) \Big| Z_{1:N} \right) \\
&= \mathbb{E}\left[ \mathbb{P}\left( \lim_{M \to +\infty} \frac{U_N^{(M)}}{M^\sigma} = \Delta_{1,h_{c,\beta}}^{-\sigma}\Gamma(1-\sigma) \Big| Z_{1:N}, \Delta_{1,h_{c,\beta}} \right) \Big| Z_{1:N} \right] \overset{(45)}{=} 1,
\end{aligned}$$

33

in other words $U_N^{(M)}/M^\sigma$ converges almost surely to the random variable $\Delta_{1,h_{c,\beta}}^{-\sigma}\Gamma(1-\sigma)$, with respect to the conditional probability $\mathbb{P}$ given $Z_{1:N}$. Note also that the posterior distribution of $\Delta_{1,h_{c,\beta}}^{-\sigma}\Gamma(1-\sigma)$ is a Gamma with parameters

$$\left(K_N + c + 1, \frac{\beta + \gamma_0^{(N)}}{\Gamma(1-\sigma)}\right),$$

thus the a.s. convergence in (17) now follows.

We proceed along the same lines as to show the validity of (18). First, we show the convergence in distribution of $U_N^{(M,r)}$ using the characteristic function, and then we show that the result also holds in an almost sure sense. From Theorem 2, the characteristic function of $U_N^{(M,r)}/M^\sigma \mid Z_{1:N}$ is given by

$$\Phi_{U_N^{(M,r)}/M^\sigma}(t) = \left(\frac{1 - p_N^{(M,r)}}{1 - p_N^{(M,r)} e^{it/M^\sigma}}\right)^{K_N+c+1}$$

where $t \in \mathbb{R}$, and $p_N^{(M)} = \rho_N^{(M,r)}/(\gamma_0^{(N)} + \rho_N^{(M,r)} + \beta)$, and $\rho_N^{(M,r)}$ was defined in the statement of Theorem 2. The expression above is equivalent to

$$\Phi_{U_N^{(M,r)}/M^\sigma}(t) = \left(\frac{\beta + \gamma_0^{(N)}}{\beta + \gamma_0^{(N)} + \rho_N^{(M,r)}(1 - e^{it/M^\sigma})}\right)^{K_N+c+1}.$$

Thanks to the well-known asymptotic relation for the ratio of gamma functions, it is easy to see that

$$\rho_N^{(M,r)} = \frac{\sigma}{r!}\Gamma(r-\sigma)(r-\sigma)\frac{\Gamma(M+1)}{\Gamma(M+1-r)}\frac{\Gamma(N+M+2-r)}{\Gamma(N+M+2-\sigma)} = \frac{\sigma}{r!}\Gamma(r-\sigma)M^\sigma(1+O(M^{-1}))$$

as $M \to +\infty$. Hence, the characteristic function under study boils down to

$$\Phi_{U_N^{(M,r)}/M^\sigma}(t) = \left(\frac{\beta + \gamma_0^{(N)}}{\beta + \gamma_0^{(N)} + \sigma\Gamma(r-\sigma)(r!)^{-1}M^{-\sigma}(1+O(M^{-1}))(1 - e^{it/M^\sigma})}\right)^{K_N+c+1}$$

$$= \left(\frac{\beta + \gamma_0^{(N)}}{\beta + \gamma_0^{(N)} - \sigma\Gamma(r-\sigma)(r!)^{-1}it + O(M^{-\sigma})}\right)^{K_N+c+1}$$

which converges, as $M \to +\infty$, to the characteristic function of a gamma random variable with parameters as in the thesis. The almost sure statement of (18) goes along similar lines, indeed one can exploit the convergence theorems proved by Masoero et al. [2021] to state that

$$\mathbb{P}\left(\lim_{M\to+\infty}\frac{U_N^{(M,r)}}{M^\sigma} = \frac{\sigma(1-\sigma)_{(r-1)}}{r!}\Delta_{1,h_{c,\beta}}^{-\sigma}\Gamma(1-\sigma)\Big|Z_{1:N}, \Delta_{1,h_{c,\beta}}\right) = 1.$$

Exactly as before, one can conclude that

$$\mathbb{P}\left(\lim_{M\to+\infty}\frac{U_N^{(M,r)}}{M^\sigma} = \frac{\sigma(1-\sigma)_{(r-1)}}{r!}\Delta_{1,h_{c,\beta}}^{-\sigma}\Gamma(1-\sigma)\Big|Z_{1:N}\right) = 1,$$

where the posterior distribution of the limiting random variable is a gamma with the same parameters as in the statement of the theorem (Equation (18)).

$\square$

# E   Multivariate extension

In the present section we discuss the multivariate version of the Bernoulli process, which we call the Bernoulli process with a condiment or the simple multinomial process, using the terminology of James [2017]. We first revise the model of James [2017] and the associated prior, called stable-Beta-Dirichlet process, then we move to introduce a new scaled prior for the model. In both the cases, we determine closed-form results to face prediction of new features with condiments. These models are extremely important in genomics to account for the presence of variants at certain genomic loci with a specific characteristic (or condiment). See, e.g., Lee et al. [2016].

## E.1   Bernoulli process with a condiment

The IBP process with a condiment has been introduced by James [2017] and we remind the definition here. For $q = 1, 2, \ldots$, we define the vector of probabilities $\boldsymbol{p} = (p_1, \ldots, p_q)$ taking values in the following set

$$S_q = \{\boldsymbol{s} := (s_1, \ldots, s_q) : \ s_j > 0 \text{ as } j = 1, \ldots, q, \ |\boldsymbol{s}| := \sum_{j=1}^{q} s_j < 1\}$$

where for a generic vector $\boldsymbol{s}$, $|\boldsymbol{s}| = \sum_{j=1}^{q} s_j$ denotes the $L^1$ norm of the vector. For a fixed vector $\boldsymbol{p} \in S_q$, we also define the *simple multinomial* distribution $\mathsf{M}(1, \boldsymbol{p})$. A vector $\boldsymbol{A} = (A_1, \ldots, A_q) \in \{0, 1\}^q$ is said to have the *simple multinomial* distribution with parameter vector $\boldsymbol{p}$ iff it has the following probability mass function

$$\mathbb{P}(\boldsymbol{A} = \boldsymbol{a}) = \mathbb{P}(A_1 = a_1, \ldots, A_q = a_q) = \begin{cases} \prod_{j=1}^{q} p_j^{a_j} \cdot (1 - |\boldsymbol{p}|)^{1 - |\boldsymbol{a}|} & \text{if } |\boldsymbol{a}| \leq 1 \\ 0 & \text{if } |\boldsymbol{a}| > 1 \end{cases}$$

and we will write $\boldsymbol{A} \sim \mathsf{M}(1, \boldsymbol{p})$. In other words $\boldsymbol{A}$ concentrates on the vectors of $\{0, 1\}^q$ for which at most one element is equal to 1 and all the other entries are zero.

The Bernoulli process with a condiment assumes that each observation $\boldsymbol{Z}$ is a multivariate $\{0, 1\}^q$-valued stochastic process

$$\boldsymbol{Z}(w) = \sum_{i \geq 1} \boldsymbol{A}_i \delta_{w_i}(w)$$

where $(w_i)_{i \geq 1}$ are features in $\mathbb{W}$ and $(\boldsymbol{A}_i)_{i \geq 1}$ are independent simple multinomial random variables with parameter vector $\boldsymbol{p}_i = (p_{i,1}, \ldots, p_{i,q})$ as $i = 1, 2, \ldots$. Here $|\boldsymbol{p}_i|$ represents the probability that an individual displays feature $w_i$, while $p_{i,j}$ is the probability that the individual exhibits feature $w_i$ with condiment $j \in \{1, \ldots, q\}$. Thus, $\boldsymbol{Z}$ is termed a *simple mutlinomial process* with parameter $\boldsymbol{\zeta} = \sum_{i \geq 1} \boldsymbol{p}_i \delta_{w_i}$, and it is denoted by $\mathsf{MP}(\boldsymbol{\zeta})$. In order to carry out BNP inference, we need to specify

a distribution for the discrete measure $\boldsymbol{\zeta}$. Thus, we obtain a multivariate version of the model (1):

$$
\begin{aligned}
\boldsymbol{Z}_n | \boldsymbol{\zeta} &\stackrel{\text{iid}}{\sim} \mathsf{MP}(\boldsymbol{\zeta}) \quad n = 1, \ldots, N \\
\boldsymbol{\zeta} &\sim \mathscr{Z}
\end{aligned}
\tag{46}
$$

where $\mathscr{Z}$ denotes the distribution of the discrete random measure $\boldsymbol{\zeta}$.

## E.2 Priors based on multivariate CRMs

In this section we consider a class of priors $\mathscr{Z}$ in (46) defined by James [2017] and based on a multivariate extension of CRMs (see Daley and Vere-Jones [2008]). In particular consider a multivariate CRM on $\mathbb{W}$:

$$
\boldsymbol{\mu} = \sum_{i \geq 1} \boldsymbol{\rho}_i \delta_{W_i}
$$

where $\boldsymbol{\rho}_i = (\rho_{i,1}, \ldots, \rho_{i,q})$ is a vector of $[0,1]$-valued random jumps with the property $\sum_{i \geq 1} |\boldsymbol{\rho}_i| < +\infty$, the $W_i$'s are i.i.d. $\mathbb{W}$-valued random locations independent of the $\boldsymbol{\rho}_i$'s. Under this nonparametric prior each observation $\boldsymbol{Z}_n$ in (46) admits the representation $\boldsymbol{Z}_n | \boldsymbol{\mu} = \sum_{i \geq 1} \boldsymbol{A}_{n,i} \delta_{W_i}$, where $\boldsymbol{A}_{n,i} = (A_{n,i,1}, \ldots, A_{n,i,q}) | \boldsymbol{\mu} \stackrel{\text{ind}}{\sim} \mathsf{M}(1, \boldsymbol{\rho}_i)$. Note that the random measure $\boldsymbol{\mu}$ equals the vector of random measures $(\mu_1, \ldots, \mu_q)$, where

$$
\mu_j = \sum_{i \geq 1} \rho_{i,j} \delta_{W_i}, \quad j = 1, \ldots, q.
$$

As a simple CRM of Section A, the multivariate extension of a CRM is characterized by its Lévy-Khintchine representation:

$$
\mathbb{E}[e^{-\int_{\mathbb{W}} f_1(w) \mu_1(\mathrm{d}w) - \cdots - \int_{\mathbb{W}} f_q(w) \mu_q(\mathrm{d}w)}]
$$

$$
= \exp \left\{ - \int_{\mathbb{W}} \int_{\mathbb{R}_+^q} (1 - e^{-s_1 f_1(w) - \cdots - s_q f_q(w)}) \lambda_{(q)}(s_1, \ldots, s_q) \mathrm{d}s_1 \cdots \mathrm{d}s_q P(\mathrm{d}w) \right\}
$$

for arbitrary measurable functions $f_1, \ldots, f_d : \mathbb{W} \to \mathbb{R}_+$, where $P$ is a probability measure on $\mathbb{W}$. The multivariate Lévy intensity $\lambda_{(q)}$ is assumed to satisfy the integral condition

$$
\int_{\mathbb{R}_+^q} \min\{1, ||\boldsymbol{s}||\} \lambda_{(q)}(s_1, \ldots, s_q) \mathrm{d}s_1 \cdots \mathrm{d}s_q < +\infty
$$

where $||\boldsymbol{s}||$ is the Euclidean norm of the vector $\boldsymbol{s}$. When $\lambda_{(q)}(s_1, \ldots, s_q)$ concentrates on $S_q$, the law of $\boldsymbol{\mu}$ may be employed as a distribution for the parameter $\boldsymbol{\zeta}$ of the simple multinomial process in (46). A possible choice indicated by James [2017] is to select a stable-Beta-Dirichlet process, which is a generalization of the Beta-Dirichlet process [Kim et al., 2012] with power law behavior. We say that a multivariate CRM $\boldsymbol{\mu} = (\mu_1, \ldots, \mu_q)$ is a stable-Beta-Dirichlet process with parameters $(\alpha, \kappa + \alpha; \boldsymbol{\gamma}; \vartheta)$, where $\boldsymbol{\gamma} = (\gamma_1, \ldots, \gamma_q)$, if it is characterized by the following Lévy intensity specification

$$
\lambda_{(q)}(\boldsymbol{s}) = \frac{\vartheta \Gamma(|\boldsymbol{\gamma}|)}{\prod_{j=1}^q \Gamma(\gamma_j)} |\boldsymbol{s}|^{-\alpha - |\boldsymbol{\gamma}|} (1 - \boldsymbol{s})^{\kappa + \alpha - 1} \prod_{j=1}^q s_j^{\gamma_j - 1} \mathbb{1}_{[0,1]}(|\boldsymbol{s}|), \ \boldsymbol{s} \in S_q
\tag{47}
$$

where $0 \leq \alpha < 1, \kappa > -\alpha, \vartheta > 0$ and $\gamma_j > 0$ for any $j = 1, \ldots, q$. We write $\boldsymbol{\mu} \sim \mathrm{mSBD}(\alpha, \kappa + \alpha; \boldsymbol{\gamma}; \vartheta)$ to denote the distribution of the stable-Beta-Dirichlet process. As emphasized by James [2017], it can be easily checked, by means of the Laplace functional, that $\sum_{j=1}^{q} \mu_j$ is a stable-Beta process of Teh and Gorur [2009], i.e. a simple CRM on $\mathbb{W}$ with Lévy intensity on $[0, 1] \times \mathbb{W}$ equal to $\vartheta s^{-\alpha-1}(1-s)^{\kappa+\alpha-1}\mathrm{d}sP(\mathrm{d}w)$.

