



---

# *Journal of Statistical Software*

June 2017, Volume 78, Issue 6.

doi: [10.18637/jss.v078.i06](https://doi.org/10.18637/jss.v078.i06)

---

## msBP: An R Package to Perform Bayesian Nonparametric Inference Using Multiscale Bernstein Polynomials Mixtures

Antonio Canale

University of Padua and Collegio Carlo Alberto

---

### Abstract

**msBP** is an R package that implements a new method to perform Bayesian multiscale nonparametric inference introduced by [Canale and Dunson \(2016\)](#). The method, based on mixtures of multiscale beta dictionary densities, overcomes the drawbacks of Pólya trees and inherits many of the advantages of Dirichlet process mixture models. The key idea is that an infinitely-deep binary tree is introduced, with a beta dictionary density assigned to each node of the tree. Using a multiscale stick-breaking characterization, stochastically decreasing weights are assigned to each node. The result is an infinite mixture model. The package **msBP** implements a series of basic functions to deal with this family of priors such as random densities and numbers generation, creation and manipulation of binary tree objects, and generic functions to plot and print the results. In addition, it implements the Gibbs samplers for posterior computation to perform multiscale density estimation and multiscale testing of group differences described in [Canale and Dunson \(2016\)](#).

*Keywords:* binary trees, density estimation, multiscale stick-breaking, multiscale testing.

---

## 1. Introduction

Multiscale methods have received abundant attention in the statistical literature, having several appealing characteristics that pushed their use in many applications. With the term “multiscale model” we refer to a model in which multiple sub-models at different scales are used simultaneously. A notable example is represented by wavelets, which are routinely used in signal and image processing, nonparametric regression, and density estimation ([Donoho, Johnstone, Kerkycharian, and Picard 1996](#)). However, from the Bayesian perspective, multiscale density estimation is surprisingly understudied. Indeed, most of the approaches rely on single-scale kernel mixtures. Among these, the Dirichlet process (DP; [Ferguson 1973, 1974](#))

mixtures of Gaussians (Lo 1984; Escobar and West 1995) are the gold standard in many applications. An exception is represented by Pólya trees (Mauldin, Sudderth, and Williams 1992; Lavine 1992a,b) that unfortunately have some unappealing characteristics. For example they tend to produce extremely spiky densities even when the true density is fairly smooth and are sensitive to the prior specification. This sensitivity can be overcome within a mixture approach (Hanson and Johnson 2002), but in this case there is a price to pay in terms of computation. Both the DP and Pólya tree mixture models are implemented in the famous **DPpackage** (Jara, Hanson, Quintana, Müller, and Rosner 2011), an R package that represents the de facto standard software for Bayesian nonparametric inference under a variety of settings.

Canale and Dunson (2016) recently proposed a Bayesian multiscale method that inherits some advantages of the DP mixture and avoids the disadvantages of Pólya trees. The key idea lies in introducing an infinitely-deep binary tree, with a beta dictionary density assigned to each node of the tree. Using a multiscale stick-breaking (Sethuraman 1994) characterization, the authors define a stochastically decreasing sequence of weights assigned to each node of the tree. This formulation implies that within a level of the tree, the densities are equivalent to Bernstein polynomials (Petroni 1999a,b). Extensions to deal with unconstrained domain data are also discussed. A similar idea appeared also in Chen, Hanson, and Zhang (2014). The DP-like characteristics are derived from the formulation of a multiscale generalization of the stick-breaking process, which can be exploited to build an efficient slice sampling algorithm. The same multiscale stick-breaking process has also been used by Wang, Canale, and Dunson (2016) to learn the joint density in massive dimensional settings, using geometric multiresolution analysis to estimate the dictionary densities over the binary tree at a first stage.

The R package **msBP** (Canale 2017) implements the multiscale stick-breaking process, and its applications to density estimation and to testing of group differences as discussed in Canale and Dunson (2016), and a series of basic R functions to deal with this family of nonparametric priors such as random density and number generation, creation and manipulation of binary trees, and generic functions to plot and print the results. The package's core is written in C++ by means of a specific 'bintree' data class and it is called from R via the `.C` function. The package is available from the Comprehensive R Archive Network (CRAN) at <https://CRAN.R-project.org/package=msBP>.

The rest of the paper is organized as follows: In the next section we outline the theoretical framework with particular emphasis on the multiscale stick-breaking process. Section 3 describes the main features of the C++ implementation while Section 4 is concerned with demonstrating the main features of the package.

## 2. A multiscale prior for densities

### 2.1. Basic formulation

Let  $x \in \mathcal{X} \subset \mathbb{R}$ , be a random variable,  $g$  be an unknown density and  $x \sim g$ . Under a Bayesian perspective  $g_0$  is assumed to be a prior guess for  $g$ , with  $G_0$  and  $G_0^{-1}$  the corresponding cumulative distribution function (CDF) and inverse CDF, respectively. A prior for  $g$  centered on  $g_0$  can be introduced through a prior for the density  $f$  of  $y = G_0(x) \in (0, 1)$ . The CDFs

$F$  and  $G$  corresponding to the densities  $f$  and  $g$ , respectively, have the following relationship

$$G(x) = F\{G_0(x)\}, x \in \mathcal{X}, \quad F(y) = G\{G_0^{-1}(y)\}, y \in (0, 1). \quad (1)$$

A similar construction also appeared in [Chen \*et al.\* \(2014\)](#). The density  $f$  is assumed to have the following structure:

$$f(y) = \sum_{s=0}^{\infty} \sum_{h=1}^{2^s} \pi_{s,h} \text{Be}(y; h, 2^s - h + 1), \quad (2)$$

where  $\text{Be}(a, b)$  denotes the beta density with mean  $a/(a+b)$ . The sequence of random weights  $\{\pi_{s,h}\}$  are constructed via the multiscale stick-breaking process described below. We will refer to the latter construction as multiscale Bernstein polynomial (msBP) model.

To build a multiscale stick-breaking process, an infinite sequence of scales  $s = 0, 1, \dots, \infty$  labelling the levels of an infinite-deep binary tree is introduced. At each scale  $s$  there will be  $2^s$  different nodes. A cartoon of the binary tree is reported in [Figure 1](#). To describe a stochastic path from the root node to the leaves, at each scale  $s$  and node  $h$  within the scale, the following independent random variables are introduced:

$$S_{s,h} \sim \text{Be}(1, a), \quad R_{s,h} \sim \text{Be}(b, b), \quad (3)$$

corresponding to the probability of stopping at node  $(s, h)$  and taking the right path after node  $(s, h)$  conditionally on not stopping, respectively. This formulation generalizes the stick-breaking process representation of [Sethuraman \(1994\)](#). Each time the stick is broken, it is then randomly divided in two parts (one for the probability of going right, the remainder for the probability of going left) before the next break. Hence, similarly to [Sethuraman \(1994\)](#), the infinite sequence of weights can be defined as

$$\pi_{s,h} = S_{s,h} \prod_{r < s} (1 - S_{r, g_{shr}}) T_{shr}, \quad (4)$$

where  $g_{shr} = \lceil h/2^{s-r} \rceil$  is the node traveled through at scale  $r$  on the way to node  $h$  at scale  $s$ ,  $T_{shr} = R_{r, g_{shr}}$  if  $(r+1, g_{shr+1})$  is the right daughter of node  $(r, g_{shr})$ , and  $T_{shr} = 1 - R_{r, g_{shr}}$  if  $(r+1, g_{shr+1})$  is the left daughter of  $(r, g_{shr})$ . Note that the general  $(s, h)$  node is related to the  $\text{Be}(h, 2^s - h + 1)$  density.

