# Prior selection for the precision parameter of Dirichlet Process Mixtures

C. Vicentini[a], I. H. Jermyn[a]

[a]*Department of Mathematical Sciences, Durham University, Upper Mountjoy, Stockton Road, DH1 3LE, Durham, UK*

**Abstract**

Consider a Dirichlet process mixture model (DPM) with random precision parameter $\alpha$, inducing $K_n$ clusters over $n$ observations through its latent random partition. Our goal is to specify the prior distribution $p(\alpha \mid \boldsymbol{\eta})$, including its fixed parameter vector $\boldsymbol{\eta}$, in a way that is meaningful.

Existing approaches can be broadly categorised into three groups. Those in the first group depend on the sample size $n$, and often rely on the linkage between $p(\alpha \mid \boldsymbol{\eta})$ and $p(K_n)$ to draw conclusions on how to best choose $\boldsymbol{\eta}$ to reflect one's prior knowledge of $K_n$; we call them *sample-size-dependent*. Those in the second and third group consist instead of using quasi-degenerate or improper priors, respectively.

In this article, we show how all three methods have limitations, especially for large $n$. Then we propose an alternative methodology which does not depend on $K_n$ or on the size of the available sample, but rather on the relationship between the largest stick lengths in the stick-breaking construction of the DPM; and which reflects those prior beliefs in $p(\alpha \mid \boldsymbol{\eta})$. We conclude with an example where existing sample-size-dependent approaches fail, while our sample-size-independent approach continues to be feasible.

*Keywords:* bayesian nonparametrics, Dirichlet process, precision parameter

## 1. Introduction

Typical usages of the Dirichlet process mixture model (DPM) are density and cluster estimation; the former is motivated by the flexibility of the DPM, and the latter by the latent random partition that the model induces on the observed data. In both cases, the precision parameter $\alpha$ of the DPM is of great significance, since it influences the smoothness of the resulting density, as well as $K_n$, the random number of clusters in the underlying partition of $n$ data points (Dorazio, 2009; Murugiah and Sweeting, 2012).

Methods have been developed to infer $\alpha$ from the data. For example, point estimates can be obtained with empirical Bayes, as outlined in Liu (1996), or alternatively a fully Bayesian approach can be used to obtain the full posterior distribution. This paper

---

*Corresponding author

Email addresses:* `carlo.vicentini@mathstat.net` (C. Vicentini), `i.h.jermyn@durham.ac.uk` (I. H. Jermyn)

focuses on the latter, which often involves Markov Chain Monte Carlo (MCMC); we address the question of how best to choose the parameter vector $\boldsymbol{\eta}$ of the prior distribution $p(\alpha \mid \boldsymbol{\eta})$, for some choice of parameterized model, and in particular what quantity of interest one's prior belief should be anchored to when eliciting $\boldsymbol{\eta}$.

Existing approaches in this domain can be broadly categorised into three groups. Those in the first group depend on the sample size $n$, and often rely on the linkage between $p(\alpha \mid \boldsymbol{\eta})$ and $p(K_n)$ to draw conclusions on how best to choose $\boldsymbol{\eta}$ to reflect one's prior knowledge of $K_n$; we call them *sample-size-dependent*. Those in the second group consist of quasi-degenerate priors, such as for example a Gamma $(a, b)$ with $a$ and $b$ both close to zero. Those in the third group consist of improper priors.

Dependence on the sample size is typically undesirable because the resulting modelling construct, including the data generating process, becomes applicable only for that particular sample size. For example, a common trait of most sample-size-dependent priors is that they try to induce a diffuse prior on $K_n$ (West and Escobar, 1993; Dorazio, 2009; Murugiah and Sweeting, 2012). However, a prior deemed diffuse or weakly informative on $K_n$ for a certain sample size $n$ would not necessarily be so for a different value of $n$. In fact, articulating one's prior beliefs through $K_n$ leads to a sequence of priors indexed by the sample size, $\{p(\alpha \mid \boldsymbol{\eta}_n)\}_{n=1}^{\infty}$. This is particularly unappealing because the DPM model is otherwise well suited for inference on streaming data, as the number of clusters that it induces is not fixed, but grows with the data size. Furthermore, by trying to be diffuse in $K_n$, such priors actually influence other important quantities of interest, in a way that is material. In our view, quasi-degenerate and improper priors do not offer viable solutions either, as the former are not compatible with the notion of multiple clusters, hence they defeat the purpose of using a DPM in the first place, while the latter can lead to an improper posterior.

We introduce a new approach to the specification of $p(\alpha)$, which is independent of the sample size and which is instead based on the appraisal of the implied joint distribution of the stick-breaking weights (either in size-biased order, or ranked); we show an example in which multiple DPMs stem from a common prior $p(\alpha; \boldsymbol{\eta})$, and where, as a result, sample-size-dependent approaches are inapplicable while our sample-size-independent method is feasible. Our approach is essentially based on the weights, or lengths, of the sticks that are broken off the initial stick of length 1, in the stick-breaking representation of the Dirichlet process; these can also be equivalently seen as the asymptotic relative cluster sizes for $n \to \infty$.

In the next sections, we proceed as follows. In section 2, we summarise some basic notions about the Dirichlet process. In section 3, we discuss the existing literature on $\alpha$ priors, and their limitations. In section 4, we introduce a novel approach to the selection of $p(\alpha \mid \boldsymbol{\eta})$, which does not depend on the sample size. In section 5, we cross-examine how sample-size-dependent and sample-size-independent approaches perform in relation to each other. In section 6, we showcase an example where existing sample-size-dependent approaches are inapplicable, while our sample-size-independent approach continues to be feasible. Section 7 outlines our conclusions.

## 2. Representations and properties of the Dirichlet process

Consider a Dirichlet process $G \sim \mathrm{DP}(\alpha, G_0)$ with precision parameter $\alpha$ and base measure $G_0$. As proved in Ferguson (1973), its posterior distribution given $n$ observations

$(\theta_1, \ldots, \theta_n)$ is in turn a DP $\left( \alpha + n, \frac{\alpha G_0 + \sum_{i=1}^n \delta_{\theta_i}}{\alpha + n} \right)$ meaning that, unconditional on $G$, we have:

$$p\left(\theta_{n+1} \mid \theta_1, \ldots, \theta_n\right) = \frac{\alpha}{\alpha + n} G_0\left(\theta_{n+1}\right) + \sum_{i=1}^n \frac{1}{\alpha + n} \delta_{\theta_i}\left(\theta_{n+1}\right). \tag{1}$$

Mixing the DP with the parametric likelihood $p\left(\cdot \mid \theta\right)$ leads to a Dirichlet process mixture (DPM) (Lo, 1984; Ferguson, 1983), which can be written as

$$\begin{aligned} y_i \mid \theta_i &\sim p\left(y_i \mid \theta_i\right), \\ \theta_i \mid G &\sim G, \\ G &\sim \mathrm{DP}\left(\alpha, G_0\right), \end{aligned} \tag{2}$$

where $\boldsymbol{y} = (y_1, \ldots, y_n)$ is a vector of $n$ observations.

We discuss two constructions of the Dirichlet process that are relevant to this article: the stick-breaking and the Poisson-Dirichlet process representations.

### 2.1. The stick-breaking representation

Sethuraman (1994) considered

$$\begin{aligned} w_1 &:= v_1, \\ w_h &:= v_h \prod_{l<h} (1 - v_l), \quad h = 2, 3, \ldots \\ v_h &\sim \mathrm{Beta}\left(1, \alpha\right), \quad h = 1, 2, \ldots, \\ m_h &\sim G_0, \quad h = 1, 2, \ldots, \end{aligned} \tag{3}$$

where $G_0$ is a probability measure, and proved that the distribution of the random measure

$$G := \sum_{h=1}^\infty w_h \delta_{m_h} \tag{4}$$

is $DP\left(\alpha, G_0\right)$.

The distribution of $\boldsymbol{w} := (w_1, w_2, \ldots)$ alone is known as the Griffiths-Engen-McCloskey distribution, and it is denoted by $\mathrm{GEM}\left(\alpha\right)$ (Arratia et al., 2003, section 4.8); its elements $w_1, w_2, \ldots$ are in size-biased order (Pitman, 1996). The probability distribution of the first $H - 1$ elements of $\boldsymbol{w}$ can be derived from the $H$-dimensional generalised Dirichlet distribution (Connor and Mosimann, 1969).

### 2.2. The Poisson-Dirichlet Process representation

Kingman (1975) introduced the Poisson-Dirichlet distribution (PDD), which he constructed through the gamma process $\xi\left(t\right)$, a stochastic process with $\xi\left(0\right) = 0$ and with increments which are independent on disjoint intervals, and gamma distributed. The PDD with parameter $\alpha$ is known to be equivalent to the distribution of the decreasing order statistics of $\boldsymbol{w}$ (Arratia et al., 2003, section 4.11)(Pitman, 1996); we denote them by $\left(w_1^\downarrow, w_2^\downarrow, \ldots\right)$. Similarly, the decreasing weight-ranked equivalent of equation 4 is known as the Poisson-Dirichlet Process (PDP).

As laid out in Watterson (1976), the conditional joint probability distribution of the first $r$ ranked weights is:

$$p\left(w_1^\downarrow, \ldots, w_r^\downarrow\right) = \alpha^r \Gamma\left(\alpha\right) e^{\gamma\alpha} \frac{w_r^{\downarrow\alpha-1}}{w_1^\downarrow \cdots w_r^\downarrow} \; g\left(\frac{1 - w_1^\downarrow - \ldots - w_r^\downarrow}{w_r^\downarrow}\right), \tag{5}$$

where $g$ is a function which can be recursively written as:

$$g\left(z\right) = z^{\alpha-1}\left[g\left(n\right) n^{1-\alpha} - \alpha \int_n^z g\left(y-1\right) y^{-\alpha} dy\right], \quad n \leq z < n+1, \; n := \lfloor z \rfloor,$$

and which is known to be particularly difficult to compute.