### E.2.1 Estimation of the unseen features with a condiment

In order to face predictive inference with the model (46) under the prior specification $\boldsymbol{\zeta} \sim \mathrm{mSBD}(\alpha, \kappa + \alpha; \boldsymbol{\gamma}; \vartheta)$, we need to characterize the predictive distribution of $\boldsymbol{Z}_{N+1}|\boldsymbol{Z}_{1:N}$ for the model (46). To this end it is worth recalling the definition of the finite-dimensional Beta-Dirichlet distribution by Kim et al. [2012]. A random vector $\boldsymbol{P} := (P_1, \ldots, P_q)$ on $S_q$ is said to follow a Beta-Dirichlet distribution with positive parameters $\alpha, \kappa$ and $\boldsymbol{\gamma} = (\gamma_1, \ldots, \gamma_q)$ if the probability density function of the random vector $(P_1, \ldots, P_q)$ has density proportional to

$$|\boldsymbol{s}|^{\alpha - |\boldsymbol{\gamma}|} \cdot (1 - |\boldsymbol{s}|)^{\kappa - 1} \prod_{j=1}^{q} s_j^{\gamma_j - 1} \cdot \mathbb{1}_{S_q}(\boldsymbol{s}) \tag{48}$$

and we write $(P_1, \ldots, P_q) \sim \mathscr{BD}(\alpha, \kappa; \boldsymbol{\gamma})$. This distribution can be characterized as follows: $|\boldsymbol{P}|$ has a Beta distribution with parameters $(\alpha, \kappa)$ and the normalized vector $(P_1/|\boldsymbol{P}|, \ldots, P_q/|\boldsymbol{P}|)$ follows a Dirichlet distribution with parameters $(\gamma_1, \ldots, \gamma_q)$.

We first characterize the distribution of $\boldsymbol{Z}_{N+1}|\boldsymbol{Z}_{1:N}$ under the prior specification $\boldsymbol{\zeta} \sim \mathrm{mSBD}(\alpha, \kappa + \alpha; \boldsymbol{\gamma}; \vartheta)$ in (46). The following result is immediate from the theory developed by James [2017].

**Theorem 4.** *For any $N \geq 1$, let $\boldsymbol{Z}_{1:N}$ be a random sample modeled as the BNP multinomial process model* (46), *with $\boldsymbol{\zeta} \sim \mathrm{mSBD}(\alpha, \kappa + \alpha; \boldsymbol{\gamma}; \vartheta)$. If $\boldsymbol{Z}_{1:N}$ displays $K_N = k$ distinct features, labeled by $W_1^*, \ldots, W_{K_N}^*$, with condiment-specific frequencies $(M_{N,1,j}, \ldots, M_{N,K_N,j}) = (m_{1,j}, \ldots, m_{k,j})$, for any $j = 1, \ldots, q$, then the conditional distribution of $\boldsymbol{Z}_{N+1}$, given $\boldsymbol{Z}_{1:N}$, coincides with the distribution of*

$$\boldsymbol{Z}_{N+1}|\boldsymbol{Z}_{1:N} \stackrel{d}{=} \boldsymbol{Z}'_{N+1} + \sum_{i=1}^{K_N} \boldsymbol{A}_{N+1,i}\delta_{W_i^*} \tag{49}$$

*where:*

i) *$\boldsymbol{Z}'_{N+1}$ is such that $\boldsymbol{Z}'_{N+1} = \sum_{i \geq 1} \boldsymbol{A}'_{N+1,i}\delta_{W_i'} \sim \mathsf{MP}(\boldsymbol{\mu}')$ and $\boldsymbol{\mu}' \sim \mathrm{mSBD}(\alpha, \kappa + M + \alpha; \boldsymbol{\gamma}; \vartheta)$;*

ii) *$\boldsymbol{A}_{N+1,1:K_N}$ is a collection of independent simple multinomial random variables with respective parameters $\boldsymbol{J}_{1:K_N}$, such that each $\boldsymbol{J}_i = (J_1, \ldots, J_q)$ has a Beta-Dirichlet distribution, i.e., $\boldsymbol{J}_i \stackrel{ind}{\sim} \mathscr{BD}(m_i - \alpha, N - m_i + \kappa + \alpha; \boldsymbol{\gamma} + \boldsymbol{m}_i)$, where we put $\boldsymbol{m}_i := (m_{i,1}, \ldots, m_{i,q})$ and $m_i = \sum_{j=1}^{q} m_{i,j} = |\boldsymbol{m}_i|$ for any $i = 1, \ldots, K_N$.*

Note that in Theorem 4 $M_{N,i,j}$ is the random number of times feature $W_i^*$ has been observed out of $\boldsymbol{Z}_{1:N}$ with condiment $j \in \{1, \ldots, q\}$, while $m_i = \sum_{j=1}^{q} m_{i,j}$ is the number of times feature $W_i^*$ has been observed out of the sample.

For any $N \geq 1$, let $\boldsymbol{Z}_{1:N}$ be an observable sample modeled as the multinomial model in (46), with $\boldsymbol{\zeta} \sim \mathrm{mSBD}(\alpha, \kappa + \alpha; \boldsymbol{\gamma}; \vartheta)$. Moreover, under the same model, for $M \geq 1$ let $\boldsymbol{Z}_{N+1:N+M} = (\boldsymbol{Z}_{N+1}, \ldots, \boldsymbol{Z}_{N+M})$ be an additional and unobserved sample. We now define the number of hitherto unobserved feature with condiment $\ell \in \{1, \ldots, q\}$ that will be recorded out of $\boldsymbol{Z}_{N+1:N+M}$ as

$$U_{N,\ell}^{(M)} := \sum_{i \geq 1} \mathbb{1}\left(\sum_{m=1}^{M} A_{m,i,\ell} > 0\right) \cdot \mathbb{1}\left(\sum_{n=1}^{N} A_{n,i,\ell} = 0\right). \tag{50}$$

Posterior inference for such a quantity could have potential interest in genomics to account for the presence of a variant with certain biological characteristics (condiment). The next theorem provides the posterior distribution of $U_{N,\ell}^{(M)}$.

**Theorem 5.** *For any $N \geq 1$, let $\boldsymbol{Z}_{1:N}$ be a random sample modeled as the BNP simple multinomial process model (46), with $\boldsymbol{\zeta} \sim \mathrm{mSBD}(\alpha, \kappa + \alpha; \boldsymbol{\gamma}; \vartheta)$. Suppose that $\boldsymbol{Z}_{1:N}$ displays $K_N = k$ distinct features, labeled by $W_1^*, \ldots, W_{K_N}^*$, with condiment-specific frequencies $(M_{N,1,j}, \ldots, M_{N,K_N,j}) = (m_{1,j}, \ldots, m_{k,j})$, for any $j = 1, \ldots, q$. Then, the posterior distribution of $U_{N,\ell}^{(M)}$, given $\boldsymbol{Z}_{1:N}$, coincides with the distribution of*

$$U_{N,\ell}^{(M)}|\boldsymbol{Z}_{1:N} \sim \mathrm{Poisson}\left(\vartheta \sum_{m=1}^{M} (-1)^{m+1} \binom{M}{m} B(m - \alpha, N + \alpha + \kappa) \frac{(\gamma_\ell)_m}{(|\boldsymbol{\gamma}|)_m}\right) \tag{51}$$

*Proof.* The proof is based on the posterior characterization provided in Theorem 4 and the evaluation of the probability generating function of the random variable $U_{N,\ell}^{(M)}$, conditionally on the sample $\boldsymbol{Z}_{1:N}$. The probability generating function is denoted as usual by $\mathcal{G}_{U_{N,\ell}^{(M)}}(\cdot)$. Thanks to the characterization (49), conditionally on $\boldsymbol{Z}_{1:N}$, the random variable $U_{N,\ell}^{(M)}$ may be written as

$$U_{N,\ell}^{(M)}|\boldsymbol{Z}_{1:N} \stackrel{d}{=} \sum_{i \geq 1} \mathbb{1}\left(\sum_{m=1}^{M} A'_{m+N,i,\ell} > 0\right).$$

Fix $t$ in a neighborhood of the origin, then one has

$$\mathcal{G}_{U_{N,\ell}^{(M)}}(t) = \mathbb{E}\left[t^{U_{N,\ell}^{(M)}} \mid \boldsymbol{Z}_{1:N}\right]. \tag{52}$$

Here, independently across $i$, $A'_{N+m,i,\ell}$ is a Bernoulli random variable with parameter $\rho'_{i,\ell}$, conditionally on the random measure $\boldsymbol{\mu}' = \sum_{i \geq 1} \boldsymbol{\rho}'_i \delta_{W'_i}$ with Lévy intensity $\lambda'_{(q)}(\boldsymbol{s})\mathrm{d}s_1 \cdots \mathrm{d}s_q P(\mathrm{d}w)$ such that

$$\lambda'_{(q)}(\boldsymbol{s}) = \frac{\vartheta \Gamma(|\boldsymbol{\gamma}|)}{\prod_{j=1}^{q} \Gamma(\gamma_j)} |\boldsymbol{s}|^{-\alpha - |\boldsymbol{\gamma}|}(1 - \boldsymbol{s})^{N + \kappa + \alpha - 1} \prod_{j=1}^{q} s_j^{\gamma_j - 1} \mathbb{1}_{[0,1]}(|\boldsymbol{s}|), \ \boldsymbol{s} \in S_q. \tag{53}$$

Thus, the expected value in (52) boils down to

$$\mathcal{G}_{U_{N,\ell}^{(M)}}(t) = \mathbb{E}\left[t^{\sum_{i\geq 1} \mathbb{1}\left(\sum_{m=1}^{M} A'_{m+N,i,\ell}>0\right)}\right] = \mathbb{E}\left[\prod_{i\geq 1} \mathbb{E}\left[t^{\mathbb{1}\left(\sum_{m=1}^{M} A'_{m+N,i,\ell}>0\right)} \mid \boldsymbol{\mu}'\right]\right]$$

$$= \mathbb{E}\left[\prod_{i\geq 1} \left(t + (1-t)\prod_{m=1}^{M} \mathbb{P}(A'_{m+N,i,\ell}=0|\boldsymbol{\mu}')\right)\right]$$

$$= \mathbb{E}\left[\prod_{i\geq 1}(t + (1-t)(1-\rho'_{i,\ell})^M)\right].$$

where we used the fact that each $A'_{m+N,i,\ell}$ is a Bernoulli random variable with parameter $\rho'_{i,\ell}$, conditionally on the random measure $\boldsymbol{\mu}'$, and in addition these random variables are conditionally independent. We now exploit the Laplace functional of the multivariate CRM $\boldsymbol{\mu}'$ to obtain

$$\mathcal{G}_{U_{N,\ell}^{(M)}}(t) = \mathbb{E}\left[\exp\left\{\sum_{i\geq 1}\log(t + (1-t)(1-\rho'_{i,\ell})^M)\right\}\right]$$

$$= \exp\left\{-(1-t)\int_{S_q}[1-(1-s_\ell)^M]\lambda'_{(q)}(\boldsymbol{s})\mathrm{d}s_1\cdots\mathrm{d}s_q\right\}$$

$$= \exp\left\{(1-t)\sum_{m=1}^{M}(-1)^m\binom{M}{m}\int_{S_q}s_\ell^m\lambda'_{(q)}(\boldsymbol{s})\mathrm{d}s_1\cdots\mathrm{d}s_q\right\} \tag{54}$$

where $\lambda'_{(q)}$ has been specified in (53) and we exploited the following formula

$$[1-(1-s_\ell)^M] = 1 - \sum_{m=0}^{M}(-1)^m\binom{M}{m}s_\ell^m = -\sum_{m=1}^{M}(-1)^m\binom{M}{m}s_\ell^m. \tag{55}$$

The integrals over $S_q$ in (54) may be easily evaluated (see, e.g., [Gradshteyn and Ryzhik, 2007, Formula 4.635.2]) to get

$$\int_{S_q}s_\ell^m\lambda'_{(q)}(\boldsymbol{s})\mathrm{d}s_1\cdots\mathrm{d}s_q = \vartheta\frac{(\gamma_\ell)_m}{(|\boldsymbol{\gamma}|)_m}\cdot B(m-\alpha, N+\alpha+\kappa).$$