The above construction leads to a meaningful sequence of weights, i.e.,  $\sum_{s=0}^{\infty} \sum_{h=1}^{2^s} \pi_{s,h} = 1$  almost surely for any  $a, b > 0$  as proved in [Lemma 1 of Canale and Dunson \(2016\)](#). An appealing aspect of this formulation is that it produces a multiscale clustering of the subjects. In particular, two subjects having similar observations may have the same cluster allocation up to some scale  $s$ , but are not clustered together on finer scales.

## 2.2. Bayesian multiscale inference on group differences

A promising feature of this multiscale stick-breaking process is its ease of generalization to more complex settings than mere density estimation. For example, the sequence of random variables defined in [Equation 3](#) can be generalized to include predictors or other forms of dependence, (e.g., spatial or temporal). Motivated by epigenetic data, [Canale and Dunson \(2016\)](#) modified model (2)–(3), to perform Bayesian multiscale inference on group differences.

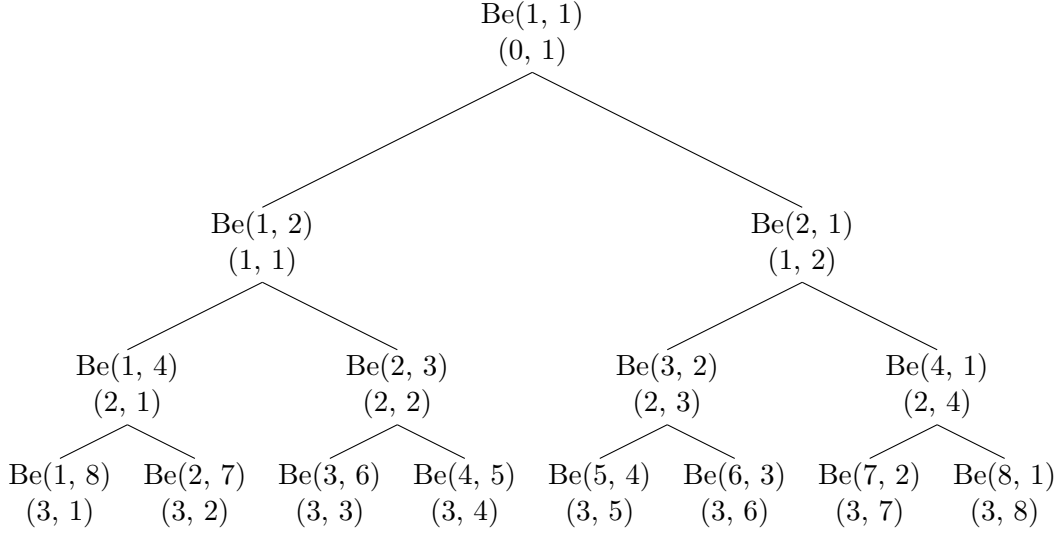


Figure 1: Binary tree with beta kernels at each node  $(s, h)$ , where  $s$  is the scale level and  $h$  is the index within the scale.

Let  $y_i$  be a bounded (between zero and one) outcome for subject  $i$  with  $y_i \sim f_{d_i}$  and  $d_i \in \{0, 1\}$ . The label  $d_i$  denotes a subject's group (e.g., cases/controls, drug/placebo). Using the msBP representation, the hypothesis  $f_0 = f_1$  is true if the groups share the same weights over the binary tree. If  $f_0 \neq f_1$ , we may have the same weights on the dictionary elements up to a given scale, so that the densities are equivalent up to that scale but not at finer scales. Thus, one can also test for  $H_0^s : f_0^s = f_1^s$ , i.e., no differences between the two groups at scale  $s$ . Clearly  $H_0^0$  is true with probability one, and thus a further modification of (3) consists to set  $S_{0,1} = 0$ .

The subjects surviving up to scale  $s$  can stop or progress to the next scale. Let  $\mathcal{N}^s$  denote these actions, with  $\mathcal{N}_{(d)}^s$  denoting the actions in group  $d$ . Conditionally on  $\mathcal{N}^s$  the posterior probability of  $H_0$  being true at scale  $s$  can be written as

$$\mathbb{P}(H_0^s | \mathcal{N}^s) = \frac{P_0^s \mathbb{P}(\mathcal{N}^s | H_0^s)}{P_0^s \mathbb{P}(\mathcal{N}^s | H_0^s) + (1 - P_0^s) \mathbb{P}(\mathcal{N}^s | H_1^s)}, \quad (5)$$

where  $P_0^s$  is our prior guess for the null being true at scale  $s$  and  $\mathbb{P}(\mathcal{N}^s | H_0^s)$  is the probability of the possible actions if  $H_0$  is true up to scale  $s$ . To compute the latter, we can use

$$\begin{aligned} \mathbb{P}(\mathcal{N}^s | H_0^s) &= \int_{\mathcal{T}} \mathbb{P}(\mathcal{N}^s | \mathcal{T}) \mathbb{P}(\mathcal{T} | a, b) d\mathcal{T} \\ &= \left\{ \frac{\Gamma(a+1) \Gamma(2b)}{\Gamma(a) \Gamma(b)^2} \right\}^{2^s} \int_{\mathcal{T}} \prod_{h=1}^{2^s} S_{s,h}^{n_{s,h}} (1 - S_{s,h})^{\hat{a}_{s,h}-1} R_{s,h}^{\hat{b}_{s,h}-1} (1 - R_{s,h})^{\hat{c}_{s,h}-1} d\mathcal{T} \\ &= \left\{ \frac{\Gamma(a+1) \Gamma(2b)}{\Gamma(a) \Gamma(b)^2} \right\}^{2^s} \prod_{h=1}^{2^s} \frac{\Gamma(1 + n_{s,h}) \Gamma(\hat{a})}{\Gamma(a + v_{s,h} + 1)} \frac{\Gamma(\hat{b}) \Gamma(\hat{c})}{\Gamma(2b + v_{s,h} - n_{s,h})}, \end{aligned} \quad (6)$$

where  $\hat{a}_{s,h} = a + v_{s,h} - n_{s,h}$ ,  $\hat{b}_{s,h} = b + r_{s,h}$ , and  $\hat{c}_{s,h} = b + v_{s,h} - n_{s,h} - r_{s,h}$ , and  $v_{s,h}$  is the number of subjects passing through node  $(s, h)$ ,  $n_{s,h}$  is the number of subjects stopping at node  $(s, h)$ , and  $r_{s,h}$  is the number of subjects that continue to the right after passing through

node  $(s, h)$ . Similarly

$$\begin{aligned} \mathbb{P}(\mathcal{N}^s | H_1^s) &= \mathbb{P}(\mathcal{N}_{(0)}^s | H_1^s) \times \mathbb{P}(\mathcal{N}_{(1)}^s | H_1^s) \\ &= \left\{ \frac{\Gamma(a+1)\Gamma(2b)}{\Gamma(a)\Gamma(b)^2} \right\}^{2^{2s}} \prod_{h=1}^{2^s} \frac{\Gamma(1+n_{s,h}^{(0)})\Gamma(\hat{a}^{(0)})}{\Gamma(a+v_{s,h}^{(0)}+1)} \frac{\Gamma(\hat{b}^{(0)})\Gamma(\hat{c}^{(0)})}{\Gamma(2b+v_{s,h}^{(0)}-n_{s,h}^{(0)})} \times \\ &\quad \prod_{h=1}^{2^s} \frac{\Gamma(1+n_{s,h}^{(1)})\Gamma(\hat{a}^{(1)})}{\Gamma(a+v_{s,h}^{(1)}+1)} \frac{\Gamma(\hat{b}^{(1)})\Gamma(\hat{c}^{(1)})}{\Gamma(2b+v_{s,h}^{(1)}-n_{s,h}^{(1)})}, \end{aligned} \quad (7)$$

where  $v_{s,h}^{(d)}$  is the number of subjects passing through node  $(s, h)$  in group  $d$ ,  $n_{s,h}^{(d)}$  is the number of subjects stopping at node  $(s, h)$  in group  $d$ , and  $r_{s,h}^{(d)}$  is the number of subjects that continue to the right after passing through node  $(s, h)$  in group  $d$ , with  $d = 0, 1$ . The global null will be the cumulative product of Equation 5 for each scale.