### 2.3. Properties of the Dirichlet process

Repeated values in $\boldsymbol{\theta}$ induce implicit clustering on $\boldsymbol{\theta}$ and, in turn, a random partition model. This can be observed in equation 1, where if $G_0$ is continuous, we obtain

$$p\left(\theta_i \notin \{\theta_1, \ldots, \theta_{i-1}\} \mid \alpha\right) = \frac{\alpha}{\alpha + i - 1}, \; i = 2, \ldots, n,$$

hence the random number of clusters $K_n \mid \alpha$ observed in a sample of $n$ observations is distributed as the sum of $n$ Bernoulli variables with parameters $p_i = \alpha/\left(\alpha + i - 1\right)$. Further,

$$\mathbb{E}\left[K_n \mid \alpha\right] = \sum_{i=1}^n \frac{\alpha}{\alpha + i - 1} \sim \alpha \log \frac{n + \alpha}{\alpha} \sim \alpha \log n, \tag{6}$$

$$\mathbb{V}\mathrm{ar}\left[K_n \mid \alpha\right] = \sum_{i=1}^n \frac{\alpha\left(i-1\right)}{\left(\alpha + i - 1\right)^2} \sim \alpha \log n, \tag{7}$$

as $n \to \infty$ (Antoniak, 1974; Arratia et al., 2003). The probability distribution of $K_n \mid \alpha$ is:

$$p\left(K_n = k \mid \alpha\right) = s_{n,k} \; \alpha^k \frac{\Gamma\left(\alpha\right)}{\Gamma\left(\alpha + n\right)}, \tag{8}$$

where $s_{n,k}$ are unsigned Stirling numbers of the first kind. The parameter $\alpha$ therefore critically controls $p\left(K_n \mid \alpha\right)$: we show in sections 3.2 and 3.3 that, for $\alpha \to 0$, the random variable $K_n \mid \alpha$ converges to the Dirac measure $\delta_1$, while for $\alpha \to \infty$, the same converges to $\delta_n$; similarly, if $\alpha$ is random, $p\left(\alpha \mid \boldsymbol{\eta}\right)$ determines the prior $p\left(K_n\right)$ through mixing, as

$$p\left(K_n = k\right) = \int_0^\infty s_{n,k} \; \alpha^k \frac{\Gamma\left(\alpha\right)}{\Gamma\left(\alpha + n\right)} \mathrm{d}p\left(\alpha \mid \boldsymbol{\eta}\right). \tag{9}$$

The following also holds (Watterson, 1974; Arratia et al., 2000):

$$d_{TV}\left(p\left(K_n \mid \alpha\right), \mathrm{Po}\left(\mathbb{E}\left[K_n \mid \alpha\right]\right)\right) = O\left(\frac{1}{\log n}\right), \tag{10}$$

where $d_{TV}$ indicates total variation distance and $\mathrm{Po}\left(\lambda\right)$ indicates the Poisson distribution with parameter $\lambda$, and which implies convergence in distribution since $1/\log n \to 0$ for $n \to \infty$. While the relationship above is unlikely to carry practical use in computations, as its rate of convergence is sublinear, it does provide conceptual insight into the limiting behaviour of $p\left(K_n \mid \alpha\right)$, as we show in section 3.1.

## 3. Existing priors for $\alpha$

In this section we discuss three approaches to prior specification. Methods in the first group (see section 3.1) are *sample-size-dependent*, and we refer to them as 'SSD'. Those in the other two groups are quasi-degenerate and improper priors (see sections 3.2 and 3.3).

Although in principle any distributional choice of $\alpha$ is allowed, all aforementioned approaches have historically been discussed by their authors in the context of the gamma distribution – i.e. $\alpha \sim \mathrm{Ga}(a,b)$. The popularity of the gamma distribution as a prior for $\alpha$ is due to reasons of computational attractiveness of the posterior, and is to be traced back to Escobar and West (1995).

### 3.1. Sample-size-dependent approaches (SSD)

Methods in this group are the $K_n$-diffuse prior of West and Escobar (1993), the DORO prior of Dorazio (2009), the SCAL prior of Murugiah and Sweeting (2012), and Jeffreys' prior (Rodríguez, 2013). The first three leverage one's prior assumptions on $K_n$ as a target to determine $\boldsymbol{\eta}$ in $p(\alpha \mid \boldsymbol{\eta})$, while Jeffreys' prior purely depends on $n$ and does not involve one's assumptions on $K_n$. All four depend on the size of the sample, $n$.

### 3.1.1. $K_n$-diffuse prior

West and Escobar (1993) use a gamma hyperprior, $\alpha \sim \mathrm{Ga}(a,b)$, "supporting a diffuse range of reasonably large values consistent with possibly large values" of $K_n$. In their article, they use $n = 74$; their prior supports "a wide range of $k$ values between about $k = 8$ and $k = 35$"[1]. We call their approach $K_n$-diffuse, to highlight that it is not necessarily diffuse in $\alpha$, but rather it is diffuse in $K_n$.

We observe that this approach, which is appealing in its simplicity, has a dependency on the size of the sample it is applied to. For example, Figure 1 shows the impact on $p(K_n)$ of a $\mathrm{Ga}(10,1)$ prior, as $n$ moves from $n = 10$ to $n = 100$:

- in the left panel ($n = 10$), $K_n$ is centred on values that are large relative to $n$, with a wide spread relative to the support of $K_n$, hence it is well-diffused;

- in the right panel ($n = 100$), $K_n$ is centred on smaller values (relative to $n$), with a smaller spread over the support of $K_n$ – meaning that it is not as well-diffused as when $n = 10$.

As a result, we conclude that if a $K_n$-diffuse prior is defined as one whose mass is well-spread around central values of $K_n$, the consequence is that, under a $\mathrm{Ga}(a,b)$ prior, $(a,b)$ needs to be updated as $n$ grows, to ensure that central values of $K_n$ continue to be well covered, and that the relative spread around central values is preserved.

---

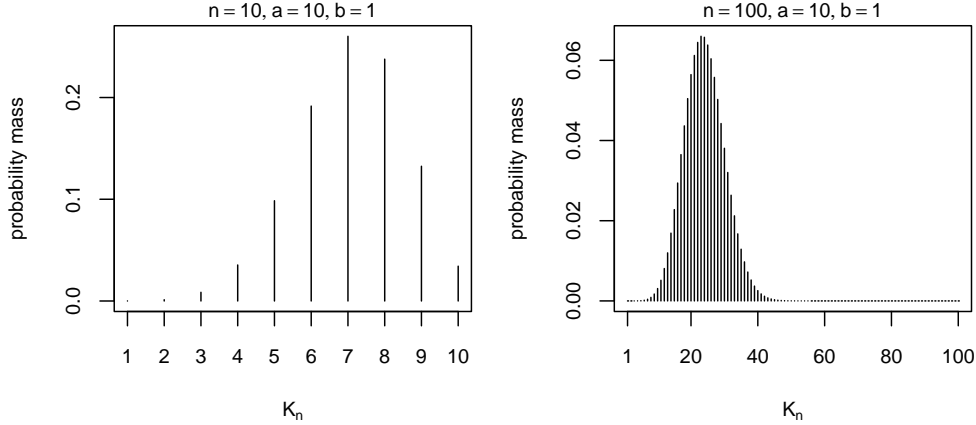[1]West and Escobar (1993) use $a = 5$ and $b = 0.5$, and $k$ in their notation is equivalent to $K_n$ in ours.

Figure 1: $K_n$-diffuse prior. While the prior probability distribution induced on $K_n$ by $\alpha \sim \mathrm{Ga}\,(10,1)$ appears reasonably diffuse for $n = 10$, it is less so for $n = 100$, as the shape and skewness of $p\,(K_n)$ change with $n$.

### 3.1.2. DORO priors

Dorazio (2009) proposes $\alpha \sim \mathrm{Ga}\,(a,b)$, with $(a,b)$ set to minimise the Kullback-Leibler distance between the prior probability distribution $p\,(K_n)$ induced by $p\,(\alpha \mid \boldsymbol{\eta})$ on $K_n$, and a target discrete distribution; in the absence of prior information about $K_n$, Dorazio (2009) uses the discrete uniform as a target. The DORO approach results in a pre-determined list of optimal values of $(a,b)$ for various values of $n$ (Table 1).

However, Figure 2 shows that the approximation of the discrete uniform resulting from the DORO prior is visibly coarse, and that it does not appear to improve as $n$ grows. In fact, asymptotic results show that, when $\alpha \sim \mathrm{Ga}\,(a,b)$, $K_n$ converges sub-linearly to a negative binomial random variable, meaning that a discrete uniform is unachievable for $n \to \infty$: from equation 10,

$$p\,(K_n \mid \alpha) \xrightarrow{d} \mathrm{Poi}\,(\alpha \log n)\,, \quad n \to \infty,$$

and $\alpha \log n \sim \mathrm{Ga}\,(a, b/\log n)$ hence $K_n \sim \mathrm{NB}\,(a, b/\,(b + \log n))$ in the limit.

Furthermore, as we will see in section 3.1.4, the shape of the prior induced by DORO on $K_n$ is also quite different from the one that is attained under Jeffreys' prior, which is one more reason why the choice by DORO of targeting a discrete uniform prior distribution on $K_n$ can be debated.

### 3.1.3. SCAL priors

Murugiah and Sweeting (2012) propose a scaling approach, where the values of $(a,b)$ are initially computed for a given $n$ according to how well they perform at recovering some known cluster structure; $(a,b)$ are subsequently rescaled to other choices of $n$, without the need to re-compute them again through the fully-fledged determination process.

Upon the initial determination of $(a,b)$, the goal of SCAL is to scale $(a,b)$ in such a way that the prior mean of $K_n$ is only affected to a small extent by the changes, while the variance is influenced to a larger extent. In particular, SCAL is based on fixing $p\,(K_n = 1)$ and $p\,(K_n \in \{c, c+1, \ldots, n-1, n\})$ to some determined values, for a
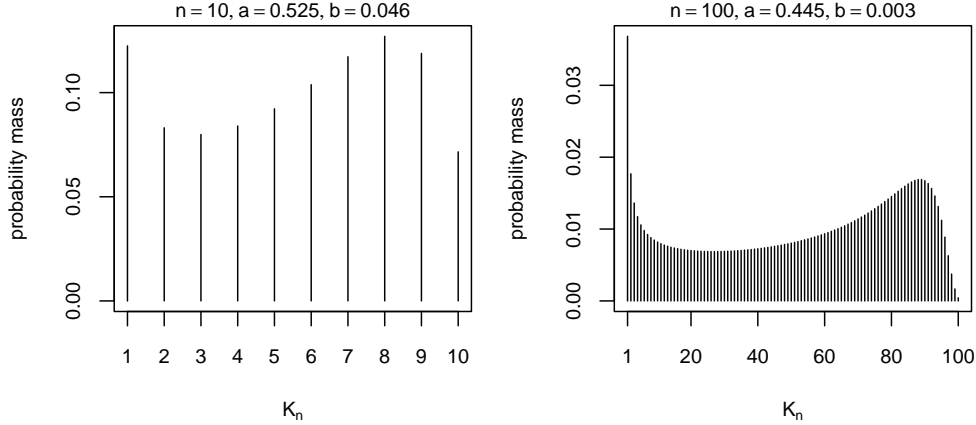
6

Figure 2: DORO prior. Prior probability distribution induced on $K_n$ by $\alpha \sim \text{Ga}(a, b)$. $K_n$ does not appear to be close to the target discrete uniform distribution, and the approximation does not appear to improve as $n$ increases.

suitably selected $c$, and deriving $(a, b)$ implicitly as $n$ varies, with equation 9. Murugiah and Sweeting (2012) suggest $c = \lceil c_0 \log n \rceil$, with $c_0 = 2$.