By substituting the previous expression in (54), we obtain

$$\mathcal{G}_{U_{N,\ell}^{(M)}}(t) = \exp\left\{(t-1)\sum_{m=1}^{M}(-1)^{m+1}\binom{M}{m}\vartheta\frac{(\gamma_\ell)_m}{(|\boldsymbol{\gamma}|)_m}\cdot B(m-\alpha, N+\alpha+\kappa)\right\}$$

which is exactly the probability generating function of a Poisson random variable with parameter

$$\sum_{m=1}^{M}(-1)^{m+1}\binom{M}{m}\vartheta\frac{(\gamma_\ell)_m}{(|\boldsymbol{\gamma}|)_m}\cdot B(m-\alpha, N+\alpha+\kappa).$$

$\square$ $\square$

39

As a consequence of Theorem 5, one can define a BNP estimator of $U_{N,\ell}^{(M)}$ with respect to a squared loss function as follows:

$$\hat{U}_{N,\ell}^{(M)} = \vartheta \sum_{m=1}^{M} (-1)^{m+1} \binom{M}{m} B(m - \alpha, N + \alpha + \kappa) \frac{(\gamma_\ell)_m}{(|\boldsymbol{\gamma}|)_m}. \tag{56}$$

We point out that for computational convenience one may write

$$\hat{U}_{N,\ell}^{(M)} = \vartheta B(1 - \alpha, N + \alpha + \kappa) \mathbb{E}_{(X,Y)} \left[ \frac{1 - (1 - XY)^M}{Y} \right] \tag{57}$$

where the expected value is taken with respect to the two independent random variables with the following beta distributions

$$X \sim \mathrm{Beta}(\gamma_\ell, |\boldsymbol{\gamma}| - \gamma_\ell), \quad Y \sim \mathrm{Beta}(1 - \alpha, N + \alpha + \kappa).$$

The equality (57) may be easily proved by observing that

$$\mathbb{E}_X[X^m] = \frac{(\gamma_\ell)_m}{(|\boldsymbol{\gamma}|)_m} \quad \text{and} \quad B(m - \alpha, N + \alpha + \kappa) = \mathbb{E}_Y[Y^{m-1}] B(1 - \alpha, N + \alpha + \kappa).$$

### E.3   Scaled stable-Beta-Dirichlet prior for multinomial processes

From Theorems 4-5, it is apparent that, under the stable-Beta-Dirichlet process, the conditional distribution of a statistic involving hitherto unobserved features, depends on the initial sample $\boldsymbol{Z}_{1:N}$ only trough the sample size $N$ and not on other sample statistics. This behavior resembles what happens for the Bernoulli process model described in the main paper when the prior $\zeta$ in (1) is a CRM. We then introduce a multivariate analogue of the stable-Beta scaled prior, that will be termed *scaled stable-Beta-Dirichlet process* with the goal to enrich the predictive structure. We introduce a discrete random measure depending on the random jump $\Delta_{1,h_{c,\beta}}$, that has been defined in the main paper as a polynomial-exponential tilting of the density function (6), whose density equals

$$f_{\Delta_{1,h_{c,\beta}}}(a) = \frac{\sigma \beta^{c+1}}{\Gamma(c+1)} a^{-\sigma(c+1)-1} \exp\left\{ -\beta a^{-\sigma} \right\} \mathbb{1}_{\mathbb{R}_+}(a) \tag{58}$$

as shown in (36). The scaled stable-Beta-Dirichlet random measure is an almost surely discrete random measure that can be represented as

$$\boldsymbol{\mu}_{\Delta_{1,h_{c,\beta}}} = \sum_{i \geq 1} \boldsymbol{\rho}_i \delta_{W_i}, \quad \boldsymbol{\rho}_i = (\rho_{i,1}, \ldots, \rho_{i,q})$$

and consisting of $q$ components

$$\mu_{\Delta_{1,h_{c,\beta}},j} = \sum_{i \geq 1} \rho_{i,j} \delta_{W_i} \quad \text{as } j = 1, \ldots, q.$$

40

Conditionally on the jump $\Delta_{1,h_{c,\beta}}$, the multivariate random measure $\boldsymbol{\mu}_{\Delta_{1,h_{c,\beta}}}$ is completely random with Lévy intensity $\lambda_{(q),\Delta_{1,h_{c,\beta}}}(\boldsymbol{s})\mathrm{d}s_1\cdots\mathrm{d}s_q P(\mathrm{d}p)$ with the specification

$$\lambda_{(q),\Delta_{1,h_{c,\beta}}}(\boldsymbol{s}) = \frac{\Gamma(|\boldsymbol{\gamma}|)}{\prod_{j=1}^{q}\Gamma(\gamma_j)}\sigma\Delta_{1,h_{c,\beta}}^{-\sigma}|\boldsymbol{s}|^{-\sigma-|\boldsymbol{\gamma}|}\prod_{j=1}^{q}s_j^{\gamma_j-1}\mathbb{1}_{[0,1]}(|\boldsymbol{s}|), \ \boldsymbol{s}\in S_q \tag{59}$$

where $0 < \sigma < 1$ and $\gamma_j > 0$ for any $j = 1,\ldots,q$. We write $\boldsymbol{\mu}_{\Delta_{1,h}} \sim \text{S-mSBD}(\sigma,\boldsymbol{\gamma};h_{c,\beta})$. A remarkable property of this model is that $\sum_{j=1}^{q}\mu_{\Delta_{1,h_{c,\beta}},j}$ is distributed as the stable-Beta scaled process prior, i.e., $|\boldsymbol{\mu}_{\Delta_{1,h_{c,\beta}}}| \sim \text{SB-SP}(\sigma,c,\beta)$. Such a property may be easily proved by means of the Laplace functionals. Note that one could potentially introduce an additional mass parameter in the model, but this is irrelevant to carry out posterior inference in the stable case.

### E.3.1 Posterior Analysis

We now provide posterior, predictive and marginal characterizations for the multivariate model (46) under the scaled stable-Beta-Dirichlet process prior specification for $\mathscr{Z}$. The results we present here may be proved by exploiting [James, 2017, Section 5], conditionally on $\Delta_{1,h_{c,\beta}}$ and then by marginalizing over the mixing distribution (58). We omit the details.

**Theorem 6.** *For any $N \geq 1$, let $\boldsymbol{Z}_{1:N}$ be a random sample modeled as the BNP simple multinomial process model* (46), *with $\boldsymbol{\zeta} \sim \text{S-mSBD}(\sigma,\boldsymbol{\gamma};h_{c,\beta})$. If $\boldsymbol{Z}_{1:N}$ displays $K_N = k$ distinct features, labeled by $W_1^*,\ldots,W_{K_N}^*$, with condiment-specific frequencies $(M_{N,1,j},\ldots,M_{N,K_N,j}) = (m_{1,j},\ldots,m_{k,j})$, for any $j = 1,\ldots,q$, then the conditional distribution of $\Delta_{1,h_{c,\beta}}$ given $\boldsymbol{Z}_{1:N}$, coincides with the distribution of*

$$\Delta_{1,h_{c,\beta}}^{-\sigma} \sim \text{Gamma}(K_N + c + 1, \beta + \gamma_0^{(N)}) \tag{60}$$

*where $\gamma_0^{(n)} = \sigma\sum_{1\leq i\leq n}B(1-\sigma,i)$. Moreover, the conditional distribution of $\boldsymbol{\zeta}$, given $\boldsymbol{Z}_{1:N},\Delta_{1,h_{c,\beta}}$, coincides with the distribution of*

$$\boldsymbol{\zeta}|(\boldsymbol{Z}_{1:N},\Delta_{1,h_{c,\beta}}) \stackrel{d}{=} \boldsymbol{\mu}'_{\Delta_{1,h_{c,\beta}}} + \sum_{i=1}^{K_N}\boldsymbol{J}_i\delta_{W_i^*} \tag{61}$$

*where:*

*i)* $\boldsymbol{\mu}'_{\Delta_{1,h_{c,\beta}}}$ *is a discrete multivariate random measure with Lévy intensity*

$$\nu'_{\Delta_{1,h_{c,\beta}}}(\mathrm{d}s_1,\ldots,\mathrm{d}s_q,\mathrm{d}w) = \frac{\Gamma(|\boldsymbol{\gamma}|)}{\prod_{j=1}^{q}\Gamma(\gamma_j)}$$
$$\times |\boldsymbol{s}|^{-\sigma-|\boldsymbol{\gamma}|}(1-\boldsymbol{s})^N\prod_{j=1}^{q}s_j^{\gamma_j-1}\mathbb{1}_{[0,1]}(|\boldsymbol{s}|)\sigma\Delta_{1,h_{c,\beta}}^{-\sigma}\mathrm{d}s_1\cdots\mathrm{d}s_q\,P(\mathrm{d}w); \tag{62}$$

*ii)* $\boldsymbol{J}_{1:K_N}$ *is a vector of independent random jumps such that each $\boldsymbol{J}_i = (J_1,\ldots,J_q)$ has a Beta-Dirichlet distribution, i.e.,*

$$\boldsymbol{J}_i|\Delta_{1,h_{c,\beta}} \sim \mathscr{BD}(m_i - \sigma, N - m_i + 1; \boldsymbol{\gamma} + \boldsymbol{m}_i) \tag{63}$$

where we put $\boldsymbol{m}_i := (m_{i,1}, \ldots, m_{i,q})$ and $m_i = \sum_{j=1}^q m_{i,j} = |\boldsymbol{m}_i|$ for any $i = 1, \ldots, K_N$.

**Theorem 7.** *For any $N \geq 1$, let $\boldsymbol{Z}_{1:N}$ be a random sample modeled as the BNP simple multinomial process model* (46), *with $\boldsymbol{\zeta} \sim$ S-mSBD$(\sigma, \boldsymbol{\gamma}; h_{c,\beta})$. If $\boldsymbol{Z}_{1:N}$ displays $K_N = k$ distinct features, labeled by $W_1^*, \ldots, W_{K_N}^*$, with condiment-specific frequencies $(M_{N,1,j}, \ldots, M_{N,K_N,j}) = (m_{1,j}, \ldots, m_{k,j})$, for any $j = 1, \ldots, q$, then the conditional distribution of $\Delta_{1,h_{c,\beta}}$ given $\boldsymbol{Z}_{1:N}$, coincides with* (60). *Moreover, the conditional distribution of $\boldsymbol{Z}_{N+1}$, given $\boldsymbol{Z}_{1:N}, \Delta_{1,h_{c,\beta}}$, coincides with the distribution of*

$$\boldsymbol{Z}_{N+1} | (\boldsymbol{Z}_{1:N}, \Delta_{1,h_{c,\beta}}) \stackrel{d}{=} \boldsymbol{Z}'_{N+1} + \sum_{i=1}^{K_N} \boldsymbol{A}_{N+1,i} \delta_{W_i^*} \tag{64}$$

*where:*

i) *$\boldsymbol{Z}'_{N+1}$ is such that $\boldsymbol{Z}'_{N+1} | \Delta_{1,h_{c,\beta}} = \sum_{i \geq 1} \boldsymbol{A}'_{N+1,i} \delta_{W'_i} \sim$ MP$(\boldsymbol{\mu}'_{\Delta_{1,h_{c,\beta}}})$ and $\boldsymbol{\mu}'_{\Delta_{1,h_{c,\beta}}} | \Delta_{1,h_{c,\beta}}$ is the completely random measure having the Lévy intensity* (62);

ii) *$\boldsymbol{A}_{N+1,1:K_N}$ is a collection of independent simple multinomial random variables with parameters $\boldsymbol{J}_{1:K_N}$, each one distributed according to Equation* (63).

**Theorem 8.** *For any $N \geq 1$, let $\boldsymbol{Z}_{1:N}$ be a random sample modeled as the BNP simple multinomial process model* (46), *with $\boldsymbol{\zeta} \sim$ S-mSBD$(\sigma, \boldsymbol{\gamma}; h_{c,\beta})$. The probability that $\boldsymbol{Z}_{1:N}$ displays a particular feature allocation of $K_N = k$ distinct features with condiment-specific frequencies $(M_{N,1,j}, \ldots, M_{N,K_N,j}) = (m_{1,j}, \ldots, m_{k,j})$, for any $j = 1, \ldots, q$, equals*

$$p_k^{(N)}(\boldsymbol{m}_1, \ldots, \boldsymbol{m}_k) = \prod_{i=1}^k \left\{ B(m_i - \sigma, N - m_i + 1) \frac{\prod_{j=1}^q (\gamma_j)_{m_{i,j}}}{(|\boldsymbol{\gamma}|)_{m_i}} \right\} \\ \times \frac{\Gamma(k+c+1)}{\Gamma(c+1)} \cdot \frac{\sigma^k \beta^{c+1}}{(\beta + \gamma_0^{(N)})^{k+c+1}}. \tag{65}$$

### E.3.2 Estimation of the unseen features with a condiment

For any $N \geq 1$, let $\boldsymbol{Z}_{1:N}$ be an observable sample modeled as the simple multinomial model in (46), with $\boldsymbol{\zeta} \sim$ S-mSBD$(\sigma, \boldsymbol{\gamma}; h_{c,\beta})$. Moreover, under the same model, for $M \geq 1$ let $\boldsymbol{Z}_{N+1:N+M} = (\boldsymbol{Z}_{N+1}, \ldots, \boldsymbol{Z}_{N+M})$ be an additional and unobserved sample. Under this model, we now determine the posterior distribution of the sample statistic $U_{N,\ell}^{(M)}$ in (50), counting the number of hitherto unobserved feature with condiment $\ell \in \{1, \ldots, q\}$ that will be recorded out of the additional sample.