Motivated by a DNA methylation arrays application, [Canale and Dunson \(2016\)](#) generalized the latter formulation in the case in which  $y_i = (y_{i1}, \dots, y_{ip})^\top$ . To deal with  $p$ -dimensional arrays, a prior for  $P_0^s$  is assumed in order to borrow information across sites and to learn the joint null probability  $P_0^s$ . This feature is not yet implemented in the **msBP** package. A similar multiscale approach to perform two-sample comparison has been proposed and successfully applied in the multivariate context in [Ma and Wong \(2011\)](#) extending the optional Pólya tree prior of [Wong and Ma \(2010\)](#). The latter approach is able to jointly perform testing of two sample difference and learn the underlying structure of the difference. Another proposal connected to Pólya trees and dealing with more than two groups, censored, and multivariate data, is discussed in [Chen and Hanson \(2014\)](#); see also [Holmes, Caron, Griffin, and Stephens \(2015\)](#) for a related approach.

### 3. The C++ implementation

All the main functions of the **msBP** package are written in C++ and most of them rely on the ‘**bintree**’ data structure, i.e.,

```
struct bintree
{
    double data;
    struct bintree *left;
    struct bintree *right;
};
```

The ‘**bintree**’ structure is composed of a root (or parent node), each of which stores data and the two links to the leaves (or daughters nodes). Clearly each leaf connects to two other leaves and it is the beginning of a new, nested binary tree. A binary tree is a well known data structure with appealing characteristics in computer science. For example, it is possible to easily access and insert data into a binary tree using search and insert functions recursively called on successive leaves. This data structure will be used to store the random variables  $S_{s,h}$  and  $R_{s,h}$ , the weights in the mixtures, and other sample statistics. Basic functions to handle the ‘**bintree**’ data structure, such as create a tree, write and extract the data on a

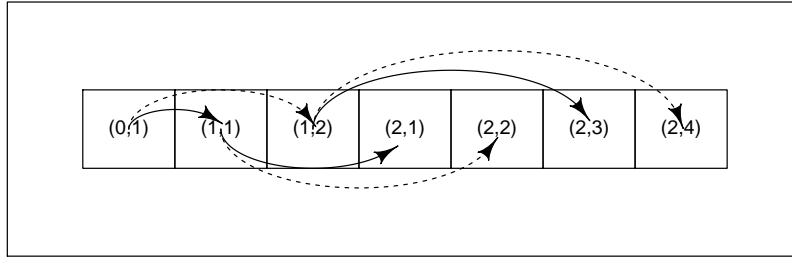


Figure 2: Behavior of the `tree2array` function. Arrows denote the branch of the original binary tree, with continuous line for the right daughter and dashed line for the left daughter. The number inside the array cells represent the original tree indexes.

given node of a tree and so forth, have also been implemented. Among them, the following functions

```
void tree2array(struct bintree *node, double *array, ...)
void array2tree(double *a, int maxScale, struct bintree *node)
```

allow for the conversion of a binary tree structure into an array and vice versa, and have been written to allow the input-output communication of R and C++ via the `.C` function. The first two arguments of the `tree2array` function are the pointers to the binary tree and to the array in which to write the values of the tree. Note that the array needs to be initialized before the use of `tree2array` and needs to have at least length  $2^s - 1$ , where  $s$  is the maximum depth of the tree. The `tree2array` function writes the array as described in Figure 2. The arguments of `array2tree`, instead, are the pointer to the array, an integer denoting the maximum scale of the binary tree, and the pointer to the binary tree structure to populate. In this latter case the binary tree structure only needs to be initialized and the function takes care of growing the tree up to the desired depth.

In addition to the basic functions already described, the **msBP** package also features more complex functions. However, most of them are then wrapped into R scripts and define the working functions of the package itself. Thus we do not further describe them here.

## 4. Usage of the **msBP** package

The main functions of **msBP** are `msBP.Gibbs`, which allows to perform nonparametric density estimation using the Gibbs sampler, and `msBP.test` which allows to perform Bayesian multi-scale testing of group differences. In this section of the article we provide examples of how to use these functions. In Section 4.1, basic and generic functions to handle the multiscale prior, to sample from an **msBP** process, and to plot the results, are described. Then, in Sections 4.2 and 4.3, the functions `msBP.Gibbs` and `msBP.test` are extensively discussed.

### 4.1. Basic and generic functions

The **msBP** package introduces two new R object classes implemented in S3. The first is the ‘`binaryTree`’ class. An object of the ‘`binaryTree`’ class represents a finite-depth binary tree.

It consists of a list containing `T` and `max.s`, the binary tree itself and an integer denoting its depth, respectively. Specifically, `T` is a list where each element is a vector containing the values of the nodes at a given scale. A binary tree of depth 3 containing the integers from 1 to 15 can be obtained with

```
R> tree <- structure(list(T = list(1, c(2, 3), c(4, 5, 6, 7),
+   c(8, 9, 10, 11, 12, 13, 14, 15)), max.s = 3), class = "binaryTree")
```

The tree structure can be converted into a vector using the `tree2vec` function

```
R> x <- tree2vec(tree)
```

while `vec2tree(x)` populates a binary tree with the values contained in the vector `x`. The latter function is ideally constructed for vectors of length  $2^n - 1$ , where  $n \in \mathbb{N}$ . However if the length  $l \neq 2^n - 1$  for any  $n$ , the function creates a tree up to scale  $\lceil \log_2(\lfloor l/2 \rfloor + 1) \rceil$  with the last leaves populated with `NA`. This object class will be largely used by other higher level functions, since the approach described in Section 2 deals with several binary trees such as Equations 3, 4 and so forth. A general `plot` function is available for the ‘`binaryTree`’ object class. The result of `plot(tree, ...)` is a cartoon of a binary tree with the root node at the top. As additional arguments, the function features: `value`, `size`, and `white`. If `value = TRUE` the numerical values of each node appear inside the node (up to the number of digits specified by `precision`); if `size = TRUE` the sizes of the nodes are proportional to their values; if `white = TRUE` the background color of the nodes is white, otherwise it is in color scale (default `gray.colors`). Figure 3 shows the output of some combinations.

The second S3 object class implemented in the `msBP` package is the ‘`msBPtree`’ class. An object of class ‘`msBPtree`’ is a list of 5 elements that represent a random draw from an `msBP(a, b)` process. The first two elements are the trees of the stopping and descending-to-the-right probabilities, described by Equation 3. Both are objects of the class ‘`binaryTree`’ with the same `max.s`. The third and fourth argument are the hyperparameters of the `msBP` prior, namely `a` and `b`. The last value is an integer with the maximum depth of both trees. To simulate a random density from an `msBP(a, b)` prior truncated at scale 3, the `msBP.rtree` function can be used as

```
R> set.seed(17012014)
R> draw <- msBP.rtree(a = 5, b = 1, max.s = 3)
```

Note that the last scale has  $S_{s,h} = 1$ . The induced trees of probabilities, calculated by means of (4) can be obtained with the `msBP.compute.prob` function as

```
R> weights <- msBP.compute.prob(draw)
```

and the results can be plotted using `plot`, as it is an object of class ‘`binaryTree`’. An additional argument `root = FALSE` sets  $S_{0,1} = 0$ . This can be used, for example, in the settings of Section 2.2. The induced random density can be drawn on a finite grid of length `n.points` of its domain with the function `msBP.pdf`, i.e.,

```
R> density <- msBP.pdf(weights, n.points = 100)
R> plot(dens ~ y, data = density, xlab = "y", ylab = "Density", type = "l")
```