Murugiah and Sweeting (2012) use simulated data sets of $n = 6$ observations to elicit $(a, b) = (1, 1)$; they then test their scaling approach up to $n = 25$, by re-calculating $(a, b)$ while keeping $p(K_n = 1) = 0.34$ and $p(K_n \in \{c, \ldots, n\}) = 0.15$. They fit a curve through the results that they obtain, to ultimately propose the following equation, as an easier approximation[2] to SCAL:

$$a = b = e^{-0.033n}. \tag{11}$$

Although Murugiah and Sweeting (2012) only obtained the optimal $(a, b)$ values for their test case for $n \in \{6, 10, 15, 20, 25\}$, we extend their results to $n \in \{50, 75, 100\}$ (see table 2). As expected, we observe dependence on $n$, albeit to a lesser extent than DORO. We also note that the proposed approximation significantly diverges from the exact values for $n > 25$; in fact, instead of the exact values $(a, b) = (0.403, 0.370)$, for $n = 100$ it yields $a = b = 0.037$, a quasi-degenerate gamma prior.

*3.1.4. Jeffreys' prior*

Jeffreys' prior was introduced by Jeffreys (1946) as the volume measure of the Riemannian metric induced on the parameter space from a Riemannian metric on the space of probability distributions containing the parameterized model, thereby ensuring that the prior is determined by the associated probability laws and not by an arbitrary labelling of those laws by parameter values. This would be true for any Riemannian metric on the space of distributions; what makes Jeffreys' prior special is that it is induced by the unique (Čencov, 1982; Campbell, 1986; Bauer et al., 2016) Riemannian metric that

---

[2]This is justified in Murugiah and Sweeting (2012) by observing that $\mathbb{E}[\alpha] = 1$, which is fixed, and that $\mathbb{V}[\alpha] = e^{0.0165n}$, which increases with $n$, as their approach originally intended.

| $n$ | $a$ | $b$ | $D_{KL}$ |
|---|---|---|---|
| 5 | 0.541 | 0.096 | 0.00458 |
| 10 | 0.525 | 0.046 | 0.01904 |
| 15 | 0.512 | 0.029 | 0.03048 |
| 20 | 0.501 | 0.021 | 0.03942 |
| 25 | 0.490 | 0.015 | 0.04660 |
| 30 | 0.486 | 0.013 | 0.05272 |
| 35 | 0.480 | 0.010 | 0.05806 |
| 40 | 0.475 | 0.009 | 0.06265 |
| 45 | 0.470 | 0.008 | 0.06684 |
| 50 | 0.467 | 0.007 | 0.07050 |
| 100 | 0.445 | 0.003 | 0.09529 |

Table 1: Optimal values of $a, b$ under the DORO approach, when $\alpha \sim \mathrm{Ga}(a, b)$ and when the target distribution $p(K_n)$ is the discrete uniform. We have enriched the original table from Dorazio (2009) with an additional entry for $n = 100$.

| $n$ | exact $(a, b)$ | approx. $(a, b)$ |
|---|---|---|
| 25 | $(0.490, 0.438)$ | $(0.438, 0.438)$ |
| 50 | $(0.466, 0.467)$ | $(0.192, 0.192)$ |
| 75 | $(0.432, 0.420)$ | $(0.084, 0.084)$ |
| 100 | $(0.403, 0.370)$ | $(0.037, 0.037)$ |

Table 2: Optimal and approximate values of $(a, b)$ under the SCAL approach, as we determined them to be according to the approach outlined in Murugiah and Sweeting (2012). Approximate values originate from equation 11. Targets are $p(K_n = 1) = 0.34$ and $p(K_n \in \{c, \ldots, n\}) = 0.15$, $c = c_0 \log n$, $c_0 = 2$.

is invariant to (appropriately behaved) mappings of the space upon which the distribution is defined. Subsequently, it was found to have other desirable properties too, and it is now widely used in statistics and best known as the most prominent example of an *objective* prior. It was first derived for the multivariate Ewens distribution by Rodríguez (2013), who then tested it on the related Dirichlet process mixture model.

For univariate cases like the one being discussed, it is obtained as:

$$p(\alpha) \propto \sqrt{I_\alpha(\alpha)} = \sqrt{\mathbb{E}_k\left[\left(\frac{\partial}{\partial \alpha} \log p(k \mid \alpha)\right)^2\right]} = \sqrt{\frac{1}{\alpha} \sum_{i=1}^{n} \frac{i-1}{(\alpha + i - 1)^2}},$$

where $I_\alpha(\alpha)$ is the Fisher information (see Rodríguez (2013)). We found the following equivalent formulation to be computationally faster:

$$p(\alpha) \propto \sqrt{\frac{1}{\alpha}\left[\psi_0(\alpha + n) - \psi_0(\alpha + 1) + \alpha(\psi_1(\alpha + n) - \psi_1(\alpha + 1))\right]},$$

which is obtained because, by definition, $\psi_0(z + 1) = \psi_0(z) + \frac{1}{z}$, and also $\psi_1(\alpha + 1) = \psi_1(\alpha) - \frac{1}{\alpha^2}$.

Its density is 0 on the entire half-line for $n = 1$, as is its integral, hence this prior carries no meaning for $n = 1$. This is natural, since when $n = 1$, $k = 1$ with certainty, and there is no dependence on $\alpha$.

For $n = 2$, Jeffreys' prior can be expressed analytically, including its normalising constant, leading to the density and cumulative probability distribution functions:

$$p(\alpha) = \frac{1}{\pi(\alpha+1)\sqrt{\alpha}}, \tag{12}$$

$$p(\alpha \leq x) = \frac{2}{\pi}\mathrm{atan}\sqrt{\alpha}. \tag{13}$$

This prior is proper and has no finite moments (Rodríguez, 2013). Figure 3 exemplifies the change in shape of $p(\alpha)$ as $n$ increases.
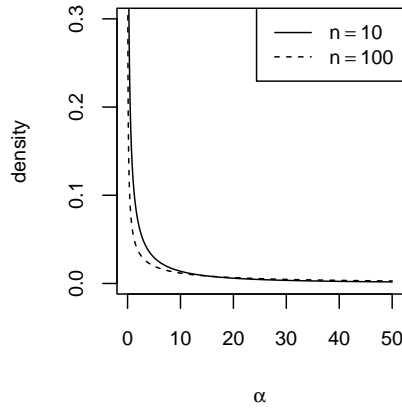


Figure 3: Density of Jeffreys' prior for $n = 10$ and $n = 100$.

Two examples of the shape of the prior distribution induced on $K_n$ by assigning Jeffreys' prior to $\alpha$ are plotted in figure 4. As noted in Rodríguez (2013), this distribution seems to be approximately symmetric around $n/2$. We include in Appendix C and Appendix D results showing that the posterior distribution induced by Jeffreys' prior on $\alpha$ is proper, as is the prior distribution it induces on $K_n$.

In their paper, after testing it on a species sampling model, Rodríguez (2013) sample 50 data points from a negative binomial $\mathrm{NB}(20, 220)$, and fit a DPM with a Gamma base measure (with mean and variance to mirror that of the negative binomial) under two priors for $\alpha$: the quasi-degenerate $\mathrm{Ga}(0.001, 0.001)$ prior, and Jeffreys' prior. They find that posterior inference on $K_n$ with the DPM leads to very similar results under these two options, with Jeffreys' prior favoring "a somewhat larger number of clusters than the Gamma prior" (Rodríguez, 2013). They also stress that Jeffreys' prior "explicitly depends on the sample size".

### 3.2. Quasi-degenerate priors

Similarly to the SCAL approximation from equation 11, other authors have independently suggested, on isolated occasions rather than as part of an attempt to design a prior elicitation framework, to use a quasi-degenerate $\mathrm{Ga}(a, b)$, with $a$ and $b$ close to zero, on the basis that it approximates the improper prior $p(\alpha) \propto 1/\alpha$ that is uniform in
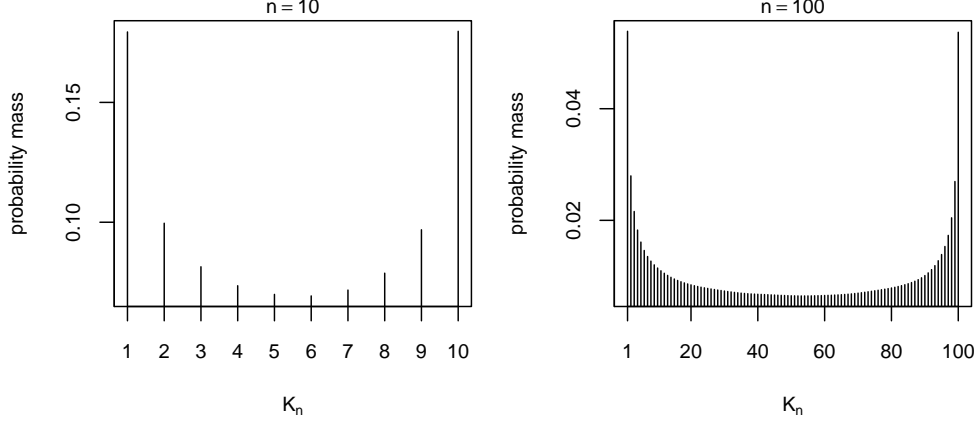
Figure 4: Prior distribution $p(K_n)$ induced by assigning Jeffreys' prior to $p(\alpha)$, for $n = 10$ and $n = 100$.

$\log(\alpha)$ (Escobar and West, 1995, 1998; Navarro et al., 2006; Lunn et al., 2012). However, gamma priors with very small $(a, b)$ are undesirable, because

$$\lim_{(a,b)\to(0,0)} p(K_n = 1) = 1,$$

hence they lead to the prior expectation of a parametric model with only one mixture component, which negates the reason for using a nonparametric prior in the first place and which is likely to overwhelm the data due to its strength; Dorazio (2009) and Murugiah and Sweeting (2012) also draw similar conclusions.

This is proven as follows. We are interested in

$$\lim_{(a,b)\to(0,0)} p(K_n = k) = s_{n,k} \lim_{(a,b)\to(0,0)} \int_0^\infty \alpha^k \frac{\Gamma(\alpha)}{\Gamma(\alpha+n)} \mathrm{d}p(\alpha \mid a, b), \qquad (14)$$

where $p(\alpha \mid a, b)$ is the gamma probability measure with parameters $(a, b)$. In what follows, we re-write equation 14 as the limit of a sequence, and we use results on the weak convergence of measures to prove that the sequence converges to 1 for $k = 1$, and to 0 for every other admissible value of $k$.