**Theorem 9.** *For any $N \geq 1$, let $\boldsymbol{Z}_{1:N}$ be a random sample modeled as the BNP simple multinomial process model* (46), *with $\boldsymbol{\zeta} \sim$ S-mSBD$(\sigma, \boldsymbol{\gamma}; h_{c,\beta})$. Suppose that $\boldsymbol{Z}_{1:N}$ displays $K_N = k$ distinct features with condiment-specific frequencies $(M_{N,1,j}, \ldots, M_{N,K_N,j}) = (m_{1,j}, \ldots, m_{k,j})$, for any $j = 1, \ldots, q$. Then, the posterior distribution of $U_{N,\ell}^{(M)}$, given $\boldsymbol{Z}_{1:N}$, coincides with the distribution of*

$$U_{N,\ell}^{(M)} | \boldsymbol{Z}_{1:N} \sim \text{NegativeBinomial} \left( K_N + c + 1, \frac{\psi_{N,\ell}^{(M)}}{\psi_{N,\ell}^{(M)} + \gamma_0^{(N)} + \beta} \right) \tag{66}$$

*where we defined*

$$\psi_{N,\ell}^{(M)} := \sigma \sum_{m=1}^{M} \binom{M}{m} (-1)^{m+1} \frac{(\gamma_\ell)_m}{(|\boldsymbol{\gamma}|)_m} B(m - \sigma, N + 1).$$

*Proof.* The proof is based on the posterior characterization in Theorem 6 and on Theorem 7. As in the proof of Theorem 5 we evaluate the probability generating function of the random variable $U_{N,\ell}^{(M)}$, conditionally on the sample $\boldsymbol{Z}_{1:N}$. The probability generating function is denoted as usual by $\mathcal{G}_{U_{N,j}^{(M)}}(\,\cdot\,)$. Thanks to the characterization (64), conditionally on $\boldsymbol{Z}_{1:N}, \Delta_{1,h_{c,\beta}}$, the random variable $U_{N,\ell}^{(M)}$ may be written as

$$U_{N,\ell}^{(M)} | (\boldsymbol{Z}_{1:N}, \Delta_{1,h_{c,\beta}}) \stackrel{d}{=} \sum_{i \geq 1} \mathbb{1}\left( \sum_{m=1}^{M} A'_{m+N,i,\ell} > 0 \right).$$

Fix $t$ in a neighborhood of the origin, then one has

$$\mathcal{G}_{U_{N,\ell}^{(M)}}(t) = \mathbb{E}\left[ t^{U_{N,\ell}^{(M)}} \mid \boldsymbol{Z}_{1:N} \right] = \mathbb{E}\left[ \mathbb{E}\left[ t^{U_{N,\ell}^{(M)}} \mid \boldsymbol{Z}_{1:N}, \Delta_{1,h_{c,\beta}} \right] \mid \boldsymbol{Z}_{1:N} \right] \tag{67}$$

by an application of the tower property. We now focus on the evaluation of the inner expected value in (67):

$$\mathbb{E}\left[ t^{U_{N,\ell}^{(M)}} \mid \boldsymbol{Z}_{1:N}, \Delta_{1,h_{c,\beta}} \right] = \mathbb{E}\left[ t^{\sum_{m=1}^{M} A'_{m+N,i,\ell}} \right]$$

$$= \mathbb{E}\left[ \prod_{i \geq 1} \mathbb{E}[1 \cdot \mathbb{P}(\sum_{m=1}^{M} A'_{m+N,i,\ell} = 0) + t \cdot \mathbb{P}(\sum_{m=1}^{M} A'_{m+N,i,\ell} > 0)] \right].$$

From Theorem 7, the $A'_{m+N,i,\ell}$s are independent random variables as $m = 1, \ldots, M$, and each one $A'_{N+m,i,\ell}$ is a Bernoulli with parameter $\rho'_{i,\ell}$, conditionally on the random measure $\boldsymbol{\mu}'_{\Delta_{1,h_{c,\beta}}} = \sum_{i \geq 1} \boldsymbol{\rho}'_i \delta_{W'_i}$ with Lévy intensity (62). As a consequence we obtain

$$\mathbb{E}\left[ t^{U_{N,\ell}^{(M)}} \mid \boldsymbol{Z}_{1:N}, \Delta_{1,h_{c,\beta}} \right] = \mathbb{E}\left[ \prod_{i \geq 1} \left[ t + (1 - t)(1 - \rho'_{i,\ell})^M \right] \right].$$

Proceeding along the same lines as in the proof of Theorem 5 we have that

$$\mathbb{E}\left[ t^{U_{N,\ell}^{(M)}} \mid \boldsymbol{Z}_{1:N}, \Delta_{1,h_{c,\beta}} \right] = \mathbb{E}\left[ \exp\left\{ \sum_{i \geq 1} \log(t + (1 - t)(1 - \rho_{i,\ell})^M) \right\} \right]$$

$$= \exp\left\{ -(1 - t) \int_{\mathbb{W}} \int_{S_q} [1 - (1 - s_\ell)^M] \nu'_{\Delta_{1,h_{c,\beta}}}(\mathrm{d}s_1, \ldots, \mathrm{d}s_q, \mathrm{d}w) \right\}.$$

Now define

$$\lambda'_{(q),\Delta_{1,h_{c,\beta}}}(\boldsymbol{s}) := \frac{\Gamma(|\boldsymbol{\gamma}|)}{\prod_{j=1}^{q} \Gamma(\gamma_j)} |\boldsymbol{s}|^{-\sigma - |\boldsymbol{\gamma}|} (1 - \boldsymbol{s})^N \prod_{j=1}^{q} s_j^{\gamma_j - 1} \mathbb{1}_{[0,1]}(|\boldsymbol{s}|) \sigma \Delta_{1,h_{c,\beta}}^{-\sigma}$$

thus, the conditional expected value under study may be written as

$$\mathbb{E}\left[t^{U_{N,\ell}^{(M)}} \mid \boldsymbol{Z}_{1:N}, \Delta_{1,h_{c,\beta}}\right] = \exp\left\{-(1-t)\int_{S_q}[1-(1-s_\ell)^M]\lambda'_{(q),\Delta_{1,h_{c,\beta}}}(\boldsymbol{s})\mathrm{d}s_1,\ldots,\mathrm{d}s_q\right\}$$

$$= \exp\left\{-(1-t)\sum_{m=1}^{M}\binom{M}{m}(-1)^{m+1}\int_{S_q}s_\ell^m\lambda'_{(q),\Delta_{1,h_{c,\beta}}}(\boldsymbol{s})\mathrm{d}s_1,\ldots,\mathrm{d}s_q\right\} \qquad (68)$$

where we applied (55). The integral over $S_q$ appearing in (68) may be evaluated resorting to [Gradshteyn and Ryzhik, 2007, Formula 4.635.2], therefore

$$\int_{S_q}s_\ell^m\lambda'_{(q),\Delta_{1,h_{c,\beta}}}(\boldsymbol{s})\mathrm{d}s_1,\ldots,\mathrm{d}s_q$$

$$= \sigma\Delta_{1,h_{c,\beta}}^{-\sigma}\frac{\Gamma(|\boldsymbol{\gamma}|)}{\prod_{j=1}^{q}\Gamma(\gamma_j)}\int_{S_q}s_\ell^m(1-|\boldsymbol{s}|)^N|\boldsymbol{s}|^{-\sigma-|\boldsymbol{\gamma}|}\prod_{j=1}^{q}s_j^{\gamma_j-1}\mathrm{d}s_1\cdots\mathrm{d}s_q$$

$$= \sigma\Delta_{1,h_{c,\beta}}^{-\sigma}\frac{(\gamma_\ell)_m}{(|\boldsymbol{\gamma}|)_m}B(m-\sigma,N+1).$$

Thus, by substituting the previous expression in (68) one obtains

$$\mathbb{E}\left[t^{U_{N,\ell}^{(M)}} \mid \boldsymbol{Z}_{1:N}, \Delta_{1,h_{c,\beta}}\right] = \exp\left\{-(1-t)\Delta_{1,h_{c,\beta}}^{-\sigma}\psi_{N,\ell}^{(M)}\right\} \qquad (69)$$

where we recall that $\psi_{N,\ell}^{(M)}$ has been defined as follows

$$\psi_{N,\ell}^{(M)} = \sigma\sum_{m=1}^{M}\binom{M}{m}(-1)^{m+1}\frac{(\gamma_\ell)_m}{(|\boldsymbol{\gamma}|)_m}B(m-\sigma,N+1).$$

As a consequence, the probability generating function in (67) equals

$$\mathcal{G}_{U_{N,\ell}^{(M)}}(t) \stackrel{(69)}{=} \mathbb{E}\left[\exp\left\{-(1-t)\Delta_{1,h_{c,\beta}}^{-\sigma}\psi_{N,\ell}^{(M)}\right\} \mid \boldsymbol{Z}_{1:N}\right].$$

The conclusion follows by a marginalization w.r.t. the posterior distribution of $\Delta_{1,h_{c,\beta}}^{-\sigma}$ which is a gamma random variable (see (60)):

$$\mathcal{G}_{U_{N,\ell}^{(M)}}(t) = \int_0^\infty e^{-(1-t)\psi_{N,\ell}^{(M)}x}\cdot\frac{(\beta+\gamma_0^{(N)})^{K_N+c+1}}{\Gamma(K_N+c+1)}x^{K_N+c}e^{-x(\gamma_0^{(N)}+\beta)}\mathrm{d}x$$

$$= \frac{(\beta+\gamma_0^{(N)})^{K_N+c+1}}{\Gamma(K_N+c+1)}\int_0^\infty e^{-[(\gamma_0^{(N)}+\beta)+(1-t)\psi_{N,\ell}^{(M)}]x}x^{K_N+c+1-1}\mathrm{d}x$$

$$= \frac{(\beta+\gamma_0^{(N)})^{K_N+c+1}}{[(\gamma_0^{(N)}+\beta)+(1-t)\psi_{N,\ell}^{(M)}]^{K_N+c+1}}$$

$$= \left(\frac{\beta+\gamma_0^{(N)}}{\gamma_0^{(N)}+\beta+\psi_{N,\ell}^{(M)}-t\psi_{N,\ell}^{(M)}}\right)^{K_N+c+1}$$

which is the probability generating function of a negative binomial distribution as in the statement.

$\square$ $\square$

As a consequence of Theorem 9, the BNP estimator of $U_{N,\ell}^{(M)}$ under a squared loss function equals

$$\hat{U}_{N,\ell}^{(M)} = (K_N + c + 1)\frac{\psi_{N,\ell}^{(M)}}{\gamma_0^{(N)} + \beta}. \tag{70}$$

For computational purposes, we finally note that the parameter $\psi_{N,\ell}^{(M)}$ in the posterior representations may be computed as

$$\psi_{N,\ell}^{(M)} = B(1 - \sigma, N + 1)\mathbb{E}_{(X,Y)}\left[\frac{1 - (1 - XY)^M}{Y}\right]$$

where the expected value is made w.r.t. two independent random variables $X$ and $Y$ having beta distributions as follows

$$X \sim \text{Beta}(\gamma_\ell, |\boldsymbol{\gamma}| - \gamma_\ell) \quad \text{and} \quad Y \sim \text{Beta}(1 - \sigma, N + 1).$$

# F  Synthetic experiments from the model

We now analyze empirically the properties of the SB-SP-Bernoulli model used in Section 3. We will use the acronym SSB for brevity in the captions. I.e., we consider the hierarchical model detailed in (1), with $\mu \sim \text{SB-SP}(\sigma, c, \beta)$. The predictive characterization detailed in Proposition 3, together with Equation (13), provides an algorithm to sample $N$ observations from the model: given $\beta > 0, \sigma \in (0,1), c > 0$,

- at every step $n = 1, \ldots, N$, conditionally on the previous $n - 1$ samples $Z_{1:n-1}$ showing $K_{n-1}$ distinct features, each feature $k = 1, \ldots, m_{K_{n-1}}$ with frequency $m_k$, sample

  – a random number of new features observed:

  $$U_{n-1}^{(1)} \mid Z_{1:n-1} \sim \text{NegativeBinonial}\left(K_{n-1} + c + 1, \frac{\gamma_{n-1}^{(1)}}{\beta + \gamma_0^{(n)}}\right);$$

  – for previously observed feature $i = 1, \ldots, K_{n-1}$:

  $$A_{n,i} \mid Z_{1:n-1} \sim \text{Bernoulli}\left(m_i - \sigma, N - m_i + 1\right);$$

In particular, for the first step, $K_0 = 0$.