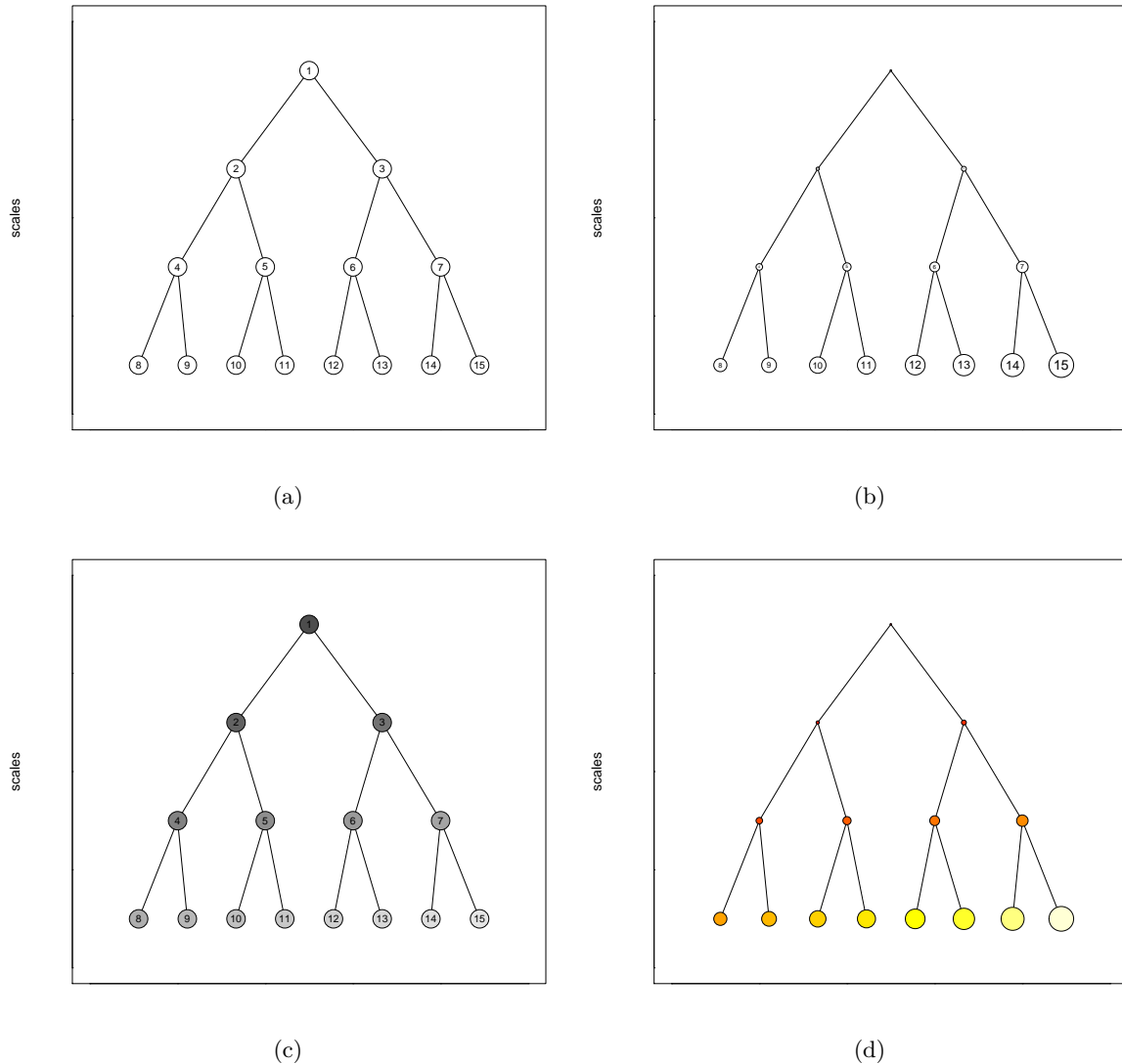


Figure 3: Output of the `plot(tree)` function with (a) default arguments values, (b) `size = TRUE`, (c) `white = FALSE`, and (d) `white = FALSE`, `size = TRUE`, `value = FALSE`, `col.grid = heat.colors(15)`.

Given a random density from an  $\text{msBP}(a, b)$  process, it is also possible to generate a sample of size  $n$  from that density. To this end we use the Algorithm 1, implemented in C++ and wrapped into R via the `msBP.rsample` function. The function `msBP.rsample` needs two parameters only: the sample size  $n$  and an object of class ‘`msBPTree`’.

## 4.2. Density estimation

Posterior density estimation under the `msBP` setup relies on a Markov chain Monte Carlo (MCMC) sampling algorithm. We briefly recall the algorithm proposed by Canale and Dunson (2016) in what follows. It basically consists of two steps: (i) multiscale cluster allocation



---

**Algorithm 1** Generating a random sample from a random density having an msBP prior.

---

```

for  $i = 1, \dots, n$  do
  loop = TRUE;
   $s_i = 0, h_i = 1$ ;
  while loop do
    let loop = FALSE with probability  $S_{s,h}$ .
    if loop then
      with probability  $R_{s_i, h_i}$ , let  $h_i = 2h_i$ .
      with probability  $1 - R_{s_i, h_i}$ , let  $h_i = 2h_i - 1$ .
    end if
  end while
  generate  $y_i \sim \text{Be}(h_i, 2^{s_i} - h_i + 1)$ .
end for

```

---

conditionally on the current values of the parameters  $\{\pi_{s,h}\}$ , and (ii) update of the probabilities parameters in the mixture, conditionally on the cluster allocation. Such a structure is typical of mixture models where first data augmentation allocates the observation to a mixture component and conditionally on such allocation the parameters of each component are updated (Bush and MacEachern 1996; MacEachern and Müller 1998; Ishwaran and James 2001).

Suppose that  $s_i$  and  $h_i$  are the scale and the node within scale labels for subject  $i$ . Conditionally on the binary tree of weights  $\{\pi_{s,h}\}$ , the posterior probability of subject  $i$  belonging to node  $(s, h)$  is simply

$$P(s_i = s, h_i = h | y_i, \pi_{s,h}) \propto \pi_{s,h} \text{Be}(y; h, 2^s - h + 1).$$

Now rewrite model (2) as

$$f(y) = \sum_{s=0}^{\infty} \pi_s \sum_{h=1}^{2^s} \bar{\pi}_{s,h} \text{Be}(y; h, 2^s - h + 1),$$

where  $\pi_s = \sum_{h=1}^{2^s} \pi_{s,h}$ , and  $\bar{\pi}_{s,h} = \pi_{s,h} / \pi_s$ . The cluster allocation uses a modification of the slice sampler of Kalli, Griffin, and Walker (2011) and is reported in Algorithm 2.

Algorithm 2 is implemented in the `msBP.postCluster` function. It requires two arguments: the sample of observations `y` and a binary tree of weights `weights`. The function makes a call to the `postCluster` C++ subroutine. The output of the function is a matrix with `length(y)` rows and two columns of cluster labels: the first denoting the scale and the second denoting the node within a given scale. The same C++ subroutine called by `msBP.postCluster` is called at each iteration of the MCMC sampler as described below.