For each path in $\mathbb{R}^+ \times \mathbb{R}^+$ where $(a, b) \to (0, 0)$, we define $\{X_l\}$, a sequence of gamma distributed random variables $X_1, X_2, \ldots$, each with parameters $(a_1, b_1), (a_2, b_2), \ldots$ identified along that path. The distribution function of $X_l$ is

$$F_{X_l}(\alpha; a_l, b_l) = \frac{\gamma(a_l, b_l\alpha)}{\Gamma(a_l)},$$

where $\gamma()$ is the lower incomplete gamma function. It is well known that

$$\frac{\gamma(s, t)}{t^s} \to \frac{1}{s},$$

as $t \to 0$. Hence

$$\lim_{(a,b)\to(0,0)} \gamma(a, b\alpha) \frac{1}{\Gamma(a)} = \lim_{(a,b)\to(0,0)} \frac{(b\alpha)^a}{a} \frac{1}{\Gamma(a)} = \lim_{l\to\infty} \frac{(b_l\alpha)^{a_l}}{a_l} \frac{1}{\Gamma(a_l)}$$

10

$$= \lim_{l \to \infty} (b_l \alpha)^{a_l} \frac{1}{a_l \Gamma(a_l)} = \lim_{l \to \infty} b_l{}^{a_l} \cdot \lim_{l \to \infty} \alpha^{a_l} = 1, \quad \forall \alpha > 0,$$

where we use the fact that

$$\lim_{(a,b) \to (0,0)} b^a = \lim_{r \to 0} (r \sin \theta)^{r \cos \theta} = \lim_{r \to 0} (r^r)^{\cos \theta} (\sin^r \theta)^{\cos \theta} = 1,$$

which holds as the paths along $a = 0$ and $b = 0$ do not belong to the function domain $\mathbb{R}^+ \times \mathbb{R}^+$ of $F_{X_l}$.

Therefore $X_l \xrightarrow{d} 0$, since the distribution function of the r.v. $X = 0$ is equal to 1 over the continuity set $(0, \infty)$ of $X$, for every path and every sequence $\{X_l\}$. A consequence is that the sequence of gamma probability measures $\{p_l\}$ induced by $\{X_l\}$ converges weakly to the Dirac measure $\delta_0$. Hence equation 14 becomes

$$\lim_{(a,b) \to (0,0)} p(K_n = k) = s_{n,k} \int_0^\infty \alpha^k \frac{\Gamma(\alpha)}{\Gamma(\alpha + n)} \mathrm{d}\delta_0(\alpha) = \begin{cases} 1, & k = 1, \\ 0, & k = 2, \ldots, n, \end{cases}$$

as $\alpha^k \Gamma(\alpha) / \Gamma(\alpha + n)$ is bounded and continuous.

### 3.3. Improper priors

In the preceding subsection, we mentioned how some studies motivate the use of quasi-degenerate gamma priors on the basis that they approximate the improper prior $p(\alpha) \propto 1/\alpha$, therefore implying that improper priors are desirable.

However, using $p(\alpha) \propto 1/\alpha$ leads to an improper posterior $p(\alpha \mid \boldsymbol{y})$ as well as to an improper implied prior $p(K_n)$, and so does $p(\alpha) \propto 1$. In fact, from equation 8 we obtain

$$p(\alpha \mid K_n = k) \propto \alpha^k \frac{\Gamma(\alpha)}{\Gamma(\alpha + n)} \cdot \frac{1}{\alpha} \sim \frac{1}{\alpha^{n-k+1}} \quad \text{as } \alpha \to \infty,$$

which is not integrable on any $(t, \infty)$ interval in the domain for $k = n$, meaning that the $\alpha$ posterior induced by the prior $1/\alpha$ cannot be normalised and is improper.

The $p(\alpha) \propto 1/\alpha$ prior also induces an improper prior on $K_n$, as

$$p(K_n = k) = s_{n,k} \int_0^\infty \alpha^k \frac{\Gamma(\alpha)}{\Gamma(\alpha + n)} \cdot \frac{1}{\alpha} \, \mathrm{d}\alpha,$$

where the integrand diverges as above.

Similar considerations apply when the $\alpha$ prior is $p(\alpha) \propto 1$, the conclusion being that the $\alpha$ posterior and the induced $K_n$ prior are improper because they are divergent for $k = n$ and $k = n - 1$.

## 4. Sample-size-independent priors for $\alpha$ (SSI)

We introduce a new prior selection approach which is independent of sample size, and which is motivated by the fact that, irrespective of how many $K_n$ clusters are observed in a sample of size $n$, there is always an underlying infinite-dimensional collection of point masses induced by the $\mathrm{DP}(\alpha, G_0)$, which is independent of $n$ and which can be made to

11

reflect one's prior beliefs. These point masses are the weights, or lengths, of the sticks that are broken off the initial stick of length 1, in the stick-breaking representation of the Dirichlet process. They also equivalently correspond to the asymptotic relative cluster sizes for $n \to \infty$. We consider two options: their size-biased random permutation, and their ranked permutation. We henceforth refer to this approach as SSI (sample-size-independent).

### 4.1. Size-biased weights

In the stick-breaking representation of equation 4, a DP is a collection of infinitely many point masses. In particular, $\{w_h\}_{h=1}^{\infty}$ is a sequence of stochastically decreasing weights; it is said to be a size-biased random permutation of the weights because the probability of each weight $w_i$ being in first position in the size-biased permutation is precisely $w_i$ (Pitman, 1996). A priori, the probability distribution of the weights as they naturally arise in the stick-breaking construction of equation 3 and 4 is the same as the probability distribution of the size-biased weights, hence for simplicity we do not distinguish between them in notation.

The conditional joint probability distribution of the first $H$ size-biased weights is a particular case of the generalised Dirichlet distribution (Connor and Mosimann, 1969):

$$p\left(w_1, \ldots, w_H \mid \alpha\right) = \alpha^H \frac{\left(1 - w_1 - \ldots - w_H\right)^{\alpha-1}}{\left(1 - w_1\right) \ldots \left(1 - w_1 - \ldots - w_{H-1}\right)}. \tag{15}$$

We plot its first two elements in Figure 5, for selected values of $\alpha$. Since $p\left(w_1, w_2, \ldots\right)$ is infinite-dimensional, we restrict our analysis to a finite number of dimensions; for practical purposes and for simplicity, we focus on the bivariate distribution of $w_1$ and $w_2$; our approach can in principle be extended to more dimensions, if necessary. Analytically, we have that:

$$p\left(w_1, w_2 \mid \alpha\right) = \frac{\alpha^2}{1 - w_1}\left(1 - w_1 - w_2\right)^{\alpha-1}.$$

The density $p\left(w_1, w_2 \mid \alpha\right)$ attains its maximum on $w_2 = 1 - w_1$ for $\alpha < 1$, at $(1, 0)$ for $1 \leq \alpha < 2$, at $w_2 = 0$ for $\alpha = 2$ and at $(0, 0)$ for $\alpha > 2$ (Figure 5). Its partial derivatives are:

$$\frac{\partial p\left(w_1, w_2 \mid \alpha\right)}{\partial w_1} = \alpha^2 \frac{\left(1 - w_1 - w_2\right)^{\alpha-2}}{\left(1 - w_1\right)^2}\left(\left(\alpha - 2\right)w_1 - w_2 - \alpha + 2\right),$$

$$\frac{\partial p\left(w_1, w_2 \mid \alpha\right)}{\partial w_2} = \alpha^2\left(1 - \alpha\right)\frac{\left(1 - w_1 - w_2\right)^{\alpha-2}}{1 - w_1},$$

and analysis of their sign leads to the considerations in Table 3.

Our prior belief in the cases displayed in Table 3 can be used to inform our choice of the parameter vector $\boldsymbol{\eta}$ of the prior distribution $p\left(\alpha \mid \boldsymbol{\eta}\right)$, since assigning a prior to $\alpha$ means mixing over those cases.

Denote by $\boldsymbol{\eta} = \left(\eta_1, \ldots, \eta_d\right)$ the $d$ parameters of the continuous probability distribution $p\left(\alpha \mid \boldsymbol{\eta}\right)$, and denote by $p_{0,t_1}, \ldots, p_{t_d,\infty}$ our prior belief associated with a partition of $(0, \infty)$ into $d+1$ nonempty subsets $\{(0, t_1], \ldots, (t_d, \infty)\}$. For example, when $d = 2$, $t_1 = 1$, $t_2 = 2$, we partition $(0, \infty)$ into $(0, 1], (1, 2]$ and $(2, \infty)$.

| $\alpha$ | $p(w_1 \mid w_2, \alpha)$ | $p(w_2 \mid w_1, \alpha)$ |
|---|---|---|
| $0 < \alpha < 1$ | increasing | increasing |
| $\alpha = 1$ | increasing | constant |
| $1 < \alpha < 2$ | concave; max. attained at $w_1 = 1 - \frac{w_2}{2-\alpha}$ | decreasing |
| $\alpha = 2, w_2 \neq 0$ | decreasing | decreasing |
| $\alpha = 2, w_2 = 0$ | constant | decreasing |
| $\alpha > 2$ | decreasing | decreasing |

Table 3: Sample-size-independent approach, size-biased. Behaviour of $w_1 \mid w_2, \alpha$ and $w_2 \mid w_1, \alpha$, for different values of $\alpha$.

We then choose $\boldsymbol{\eta}$ so that $p(\alpha \mid \boldsymbol{\eta})$ reflects our prior belief by solving:

$$\begin{cases} p(0 < \alpha \leq t_1 \mid \boldsymbol{\eta}) = p_{0,t_1}, \\ p(t_1 < \alpha \leq t_2 \mid \boldsymbol{\eta}) = p_{t_1,t_2}, \\ \dots \\ p(\alpha > t_d \mid \boldsymbol{\eta}) = p_{t_d,\infty}. \end{cases}$$

The resulting system of equations can be either solved analytically, if the cumulative probability distribution admits an explicit representation of its inverse, or numerically. For example, this is analytically feasible when $\alpha \sim \mathrm{Exp}(\eta)$, and we obtain:

$$\eta = \log\left((1 - p_{0,t_1})^{-\frac{1}{t_1}}\right), \tag{16}$$

and clearly, when $d = 1$, $p_{t_1,\infty} = 1 - p_{0,t_1}$.

When instead $\alpha \sim \mathrm{Ga}(\eta_1, \eta_2)$, we obtain:

$$\begin{cases} \dfrac{\gamma(\eta_1, \eta_2 t_1)}{\Gamma(\eta_1)} = p_{0,t_1}, \\ \dfrac{\gamma(\eta_1, \eta_2 t_2) - \gamma(\eta_1, \eta_2 t_1)}{\Gamma(\eta_1)} = p_{t_1,t_2}. \end{cases}$$

For values of $t_1, t_2$ that mirror those from table 3, and for some arbitrary choices of the underlying probabilities, we obtain the results in table 4. These results are purely for exemplification, and the approach that we outline in this section can be used to calculate $\boldsymbol{\eta}$ for any partition of $(0, \infty)$, any associated probabilities, and any distributional choice of $p(\alpha \mid \boldsymbol{\eta})$.