## F.1  Predictive behavior of the number of new features from the prior

First, we investigate the predictive behavior of the model as we vary the hyperparameters of the process — $\beta, \sigma, c$. Because our interest is in understanding the coverage properties of the posterior predictive distribution induced by the model, we report, together with the predictive mean, also posterior predictive credible intervals. In this first set of simulations reported in Appendices F.1.1 to F.1.3, we assume the hyperparameter $\sigma, c, \beta$ to be known.
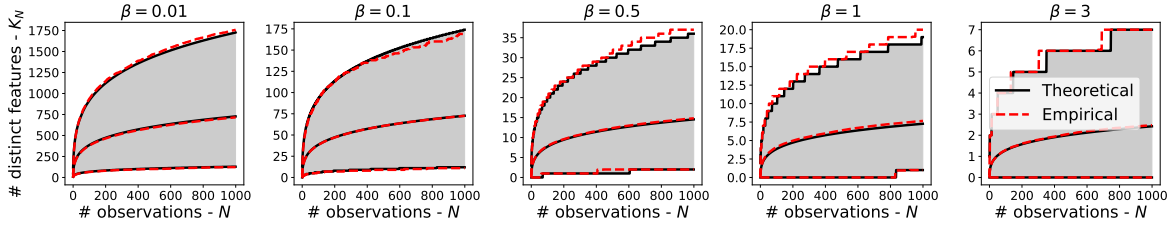
Figure 6: 90% centered credible interval for the number of distinct features $K_N$ ($y$-axis) as a function of the sample size $N$ ($x$-axis). We fix $\beta = 1$, $c = 5$, and vary $\sigma$ across subplots. For the $5\%, 50\%, 95\%$ quantiles, we compare the theoretical value (solid black lines) to empirical result (dashed red lines), obtained by drawing $N_{MC} = 1000$ different datasets with the same parameter specification.



Figure 7: 90% centered credible interval for the number of distinct features $K_N$ ($y$-axis) as a function of the sample size $N$ ($x$-axis). We repeat the same experiments as in Figure 6, but now fix $\beta = 1$, $\sigma = 0.2$, and vary $c$ across subplots.

### F.1.1    The role of $\sigma$

We start by analyzing the role of $\sigma$ in Figure 6. As suggested by the asymptotic behavior analyzed in Theorem 3, $\sigma$ directly controls the asymptotic rate of growth of the number of distinct features: as $\sigma$ increases, the expected number of variants increases, approaching a linear behavior as $\sigma \to 1$. We notice that this behavior is reminiscent of the tail parameter of the stable beta-Bernoulli process [Teh and Gorur, 2009, Broderick et al., 2012].

### F.1.2    The role of c

We now move to the analysis of the polynomial tilting parameter, $c$. As suggested by the predictive distribution given in Equation (13), $c$ acts as a "prior" number of features. That is, in the prior, the expected number of features to be observed from $N$ samples is a Negative Binomial random variable with parameters $c + 1, \gamma_0^{(N)}/(\beta + \gamma_0^{(N)})$, i.e. with expectation given by

$$\mathbb{E}[U_0^N] = (c + 1)\left(\frac{\gamma_0^{(N)}}{\beta}\right).$$

Again, larger values of $c$ induce a higher rate of growth in the number of features, as showed in Figure 7.

### F.1.3    The role of $\beta$

Last, we analyze the role of the exponential tilting parameter, $\beta$. Inspecting again the predictive distribution Equation (13), $\beta$ affects the number of new variants thorugh the success probability of

46

Figure 8: 90% centered credible interval for the number of distinct features $K_N$ ($y$-axis) as a function of the sample size $N$ ($x$-axis). We repeat the same exercise as in Figures 6-7 but now $c = 1$, $\sigma = 0.2$, and vary $\beta$ across subplots.
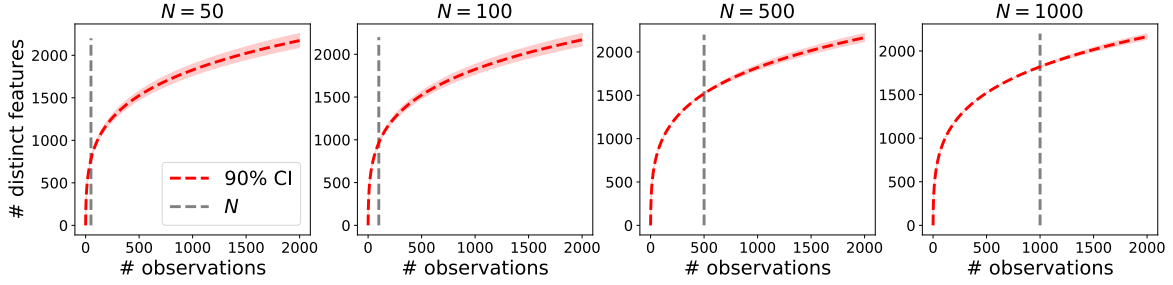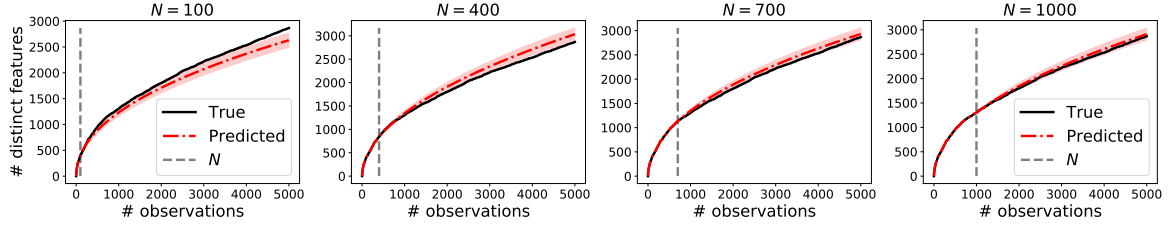


Figure 9: 90% centered credible interval for the expected number of distinct features $\mathbb{E}[U_N^{(M)} \mid Z_{1:N}]$ ($y$-axis) as a function of the sample size $N$ ($x$-axis). We fix $\beta = 1$, $c = 5$, $\sigma = 0.5$, and total sequencing capacity $L = 2000$. In different subplots, we show $\mathbb{E}[U_N^{(M)} \mid Z_{1:N}]$ for different values of $N$. Here, the first $N$ samples display exactly $K_N = \mathbb{E}[U_0^{(N)}]$ distinct features.

the negative binomial — for fixed $c, \sigma, N, M, Z_{1:N}$, the expected number of new variants $U_N^{(M)} \mid Z_{1:N}$ depends inversely on the parameter $\beta$. We verify this empirically in Figure 8.

## F.2  Predictive behavior of the number of new features from the posterior

Next, we perform a slightly different exercise from the one described above. We still assume the parameters to be known, and we investigate how the posterior predictive behavior varies as we change the number of training samples $N$ with respect to a total sampling "capacity" $L$. Intuitively, for a fixed value of this "sampling capacity", $N + M = L$, the expected number of observed features from the model should be the independent of the choice of $N, M$. However, we expect the distribution (e.g., the posterior variance), to concentrate as $N$ increases relative to $M$. To perform this experiment, we do as follow: we fix $\beta, c, \sigma$ and, for each $\ell = 1, \ldots, 2000$, we let $K_\ell = U_0^{(\ell)}$. Next, for $N \in \{50, 100, 500, 1000\}$, we compute $U_N^{(M)} \mid Z_{1:N}$, where we condition on the number of observed variants as given by the curve $\{K_\ell\}_{\ell=1,\ldots,2000}$. As displayed in Figure 9 and Figure 10, the width of the credible intervals shrinks with increasing training sizes $N$.

## F.3  Estimation of the parameters

Next, we move to the more interesting scenario in which the parameters are unknown and need to be inferred from the data. The natural way to estimate the unknown parameters is to maximize a likelihood criterion, such as the marginal distribution of the feature counts $m_1, \ldots, m_K$, given in

Figure 10: 90% centered credible interval for the expected number of distinct features $\mathbb{E}[U_N^{(M)} \mid Z_{1:N}]$ ($y$-axis) as a function of the sample size $N$ ($x$-axis). We repeat the same exercise as in Figure 9 but now fix $\beta = 2$, $c = 1000$, $\sigma = 0.2$.
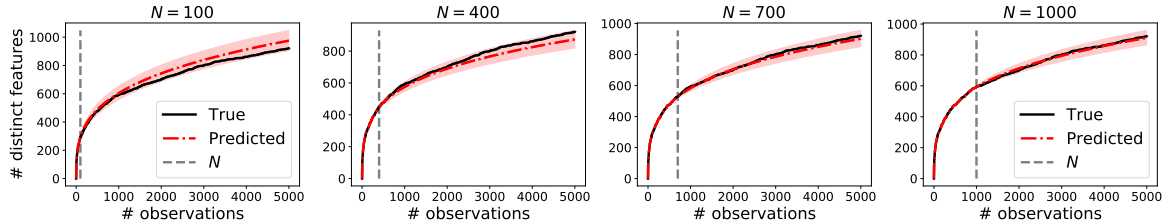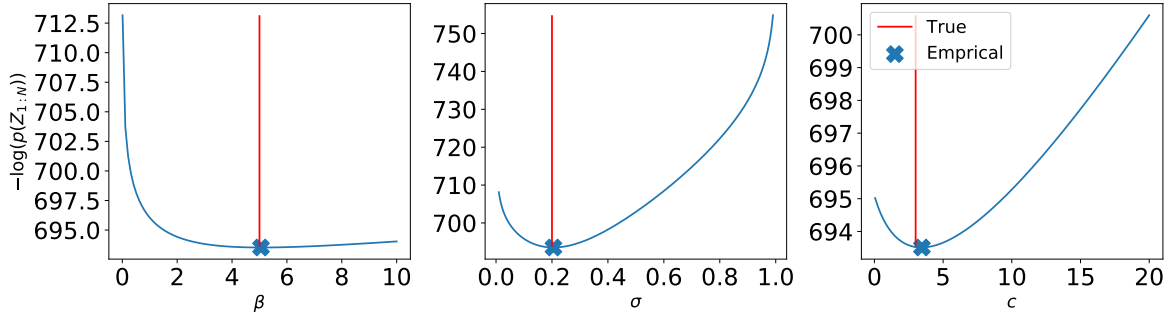


Figure 11: 99% credible interval centered around the posterior predictive mean (dashed red line) of the number of distinct features $U_N^{(M)} \mid Z_{1:N}$ ($y$-axis) as a function of the sample size $N$ ($x$-axis). We fix $\beta = 1$, $c = 20$, $\sigma = 0.5$, and learn the parameters for different training size $N \in \{100, 400, 700, 1000\}$ across subplots for a total sequencing capacity $L = 5000$.

Equation (10). We found this method to work well both on real data, as displayed in Section 4, and on synthetic data. We here report some results in Figures 11 and 12. In general, and not surprisingly, the precision of our estimates increases with larger sample sizes.

In our synthetic experiments, as expected, the values maximizing the marginal likelihood converge to the underlying true values of the data generating process as the sample size $N \to \infty$. By performing a visual investigation, we find that indeed the negativd marginal likelihood is a convex function in each argument, with a unique, well-defined minimum (see Figures 13 to 15).

When most of the features are very rare (e.g., they appear once or twice in the sample), we found that an alternative empirical Bayes approach, akin to the one adopted in Masoero et al. [2021], worked better, as further discussed in Appendix G.



Figure 12: 99% credible interval centered around the posterior predictive mean (dashed red line) of the number of distinct features $U_N^{(M)} \mid Z_{1:N}$ ($y$-axis) as a function of the sample size $N$ ($x$-axis). We repeat the same experiment as in Figure 11, but now for $\beta = 1$, $c = 100$, $\sigma = 0.25$.

Figure 13: We draw a synthetic dataset of size $N = 10'000$ from a SSB with parameters $\beta = 5, \sigma = 0.1, c = 3$. In the left subplot, we plot the value of the negative marginal likelihood (vertical axis) as we vary the value of $\beta$ (horizontal axis), keeping $\sigma = 0.1$ and $c = 3$ fixed at the true value. We repeat the same procedure, now varying $\sigma$ and keeping $\beta = 5, c = 3$ fixed at their true value in the central subplot. Last, in the right subplot, we inspect the marginal likelihood as we vary the value of $c$, keeping $\beta = 5, \sigma = 0.1$. We then minimize numerically the negative log-likelihood, and report in each subplot with a blue cross the numerical value of the corresponding hyperparameter (horizontal axis) together with the corresponding marginal likelihood value (vertical axis).



Figure 14: We repeat the same exercise of Figure 13 for $N = 1'000$, $\beta = 10$, $\sigma = 0.7$, $c = 20$.



Figure 15: We repeat the same exercise as in Figures 13 and 14 for $N = 100$, $\beta = 5$, $\sigma = 0.35$ and $c = 102$.

49

Figure 16: Frequencies distribution for different choices of the parameter $\xi$.



Figure 17: Estimates for the number of new features for the SB-SP-Bernoulli (red) and the stable beta-Bernoulli (blue) processes as the exponent $\xi$ varies across subplots. Shaded regions cover a 95% credible interval around the predictive mean. The solid black line represents the true counts. Here, the training is done using the first $N = 100$ observations, and extrapolating up to the remaining $M = 1900$ observations.