Conditionally on the cluster allocations, the stopping and descending-right probabilities can be updated from their full conditional posteriors, namely:

$$S_{s,h} \sim \text{Be}(1 + n_{s,h}, a + v_{s,h} - n_{s,h}), \quad R_{s,h} \sim \text{Be}(b + r_{s,h}, b + v_{s,h} - n_{s,h} - r_{s,h}). \quad (8)$$

Calculation of  $v_{s,h}$  and  $r_{s,h}$  can be performed via the `msBP.nrvTrees` function, a wrapper calling the `allTree` C++ subroutine. The input of `msBP.nrvTrees` is the output of

---

**Algorithm 2** Multiscale cluster posterior allocation for  $i$ -th subject.

---

```

for each scale  $s$  do
  calculate  $\pi_s = \sum_{h=1}^{2^s} \pi_{s,h}$ .
end for
simulate  $u_i | y_i, s_i \sim U(0, \pi_{s_i})$ .
for each scale  $s$  do
  if  $\pi_s > u_i$  then
    for  $h = 1, \dots, 2^s$  do
      compute  $\bar{\pi}_{s,h} = \pi_{s,h} / \pi_s$ .
    end for
    compute  $P(s_i = s | u_i, y_i) \propto \sum_{h=1}^{2^s} \bar{\pi}_{s,h} \text{Be}(y_i; h, 2^s - h + 1)$ .
  else
     $P(s_i = s | u_i, y_i) = 0$ .
  end if
end for
sample  $s_i$  with probability  $P(s_i = s | u_i, y_i) \propto \mathbb{I}(s : \pi_s > u_i) \sum_{h=1}^{2^s} \bar{\pi}_{s,h} \text{Be}(y_i; h, 2^s - h + 1)$ .
sample  $h_i$  with probability  $P(h_i = h | y_i, s_i) \propto \bar{\pi}_{s_i,h} \text{Be}(y_i; h, 2^{s_i} - h + 1)$ .

```

---

`msBP.postCluster`, i.e., a matrix with 2 columns and number of rows equal to the sample size. The output of `msBP.nrvTrees` is a list containing tree objects of the class ‘`binaryTree`’. The main function implemented in the **msBP** package is the `msBP.Gibbs` function, performing the actual MCMC simulation from the posterior. The function basically iterates between cluster allocation, using the `postCluster` C++ subroutine and parameter updating, calculating first the elements  $n_{s,h}$ ,  $r_{s,h}$ , and  $v_{s,h}$  by means of the `allTree` C++ subroutine, and then using Equation 8. The Markov chain sampling is written in C++ but additional R language is used to initialize the function.

To describe the use of `msBP.Gibbs`, we will now walk the reader through the entire process of density estimation under the msBP setup. We will start showing how to elicit prior information, how to run the sampler, and how to analyze the output of the analysis. We do this using the famous Galaxy dataset (Roeder 1990). The dataset contains the velocity of 82 galaxies. The histogram of the speeds reveals that the data are clearly multimodal. This feature supports the Big Bang theory, as the different modes of density can be thought as clusters of galaxies moving at different speed.

```

R> data("galaxy", package = "DPpackage")
R> x <- galaxy$speed / 1000

```

We start by discussing prior elicitation. In Section 2.1 we assumed that  $g_0$  is our prior guess for  $g$ , the density of  $x$  and we want to center our msBP process in such a prior. We discussed that if  $G_0$  and  $G_0^{-1}$  are the corresponding CDF and inverse CDF, we can first transform the data with  $y = G_0(x) \in (0, 1)$ , and then estimate  $f \sim \text{msBP}(a, b)$ . It can be shown (see Canale and Dunson 2016) that the expectation  $\mathbb{E}\{F(A)\} = \lambda(A)$ , where  $\lambda(A)$  is the Lebesgue measure over the set  $A$  and  $F(A) = \int_A f$ . Since the prior for  $f$  is centered about the uniform, the prior on  $g$  is automatically centered in  $g_0$ . To allow this from a practical viewpoint we can use the argument `g0` of the `msBP.Gibbs` function. The package features four different prior guesses for  $g_0$ : `g0 = c("uniform", "normal", "gamma", "empirical")` for uniform,

normal, gamma and, following an empirical approach, the kernel density estimate. As default choice the function implements the "empirical" specification. For "normal" and "gamma", the parameters can be fixed or an additional prior distribution can be assumed. The former approach is adopted using `g0par` a vector of size two corresponding to mean and standard deviation of the normal, or shape and rate parameters for the gamma, respectively. The latter approach is adopted using `hyper$hyperpriors$g0 = TRUE`. In this case the model becomes

$$y = G_0(x; \theta), \quad \theta \sim P(\theta),$$

and thus an additional step of Gibbs sampling to simulate  $\theta$  is necessary. The full conditional posterior of  $\theta$  is simply

$$P(\theta | -) \propto P(\theta) \prod_{i=1}^n f(G_0(x_i); \theta) g_0(x_i; \theta),$$

and to sample from the latter full conditional posterior distribution the package uses a Metropolis-Hastings step (Hastings 1970) with proposal equal to the prior. Currently only `g0 = "normal"` is allowed with normal-inverse-gamma prior.

Then one has to specify the values of  $a$  and  $b$ , the hyperparameters of the msBP prior. The hyperparameter  $a$  controls the decline in probabilities over scales. Let  $S^{(i)}$  denotes the scale of  $i$ -th observation. It can be shown that  $E(S^{(i)}) = a$  which means that for small  $a$ , high probability is placed on coarse scales, leading to smoother densities and as  $a$  increases, finer scale densities will be weighted more, leading to spiker realizations. Additional hyperpriors for  $a$  and  $b$  can be assumed. Clearly, this will lead to additional sampling steps in the Gibbs sampling algorithm. In the `msBP.Gibbs` function this can be achieved by setting `hyper$hyperpriors$a = TRUE` and `hyper$hyperpriors$b = TRUE`, respectively. Specifically the algorithm implements  $a \sim \text{Ga}(\beta, \gamma)$  and  $b \sim \text{Ga}(\delta, \lambda)$ . This leads to the following conditional posterior distributions:

$$a| - \sim \text{Ga} \left( \beta + 2^{s'+1} - 1, \gamma - \sum_{s=0}^{s'} \sum_{h=1}^{2^s} \log(1 - S_{s,h}) \right), \quad (9)$$

and

$$P(b| -) \propto \frac{b^{\delta-1}}{B(b, b)^{2^{s'+1}-1}} \exp \left\{ b \left( \sum_{s=0}^{s'} \sum_{h=1}^{2^s} \log\{R_{s,h}(1 - R_{s,h})\} - \lambda \right) \right\}, \quad (10)$$

where  $s'$  is the maximum occupied scale and  $B(p, q)$  is the Beta function. To sample from the conditional posterior distribution of  $b$ , a Griddy-Gibbs approach over the grid defined by `hyper$hyperpar$gridB` is used (see Ritter and Tanner 1992). For sake of illustration we run and discuss `msBP.Gibbs` under the following different prior specifications:

```
R> hyper1 <- list(hyperprior = list(a = TRUE, b = TRUE, g0 = FALSE),
+   hyperpar = list(beta = 50, gamma = 5, delta = 10, lambda = 1,
+   gridB = seq(0, 20, length = 30)))
R> g0_1 <- "empirical"
R> hyper2 <- list(hyperprior = list(a = TRUE, b = TRUE, g0 = FALSE),
+   hyperpar = list(beta = 50, gamma = 5, delta = 10, lambda = 1,
+   gridB = seq(0, 20, length = 30)))
```

```
R> g0_2 <- "normal"
R> g0par_2 <- c(21, 2.5)
R> hyper2 <- list(hyperprior = list(a = TRUE, b = TRUE, g0 = TRUE),
+   hyperpar = list(beta = 50, gamma = 5, delta = 10, lambda = 1,
+     gridB = seq(0, 20, length = 30), mu0 = 21, kappa0 = 0.1,
+     alpha0 = 1, beta0 = 20))
R> g0_3 <- "normal"
R> g0par_3 <- c(21, 2.5)
```

which correspond to: (i)  $g_0$  is assumed to be equal to the empirical kernel density estimate, (ii)  $g_0$  is assumed to be normal with mean 21 and variance 2.5, and (iii)  $g_0$  is assumed to be normal with random mean and variance with prior  $(\mu, \sigma^2) \sim N(\mu, \mu_0, \kappa_0 \sigma^2) \text{I-Ga}(\sigma^2; \alpha_0, \beta_0)$ . In all cases the parameters of the msBP prior are assumed to be random with hyperprior distributions  $a \sim \text{Ga}(50, 5)$ , and  $b \sim \text{Ga}(10, 1)$ , with the prior for  $b$  evaluated on a grid from 0 to 20 of length 30.