Unconditionally, we have that:

$$p(w_1, w_2) = \int_0^\infty \alpha^2 \, \frac{(1 - w_1 - w_2)^{\alpha-1}}{1 - w_1} \, \mathrm{d}p(\alpha \mid \boldsymbol{\eta}),$$

which, when $\alpha \sim \mathrm{Ga}(a, b)$, leads to the following analytic expressions (see Appendix A):

$$p(w_1, w_2) = a(a+1)b^a \, \frac{(b - \log(1 - w_1 - w_2))^{-a-2}}{(1 - w_1)(1 - w_1 - w_2)},$$

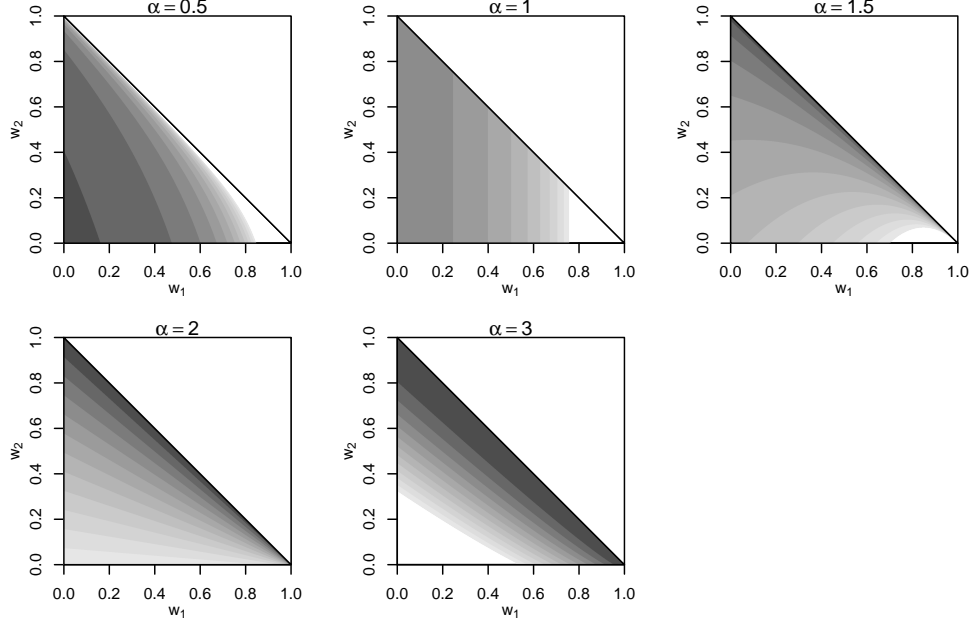$$p(w_1) = ab^a \, \frac{(b - \log(1 - w_1))^{-a-1}}{(1 - w_1)},$$

13

Figure 5: Sample-size-independent approach, size-biased. Conditional joint probability distribution $p(w_1, w_2 \mid \alpha)$, for different values of $\alpha$. Darker colours indicate smaller values.

| Distribution | $p(0 < \alpha < 1)$ | $p(1 < \alpha < 2)$ | $p(\alpha > 2)$ | $\eta_1$ | $\eta_2$ |
|---|---|---|---|---|---|
| Gamma | 1/3 | 1/3 | 1/3 | 1.814 | 1.036 |
| Gamma | 1/2 | 1/4 | 1/4 | 1.000 | 0.693 |
| Lognormal | 1/3 | 1/3 | 1/3 | 0.347 | 0.805 |
| Lognormal | 1/2 | 1/4 | 1/4 | 0.000 | 1.028 |
| Half-Cauchy | 1/2 | $p(\alpha > 1) = 1/2$ | | 0.000 | 1.000 |

Table 4: Sample-size-independent approach, size-biased. A selection of results for different distributional choices of $\alpha$, and for different choices of $p(0 < \alpha < 1)$, $p(1 < \alpha < 2)$, $p(\alpha > 2)$.

which can potentially be used for further analytical considerations, when setting $(a, b)$.

The plots in Figure 6 confirm that the joint unconditional distribution reflects a mix of the characteristics of the conditional joint distributions, in that the probability is amassed at the vertices $(0, 1)$ and $(1, 1)$, and along the edge that connects $(1, 0)$ and $(0, 0)$.

### 4.2. Ranked weights

The same approach from section 4.1 can also be applied to ranked stick weights, in the Poisson-Dirichlet distribution representation (see section 2.2). We re-write equation 5 as follows:

$$p\left(w_1^{\downarrow}, \ldots, w_r^{\downarrow} \mid \alpha\right) = \alpha^r \frac{\left(1 - w_1^{\downarrow} - \ldots - w_r^{\downarrow}\right)^{\alpha - 1}}{w_1^{\downarrow} \cdots w_r^{\downarrow}} F_\alpha\left(\frac{w_r^{\downarrow}}{1 - w_1^{\downarrow} - \ldots - w_r^{\downarrow}}\right), \qquad (17)$$

14
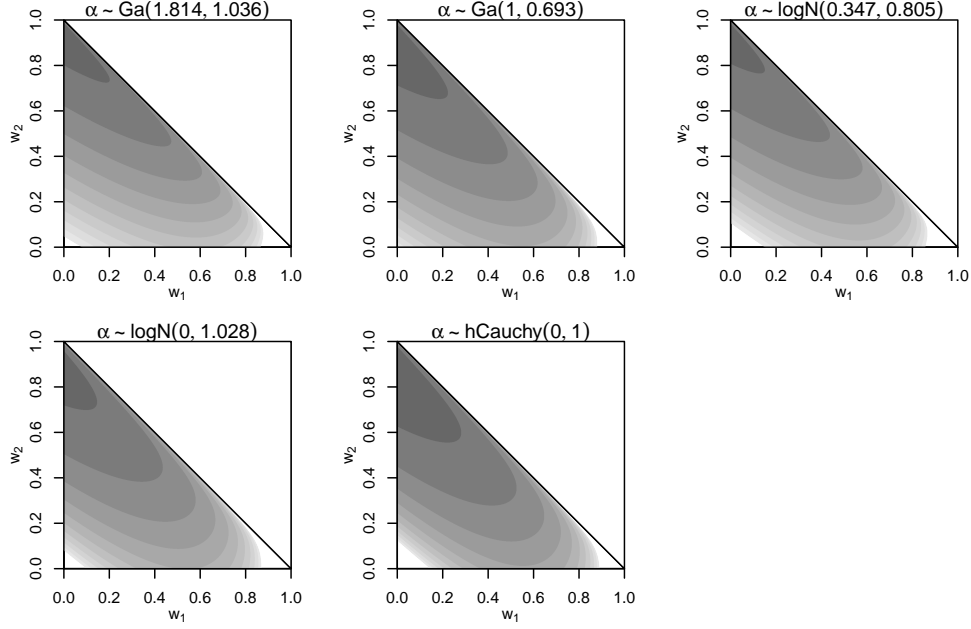
Figure 6: Sample-size-independent approach, size-biased. Joint probability distribution $p(w_1, w_2)$, for the distributional choices of $\alpha$ identified in table 4. Darker colours indicate smaller values.

with $F_\alpha(x) := p\left(w_1^\downarrow \leq x \mid \alpha\right)$, which can be obtained by simulation. In the bivariate case, equation 17 becomes:

$$p\left(w_1^\downarrow, w_2^\downarrow \mid \alpha\right) = \alpha^2 \frac{\left(1 - w_1^\downarrow - w_2^\downarrow\right)^{\alpha-1}}{w_1^\downarrow w_2^\downarrow} \, F_\alpha\left(\frac{w_2^\downarrow}{1 - w_1^\downarrow - w_2^\downarrow}\right), \tag{18}$$

which is defined on

$$E = \left\{\left(w_1^\downarrow, w_2^\downarrow\right) : \left(w_1^\downarrow + w_2^\downarrow < 1\right) \, \wedge \, \left(w_2^\downarrow < w_1^\downarrow\right)\right\}.$$

Since $F_\alpha$ is a cumulative probability distribution function, we have that

$$F_\alpha\left(\frac{w_2^\downarrow}{1 - w_1^\downarrow - w_2^\downarrow}\right) = 1, \quad \text{for } w_2^\downarrow \geq \frac{1}{2} - \frac{w_1^\downarrow}{2}.$$

As such, the restriction of equation 18 to

$$A = \left\{\left(w_1^\downarrow, w_2^\downarrow\right) : \left(w_1^\downarrow + w_2^\downarrow < 1\right) \, \wedge \, \left(w_2^\downarrow < w_1^\downarrow\right) \, \wedge \, \left(w_2^\downarrow \geq \frac{1}{2} - \frac{1}{2}w_1^\downarrow\right)\right\}$$

yields

$$p\,|_A\left(w_1^\downarrow, w_2^\downarrow \mid \alpha\right) = \alpha^2 \frac{\left(1 - w_1^\downarrow - w_2^\downarrow\right)^{\alpha-1}}{w_1^\downarrow w_2^\downarrow}, \tag{19}$$

15

| $\alpha$ | $p\mid_A \left(w_1^\downarrow \mid w_2^\downarrow, \alpha\right)$ | $p\mid_A \left(w_2^\downarrow \mid w_1^\downarrow, \alpha\right)$ |
|---|---|---|
| $0 < \alpha < 1$ | convex; min. at $w_1^\downarrow = \frac{1-w_2^\downarrow}{2-\alpha}$ | convex; min. at $w_2^\downarrow = \frac{1-w_1^\downarrow}{2-\alpha}$ |
| $\alpha \geq 1$ | decreasing | decreasing |

Table 5: Sample-size-independent approach, ranked. Behaviour of $w_1^\downarrow \mid w_2^\downarrow, \alpha$ and $w_2^\downarrow \mid w_1^\downarrow, \alpha$, for different values of $\alpha$, over $A$.

whose partial derivatives are easier to study than those of equation 18, and which bring some insights. In particular:

$$\frac{\partial p\mid_A \left(w_1^\downarrow, w_2^\downarrow \mid \alpha\right)}{\partial w_1^\downarrow} = \alpha^2 \frac{\left(1 - w_1^\downarrow - w_2^\downarrow\right)^{\alpha-2}}{w_1^{\downarrow 2} w_2^\downarrow} \left((2-\alpha) w_1^\downarrow + w_2^\downarrow - 1\right),$$

$$\frac{\partial p\mid_A \left(w_1^\downarrow, w_2^\downarrow \mid \alpha\right)}{\partial w_2^\downarrow} = \alpha^2 \frac{\left(1 - w_1^\downarrow - w_2^\downarrow\right)^{\alpha-2}}{w_1^\downarrow w_2^{\downarrow 2}} \left((2-\alpha) w_2^\downarrow + w_1^\downarrow - 1\right),$$

which leads to the considerations in table 5. Following the same approach we used in section 4.1, we also plot in Figure 7 the joint bivariate distribution of $p\left(w_1^\downarrow, w_2^\downarrow \mid \alpha\right)$ for key values of $\alpha$, to spot any further visible patterns in the behaviour of $p\left(w_1^\downarrow, w_2^\downarrow \mid \alpha\right)$ as $\alpha$ varies.