# G Synthetic experiments from Zipf distributions

To compare the predictive performance of the SB-SP-Bernoulli process proposed in Section 3.1 to existing competing methods, we also consider synthetic data from Zipf-distributed frequencies (see Figure 16). That is to say, we imagine that there exists a countable number of features in the population, and that, for some $\xi > 0$, feature $k$ is observed independently of any other feature with probability $\pi_k = (k + 1)^{-\xi}$. An observation $X_\ell$ then a binary vector, in which, conditionally on the frequencies $\pi = (\pi_1, \pi_2, \ldots)$ the $k$-th coordinate is a Bernoulli random variable:

$$X_{\ell,k} \mid \pi \overset{i.i.d.}{\sim} \text{Bernoulli}(\pi_k), \quad \pi = \{(k+1)^{-\xi}\}_{k \geq 1}. \tag{71}$$

We perform our simulations as follows: we fix a total sequencing capacity of $L = 2000$, and draw $L$ i.i.d. samples from the model, following the recipe given in Equation (71). For simulation purposes, we only consider the first $K = 10^5$ features to have non-zero probability, i.e. $\pi_k = 0$, for all $k > K$. We compare the estimates of our proposed SB-SP-Bernoulli model (Section 3.1), to the stable beta-Bernoulli process [3BP], the linear program of Zou et al. [2016], the first four orders of the Jackknife estimator originally proposed in Burnham and Overton [1978] and recently employed in the genomics context by Gravel [2014], and the Good-Toulmin estimator, recently used in Chakraborty et al. [2019], with the two alternative smoothing choices described in Orlitsky et al. [2016]. Estimates for Bayesian methods are obtained by using the posterior predictive mean for the number of new variants conditionally on the observed sample, with hyper-parameters learned by numerically maximizing the marginal distribution (EFPF) of the features counts, as described in Section 4.1.

As expected, we find the nonparametric Bayesian estimators to do particularly well for larger values of the exponent $\xi$ — that is when most features are exceedingly rare. The SB-SP-Bernoulli and the SB-SP-Bernoulli-parameter beta-Bernoulli processes performed comparably on these datasets, both in terms of estimation accuracy and uncertainty quantification, as displayed in Figure 17.
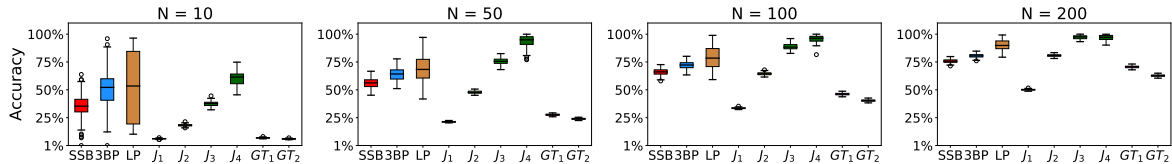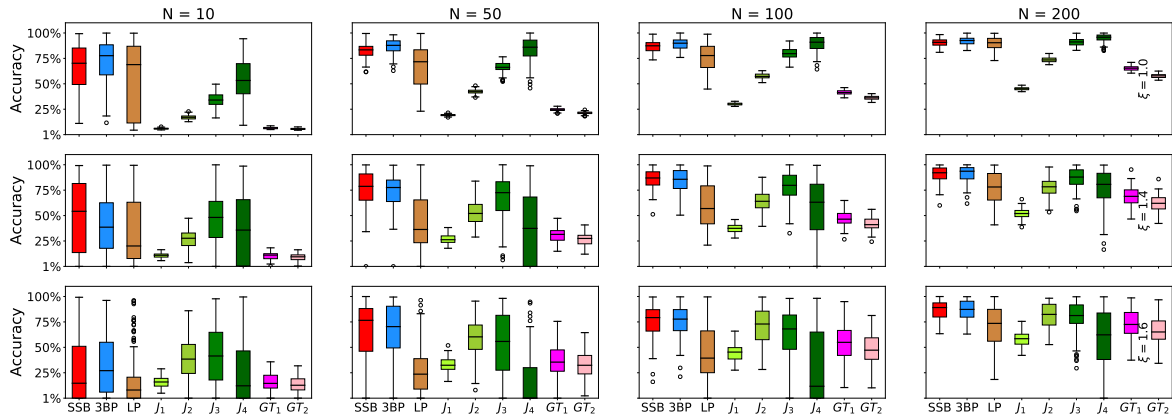
50

Figure 18: Accuracy of the competing methods (SB-SP-Bernoulli [SSP], stable beta-Bernoulli [3BP], Jackknife [J], linear program [LP], Good-Toulmin [GT]) on simulated data from a Zipf model (Equation (71)) with parameter $\xi = 1.2$. For $L = 2000$, we report $v_{N,a}^{(M)}$ as $N$ increases, for $M = L - N$. For each $N$, results across $S = 100$ datasets are reported in the boxplots.



Figure 19: Accuracy of the competing methods (SB-SP-Bernoulli [SSP], stable beta-Bernoulli [3BP], Jackknife [J], linear program [LP], Good-Toulmin [GT]) on simulated data from a Zipf model (Equation (71)) with parameter $\xi = 0.8$. For $L = 2000$, we report $v_{N,a}^{(M)}$ as $N$ increases, for $M = L - N$. For each $N$, results across $S = 100$ datasets are reported in the boxplots.

To better asses the predictive quality of the different methods, we ran extensive simulation experiments; for each value of $\xi \in \{0.8, 1, 1.2, 1.4, 1.6\}$, we generated $S = 100$ datasets of size $L = 2000$, and for each value of $N \in \{10, 50, 100, 200\}$ we trained each method, and extrapolated to predict the number of new variants to be observed up to $M = L - N \in \{1990, 1950, 1900, 1800\}$ remaining samples. We report as measure of accuracy the percentage accuracy incurred by each estimation method $v_{N,a}^{(M)}$, defined in Equation (19), at the largest extrapolation level $M = L - N$, across different values of $N$ and all $S = 100$ simulation studies. Results are reported via boxplots in Figures 18 to 20. While all methods improve their performance with larger sample sizes, we find that the BNP estimators (SSP, 3BP) provide relatively more accurate results for smaller samples sizes (e.g., $N = 10, N = 50$ in Figures 18 and 19). The performance of the BNP methods exceed those of competing methods for larger values of the exponent ($\xi \in \{1.2, 1.4, 1.6\}$), while higher order Jackknife and linear programs tend to do better for smaller values of the exponent ($\xi \in \{0.8, 1\}$).

# H  Additional experiments on the gnomAD dataset

## H.1  Experimental setup

In order to run our experiments, we use data from the gnomAD (genome aggregation dataset) discovery project [Karczewski et al., 2020], the largest and most comprehensive publicly available human genome dataset. We follow the same experimental setup adopted in Masoero et al. [2021]. We briefly summarize this setup in this section. The gnomAD dataset contains 125'748 exomes sequences (i.e. protein-coding regions of the genome), from 8 main populations. Sample size varies widely across sub populations, e.g. the "Other" subgroup counts about 3'000 observations, while "South Easy Asian" contains almost 16'000 individuals (see Karczewski et al. [2020] for additional details).

Figure 20: Accuracy of the competing methods (SB-SP-Bernoulli [SSP], stable beta-Bernoulli [3BP], Jackknife [J], linear program [LP], Good-Toulmin [GT]) on simulated data from a Zipf model (Equation (71)) with parameter $\xi \in \{1, 1.4, 1.6\}$ (top row, center row, bottom row). For $L = 2000$, we report $v_{N,a}^{(M)}$ as $N$ increases, for $M = L - N$. For each $N$, results across $S = 100$ datasets are reported in the boxplots.

For privacy reasons not all individual sequences are accessible. Hence, in order to run our analysis we generate synthetic data which closely resembles the true data as follows. For every subpopulation with $N$ individuals and every position $j = 1, \ldots, K$ in the exome, we have access to the total number of individuals $N_j$ showing variation at position $j$. We compute the empirical frequency of variation at site $j$, $\hat{\theta}_j := N_j/N$ for all $j = 1, \ldots, K$. Our data is then generated by sampling independent Bernoulli random vectors $X_1, \ldots, X_N$, with $X_n = [x_{n,1}, \ldots, x_{n,K}]$. The entries in the vector are independent Bernoulli random variables, $x_{n,j} \sim \text{Bernoulli}(\hat{\theta}_j)$.

## H.2 Results from the gnomAD data

For each of eight subpopulations in the data, we performed the following experiment. Let $\hat{\theta} = [\hat{\theta}_1, \ldots, \hat{\theta}_{K_{\max}}] \subseteq [0, 1]$ denote the "genetic signature" of the population, with $\hat{\theta}_k = N_k/N$, with $N_{tot}$ the total number of individuals in the population and $N_k$ the number of individuals in the population displaying such variant, $1 \leq N_k \leq N_{tot}$. Then, for each population, we generate $S = 50$ datasets by drawing $N_{tot}$ i.i.d. binary random vectors of length $K_{\max}$ as described above, with biases given by $\hat{\theta}$. We then retain for each dataset $N \in \{50, 100\}$ observations for training, and try to predict the number of new variants that are going to be observed if we were to sample additional $M = N_{tot} - N$ observations.

In a nutshell, also on this data, the findings are similar to the results obtained on the MSK-IMPACT cancer data. In particular, we find that when the sample size $N$ is small, the proposed SB-SP Bernoulli model leads to predictions that are often comparable or more accurate than competing methods.