The number of iterations to perform in the MCMC chain can be set via the function argument `mcmc`, a list including `nb`, the number of burn-in iterations, `nrep` the total number of iterations (including `nb`), and `ndisplay` the multiple of iterations to be displayed on screen while the C++ routine is running:

```
R> mcmc <- list(nrep = 10000, nb = 5000, ndisplay = 1000)
```

To obtain a posterior estimate of the density, the `grid` argument needs to be fixed. It consists of a named list giving the parameters for plotting the posterior mean density over a finite grid of points. It must include `low` and `upp` giving the lower and upper bound respectively of the grid and `n.points`, an integer giving the number of evaluation points.

```
R> grid <- list(n.points = 150, low = 5, upp = 38)
```

Additional arguments to be set are `maxS` and `printing`. The former is an upper bound for the binary trees involved in the MCMC sampling, and the latter is a control argument. If `printing = TRUE` the C++ routine prints on standard output what it is doing every `ndisplay` iterations. The default choice is `printing = FALSE`. With the following code we run the MCMC algorithm:

```
R> set.seed(17012014)
R> fit.msbp.1 <- msBP.Gibbs(speeds, a = 10, b = 10, g0 = g0_1,
+   mcmc = mcmc, hyper = hyper1, maxS = 5, grid = grid)
R> fit.msbp.2 <- msBP.Gibbs(speeds, a = 10, b = 10, g0 = g0_2,
+   g0par = g0par_2, mcmc = mcmc, hyper = hyper2, maxS = 5, grid = grid)
R> fit.msbp.3 <- msBP.Gibbs(speeds, a = 10, b = 10, g0 = g0_3,
+   g0par = g0par_3, mcmc = mcmc, hyper = hyper3, maxS = 5, grid = grid)
```

The function output is a named list containing four objects:

- `density`: a named list containing `postMeanDens`, the posterior mean density estimate evaluated over `xDens` and the related lower and upper pointwise 95% credible bands (`postLowDens` and `postUppDens`).

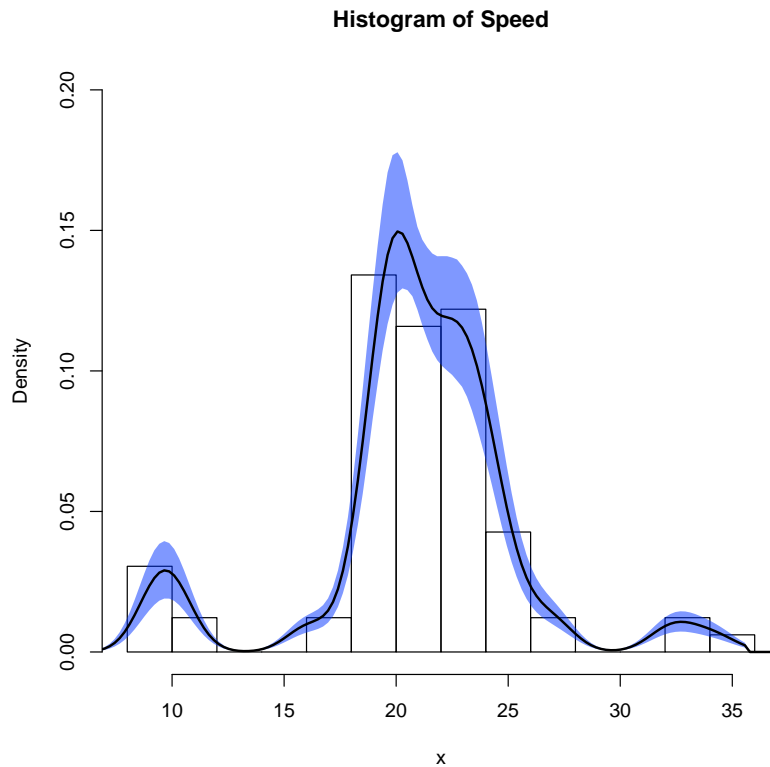


Figure 4: Posterior mean density for the Galaxy dataset.

- **mcmc**: a named list containing the MCMC chains: **dens** is a matrix  $(\text{nrep}-\text{nb}) \times \text{n.grid}$  containing the values of the density for each MCMC iteration, **a** and **b** are vectors containing the MCMC replicates for the two msBP parameters (if **hyperprior\$a** or **hyperprior\$b** are set to **TRUE**), **scale** is a matrix where each column is an MCMC sample of the total mass for each scale, **R** and **S** are matrices where each column is the **tree2vec** form of the corresponding trees, **weights** is a matrix where each column is the **tree2vec** form of the corresponding tree of weights, **s** and **h** are matrices where each column is the MCMC chain for the node labels for a subject, **mu** and **sigma** are vectors containing the MCMC replicates for the two parameters of the normal transformation of the data (if **hyper\$hyperprior\$g0** was set to **TRUE**).
- **postmean**: a named list containing posterior means over the MCMC samples of *a*, *b*, and of all binary trees. If **hyper\$hyperprior\$g0** was set to **TRUE**, the named list contains also the posterior means of the two parameters of the normal transformation of the data.
- **fit**: a named list containing the LPML, mean, and median of the log CPO (conditional predictive ordinate).

The histogram of the raw data and the plot of the posterior mean density and the related 95% credible bands for the first specification are reported in Figure 4.

To assess the convergence of the MCMC, one can visually inspect the traceplots of the chains for some parameter of interest. In general `fit.msbp.1$mcmc` contains the MCMC chains of

all the model's parameters. For example, if hyperpriors on the msBP prior parameters have been assumed, one can monitor the convergence of the chains for  $a$  and  $b$  with

```
R> plot(fit.msbp.1$mcmc$a, type = "l", main = "Traceplot for a", ylab = "")
R> plot(fit.msbp.1$mcmc$b, type = "l", main = "Traceplot for b", ylab = "")
```

and test for convergence using, for example, the [Geweke \(1992\)](#) diagnostics implemented in the `coda` package ([Plummer, Best, Cowles, and Vines 2006](#)), i.e.,

```
R> library("coda")
R> geweke.diag(fit.msbp.1$mcmc$a)
```

```
Fraction in 1st window = 0.1
Fraction in 2nd window = 0.5
```

```
var1
-0.1161
```

```
R> geweke.diag(fit.msbp.1$mcmc$b)
```

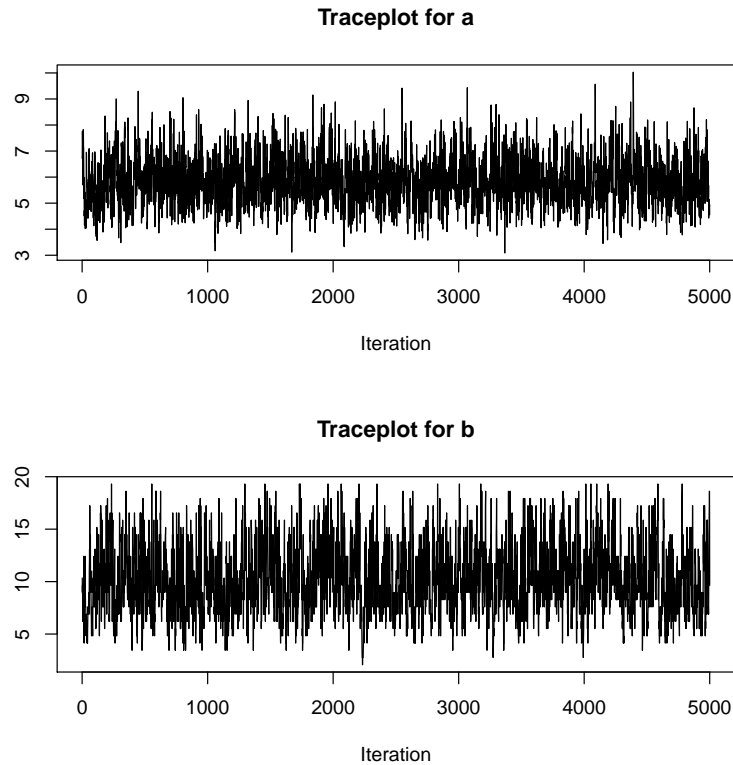
```
Fraction in 1st window = 0.1
Fraction in 2nd window = 0.5
```

```
var1
-1.769
```