The process for choosing the probability distribution $p\left(\alpha \mid \boldsymbol{\eta}\right)$ and its parameters is analogous to the one in section 4.1. In Figure 8, we plot $p\left(w_1^\downarrow, w_2^\downarrow\right)$ for different distributional choices of $p\left(\alpha \mid \boldsymbol{\eta}\right)$, to show the impact of marginalising $\alpha$ out.

## 5. Discussion

In section 3 we studied the behaviour of SSD, quasi-degenerate and improper priors in view of the distribution that they induce on $K_n$, while in 4 we assessed SSI priors with respect to their implied stick-breaking distribution. Here we do the opposite: we cross-check how SSD and SSI priors behave in relation to each other's driving metric. For the sake of brevity, we only discuss size-biased stick breaking weights; conclusions with respect to ranked weights are similar.

We observe that the behaviour of the $K_n$-diffuse, the DORO, the quasi-degenerate[3] and Jeffreys' priors with respect to $p\left(w_1, w_2\right)$ (Figure 9) is markedly different from the behaviour of SSI (Figure 6), and more extreme, because probability in SSD is mostly concentrated at $(0,0)$ or $(1,0)$, while in SSI it is more spread out. SCAL is closer to SSI in this respect. By looking at 10 too, we conclude that the $K_n$-diffuse, the DORO and Jeffreys' priors would likely attract posterior estimates of $(w_1, w_2)$ towards $(0,0)$, while the quasi-degenerate prior would attract them towards $(1,0)$. We also note that neither SCAL nor the quasi-degenerate priors are diffuse in $K_n$ (Figure 11).

Conversely, with respect to SSI, we consider the five test cases identified in Table 4, and we plot their implied distribution of $K_n$ in Figure 12. We observe that their implied

---

[3]We parameterized the quasi-degenerate prior as a $\mathrm{Ga}\left(0.403, 0.370\right)$, to mirror the result from the approximation method of SCAL, although we could have used any smaller value of $(a, b)$. Conclusions would not materially change.
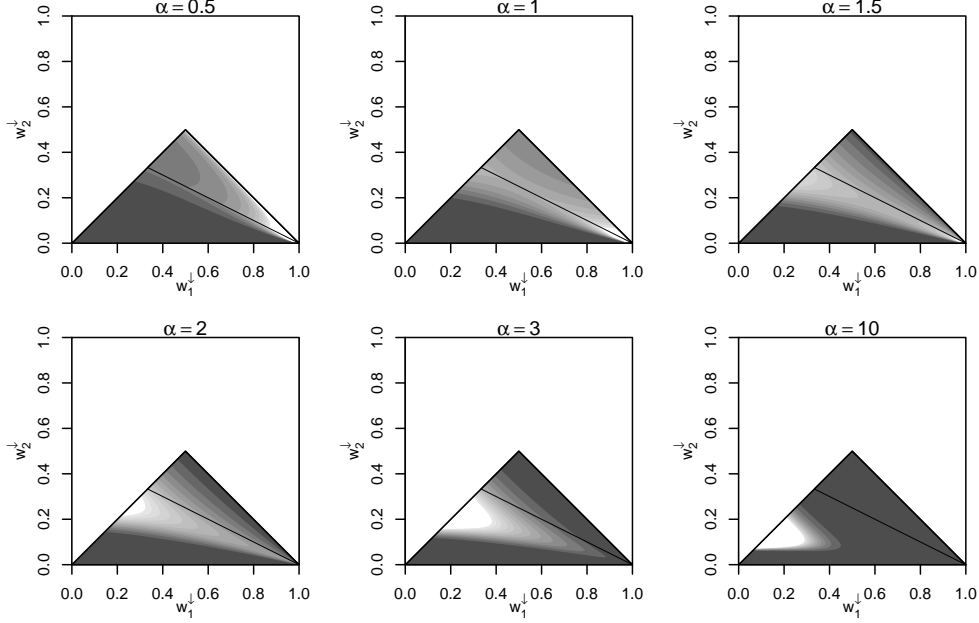
Figure 7: Sample-size-independent approach, ranked. Joint probability distribution $p\left(w_1^\downarrow, w_2^\downarrow \mid \alpha\right)$, for different values of $\alpha$. The top right triangle identifies $A \subset E$ (see equation 18). Darker colours indicate smaller values.

$p\left(K_n\right)$ is concentrated over values that are small, relative to $n$. We highlight that these five cases are just examples, as one's prior information to be reflected with SSI may well be entirely different from the cases in table 4.

We more generally conclude that choices of $p\left(\alpha \mid \boldsymbol{\eta}\right)$ that are diffuse in $K_n$ are not necessarily diffuse in $(w_1, w_2)$, and vice-versa. As such, users should assess whether the SSD or SSI prior approach is more suitable with respect to their particular problem at hand; using both $p\left(K_n\right)$ and $p\left(w_1, w_2\right)$ as a reference to set $p\left(\alpha \mid \boldsymbol{\eta}\right)$ could also be an option, if $p\left(K_n\right)$ is relevant to the application. We argue that SSI priors still hold a number of advantages over SSD, as we summarize in Section 7.

## 6. Case study: multiple DPMs

In this section we point to a setting where, by construction, the sample-size-dependent approach to choosing $p\left(\alpha \mid \boldsymbol{\eta}\right)$ is inapplicable, while the sample-size-independent approach can still be used. This was inspired by (Müller and Rodriguez, 2013, figure 5.1).

Consider a partially exchangeable setting of $J$ groups with $n_j$ observations in group $j$, where $j = 1, \ldots, J$ and the groups share the same common underlying precision parameter $\alpha$. This framework, while extremely simple, allows information to be borrowed between groups through their dependence on a common $\alpha$.

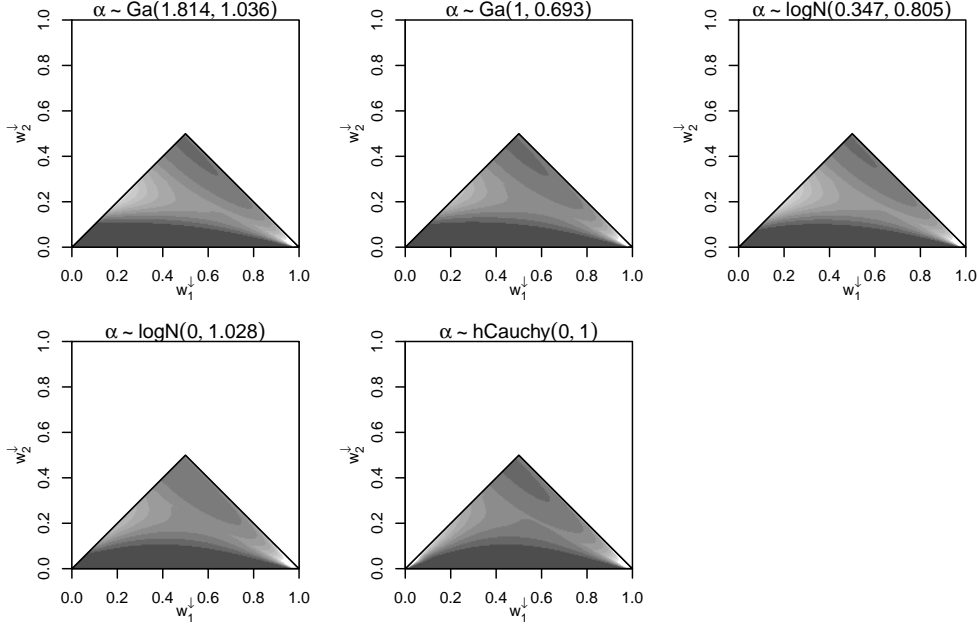We denote the observations by $y_{i,j}$, $i = 1, \ldots, n_j$, $j = 1, \ldots, J$, and the random

17

Figure 8: Sample-size-independent approach, ranked. Joint probability distribution $p\left(w_1^\downarrow, w_2^\downarrow\right)$, for the distributional choices of $\alpha$ identified in Table 4. Darker colours indicate smaller values.

number of clusters in group $j$ by $K_{n,j}$. The model is:

$$
\begin{aligned}
y_{i,j} \mid \theta_{i,j} &\sim p\left(y_{i,j} \mid \theta_{i,j}\right) \\
\theta_{i,j} \mid G_j &\sim G_j \\
G_j \mid \alpha &\sim \mathrm{DP}\left(\alpha, G_{0j}\right) \\
\alpha &\sim p\left(\alpha \mid \boldsymbol{\eta}\right).
\end{aligned}
$$

It is clear that, by construction, multiple $K_{n,j}$ are involved, hence there is no unique probability distribution for $K_n$ and no unique sample size $n$ that can be used as a target for the sample-size-dependent approaches that we have described. Conversely, our sample-size-independent approach is still applicable.

A different example of a set-up that is similarly unsuitable for SSD priors due to the lack of a unique sample size is that of online inferential algorithms for streaming data, because the sample size changes continuously as more data is collected.

## 7. Conclusions

In section 3 we highlighted the limitations of previous approaches. Their limitations include:

- Dependence on $n$. All SSD priors are only optimal for one particular value of $n$, and they need to be updated as $n$ varies;

- Unmet assumptions. DORO minimises the Kullback-Leibler distance between $p(K_n)$ and the discrete uniform distribution, but the discrete uniform target can never be attained, as $p(K_n)$ converges to a negative binomial for $n \to \infty$ (see section 3.1.2). The assumption that a non-informative prior would be reflected in a discrete uniform distribution induced on $K_n$ is not supported by Jeffreys' prior either, which has a very different shape from discrete uniform (see Figure 4);

- Asymptotic breakdown. The approximation suggested by SCAL in relation to their test case, as an alternative to their fully-fledged approach, is only valid over a narrow range of values of $n$; for moderately large $n$, it results in $\alpha \sim \mathrm{Ga}(a, b)$, $(a, b) \approx (0, 0)$. This is undesirable, as it implies $p(K_n = 1) \approx 1$, which negates the reason for using a nonparametric model in the first place (see sections 3.1.3, 3.2);

- Over-informativeness. Although $K_n$-diffuse and DORO priors are diffuse in $K_n$, they are very informative with respect to $(w_1, w_2)$ and $\left(w_1^{\downarrow}, w_2^{\downarrow}\right)$ (see section 5);

- Inapplicability. SSD priors are inapplicable when multiple DPMs with different sample sizes share the same $\alpha$ (see section 6); SSD priors are not well suited either to cases where the sample size is allowed to grow.

- Incompatibility. Quasi-degenerate priors are not compatible with the notion of multiple clusters, hence they defeat the purpose of using a DPM in the first place (see 3.2);

- Impropriety. Improper priors can lead to an improper posterior $p(\alpha \mid \boldsymbol{y})$ as well as to an improper prior induced on $K_n$ (see 3.3).