First, we report the accuracy metric $v_N^{(M)}$ for eight subpopulations in gnomAD, Afroamerican (Amr.), South East Asian (SE. As.), Other East Asian (Ot. E. As.), Finnish (Fin.), South European (S. Eu.), Swedish (Swe.), South Asian (S. As.) and the remaining Other. In Figure 21 we show results (over $S = 50$ Monte-Carlo re-draws of the data from the estimated frequencies $\hat{\theta}$) of retaining $N = 50$ datapoints for training, and extrapolating to the largest available sample size $M$. In Figure 22 we report results for the same metric, with training performed by retaining $N = 100$ datapoints.
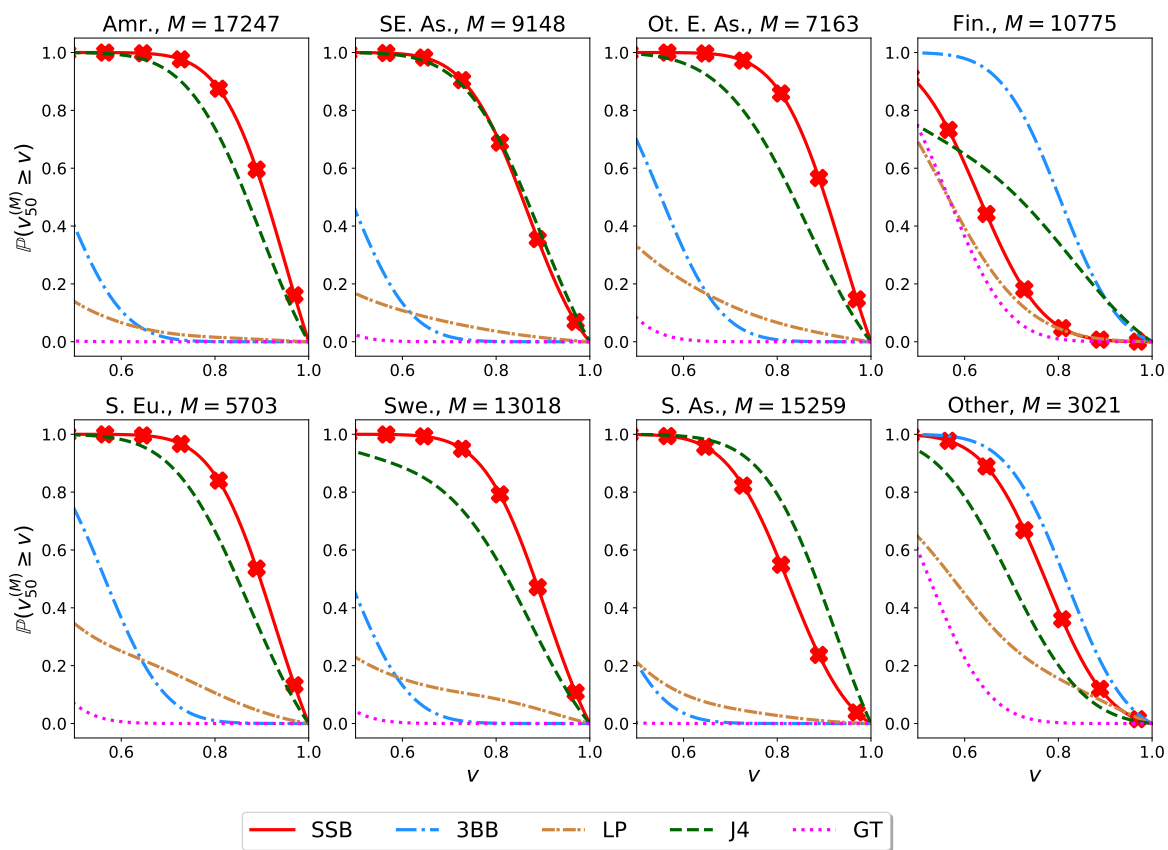
Figure 21: Accuracy metric $v_{50}^{(M)}$ for eight subpopulations in the gnomAD dataset. For each subpopulation we retain $N = 50$ observations for training, and extrapolate to the largest possible value $M$. Results are over $S = 50$ Monte-Carlo draws of the data, as described in Appendix H.1
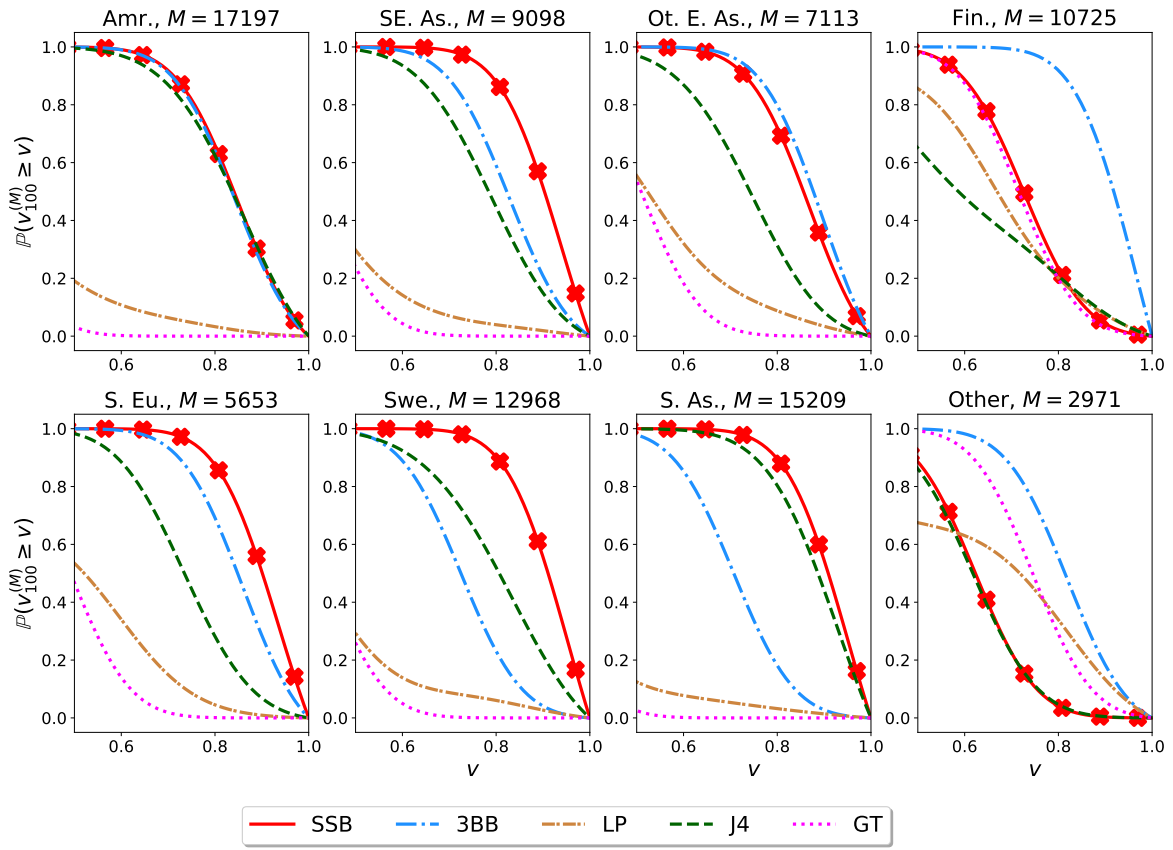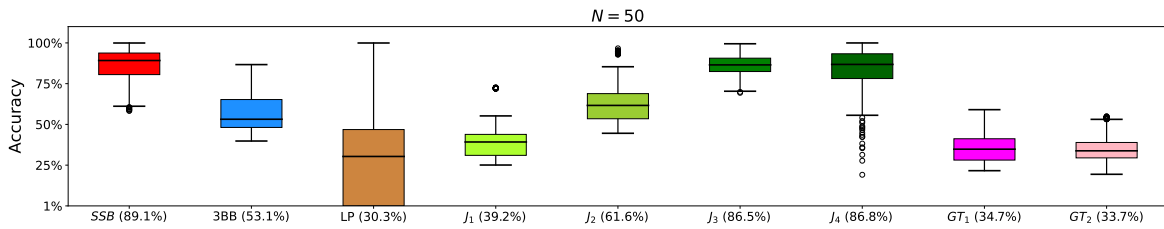
Figure 22: Same setup as in Figure 21, now for $N = 100$.

Figure 23: Accuracy of the compared methods, now over all the eight subpopulations and over 50 Monte Carlo draws for each population. $N = 50$, and $M$ is set to be the largest possible extrapolation size for each subpopulation.



Figure 24: Same setup as in Figure 23, now for $N = 100$.

Next, we provider boxplots that report the (aggregated) accuracy of the metric $v_N^{(M)}$ across all the eight populations, and all the $S = 50$ Monte-Carlo draws (so that each boxplot reports the accuracy of a total of $50 \times 8 = 400$ accuracy values), for $N = 50$ (Figure 23) as well as $N = 100$ (Figure 24).

Figure 25: Same setup as in Figure 23, but only for the American subpopulation.



Figure 26: Same setup as in Figure 23, but only for the Other East Asian subpopulation.
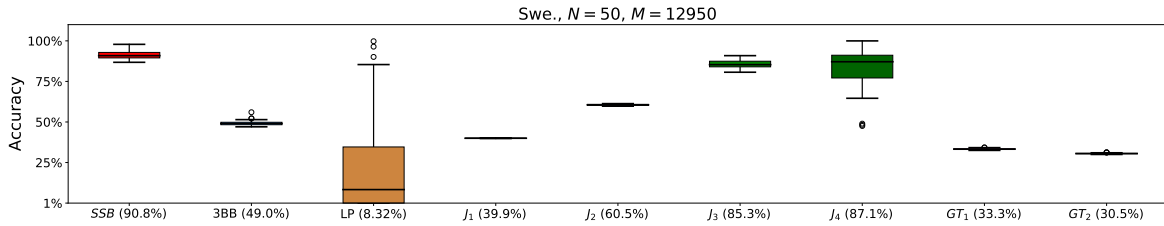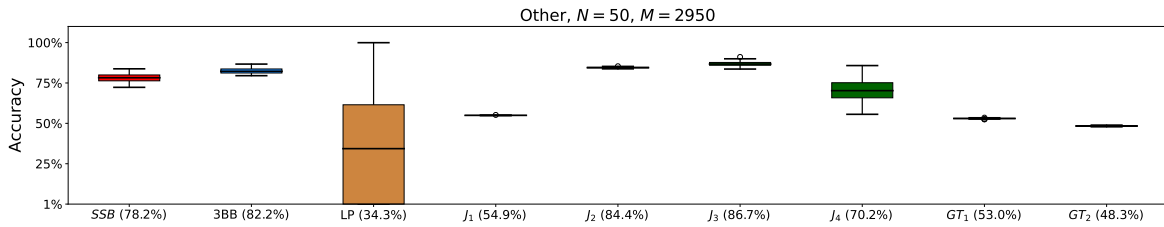
## H.3 Additional boxplots

Since in Figure 23 and Figure 24 we are aggregating result in which $N$ is consistent for all populations, but $M$ differs, we also report boxplots of each subpopulation individually.



Figure 27: Same setup as in Figure 23, but only for the East Asian subpopulation.

Figure 28: Same setup as in Figure 23, but only for the Finnish subpopulation.



Figure 29: Same setup as in Figure 23, but only for the Southern European subpopulation.



Figure 30: Same setup as in Figure 23, but only for the Swedish subpopulation.



Figure 31: Same setup as in Figure 23, but only for the "Other" subpopulation.
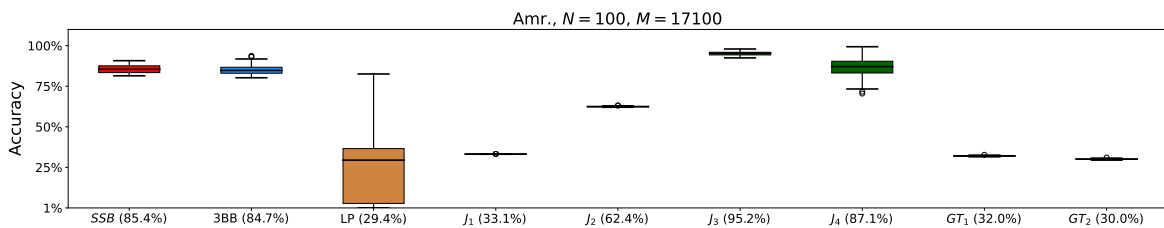


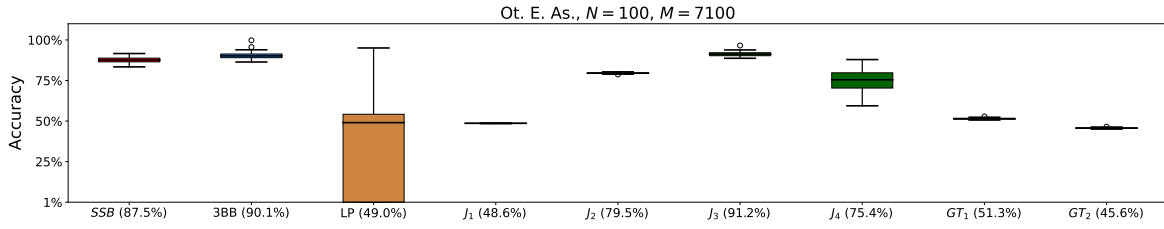Figure 32: Same setup as in Figure 24, but only for the American subpopulation.

**Figure 33:** Same setup as in Figure 24, but only for the Other East Asian subpopulation.
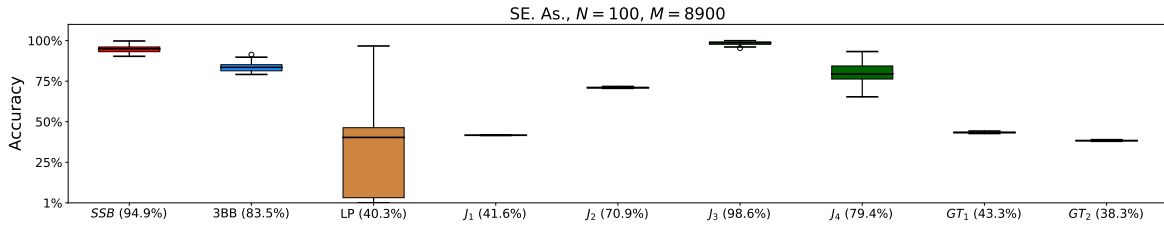


**Figure 34:** Same setup as in Figure 24, but only for the East Asian subpopulation.
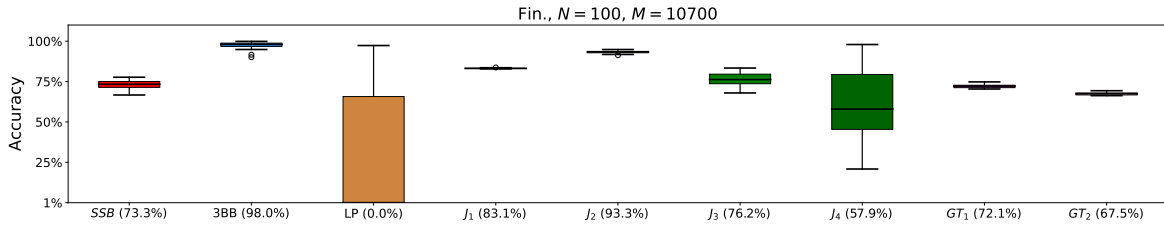


**Figure 35:** Same setup as in Figure 24, but only for the Finnish subpopulation.
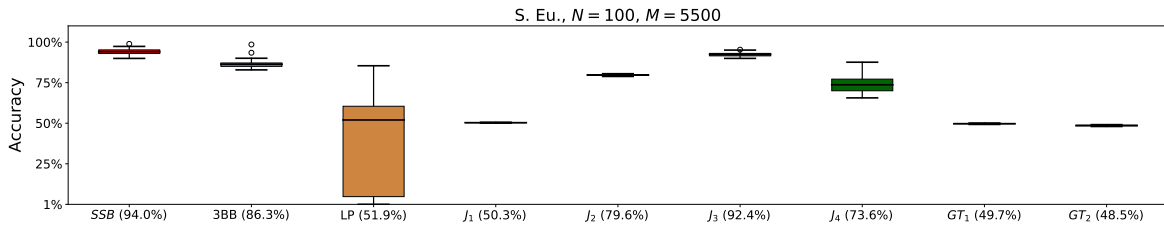


**Figure 36:** Same setup as in Figure 24, but only for the Southern European subpopulation.
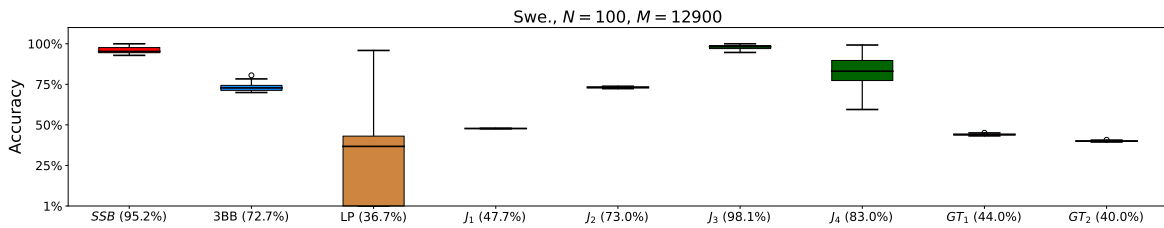


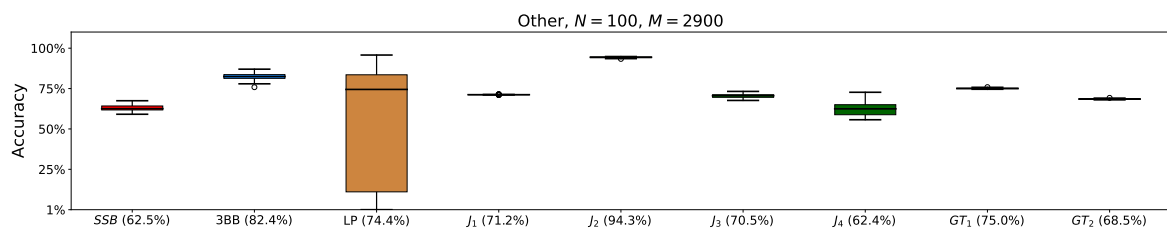**Figure 37:** Same setup as in Figure 24, but only for the Swedish subpopulation.