We finally compare the fit obtained with the different prior specifications with the fit obtained running `DPdensity` and `PTdensity` from package **DPpackage**. As prior specification for the latter models we rely on the specifications described in the documentation of the package, reported for sake of completeness in the online supplements of this article. The comparison is based on the log pseudo-marginal likelihood (LPML) criterion. The LPML is a leave-one-out cross validatory measure based on the predictive densities, see [Green and Richardson \(2001\)](#) and [Gelfand and Dey \(1994\)](#) for details. The three msBP specifications have a LPML of  $-215$ ,  $-298$ , and  $-265$ , respectively. The best performance is obtained using the first prior specification which centers the prior expected density in the kernel density estimate of the raw data. The latter is practically equivalent to the best fits obtained with `DPdensity` and `PTdensity`, where the LPML is equal to  $-210$  and to  $-216$ , respectively.

### 4.3. Inference in group differences

The function `msBP.test` performs multiscale hypothesis testing of difference in the distribution of two groups. It exploits the closed form expression for the conditional posterior probability for  $H_0^s$  in Equation 5. However, since it cannot be directly used due to the dependence on the unknown multiscale clustering structure, the function relies on a Gibbs sampling algorithm. Again, the algorithm is made of two steps: multiscale cluster allocation, and update of the tree of weights. For node  $h$  at scale  $s$ , let  $\pi_{s,h}^{(0)}$  denote the weight under  $H_0^s$  and

Figure 5: Posterior draw for  $a$  and  $b$ .

$\pi_{s,h}^{(1,d)}$  for  $d = 0, 1$  denote the group-specific weights under  $H_1^s$ . The allocation of subject  $i$ , at each iteration, will be made via `msBP.postCluster` using the following set of weights:

$$\pi_{s,h}^{(d_i)} = P(H_0^s | \mathcal{N}_{(0)}^s, \mathcal{N}_{(1)}^s) \pi_{s,h}^{(0)} + \{1 - P(H_0^s | \mathcal{N}_{(0)}^s, \mathcal{N}_{(1)}^s)\} \pi_{s,h}^{(1,d_i)}. \quad (11)$$

Then, at a given iteration the quantities in Equations 6–7 can be calculated explicitly, and used to update the stopping and descending probabilities.

We describe the parameters and the behavior of the function via the Indian Liver dataset, available at the UCI Machine Learning repository (Bache and Lichman 2013). This dataset contains data on 580 subjects where 413 are liver patients and 167 are non-liver patients. Subjects with liver problems typically register higher levels of bilirubin in their blood and thus we want to test if there is a difference in the distribution of the relative direct bilirubin, calculated as the ratio of the direct bilirubin over the total bilirubin. An histogram of the raw data is reported in Figure 6(a).

The `msBP.test` function, in addition to a vector of observations and to a vector of group labels, requires prior values for  $a$ ,  $b$ , and for the probability of  $H_0$ . The choice of the hyperparameters  $a$  and  $b$  can be made using prior information. In what follows, however, the choice is done with a two-step procedure. First, the density of the pooled dataset is fitted with the `msBP.Gibbs` function assuming hyperpriors for  $a$  and  $b$

```
R> mcmc.test <- list(nrep = 8000, nb = 4000, ndisplay = 1000)
R> hyper.test <- list(hyperprior = list(a = TRUE, b = TRUE),
+   hyperpar = list(beta = 5, gamma = 0.5, delta = 1, lambda = 1))
```

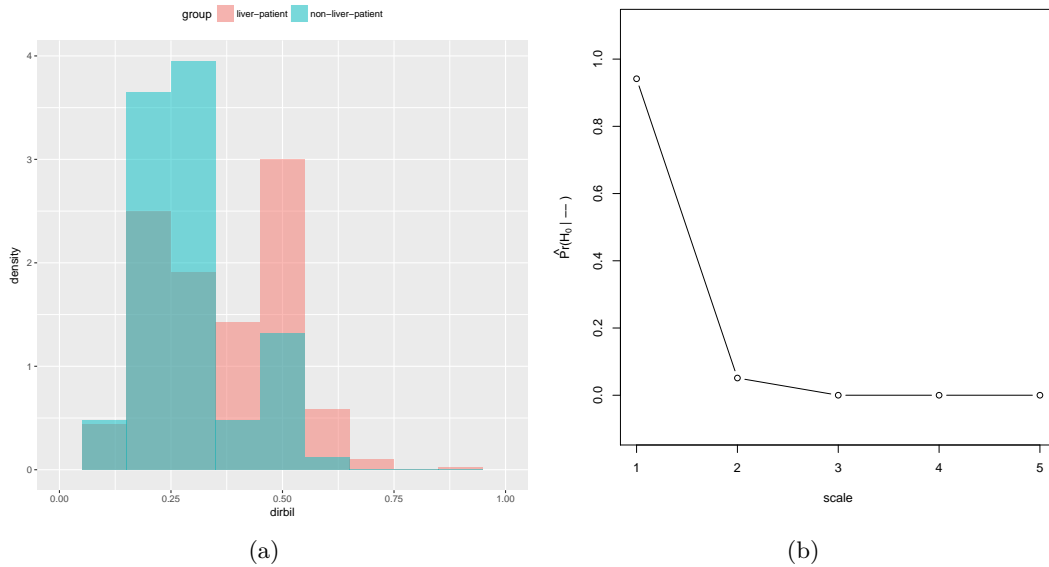


Figure 6: Histogram of the raw data (a) and posterior mean of  $P(H_{0s} | -)$  as function of  $s$  for the Indian liver dataset.

```
R> set.seed(17012014)
R> dens.liver <- msBP.Gibbs(liver$dirbil, a = 10, b = 10, g0 = "unif",
+   mcmc = mcmc.test, hyper = hyper.test, maxScale = 5)
```

Then the posterior mean of  $a$  and  $b$  are used as plug-in estimates for the testing. We fix the prior probability of  $H_0$  to 0.5 in order to put equal weights to the null and the alternative, and we fix the upper bound for the scales to 5. The function can be executed via

```
R> test.liver <- msBP.test(liver$dirbil, a = dens.liver$postmean$a,
+   b = dens.liver$postmean$b, group = liver$group, priorH0 = 0.5,
+   mcmc = mcmc.test, plot.it = TRUE, maxScale = 5)
```

The function's output is a list containing all the MCMC replicates for  $P(H_0^s | -)$  along with their posterior means and the global Bayes factor

$$\text{BF} = \frac{P(H_0 | -)}{P(H_1 | -)}.$$

Figure 6(b) reports the posterior mean of the global null hypothesis, as function of the scale. The differences between the two groups are minimal at the first scale but start to become evident for increasing scales.

## 5. Conclusions

We have presented a detailed introduction to the R package **msBP**, which implements a recently introduced multiscale stick-breaking prior and allows to perform density estimation and to test for differences in the distribution of two groups. The package implements also basic and generic functions to handle the involved multiscale trees structures.