In 4, we introduced a new approach, based on the appraisal of the implied joint distribution of the stick-breaking weights (either in size-biased order, or ranked). This is equivalent to thinking in terms of the asymptotic relative cluster sizes for $n \to \infty$. Our approach does not suffer from any of the aforementioned limitations, although we acknowledge that SSD priors still have a place, depending on the specific nature of the problem at hand, and that potentially both $p(K_n)$ and the stick-breaking weights could be concurrently leveraged to inform one's choice of $p(\alpha \mid \boldsymbol{\eta})$.

Out of the two alternatives that we propose (size-biased weights, and ranked weights), the one that is based on $p(w_1, w_2)$ appears to be the easiest to interpret, compute, and evaluate, due to the distinctive behaviour of $p(w_1, w_2 \mid \alpha)$ for various levels of $\alpha$ (Figure 5) and to the availability of simpler analytical formulae. Quantitative measurements can be carried out, as exemplified in Table 4, to determine how to mix precisely over those behaviours. A simple exact formula is also available when $\alpha \sim \mathrm{Exp}(\eta)$ (see equation 16). We envisage use of SSI priors by practitioners as a valid alternative to SSD priors, not only because of the principle that underlies them (which has general appeal in all applications) but also because they allow for two applications that would not otherwise be feasible with SSD priors. One is the area of online learning algorithms, where the focus is to obtain an end-to-end inferential algorithm for streaming data that allows updates to the posterior estimates as the size of the data grows: in such applications, a prior that does not depend on the sample size is clearly very desirable. DPM models are also particularly well suited to streaming data as they allow the number of clusters to grow as the size of the data set increases. The second is the case of multiple DPM models that

19

share a common underlying precision parameter, which allows the borrowing of strength between models and which is not feasible with traditional sample-size-dependent priors because there is no single sample size with which to parameterize the prior (see section 6).

## Appendix A. Unconditional size-biased weight distribution

When $\alpha \sim \text{Ga}(a, b)$, the joint distribution of the first $H$ size-biased weights of a stick-breaking process has a simple closed analytical formula, once $\alpha$ is integrated out of equation 15:

$$
\begin{aligned}
p(w_1, \ldots, w_H) &= \int_0^\infty \frac{\alpha^H (1 - w_1 - \ldots - w_H)^{\alpha - 1}}{(1 - w_1) \ldots (1 - w_1 - \ldots - w_{H-1})} dp(\alpha) \\
&= \frac{\Gamma(a + H)}{\Gamma(a)} b^a \frac{[b - \log(1 - w_1 - \ldots - w_H)]^{-a-H}}{(1 - w_1) \ldots (1 - w_1 - \ldots - w_H)} \qquad \text{(A.1)} \\
&\propto \frac{[b - \log(1 - w_1 - \ldots - w_H)]^{-a-H}}{(1 - w_1)(1 - w_1 - \ldots - w_H)}. \qquad \text{(A.2)}
\end{aligned}
$$

The cumulative probability distribution of $w_1$ also has a simple explicit analytical representation:

$$
p(w_1 \leq x) = \int_0^x p(w_1)\, dw_1 = 1 - \left( \frac{b}{b - \log(1 - x)} \right)^a.
$$

## Appendix B. Sampling from Jeffreys' prior

As Jeffreys' prior for the Dirichlet process is a non-standard distribution, there is value in discussing how to sample from it. Like any other distribution for which there are no widely known sampling methods, one can resort to any of the general-purpose approaches described, for example, in Robert and Casella (2004), such as the accept-reject method, importance sampling, the 2-d slice sampler or Metropolis-Hastings algorithms.

The accept-reject algorithm requires that a constant $M$ exists such that the target distribution $f$ and the proposal distribution $g$ meet the condition $\frac{f(\alpha)}{g(\alpha)} \leq M$, for all values of $\alpha$. This can be achieved by using Jeffreys' prior for $n = 2$ (which has an analytical solution, as per equations 12 and 13) as the proposal distribution $g$; easy calculations lead then to

$$
\frac{f(\alpha)}{g(\alpha)} = \sqrt{1 + 2 \frac{(\alpha + 1)^2}{(\alpha + 2)^2} + \ldots + (n - 1) \frac{(\alpha + 1)^2}{(\alpha + n - 1)^2}} \leq M = \sqrt{\sum_{i=1}^{n-1} i}.
$$

Once $M$ is determined, the steps to the algorithm are:

- generate $X$ from the proposal distribution $g$, and generate $U \sim \text{Unif}(0, 1)$;

- accept the value above if $U \leq f(X) / (Mg(X))$; return to the step above if otherwise.

We also implemented, to compare:

- the 2d slice sampler

- the independence Metropolis-Hastings algorithm with a Jeffreys' prior distribution ($n = 2$);

- the random-walk Metropolis algorithm with half-Cauchy and normal proposal distributions;

and found the convergence speed to be fastest for the 2d slice sampler, followed by the independence Metropolis-Hastings algorithm with Jeffreys' prior distribution and $n = 2$, followed again by the random-walk Metropolis with half-Cauchy proposal. By comparison, Rodríguez (2013) use random-walk Metropolis-Hastings with a normal proposal distribution.

### Appendix  C. Propriety of the posterior induced by Jeffreys' prior on $\alpha$

Jeffreys' prior induces a proper posterior $p\left(\alpha \mid K_n = k\right)$. Given

$$p\left(\alpha \mid K_n = k\right) \propto p\left(K_n = k \mid \alpha\right) p\left(\alpha\right)$$

$$\propto \alpha^k \frac{\Gamma\left(\alpha\right)}{\Gamma\left(\alpha + n\right)} \sqrt{\frac{1}{\alpha} \sum_{i=1}^{n} \frac{i-1}{\left(\alpha + i - 1\right)^2}},$$

its limit for $\alpha \to 0$ is

$$\lim_{\alpha \to 0} p\left(\alpha \mid K_n = k\right) = \alpha^{k-1} \frac{\Gamma\left(1\right)}{\Gamma\left(n\right)} \alpha^{-\frac{1}{2}} \sqrt{\sum_{i=1}^{n} \frac{1}{i-1}} \approx \alpha^{k-\frac{3}{2}},$$

which is convergent over $(0, 1)$ (limit comparison test). Its limit for $\alpha \to \infty$ is

$$\lim_{\alpha \to \infty} p\left(\alpha \mid K_n = k\right) = \alpha^{k-n-\frac{3}{2}},$$

which is proper over $(t, \infty)$, for $t \geq 1$.

### Appendix  D. Propriety of the prior induced by Jeffreys' prior on $K_n$

Recall from equation 8 the expression for $p\left(K_n = k \mid \alpha\right)$. The prior it induces on $K_n$ is

$$p\left(K_n = k\right) = s_{n,k} \int_0^\infty \alpha^k \frac{\Gamma\left(\alpha\right)}{\Gamma\left(\alpha + n\right)} \sqrt{\frac{1}{\alpha} \sum_{i=1}^{n} \frac{i-1}{\left(\alpha + i - 1\right)^2}} \, \mathrm{d}\alpha$$

$$\propto s_{n,k} \int_0^\infty \frac{\Gamma\left(\alpha\right)}{\Gamma\left(\alpha + n\right)} \alpha^{k-\frac{1}{2}} \sqrt{\sum_{i=1}^{n} \frac{i-1}{\left(\alpha + i - 1\right)^2}} \, \mathrm{d}\alpha.$$

21

Denoting the integrand by $g(\alpha)$,

$$\lim_{\alpha \to 0} g(\alpha) = \lim_{\alpha \to 0} \frac{\Gamma(\alpha+1)}{\Gamma(\alpha+n)} \alpha^{k-\frac{3}{2}} \sqrt{\sum_{i=1}^{n} \frac{i-1}{(\alpha+i-1)^2}}$$

$$= \lim_{\alpha \to 0} \frac{\alpha^{k-\frac{3}{2}}}{\Gamma(n)} \sqrt{\sum_{i=1}^{n} \frac{1}{i-1}}$$

$$\propto \alpha^{k-\frac{3}{2}},$$

which converges over $(0, 1)$, and

$$\lim_{\alpha \to \infty} g(\alpha) \propto \alpha^{k-\frac{1}{2}-n-1},$$

which converges over $(1, \infty)$.

## References

Antoniak, C.E., 1974. Mixtures of Dirichlet processes with applications to Bayesian nonparametric problems. The Annals of Statistics 2, 1152–1174. URL: https://www.jstor.org/stable/2958336.

Arratia, R., Barbour, A.D., Tavaré, S., 2000. The number of components in a logarithmic combinatorial structure. The Annals of Applied Probability 10, 331–361. doi:10.1214/aoap/1019487347.

Arratia, R., Barbour, A.D., Tavaré, S., 2003. Logarithmic Combinatorial Structures: a Probabilistic Approach. European Mathematical Society. doi:10.4171/000.

Bauer, M., Bruveris, M., Michor, P.W., 2016. Uniqueness of the Fisher–Rao metric on the space of smooth densities. Bull. Lond. Math. Soc. 48, 499–506. doi:10.1112/blms/bdw020.

Campbell, L.L., 1986. An extended Čencov characterization of the information metric. Proc. Am. Math. Soc. 98, 135–141. doi:10.1090/s0002-9939-1986-0848890-5.

Čencov, N.N., 1982. Statistical Decision Rules and Optimal Inferences. volume 53. AMS.

Connor, R.J., Mosimann, J.E., 1969. Concepts of independence for proportions with a generalization of the Dirichlet distribution. Journal of the American Statistical Association 64, 194–206. doi:10.1080/01621459.1969.10500963.

Dorazio, R.M., 2009. On selecting a prior for the precision parameter of Dirichlet process mixture models. Journal of Statistical Planning and Inference 139, 3384–3390. URL: 10.1016/j.jspi.2009.03.009.

Escobar, M.D., West, M., 1995. Bayesian density estimation and inference using mixtures. Journal of the American Statistical Association 90, 577–588. doi:10.1080/01621459.1995.10476550.

Escobar, M.D., West, M., 1998. Computing nonparametric hierarchical models, in: Practical Nonparametric and Semiparametric Bayesian Statistics. Springer, pp. 1–22. doi:10.1007/978-1-4612-1732-9_1.

Ferguson, T.S., 1973. A Bayesian analysis of some nonparametric problems. The Annals of Statistics 1, 209–230. URL: https://www.jstor.org/stable/2958008.

Ferguson, T.S., 1983. Bayesian density estimation by mixtures of normal distributions, in: Recent Advances in Statistics: Papers in Honor of Herman Chernoff on his Sixtieth Birthday. Academic Press. volume 24, pp. 287–302. doi:10.1016/B978-0-12-589320-6.50018-6.