Figure 38: Same setup as in Figure 24, but only for the "Other" subpopulation.

# References

D. Aldous. *Exchangeability and related topics*, volume 1117 of *Lecture Notes in Mathematics.* Springer-Verlag, Berlin, 1983.

F. Ayed and F. Caron. Nonnegative Bayesian nonparametric factor models with completely random measures. *Stat. Comput.*, 31(5):Paper No. 63, 24, 2021. ISSN 0960-3174.

S. Bacallado, M. Battiston, S. Favaro, and L. Trippa. Sufficientness postulates for Gibbs-type priors and hierarchical generalizations. *Statist. Sci.*, 32:487–500, 2017.

T. Broderick, M. I. Jordan, and J. Pitman. Beta processes, stick-breaking and power laws. *Bayesian analysis*, 7:439–476, 2012.

T. Broderick, J. Pitman, and M. I. Jordan. Feature allocations, probability functions, and paintboxes. *Bayesian Analysis*, 8:801–836, 2013.

T. Broderick, A. C. Wilson, and M. I. Jordan. Posteriors, conjugacy, and exponential families for completely random measures. *Bernoulli*, 24:3181–3221, 2018.

K. P. Burnham and W. S. Overton. Estimation of the size of a closed population when capture probabilities vary among animals. *Biometrika*, 65:625–633, 1978.

S. Chakraborty, A. Arora, C. B. Begg, and R. Shen. Using somatic variant richness to mine signals from rare variants in the cancer genome. *Nature Communications*, 10:5506, 2019.

D. J. Daley and D. Vere-Jones. *An introduction to the theory of point processes. Vol. II.* Probability and its Applications (New York). Springer, New York, second edition, 2008. General theory and structure.

T. S. Ferguson. A Bayesian analysis of some nonparametric problems. *Ann. Statist.*, 1:209–230, 1973.

T. S. Ferguson and M. J. Klass. A representation of independent increment processes without Gaussian components. *Ann. Math. Statist.*, 43:1634–1643, 1972.

A. M. Goldstein, Y. Xiao, J. Sampson, B. Zhu, M. Rotunno, H. Bennett, Y. Wen, K. Jones, A. Vogt, and L. Burdette. Rare germline variants in known melanoma susceptibility genes in familial melanoma. *Human molecular genetics*, 26:4886–4895, 2017.

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

S. Gravel. Predicting discovery rates of genomic features. *Genetics*, 197:601–610, 2014.

T. L. Griffiths and Z. Ghahramani. The Indian buffet process: an introduction and review. *J. Mach. Learn. Res.*, 12:1185–1224, 2011.

R. D. Hernandez, L. H. Uricchio, K. Hartman, C. Ye, A. Dahl, and N. Zaitlen. Ultrarare variants drive substantial cis heritability of human gene expression. *Nature genetics*, 51:1349–1355, 2019.

J. R. Huyghe, S. A. Bien, T. A. Harrison, H. M. Kang, S. Chen, S. L. Schmit, D. V. Conti, C. Qu, J. Jeon, and C. K. Edlund. Discovery of common and rare genetic risk variants for colorectal cancer. *Nature Genetics*, 51:76–87, 2019.

I. Ionita-Laza and N. M. Laird. On the optimal design of genetic variant discovery studies. *Statistical Applications in Genetics and Molecular Biology*, 9, 2010.

I. Ionita-Laza, C. Lange, and N. M. Laird. Estimating the number of unseen variants in the human genome. *Proceedings of the National Academy of Sciences*, 106:5008–5013, 2009.

L. F. James. Bayesian Poisson calculus for latent feature modeling via generalized Indian buffet process priors. *Ann. Statist.*, 45:2016–2045, 2017.

L. F. James, P. Orbanz, and Y. W. Teh. Scaled subordinators and generalizations of the Indian buffet process. *arXiv preprint arXiv:1510.07309*, 2015.

Ö. Johansson, G. Samelius, E. Wikberg, G. Chapron, C. Mishra, and M. Low. Identification errors in camera-trap studies result in systematic population overestimation. *Scientific Reports*, 10:1–10, 2020.

O. Kallenberg. Commutativity properties of conditional distributions and Palm measures. *Commun. Stoch. Anal.*, 4:21–34, 2010.

O. Kallenberg. *Random measures, theory and applications*. Springer, Cham, 2017.

K. J. Karczewski, L. C. Francioli, G. Tiao, B. B. Cummings, J. Alföldi, Q. Wang, R. L. Collins, K. M. Laricchia, A. Ganna, and D. P. Birnbaum. The mutational constraint spectrum quantified from variation in 141,456 humans. *Nature*, 581(7809):434–443, 2020.

Y. Kim, L. James, and R. Weissbach. Bayesian analysis of multistate event history data: beta-dirichlet process prior. *Biometrika*, 99:127–140, 2012.

J. Kingman. Completely random measures. *Pacific Journal of Mathematics*, 21:59–78, 1967.

J. Kingman. *Poisson Processes*. Oxford Studies in Probability. Clarendon Press, 1992.

J. F. Kingman. Random discrete distributions. *J. R. Stat. Soc. Ser. B Stat. Methodol.*, 37:1–15, 1975.

D. Knowles and Z. Ghahramani. Nonparametric Bayesian sparse factor models with application to gene expression modeling. *Ann. Appl. Stat.*, 5(2B):1534–1552, 2011. ISSN 1932-6157.

K. Lawrenson, S. Kar, K. McCue, K. Kuchenbaeker, K. Michailidou, J. Tyrer, J. Beesley, S. J. Ramus, Q. Li, and M. K. Delgado. Functional mechanisms underlying pleiotropic risk alleles at the 19p13. 1 breast–ovarian cancer susceptibility locus. *Nature Communications*, 7:1–22, 2016.

A. Lee, N. Mavaddat, A. N. Wilcox, A. P. Cunningham, T. Carver, S. Hartley, C. B. de Villiers, A. Izquierdo, J. Simard, and M. K. Schmidt. Boadicea: a comprehensive breast cancer risk prediction model incorporating genetic and nongenetic risk factors. *Genetics in Medicine*, 21:1708–1718, 2019.

J. Lee, P. Müller, S. Sengupta, K. Gulukota, and Y. Ji. Bayesian inference for intratumour heterogeneity in mutations and copy number variation. *J. R. Stat. Soc. Ser. C. Appl. Stat.*, 65:547–563, 2016.

A. Lijoi and I. Prünster. Models beyond the Dirichlet process. In N. L. Hjort, C. Holmes, P. Müller, and S. Walker, editors, *Bayesian Nonparametrics*, pages 80–136. Cambridge University Press, 2010.

A. Lijoi, R. H. Mena, and I. Prünster. Bayesian nonparametric estimation of the probability of discovering new species. *Biometrika*, 94:769–786, 2007.

Y. Liu, J. Xia, J. McKay, S. Tsavachidis, X. Xiao, M. R. Spitz, C. Cheng, J. Byun, W. Hong, and Y. Li. Rare deleterious germline variants and risk of lung cancer. *NPJ precision oncology*, 5:1–12, 2021.

L. Masoero, F. Camerlenghi, S. Favaro, and T. Broderick. More for less: predicting and maximizing genomic variant discovery via Bayesian nonparametrics. *Biometrika*, 2021. doi: 10.1093/biomet/asab012.

Y. Momozawa and K. Mizukami. Unique roles of rare variants in the genetics of complex diseases in humans. *Journal of Human Genetics*, pages 1–13, 2020.

T. Nguyen-Dumont, R. J. MacInnis, J. A. Steen, D. Theys, H. Tsimiklis, F. Hammet, M. Mahmoodi, B. J. Pope, D. J. Park, and K. Mahmood. Rare germline genetic variants and risk of aggressive prostate cancer. *International journal of cancer*, 147:2142–2149, 2020.

A. Orlitsky, A. T. Suresh, and Y. Wu. Optimal prediction of the number of unseen species. *Proceedings of the National Academy of Sciences*, 113:13283–13288, 2016.

C. M. Phelan, K. B. Kuchenbaecker, J. P. Tyrer, S. P. Kar, K. Lawrenson, S. J. Winham, J. Dennis, A. Pirie, M. J. Riggan, and G. Chornokur. Identification of 12 new susceptibility loci for different histotypes of epithelial ovarian cancer. *Nature Genetics*, 49:680–691, 2017.

J. Pitman. Some developments of the Blackwell-MacQueen urn scheme. In *Statistics, probability and game theory*, volume 30 of *IMS Lecture Notes Monogr. Ser.*, pages 245–267. Inst. Math. Statist., Hayward, CA, 1996.

J. Pitman. *Combinatorial stochastic processes*, volume 1875 of *Lecture Notes in Mathematics*. Springer-Verlag, Berlin, 2006.

J. Pitman and M. Yor. The two-parameter Poisson-Dirichlet distribution derived from a stable subordinator. *Ann. Probab.*, 25:855–900, 1997.

S. Rashkin, G. Jun, S. Chen, G. R. Abecasis, Genetics, and E. of Colorectal Cancer Consortium. Optimal sequencing strategies for identifying disease-associated singletons. *PLoS Genetics*, 13, 2017.

R. Rasnic, N. Linial, and M. Linial. Expanding cancer predisposition genes with ultra-rare cancer-exclusive human variations. *Scientific reports*, 10:1–9, 2020.

E. Regazzini. Intorno ad alcune questioni relative alla definizione del premio secondo la teoria della credibilià. *Giornale dell'Istituto Italiano degli Attuari*, 41:77–89, 1978.

J. G. Sanders, S. Nurk, R. A. Salido, J. Minich, Z. Z. Xu, Q. Zhu, C. Martino, M. Fedarko, T. D. Arthur, and F. Chen. Optimizing sequencing protocols for leaderboard metagenomics by combining long and short reads. *Genome Biology*, 20:1–14, 2019.

K. Schwarze, J. Buchanan, J. M. Fermont, H. Dreau, M. W. Tilley, J. M. Taylor, P. Antoniou, S. J. Knight, C. Camps, and M. M. Pentony. The complete costs of genome sequencing: a microcosting study in cancer and rare diseases from a single center in the united kingdom. *Genetics in Medicine*, 22:85–94, 2020.

C. A. Souza, N. Murphy, C. Villacorta-Rath, L. N. Woodings, I. Ilyushkina, C. E. Hernandez, B. S. Green, J. J. Bell, and J. M. Strugnell. Efficiency of ddRAD target enriched sequencing across spiny rock lobster species (palinuridae: Jasus). *Scientific reports*, 7:1–14, 2017.

Y. Teh and D. Gorur. Indian buffet processes with power-law behavior. *Advances in neural information processing systems*, 22:1838–1846, 2009.

C. Wendt and S. Margolin. Identifying breast cancer susceptibility genes–a review of the genetic background in familial breast cancer. *Acta Oncologica*, 58:135–146, 2019.

S. Zabell. *The continuum of inductive methods revisited*. Cambridge University Press, 2005.

M. J. Zhang, V. Ntranos, and D. Tse. Determining sequencing depth in a single-cell RNA-seq experiment. *Nature Communications*, 11:1–11, 2020.

J. Zou, G. Valiant, P. Valiant, K. Karczewski, S. O. Chan, K. Samocha, M. Lek, S. Sunyaev, M. Daly, and D. G. MacArthur. Quantifying unobserved protein-coding variants in human populations provides a roadmap for large-scale sequencing projects. *Nature Communications*, 7:13293, 2016.