## Acknowledgements

The author thanks David Dunson and Giovanna Menardi for comments on early versions of the manuscript. Comments on the package implementation by Roberto Vigo, Marco Pischedda, Brian Ripley, and the R-package-devel mailing list are gratefully acknowledged. This work has been conducted while the author was affiliated to the University of Turin.

## References

- Bache K, Lichman M (2013). “UCI Machine Learning Repository.” URL <http://archive.ics.uci.edu/ml>.
- Bush CA, MacEachern SN (1996). “A Semiparametric Bayesian Model for Randomised Block Designs.” *Biometrika*, **83**(2), 275–285. doi:10.1093/biomet/83.2.275.
- Canale A (2017). *msBP: Multiscale Bernstein Polynomials for Densities*. R package version 1.3, URL <https://CRAN.R-project.org/package=msBP>.
- Canale A, Dunson D (2016). “Multiscale Bernstein Polynomial for Densities.” *Statistica Sinica*, **26**(3), 1175–1195. doi:10.5705/ss.202015.0163.
- Chen Y, Hanson T (2014). “Bayesian Nonparametric  $k$ -Sample Tests for Censored and Uncensored Data.” *Computational Statistics & Data Analysis*, **71**, 335–346. doi:10.1016/j.csda.2012.11.003.
- Chen Y, Hanson T, Zhang J (2014). “Accelerated Hazards Model Based on Parametric Families Generalized with Bernstein Polynomials.” *Biometrics*, **70**(1), 192–201. doi:10.1111/biom.12104.
- Donoho DL, Johnstone IM, Kerkycharian G, Picard D (1996). “Density Estimation by Wavelet Thresholding.” *The Annals of Statistics*, **24**(2), 508–539. doi:10.1214/aos/1032894451.
- Escobar MD, West M (1995). “Bayesian Density Estimation and Inference Using Mixtures.” *Journal of the American Statistical Association*, **90**(430), 577–588. doi:10.1080/01621459.1995.10476550.
- Ferguson TS (1973). “A Bayesian Analysis of Some Nonparametric Problems.” *The Annals of Statistics*, **1**(2), 209–230. doi:10.1214/aos/1176342360.
- Ferguson TS (1974). “Prior Distributions on Spaces of Probability Measures.” *The Annals of Statistics*, **2**(4), 615–629. doi:10.1214/aos/1176342752.
- Gelfand AE, Dey DK (1994). “Bayesian Model Choice: Asymptotics and Exact Calculations.” *Journal of the Royal Statistical Society B*, **56**(3), 501–514.
- Geweke J (1992). “Evaluating the Accuracy of Sampling-Based Approaches to the Calculation of Posterior Moments.” In JM Bernardo, JO Berger, AP Dawid, AFM Smith (eds.), *Bayesian Statistics 4*, pp. 169–188. Oxford University Press, Oxford.

- Green P, Richardson S (2001). “Modelling Heterogeneity with and without the Dirichlet Process.” *Scandinavian Journal of Statistics*, **28**(2), 355–375. doi:10.1111/1467-9469.00242.
- Hanson T, Johnson WO (2002). “Modeling Regression Error with a Mixture of Polya Trees.” *Journal of the American Statistical Association*, **97**(460), 1020–1033. doi:10.1198/016214502388618843.
- Hastings WK (1970). “Monte Carlo Sampling Methods Using Markov Chains and Applications.” *Biometrika*, **57**(1), 97–109. doi:10.2307/2334940.
- Holmes CC, Caron C, Griffin JE, Stephens DA (2015). “Two-Sample Bayesian Nonparametric Hypothesis Testing.” *Bayesian Analysis*, **10**(2), 297–320. doi:10.1214/14-ba914.
- Ishwaran H, James LF (2001). “Gibbs Sampling Methods for Stick Breaking Priors.” *Journal of the American Statistical Association*, **96**(453), 161–173. doi:10.1198/016214501750332758.
- Jara A, Hanson T, Quintana FA, Müller P, Rosner GL (2011). “**DPpackage**: Bayesian Semi-And Nonparametric Modeling in R.” *Journal of Statistical Software*, **40**(5), 1–30. doi:10.18637/jss.v040.i05.
- Kalli M, Griffin J, Walker SG (2011). “Slice Sampling Mixture Models.” *Statistics and Computing*, **21**(1), 93–105. doi:10.1007/s11222-009-9150-y.
- Lavine M (1992a). “Some Aspects of Polya Tree Distributions for Statistical Modelling.” *The Annals of Statistics*, **20**(3), 1222–1235. doi:10.1214/aos/1176348767.
- Lavine M (1992b). “More Aspects of Polya Tree Distributions for Statistical Modelling.” *The Annals of Statistics*, **22**(3), 1161–1176. doi:10.1214/aos/1176325623.
- Lo AY (1984). “On a Class of Bayesian Nonparametric Estimates: I. Density Estimates.” *The Annals of Statistics*, **12**(1), 351–357. doi:10.1214/aos/1176346412.
- Ma L, Wong WH (2011). “Coupling Optional Pólya Trees and the Two Sample Problem.” *Journal of the American Statistical Association*, **106**(496), 1553–1565. doi:10.1198/jasa.2011.tm10003.
- MacEachern SN, Müller P (1998). “Estimating Mixture of Dirichlet Process Models.” *Journal of Computational and Graphical Statistics*, **7**(2), 223–238. doi:10.1080/10618600.1998.10474772.
- Mauldin D, Sudderth WD, Williams SC (1992). “Polya Trees and Random Distributions.” *The Annals of Statistics*, **20**(3), 1203–1203. doi:10.1214/aos/1176348766.
- Petrone S (1999a). “Bayesian Density Estimation Using Bernstein Polynomials.” *Canadian Journal of Statistics*, **27**(1), 105–126. doi:10.2307/3315494.
- Petrone S (1999b). “Random Bernstein Polynomials.” *Scandinavian Journal of Statistics*, **26**(3), 373–393. doi:10.1111/1467-9469.00155.

- Plummer M, Best N, Cowles K, Vines K (2006). “**coda**: Convergence Diagnosis and Output Analysis for MCMC.” *R News*, **6**(1), 7–11. URL [https://www.R-project.org/doc/Rnews/Rnews\\_2006-1.pdf](https://www.R-project.org/doc/Rnews/Rnews_2006-1.pdf).
- Ritter C, Tanner MA (1992). “Facilitating the Gibbs Sampler: The Gibbs Stopper and the Griddy-Gibbs Sampler.” *Journal of the American Statistical Association*, **87**(419), 861–868. doi:10.1080/01621459.1992.10475289.
- Roeder K (1990). “Density Estimation with Confidence Sets Emplified by Superclusters and Voids in Galaxies.” *Journal of the American Statistical Association*, **85**(411), 617–624. doi:10.1080/01621459.1990.10474918.
- Sethuraman J (1994). “A Constructive Definition of Dirichlet Priors.” *Statistica Sinica*, **4**(2), 639–650. doi:10.5705/ss.2012.047.
- Wang Y, Canale A, Dunson DB (2016). “Scalable Geometric Density Estimation.” In A Gretton, CC Robert (eds.), *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics*, volume 51 of *Proceedings of Machine Learning Research*, pp. 857–865. PMLR, Cadiz, Spain. URL <http://proceedings.mlr.press/v51/wang16e.html>.
- Wong WH, Ma L (2010). “Optional Pólya Tree and Bayesian Inference.” *The Annals of Statistics*, **38**(3), 1433–1459. doi:10.1214/09-aos755.

**Affiliation:**

Antonio Canale  
Department of Statistical Sciences  
University of Padua and  
Collegio Carlo Alberto  
E-mail: [canale@stat.unipd.it](mailto:canale@stat.unipd.it)  
URL: [homes.stat.unipd.it/antoniocanale](http://homes.stat.unipd.it/antoniocanale)