Jeffreys, H., 1946. An invariant form for the prior probability in estimation problems. Proceedings of the Royal Society of London. Series A. Mathematical and Physical Sciences 186, 453–461. doi:10.1098/rspa.1946.0056.

Kingman, J.F.C., 1975. Random Discrete Distributions. Journal of the Royal Statistical Society. Series B (Methodological) 37, 1–22. doi:10.1111/j.2517-6161.1975.tb01024.x.

Liu, J.S., 1996. Nonparametric hierarchical Bayes via sequential imputations. The Annals of Statistics 24, 911–930. doi:10.1214/aos/1032526949.

Lo, A.Y., 1984. On a class of Bayesian nonparametric estimates: I. density estimates. The Annals of Statistics 12, 351–357. URL: https://www.jstor.org/stable/2241054.

Lunn, D., Jackson, C., Best, N., Thomas, A., Spiegelhalter, D., 2012. The BUGS Book: A Practical Introduction to Bayesian Analysis. Chapman & Hall/CRC. doi:10.1201/b13613.

Müller, P., Rodriguez, A., 2013. Nonparametric Bayesian Inference. Institute of Mathematical Statistics. doi:10.1214/cbms/1362163742.

Murugiah, S., Sweeting, T., 2012. Selecting the precision parameter prior in Dirichlet process mixture models. Journal of Statistical Planning and Inference 142, 1947–1959. doi:10.1016/j.jspi.2012.02.013.

Navarro, D.J., Griffiths, T.L., Steyvers, M., Lee, M.D., 2006. Modeling individual differences using Dirichlet processes. Journal of Mathematical Psychology 50, 101–122. doi:10.1016/j.jmp.2005.11.006.

Pitman, J., 1996. Random discrete distributions invariant under size-biased permutation. Advances in Applied Probability 28, 525–539. doi:10.2307/1428070.

Robert, C.P., Casella, G., 2004. Monte Carlo Statistical Methods. Springer. doi:10.1007/978-1-4757-4145-2.

Rodríguez, A., 2013. On the Jeffreys prior for the multivariate Ewens distribution. Statistics & Probability Letters 83, 1539–1546. doi:10.1016/j.spl.2013.02.014.

Sethuraman, J., 1994. A constructive definition of Dirichlet priors. Statistica Sinica , 639–650URL: https://www.jstor.org/stable/24305538.

Watterson, G.A., 1974. The sampling theory of selectively neutral alleles. Advances in Applied Probability 6, 463–488. doi:10.2307/1426228.

Watterson, G.A., 1976. The stationary distribution of the infinitely-many neutral alleles diffusion model. Journal of Applied Probability 13, 639–651. doi:10.2307/3212519.

West, M., Escobar, M.D., 1993. Hierarchical Priors and Mixture Models, with Application in Regression and Density Estimation. Institute of Statistics and Decision Sciences, Duke University. doi:10.1057/jors.1995.91.
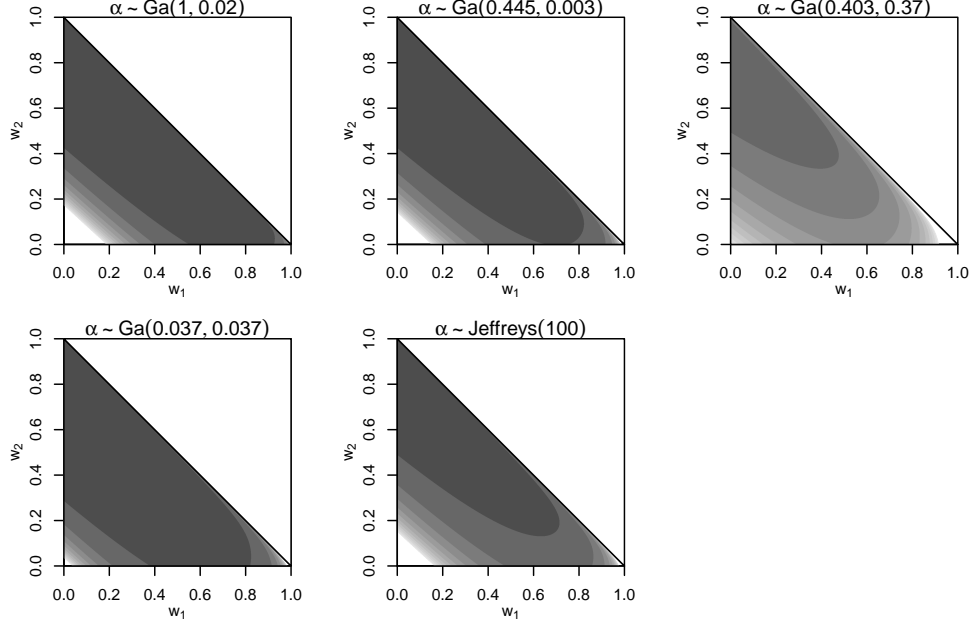
Figure 9: Size-biased joint probability distribution $p(w_1, w_2)$ underlying the $K_n$-diffuse, DORO, SCAL, quasi-degenerate and Jeffreys' priors for $n = 100$. Darker colours indicate smaller values.
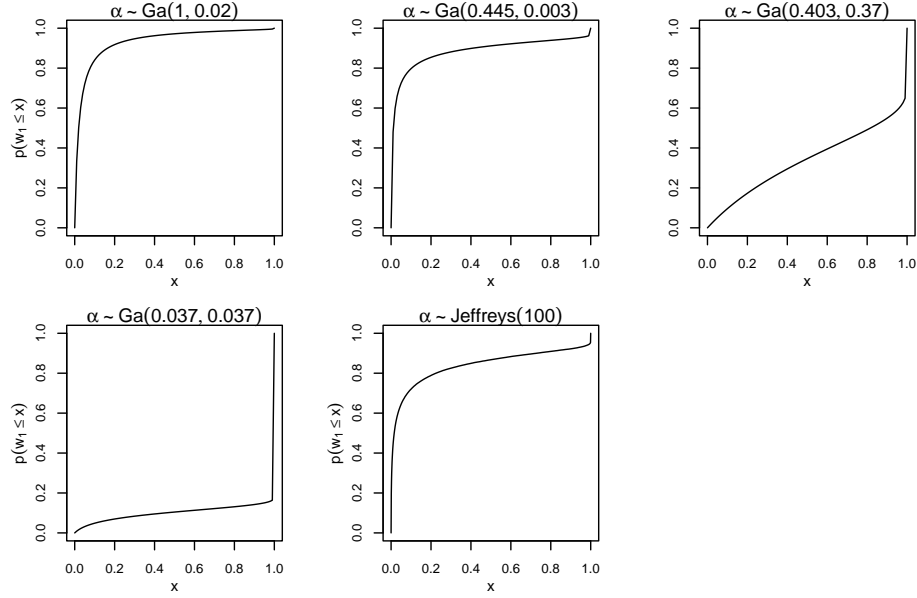


Figure 10: Size-biased cumulative probability distributions of $w_1$ underlying the $K_n$-diffuse, DORO, SCAL, quasi-degenerate and Jeffreys' priors for $n = 100$.
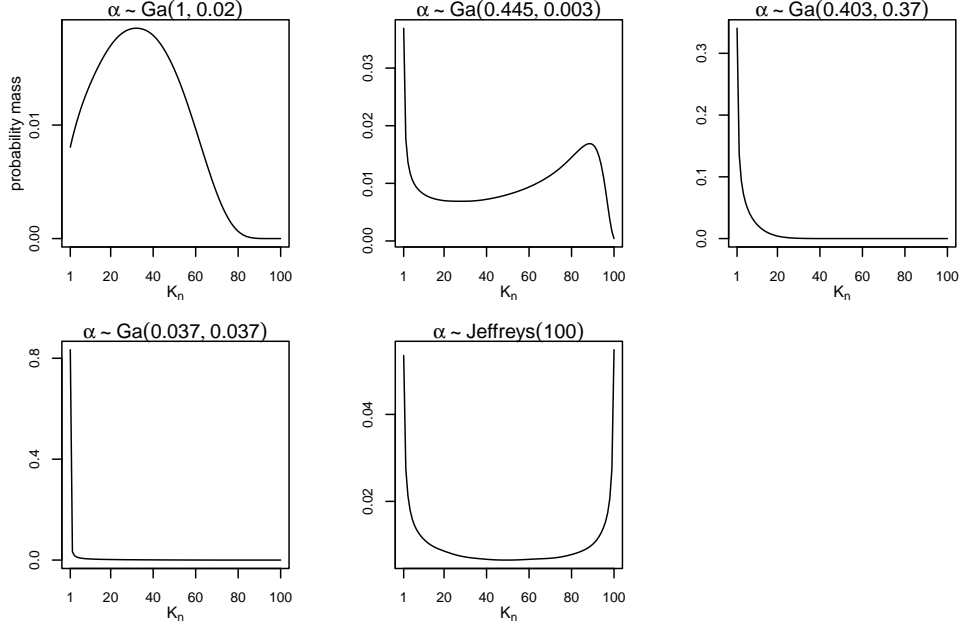
Figure 11: Prior distributions $p(K_n)$ induced by $p(\alpha)$, for the $K_n$-diffuse, DORO, SCAL, quasi-degenerate and Jeffreys' priors for $n = 100$.
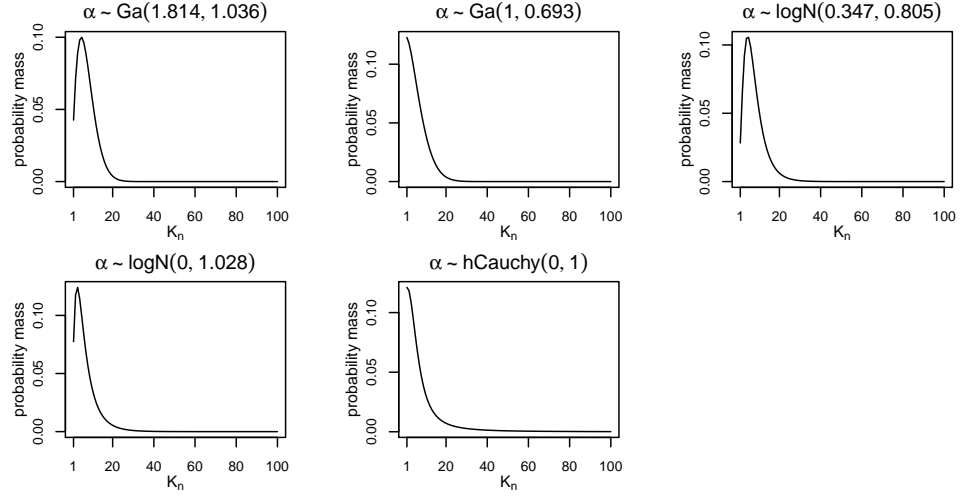


Figure 12: Prior distributions $p(K_n)$ induced by $p(\alpha)$, for $n = 100$, for the distributional choices of $\alpha$ identified in Table 4.

25